

PCA analysis

Rachel Richardson

November 18, 2018

In order to determine if the lipid profiles are different between experimental groups, we've run a PCA across lipid data. (Also to check for batch effects, not shown here.)

Step 1: Download data into R and make workable dataframes for analysis with all data and maternal diet-based subsets.

```
DGs1.2 <- read.csv("1.2DGs.changed LOQ.csv", head = TRUE, row.names = 1)
AC <- read.csv("AC.formatted.noloq.csv", head = TRUE, row.names = 1)
Cer <- read.csv("Cerimides.changedLOQ.csv", head = TRUE, row.names = 1)
DGs1.3 <- read.csv("CM033018 1-3DGs.changedLOQ.csv", head = TRUE, row.names = 1)
dh <- read.csv("dhCer.changedLOQ.csv", head = TRUE, row.names = 1)
Glu <- read.csv("GluCer.changedLOQ.csv", head = TRUE, row.names = 1)
hex <- read.csv("hexosylCer.changedLOQ.csv", head = TRUE, row.names = 1)
Lac <- read.csv("LacCer.changedLOQ.csv", head = TRUE, row.names = 1)
mye <- read.csv("Sphingomyelins.formatted.noloq.csv", head = TRUE, row.names = 1)
sine <- read.csv("Sphingosine.formatted.noloq.editfordismatrix.csv", header = TRUE, row.names = 1)
TAG <- read.csv("TAG.changedLOQ.csv", head = TRUE, row.names = 1)

#Create a dataframe with all PCA
ALL <- cbind(AC, Cer, DGs1.2, DGs1.3, dh, Glu, hex, Lac, mye, sine, TAG)

#Create a list of all data structures and the names of each.
plotlist <- list(DGs1.2, DGs1.3, TAG, AC, dh, Glu, hex, Lac, sine, Cer, mye, ALL)
plotname <- c("DGs1.2", "DGs1.3", "TAG", "AC", "dh", "Glu", "hex", "Lac", "sine", "Cer", "mye", "ALL")

#Clean up labels in ALL
Labels <- gsub("..pmol.", "", colnames(ALL))
Labels <- gsub("X", "", Labels)
colnames(ALL) <- Labels

#Create subsets
CTR <- rbind(ALL[1:8,], ALL[14:21,], ALL[30:36,])
HFD <- rbind(ALL[9:13,], ALL[22:29,])
```

Next, we calculate the principal components with `prcomp()`. Scree plots show how each principal component describes the data and how many of these components we should reasonably consider to describe data trends.

```
#Principal component
prin_comp <- prcomp(ALL, scale. = T, center = T)

#Variance calculation
Totalvar <- (prin_comp$sdev)^2
Proportionvar <- Totalvar/sum(Totalvar)

#plotting
par(mfrow=c(1,2))

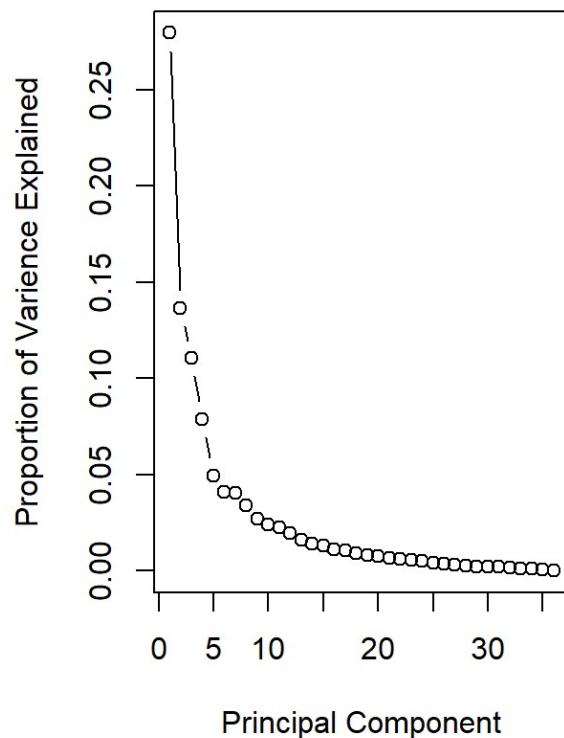
#scree plot

plot(Proportionvar, xlab = "Principal Component", ylab = "Proportion of Variance Explained", type = "b", main = "Proportion for principal component")

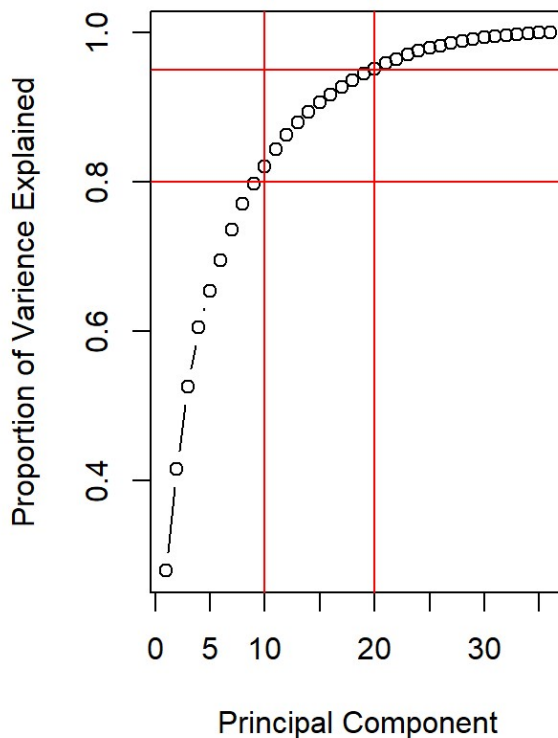
#cumulative scree plot

plot(cumsum(Proportionvar),xlab = "Principal Component", ylab = "Proportion of Variance Explained", type = "b", main = "Cumulative, lines at 80 and 95%")
abline(h=0.95, v=20, col="red")
abline(h=0.80, v=10, col="red")
```

Proportion for principal componer



Cumulative, lines at 80 and 95%



We've elected to use several additional packages to help us visualize our PCA:

```
library(ggpubr) #For arranging biplots
library(devtools) #For arranging biplots
library(ggbiplot) #For arranging all plots

library("FactoMineR") #Specific for PCA visualizing
library("factoextra") #Specific for PCA visualizing
```

```

#Display PCA; png saving options anf format adjustments are commented out.

###png("PCA1.png", height = 600, width = 800)

#Define sample groups (update with metadata to observe possible batch effects)
groups <- c(rep("Ln/CTR-CTR",8),rep("Ln/CTR-HFD",5),rep("Ln/HFD-CTR",8),rep("Ln/HFD-
HFD",8),rep("Ob/HFD-CTR",7))

#Plot with circles
g <- ggbiplot(prin_comp, obs.scale = 1, var.scale = 1, groups = groups, ellipse = TRU
E, circle = TRUE, var.axes = FALSE, varname.size = 0) +

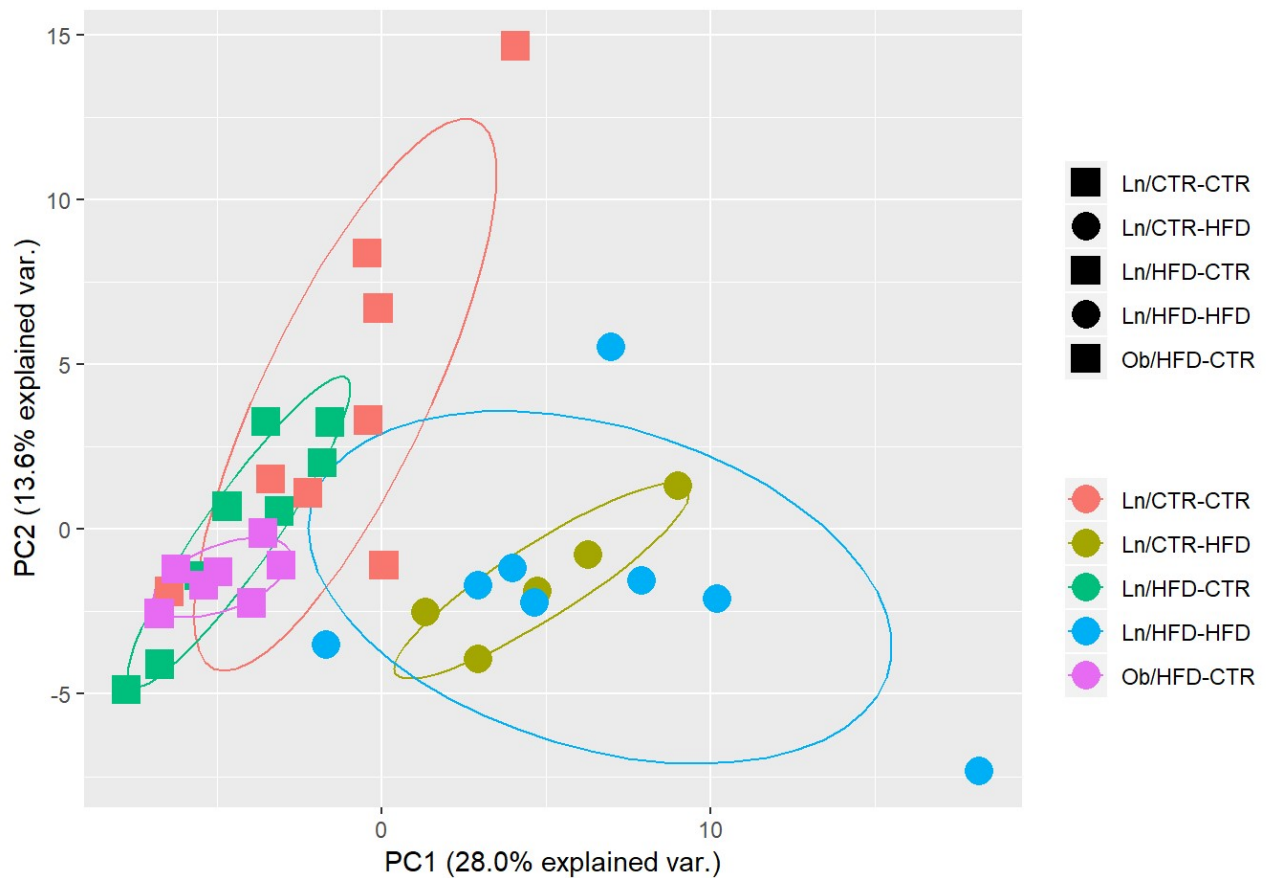
  geom_point(aes(colour=groups, shape = groups), size = 5)+ #Plot points

  scale_shape_manual(name= "", values = c(15,16,15,16,15))+ #Manually adjusted shape b
ased on Juvinille diet update with metadata to observe possible batch effects)
  scale_color_manual(name="", values=c("#F8766D", "#A3A500", "#00BF7D", "#00B0F6", "#E
76BF3"))+ #Experimental group colors
  guides(shape = guide_legend(order = 1), color = guide_legend(order = 2)) # Order leg
ends

#g <- g + theme_light() +
#       theme(panel.grid = element_line(colour = "white"),
#             legend.direction = 'vertical',
#             legend.position = c(0.85, 0.75),
#             axis.text.x = element_text(size = 26), axis.text.y = element_text(size
# = 26),
#             axis.title.x = element_text(size = 30), axis.title.y = element_text(siz
# e = 30),
#             legend.text = element_text(size = 24), legend.title = element_text(siz
# e = 0))

print(g) #Show plot

```

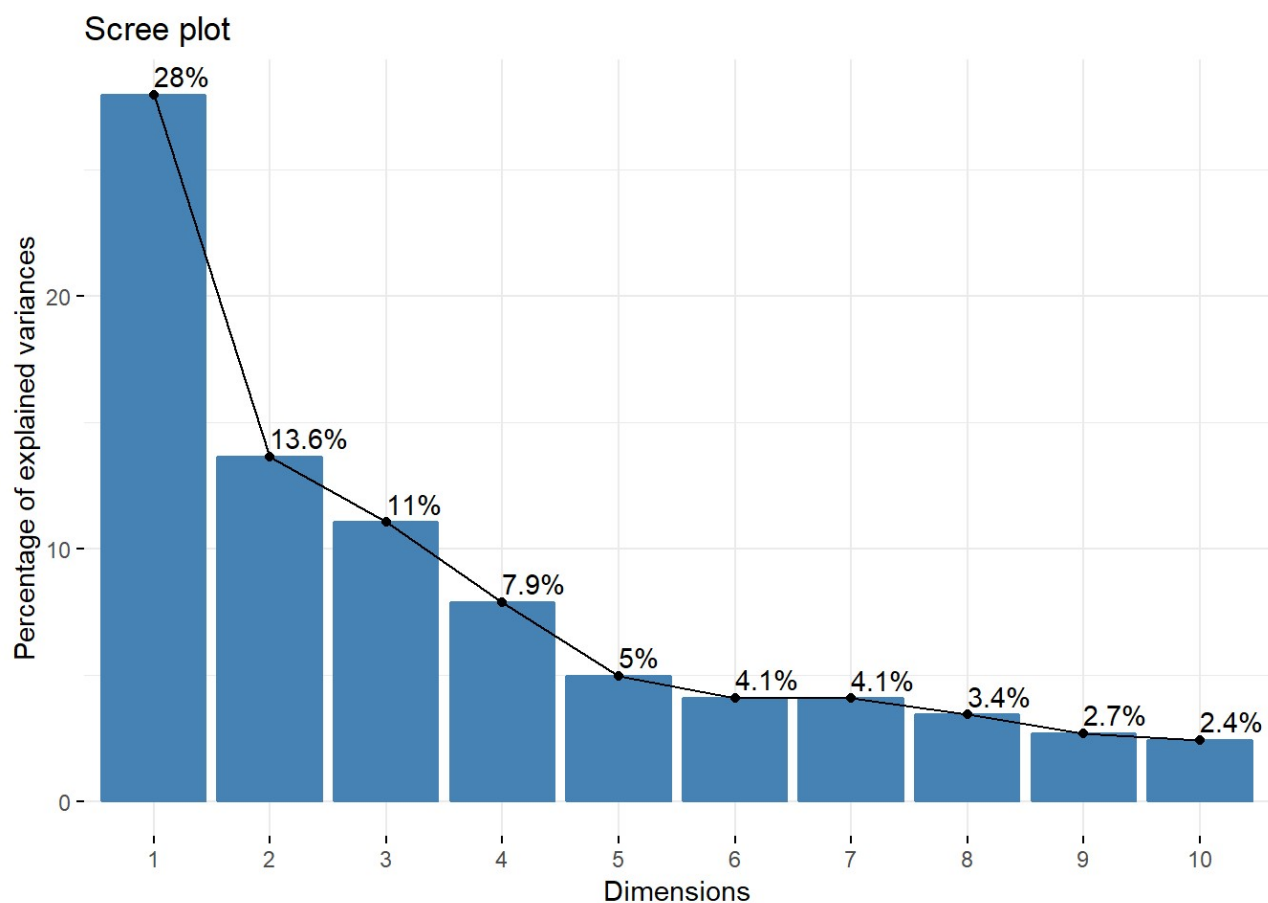


```
###dev.off()
```

Using these packages, we can recreate our scree plot and investigate the variance explained by the calculated principal components and which lipids are contributing the most to each component.

```
PCA_36 <- PCA(ALL, scale.unit = TRUE, ncp = 36, graph = FALSE) #PCA used in the PCA specific package

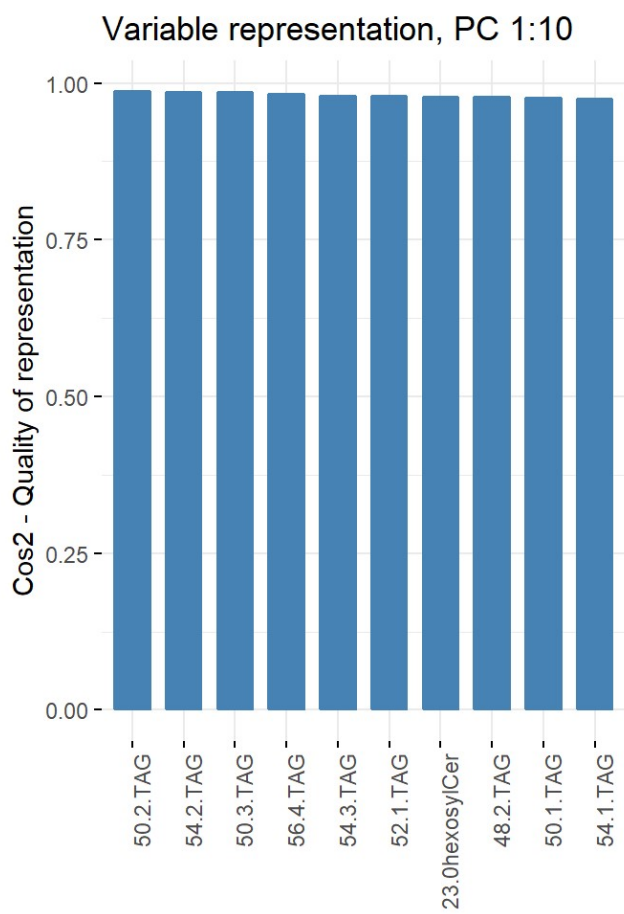
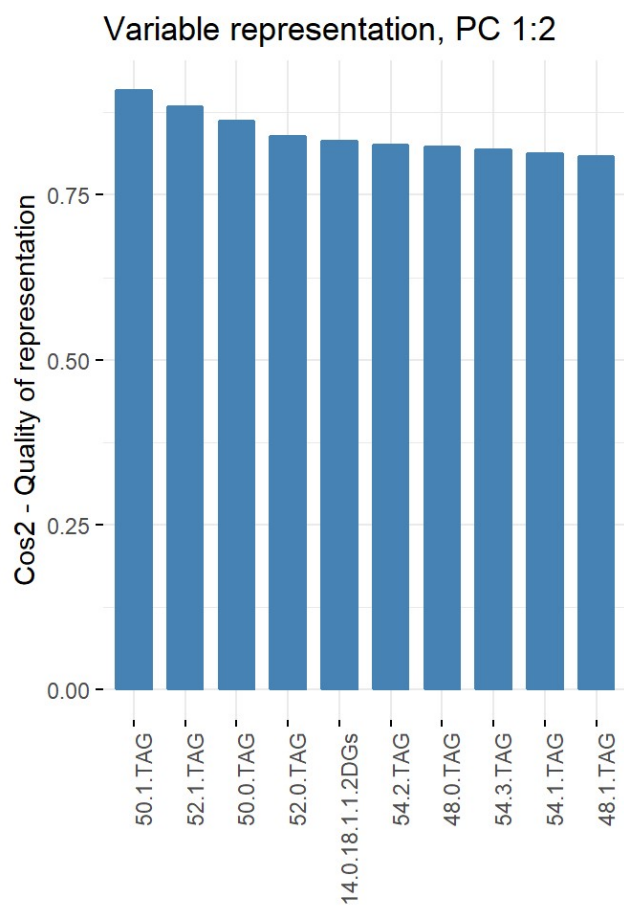
#scree plot
fviz_eig(PCA_36, addlabels = TRUE)
```



#Variable represation for the first two prinicpal components (as used in the PCA), and out to 10 principal components (Noted to account for 80% of variation in the data base d on earlier scree plots.)

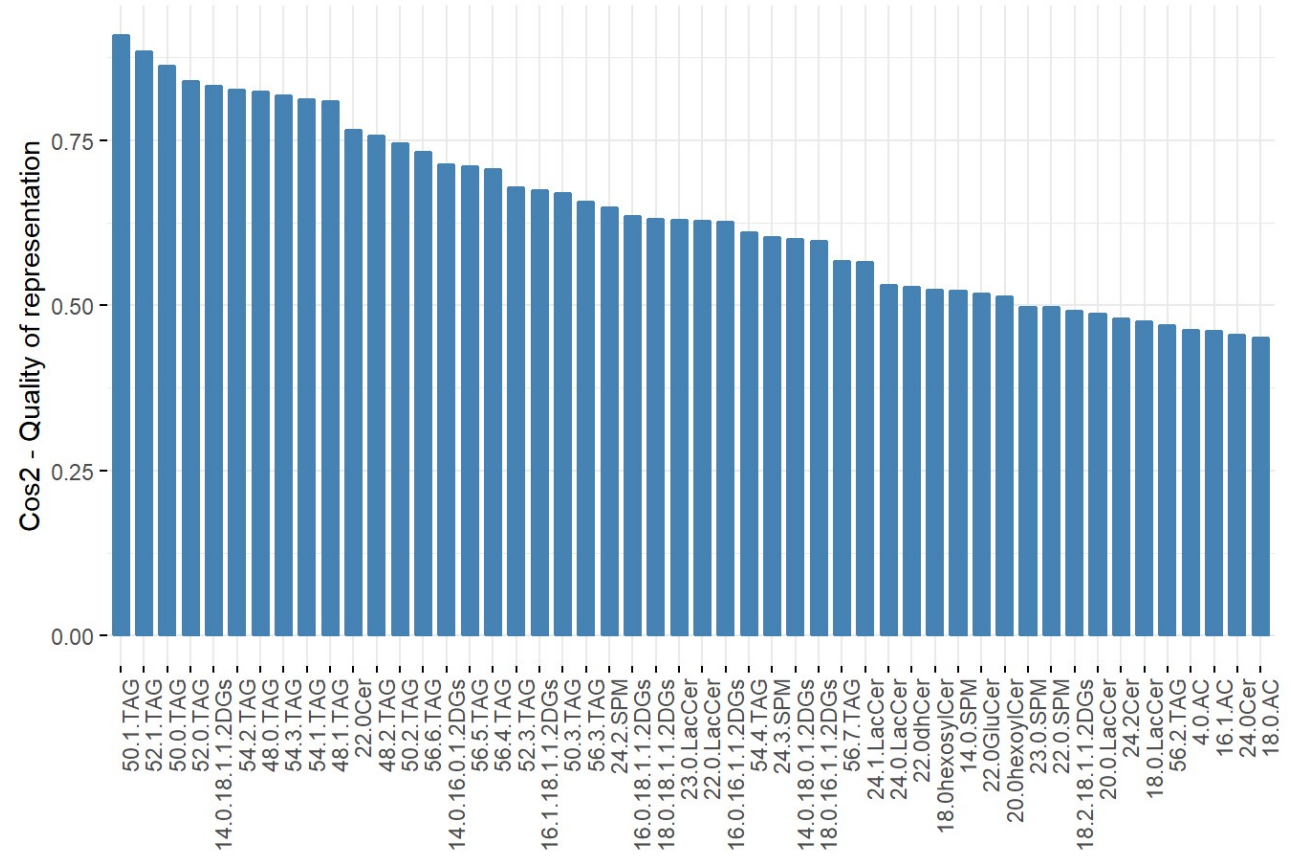
```
qual1 <- fviz_cos2(PCA_36, choice = "var", axes = 1:2, top = 10, xtickslab.rt = 90)+ 1
abs(title = "Variable representation, PC 1:2")
qual2 <- fviz_cos2(PCA_36, choice = "var", axes = 1:10, top = 10, xtickslab.rt = 90)+
labs(title = "Variable representation, PC 1:10")
qual3 <- fviz_cos2(PCA_36, choice = "var", axes = 1:2, top = 50, xtickslab.rt = 90)+ 1
abs(title = "Variable representation, PC 1:2")
qual4 <- fviz_cos2(PCA_36, choice = "var", axes = 1:10, top = 50, xtickslab.rt = 90)+
labs(title = "Variable representation, PC 1:10")
```

```
ggarrange(qual1, qual2, ncol = 2, nrow = 1)
```

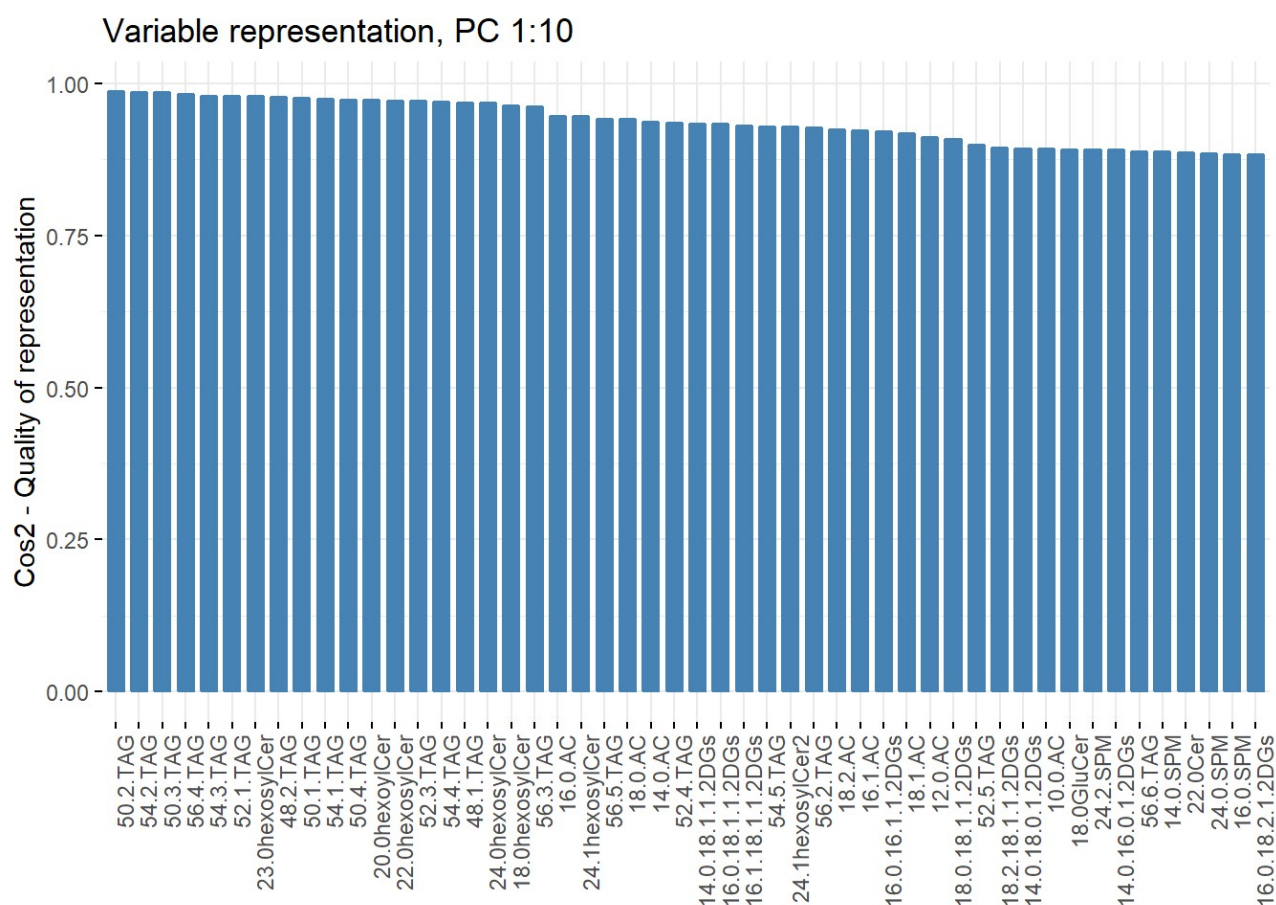


qual3

Variable representation, PC 1:2



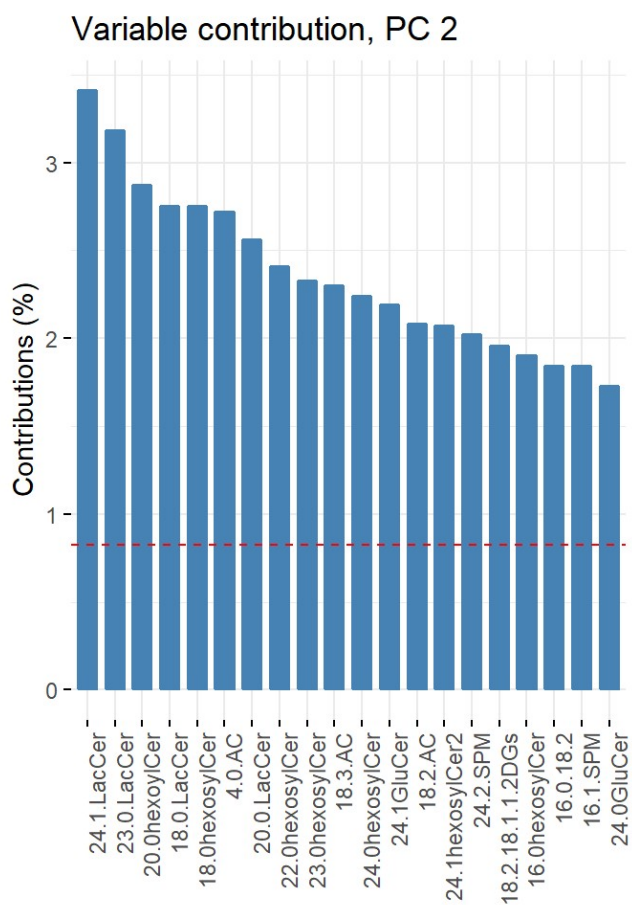
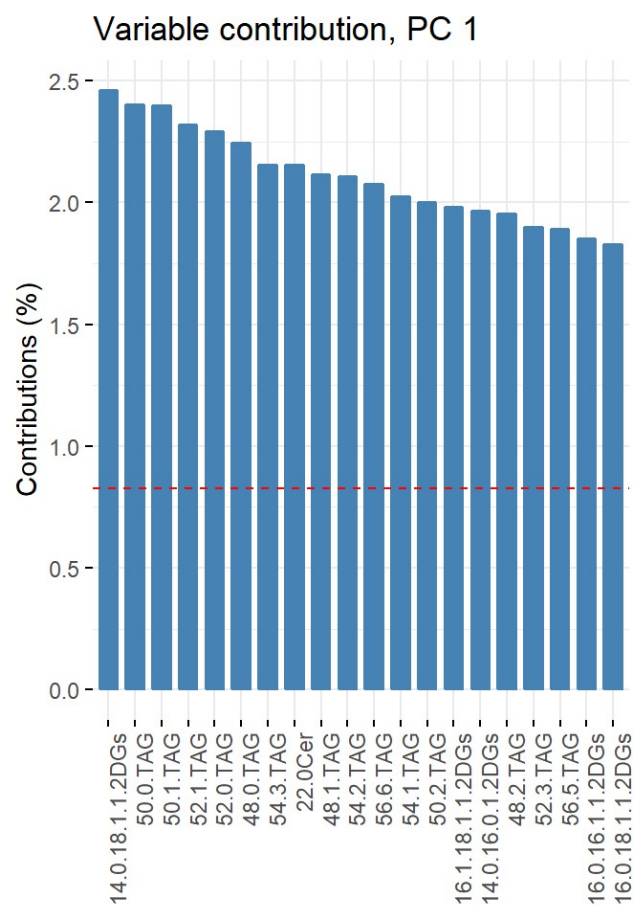
qual14



#Lipid contribution to the first four principal components (as used in the PCA), and out to 10 and 20 principal components (Noted to account for 80%, 95% of variation in the data based on earlier scree plots.)

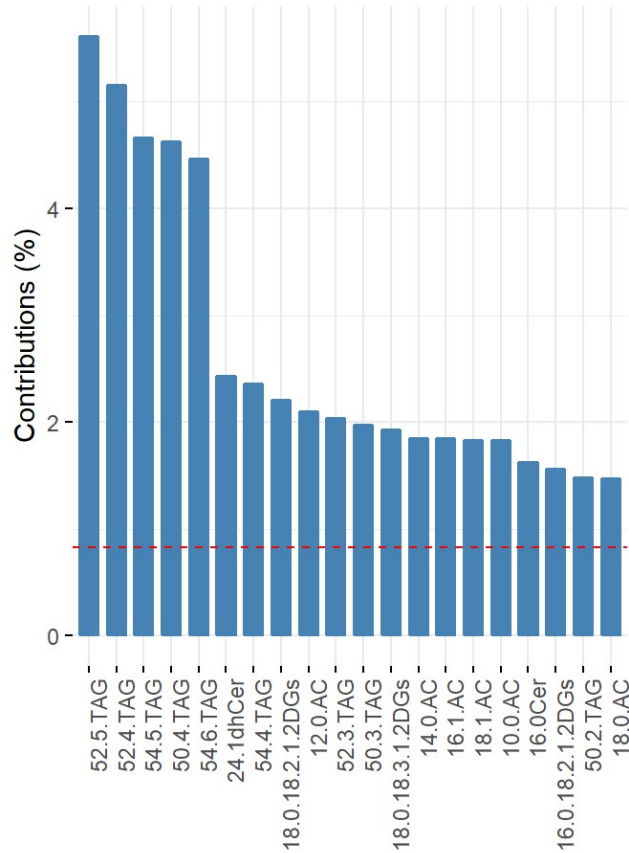
```
con1 <- fviz_contrib(PCA_36, choice = "var", axes = 1, top = 20, xtickslab.rt = 90)+ 1
abs(title = "Variable contribution, PC 1")
con2 <- fviz_contrib(PCA_36, choice = "var", axes = 2, top = 20, xtickslab.rt = 90)+ 1
abs(title = "Variable contribution, PC 2")
con3 <- fviz_contrib(PCA_36, choice = "var", axes = 3, top = 20, xtickslab.rt = 90)+ 1
abs(title = "Variable contribution, PC 3")
con4 <- fviz_contrib(PCA_36, choice = "var", axes = 4, top = 20, xtickslab.rt = 90)+ 1
abs(title = "Variable contribution, PC 4")

ggarrange(con1, con2, ncol = 2, nrow = 1)
```

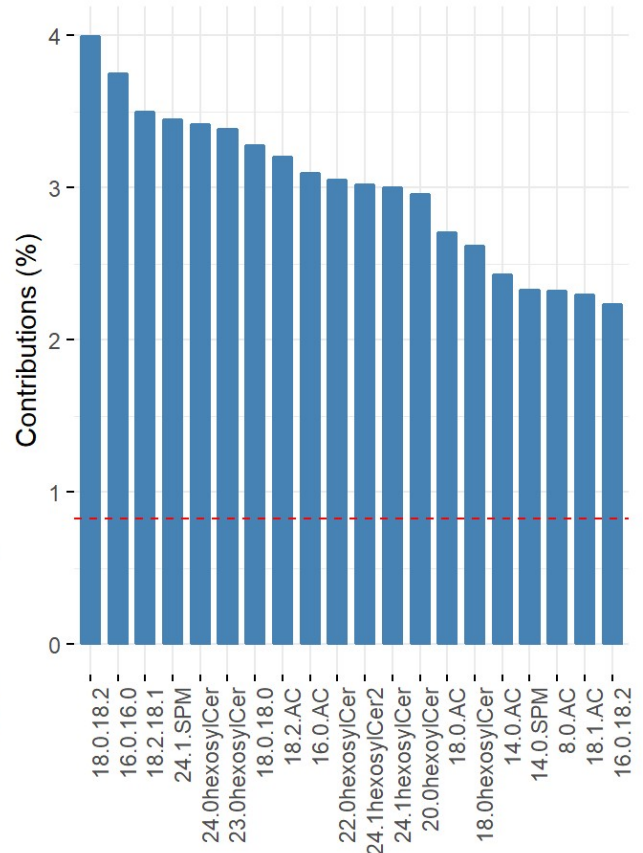


```
ggarrange(con3,con4, ncol = 2, nrow = 1)
```

Variable contribution, PC 3

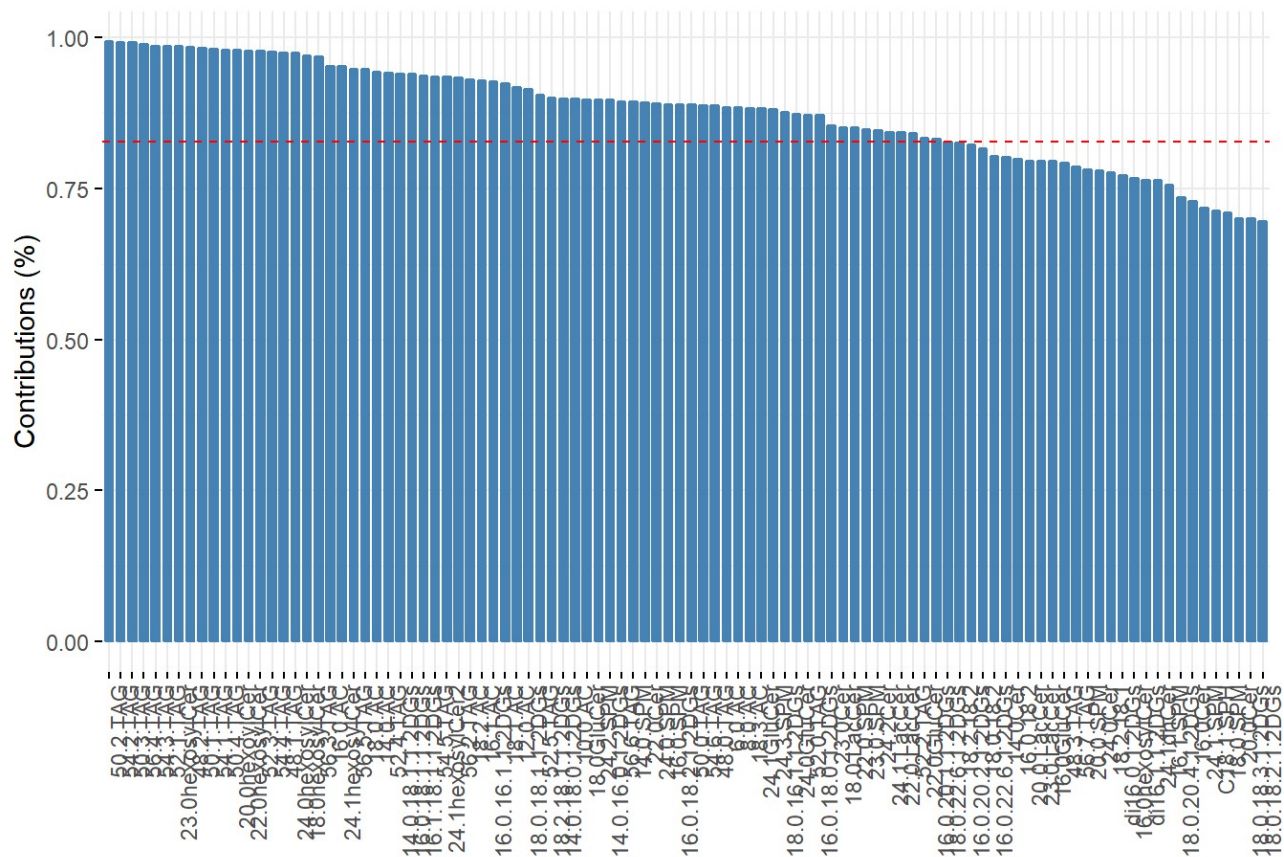


Variable contribution, PC 4



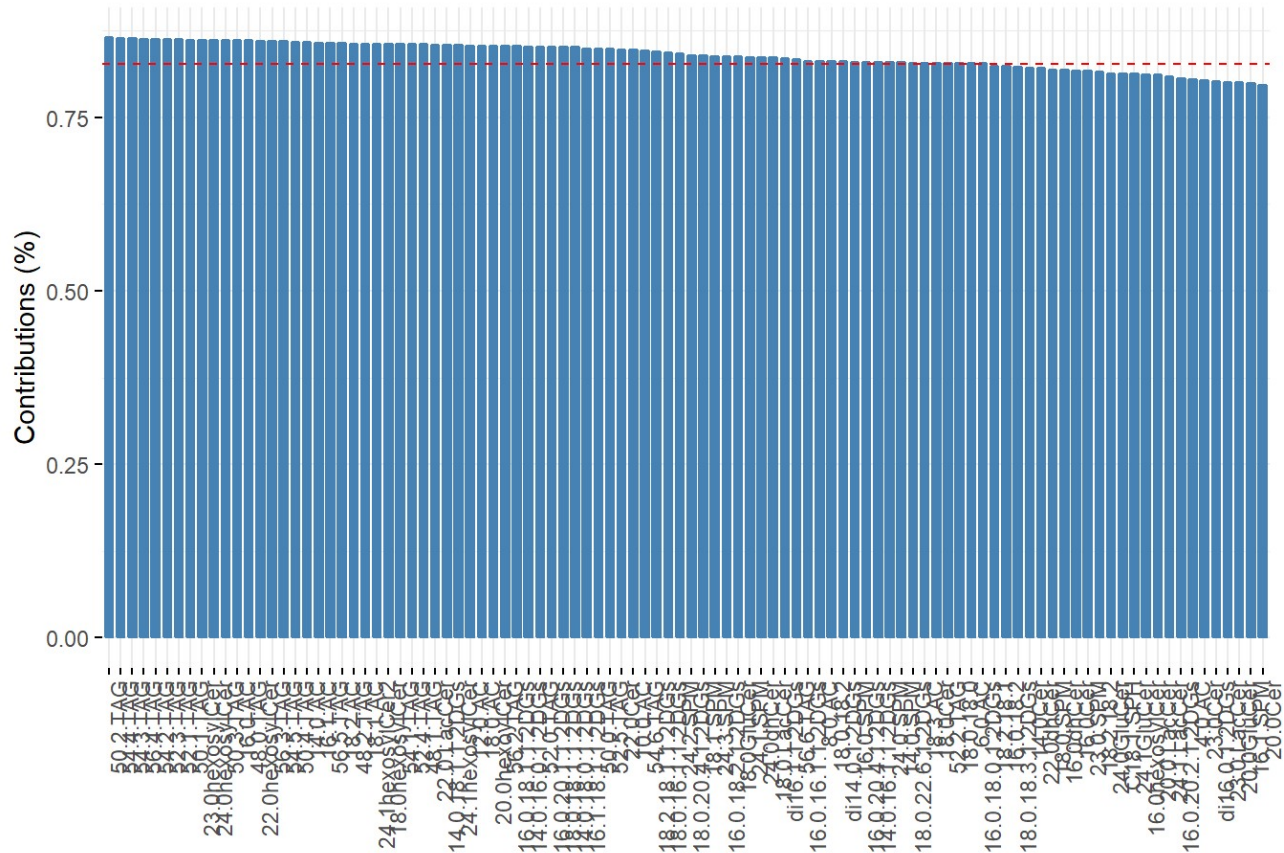
```
fviz_contrib(PCA_36, choice = "var", axes = 1:10, top = 100, xtickslab.rt = 90)+ labs
(title = "Variable contribution, PC 1:10")
```

Variable contribution, PC 1:10



```
fviz_contrib(PCA_36, choice = "var", axes = 1:20, top = 100, xtickslab.rt = 90)+ labs
(title = "Variable contribution, PC 1:20")
```

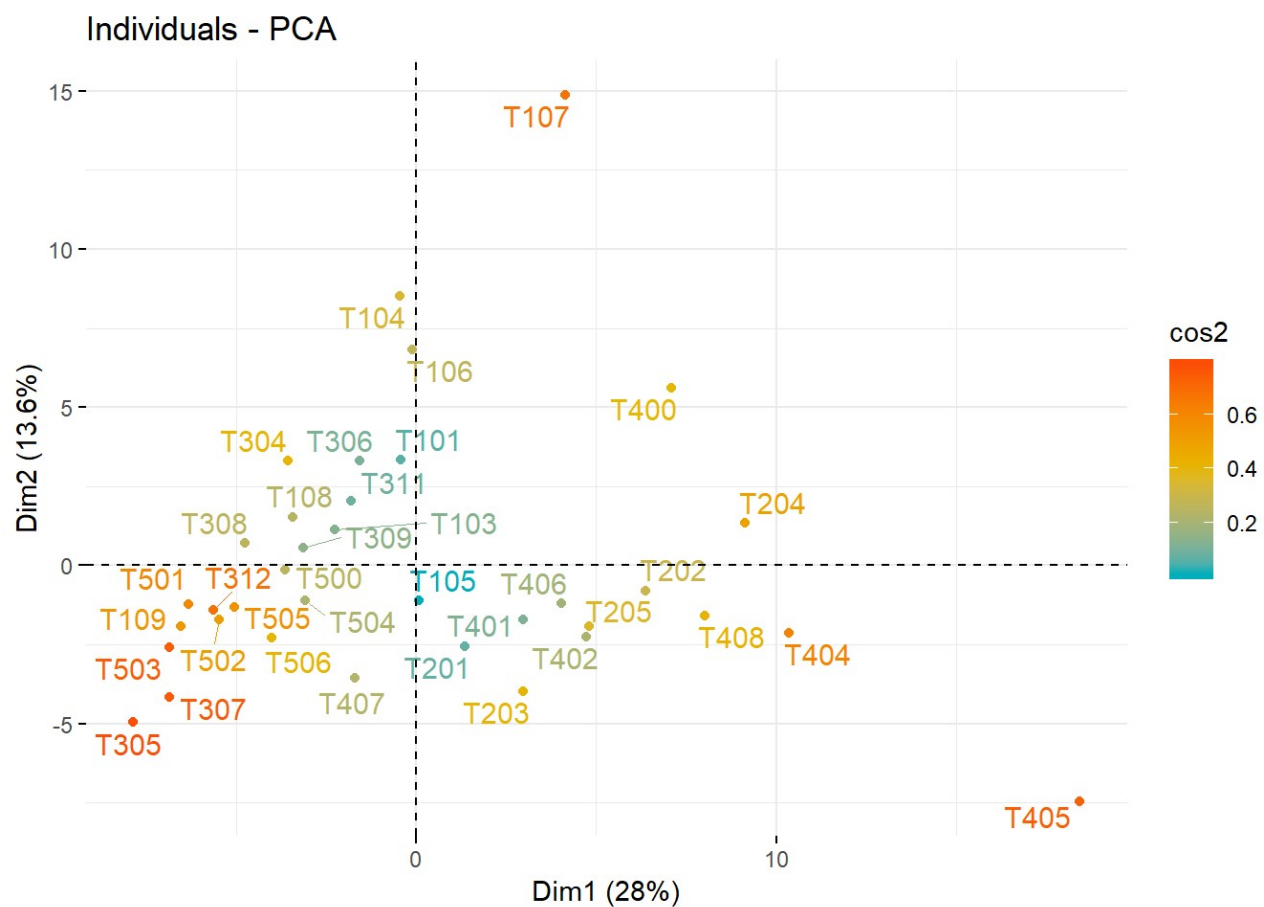
Variable contribution, PC 1:20



#Redlines represent overall mean

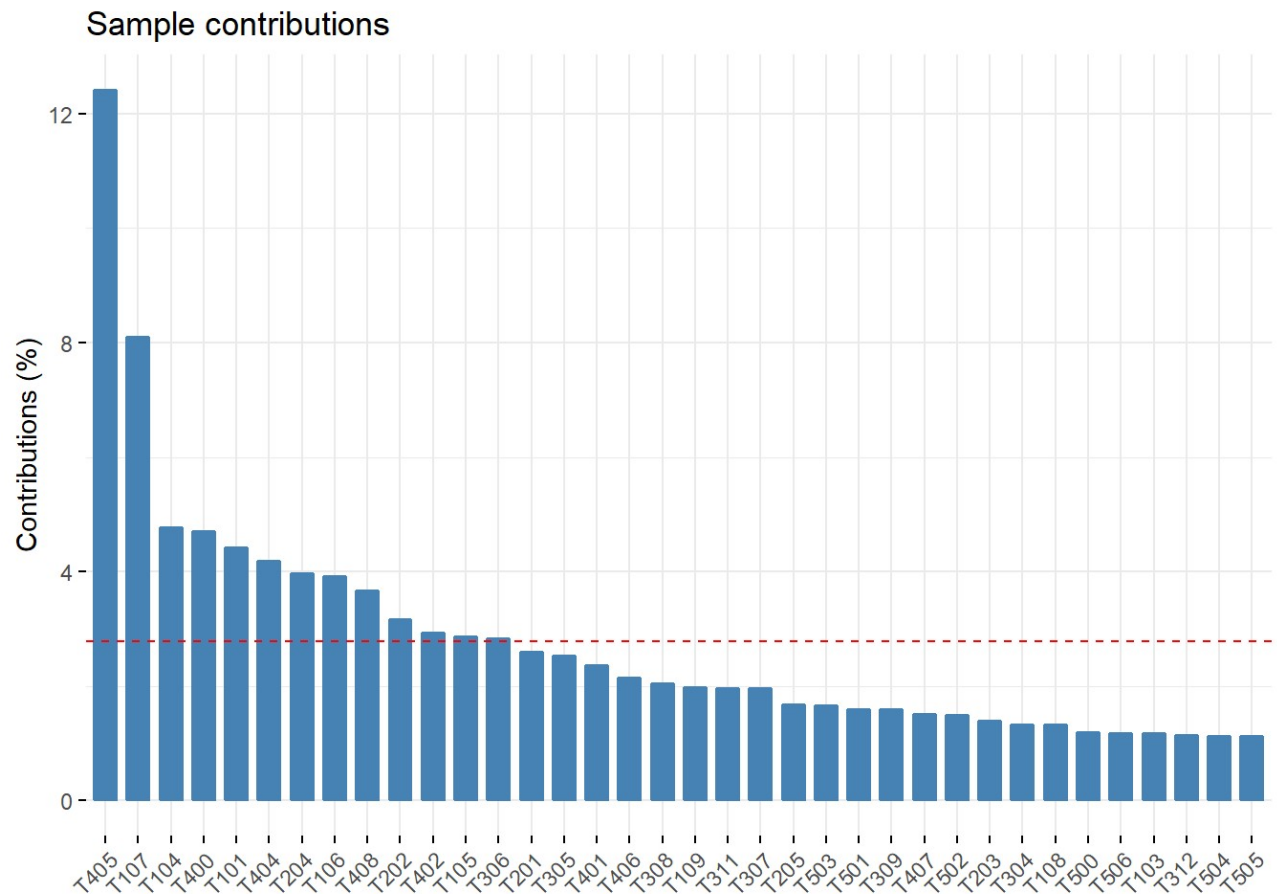
#Labeled PCA with sample names and contribution to PCA from lipids in each sample. Notice that extreme points have the most significant contributions.

```
fviz_pca_ind(PCA_36, col.ind = "cos2",
             gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),
             repel = TRUE # Avoid text overlapping (slow if many points)
             )
```



#Associated sample contribution bar graph

```
fviz_contrib(PCA_36, choice = "ind", axes = 1:35)+ labs(title = "Sample contribution
s")
```



*#var\$cos2: represents the quality of representation for variables on the factor map. It's calculated as the squared coordinates: $\text{var.cos2} = \text{var.coord} * \text{var.coord}$.*

Swapping out the All dataframe with the subsets of HFD and CTR, isolated trends within maternal diet can be discerned as well.