

Preliminary Outlier Detection in Missing Values Free data

December 19, 2021

Problems : Valid estimation of a composition of rocks.

1. composition of major elements with outliers
2. composition of major elements with outliers and missing values
3. : composition of major elements with outliers and missing values
4. : composition of all elements with outliers and missing values

An initial goal of this work is to provide robust estimates of the composition of each rock. This problem is made difficult by the presence of both missing values (laboratories did not measure any composition) and outliers (the measures reported are extreme).

Setting aside the problem of missing values for now, one may look at a subcomposition $\mathbf{x} \in \mathbb{S}^D$ of oxides of major elements.

One begin by importing a dataset, for example GeoPT48, a monzonite.

```
setwd("/home/max/Documents/MStatistics/MA2/Thesis/Repository/")
data <- read_csv("data/raw/GeoPT48 -84Ra.csv")

##
## - Column specification -----
## cols(
##   .default = col_double(),
##   Laboratory = col_character(),
##   Au = col_logical(),
##   N = col_logical(),
##   Os = col_logical()
## )
## i Use 'spec()' for the full column specifications.

data2 <- read_csv("data/raw/GeoPT46 -84Ra.csv")

## Warning: Missing column names filled in: 'X1' [1]
##
## - Column specification -----
## cols(
##   .default = col_double(),
##   Laboratory = col_character(),
##   Ir = col_logical(),
##   N = col_logical(),
##   Os = col_logical(),
##   Rh = col_logical(),
##   Ru = col_logical()
## )
## i Use 'spec()' for the full column specifications.
```

Then, one looks at the subcomposition of major elements :

```

sel <-c("SiO2","TiO2","Al2O3","Fe2O3T","MnO","MgO","CaO","Na2O",
        "K2O","P2O5")

df.majors.raw <- select(data,all_of(sel))
df.majors2.raw <- select(data2,all_of(sel))

# data cleaning we remove all columns filled with NA
df.majors <- df.majors.raw[rowSums(is.na(df.majors.raw)) != ncol(df.majors.raw),]
df.majors2 <- data.frame(clo(df.majors2.raw[rowSums(is.na(df.majors2.raw)) != ncol(df

# closure operation
df.majors <- data.frame(clo(df.majors))
dim(df.majors.raw) # 97 observations

## [1] 97 10

dim(df.majors) # 86 observations

## [1] 86 10

# remove columns containing even only one missing value
df.majors.nafree <- df.majors.raw[rowSums(is.na(df.majors.raw)) == 0,] # 62 observati
df.majors2.nafree <- df.majors2.raw[rowSums(is.na(df.majors2.raw)) == 0,] # 78 observ

```

Then missing values are imputed by the column geometric mean

```

geomean.v <- sapply(rbind(df.majors),geomean)
for (i in 1:ncol(df.majors)){
  df.majors[,i][is.na(df.majors[,i])] <- geomean.v[i]
}

```

One then transform the dataset using the CLR transformation which is a mapping $\mathbb{S}^D \rightarrow \mathbb{U}^D$:

$$\mathbf{z} = \text{clr}(\mathbf{x} = [\log(x_1/g(x)), \dots, \log(x_D/g(x))]) \quad (1)$$

Where \mathbb{U}^D is an hyperplane of \mathbb{R}^D defined as :

$$U^D = \left\{ [u_1, \dots, u_D] : \sum_{i=1}^D u_i = 0 \right\}$$

```

clr.df <- clr(df.majors)
clr.df.nafree <- clr(df.majors.nafree)
clr.df.nafree2 <- clr(df.majors2.nafree)

```

1 Outlier Detection

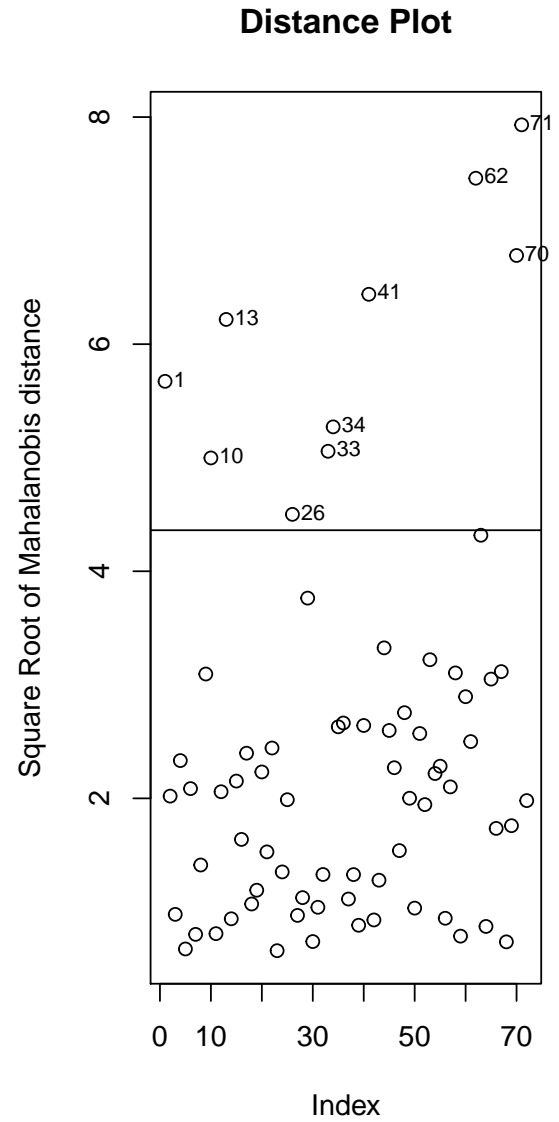
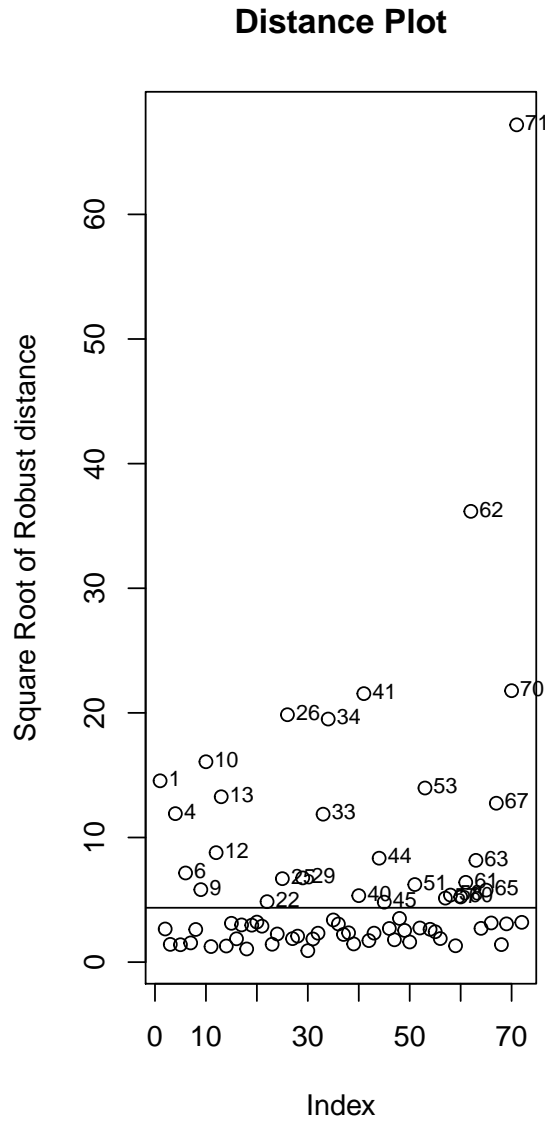
```
# outlier detection on the dataframe where missing values entries are replaced with c
ilrdf.mcd <- covMcd(clr2ilr(clr.df))
tolEllipsePlot(clr.df,classic = T,m.cov = clrdf.mcd)

## Error in tolEllipsePlot(clr.df, classic = T, m.cov = clrdf.mcd): Dimension
{= ncol(x)} must be 2!

1-length(ilrdf.mcd$best)/nrow(clr.df) # 44 % of outliers in this approach

## [1] 0.4418605

# outlier detection on the dataframe where missing values are not taken into account
ilrdf.nafree.mcd <- covMcd(clr2ilr(clr.df.nafree))
plot(ilrdf.nafree.mcd,which = c("distance"),classic = TRUE)
```



```
1-length(ilrdf.nafree.mcd$best)/nrow(clr.df.nafree) # 43 % of outliers in this approach
## [1] 0.4305556
```

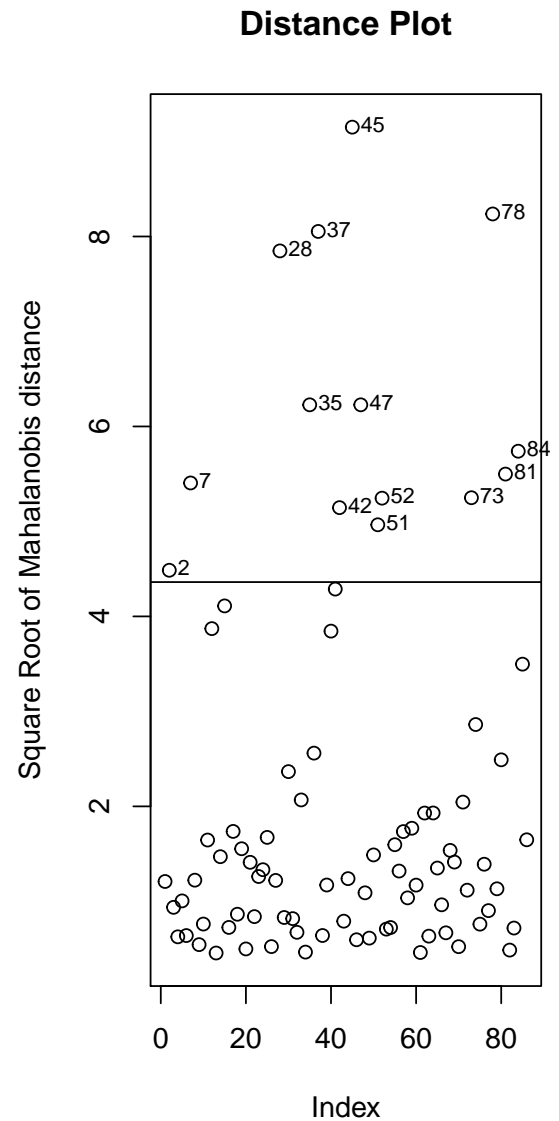
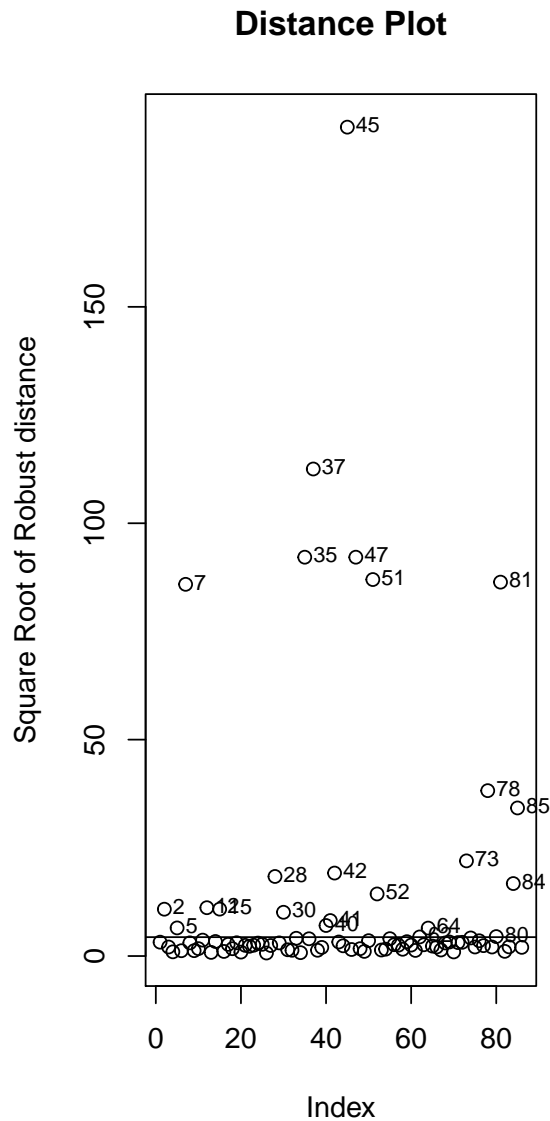
Is this proportion of outliers so high a particularity of the rock ? No, one sees that with the dataframe GeoPT46, again we flag almost 43% of observations as outliers, this is too much for practical reasons. `cov.mcd` has a default parameter `alpha` which is the proportion of the total sample size to take a subset of through the relationship :

$$\alpha \times n = h \quad (2)$$

We can take $\alpha = 0.75$ so this means the breakdown value of our estimator would be 25 %.

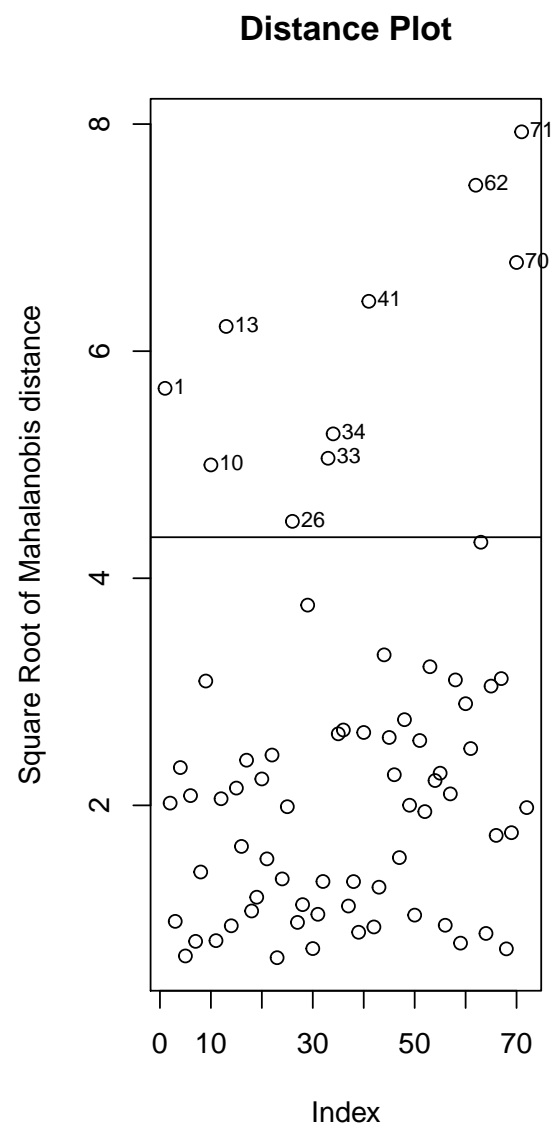
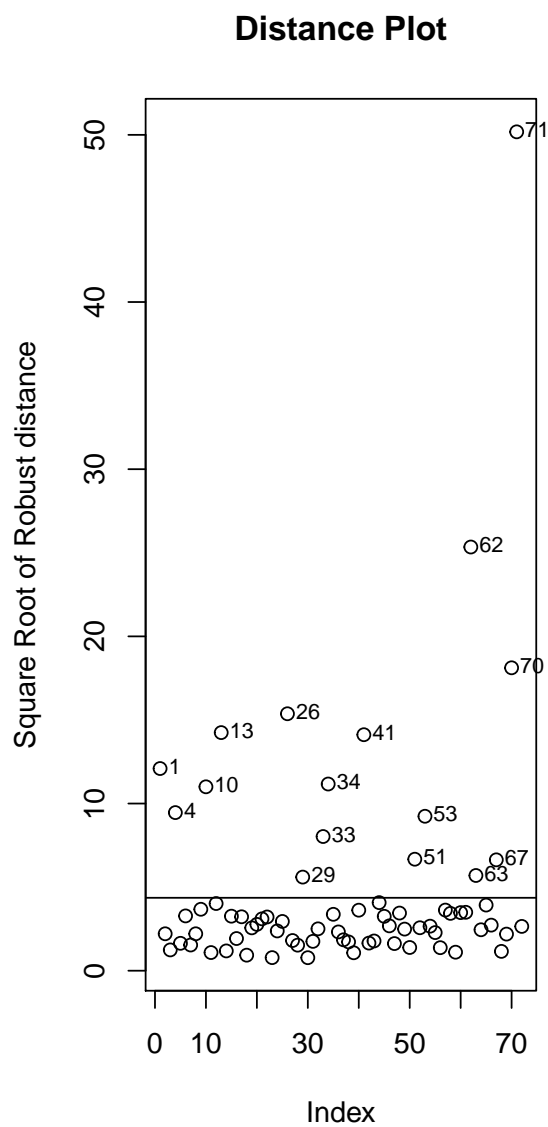
```
# outlier detection on the dataframe where missing values entries are replaced with c
ilrddf.mcd <- covMcd(clr2ilr(clr.df),alpha = .75)

plot(ilrddf.mcd,which = c("distance"),classic = TRUE)
```



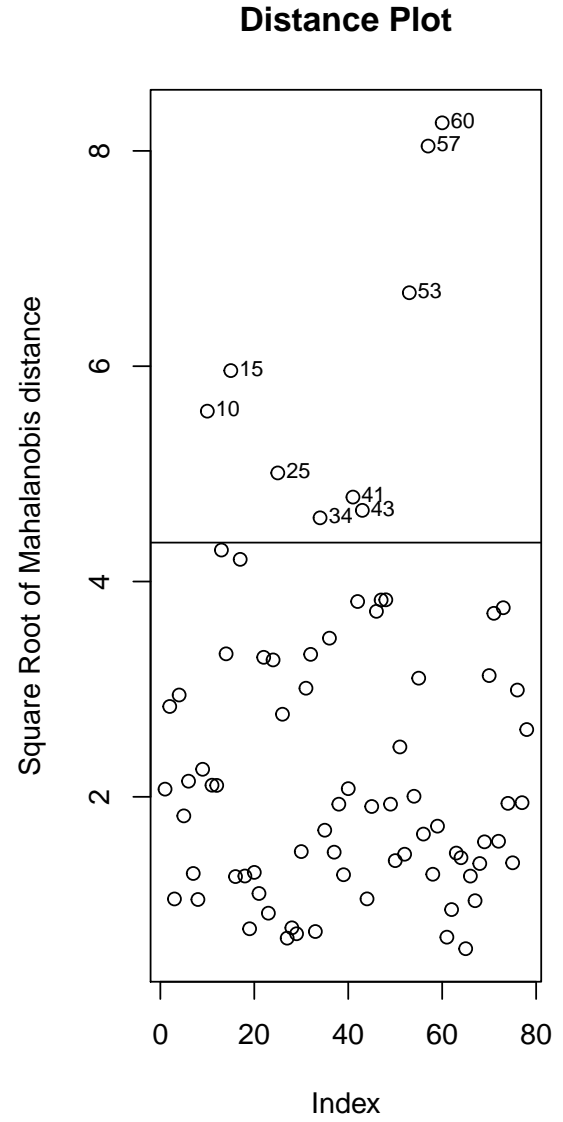
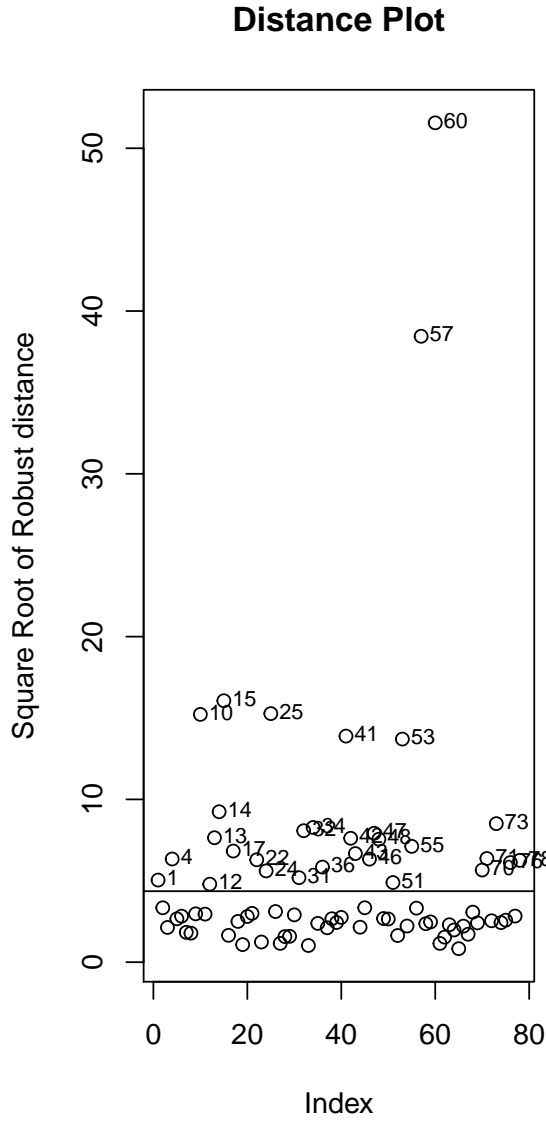
```
1-length(ilrddf.mcd$best)/nrow(clr.df) # 22 % of outliers in this approach
## [1] 0.2209302
```

```
# outlier detection on the dataframe where missing values are not taken into account
ilrddf.nafree.mcd <- covMcd(clr2ilr(clr.df.nafree),alpha=.75)
plot(ilrddf.nafree.mcd,which = c("distance"),classic = TRUE)
```



```
1-length(ilrdf.nafree.mcd$best)/nrow(clr.df.nafree) # 22 % of outliers in this approach
## [1] 0.2222222
```

```
ilrdf2.mcd <- covMcd(clr2ilr(clr.df.nafree2))
plot(ilrdf2.mcd,which = c("distance"),classic = TRUE)
```



```
1-length(ilrdf2.mcd$best)/nrow(clr.df.nafree2) # 43 % outliers in this approach, again
## [1] 0.4358974
```

Then principal component analysis is conducted. Z denotes the mean-centered data matrix X :

$$z_{ij} = x_{ij} - \mu_j$$

Where μ_j denotes the arithmetic mean of the j -th column. Recall that here using the arithmetic mean is justified because X now lives in a subspace of \mathbb{R}^D which is no longer constrained by the unit sum.

Perform Singular Value Decomposition on clr.df :

$$Z = UDW^T = (UD)W^T = Z^*W^T, Z \in \mathbb{R}^{n \times (d-1)}, U \in \mathbb{R}^{n \times p}, D \in \mathbb{R}^{p \times p}, W \in \mathbb{R}^{(d-1) \times p} \quad (3)$$

Z^* denotes the projection of the mean-centered data matrix on a space of dimension p . Hopefully, Z^* contains enough meaningful information about Z while having a much lower number of dimensions. To evaluate how good Z^* approximates Z , one looks at the proportion of variance explained by each of the components and more specifically, the cumulative proportion of variance explained by each of the components when these components are ranked from most to less important.

```
pca.clr <- prcomp(clr.df, scale = T, rank. = ncol(clr.df)-1 )
summary(pca.clr)

## Importance of first k=9 (out of 10) components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  2.6562 0.97703 0.91318 0.7470 0.56855 0.34083 0.25171
## Proportion of Variance 0.7055 0.09546 0.08339 0.0558 0.03233 0.01162 0.00634
## Cumulative Proportion 0.7055 0.80098 0.88437 0.9402 0.97249 0.98411 0.99045
##
##          PC8      PC9
## Standard deviation  0.24623 0.18687
## Proportion of Variance 0.00606 0.00349
## Cumulative Proportion 0.99651 1.00000

# repeat for nafree df
pca.clr.nafree <- prcomp(clr.df.nafree, scale = T, rank. = ncol(clr.df)-1 )
summary(pca.clr.nafree)

## Importance of first k=9 (out of 10) components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  2.2003 1.4357 1.0267 0.77238 0.73015 0.66663 0.51892
## Proportion of Variance 0.4841 0.2061 0.1054 0.05966 0.05331 0.04444 0.02693
## Cumulative Proportion 0.4841 0.6902 0.7957 0.85531 0.90863 0.95307 0.97999
##
##          PC8      PC9
## Standard deviation  0.37047 0.25063
## Proportion of Variance 0.01372 0.00628
## Cumulative Proportion 0.99372 1.00000
```

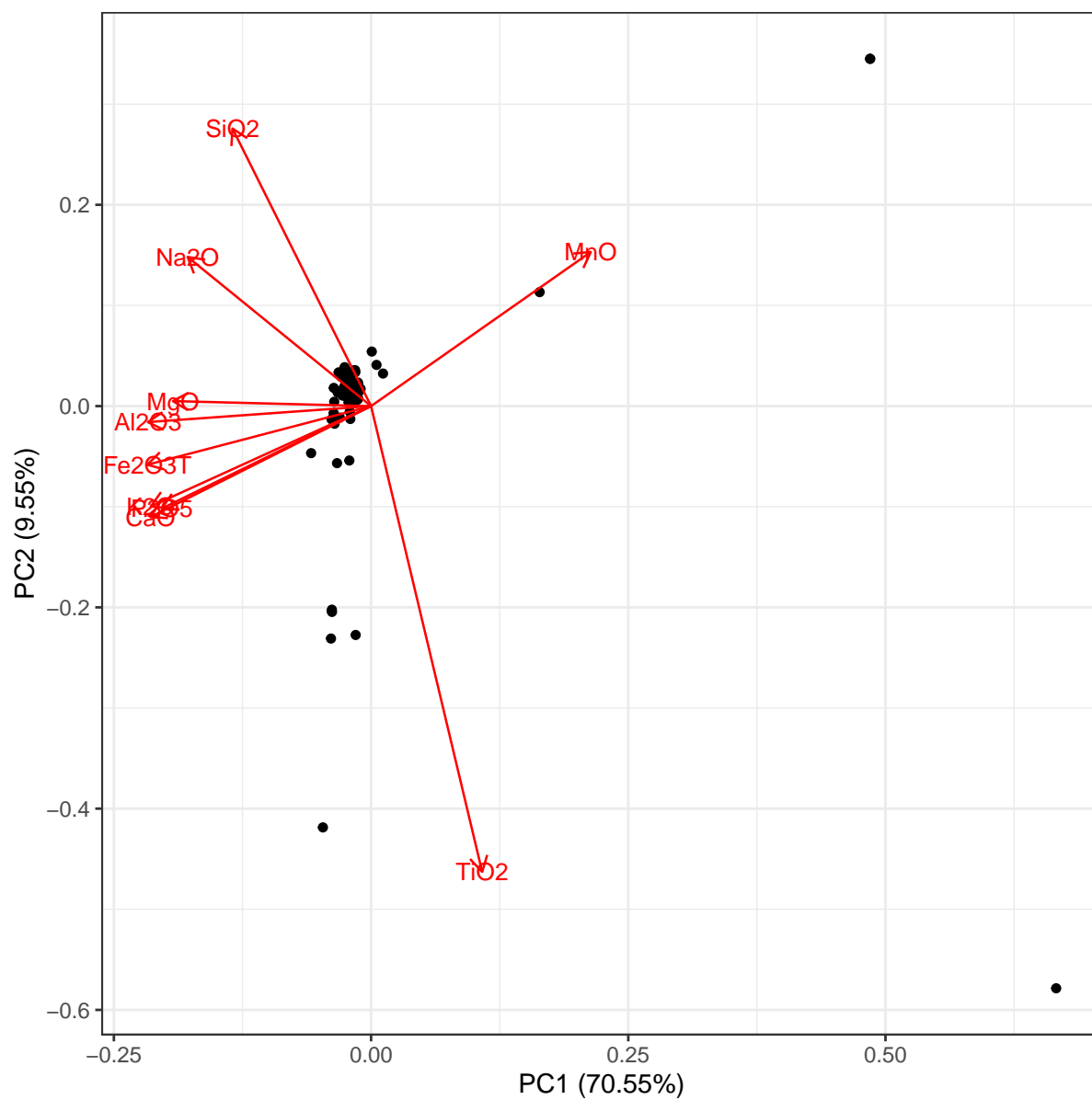
The importance of a component, in terms of proportion of variance explained, is directly related to the D matrix whose eigenvalues squared are directly related to the proportion of variance explained through the relationship :

$$\lambda_i = d_i^2 / (n - 1) \quad (4)$$

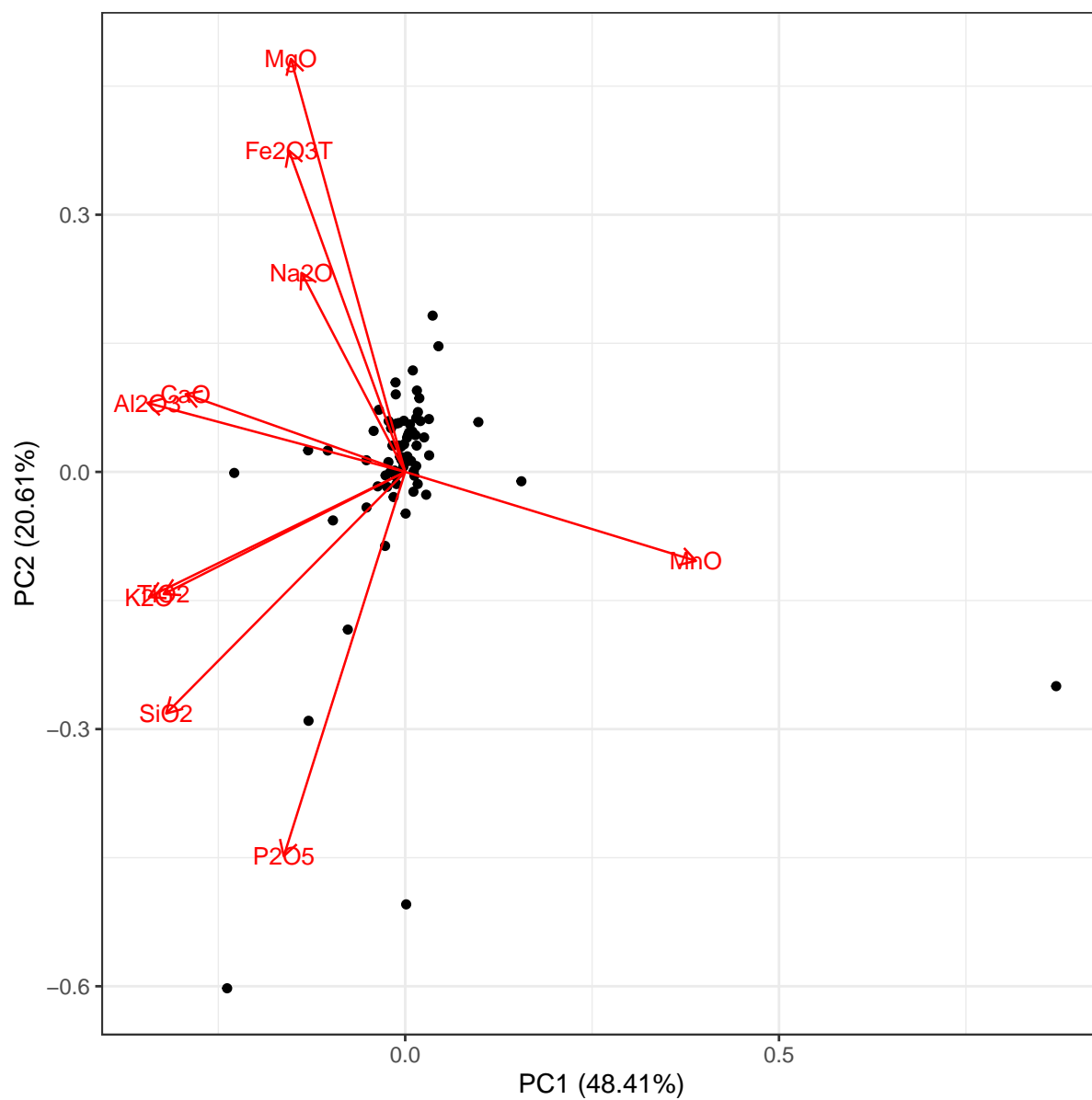
Where d_i denotes the i -th diagonal element in the square matrix D of the singular values and n is the number of principal components which is equal to the dimensionality of the original data matrix X .

Now, one looks at the rank-two approximation of Z (as the biplot does) :

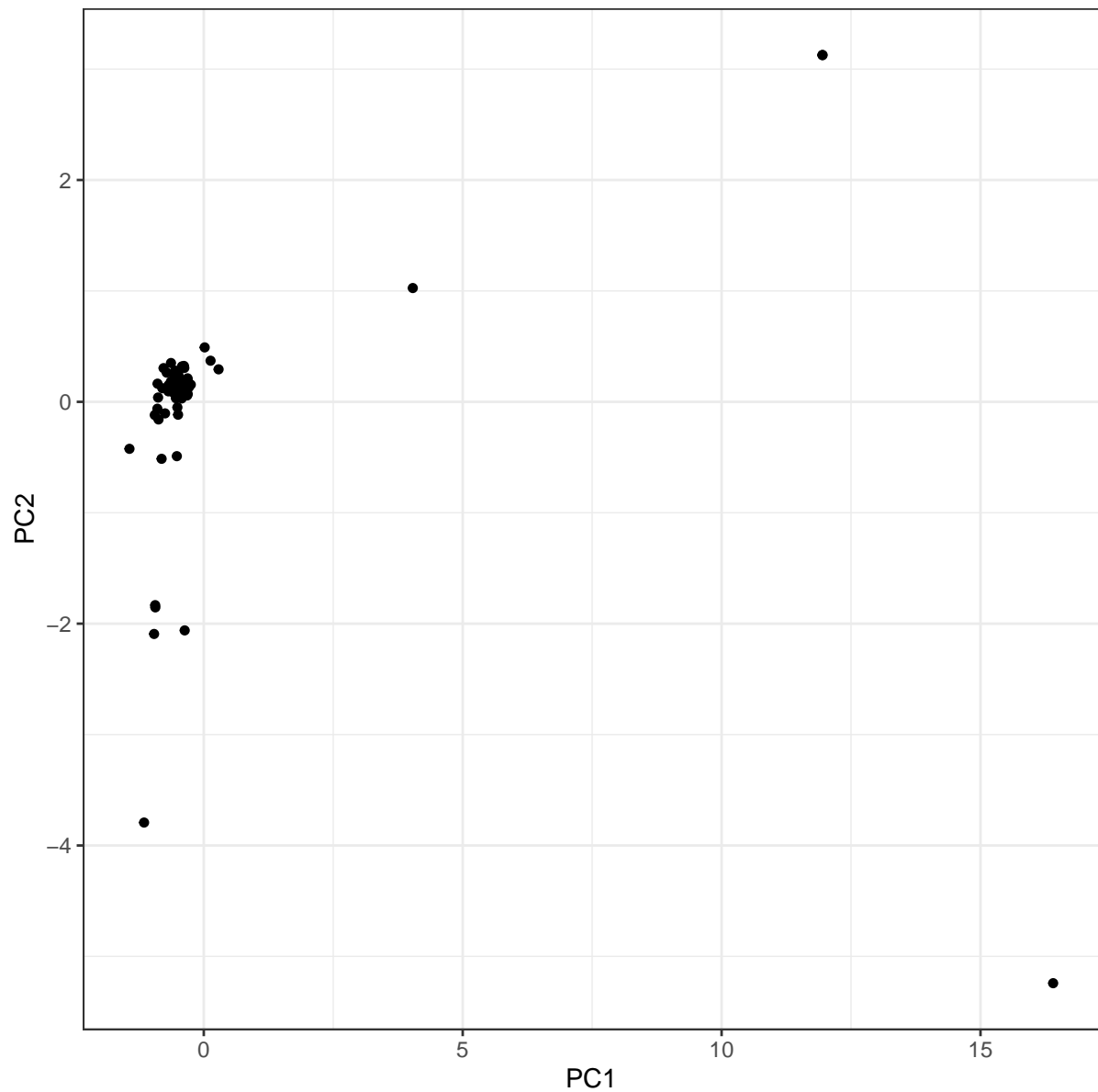
```
# autoplot
autoplot(pca.clr, loadings=T, loadings.label=T)+theme_bw()
```



```
autoplot(pca.clr.nafree,loadings=T,loadings.label=T)+theme_bw()
```

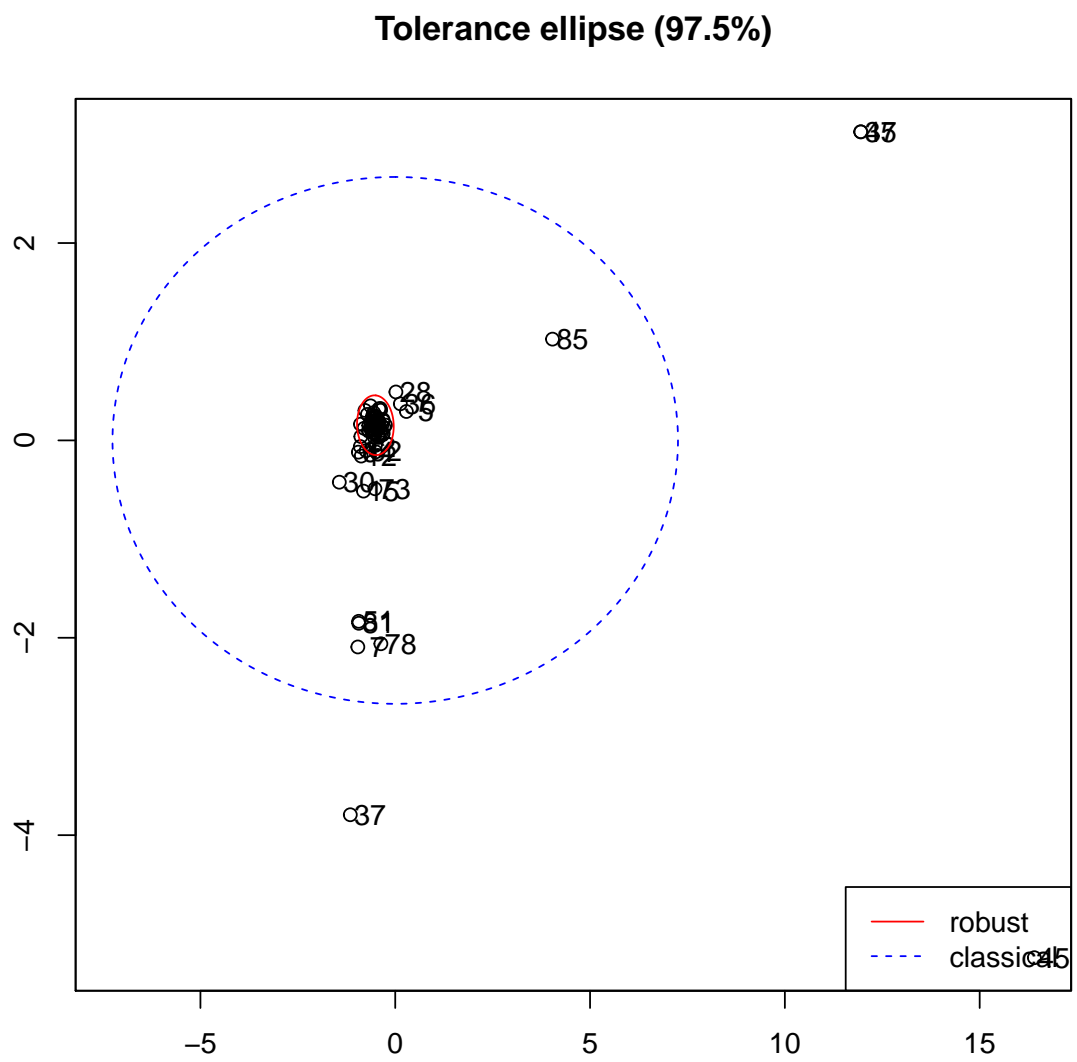


```
# manually extracting the rank 2 approximation of Z
Z.approx <- data.frame(pca.clr$x[,c(1,2)])
colnames(Z.approx) = c("PC1", "PC2")
ggplot(aes(x=PC1, y=PC2), data=Z.approx) + theme_bw() + geom_point()
```

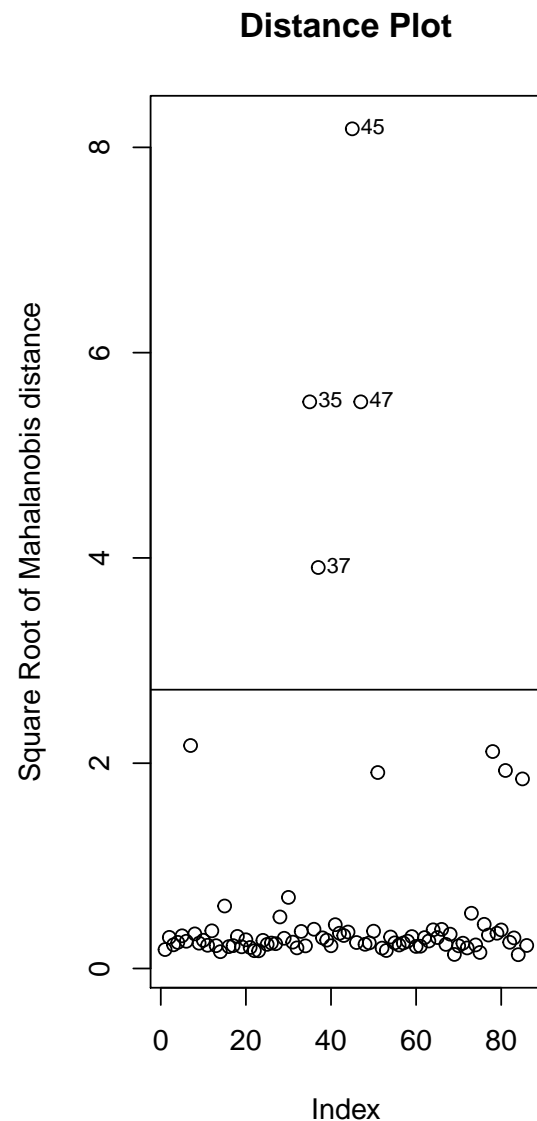
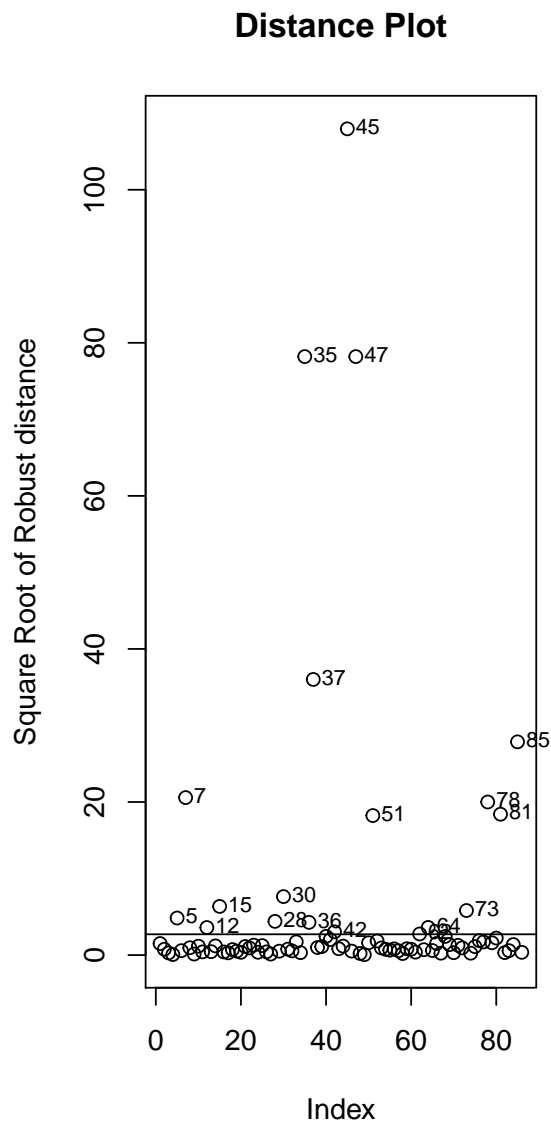


Flagging the outliers, remarkable discrepancy between robust method (MCD estimator of the location, MCD is affine equivariant and has a high breakdown value) vs classical method (sample covariance matrix).

```
Z2.mcd <- covMcd(Z.approx,alpha = 0.75)
tolEllipsePlot(Z.approx,classic = T,m.cov = Z2.mcd)
```



```
plot(Z2.mcd, which = c("distance"), classic = TRUE)
```



```
Z2.outfree <- Z.approx[Z2.mcd$best,]
nrow(Z2.outfree)/nrow(Z.approx)

## [1] 0.755814
```

This approach has a pitfall. The estimation of the PC is itself not robust. Robcompositions package provides a robust PCA estimation method :

The compositional data set is expressed in isometric logratio coordinates. Afterwards, robust principal component analysis is performed. Resulting loadings and scores are back-transformed to the clr space where the compositional biplot can be shown. CITE ROBCOMP R package

```
rob.pca.clr <- pcaCoDa(df.majors)
summary(rob.pca.clr)
```

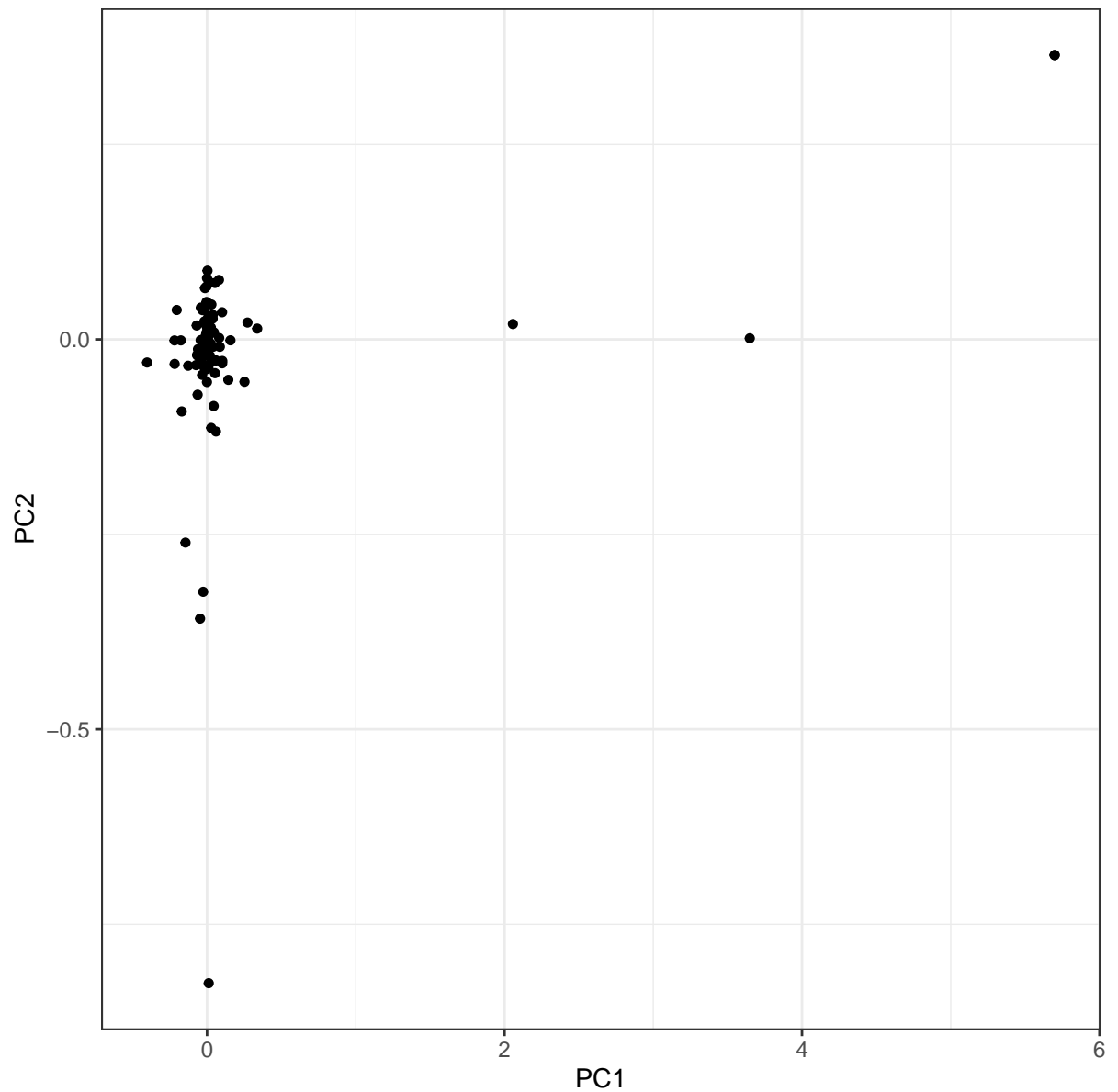
Importance of components:

##	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5
## Standard deviation	0.0730207	0.03302854	0.03242246	0.02468697	0.0205167
## Proportion of Variance	0.5873237	0.12016122	0.11579171	0.06713070	0.0463661
## Cumulative Proportion	0.5873237	0.70748494	0.82327665	0.89040735	0.9367734

##	Comp.6	Comp.7	Comp.8	Comp.9
## Standard deviation	0.01848889	0.01108852	0.009229454	0.004901582
## Proportion of Variance	0.03765368	0.01354355	0.009382911	0.002646416
## Cumulative Proportion	0.97442713	0.98797067	0.997353584	1.000000000

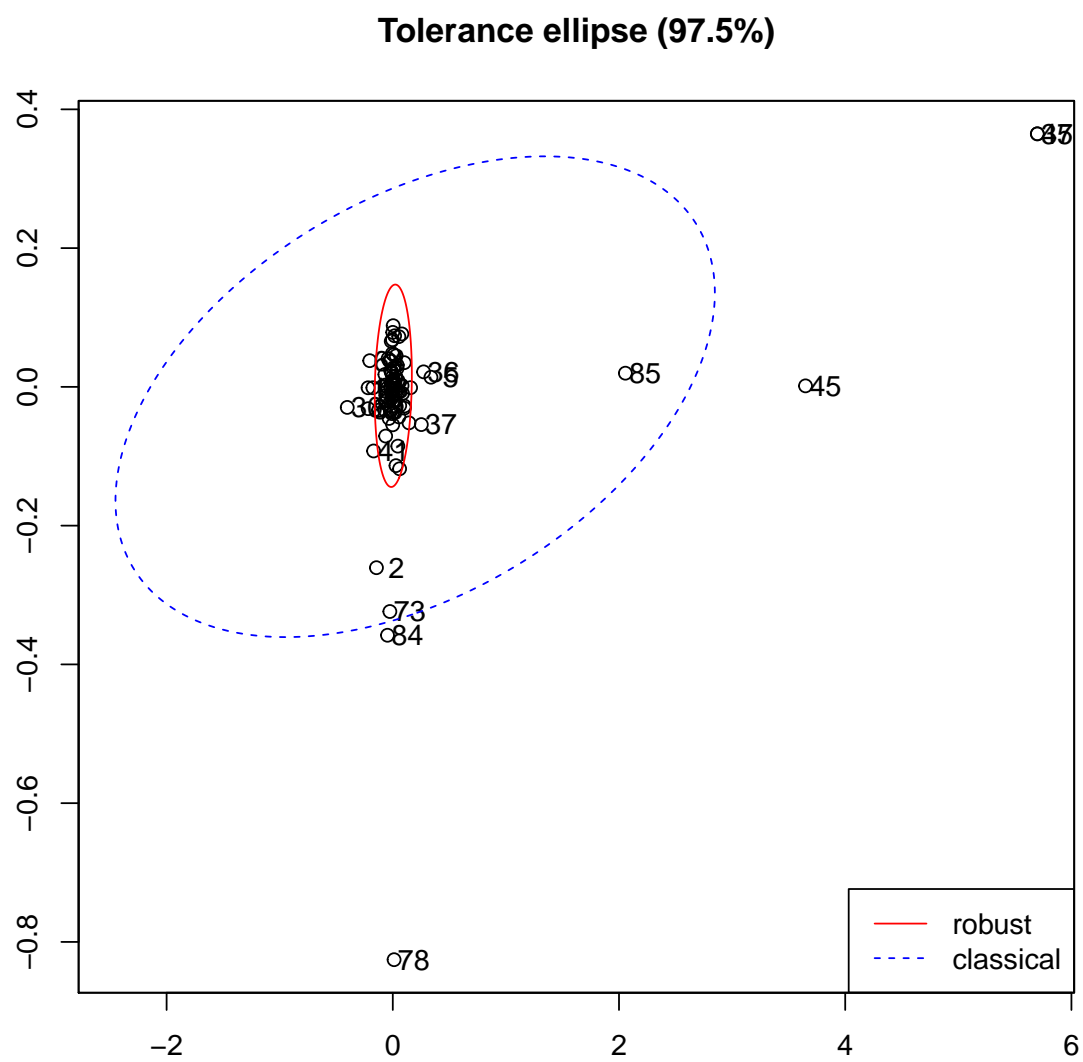
One sees that the first two PC's found by using the robust method explains (89 %) whereas in the classical method, the first two PC's explained only 76 %. The outlier detection is repeated in the first 2 robust PC's subspace.

```
Z.approx.rob <- data.frame(rob.pca.clr$scores[,c(1,2)])
colnames(Z.approx.rob) = c("PC1", "PC2")
ggplot(aes(x=PC1, y=PC2), data=Z.approx.rob) + theme_bw() + geom_point()
```

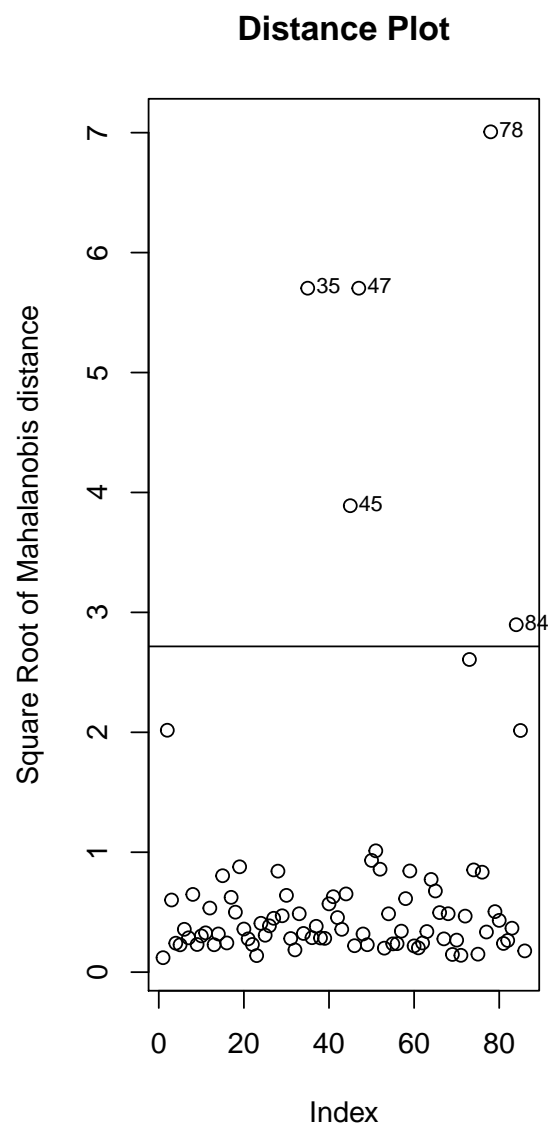
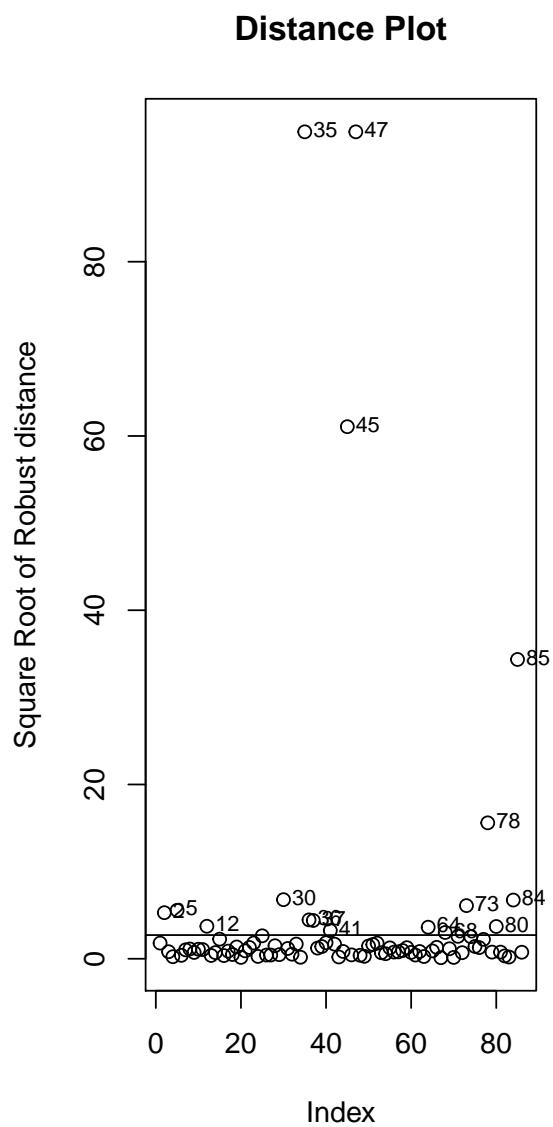


At first glance, there seem to be less outliers in the robust first two PC's space. This is checked by using diagnostic plots :

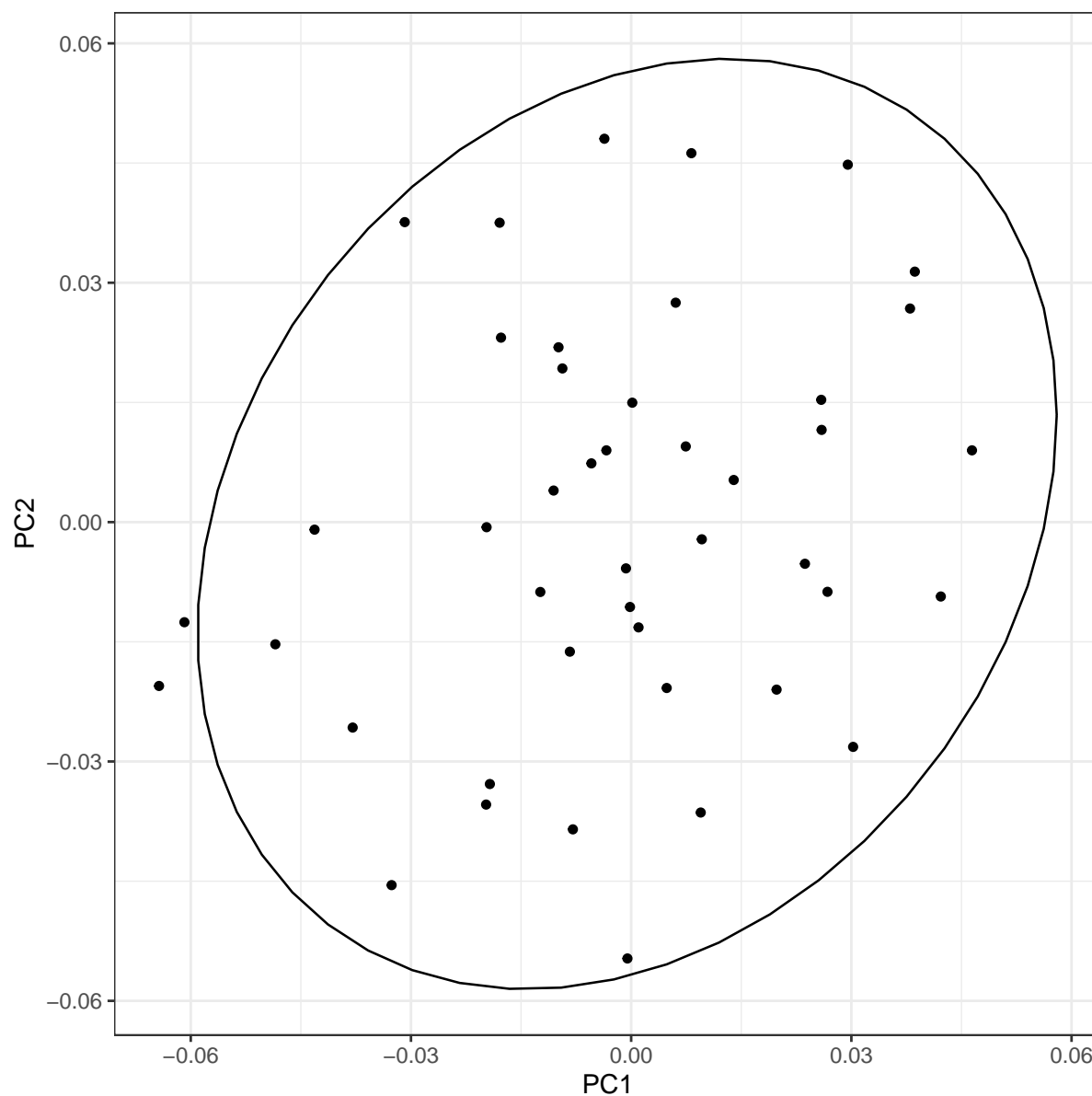
```
Z2.rob.mcd <- covMcd(Z.approx.rob,)  
tolEllipsePlot(Z.approx.rob, classic = T)
```

```
plot(Z2.rob.mcd,which = c("distance"),classic = TRUE)
```



```
Z2.rob.outfree <- Z.approx.rob[Z2.rob.mcd$best,]
ggplot(aes(x=PC1,y=PC2),data=Z2.rob.outfree)+theme_bw()+geom_point()+stat_ellipse()
```



There seems to be no association between PC1 and PC2, which makes much more sense. Eventually, an estimation for rock 1 composition is

```
df.majors.out.free <- df.majors[Z2.rob.mcd$best,]
df.majors.out.free.clr <- data.frame(clr(df.majors.out.free))
clr.mean <- colMeans(df.majors.out.free.clr)
clr.cov <- Cov(df.majors.out.free.clr)

mean.estimate <- as.vector(clrInv(clr.mean))
cov.estimate <- clrInv(as.matrix(clr.cov$cov))
cov.estimate
```

##	SiO2	TiO2	Al2O3	Fe2O3T	MnO
## SiO2	"0.10289487"	"0.09962841"	"0.09987669"	"0.09979176"	"0.0996833"
## TiO2	"0.09963294"	"0.10007568"	"0.10001475"	"0.10002610"	"0.1000330"
## Al2O3	"0.09988130"	"0.10001482"	"0.10002404"	"0.10001191"	"0.1000105"

```
## Fe2O3T "0.09979635" "0.10002615" "0.10001190" "0.10005504" "0.1000140"
## MnO    "0.09968785" "0.10003302" "0.10001042" "0.10001398" "0.1000981"
## MgO    "0.09971372" "0.10001797" "0.10000860" "0.10001880" "0.1000277"
## CaO    "0.09959408" "0.10005293" "0.10000200" "0.10002234" "0.1000343"
## Na2O   "0.09966646" "0.10004172" "0.10001381" "0.09999442" "0.1000451"
## K2O    "0.09967389" "0.10004639" "0.10001272" "0.10003256" "0.1000178"
## P2O5   "0.09949945" "0.10005849" "0.10001994" "0.10002814" "0.1000317"
##      MgO      CaO      Na2O      K2O      P2O5
## SiO2   "0.09970915" "0.09958957" "0.09966191" "0.09966934" "0.09949499"
## TiO2   "0.10001794" "0.10005295" "0.10004171" "0.10004637" "0.10005855"
## Al2O3  "0.10000864" "0.10000209" "0.10001387" "0.10001277" "0.10002008"
## Fe2O3T "0.10001883" "0.10002242" "0.09999446" "0.10003259" "0.10002826"
## MnO    "0.10002772" "0.10003430" "0.10004514" "0.10001776" "0.10003176"
## MgO    "0.10008386" "0.10003285" "0.10002368" "0.10002052" "0.10005227"
## CaO    "0.10003280" "0.10009855" "0.10004622" "0.10004775" "0.10006905"
## Na2O   "0.10002366" "0.10004625" "0.10009211" "0.10002929" "0.10004714"
## K2O    "0.10002050" "0.10004779" "0.10002930" "0.10005981" "0.10005929"
## P2O5   "0.10005218" "0.10006900" "0.10004706" "0.10005921" "0.10013487"
## attr(,"class")
## [1] "acomp"
```

How does it compare with the naive estimates of GeoPT ?

```
mean.naive <- colMeans(clo(df.majors))
od <- outCoDa(df.majors, quantile = 0.975, method = "robust", alpha = 0.9, coda = TRUE)
df.major.wo.outliers <- df.majors[od$outlierIndex,]
mean.naive.wo.outlier <- colMeans(clo(df.major.wo.outliers))
```

Not much difference for SiO₂ *but* TiO₂ concentration is twice as high in the naive way, 5 times higher in the naive way when removing outliers. MnO concentration is 10 times higher in the naive way.

```
df <- data.frame(mean.estimate,mean.naive,mean.naive.wo.outlier)
df

##      mean.estimate mean.naive mean.naive.wo.outlier
## SiO2      0.572280369 0.562480850      0.509635270
## TiO2      0.011132233 0.016099834      0.030527015
## Al2O3     0.201225261 0.197964339      0.200463923
## Fe2O3T    0.049229936 0.048373292      0.048931757
## MnO       0.001432903 0.013702982      0.047356073
## MgO       0.011614807 0.011273968      0.011162882
## CaO       0.038608273 0.038366921      0.040530951
## Na2O      0.069131939 0.066717132      0.064349912
## K2O       0.039169507 0.038844793      0.040388171
## P2O5      0.006174773 0.006175889      0.006654047
```

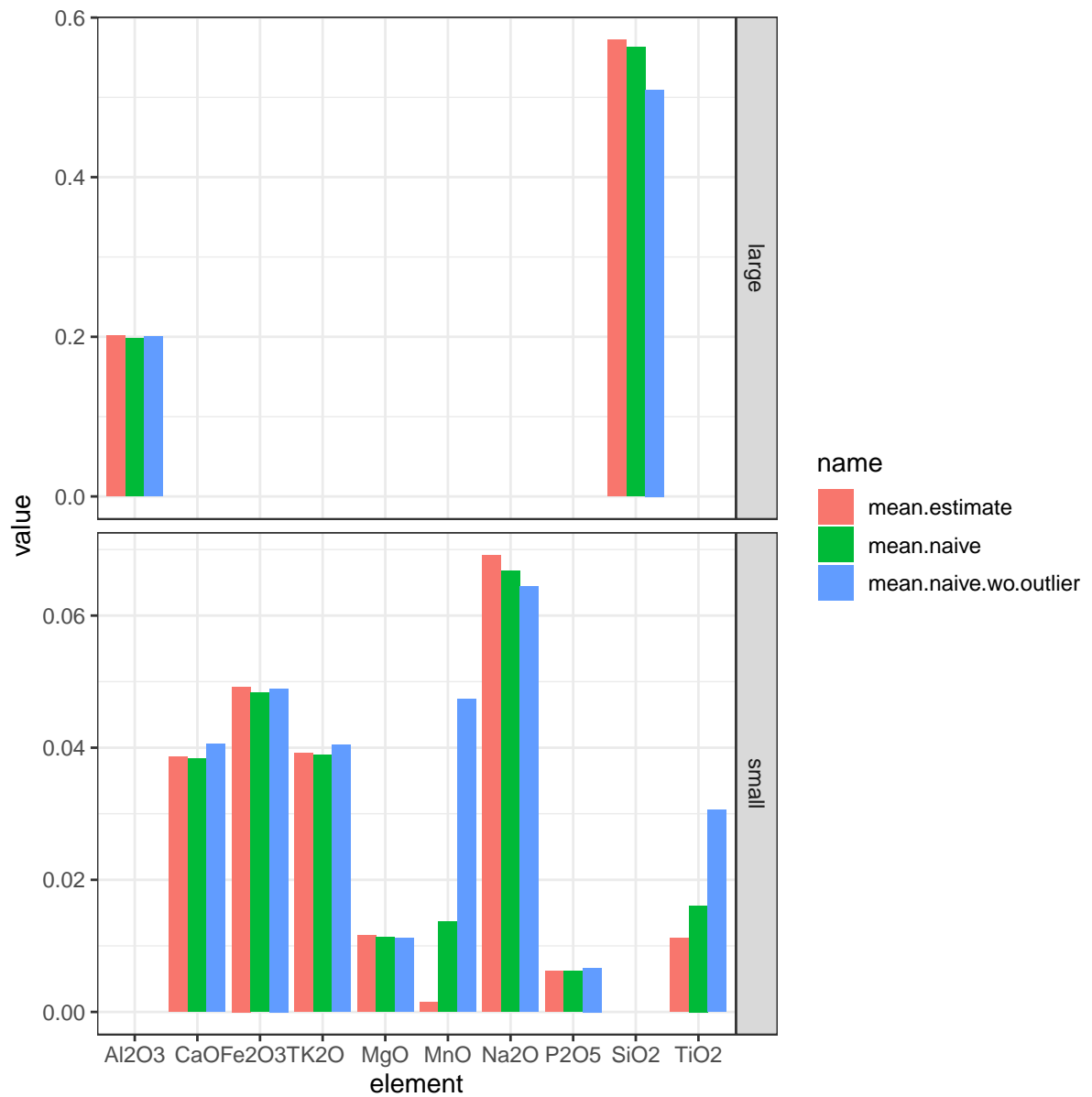
```

df$element <- row.names(df)
df.plot <- df %>% pivot_longer(cols=starts_with("mean"))
df.plot

## # A tibble: 30 x 3
##   element name          value
##   <chr>   <chr>         <dbl>
## 1 SiO2    mean.estimate      0.572
## 2 SiO2    mean.naive         0.562
## 3 SiO2    mean.naive.wo.outlier 0.510
## 4 TiO2    mean.estimate      0.0111
## 5 TiO2    mean.naive         0.0161
## 6 TiO2    mean.naive.wo.outlier 0.0305
## 7 Al2O3    mean.estimate      0.201
## 8 Al2O3    mean.naive         0.198
## 9 Al2O3    mean.naive.wo.outlier 0.200
## 10 Fe2O3T mean.estimate      0.0492
## # ... with 20 more rows

df.plot$group <- ifelse(df.plot$element %in% c("SiO2","Al2O3"),yes = "large","small")
ggplot(aes(x=element,y=value),data = df.plot)+geom_col(aes(fill=name),position = "dodge")

```



Preliminary conclusion : naive without outlier downplay the concentration of elements present in large quantities and blow the concentration of elements present in small quantities.