

Perform Differential Expression analysis using DESeq2 package

Load Libraries

```
library(DESeq2)
```

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

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

```
library(pheatmap)
```

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

```
library(ggplot2)
```

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

```
library(tinytex)
```

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

Read counts data

```
read_count_data <- function(file_path){  
  counts_data <- read.csv(file_path, row.names = 1)  
  return (head(counts_data))  
}
```

```
read_count_data("pasilla_gene_exp.csv")
```

```
##           treated1 treated2 treated3 untreated1 untreated2 untreated3  
## FBgn0000003         0         0         1         0         0         0  
## FBgn0000008        140         88         70         92        161        76  
## FBgn0000014         4         0         0         5         1         0  
## FBgn0000015         1         0         0         0         2         1  
## FBgn0000017       6205       3072       3334       4664       8714       3564  
## FBgn0000018        722        299        308        583        761        245  
##           untreated4  
## FBgn0000003         0  
## FBgn0000008        70  
## FBgn0000014         0  
## FBgn0000015         2  
## FBgn0000017       3150  
## FBgn0000018       310
```

Read Sample metadata

```
read_metadata <- function(file_path){  
  coldata <- read.csv(file_path, row.names = 1)  
  return (coldata)  
}  
  
read_metadata("pasilla_meta.data.csv")
```

```
##           condition      type  
## treated1      treated single-read  
## treated2      treated paired-end  
## treated3      treated paired-end  
## untreated1 untreated single-read  
## untreated2 untreated single-read  
## untreated3 untreated paired-end  
## untreated4 untreated paired-end
```

convert condition and types columns in coldata object to factor

```
convert_chr_to_factor <- function(){  
  coldata$condition <- factor(coldata$condition)  
  coldata$type <- factor(coldata$type)  
}  
  
convert_chr_to_factor()
```

make sure the row names in colData matches to the column names in counts_data

```
all(rownames(coldata) %in% colnames(counts_data))
```

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

IS the columns of the count matrix and the rows of the colData (information about samples) are in the same order?

```
all(rownames(coldata) == colnames(counts_data))
```

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

if Not, make them in the same order.

```
counts_data <- counts_data[, rownames(coldata)]
all(rownames(coldata) == colnames(counts_data))
```

Construct a DESeqDataSet.

```
deseqdataset <- DESeqDataSetFromMatrix(countData = counts_data,
                                       colData = coldata,
                                       design = ~ condition)
```

```
## Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
## design formula are characters, converting to factors
```

```
deseqdataset
```

```
## class: DESeqDataSet
## dim: 14599 7
## metadata(1): version
## assays(1): counts
## rownames(14599): FBgn0000003 FBgn0000008 ... FBgn0261574 FBgn0261575
## rowData names(0):
## colnames(7): treated1 treated2 ... untreated3 untreated4
## colData names(2): condition type
```

Pre-filtering: removing rows with low gene counts

keep rows that have at least 10 reads total

```
pre_filter <- function(){
  keep <- rowSums(counts(deseqdataset)) >=10
  deseqdataset <- deseqdataset[keep,]

  return (deseqdataset)
}

pre_filter()
```

```
## class: DESeqDataSet
## dim: 9921 7
## metadata(1): version
## assays(1): counts
## rownames(9921): FBgn0000008 FBgn0000014 ... FBgn0261574 FBgn0261575
## rowData names(0):
## colnames(7): treated1 treated2 ... untreated3 untreated4
## colData names(2): condition type
```

Set the factor level

```
factor_level <- function(){  
  deseqdataset$condition <- relevel(deseqdataset$condition, ref = "untreated")  
}  
  
factor_level()
```

Differential expression analysis

```
diff_expr_analysis <- function(){  
  deseqdataset <- DESeq(deseqdataset)  
  result <- results(deseqdataset)  
  result01 <- results(deseqdataset, alpha = 0.01 , lfcThreshold = 1.5)  
  return (result01)  
}  
  
diff_expr_analysis()
```

```
## estimating size factors
```

```
## estimating dispersions
```

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

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

```
## final dispersion estimates
```

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

```
## log2 fold change (MLE): condition untreated vs treated
```

```
## Wald test p-value: condition untreated vs treated
```

```
## DataFrame with 14599 rows and 6 columns
```

```
##           baseMean log2FoldChange    lfcSE      stat    pvalue    padj  
##           <numeric>      <numeric> <numeric> <numeric> <numeric> <numeric>  
## FBgn0000003    0.171569   -1.02604541  3.805503  0.000000  1.000000    1  
## FBgn0000008   95.144079   -0.00215142  0.223884  0.000000  1.000000    1  
## FBgn0000014    1.056572    0.49673557  2.160264  0.000000  1.000000    1  
## FBgn0000015    0.846723    1.88276170  2.106432  0.181711  0.85581    1  
## FBgn0000017  4352.592899    0.24002523  0.126024  0.000000  1.000000    1  
## ...           ...           ...           ...           ...           ...  
## FBgn0261571  8.73437e-02   -0.9002942  3.810165    0          1          1  
## FBgn0261572  6.19714e+00    0.9591315  0.777017    0          1          1  
## FBgn0261573  2.24098e+03   -0.0126158  0.112701    0          1          1  
## FBgn0261574  4.85774e+03   -0.0152569  0.193148    0          1          1  
## FBgn0261575  1.06836e+01   -0.1635594  0.938909    0          1          1
```

Top 10 differentially expression genes

```
ordered_result <- result01[order(result01$pvalue, decreasing = TRUE), ]
top10 <- head(ordered_result, n=10)
top10
```

```
## log2 fold change (MLE): condition untreated vs treated
## Wald test p-value: condition untreated vs treated
## DataFrame with 10 rows and 6 columns
##
```

	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>
## FBgn00000003	1.71569e-01	-1.02604541	3.805503	0	1	1
## FBgn00000008	9.51441e+01	-0.00215142	0.223884	0	1	1
## FBgn00000014	1.05657e+00	0.49673557	2.160264	0	1	1
## FBgn00000017	4.35259e+03	0.24002523	0.126024	0	1	1
## FBgn00000018	4.18615e+02	0.10479911	0.148280	0	1	1
## FBgn00000022	7.96749e-02	0.71976369	3.814589	0	1	1
## FBgn00000024	6.40629e+00	-0.21080708	0.690349	0	1	1
## FBgn00000028	4.38900e-01	-1.41421471	2.779524	0	1	1
## FBgn00000032	9.89730e+02	0.09190505	0.147697	0	1	1
## FBgn00000036	5.89596e-01	0.06544922	2.304771	0	1	1

Explore results

```
summary(result)
```

```
##
## out of 12359 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)      : 540, 4.4%
## LFC < 0 (down)    : 521, 4.2%
## outliers [1]      : 1, 0.0081%
## low counts [2]     : 4035, 33%
## (mean count < 7)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

```
summary(result01)
```

```
##
## out of 12359 with nonzero total read count
## adjusted p-value < 0.01
## LFC > 1.50 (up)    : 11, 0.089%
## LFC < -1.50 (down) : 7, 0.057%
## outliers [1]       : 1, 0.0081%
## low counts [2]      : 0, 0%
## (mean count < 0)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

How many adjusted p-values less than 0.01?

```
sum(result01$padj < 0.01 , na.rm = TRUE)
```

```
## [1] 18
```

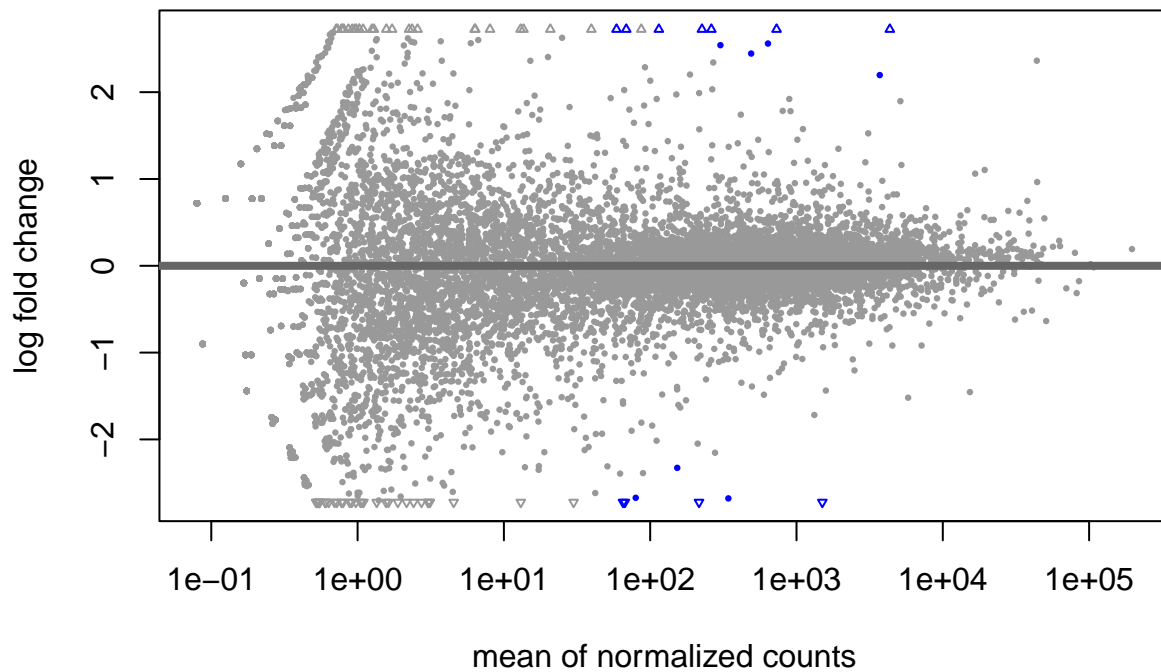
Write results to CSV file

```
write_sig_genes <- function(out_path){  
  write.csv(ordered_result, file = out_path)  
}  
  
write_sig_genes("Significant_genes.csv")
```

Visualizing the results

MA-plot

```
ma_plot <- function(){  
  plotMA(result01)  
}  
  
ma_plot()
```



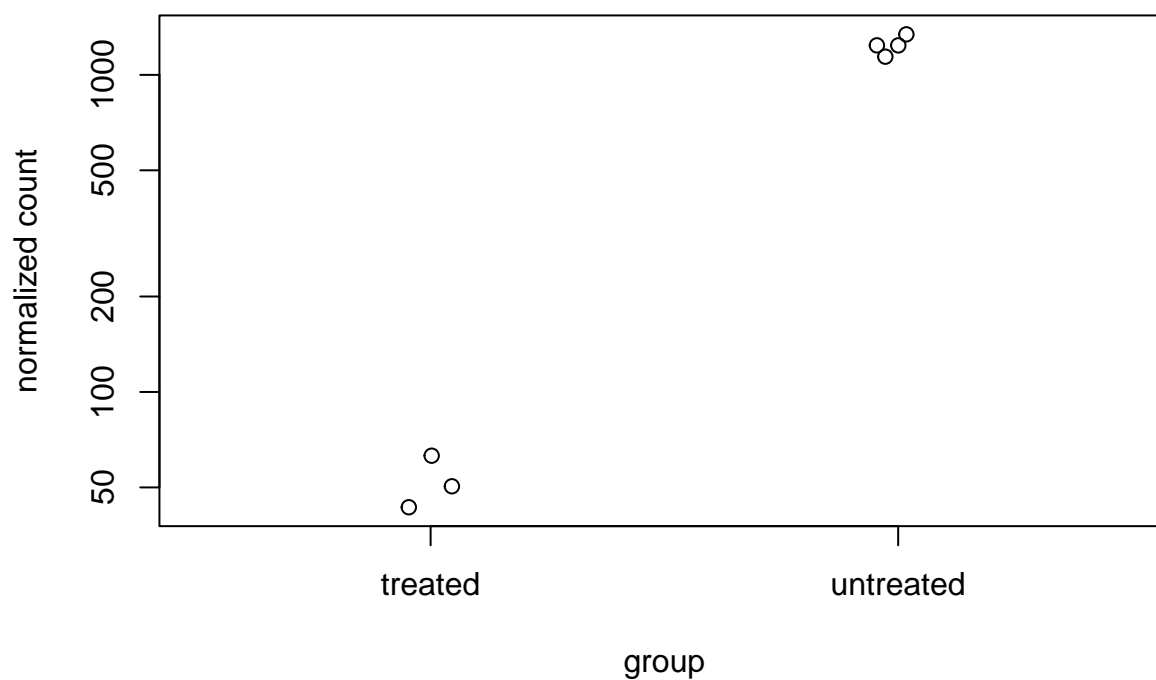
Plot counts

Here we specify the gene which had the smallest p value from the results table

```
plot_counts <- function(){
  plotCounts(deseqdataset, gene = which.min(result01$padj), intgroup = "condition")
}

plot_counts()
```

FBgn0039155

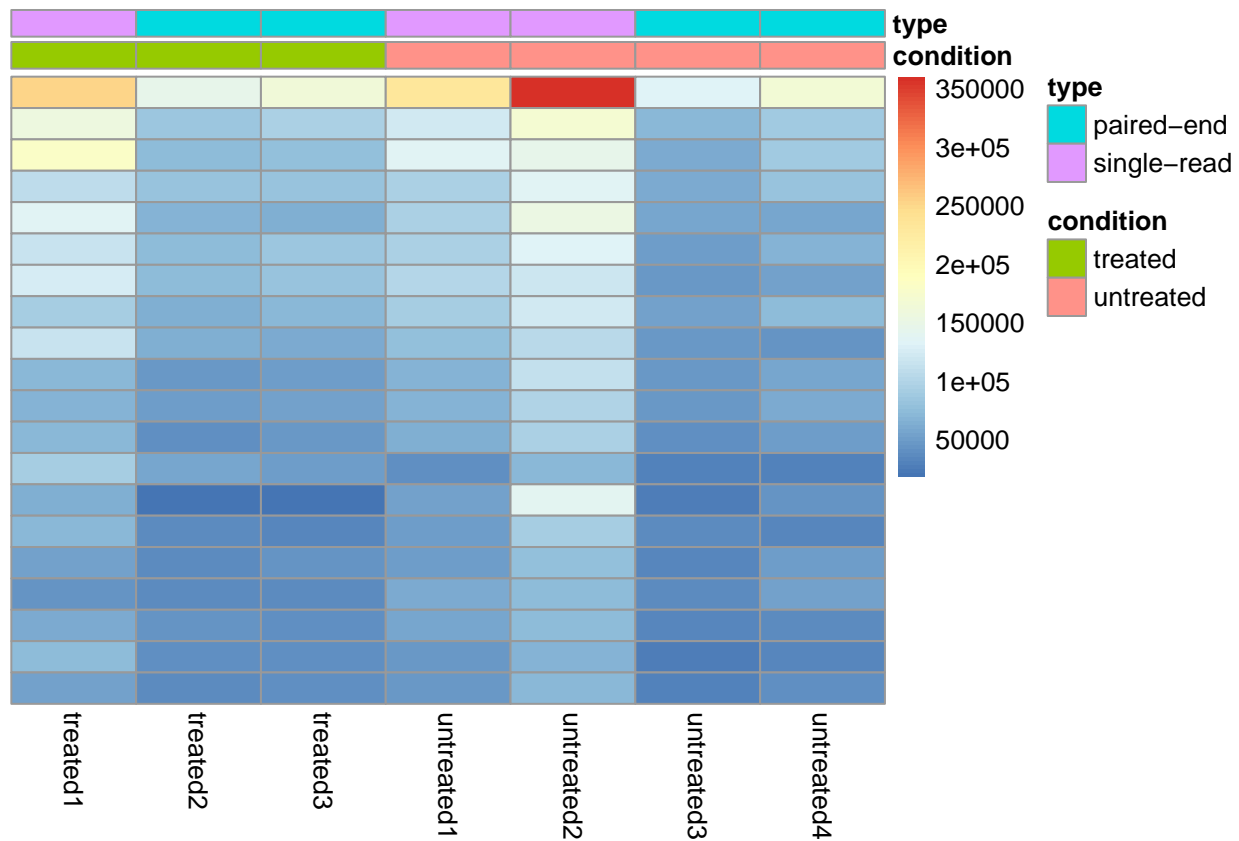


Heatmap

```
heatmap <- function(){
  select <- order(rowMeans(counts(deseqdataset)),
                  decreasing = TRUE)[1:20]

  df <- as.data.frame(colData(deseqdataset)[,c("condition", "type")])
  pheatmap::pheatmap(assay(deseqdataset)[select,], cluster_rows = FALSE,
                      show_rownames = FALSE, cluster_cols = FALSE,
                      annotation_col = df)
}

heatmap()
```

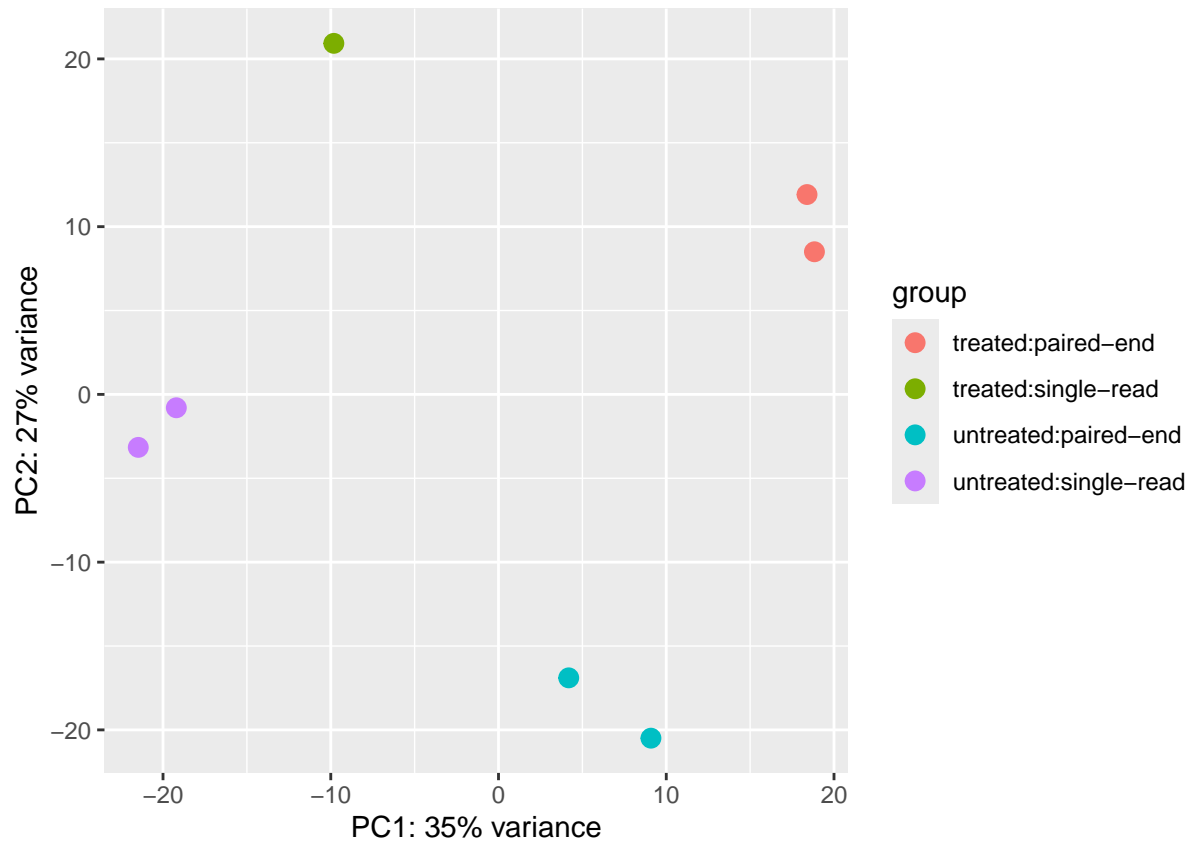



PCA plot

```
pca_plot <- function(){
  normalized = normTransform(deseqdataset)
  plotPCA(normalized, intgroup=c("condition","type"))
}
```

```
pca_plot()
```

```
## using ntop=500 top features by variance
```



Volcano plot

```
vpcano_plot <- function(){
  result.df <- as.data.frame(result)

  result.df$diffexpressed <- "NO"
  result.df$diffexpressed[result.df$log2FoldChange > 1.5 &
    result.df$padj < 0.01] <- "UP"
  result.df$diffexpressed[result.df$log2FoldChange < -1.5 &
    result.df$padj < 0.01] <- "DOWN"

  ggplot(data = result.df, aes(x = log2FoldChange, y = -log10(pvalue),
    col = diffexpressed))+
    geom_point()+ theme_minimal()+
    geom_vline(xintercept = c(-1.5, 1.5), col = "black", linetype = 'dashed') +
    geom_hline(yintercept = -log10(0.01), col = "black", linetype = 'dashed') +
    scale_color_manual(values = c("#00AFBB", "grey", "#FFDB6D"),
      labels = c("Downregulated", "Not significant", "Upregulated"))
}

vpcano_plot()
```

```
## Warning: Removed 2241 rows containing missing values or values outside the scale range
## ('geom_point()').
```

