

Project02 - Group01

Eva, Tobi, Kathi, Laura

14 Juni 2019

data loading

```
wd = getwd()

NCI_TPW_gep_treated = readRDS(paste0(wd, "/Data/NCI_TPW_gep_treated.rds"))
NCI_TPW_gep_untreated = readRDS(paste0(wd, "/Data/NCI_TPW_gep_untreated.rds"))
Metadata = read.delim(paste0(wd, "/Data/NCI_TPW_metadata.tsv"), header = TRUE, sep = "\t", stringsAsFactors = TRUE)
Cellline_Annotation = read.delim(paste0(wd, "/Data/cellline_annotation.tsv"), header = TRUE, sep = "\t", stringsAsFactors = TRUE)
Drug_Annotation = read.delim(paste0(wd, "/Data/drug_annotation.tsv"), header = TRUE, sep = "\t", stringsAsFactors = TRUE)
CCLE_mutations = readRDS(paste0(wd, "/Data/CCLE_mutations.rds"))
CCLE_copynumber = readRDS(paste0(wd, "/Data/CCLE_copynumber.rds"))
CCLE_basalexpression = readRDS(paste0(wd, "/Data/CCLE_basalexpression.rds"))
NegLogGI50 = as.data.frame(readRDS(paste0(wd, "/Data/NegLogGI50.rds")))
Treated = data.frame(NCI_TPW_gep_treated)
Untreated = data.frame(NCI_TPW_gep_untreated)
```

data scaling

After checking for normalization, we scaled our data in the first place to provide the scaled data for further analysis.

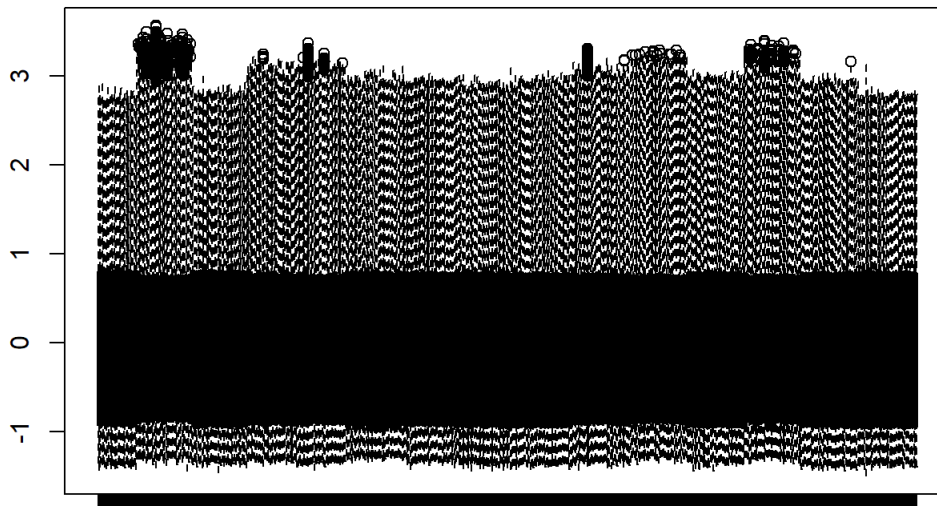
```
list = list(Treated, Untreated)
nlist = lapply(list, scale)
Treated = as.data.frame(nlist[[1]])
Untreated = as.data.frame(nlist[[2]])
Fold_Change = Treated - Untreated
Fold_Change = data.frame(Fold_Change)
rm(NCI_TPW_gep_treated, NCI_TPW_gep_untreated, list, nlist)
```

1. Broad analysis

Boxplots (already normalized)

This step was done before scaling the data. The boxplots showed a deviation which is the reason for scaling the data.

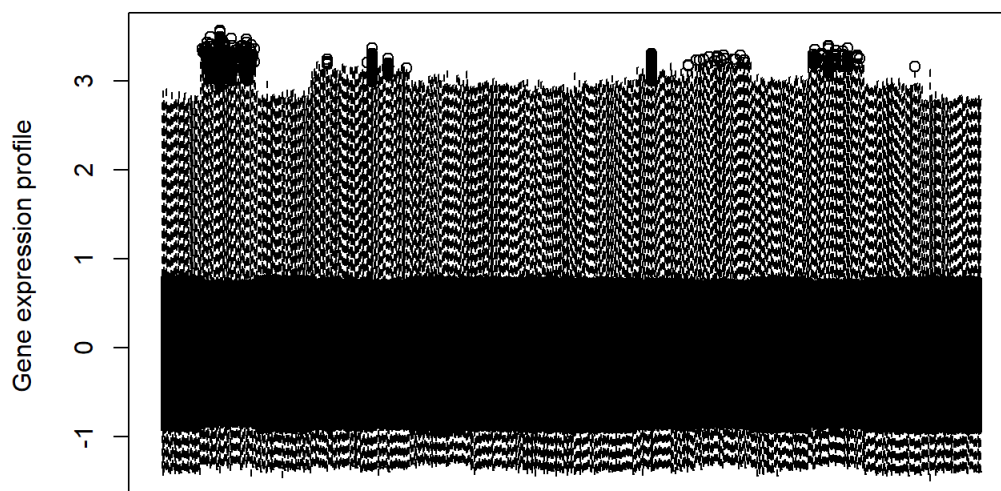
```
boxplot(Treated)
```



6.0_5.Azacytidine_5000nM_24h SR_gemcitabine_2000nM_24h LOX_vorinostat_5000nM_2

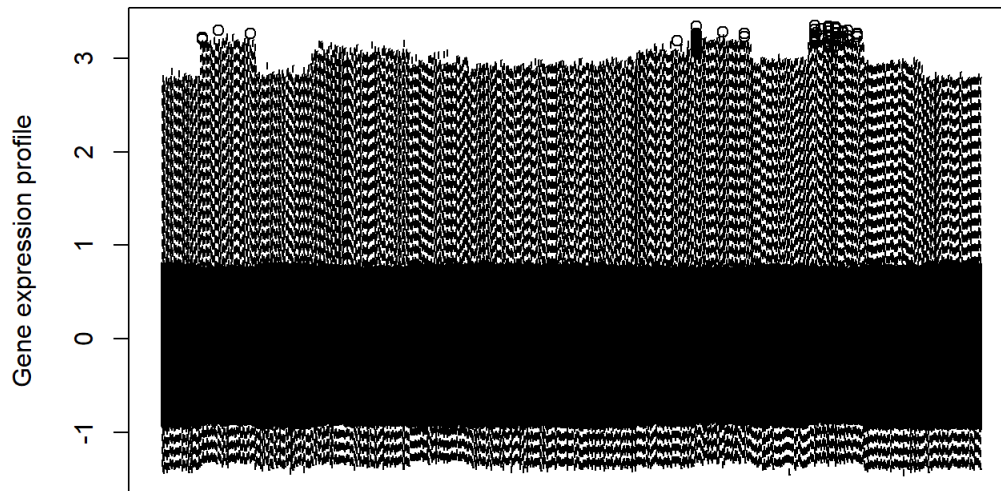
```
boxplot(Treated, ylab = "Gene expression profile", main = "Treated genexpressionprofiles", xaxt = "n")
```

Treated genexpressionprofiles



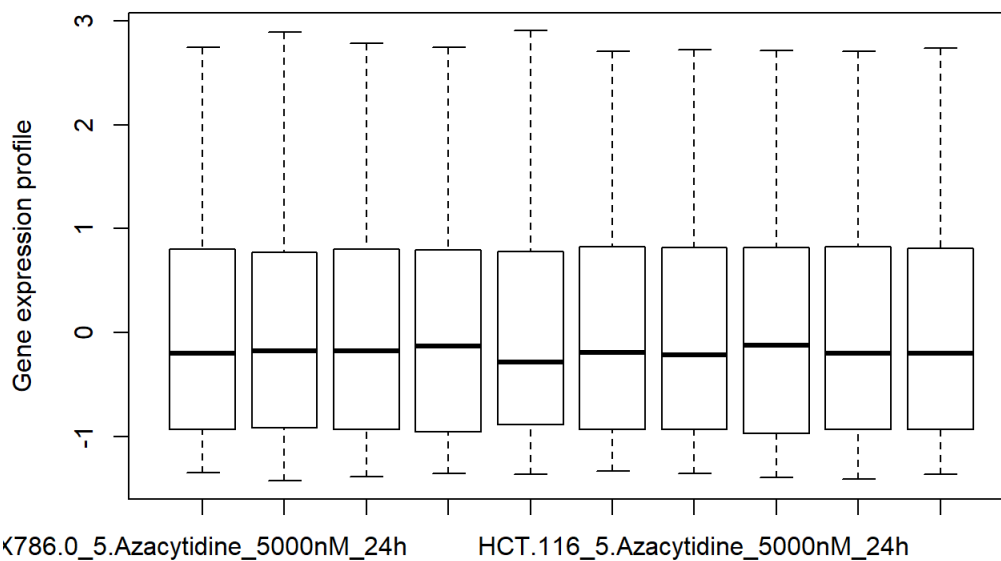
```
boxplot(Untreated, ylab = "Gene expression profile", main = "Untreated genexpressionprofiles", xaxt = "n")
```

Untreated genexpressionprofiles



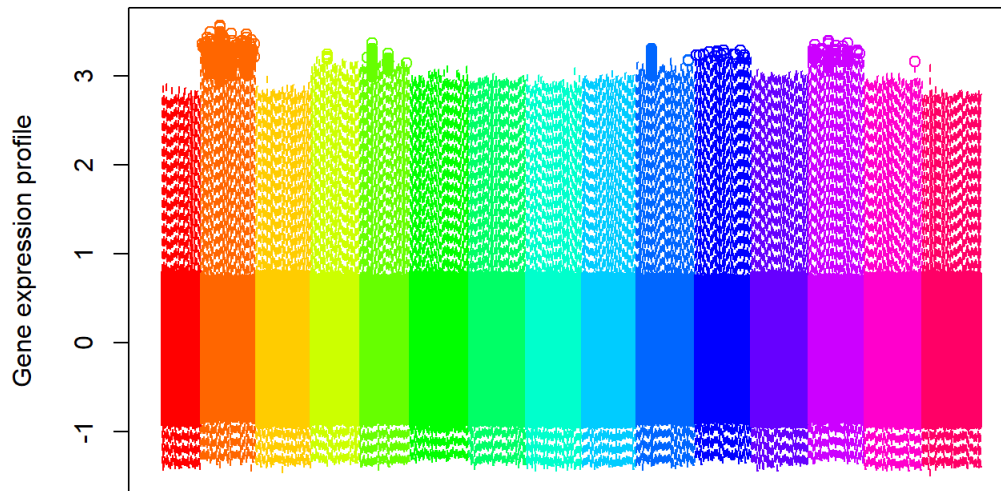
```
boxplot(Treated[,1:10], ylab = "Gene expression profile", main = "First 10 reated genexpressionprofiles")
```

First 10 reated genexpressionprofiles



```
Treated1 = readRDS(paste0(wd, "/Data/NCI_TPW_gep_treated.rds"))
df = data.frame(t(Treated1))
df.data <- data.frame(sample = rownames(df))
adjustedMeda = subset(Metadata, sample %in% intersect(Metadata$sample, df.data$sample))
rm(df,df.data, Treated1)
palette(rainbow(15))
boxplot(Treated, border=adjustedMeda$drug,xlab= "Different Drugs",ylab = "Gene expression profile", main = "Teated genexpressionprofiles",xaxt ="n")
```

Teated genexpressionprofiles



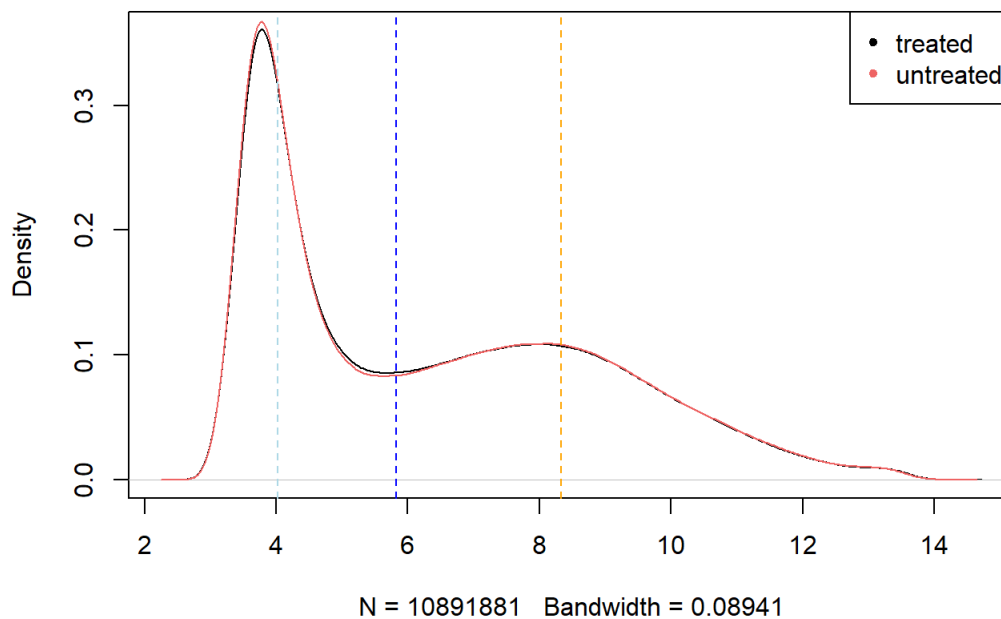
Different Drugs

Densityplot

The abline shows the 3 quantiles (25% 50% 75%)

```
NCI_TPW_gep_treated = readRDS(paste0(wd, "/Data/NCI_TPW_gep_treated.rds"))
NCI_TPW_gep_untreated = readRDS(paste0(wd, "/Data/NCI_TPW_gep_untreated.rds"))
plot(density(NCI_TPW_gep_treated), "Densityplot Treated vs Untreated")
lines(density(NCI_TPW_gep_untreated), col = "indianred2")
legend("topright", legend = c("treated", "untreated"), col = c("black", "indianred2"), pch = 20)
abline(v = quantile(NCI_TPW_gep_treated)[2:4], col = c("lightblue", "blue", "orange"), lty = 2)
```

Densityplot Treated vs Untreated



k-means clustering

To look for clusters in the raw data we performed a k-means clustering and searched for potentially clusters.

```
# Performing a k-means on Treated
#Determining the number of clusters
topVarTreated = apply(Treated, 1, var)
summary(topVarTreated)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.002893 0.029461 0.069002 0.124300 0.135476 2.138284
```

```
# Using the most variable, thus informative genes
topVarTreated75 = Treated[topVarTreated > quantile(topVarTreated, probs = 0.75), ]
dim(topVarTreated75)
```

```
## [1] 3325 819
```

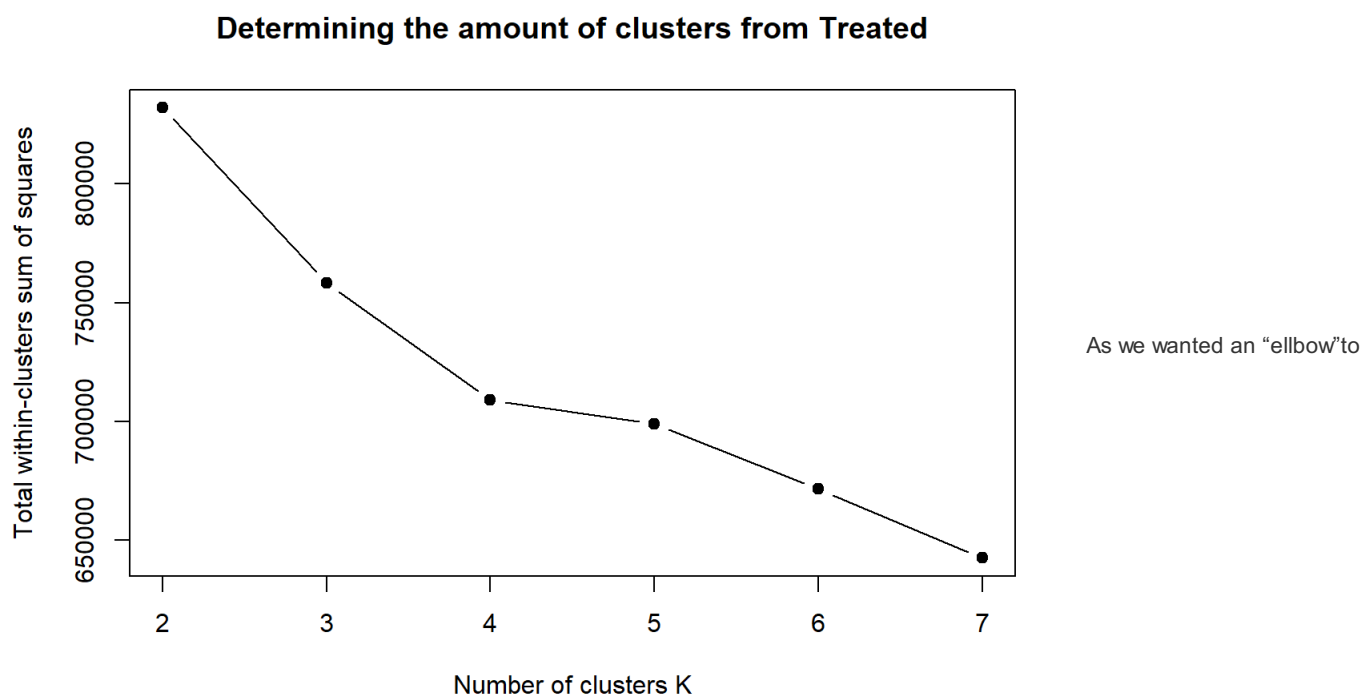
```
km = kmeans(x = t(topVarTreated75), centers = 3, nstart = 10)
km$tot.withinss
```

```
## [1] 758323.6
```

```
km = kmeans(x = t(topVarTreated75), centers = 2, nstart = 10)
km$tot.withinss
```

```
## [1] 832093.5
```

```
#running a loop for the best n (searching for "ellbow")
wss = sapply(2:7, function(k) {
  kmeans(x = t(topVarTreated75), centers = k)$tot.withinss})
plot(2:7, wss, type = "b", pch = 19, xlab = "Number of clusters K", ylab = "Total within-clusters sum of squares", main = "Determining the amount of clusters from Treated")
```



get a good result we can say in a way that our data are not really good to cluster. To look in a other way, we also provided the clusters by

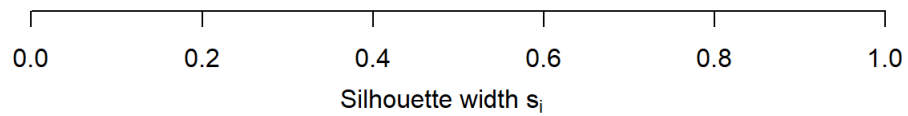
Silhouette plot of (x = km\$cluster, dist = D)

n = 819

10 clusters C_j

j	n_j	$\text{ave}_{i \in C_j} s_i$
1	101	0.08
2	44	0.19
3	102	0.19
4	110	0.05
5	27	0.36
6	108	0.23
7	62	0.20
8	130	0.06
9	54	0.13
10	81	0.10

the silhouette-method.



Average silhouette width : 0.14

PCA

```
pca <- prcomp(t(Fold_Change), scale = TRUE)
```

```
# sdev calculates variation each PC accounts for
pca.var <- pca$sdev^2
# since percentages make more sense than normal variation values
# calculate % or variation, which is much more interesting
pca.var.per <- round(pca.var/sum(pca.var)*100, 1)

barplot(pca.var.per, main = "Scree plot", xlab = "Principal Components", ylab = "% variation")
```

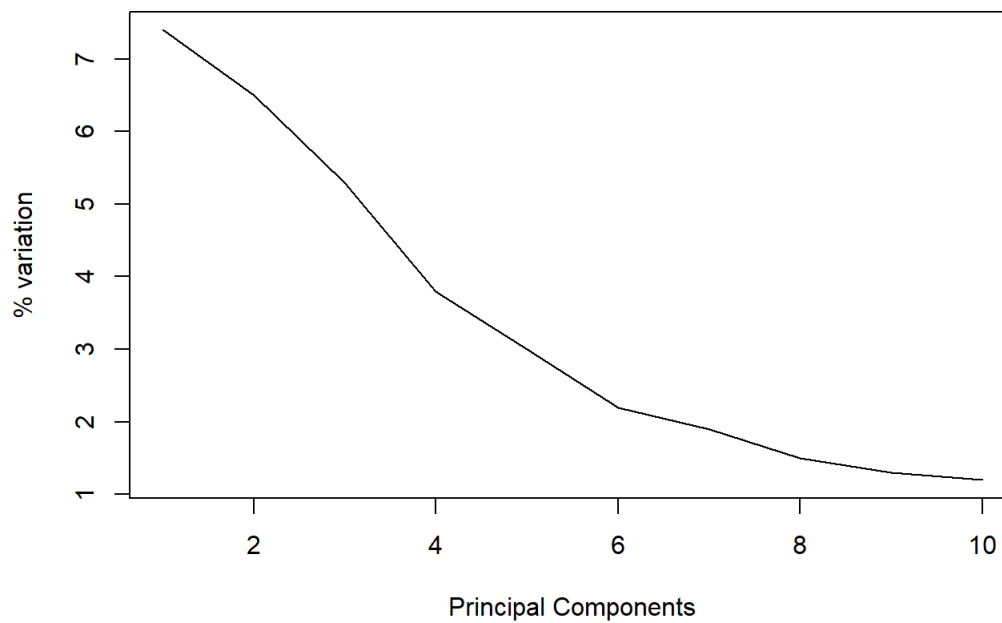
Scree plot



Principal Components

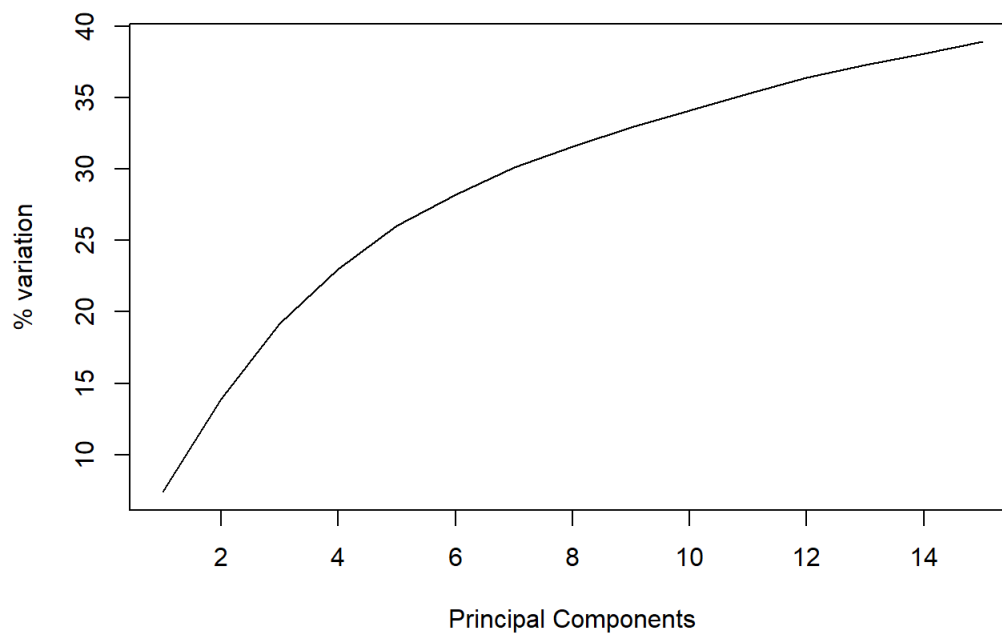
```
plot(pca.var.per[1:10], main = "Elbow plot", type = "l", xlab = "Principal Components", ylab = "% variation")
```

Elbow plot



```
plot(cumsum(pca.var.per[1:15]), main = "cumulative variation", type = "l", xlab = "Principal Components",  
ylab = "% variation")
```

cumulative variation



```
#creating data frame with all pcs  
#cleaning up sample names as they differed between matrices  
pca.data <- data.frame(pca$x)  
rownames(pca.data) <- gsub(x = rownames(pca.data), pattern = "x786", replacement = "786")  
pca.data <- cbind(sample = rownames(pca.data), pca.data)
```

```
## get names of top 10 genes that contribute most to pc1
loading_scores_1 <- pca$rotation[,1]
gene_score <- abs(loading_scores_1) ## sort magnitude
gene_score_ranked <- sort(gene_score, decreasing = TRUE)

top_10_genes <- names(gene_score_ranked[1:10])
top_10_genes # show names of top 10 genes
```

```
## [1] "DNAJC2" "NGDN" "GTPBP4" "CCDC59" "DNTTIP2" "AKAP8" "PAPSS1"
## [8] "TRMT1" "BRF2" "YRDC"
```

```
### Metadata color matrix for coloring
Metadata$sample <- gsub(x = Metadata$sample, pattern = "-", replacement = ".")

metad.cl <- subset(Metadata, Metadata$sample %in% pca.data$sample)
## adjust row length of metadata to pca.data

metad.cl$mechanism <- Drug_Annotation$Mechanism[match(metad.cl$drug, Drug_Annotation$Drug)]
metad.cl$msi <- Cellline_Annotation$Microsatellite_instability_status[match(metad.cl$cell, Cellline_Annotation$Cell_Line_Name)]
```

```
# plotting all informative PCs
#color vectors for coloring by drug and tissue
viridis <- viridis(9)
color_tissue = viridis[metad.cl$tissue]
tissue <- levels(metad.cl$tissue)

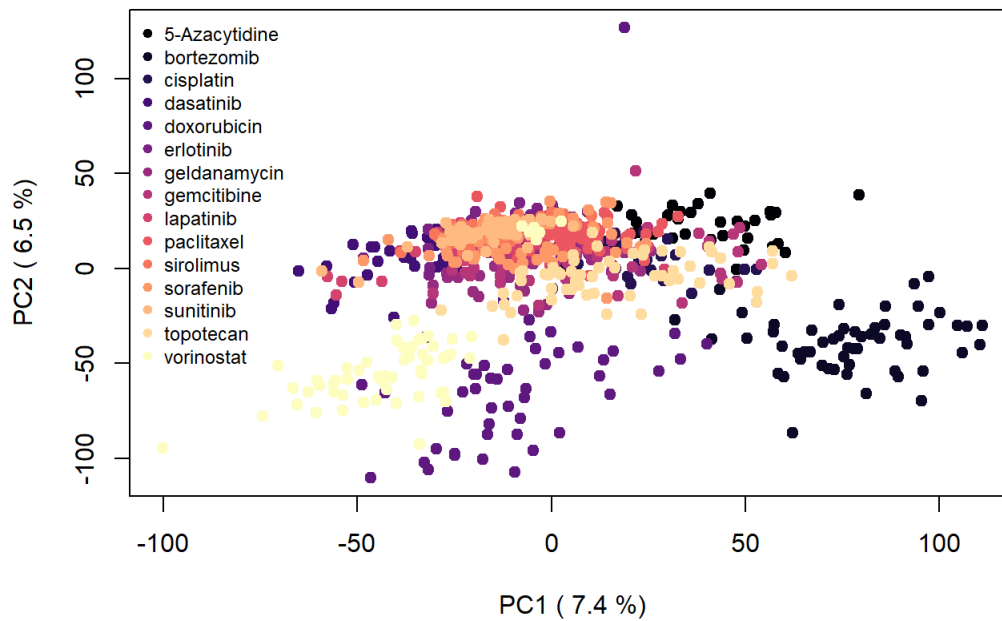
magma <- magma(15)
color_drug = magma[metad.cl$drug]
drug <- levels(metad.cl$drug)

## colored by drug
#plot PC1 and PC2
plot(pca$x[,1],
     pca$x[,2],
     col = color_drug,
     pch = 19,
     xlab = paste("PC1 (",pca.var.per[1],"%)" ),
     ylab = paste("PC2 (",pca.var.per[2],"%)" ) )
#create legend
legend("topleft",
      legend = drug,
      col = magma,
      pch = 19,
      xpd = "TRUE",
      bty = "n",
      cex = 0.75
)
```

```
## Warning in par(xpd = xpd): NAs durch Umwandlung erzeugt
```

```
#create title
mtext("PCA of Fold Change colored by drug",
      side = 3,
      line = -2,
      cex = 1.2,
      font = 2,
      outer = TRUE)
```


PCA of Fold Change colored by drug

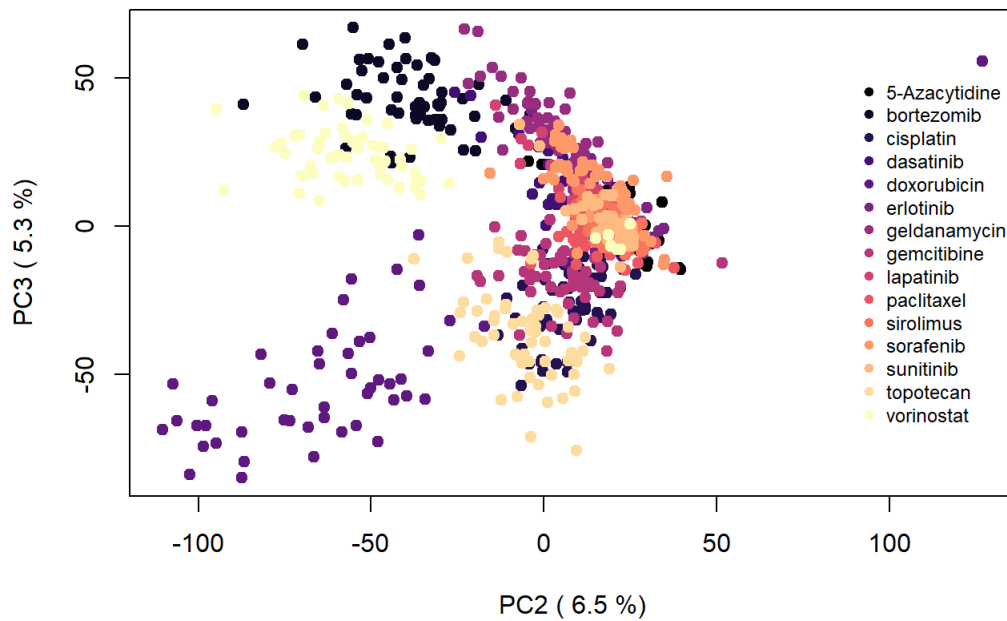


```
#plot PC2 and PC3
plot(pca$x[,2],
     pca$x[,3],
     col = color_drug,
     pch = 19,
     xlab = paste("PC2 (",pca.var.per[2],"%)" ),
     ylab = paste("PC3 (",pca.var.per[3],"%)" ))
#create legend
legend("right",
      legend = drug,
      col = magma,
      pch = 19,
      xpd = "TRUE",
      bty = "n",
      cex = 0.75,
      inset = c(0, 2)
)
```

```
## Warning in par(xpd = xpd): NAs durch Umwandlung erzeugt
```

```
#create title
mtext("PCA of Fold Change colored by drug",
      side = 3,
      line = -2,
      cex = 1.2,
      font = 2,
      outer = TRUE)
```

PCA of Fold Change colored by drug

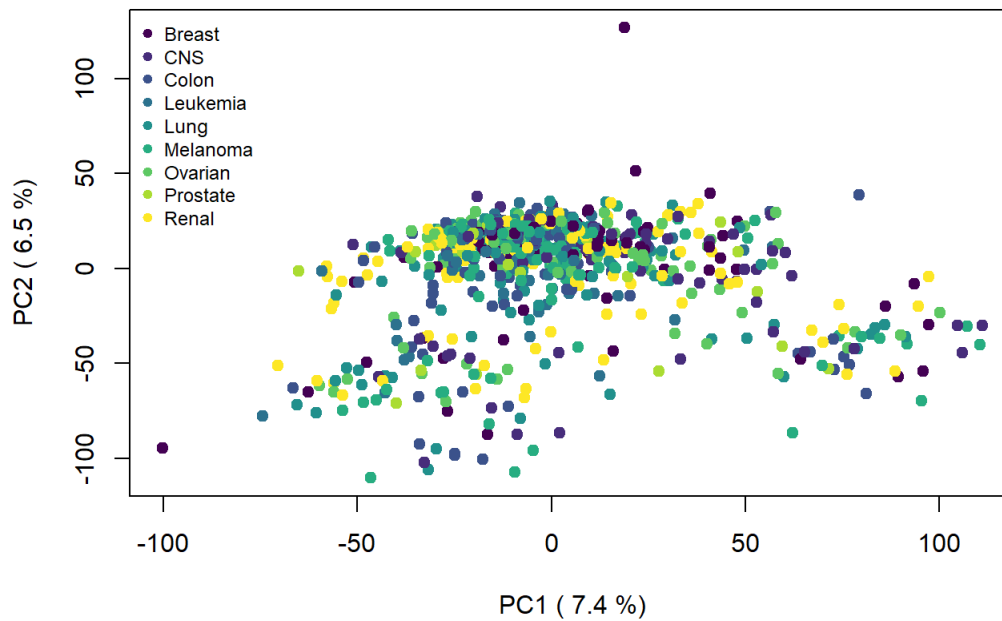


```
## colored by tissue
#plot PC1 and PC2
plot(pca$x[,1],
     pca$x[,2],
     col = color_tissue,
     pch = 19,
     xlab = paste("PC1 (",pca.var.per[1],"%)" ),
     ylab = paste("PC2 (",pca.var.per[2],"%)" ))
#create legend
legend("topleft",
      legend = tissue,
      col = viridis,
      pch = 19,
      xpd = "TRUE",
      bty = "n",
      cex = 0.75
)
```

```
## Warning in par(xpd = xpd): NAs durch Umwandlung erzeugt
```

```
#create title
mtext("PCA of Fold Change colored by tissue",
      side = 3,
      line = -2,
      cex = 1.2,
      font = 2,
      outer = TRUE)
```

PCA of Fold Change colored by tissue



```
#plot PC2 and PC3
plot(pca$x[,2],
     pca$x[,3],
     col = color_tissue,
     pch = 19,
     xlab = paste("PC2 (",pca.var.per[2],"%)",
     ylab = paste("PC3 (",pca.var.per[3],"%)",
#create legend
legend("right",
      legend = tissue,
      col = viridis,
      pch = 19,
      xpd = "TRUE",
      bty = "n",
      cex = 0.75,
      inset = c(0, 2)
)
```

```
## Warning in par(xpd = xpd): NAs durch Umwandlung erzeugt
```

```
#create title
mtext("PCA of Fold Change colored by tissue",
      side = 3,
      line = -2,
      cex = 1.2,
      font = 2,
      outer = TRUE)
```

PCA of Fold Change colored by tissue

