

Geometrical Unsupervised Learning of Breast Cancer SubTyping

2024-05-09

Abstract

Molecular subtyping of cancer is recognized as a critical and challenging step towards individualized therapy. Breast cancer molecular subtyping has revolutionized the treatment and lowered the mortality rate.

In the study, we further analyse breast cancer's subtyping using geometrical unsupervised learning. We aim to find a new way to examine the gene expression data by applying on it a different topology. The topology is applied by using different distance function on the sample space.

In the study we used 5 metrical functions - Pairwise Correlation, Cosine Similarity, Kernel Distance, PCA and Euclidean distance for the control analyses.

- Pairwise correlation measures the linear relationship between two variables.
- Cosine similarity measures the cosine of the angle between two non-zero vectors.
- Kernel distance measures nonlinear relationship using a kernel function.
- Principal Component Analysis (PCA) is a dimensionality reduction technique, preserving as much variance as possible. After the reduction of the samples to one dimensional variable we measured the distance between the samples using the Euclidean distance.

The research contains RNA sequence samples taken from TCGA dataset. The data contains 313 samples. The samples are divided by PAM50 classification to 35 “Normal” samples, 68 “Her2” samples and 70 “Luminal A”, “Luminal B” and “Basal” samples.

The study pipeline is Desq2 analyses, heatmap and unsupervised hierarchical clustering, survival analyses and pathways enrichments analyses.

Each of the 5 metrics processed through the pipeline. The Euclidean distance, Pairwise Correlation, Cosine Similarity had similar results. Those results are consistent with the PAM50 common subtyping. The kernel distance did not form clusters.

The PCA had different clustering results.

The PCA resulted 3 clusters. Those clusters received the lowest p value in the survival analyses, yet nonsignificant.

The pathway enrichment analyses didn’t contribute difference to any of the clustering of all metrics.

A deeper understanding of the PCA clustering can be done to understand if the results are significant.

Possible forward research can be to repeat the analyses with different datasets and to further understand the clusters division.

```
setwd('C:/technion/Bio/project/')  
library(TCGAbiolinks) # download the TCGA data  
  
library(DESeq2)
```

```
## Warning: package 'DESeq2' was built under R version 4.3.3
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
## Loading required package: BiocGenerics
```

```
##  
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:stats':  
##  
## IQR, mad, sd, var, xtabs
```

```
## The following objects are masked from 'package:base':  
##  
## anyDuplicated, aperm, append, as.data.frame, basename, cbind,  
## colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,  
## get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,  
## match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,  
## Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,  
## table, tapply, union, unique, unsplit, which.max, which.min
```

```
##  
## Attaching package: 'S4Vectors'
```

```
## The following object is masked from 'package:utils':  
##  
## findMatches
```

```
## The following objects are masked from 'package:base':  
##  
## expand.grid, I, unname
```

```
## Loading required package: IRanges
```

```
##  
## Attaching package: 'IRanges'
```

```
## The following object is masked from 'package:grDevices':  
##  
## windows
```

```
## Loading required package: GenomicRanges
```

```
## Loading required package: GenomeInfoDb
```

```
## Loading required package: SummarizedExperiment
```

```
## Loading required package: MatrixGenerics
```

```
## Loading required package: matrixStats
```

```
##  
## Attaching package: 'MatrixGenerics'
```

```
## The following objects are masked from 'package:matrixStats':
##
##   colAlls, colAnyNAs, colAnys, colAvgPerRowSet, colCollapse,
##   colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
##   colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
##   colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
##   colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
##   colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
##   colWeightedMeans, colWeightedMedians, colWeightedSds,
##   colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgPerColSet,
##   rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
##   rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
##   rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
##   rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
##   rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
##   rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
##   rowWeightedSds, rowWeightedVars
```

```
## Loading required package: Biobase
```

```
## Welcome to Bioconductor
##
##   Vignettes contain introductory material; view with
##   'browseVignettes()'. To cite Bioconductor, see
##   'citation("Biobase")', and for packages 'citation("pkgname")'.
```

```
##
## Attaching package: 'Biobase'
```

```
## The following object is masked from 'package:MatrixGenerics':
##
##   rowMedians
```

```
## The following objects are masked from 'package:matrixStats':
##
##   anyMissing, rowMedians
```

```
library(SummarizedExperiment)
```

```
library(ComplexHeatmap) # Heatmap (works with personalize distance function)
```

```
## Loading required package: grid
```

```
## =====
## ComplexHeatmap version 2.18.0
## Bioconductor page: http://bioconductor.org/packages/ComplexHeatmap/
## Github page: https://github.com/jokergoo/ComplexHeatmap
## Documentation: http://jokergoo.github.io/ComplexHeatmap-reference
##
## If you use it in published research, please cite either one:
## - Gu, Z. Complex Heatmap Visualization. iMeta 2022.
## - Gu, Z. Complex heatmaps reveal patterns and correlations in multidimensional
##   genomic data. Bioinformatics 2016.
##
##
## The new InteractiveComplexHeatmap package can directly export static
## complex heatmaps into an interactive Shiny app with zero effort. Have a try!
##
## This message can be suppressed by:
##   suppressPackageStartupMessages(library(ComplexHeatmap))
## =====
```

```
library(dendsort) # heatmap reordering columns
```

```
## Warning: package 'dendsort' was built under R version 4.3.3
```

```
library(factoextra) # elbow method
```

```
## Loading required package: ggplot2
```

```
## Warning: package 'ggplot2' was built under R version 4.3.3
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
library(survival) # survival analyses  
library(survminer) # survival analyses
```

```
## Loading required package: ggpubr
```

```
##  
## Attaching package: 'survminer'
```

```
## The following object is masked from 'package:survival':  
##  
##      myeloma
```

```
library(msigdb) # enrichment analysis  
library(clusterProfiler) # enrichment analysis
```

```
## Warning: package 'clusterProfiler' was built under R version 4.3.3
```

```
##
```

```
## clusterProfiler v4.10.1 For help: https://yulab-smu.top/biomedical-knowledge-mining-book/  
##  
## If you use clusterProfiler in published research, please cite:  
## T Wu, E Hu, S Xu, M Chen, P Guo, Z Dai, T Feng, L Zhou, W Tang, L Zhan, X Fu, S Liu, X Bo, and G Yu. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. The Innovation. 2021, 2(3):100141
```

```
##  
## Attaching package: 'clusterProfiler'
```

```
## The following object is masked from 'package:IRanges':  
##  
##      slice
```

```
## The following object is masked from 'package:S4Vectors':  
##  
##      rename
```

```
## The following object is masked from 'package:stats':  
##  
##      filter
```

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following object is masked from 'package:Biobase':
##
##      combine
```

```
## The following object is masked from 'package:matrixStats':
##
##      count
```

```
## The following objects are masked from 'package:GenomicRanges':
##
##      intersect, setdiff, union
```

```
## The following object is masked from 'package:GenomeInfoDb':
##
##      intersect
```

```
## The following objects are masked from 'package:IRanges':
##
##      collapse, desc, intersect, setdiff, slice, union
```

```
## The following objects are masked from 'package:S4Vectors':
##
##      first, intersect, rename, setdiff, setequal, union
```

```
## The following objects are masked from 'package:BiocGenerics':
##
##      combine, intersect, setdiff, union
```

```
## The following objects are masked from 'package:stats':
##
##      filter, lag
```

```
## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

```
library(matrixStats)
library(circlize) # Heatmaps color
```

```
## Warning: package 'circlize' was built under R version 4.3.3
```

```
## =====
## circlize version 0.4.16
## CRAN page: https://cran.r-project.org/package=circlize
## Github page: https://github.com/jokergoo/circlize
## Documentation: https://jokergoo.github.io/circlize\_book/book/
##
## If you use it in published research, please cite:
## Gu, Z. circlize implements and enhances circular visualization
##   in R. Bioinformatics 2014.
##
## This message can be suppressed by:
##   suppressPackageStartupMessages(library(circlize))
## =====
```

```
library(RColorBrewer)
library(scales)
```

```
## Warning: package 'scales' was built under R version 4.3.3
```

Download TCGA Data of BRCA RNA-seq

```
query_TCGA = GDCquery(
  project = "TCGA-BRCA",
  data.category = "Transcriptome Profiling",
  experimental.strategy = "RNA-Seq",
  workflow.type = "STAR - Counts", #raw counts
  data.type = 'Gene Expression Quantification',
  sample.type = c("Primary Tumor", "Solid Tissue Normal"))#removed metastatic tumor
```

```
## -----
```

```
## o GDCquery: Searching in GDC database
```

```
## -----
```

```
## Genome of reference: hg38
```

```
## -----
```

```
## oo Accessing GDC. This might take a while...
```

```
## -----
```

```
## ooo Project: TCGA-BRCA
```

```
## -----
```

```
## oo Filtering results
```

```
## -----
```

```
## ooo By experimental.strategy
```

```
## ooo By data.type
```

```
## ooo By workflow.type
```

```
## ooo By sample.type
```

```
## -----
```

```
## oo Checking data
```

```
## -----
```

```
## ooo Checking if there are duplicated cases
```

```
## ooo Checking if there are results for the query
```

```
## -----
```

```
## o Preparing output
```

```
## -----
```

```
GDCdownload(query = query_TCGA)
```

```
## Downloading data for project TCGA-BRCA
```

```
## Of the 1224 files for download 1224 already exist.
```

```
## All samples have been already downloaded
```

```
data = GDCprepare(query_TCGA)
```

```
## |=====|100% Completed after 2 m
```

```
## Starting to add information to samples
```

```
## => Add clinical information to samples
```

```
## => Adding TCGA molecular information from marker papers
```

```
## => Information will have prefix 'paper_'
```

```
## brca subtype information from:doi.org/10.1016/j.ccell.2018.03.014
```

```
## Available assays in SummarizedExperiment :  
##   => unstranded  
##   => stranded_first  
##   => stranded_second  
##   => tpm_unstrand  
##   => fpkm_unstrand  
##   => fpkm_uq_unstrand
```

Extract metadata and Counts matrix

```
# ----- metadata  
metadata <- as.data.frame(colData(data)[, c("gender", "vital_status", "age_at_index", "days_to_l  
ast_follow_up", "paper_pathologic_stage", "paper_BRCA_Subtype_PAM50", "tissue_type")])  
  
colnames(metadata) <- c("Gender", "SurvivalStatus", "Age", "DaysToLastFollowUp", "PathologicStag  
e", "PAM50", "TissueType")  
metadata <- na.omit(metadata)  
metadata <- metadata %>% mutate(SurvivalStatus = ifelse(SurvivalStatus == "Alive", FALSE, Surviv  
alStatus)) %>%  
  mutate(SurvivalStatus = ifelse(SurvivalStatus == "Dead", TRUE, SurvivalStatus))
```

The metadata is biased to some Cancer sub types: Basal 176, Her2 68, LumA 529, LumB 186, Normal 35. For having a more reliable dataset we will choose an unbiased subset dataset. in total the research contains 313 samples: 35 Normal, 68 Her2 and 70 LumA, LumB and Basal.

```
normal_data <- metadata[metadata$PAM50 == "Normal", ]
her2_data <- metadata[metadata$PAM50 == "Her2", ]
basal_data <- metadata[metadata$PAM50 == "Basal", ]
luma_data <- metadata[metadata$PAM50 == "LumA", ]
lumb_data <- metadata[metadata$PAM50 == "LumB", ]

metadata <- rbind(normal_data, her2_data, basal_data[1:70, ], luma_data[1:70, ], lumb_data[1:70,
])
```

```
# ----- gene names
gene_ids_to_names <- rowData(data)[, c("gene_id", "gene_name")]

# ----- counts matrix
counts <- as.data.frame(assay(data, "unstranded"))
counts <- counts[, colnames(counts) %in% rownames(metadata)] #filter
counts <- counts[, rownames(metadata)] # reorder

# Make sure the counts matrix corresponds to the metadata
all(rownames(metadata) == colnames(counts))
```

```
## [1] TRUE
```

Prepare and Run DESeq2

- Constructs a DESeq2 object

```
# Set Factors for dds analyses
metadata$PathologicStage <- as.factor(metadata$PathologicStage)
metadata$SurvivalStatus <- as.factor(metadata$SurvivalStatus)

dds <- DESeqDataSetFromMatrix(countData=counts,
                              colData=metadata,
                              design= ~ PathologicStage + SurvivalStatus)

dds <- DESeq(dds)
```

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

```
## -- replacing outliers and refitting for 11573 genes
## -- DESeq argument 'minReplicatesForReplace' = 7
## -- original counts are preserved in counts(dds)
```

```
## estimating dispersions
```

```
## fitting model and testing
```


Normalize Results

```
res <- results(dds)
significant_genes <- length(which(res$padj < 0.05)) # number of significant genes
print(paste(significant_genes, " significant genes were found."))

## [1] "432  significant genes were found."

counts.vst <- vst(dds) # normal the data
counts.vst <- assay(counts.vst)

var_per_gene <- apply(counts.vst, 1, var) # Calculate the variance per gene
selectedGenes <- names(var_per_gene[order(var_per_gene, decreasing = T)][1:50]) # Take the first 50 genes

counts.vst.significant <- counts.vst[selectedGenes,] # Construct a new matrix only for the significant genes

# Make names informative names
genes_names <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% selectedGenes]
rownames(counts.vst.significant) <- genes_names
colnames(counts.vst.significant) <- paste0(metadata$Gender, "_", metadata$Age, "_", metadata$PathologicStage, "_seq", seq_along(metadata$Gender))
```

Heatmaps by Distance functions

In this section we will plot for each distance function the Heatmap graph. The data is the normalized counts results. For each plot we added annotation for the PAM50 BRCA subtypes.

```

# Heatmap Colors
col_fun = colorRamp2(c(0, 5, 10, 15, 20), brewer.pal(n = 5, name = "RdBu"))
col_PAM50 = c("Basal" = "#FB9A99", "Her2"="#B2DF8A", "LumA"="#A6CEE3", "LumB"="#1F78B4", "Normal"="#FDBF6F")

# ----- Heat Maps -----

# ----- Euclidean Distance

euc_dist <- dist(counts.vst.significant, method = 'euclidean') # Create euclidean distance matrix for heatmap
euc_row_dend = dendsort(hclust(dist(counts.vst.significant))) # for row ordering
euc_col_dend = dendsort(hclust(dist(t(counts.vst.significant)))) # for column ordering

# the heatmap
ht_euc = ComplexHeatmap::Heatmap(counts.vst.significant,
                                column_title = "Euclidean Distance",
                                name = " ",
                                show_row_names = TRUE, show_column_names = FALSE,
                                col = col_fun,
                                clustering_distance_rows = function(x) euc_dist, # calculates the heatmap
                                by_distance_matrix
                                clustering_distance_columns = function(x) euc_dist,
                                clustering_method_rows = "complete",
                                clustering_method_columns = "complete",
                                column_dend_height = unit(2, "cm"),
                                cluster_rows = euc_row_dend, cluster_columns = euc_col_dend,
                                row_names_rot = -45,
                                row_names_gp = grid::gpar(fontsize = 4),
                                top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM50), col=list(PAM50 = col_PAM50)),
                                heatmap_width = unit(12, "cm")
                                )

# ----- Pair wise Correlation

# A function that calculates pairwise distance for matrix
cor_distance <- function(mat) {
  as.dist(1 - cor(t(mat)))
}

cor_row_dend = dendsort(hclust(cor_distance(counts.vst.significant))) # for row ordering
cor_col_dend = dendsort(hclust(cor_distance(t(counts.vst.significant)))) # for column ordering

# the heatmap
ht_cor = ComplexHeatmap::Heatmap(counts.vst.significant,
                                column_title = "Correlation Pairwise Distance",
                                name = " ",
                                show_row_names = TRUE, show_column_names = FALSE,
                                col = col_fun,
                                clustering_distance_rows = function(x, y) 1 - cor(x, y), # calculates the
                                heatmap by function for each sample
                                clustering_distance_columns = function(x, y ) 1 - cor(x, y),
                                clustering_method_rows = "complete",
                                clustering_method_columns = "complete",
                                column_dend_height = unit(2, "cm"),
                                cluster_rows = cor_row_dend, cluster_columns = cor_col_dend,
                                row_names_rot = -45,
                                row_names_gp = grid::gpar(fontsize = 4),
                                top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM50), col=list(PAM50 = col_PAM50)),
                                heatmap_width = unit(12, "cm")
                                )

# ----- Cosine Similarity
library(lsa) # for cosine similarity function

```

Warning: package 'lsa' was built under R version 4.3.3

Loading required package: SnowballC

```

# Convert similarity to distance so it can be used to create the clustering
cos_dist <- as.dist(1 - cosine(counts.vst.significant))

# the heatmap
ht_cos = ComplexHeatmap::Heatmap(counts.vst.significant,
                                  column_title = "Cosine Simalarity",
                                  name = " ",
                                  show_row_names = TRUE, show_column_names = FALSE,
                                  col = col_fun,
                                  clustering_distance_rows = function(x,y) 1- cosine(x, y), # calculates the
heatmap by function for each sample
                                  clustering_distance_columns = function(x,y) 1- cosine(x, y),
                                  clustering_method_rows = "complete",
                                  clustering_method_columns = "complete",
                                  column_dend_height = unit(2, "cm"),
                                  row_names_rot = -45,
                                  row_names_gp = grid::gpar(fontsize = 4),
                                  top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM5
0), col=list(PAM50 = col_PAM50)),
                                  heatmap_width = unit(12, "cm")
                                  )

# ----- Kernel Similarity
# A function that calculates kernel similarity for matrix
ker_dis_matrix <- function(mat) {
  product_matrix <- t(mat) %*% mat
  result_matrix <- (product_matrix + 1)^2
  return(result_matrix)
}

ker_row_dend = dendsort(hclust(cos_dist)) # for row ordering
ker_col_dend = dendsort(hclust(t(cos_dist))) # for column ordering

# the heatmap
ht_ker = ComplexHeatmap::Heatmap(counts.vst.significant,
                                  column_title = "Kernel based distance (for nonlinear distance)",
                                  name = " ",
                                  show_row_names = TRUE, show_column_names = FALSE,
                                  col = col_fun,
                                  clustering_distance_rows = function(x, y) (t(x)%*%y + 1)^2, # calculates t
he heatmap by function for each sample
                                  clustering_distance_columns = function(x, y) (t(x)%*%y + 1)^2,
                                  clustering_method_rows = "complete",
                                  clustering_method_columns = "complete",
                                  column_dend_height = unit(2, "cm"),
                                  row_names_rot = -45,
                                  row_names_gp = grid::gpar(fontsize = 4),
                                  top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM5
0), col=list(PAM50 = col_PAM50)),
                                  heatmap_width = unit(12, "cm")
                                  )

# ----- Mahalanobis Distance

rescaled_data <- rescale(counts.vst.significant) # normalize for smaller values for the mahalano
bis distance function, so it wont exceed the max int value of the computer

# data for the mahalanobis_distance function
means <- colMeans(rescaled_data)
cov_matrix <- cov(rescaled_data) + diag(1e-4, ncol(rescaled_data)) # add a very small number so
the matrix will be diagonalizable
n_features <- ncol(rescaled_data)

# A function that calculates mahalanobis distance for samples
mahalanobis_distance <- function(x, y) {

```

```

x <- matrix(x, nrow = 1, ncol = n_features)
y <- matrix(y, nrow = 1, ncol = n_features)

# Check if x or y need padding
if (length(x) < n_features) {
  x <- x + rep(0, n_features - length(x))
}
if (length(y) < n_features) {
  y <- y + rep(0, n_features - length(y))
}
# Calculate distance
return(as.numeric(sqrt((x - y)%*% solve(cov_matrix) %*% t(x - y))))
}

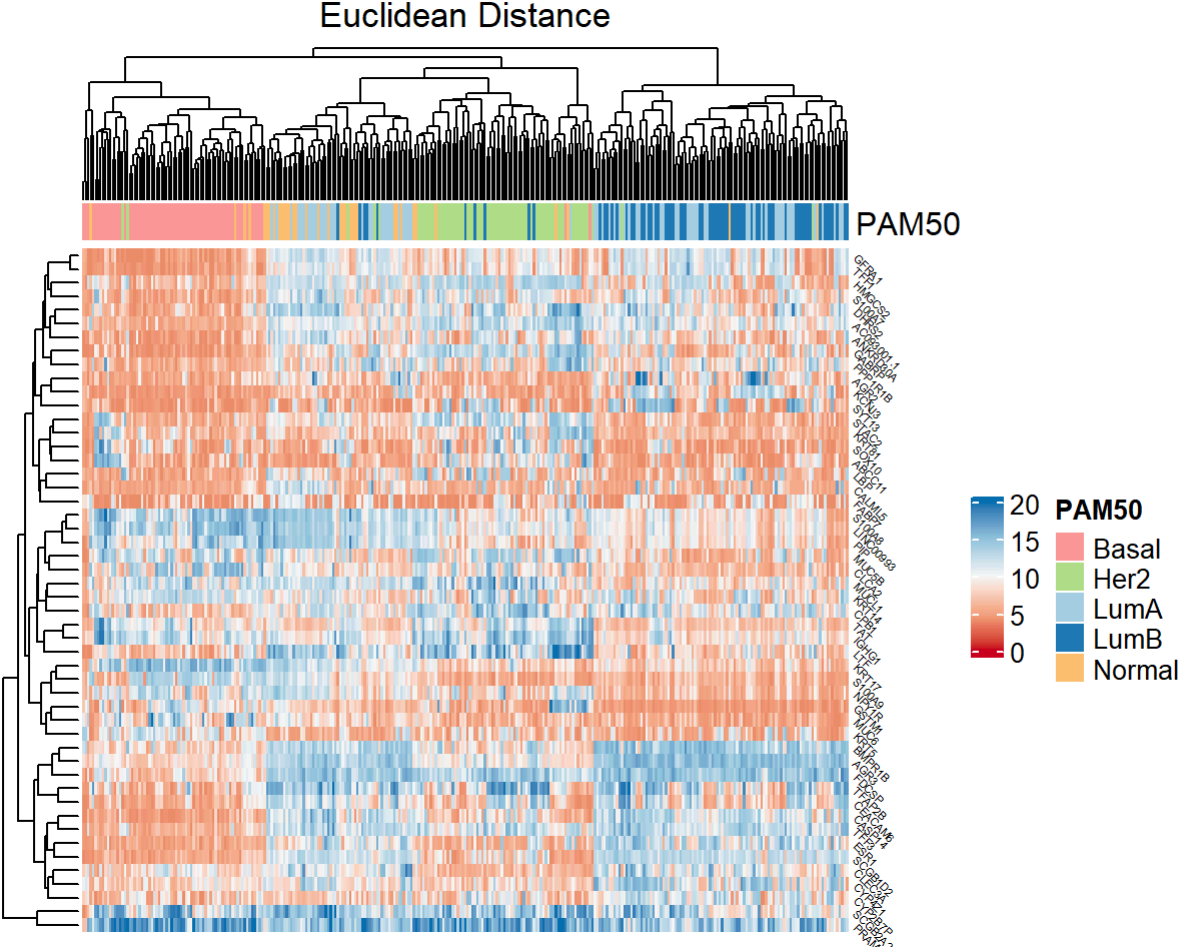
# the heatmap
ht_mah = ComplexHeatmap::Heatmap(counts.vst.significant,
  column_title = "Mahalanobis distance",
  name = " ",
  show_row_names = TRUE, show_column_names = FALSE,
  col = col_fun,
  clustering_distance_rows = function(x,y) mahalanobis_distance(x,y), # calculates the heatmap by function for each sample
  clustering_distance_columns = function(x,y) mahalanobis_distance(x,y),
  clustering_method_rows = "complete",
  clustering_method_columns = "complete",
  column_dend_height = unit(2, "cm"),
  row_names_rot = -45,
  row_names_gp = grid::gpar(fontsize = 4),
  top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM50), col=list(PAM50 = col_PAM50)),
  heatmap_width = unit(12, "cm")
)

# ----- PCA

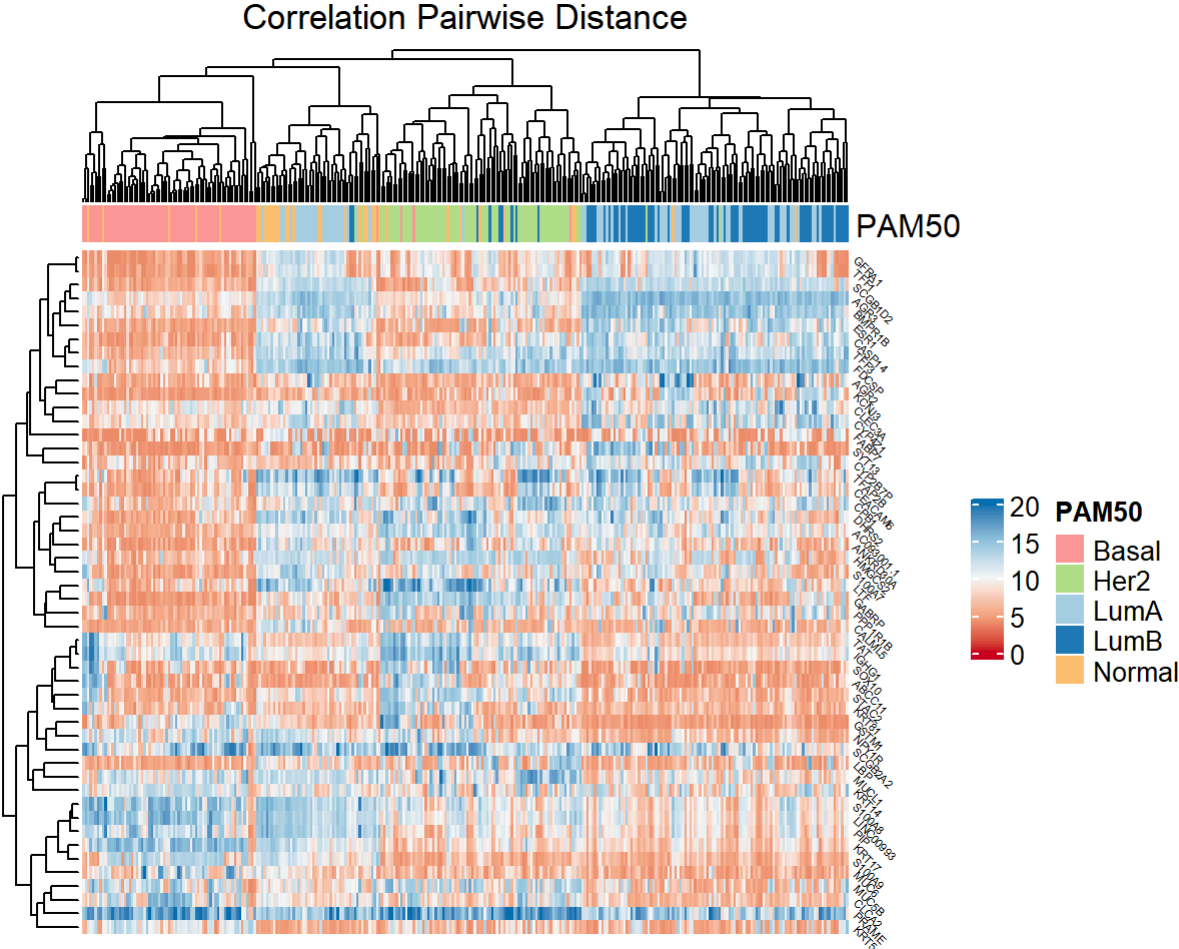
pcaResults = prcomp(t(counts.vst.significant))
pca_hc = ComplexHeatmap::Heatmap(counts.vst.significant,
  column_title = "distance by PCA",
  name = " ",
  show_row_names = TRUE, show_column_names = FALSE,
  col = col_fun,
  clustering_distance_rows = function(x, y) sqrt(sum((pcaResults$x[x, ] - pcaResults$x[y, ])^2)), # calculates the heatmap by function for each sample
  clustering_distance_columns = function(x, y) sqrt(sum((pcaResults$x[x, ] - pcaResults$x[y, ])^2)),
  clustering_method_rows = "complete",
  clustering_method_columns = "complete",
  column_dend_height = unit(2, "cm"),
  row_names_rot = -45,
  row_names_gp = grid::gpar(fontsize = 4),
  top_annotation = ComplexHeatmap::HeatmapAnnotation(PAM50 = c(metadata$PAM50), col=list(PAM50 = col_PAM50)),
  heatmap_width = unit(12, "cm")
)

# prints the plots
ht_euc

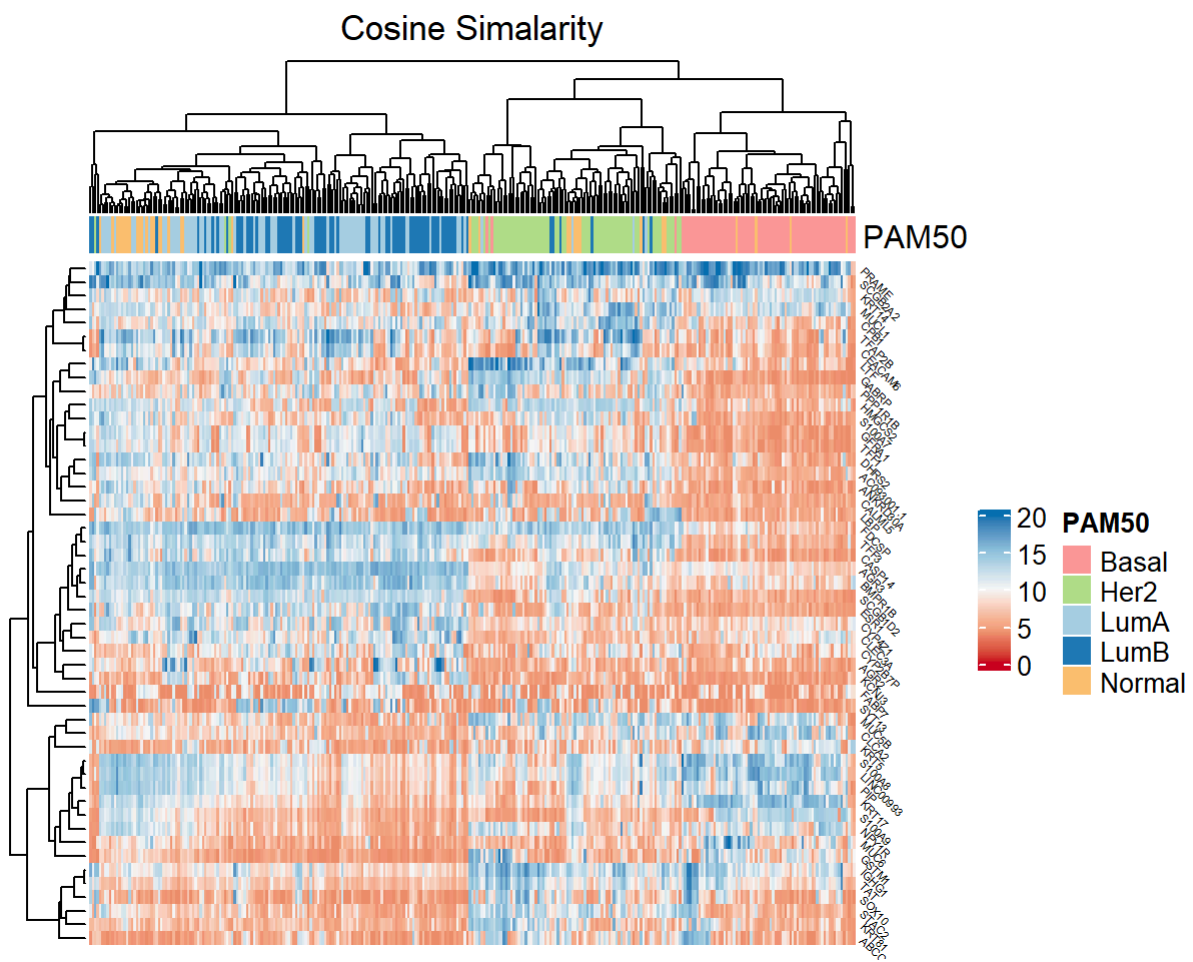
```



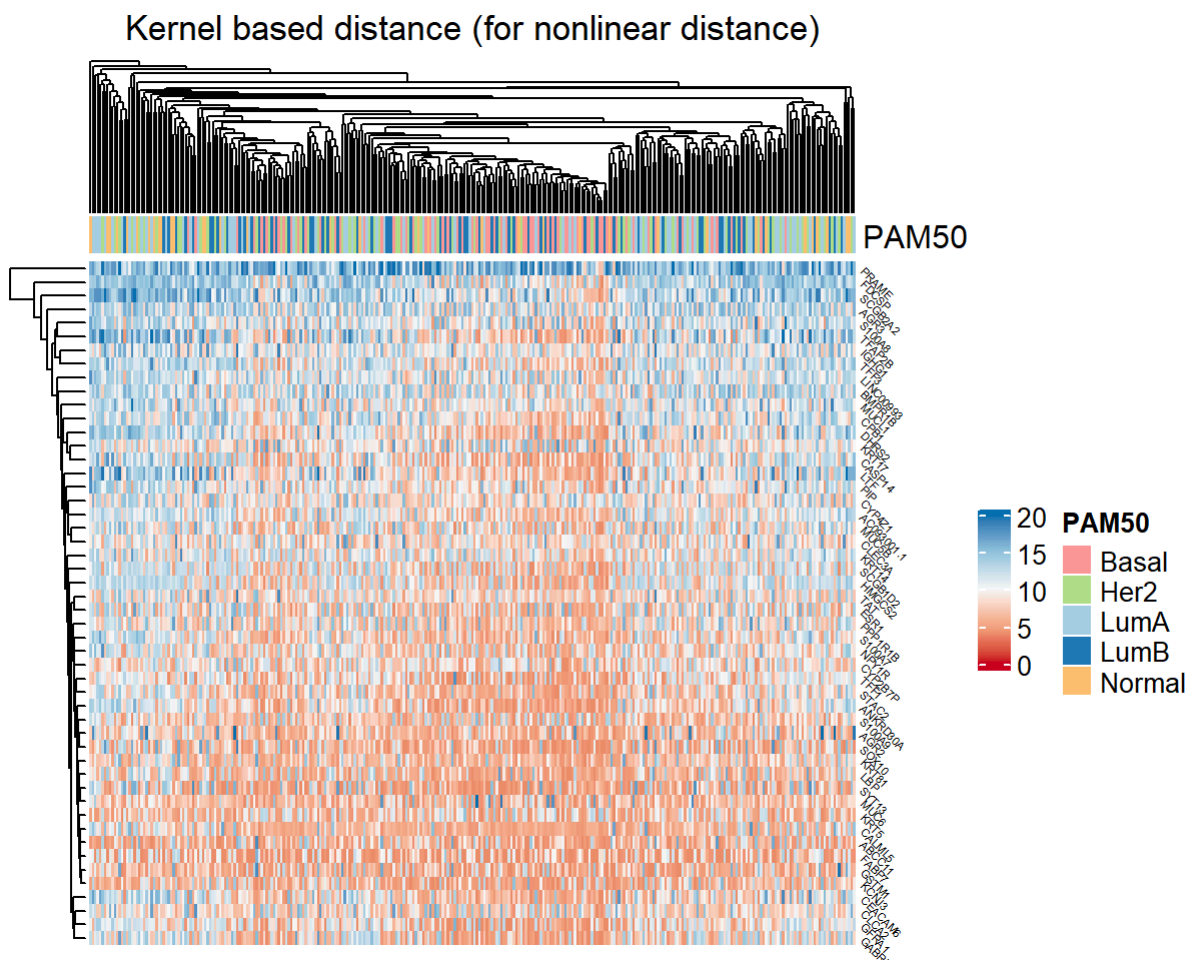
ht_cor



ht_cos



ht_ker



pca_hc


```
k <- 5
```

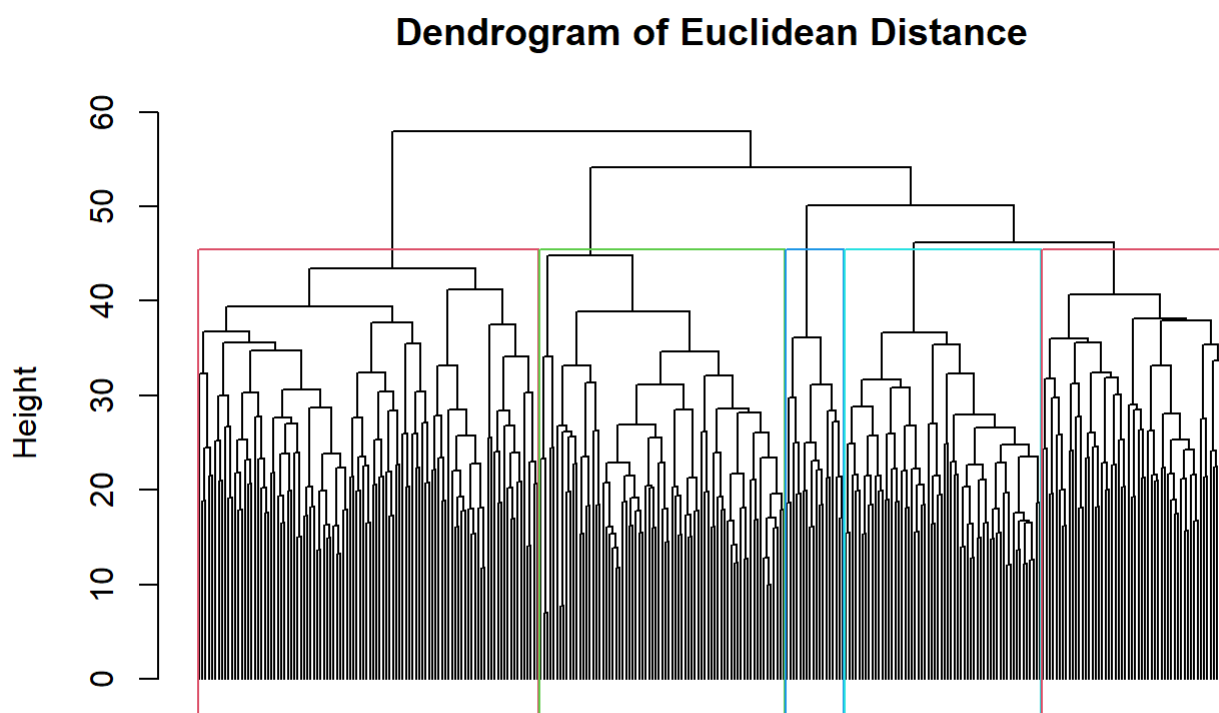
```
# Clustering By distance functions and add to metadata
```

```
# ----- Euclidean Distance
```

```
hclust_res_euc <- hclust(dist(t(counts.vst.significant))) # clustering
```

```
plot(hclust_res_euc, cex = 0.6, hang = -1, main = "Dendrogram of Euclidean Distance", labels=FALSE, xlab="", sub="") # plot
```

```
rect.hclust(hclust_res_euc, k = k, border = 2:5)
```



```
clusters_euc <- cutree(hclust_res_euc, k =k) # add to metadata
```

```
metadata$ClusterEuc <- clusters_euc
```

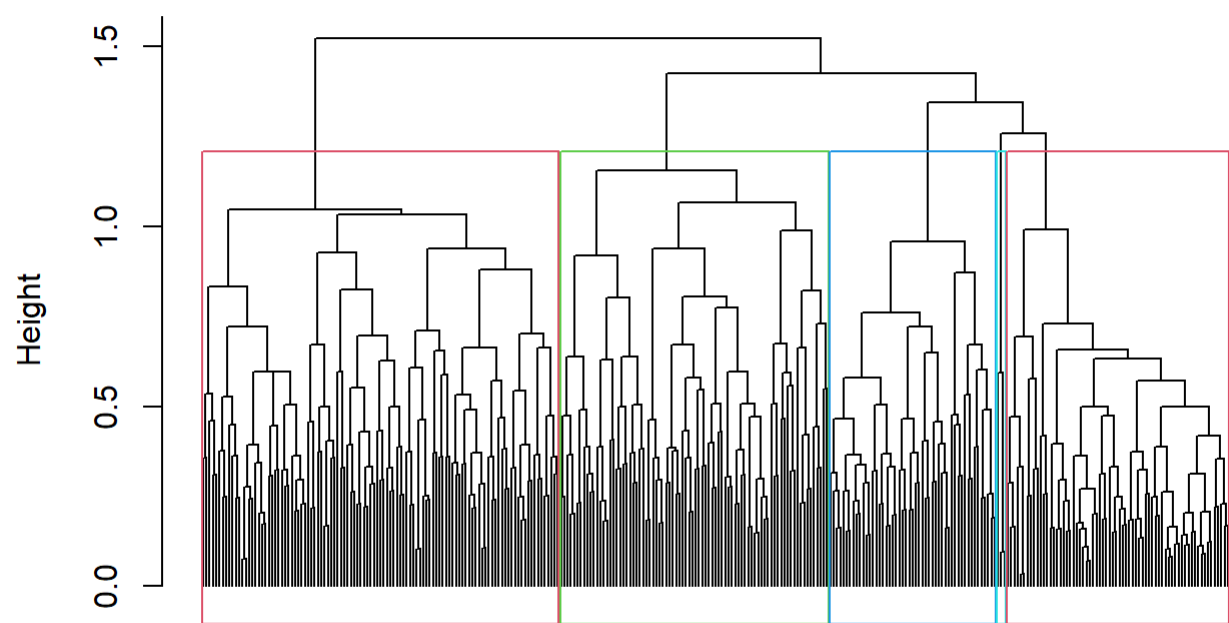
```
# ----- Pair wise Correlation
```

```
hclust_res_cor <- hclust(cor_distance(t(counts.vst.significant))) # clustering
```

```
plot(hclust_res_cor, cex = 0.6, hang = -1, main = "Dendrogram of Pair wise Correlation", labels=FALSE, xlab="", sub="") # plot
```

```
rect.hclust(hclust_res_cor, k = k, border = 2:5)
```

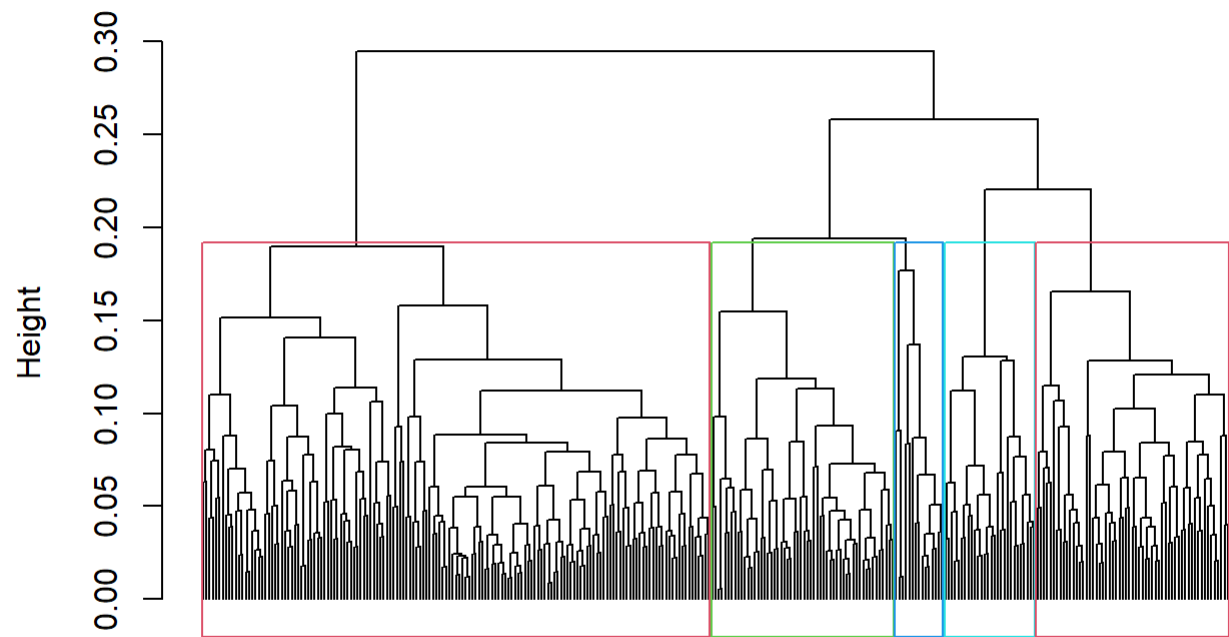
Dendrogram of Pair wise Correlation



```
clusters_cor <- cutree(hclust_res_cor, k =k) # add to metadata
metadata$ClusterCor <- clusters_cor

# ----- Cosine Similarity
hclust_res_cos <- hclust(t(cos_dist)) # clustering
plot(hclust_res_cos, cex = 0.6, hang = -1, main = "Dendrogram of Cosine Similarity", labels=FALS
E, xlab="", sub="") # plot
rect.hclust(hclust_res_cos, k = k, border = 2:5)
```

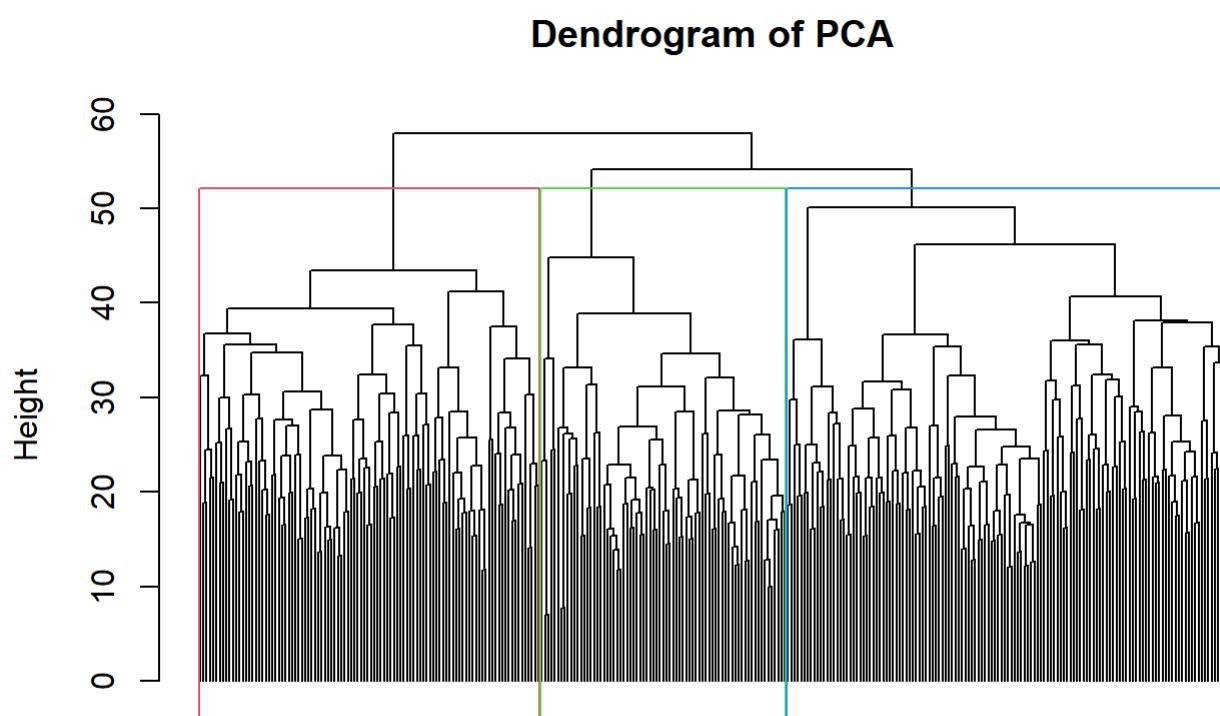
Dendrogram of Cosine Similarity



```
clusters_cos <- cutree(hclust_res_cos, k =k) # add to metadata
metadata$ClusterCos <- clusters_cos

# ----- Kernel Similarity
#not clustering where created.

# ----- PCA
pcaResults <- prcomp(t(counts.vst.significant))
hclust_res_pca <- hclust(dist(pcaResults$x)) # clustering
plot(hclust_res_pca, cex = 0.6, hang = -1, main = "Dendrogram of PCA", labels=FALSE, xlab="", sub="") # plot
rect.hclust(hclust_res_pca, k = 3, border = 2:5)
```



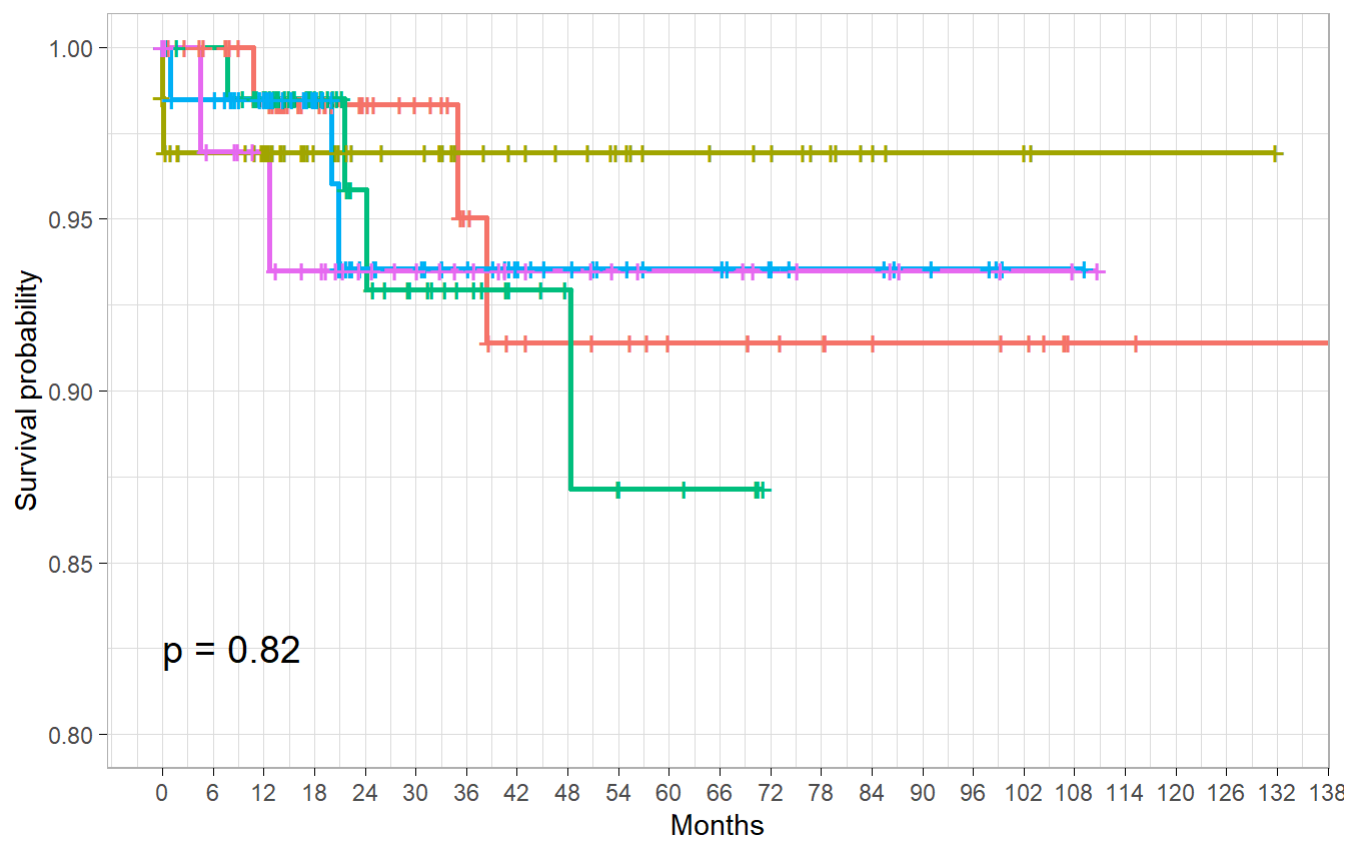
```
clusters_pca <- cutree(hclust_res_pca, k =3) # add to metadata
metadata$ClusterPCA <- clusters_pca
```

Survival Analyses

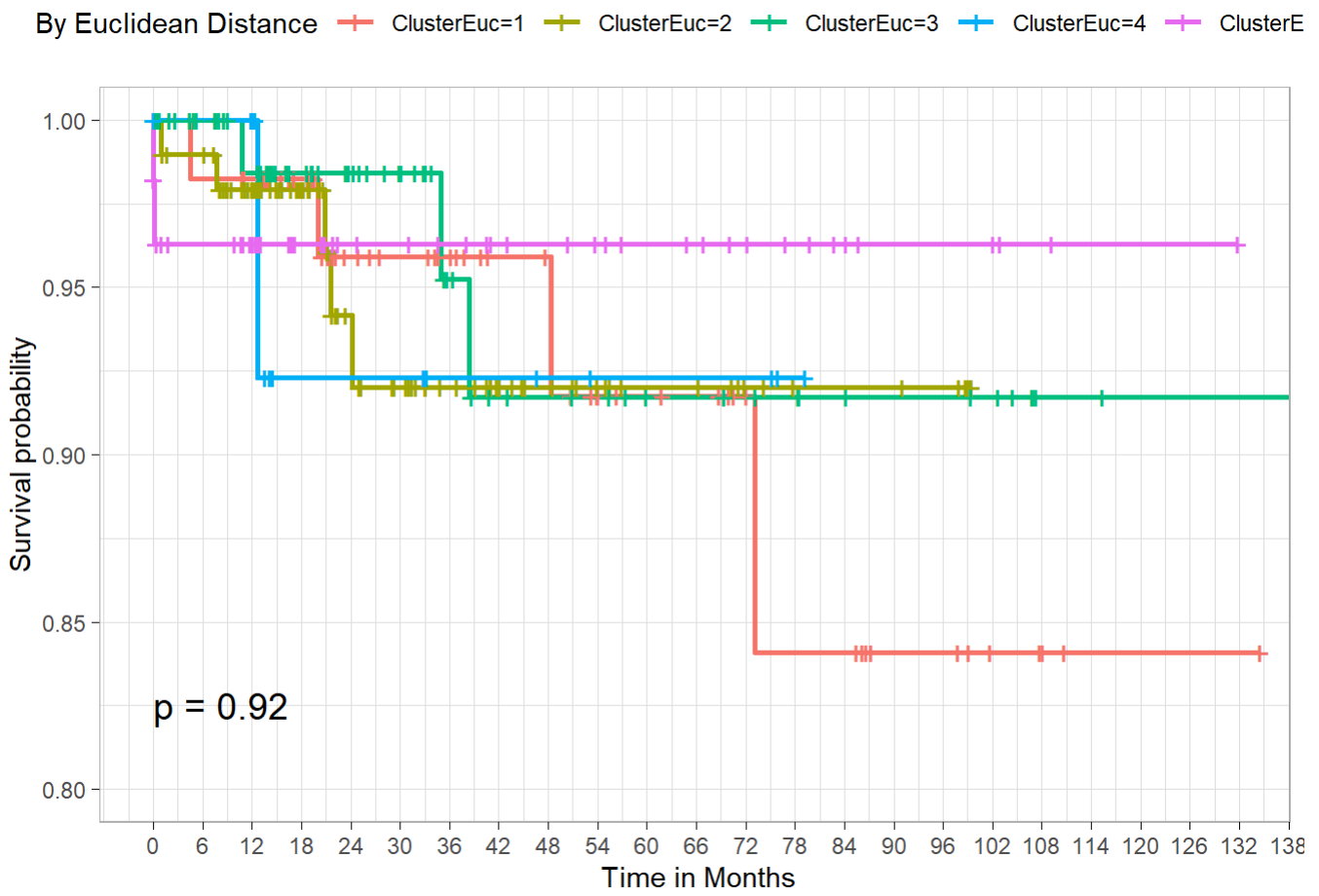
In this section we will do survival analyses with “Kaplan-Meier” curve.

```
metadata$SurvivalStatus <- as.logical(metadata$SurvivalStatus)
# ----- PAM50
ggsurvplot(
  survfit(Surv(metadata$DaysToLastFollowUp, metadata$SurvivalStatus) ~ PAM50, data = metadata),
  data = metadata,
  xlab = "Months",
  xscale = 30.4,
  break.x.by = 182.4,
  surv.median.line = "hv",
  legend.title = "By PAM50",
  pval = TRUE,
  pval.coord = c(0.1, 0.825),
  xlim = c(0,4000),
  ylim=c(0.8, 1),
  ggtheme = theme_light()
)
```

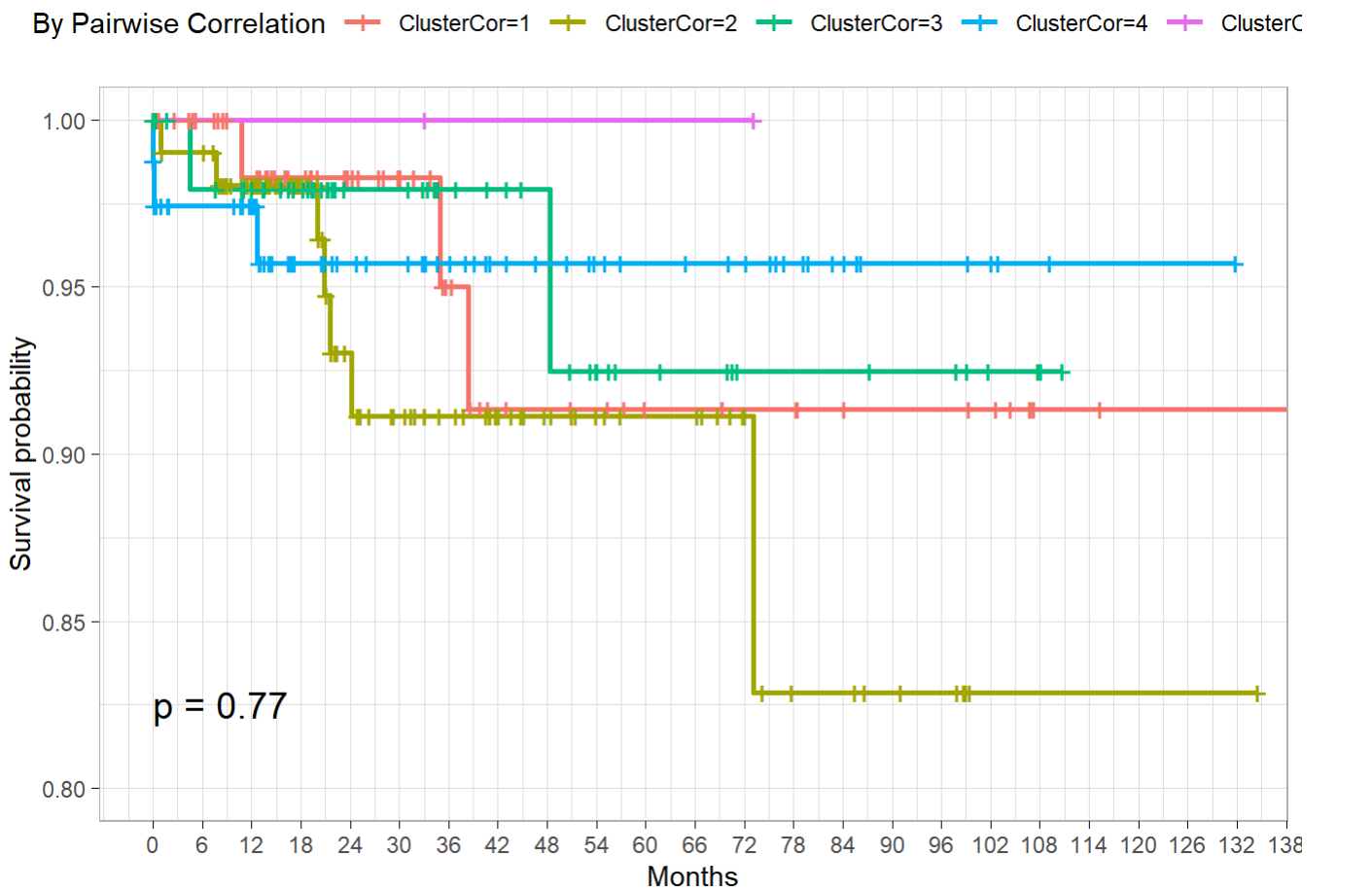
By PAM50 PAM50=Basal PAM50=Her2 PAM50=LumA PAM50=LumB PAM50=Normal



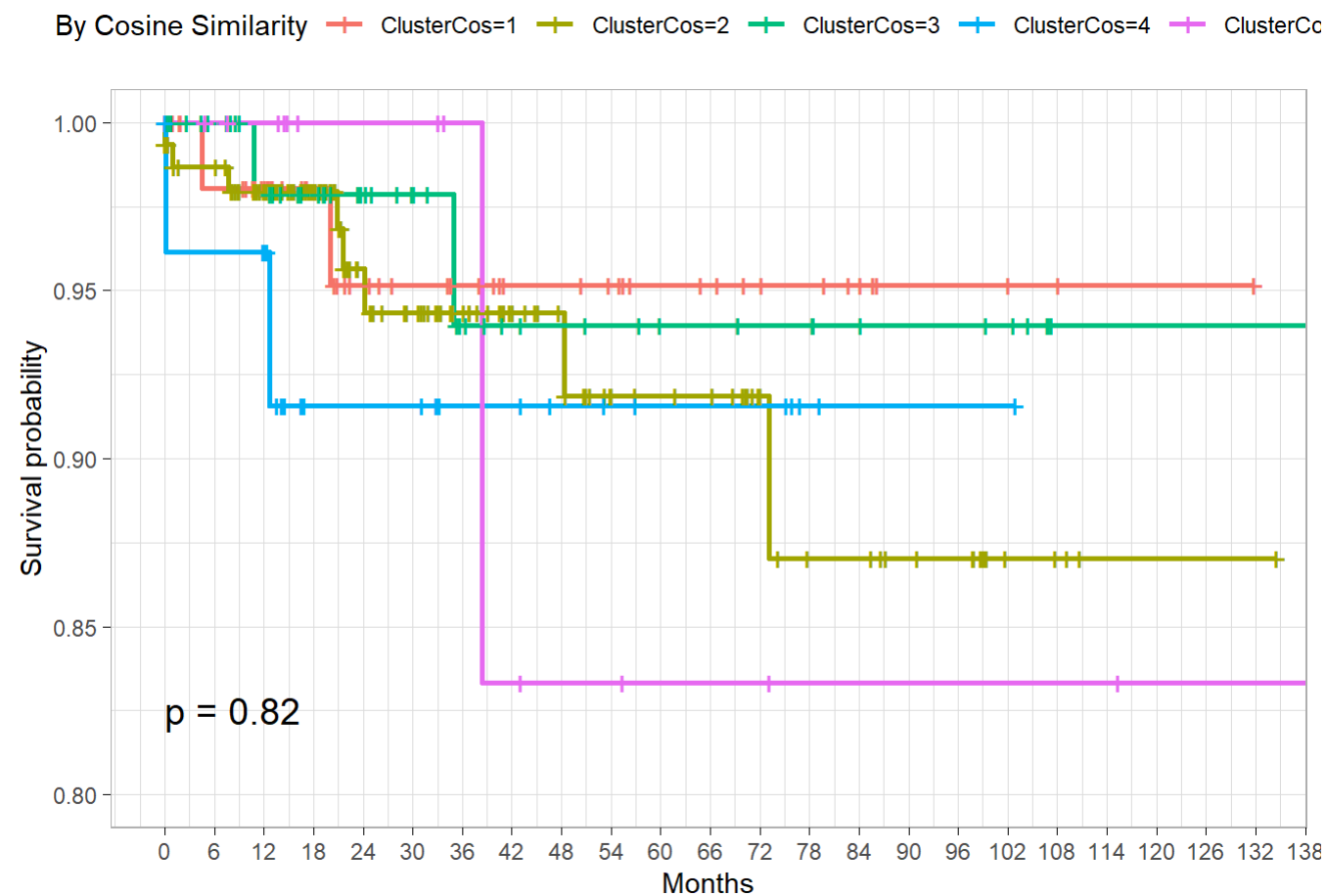
```
# ----- Euclidean Distance
ggsurvplot(
  survfit(Surv(metadata$DaysToLastFollowUp, metadata$SurvivalStatus) ~ ClusterEuc, data = meta
data),
  data = metadata,
  xlab = " Time in Months",
  xscale = 30.4,
  xlim = c(0,4000),
  ylim=c(0.8, 1),
  break.x.by = 182.4,
  pval = TRUE,
  pval.coord = c(0.1, 0.825),
  surv.median.line = "hv",
  legend.title = "By Euclidean Distance",
  ggtheme = theme_light()
)
```



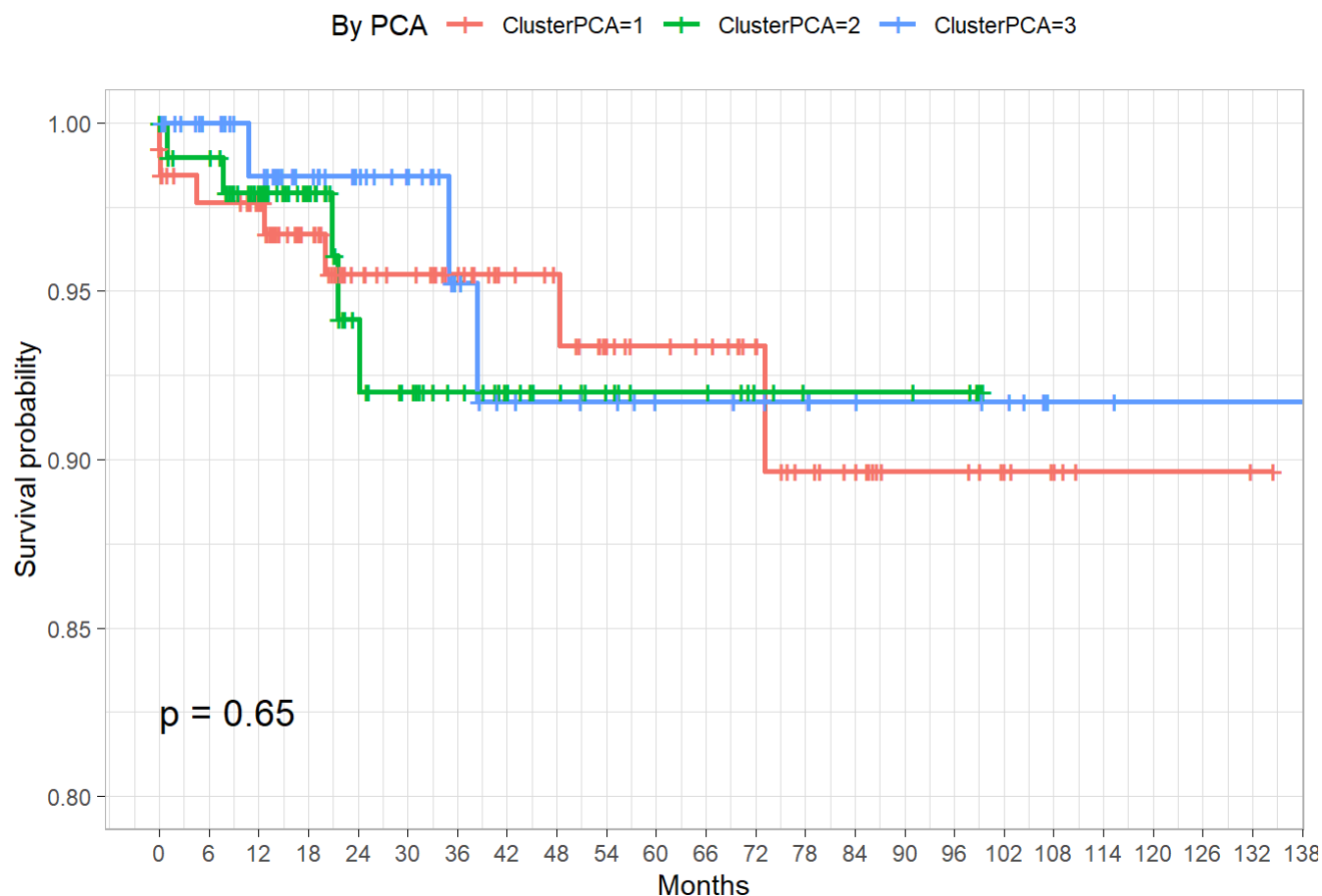
```
# ----- Pairwise Correlation
ggsurvplot(
  survfit(Surv(metadata$DaysToLastFollowUp, metadata$SurvivalStatus) ~ ClusterCor, data = meta
data),
  data = metadata,
  xlab = "Months",
  xscale = 30.4,
  break.x.by = 182.4,
  surv.median.line = "hv",
  legend.title = "By Pairwise Correlation",
  xlim = c(0,4000),
  ylim=c(0.8, 1),
  pval = TRUE,
  pval.coord = c(0.1, 0.825),
  ggtheme = theme_light()
)
```



```
# ----- Cosine Similarity
ggsurvplot(
  survfit(Surv(metadata$DaysToLastFollowUp, metadata$SurvivalStatus) ~ ClusterCos, data = meta
data),
  data = metadata,
  xlab = "Months",
  xscale = 30.4,
  break.x.by = 182.4,
  surv.median.line = "hv",
  legend.title = "By Cosine Similarity",
  xlim = c(0,4000),
  ylim=c(0.8, 1),
  pval = TRUE,
  pval.coord = c(0.1, 0.825),
  ggtheme = theme_light()
)
```



```
# ----- PCA
ggsurvplot(
  survfit(Surv(metadata$DaysToLastFollowUp, metadata$SurvivalStatus) ~ ClusterPCA, data = meta
data),
  data = metadata,
  xlab = "Months",
  xscale = 30.4,
  break.x.by = 182.4,
  surv.median.line = "hv",
  legend.title = "By PCA",
  xlim = c(0,4000),
  ylim=c(0.8, 1),
  pval = TRUE,
  pval.coord = c(0.1, 0.825),
  ggtheme = theme_light()
)
```



Enrichment Analysis

This section is divided to 2 parts. The first part is GSEA on the counts.vst results, where the results are analyses by cluster. In this section no significant pathways for any cluster of the distance functions. Because of the similarity between the results, we only presents the euclidean distance and the PCA.

```
# prepare hallmarks
hallmarks <- msigdbr(species = "Homo sapiens", category = "H") %>%
  dplyr::select(gs_name, gene_symbol)
hallmarks <- msigdbr(species = "Homo sapiens", category = "H")
hallmarks <- hallmarks[,c('gs_name', 'gene_symbol')]

# ----- Euclidean Distance

results_list <- list()
for (i in unique(metadata$ClusterEuc)) {
  cluster_samples <- which(metadata$ClusterEuc == i)
  gene_list <- apply(counts.vst.significant[, cluster_samples], 1, var)
  names(gene_list) <- rownames(counts.vst.significant)
  gene_list <- sort(gene_list, decreasing = TRUE)

  gsea_results <- GSEA(gene_list, TERM2GENE = hallmarks, verbose = TRUE, scoreType = "pos")
  results_list[[paste("Cluster", i)]] <- gsea_results
}
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```



```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
# Assuming the GSEA results are stored in 'results_list'
for (i in names(results_list)) {
  gsea_results <- results_list[[i]]

  # Check if the results are empty
  if (length(gsea_results) == 0 || nrow(gsea_results@result) == 0) {
    print(paste("No significant pathways for Euclidean", i))
  } else {
    print(paste("Results for", i))
    dotplot(gsea_results)
  }
}
```

```
## [1] "No significant pathways for Euclidean Cluster 1"
## [1] "No significant pathways for Euclidean Cluster 2"
## [1] "No significant pathways for Euclidean Cluster 3"
## [1] "No significant pathways for Euclidean Cluster 4"
## [1] "No significant pathways for Euclidean Cluster 5"
```

```
# ----- PCA

results_list <- list()
for (i in unique(metadata$ClusterPCA)) {
  cluster_samples <- which(metadata$ClusterPCA == i)
  gene_list <- apply(counts.vst.significant[, cluster_samples], 1, var)
  names(gene_list) <- rownames(counts.vst.significant)
  gene_list <- sort(gene_list, decreasing = TRUE)

  gsea_results <- GSEA(gene_list, TERM2GENE = hallmarks, verbose = TRUE, scoreType = "pos")
  results_list[[paste("Cluster", i)]] <- gsea_results
}
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## no term enriched under specific pvalueCutoff...
```

```
# Assuming the GSEA results are stored in 'results_list'
for (i in names(results_list)) {
  gsea_results <- results_list[[i]]

  # Check if the results are empty
  if (length(gsea_results) == 0 || nrow(gsea_results@result) == 0) {
    print(paste("No significant pathways for PCA Cluster", i))
  } else {
    print(paste("Results for", i))
    dotplot(gsea_results)
  }
}
```

```
## [1] "No significant pathways for PCA Cluster Cluster 1"
## [1] "No significant pathways for PCA Cluster Cluster 2"
## [1] "No significant pathways for PCA Cluster Cluster 3"
```

The second part we repeated the DESeq analyses design by each distance function Clustering and did the GSEA by the new dds results.

```
# ----- Euclidean Distance

metadata$ClusterEuc <- as.factor(metadata$ClusterEuc)
euclidean_dds <- DESeqDataSetFromMatrix(countData=counts,
                                       colData=metadata,
                                       design= ~ ClusterEuc)

euclidean_dds <- DESeq(euclidean_dds)
```

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

```
## -- replacing outliers and refitting for 9868 genes
## -- DESeq argument 'minReplicatesForReplace' = 7
## -- original counts are preserved in counts(dds)
```

```
## estimating dispersions
```

```
## fitting model and testing
```

```
euc.counts.vst <- vst(euclidean_dds)
euc.counts.vst <- assay(euc.counts.vst)

rownames(euc.counts.vst) <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% rownames
(euc.counts.vst)]
var_per_gene <- apply(euc.counts.vst, 1, var) # Calculate the variance per gene
var_per_gene <- var_per_gene[order(var_per_gene, decreasing = T)]
selectedGenes <- names(var_per_gene[order(var_per_gene, decreasing = T)])
selectedGenes <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% selectedGenes]

hm <- GSEA(var_per_gene, TERM2GENE = hallmarks, pvalueCutoff = 0.05, eps = 0)
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

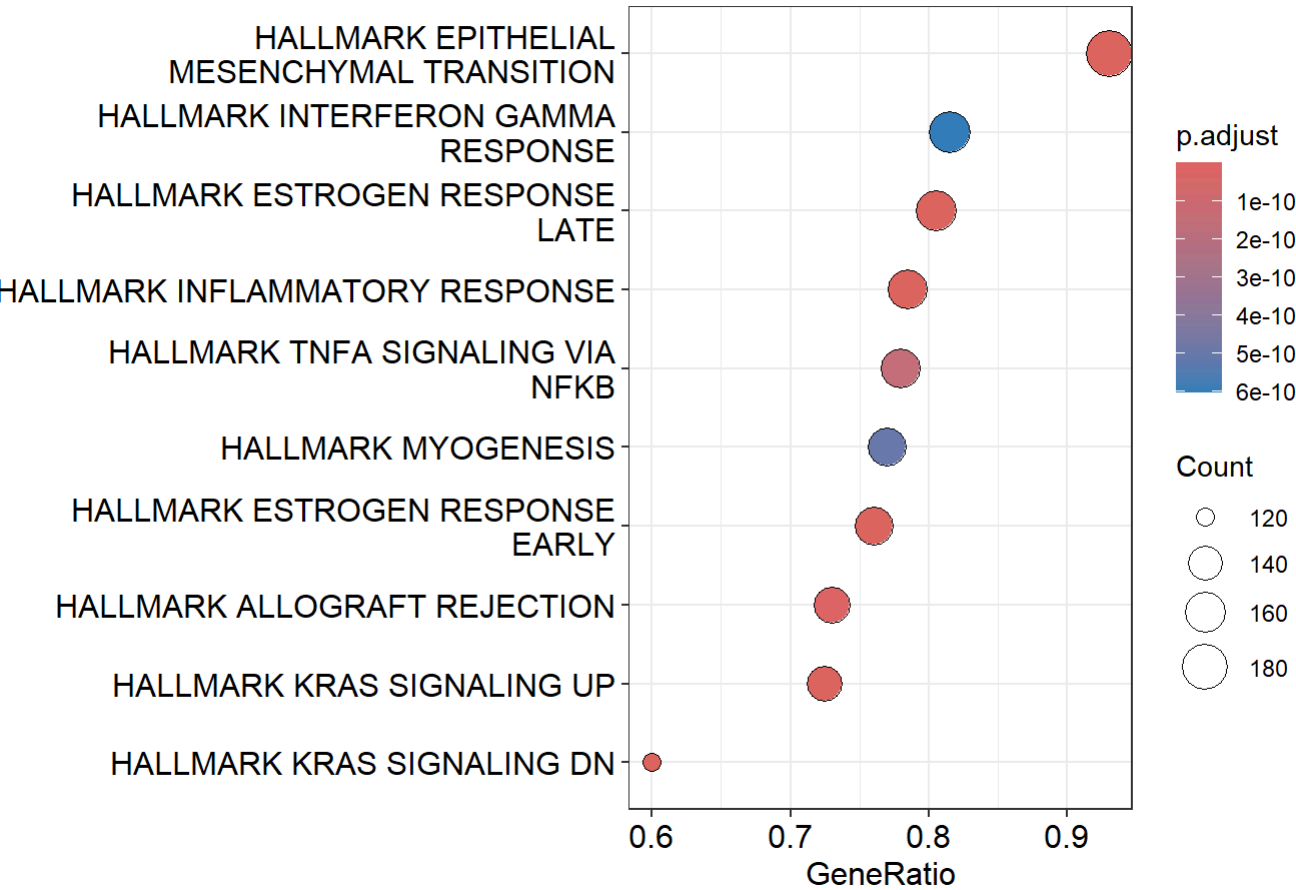
```
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are
ties in the preranked stats (1.02% of the list).
## The order of those tied genes will be arbitrary, which may produce unexpected results.
```

```
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize,
## gseaParam, : There are duplicate gene names, fgsea may produce unexpected
## results.
```

```
## leading edge analysis...
```

```
## done...
```

```
dotplot(hm)
```



```
# ----- Pairwise Correlation
metadata$ClusterCor <- as.factor(metadata$ClusterCor)
correlation_dds <- DESeqDataSetFromMatrix(countData=counts,
                                         colData=metadata,
                                         design= ~ ClusterCor)
correlation_dds <- DESeq(correlation_dds)
```

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

```
## -- replacing outliers and refitting for 4864 genes
## -- DESeq argument 'minReplicatesForReplace' = 7
## -- original counts are preserved in counts(dds)
```

```
## estimating dispersions
```

```
## fitting model and testing
```

```
cor.counts.vst <- vst(correlation_dds)
cor.counts.vst <- assay(cor.counts.vst)

rownames(cor.counts.vst) <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% rownames
(cor.counts.vst)]
var_per_gene <- apply(cor.counts.vst, 1, var) # Calculate the variance per gene
var_per_gene <- var_per_gene[order(var_per_gene, decreasing = T)]
selectedGenes <- names(var_per_gene[order(var_per_gene, decreasing = T)])
selectedGenes <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% selectedGenes]

hallmarks <- msigdb(species = "Homo sapiens", category = "H")
hallmarks <- hallmarks[,c('gs_name', 'gene_symbol')]
hm <- GSEA(var_per_gene, TERM2GENE = hallmarks, pvalueCutoff = 0.05, eps = 0)
```

```
## preparing geneSet collections...
```

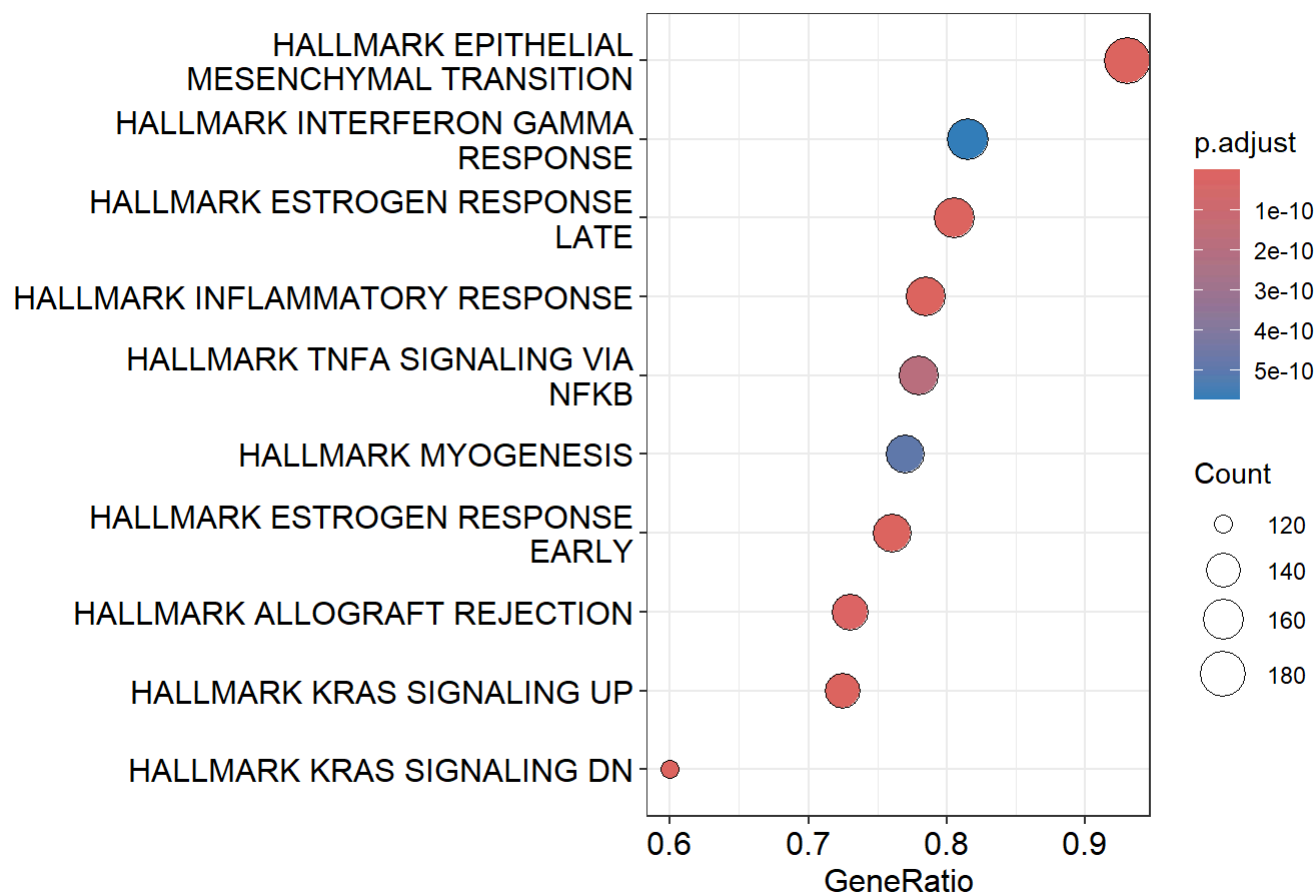
```
## GSEA analysis...
```

```
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are
ties in the preranked stats (1.02% of the list).
## The order of those tied genes will be arbitrary, which may produce unexpected results.
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are
duplicate gene names, fgsea may produce unexpected results.
```

```
## leading edge analysis...
```

```
## done...
```

```
dotplot(hm)
```



```
# ----- Cosine Similarity
metadata$ClusterCos <- as.factor(metadata$ClusterCos)
cosine_dds <- DESeqDataSetFromMatrix(countData=counts,
                                     colData=metadata,
                                     design= ~ ClusterCos)
cosine_dds <- DESeq(cosine_dds)
```

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

```
## -- replacing outliers and refitting for 9778 genes
## -- DESeq argument 'minReplicatesForReplace' = 7
## -- original counts are preserved in counts(dds)
```

```
## estimating dispersions
```

```
## fitting model and testing
```

```
cos.counts.vst <- vst(cosine_dds)
cos.counts.vst <- assay(cos.counts.vst)

rownames(cos.counts.vst) <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% rownames
(cos.counts.vst)]
var_per_gene <- apply(cos.counts.vst, 1, var) # Calculate the variance per gene
var_per_gene <- var_per_gene[order(var_per_gene, decreasing = T)]
selectedGenes <- names(var_per_gene[order(var_per_gene, decreasing = T)])
selectedGenes <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% selectedGenes]

hallmarks <- msigdb(species = "Homo sapiens", category = "H")
hallmarks <- hallmarks[,c('gs_name', 'gene_symbol')]
hm <- GSEA(var_per_gene, TERM2GENE = hallmarks, pvalueCutoff = 0.05, eps = 0)
```

```
## preparing geneSet collections...
```

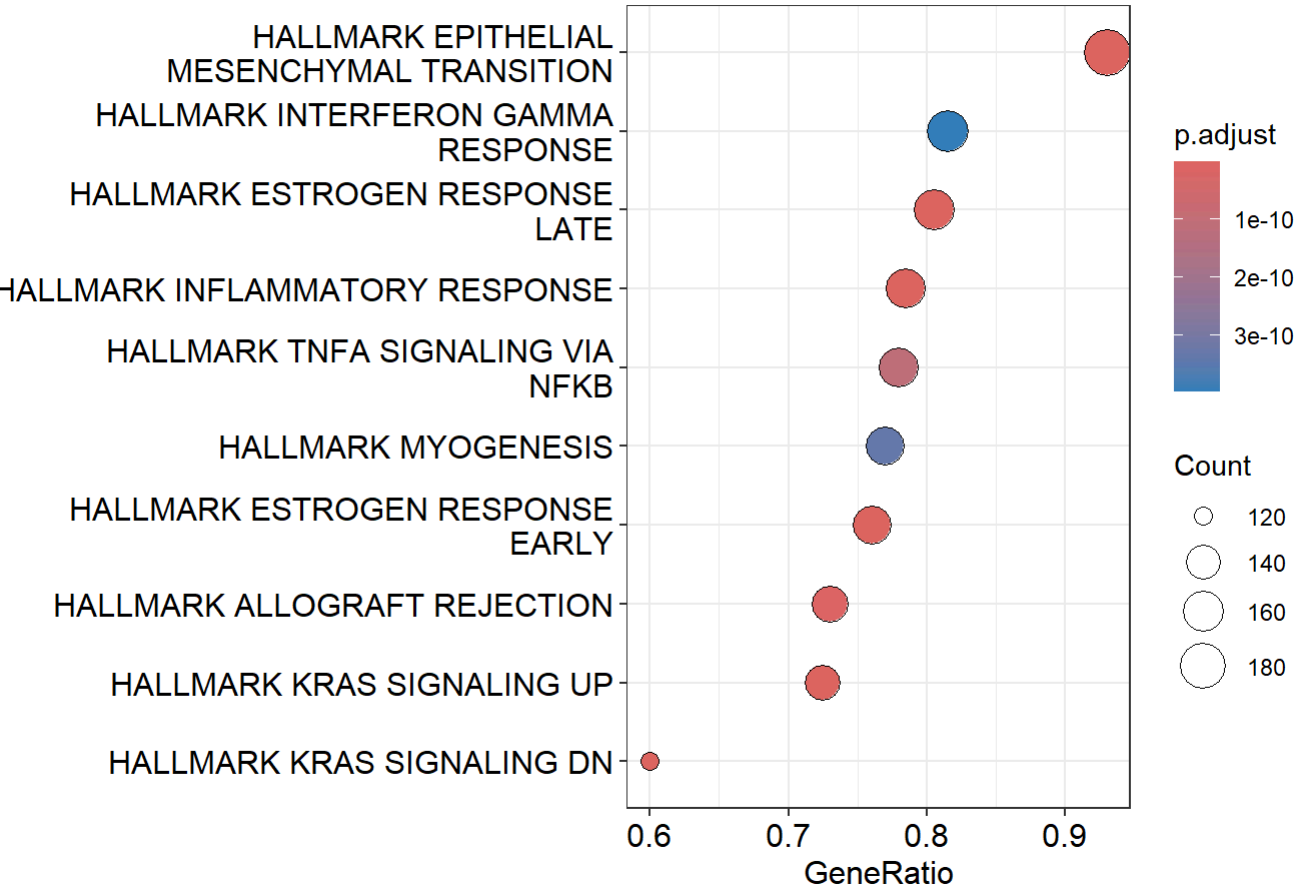
```
## GSEA analysis...
```

```
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are
ties in the preranked stats (1.02% of the list).
## The order of those tied genes will be arbitrary, which may produce unexpected results.
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are
duplicate gene names, fgsea may produce unexpected results.
```

```
## leading edge analysis...
```

```
## done...
```

```
dotplot(hm)
```



```
# ----- PCA
metadata$ClusterPCA <- as.factor(metadata$ClusterPCA)
pca_dds <- DESeqDataSetFromMatrix(countData=counts,
                                  colData=metadata,
                                  design= ~ ClusterPCA)
pca_dds <- DESeq(pca_dds)
```

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

```
## -- replacing outliers and refitting for 12224 genes  
## -- DESeq argument 'minReplicatesForReplace' = 7  
## -- original counts are preserved in counts(dds)
```

```
## estimating dispersions
```

```
## fitting model and testing
```

```
pca.counts.vst <- vst(pca_dds)  
pca.counts.vst <- assay(pca.counts.vst)  
  
rownames(pca.counts.vst) <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% rownames  
(pca.counts.vst)]  
var_per_gene <- apply(pca.counts.vst, 1, var) # Calculate the variance per gene  
var_per_gene <- var_per_gene[order(var_per_gene, decreasing = T)]  
selectedGenes <- names(var_per_gene[order(var_per_gene, decreasing = T)])  
selectedGenes <- gene_ids_to_names$gene_name[gene_ids_to_names$gene_id %in% selectedGenes]  
  
hallmarks <- msigdb(species = "Homo sapiens", category = "H")  
hallmarks <- hallmarks[,c('gs_name', 'gene_symbol')]  
hm <- GSEA(var_per_gene, TERM2GENE = hallmarks, pvalueCutoff = 0.05, eps = 0)
```

```
## preparing geneSet collections...
```

```
## GSEA analysis...
```

```
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are  
ties in the preranked stats (1.02% of the list).  
## The order of those tied genes will be arbitrary, which may produce unexpected results.  
## Warning in preparePathwaysAndStats(pathways, stats, minSize, maxSize, gseaParam, : There are  
duplicate gene names, fgsea may produce unexpected results.
```

```
## leading edge analysis...
```

```
## done...
```

```
dotplot(hm)
```

