

Ce notebook s'agit d'une ACP pour swh

Charger les données :

on change l'index pour qu'il soit les noms puis on supprime cette variable:

```
swh<-read.csv('data_acp_time-series-swh.csv',header = T)
rownames(swh)<-swh$index
swh$index<-NULL
#str(swh) # donne 518 obs et 288 variables
```

Charger les libraires :

```
library("FactoMineR")

## Warning: package 'FactoMineR' was built under R version 3.4.2

library("factoextra")

## Warning: package 'factoextra' was built under R version 3.4.2

## Loading required package: ggplot2

## Welcome! Related Books: `Practical Guide To Cluster Analysis in R` at
https://goo.gl/13EFCZ

library("qtlcharts")

## Warning: package 'qtlcharts' was built under R version 3.4.2
```

Corrélation entre les variables : on choisit que les variables numériques

```
cor(swh)[1:5,1:5]

##           X1993_1 X1993_10 X1993_11 X1993_12 X1993_2
## X1993_1  1.0000000 0.6509012 0.8260983 0.7950397 0.8057792
## X1993_10 0.6509012 1.0000000 0.7245957 0.6597373 0.6493819
## X1993_11 0.8260983 0.7245957 1.0000000 0.7944981 0.8159978
## X1993_12 0.7950397 0.6597373 0.7944981 1.0000000 0.6464883
## X1993_2  0.8057792 0.6493819 0.8159978 0.6464883 1.0000000
```

La table de corrélation montre l'existence de plusieurs variables qui sont corrélées entre eux. Ceci est une bonne nouvelle dans le sens où l'ACP nécessite que les variables soient corrélées afin d'extraire de l'information contenue dans l'inertie totale.

On peut donc commencer l'analyse en composantes principales.

ACP

on va utiliser toutes les variables quantitatives (288 variables) pour effectuer notre ACP. Ainsi que tous les individus (518 individus).

```
res.pca=PCA(swh,graph=FALSE,scale.unit = FALSE)
print(res.pca)

## **Results for the Principal Component Analysis (PCA)**
## The analysis was performed on 518 individuals, described by 288 variables
## *The results are available in the following objects:
##
##      name                description
## 1  "$eig"                "eigenvalues"
## 2  "$var"                "results for the variables"
## 3  "$var$coord"          "coord. for the variables"
## 4  "$var$cor"            "correlations variables - dimensions"
## 5  "$var$cos2"           "cos2 for the variables"
## 6  "$var$contrib"        "contributions of the variables"
## 7  "$ind"                "results for the individuals"
## 8  "$ind$coord"          "coord. for the individuals"
## 9  "$ind$cos2"           "cos2 for the individuals"
## 10 "$ind$contrib"        "contributions of the individuals"
## 11 "$call"               "summary statistics"
## 12 "$call$centre"        "mean of the variables"
## 13 "$call$ecart.type"    "standard error of the variables"
## 14 "$call$row.w"         "weights for the individuals"
## 15 "$call$col.w"         "weights for the variables"
```

res.pca est un objet qui contient plusieurs variables à analyser. Dans ce qui suit on va analyser chaque attribut de cet objet.

```
eig.val<-get_eigenvalue(res.pca)
head(eig.val)

##      eigenvalue variance.percent cumulative.variance.percent
## Dim.1 197.316044      77.4740455      77.47405
## Dim.2  15.351697       6.0276806      83.50173
## Dim.3   4.153028       1.6306422      85.13237
## Dim.4   2.403778       0.9438177      86.07619
## Dim.5   1.577190       0.6192667      86.69545
## Dim.6   1.454910       0.5712548      87.26671
```

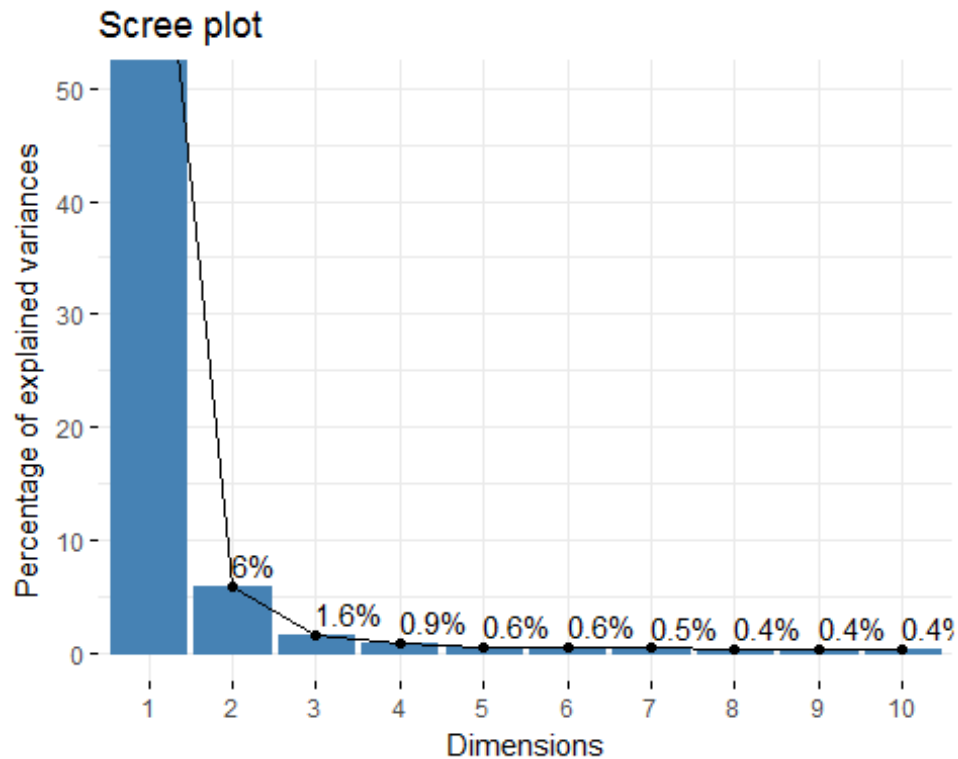
D'après les résultats ci-dessus, on peut conserver 2 axes principales vu qu'ils expliquent 83.50% de l'inertie totale contenue dans notre jeu de données. A noter qu'ici on n'a pas le droit d'utiliser le critère de Kaïser pour choisir le nombre d'axe à conserver, car les données ne sont pas normalisées. On préfère la normalisation dans le cas où les variables sont hétérogènes.

```
=====
=====
```

Critère de Kaiser (1961) : est utilisé pour déterminer le nombre d'axes principaux à garder après l'ACP. Une valeur propre > 1 \Rightarrow la CP en question représente plus de variance par rapport à une seule variable d'origine. Ceci est valable que si les données sont normalisées.

Afin de bien justifier notre choix on peut fournir les diagrammes suivants :

```
fviz_eig(res.pca, addlabels=T, ylim=c(0,50))
```



On peut, à partir du graph, se limiter à 2 composantes principales, soit 83.50% de la variance totale. Après ces deux axes, la variance cumulée ne change pas beaucoup.

Etude des variables :

Corrélation avec les CPs, Qualité de représentation, contributions aux CPs et cercles de corrélations.

Decription des dimensions :

dimdesc: est utilisée pour identifier les variables les plus significativement associées avec une CP donnée.

```
res.desc<-dimdesc(res.pca, axes=c(1,2), proba=0.05)  
print('====pour Dim.1====')
```

```
## [1] "====pour Dim.1===="
```

```

res.desc$Dim.1$quanti[1:5,]

##           correlation      p.value
## X2011_2      0.9594142 7.803638e-286
## X2003_3      0.9587245 5.512818e-284
## X2002_3      0.9532591 2.311457e-270
## X2011_11     0.9521820 7.167916e-268
## X2012_11     0.9490431 6.313158e-261

print('====pour Dim.2====')

## [1] "====pour Dim.2===="

res.desc$Dim.2$quanti[1:5,]

##           correlation      p.value
## X1998_4      0.5895354 8.436823e-50
## X2014_2      0.5733113 1.355737e-46
## X1996_1      0.5009079 2.972598e-34
## X2001_3      0.5000036 4.066459e-34
## X1997_11     0.4935803 3.665965e-33

tail(res.desc$Dim.2$quanti)

##           correlation      p.value
## X1999_2     -0.4742463 2.086747e-30
## X2007_11    -0.4961653 1.521613e-33
## X2012_2     -0.5246527 5.627249e-38
## X2004_11    -0.5320974 3.321372e-39
## X2000_3     -0.5360614 7.150779e-40
## X2001_11    -0.5438540 3.289728e-41

```

Conclusion : - Les variables les plus corrélées avec la Dim.1 (CP1) sont les mois suivantes :2011_2,2003_3,2002_3,2011_11,2012_11,2005_10,2013_12,2015_3,2002_1,2003_2==> Dim1 peut être interprété comme étant la moyenne de swh sur ces mois - Les variables les plus corrélées positivement avec Dim.2 (CP2) sont les mois suivantes : 1998_4,2014_2,1996_1 et 2001_3. - Les variables les plus corrélées négativement avec Dim.2 sont les mois suivantes : 2007_11,2012_2,2004_11,2000_3 et 2001_11

Rq: ce n'est pas une vraie moyenne (on verra pourquoi à l'aide des résultats suivants).

a.2 information générales sur les variables :

```

var<- get_pca_var(res.pca)
print('===== Coordonnées des dim :=====')

## [1] "===== Coordonnées des dim :===== "

head(var$coord)

##           Dim.1      Dim.2      Dim.3      Dim.4      Dim.5
## X1993_1  1.5639411 -0.4782969 -0.2993860  0.25021321 -0.08612173
## X1993_10 0.6515051  0.1053498  0.2658732  0.16224008  0.10395235

```

```
## X1993_11 1.1682282 -0.1425860 0.1300491 -0.11395730 -0.07404538
## X1993_12 1.2664339 0.0870997 -0.3859654 -0.01920919 0.22685773
## X1993_2 1.0093550 -0.6001314 0.1982075 0.08112487 0.06062967
## X1993_3 1.1173607 -0.2795167 0.1223889 -0.01174728 -0.01571461

print('===== Qualité de représentation : =====')

## [1] "===== Qualité de représentation : ====="

head(var$cos2)

##           Dim.1      Dim.2      Dim.3      Dim.4      Dim.5
## X1993_1  0.7800245 0.072956259 0.02858448 1.996585e-02 0.0023653363
## X1993_10 0.6368387 0.016651788 0.10605786 3.949209e-02 0.0162129639
## X1993_11 0.8635915 0.012864910 0.01070207 8.217454e-03 0.0034693561
## X1993_12 0.7838553 0.003707696 0.07280606 1.803388e-04 0.0251523183
## X1993_2  0.6354128 0.224626390 0.02450239 4.104649e-03 0.0022926538
## X1993_3  0.8516311 0.053294231 0.01021759 9.413244e-05 0.0001684504

print('===== Contributions des variables :=====')

## [1] "===== Contributions des variables :===== "

head(var$contrib)

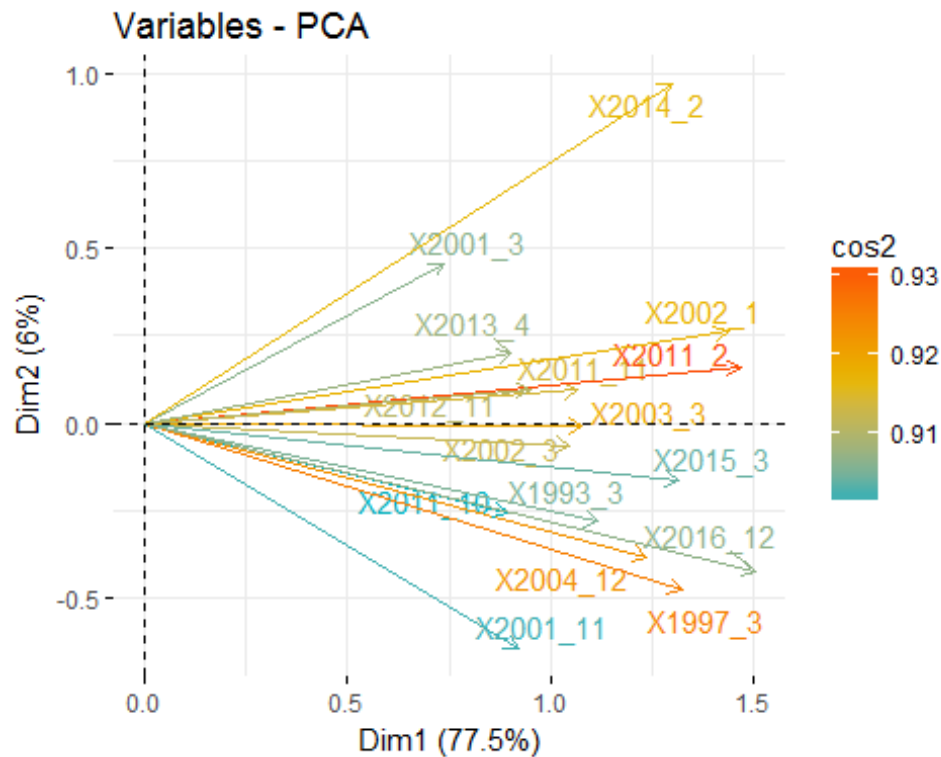
##           Dim.1      Dim.2      Dim.3      Dim.4      Dim.5
## X1993_1  1.2395910 1.49018003 2.1582318 2.604510932 0.47026385
## X1993_10 0.2151163 0.07229541 1.7020968 1.095020006 0.68514857
## X1993_11 0.6916605 0.13243336 0.4072398 0.540244092 0.34762582
## X1993_12 0.8128355 0.04941706 3.5870038 0.015350550 3.26304649
## X1993_2  0.5163278 2.34604467 0.9459659 0.273787558 0.23307007
## X1993_3  0.6327387 0.50893125 0.3606776 0.005740905 0.01565753
```

Quelques remarques à partir des tableaux ci-dessus : Dim.1= $1.56X_{1993_1} + 0.65X_{1993_10}$ +... (voilà pourquoi ce n'est pas une moyenne) de même pour les autres.

a.3 : Les cercles de corrélations :

ci-dessous on fournit un cercle avec un gradient de couleurs.

```
fviz_pca_var(res.pca,
col.var="cos2",gradient.cols=c("#00AFBB","#E7B800","#FC4E07"),repel=T,select.
var = list(cos2=0.90))
```



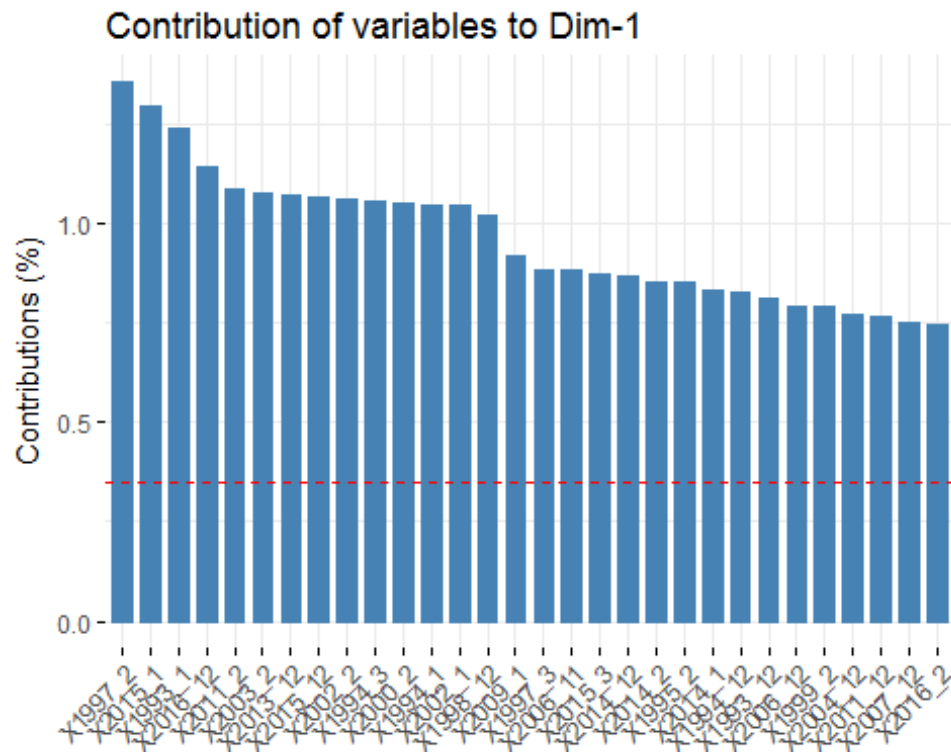
->on représente les variables par leurs corrélation sur un cercle : -> Les variables positivement corrélées sont regroupé par quart de cercle -> Les variables négativement corrélées sont positionnées sur les côtés opposés de l'origine du cercle(quadrant opposés)
-> La distance entre les variables et l'origine mesure la qualité de représentation des variables. -> Et donc les variables les plus loin de l'origine sont les bien représentés par l'ACP

Rq: Comme représenté sur la figure ci-dessus, les mois les plus bien représenté sont : 2011_2,1997_3,2004_12,2002_1 et 2014_2.

a.4 : Contributions des variables aux axes principaux

Pour la dim.1 :

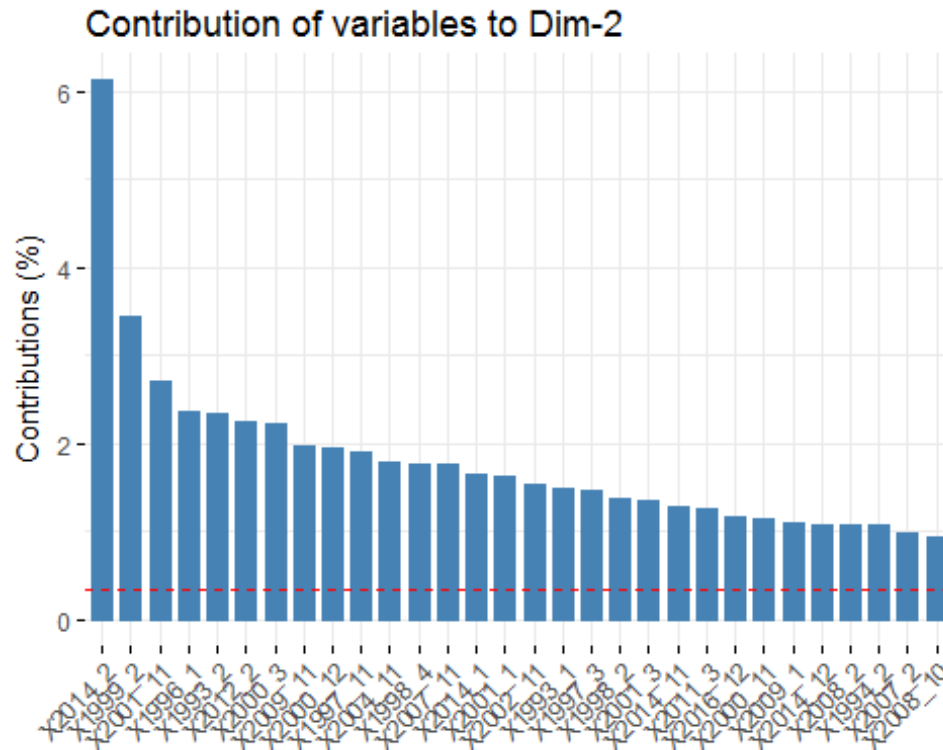
```
fviz_contrib(res.pca,choice="var", axes=1,top=30)
```



On voit que les mois 1997_2, 2015_1, 1993_1, 2011_2 contribuent le plus à la dimension 1.

Pour la dim.2 :

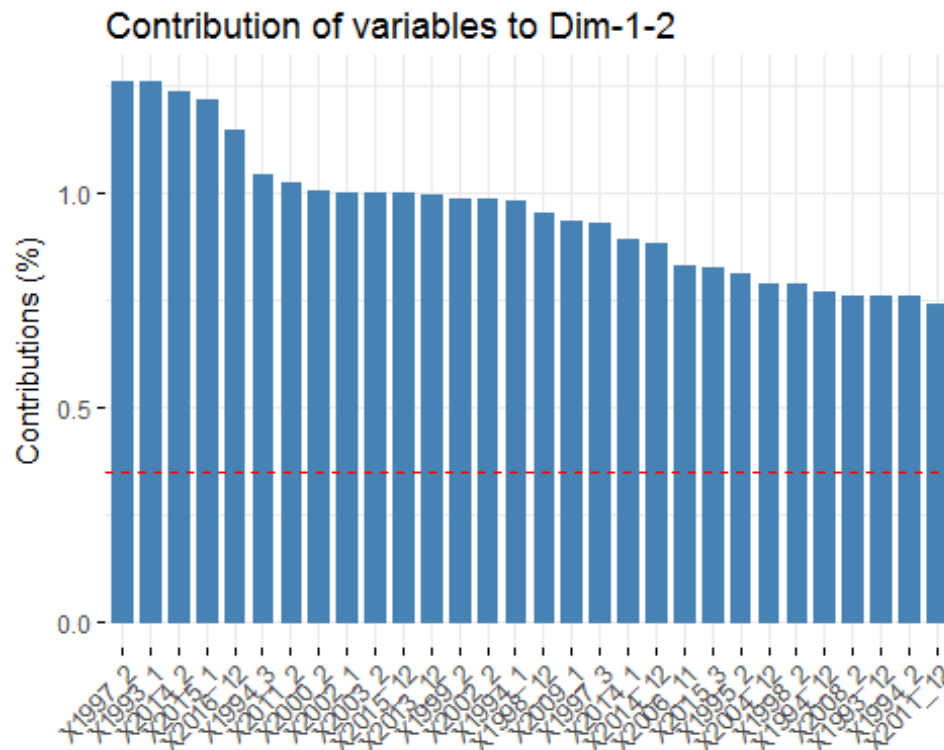
```
fviz_contrib(res.pca,choice="var", axes=2,top=30)
```



On voit que 2014_2, 1999_2, 2001_11, 1996_1 et 1993_2 sont les mois les plus contributifs à la dimension 2. RQ : La ligne en pointillé rouge indique la contribution en moyenne attendue. si une variable dépasse ce seuil ==> elle est importante pour contribuer à la composante.

Notez que la contribution totale à PC1 et PC2 peut être obtenue avec le code R suivant:

```
fviz_contrib(res.pca, choice="var", axes=1:2, top=30)
```

Généralement, les mois les plus contributifs sont les mois de l'hiver. ###b. Etude des individus :

```
ind<-get_pca_ind(res.pca)
print('==== Coordonnées des individus :====')

## [1] "==== Coordonnées des individus :===="

head(ind$coord)

##           Dim.1    Dim.2    Dim.3    Dim.4    Dim.5
## -10_38 -3.8635902  4.679168  2.794974 -1.02209329 -1.7022777
## -10_39 -2.9554985  5.015904  2.481558 -1.42527414  0.9583765
## -10_40 -1.6657553  5.229145  2.115656 -1.27771384  1.6952736
## -10_41 -0.5442595  5.382409  2.182068 -0.61515267  0.7080205
## -10_42  0.9791744  4.980011  1.812504  0.08561013 -0.3201328
## -10_43  2.9609207  4.944861  1.715110  1.14779967  1.5248697

print('==== Qualité de représentation des individus:====')

## [1] "==== Qualité de représentation des individus:===="

head(ind$cos2)

##           Dim.1    Dim.2    Dim.3    Dim.4    Dim.5
## -10_38 0.187994444  0.2757403  0.09838263  0.0131566090  0.036494187
## -10_39 0.119823518  0.3451270  0.08447522  0.0278661762  0.012599487
## -10_40 0.039284839  0.3871363  0.06337131  0.0231137158  0.040689482
## -10_41 0.003813061  0.3729194  0.06129121  0.0048711044  0.006452875
```

```
## -10_42 0.014617827 0.3781145 0.05008651 0.0001117411 0.001562511
## -10_43 0.117861315 0.3287198 0.03954594 0.0177112784 0.031259583

print('==== Contributions des individus :====')

## [1] "==== Contributions des individus :===="

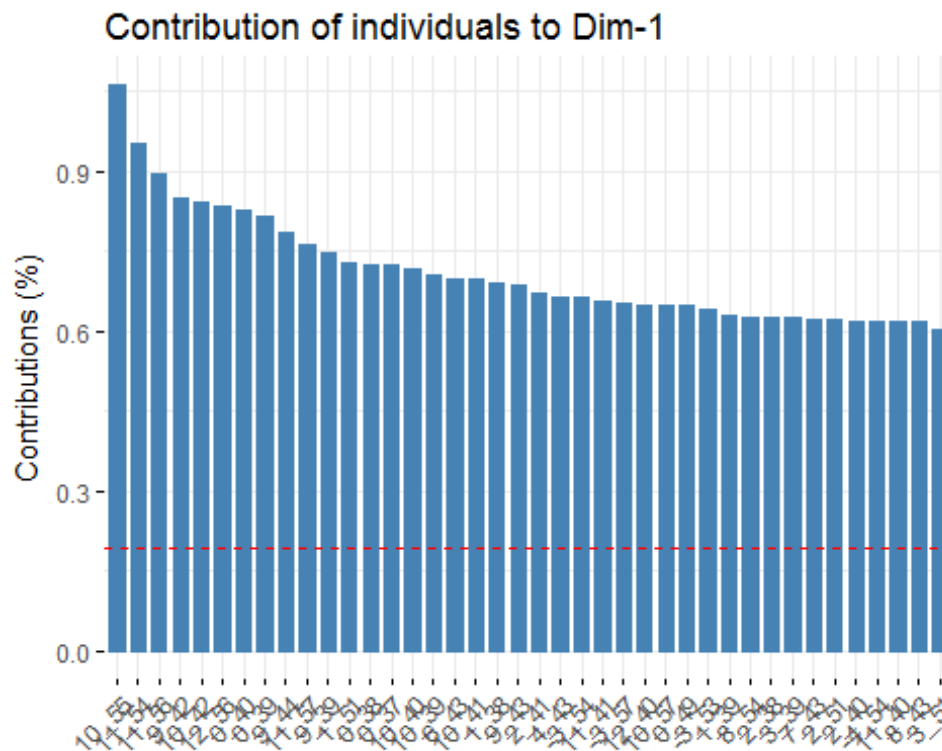
head(ind$contrib)

##           Dim.1      Dim.2      Dim.3      Dim.4      Dim.5
## -10_38 0.0146046098 0.2753285 0.3631290 0.0838990503 0.35468854
## -10_39 0.0085461268 0.3163824 0.2862557 0.1631446296 0.11242390
## -10_40 0.0027147525 0.3438548 0.2080633 0.1311122195 0.35177579
## -10_41 0.0002898144 0.3643067 0.2213309 0.0303907767 0.06135896
## -10_42 0.0009380543 0.3118706 0.1527086 0.0005886082 0.01254431
## -10_43 0.0085775133 0.3074837 0.1367381 0.1058054806 0.28461109
```

On va utiliser des graphiques pour rapidement identifier les individus qui contribuent très bien à une telle Dimension principale.

->Contribution totale des individus sur PC1 :

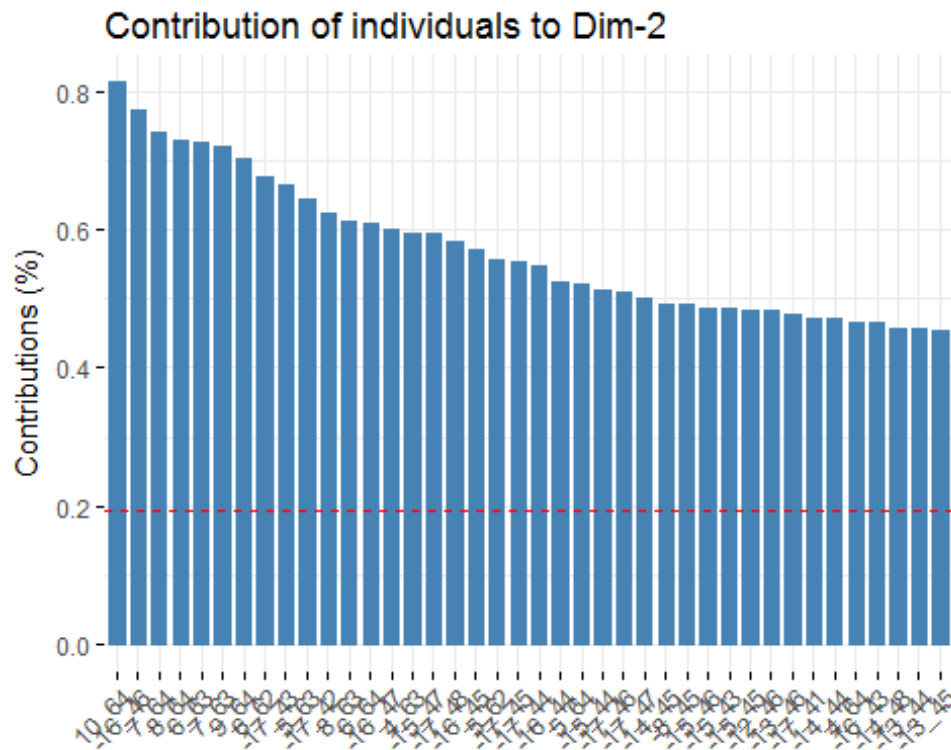
```
fviz_contrib(res.pca,choice="ind", axes=1,top = 40)
```



Donc la région qui contribue le plus à la définition du premier axe principale est la région comprise entre : -4=<lon<=12 et 37=<lat<57.

Contribution totale des individus sur PC2 :

```
fviz_contrib(res.pca,choice="ind", axes=2,top=40)
```

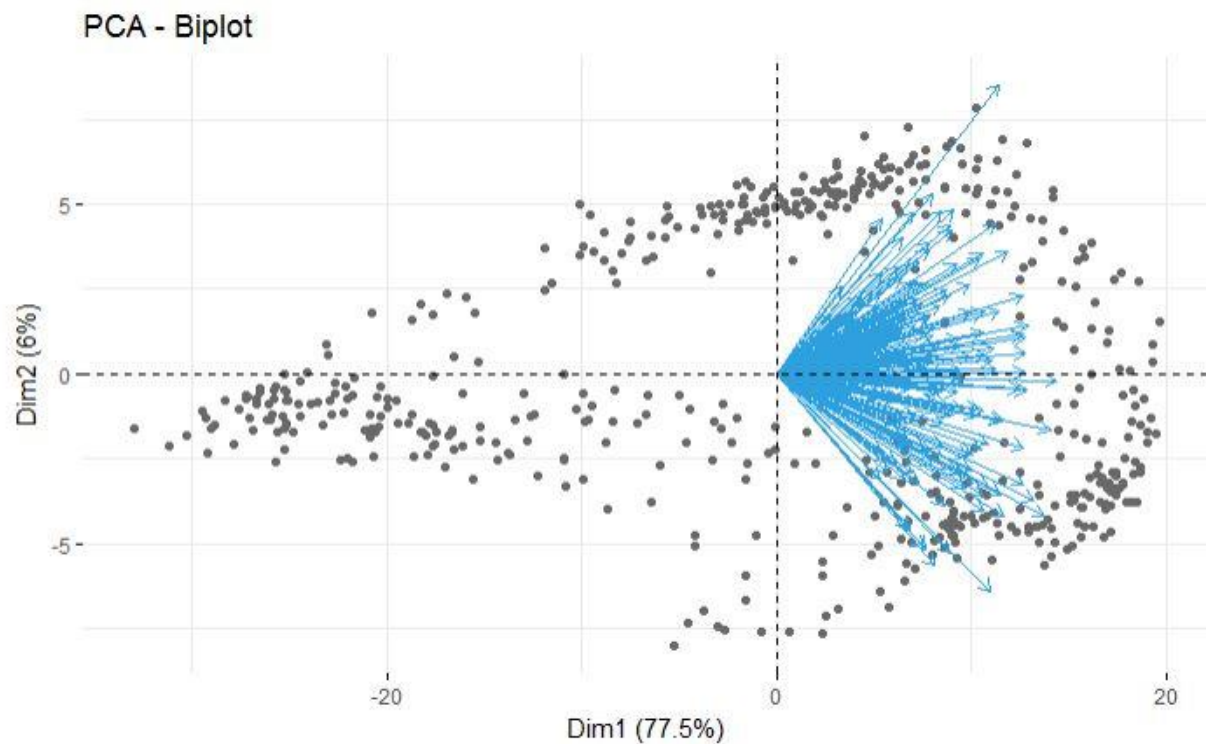


Donc la région qui contribue le plus à la définition du premier axe principale est la région comprise entre :

-17=<lon<=9 et 42=<lat<=64

Biplot :

Pour raison de lisibilité, j'ai supprimé tous les texts. Donc la figure ne contient que la représentation des variables(les flèches en bleu) ainsi que la nuage des individus(les points en noires).



Globalement un biplot peut être interprété comme suit : ->un individu qui se trouve du même côté d'une variable donnée a une valeur élevée pour cette variable ->un individu qui se trouve sur le côté opposé d'une variable donnée a une faible valeur pour cette variable.

NB: il faut se méfier des individus proches de l'origine : mal représentés, ou proches de la moyenne car ils sont mal représentés.

Commentaire pour notre jeu de données :

- Pour le quart du cercle en haut à droite : ces couple lon/lat (points noirs) ont généralement de grandes valeurs de swl pendant ces mois qui sont sur le même quart.
- Pour le quart du cercle en bas à gauche : ces couple lon/lat ont généralement des petites valeurs de swl pendant les mêmes mois que le cas précédent.
- Pour le quart du cercle en bas à droite : ces couple lon/lat (points noirs) ont généralement de grandes valeurs de swl pendant ces mois qui sont sur le même quart.
- Pour le quart du cercle en haut à gauche : ces couple lon/lat ont généralement des petites valeurs de swl pendant les mêmes mois que le cas précédent.

Malheureusement, à cause de la lisibilité du graph, on ne peut pas voir quels sont ces mois et les régions correspondantes.