

# EW Sarcoma

2025-12-02

## R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
rm(list = ls())
```

```
# Load libraries
library(pheatmap)
library(recount)
library(DESeq2)
library(tidyverse)
library(ggpubr)
library(EnhancedVolcano)
library(EnsDb.Hsapiens.v86)
library(msigdbr)
library(clusterProfiler)
library(enrichplot)
library(cowplot)
library(org.Hs.eg.db)
library(gridExtra)
library(ggridges)
```

```
#The same as going the recount2 and download them manually
project_info <- abstract_search(query = "Ewing sarcoma")
```

```
# Download the study data
download_study("SRP015989")
```

```
## 2026-01-19 11:56:30.893008 downloading file rse_gene.Rdata to SRP015989
```

```
# Load the data
load("SRP015989/rse_gene.Rdata")
```

```
# Fix the colData to give a column with the appropriate groups
rse_gene$condition <- c(rep("shCTR", 3), rep("shEF1", 4))
```

```
rownames(rse_gene) <- gsub("\\..*$", "", rownames(rse_gene))
head(rownames(rse_gene))
```

```

## [1] "ENSG00000000003" "ENSG00000000005" "ENSG000000000419" "ENSG000000000457"
## [5] "ENSG000000000460" "ENSG000000000938"

# Create DESeq2 data
dds <- DESeqDataSet(rse_gene, design = ~condition)

## converting counts to integer mode

## Warning in DESeqDataSet(rse_gene, design = ~condition): 45 duplicate rownames
## were renamed by adding numbers

## Warning in DESeqDataSet(rse_gene, design = ~condition): some variables in
## design formula are characters, converting to factors

#filter out low count genes across all samples
dds <- dds[rowSums(counts(dds)) > 10, ]

# relevel so that the control experiment is the base condition
dds$condition <- relevel(dds$condition, ref = "shCTR")

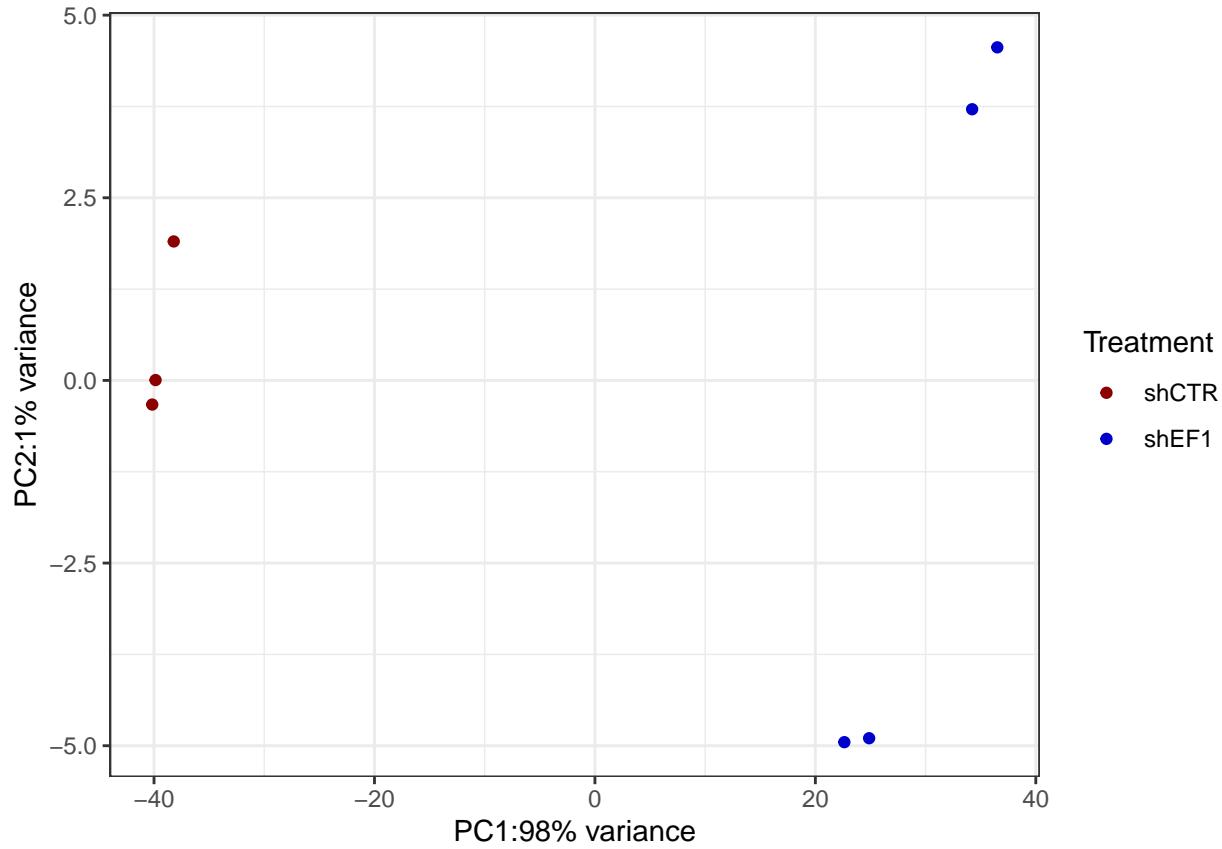
# transform data into log scale to make the variance more uniform
logdds <- rlog(dds)
logdata <- as.data.frame(assay(logdds))

#PCA is performed to confirm replicates cluster correctly, detect batch effects or outliers, and ensure
PCA_data <- plotPCA(logdds, intgroup = "condition", returnData = TRUE)

## using ntop=500 top features by variance

percentVar <- round(100 * attr(PCA_data, "percentVar"))
ggplot(PCA_data, aes(x = PC1, y = PC2, color = condition)) + geom_point() + theme_bw() +
  xlab(paste0("PC1:", percentVar[1], "% variance")) +
  ylab(paste0("PC2:", percentVar[2], "% variance")) +
  scale_color_manual(name = "Treatment", values = c("red4", "blue3"))

```



This PCA plot shows a very strong separation between shCTR and shEF1 along PC1, indicating that EF1 knockdown causes a dominant transcriptional change. Replicates cluster tightly within each group, suggesting good sample quality and no major batch effects.

```
# Perform DESeq2 analysis
dds <- DESeq(dds)

## estimating size factors

## estimating dispersions

## gene-wise dispersion estimates

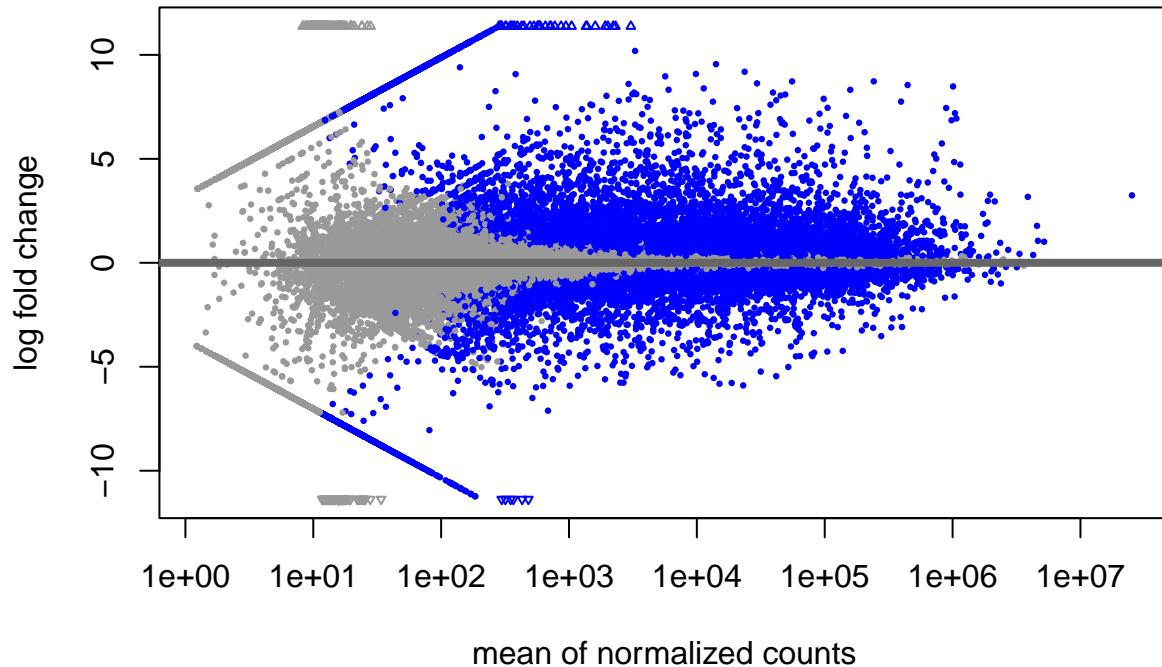
## mean-dispersion relationship

## final dispersion estimates

## fitting model and testing

res <- results(dds)

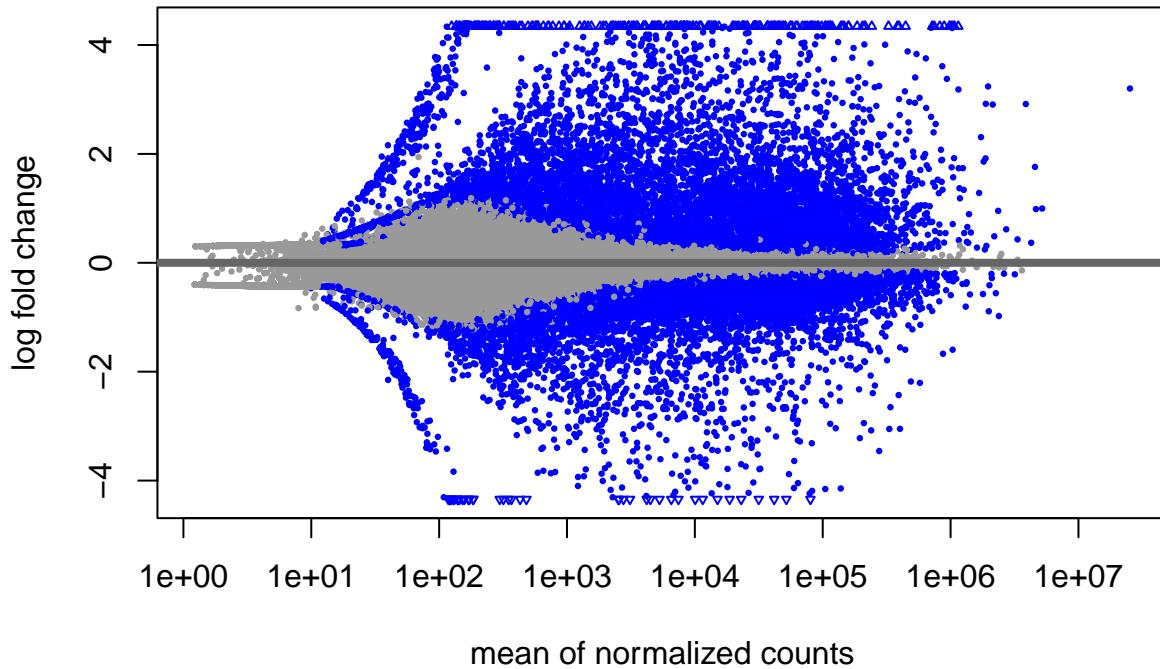
# -- plotMA
plotMA(res)
```



```
# LFC shrink
#find the fold change is exaggerated and shrink into right
resNorm <- lfcShrink(dds = dds, res = res, type = "normal", coef = 2)
```

```
## using 'normal' for LFC shrinkage, the Normal prior from Love et al (2014).
##
## Note that type='apeglm' and type='ashr' have shown to have less bias than type='normal'.
## See ?lfcShrink for more details on shrinkage type, and the DESeq2 vignette.
## Reference: https://doi.org/10.1093/bioinformatics/bty895
```

```
#coef = 2 means change the second results
# -- plotMA with resNorm
plotMA(resNorm)
```



```

#Create a result dataframe called resdf
resdf <- as.data.frame(resNorm) %>%
  rownames_to_column() %>%
  rename(ENSEMBL = rowname)

#extract gene lists based on ENSEMBL
genelist <- AnnotationDbi::select(org.Hs.eg.db, keys = rownames(resNorm), keytype = "ENSEMBL",
                                    columns = c("ENTREZID", "SYMBOL", "GENETYPE", "GENENAME"))

## 'select()' returned 1:many mapping between keys and columns

#filter genes with no padj, no symbol, or duplicated symbol
resdf <- resdf %>%
  left_join(genelist, by = "ENSEMBL") %>%
  filter(!is.na(padj)) %>%
  filter(!is.na(SYMBOL)) %>%
  filter(!duplicated(SYMBOL))

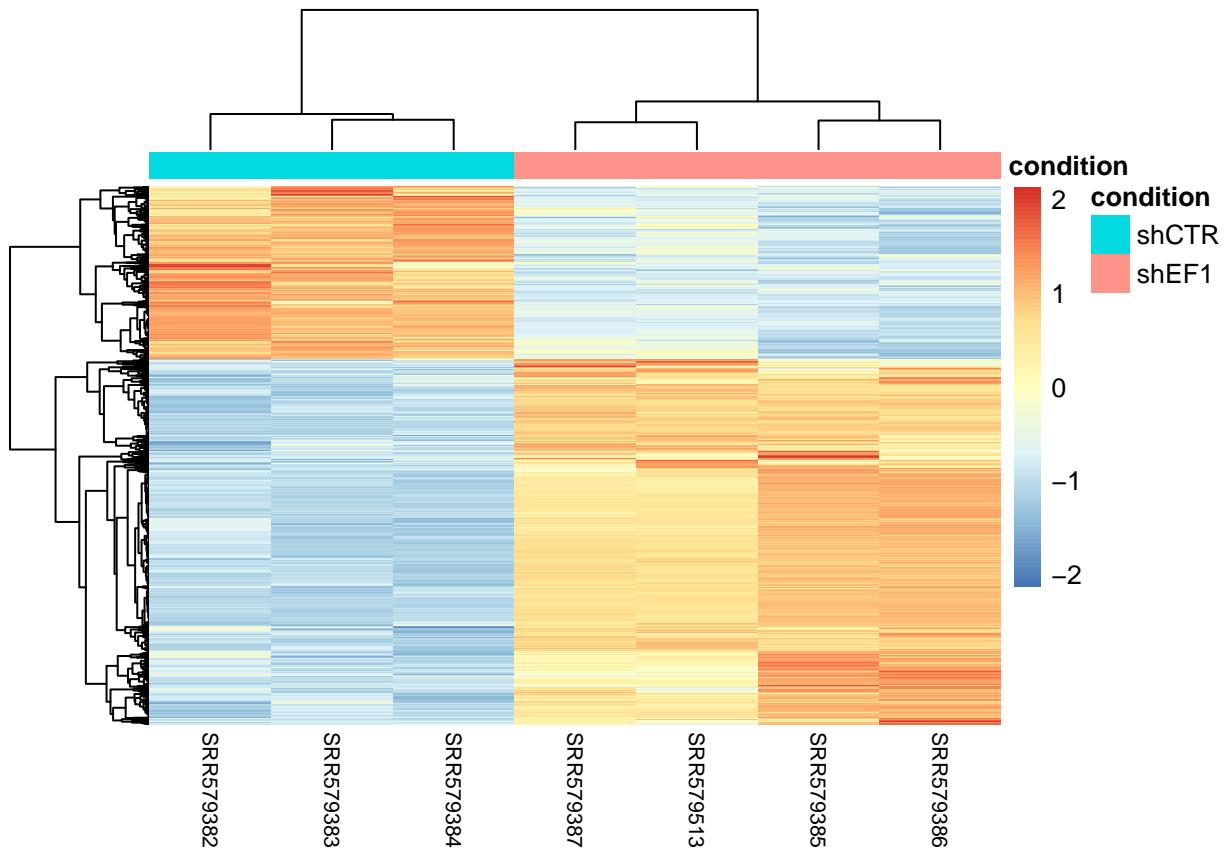
#filter out significant genes
sig_gene_resdata <- resdf[resdf$baseMean >= 20 & abs(resdf$log2FoldChange) >= 1 &
                           resdf$padj <= 0.05, ]
sig_gene <- sig_gene_resdata$ENSEMBL

```

```
#heatmap for significant genes
sig_gene_logdata <- logdata[sig_gene,]

heatmap_anno <- as.data.frame(colData(dds)) %>%
  dplyr::select(condition)

pheatmap(sig_gene_logdata, scale = "row", clustering_distance_rows = "correlation",
         annotation_col = heatmap_anno, title = "Differently expressed gene",
         show_rownames = FALSE, fontsize_col = 8)
```



```
#count the number of up-regulated genes
```

```
nrow(resdf[resdf$baseMean >= 20 & resdf$log2FoldChange >= 1 & resdf$padj <= 0.05, ])
```

```
## [1] 3086
```

```
#count the number of down-regulated genes
```

```
nrow(resdf[resdf$baseMean >= 20 & resdf$log2FoldChange <=-1 & resdf$padj <= 0.05, ])
```

```
## [1] 1455
```

```
# extract data for the top 10-upregulated genes
```

```
top10_ensembl <- sig_gene_resdata %>%
  arrange(desc(log2FoldChange)) %>%
  head(10) %>%
```

```

pull(ENSEMBL)

top10_geneid <- sig_gene_resdata %>%
  arrange(desc(log2FoldChange)) %>%
  head(10) %>%
  pull(SYMBOL)

upregulated_heatdata <- sig_gene_logdata[top10_ensembl,]
rownames(upregulated_heatdata) <- top10_geneid

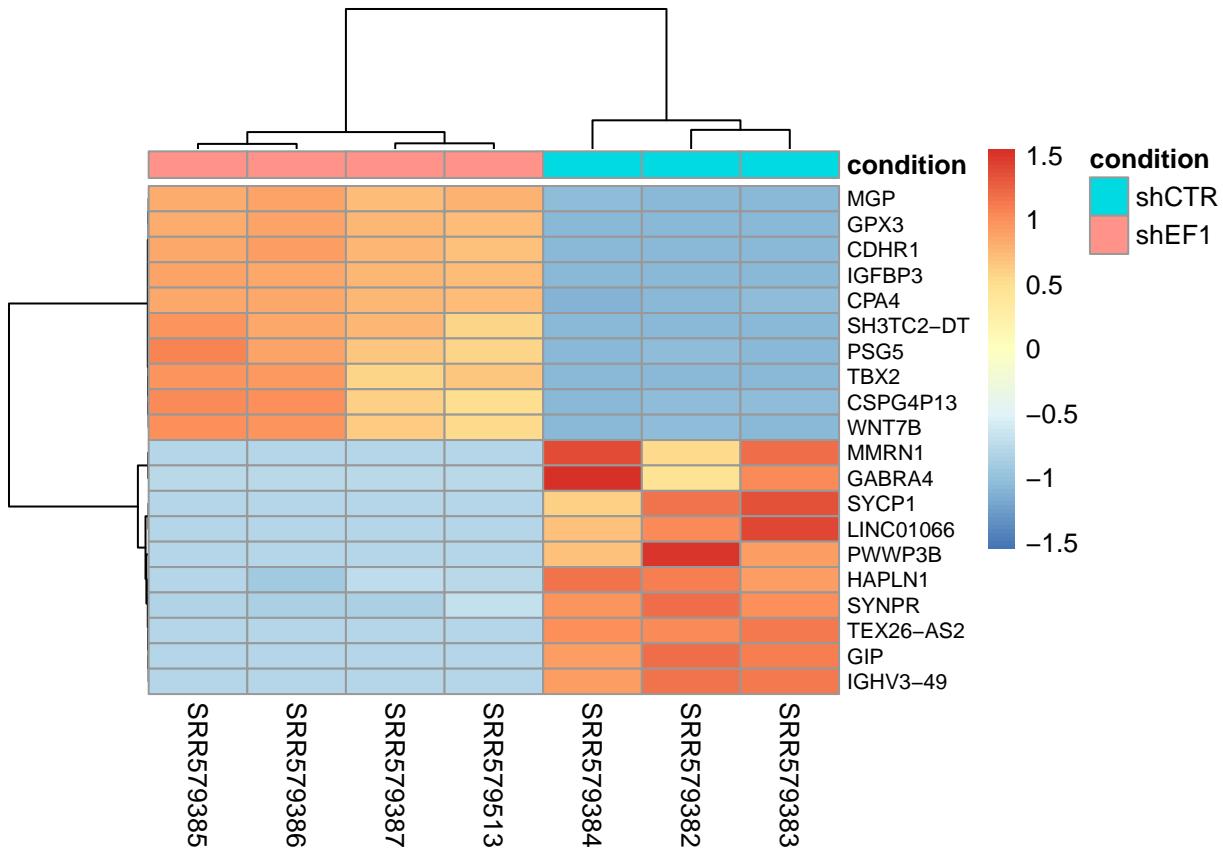
# extract data for the top 10-downregulated genes
down10_ensembl <- sig_gene_resdata %>%
  arrange(log2FoldChange) %>%
  head(10) %>%
  pull(ENSEMBL)

down10_geneid <- sig_gene_resdata %>%
  arrange(log2FoldChange) %>%
  head(10) %>%
  pull(SYMBOL)

downregulated_heatdata <- sig_gene_logdata[down10_ensembl,]
rownames(downregulated_heatdata) <- down10_geneid

#Heatmap of the top and down regulated genes
mergeheatdata <- rbind(downregulated_heatdata, upregulated_heatdata)
pheatmap(mergeheatdata, scale = "row", clustering_distance_rows = "correlation",
         annotation_col = heatmap_anno, fontsize_row = 8)

```



```

# Volcano plot
EnhancedVolcano(resdf, x = "log2FoldChange", y = "padj", lab = resdf$SYMBOL,
                 FCcutoff = 1, pCutoff = 0.05, pointSize = 1, labSize = 3,
                 title = "Volcano Plot", titleLabSize = 20, legendPosition = "right",
                 legendIconSize = 8, legendLabSize = 8, gridlines.major = FALSE,
                 gridlines.minor = FALSE, border = "full")

## Warning: One or more p-values is 0. Converting to 10^-1 * current lowest
## non-zero p-value...

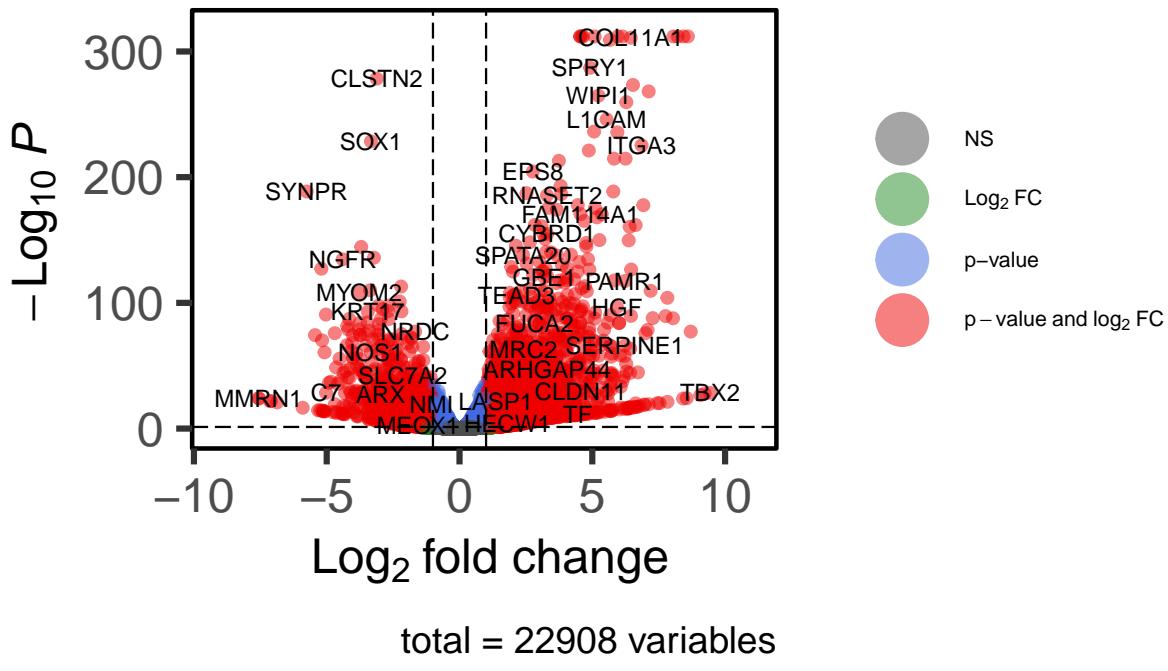
## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## i The deprecated feature was likely used in the EnhancedVolcano package.
##   Please report the issue to the authors.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

## Warning: The 'size' argument of 'element_rect()' is deprecated as of ggplot2 3.4.0.
## i Please use the 'linewidth' argument instead.
## i The deprecated feature was likely used in the EnhancedVolcano package.
##   Please report the issue to the authors.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

```

# Volcano Plot

## Enhanced Volcano



```
#Prepare gene list with decreasing log2FC for GSEA analysis (both GO and KEGG)
GSEAlist <- resdf$log2FoldChange
names(GSEAlist) <- resdf$ENTREZID

#remove NAs & remove duplicated IDs
GSEAlist <- GSEAlist[!is.na(GSEAlist)]
GSEAlist <- GSEAlist[!duplicated(names(GSEAlist))]

# sort for GSEA
GSEAlist <- sort(GSEAlist, decreasing = TRUE)

#Gene Set Enrichment Analysis for GO/KEGG
gseaGO <- gseGO(geneList = GSEAlist, ont = "ALL", minGSSize = 30, keyType = "ENTREZID",
                  pvalueCutoff = 0.05, pAdjustMethod = "BH", verbose = TRUE,
                  OrgDb = org.Hs.eg.db)

## using 'fgsea' for GSEA analysis, please cite Korotkevich et al (2019).

## preparing geneSet collections...

## GSEA analysis...

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

```

## Warning in fgseaMultilevel(pathways = pathways, stats = stats, minSize =
## minSize, : There were 8 pathways for which P-values were not calculated
## properly due to unbalanced (positive and negative) gene-level statistic values.
## For such pathways pval, padj, NES, log2err are set to NA. You can try to
## increase the value of the argument nPermSimple (for example set it nPermSimple
## = 10000)

## Warning in fgseaMultilevel(pathways = pathways, stats = stats, minSize =
## minSize, : For some of the pathways the P-values were likely overestimated. For
## such pathways log2err is set to NA.

## Warning in fgseaMultilevel(pathways = pathways, stats = stats, minSize =
## minSize, : For some pathways, in reality P-values are less than 1e-10. You can
## set the 'eps' argument to zero for better estimation.

## leading edge analysis...

## done...

gseaKEGG <- gseKEGG(geneList = GSEAlist, organism = "hsa", minGSSize = 30,
                      pvalueCutoff = 0.05, pAdjustMethod = "BH", verbose = TRUE)

## Reading KEGG annotation online: "https://rest.kegg.jp/link/hsa/pathway"...

## Reading KEGG annotation online: "https://rest.kegg.jp/list/pathway/hsa"...

## using 'fgsea' for GSEA analysis, please cite Korotkevich et al (2019).

## preparing geneSet collections...

## GSEA analysis...

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

## Warning in fgseaMultilevel(pathways = pathways, stats = stats, minSize =
## minSize, : For some of the pathways the P-values were likely overestimated. For
## such pathways log2err is set to NA.

## Warning in fgseaMultilevel(pathways = pathways, stats = stats, minSize =
## minSize, : For some pathways, in reality P-values are less than 1e-10. You can
## set the 'eps' argument to zero for better estimation.

## leading edge analysis...

## done...

```

```

#Make gene IDs readable (converts Entrez IDs to gene symbols)
gseaGO <- setReadable(gseaGO, OrgDb = org.Hs.eg.db, keyType = 'ENTREZID')
gseaKEGG <- setReadable(gseaKEGG, OrgDb = org.Hs.eg.db, keyType = 'ENTREZID')

#function to find the associated pathways with gene
gene_pathway_gsea <- function(gene, pathwaylist) {
  down_list <- list()
  for (i in seq(1, nrow(pathwaylist), by = 1)) {
    seq = pathwaylist$core_enrichment[i]
    genes <- strsplit(seq, "/")[[1]] # split string into vector
    down_list[[i]] <- genes
  }

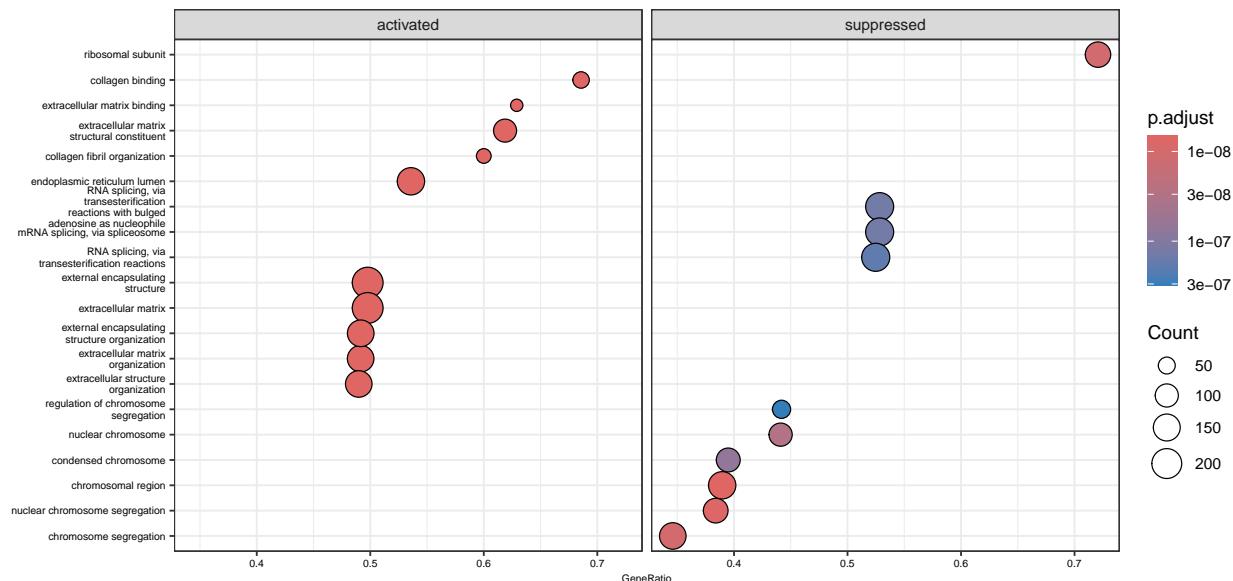
  for (a in seq_len(nrow(pathwaylist))) {
    list_gene = down_list[[a]]
    for (b in seq_len(length(list_gene))) {
      if (list_gene[b] == gene) {
        print(paste0("the pathway called ", pathwaylist$Description[[a]],
                     " contains gene ", gene, " with the NES score: ", pathwaylist$NES[[a]]))
      }
    }
  }
}

```

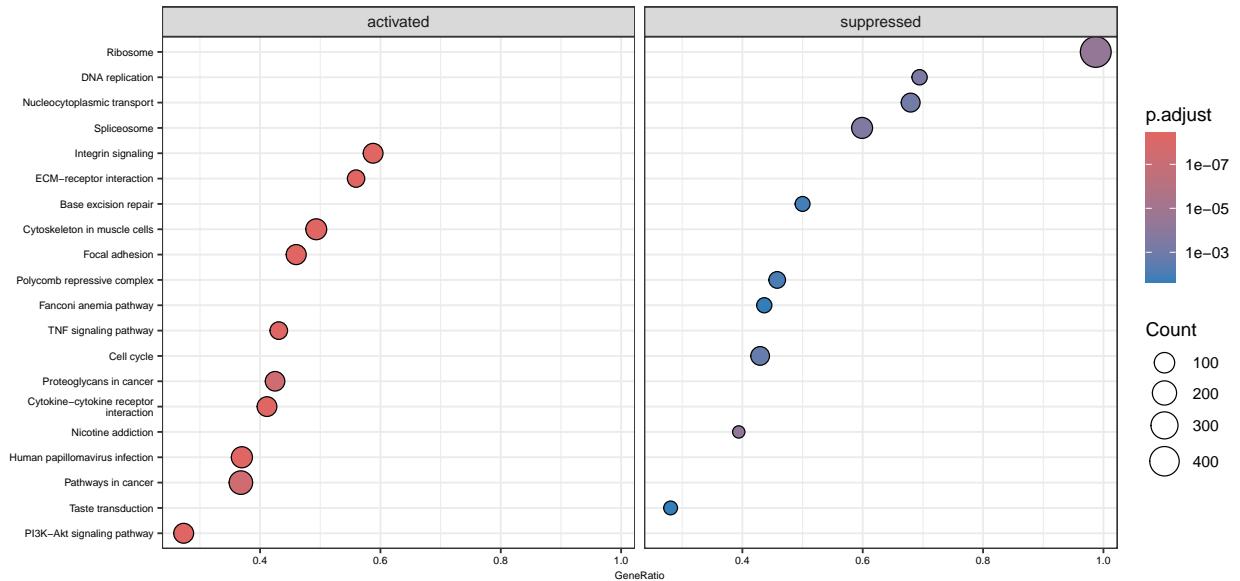
```

#Perform dotplot using GO and KEGG
dotplot(gseaGO, showCategory = 10, font.size = 6, split = ".sign") + facet_grid(.~.sign)

```

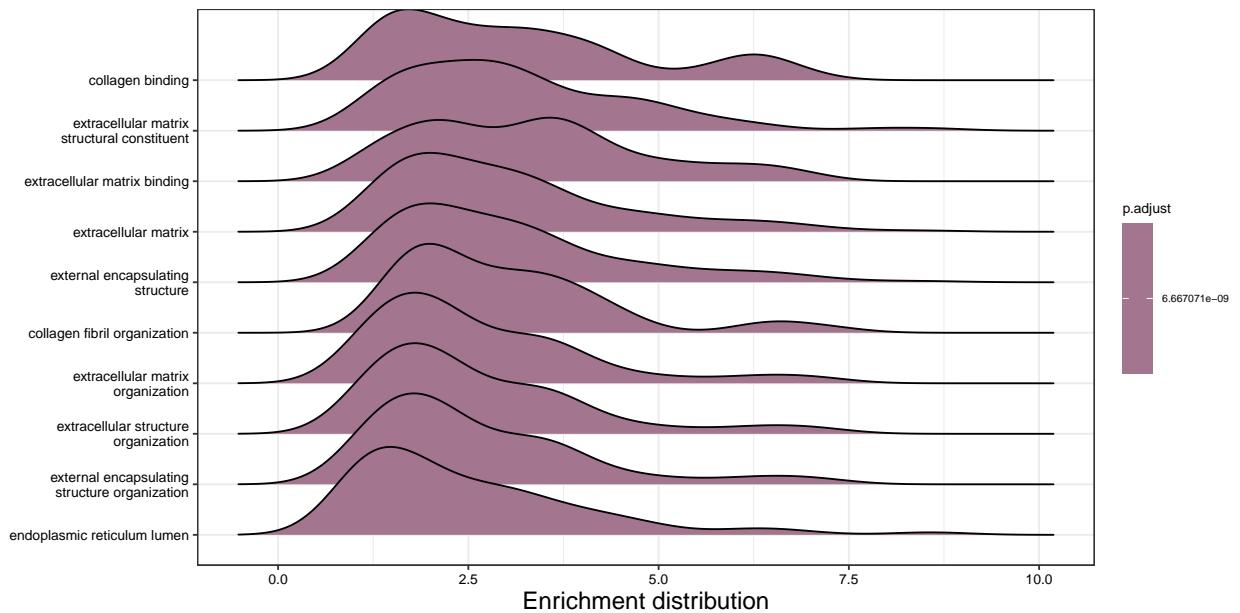


```
dotplot(gseaKEGG, showCategory = 10, font.size = 6, split = ".sign") + facet_grid(.~.sign)
```



```
#Perform ridgeplot using GO and KEGG
ridgeplot(gseaGO, showCategory = 10) + labs(x = "Enrichment distribution") +
  theme(text = element_text(size = 8),
        axis.text.x = element_text(size = 8),
        axis.text.y = element_text(size = 8))
```

## Picking joint bandwidth of 0.494

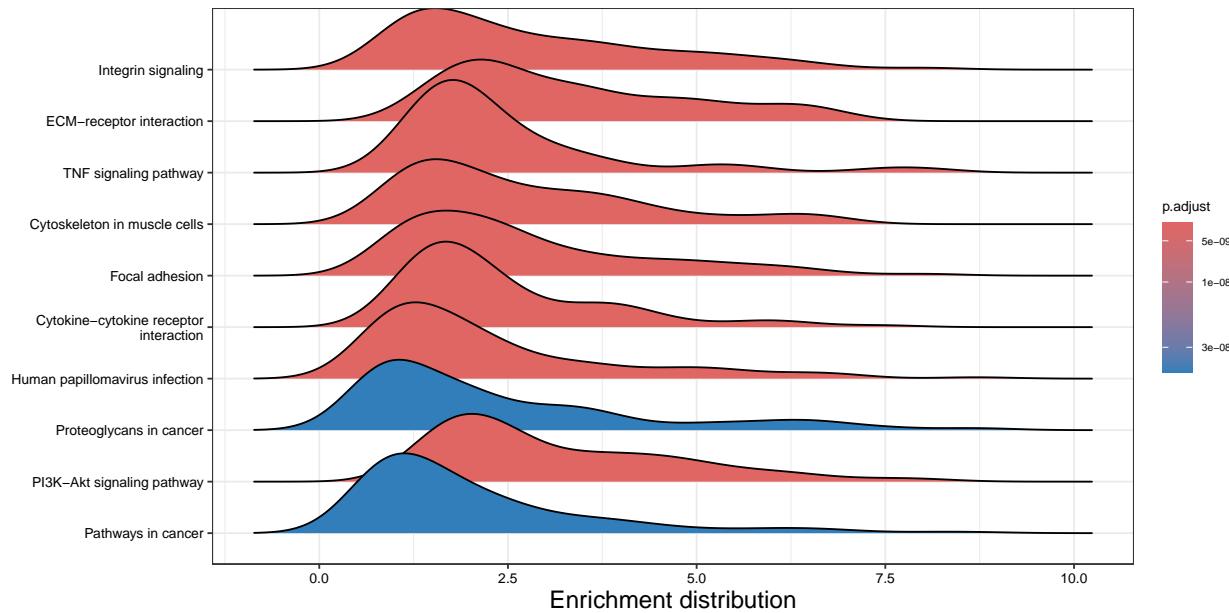


```

ridgeplot(gseaKEGG, showCategory = 10) + labs(x = "Enrichment distribution") +
  theme(text = element_text(size = 8),
        axis.text.x = element_text(size = 8),
        axis.text.y = element_text(size = 8))

```

## Picking joint bandwidth of 0.51

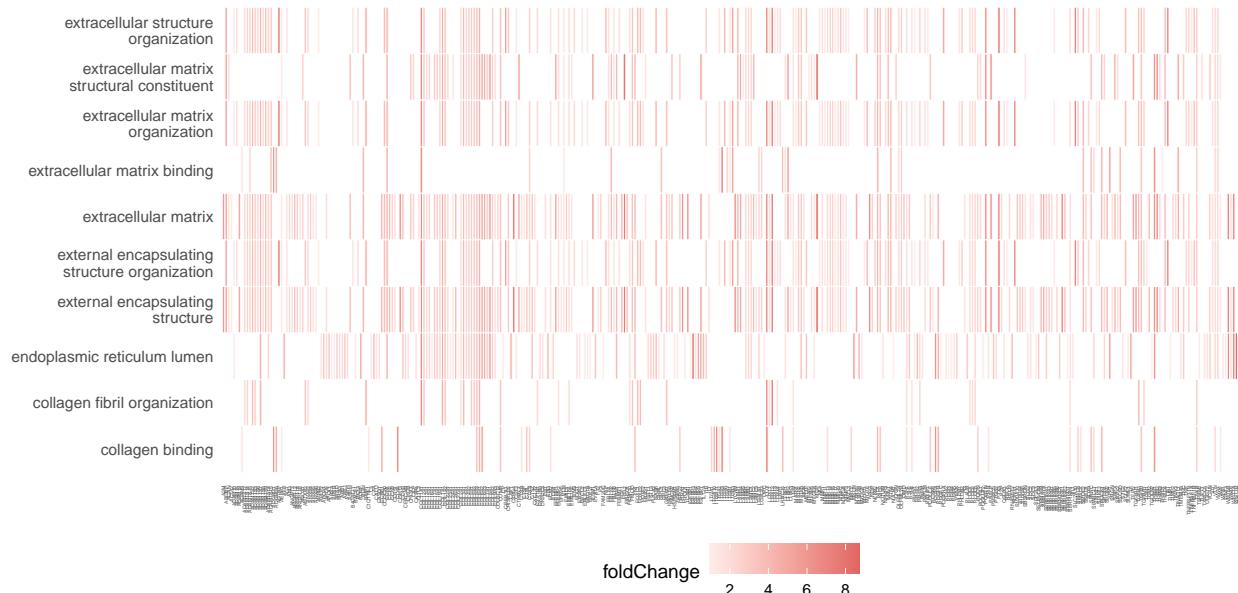


#Perform heatplot using GO and KEGG

```

heatplot(gseaGO, showCategory = 10, foldChange = GSEAlist) +
  theme(axis.text.x = element_text(size = 3, angle = 90, hjust = 1, vjust = 0.5)) +
  theme(legend.position = "bottom")

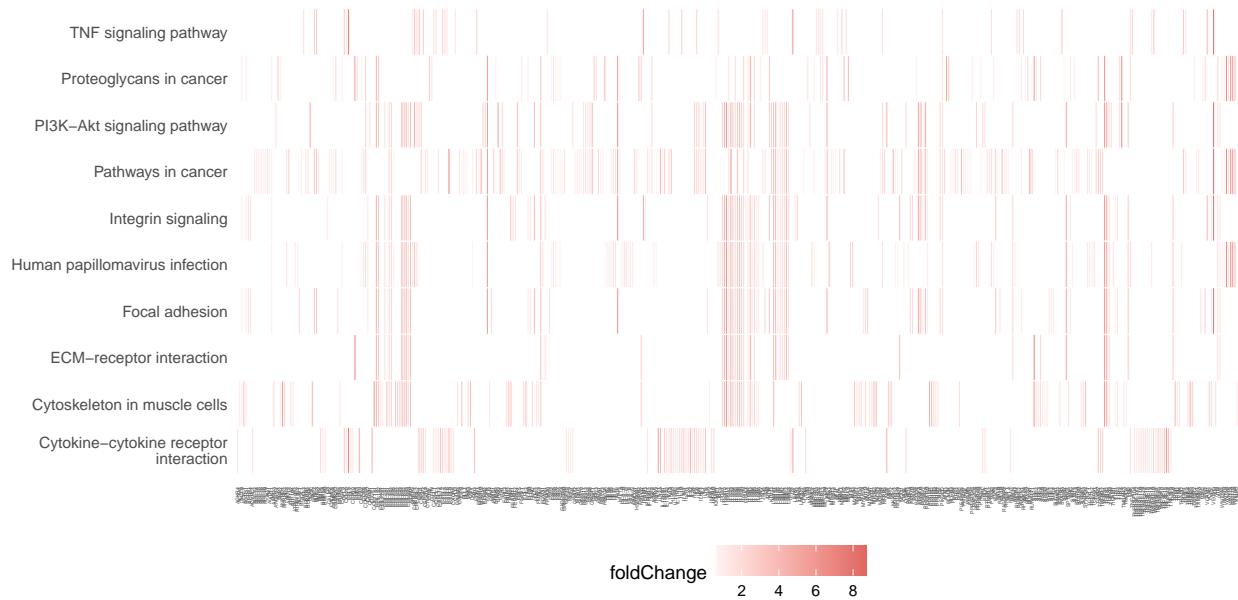
```



```

heatplot(gseaKEGG, showCategory = 10, foldChange = GSEAlist) +
  theme(axis.text.x = element_text(size = 3, angle = 90, hjust = 1, vjust = 0.5)) +
  theme(legend.position = "bottom")

```



```

#Perform category plot using GO and KEGG
#categorySize can be either 'pvalue' or 'geneNum'
#GO_cnet <- cnetplot(gseaGO,font.size=4, categorySize="geneNum", foldChange=GSEAlist,max.overlaps=50,no
#KEGG_cnet <- cnetplot(gseaKEGG,font.size=4, categorySize="geneNum", foldChange=GSEAlist,max.overlaps=5
#cowplot::plot_grid(GO_cnet, KEGG_cnet)

```