



**Master 2 BioInformatique Modélisation et Statistique (BIMS), Université de Rouen Normandie**

**Functional Enrichment Analysis - Development of a Shiny application**



**Differential expression Analysis In ShinY for RNAseq**

**Alizée Bardon, Fiona Bottin et Sara Bencheikh**

## Entrée

## Méthode

## Visualisation

Whole data  
inspection

MAplot  
VolcanoPlot

Gènes  
Différentiellement  
Exprimés

Go term  
enrichement

ORA : EnrichGo

GSEA : gseGo

ORA only  
barplot  
goplot

Pathway  
enrichement

ORA : enrichKEGG

GSEA : gseKEGG

ORA & GSEA  
dotplot  
pathview

Liste de  
Gènes  
ordonnées

Protein  
domain  
enrichement

ORA : coder + enricher

GSEA : GSEA

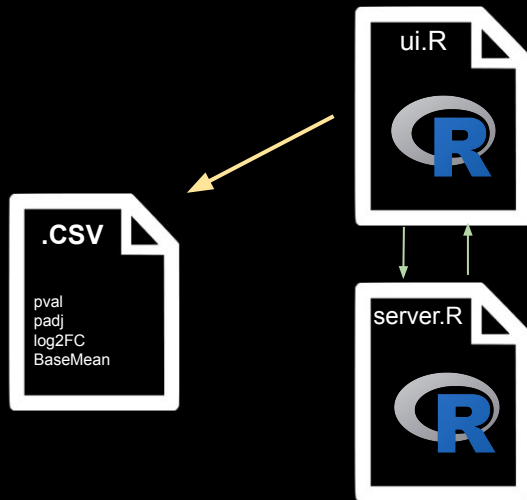
GSEA only  
gseplot



Entrée

DAISy

Interface



WGI



GO



Pathway



Protein Domain

Traduction des  
organismes

Style CSS



Images

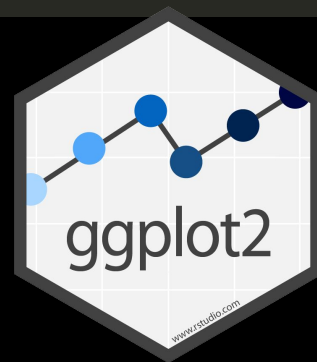
## Lancement de la démonstration



## Whole genome inspection

```
re <- reactive({  
  file <- input$file1  
  ext <- tools::file_ext(file$datapath)  
  req(file)  
  validate(need(ext == "csv", "Invalid file. Please upload a .csv file"))  
  data <- read.csv(file$datapath, header = TRUE, sep = ";")  
  data <- na.omit(data)  
  required_columns <- c("ID", "baseMean", "log2FC", "pval", "padj")  
  column_names <- colnames(data)  
  shiny::validate(need(all(required_columns %in% column_names), "Missing required columns"))  
  data  
})
```

```
p <- ggplot(  
  data = data_plot_plotly,  
  mapping = aes(x = log2FC, y = -log10(padj), col = diffexpressed, key = key)) +  
  geom_point(size = 0.45) +  
  theme_bw() +  
  scale_color_manual(values = c("red", "black", "green")) +  
  xlab("log2 fold change") +  
  ylab("-log10(p-value)")
```



# Gene Ontology enrichment

```

goGse_annot <- eventReactive(input$Run_Annotation_go, if(input$method_go == 2){
  espece_id <- espece_id()
  orga_translate_table <- orga_translate_table()
  go <- re()
  original_gene_list <- go$log2FC
  names(original_gene_list) <- go$ID
  gene_list <- na.omit(original_gene_list)
  gene_list <- sort(gene_list, decreasing=TRUE)

  gse <- gseGO(geneList=gene_list,
               ont = input$Ontology,
               keyType = 'ENSEMBL',
               pvalueCutoff = input$pvalue_go,
               minGSSize = 3,
               maxGSSize = 800,
               verbose = TRUE,
               OrgDb = orga_translate_table[espece_id,2],
               pAdjustMethod = input$Ajustement_go)
})

```

GSEA

```

goGse_enrich <- eventReactive(input$Run_Annotation_go, if(input$method_go == 1){
  espece_id <- espece_id()
  orga_translate_table <- orga_translate_table()
  go <- re()
  original_gene_list <- go$log2FC
  names(original_gene_list) <- go$ID
  gene_list <- na.omit(original_gene_list)
  gene_list = sort(gene_list, decreasing = TRUE)
  sig_genes_df = subset(go, padj < input$pvalue_go)
  genes <- sig_genes_df$log2FC
  names(genes) <- sig_genes_df$ID
  genes <- na.omit(genes)
  thresholdLog2FoldChange <- thresholdLog2FoldChange()

  if (input$type_go == "over"){
    genes <- names(genes)[genes > thresholdLog2FoldChange]
  }
  else if(input$type_go == "under") {
    genes <- names(genes)[genes < -thresholdLog2FoldChange]
  }
  else if(input$type_go == "both") {
    genes <- names(genes)[abs(genes) > thresholdLog2FoldChange]
  }

  go_enrich <- enrichGO(gene = genes,
                        universe = names(gene_list),
                        OrgDb = orga_translate_table[espece_id,2],
                        keyType = 'ENSEMBL',
                        readable = T,
                        ont = input$Ontology,
                        pvalueCutoff = input$pvalue_go,
                        pAdjustMethod = input$Ajustement_go
                        )
})

```

ORA



# KEGG, Pathway enrichment

```

kegg_data <- eventReactive(input$Run_Pathway, {
  data <- re()
  espece_id <- espece_id()
  orga_translate_table <- orga_translate_table()
  adj_method <- input$kegg_adj_method

  ids = bitr(data$ID, fromType = "ENSEMBL", toType = "ENTREZID",
             OrgDb=orga_translate_table[espece_id,2])

  dedup_ids = ids[!duplicated(ids[c("ENSEMBL")]),]
  df2 = data[data$ID %in% dedup_ids$ENSEMBL,]

  df2$Y = dedup_ids$ENTREZID

  kegg_gene_list <- df2$log2FC
  names(kegg_gene_list) <- df2$Y

  kegg_gene_list<-na.omit(kegg_gene_list)

  kegg_gene_list = sort(kegg_gene_list, decreasing = TRUE)

  if(input$method == 2){
    res <- gse_kegg(kegg_gene_list, adj_method, orga_translate_table[espece_id,])
  }
  if(input$method == 1){
    res <- ora_kegg(df2, kegg_gene_list, adj_method, orga_translate_table[espece_id,])
  }
  list(res =res, kegg_gene_list=kegg_gene_list)
})

```

GSEA

ORA

```

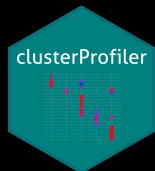
gse_kegg <- function(kegg_gene_list, adj_method, orga_translate_table) {
  if (input$db == 1){
    gse_result <- gseKEGG(geneList = kegg_gene_list,
                          organism = orga_translate_table[1,4],
                          nPerm = 10000,
                          minGSSize = 3,
                          maxGSSize = 800,
                          pvalueCutoff = input$pvalue_gsea,
                          pAdjustMethod = adj_method,
                          keyType = "ncbi-geneid")
  }
  else {
    gse_result <- gsePathway(geneList = kegg_gene_list,
                             organism = orga_translate_table[1,5],
                             nPerm = 10000,
                             minGSSize = 3,
                             maxGSSize = 800,
                             pvalueCutoff = input$pvalue_gsea,
                             pAdjustMethod = adj_method)
  }
  return(gse_result)
}

```

KEGG

Reactome

GSEA



# KEGG, Pathway enrichment

```
kegg_data <- eventReactive(input$Run_Pathway, {
  data <- re()
  espece_id <- espece_id()
  orga_translate_table <- orga_translate_table()
  adj_method <- input$kegg_adj_method

  ids = bitr(data$ID, fromType = "ENSEMBL", toType = "ENTREZID",
            OrgDb=orga_translate_table[espece_id,2])

  dedup_ids = ids[!duplicated(ids[c("ENSEMBL")]),]
  df2 = data[data$ID %in% dedup_ids$ENSEMBL,]

  df2$Y = dedup_ids$ENTREZID

  kegg_gene_list <- df2$log2FC
  names(kegg_gene_list) <- df2$Y

  kegg_gene_list<-na.omit(kegg_gene_list)

  kegg_gene_list = sort(kegg_gene_list, decreasing = TRUE)

  if(input$method == 2){
    res <- gse_kegg(kegg_gene_list, adj_method, orga_translate_table[espece_id,])
  }
  if(input$method == 1){
    res <- ora_kegg(df2, kegg_gene_list, adj_method, orga_translate_table[espece_id,])
  }
  list(res =res, kegg_gene_list=kegg_gene_list)
})
```

GSEA

ORA

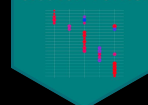
```
ora_kegg <- function(df2, kegg_gene_list, adj_method, orga_translate_table) {
  # Extract significant results from previous pvalue adjusted threshold
  kegg_sig_genes_df = subset(df2, padj < input$pvalue)
  # Vector of log2foldChange with ENTREZ IDs as names
  kegg_genes <- kegg_sig_genes_df$log2FC
  names(kegg_genes) <- kegg_sig_genes_df$Y
  kegg_genes <- na.omit(kegg_genes)
  # filter on log2fold change (under or over expressed DEG)
  if (input$type == 1){
    kegg_genes <- names(kegg_genes)[kegg_genes > input$thresholdLog2FoldChange]
  }
  else if(input$type == 2) {
    kegg_genes <- names(kegg_genes)[kegg_genes < - input$thresholdLog2FoldChange]
  }
  else {
    kegg_genes <- names(kegg_genes)[kegg_genes < - input$thresholdLog2FoldChange
                                || kegg_genes > input$thresholdLog2FoldChange]
  }
  #Database to use for annotation
  if (input$db == 1){
    ora_result <- enrichKEGG(gene=kegg_genes, universe=names(kegg_gene_list),
                           organism=orga_translate_table[1,4],
                           pvalueCutoff = input$pvalue_gsea,
                           keyType = "ncbi-geneid",
                           pAdjustMethod = adj_method)
  }
  else {
    ora_result <- enrichPathway(gene=kegg_genes, universe=names(kegg_gene_list),
                              organism=orga_translate_table[1,5],
                              pvalueCutoff = input$pvalue_gsea,
                              pAdjustMethod = adj_method)
  }
  return(ora_result)
}
```

KEGG

Reactome

ORA

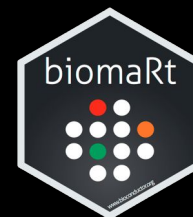
clusterProfiler





## Protein domain enrichment

```
interpro_id_raw <- getBM(  
  attributes=c('interpro', 'interpro_description', 'ensembl_transcript_id', 'ensembl_gene_id'),  
  filters = 'ensembl_gene_id',  
  values = resOrdered$ID,  
  uniqueRows = FALSE,  
  mart = ensembl)  
  
create_table_enrichment = function(GeneList, GeneRef){  
  
  Gene.Bg.ratio = get_Gene_and_Bg_ratio(GeneList = GeneList, GeneRef = GeneRef)  
  Bg.ratio = signif(100 * Gene.Bg.ratio$m / (Gene.Bg.ratio$m + Gene.Bg.ratio$n), 3)  
  Gene.ratio = signif(100 * Gene.Bg.ratio$x / Gene.Bg.ratio$k, 3)  
  
  test = hypergeom_test(x = Gene.Bg.ratio$x, k = Gene.Bg.ratio$k,  
                        m = Gene.Bg.ratio$m, n = Gene.Bg.ratio$n)  
  
  table.enrich = data.frame(interpro_ID = Gene.Bg.ratio$Term,  
                            pvalue = signif(test$pvalue_fisher.test, 3),  
                            padj = signif(test$padj, 3),  
                            BgRatio = Bg.ratio,  
                            GeneRatio = Gene.ratio,  
                            count = Gene.Bg.ratio$x,  
                            )  
  
  return (table.enrich[order(table.enrich$pval), ])
```

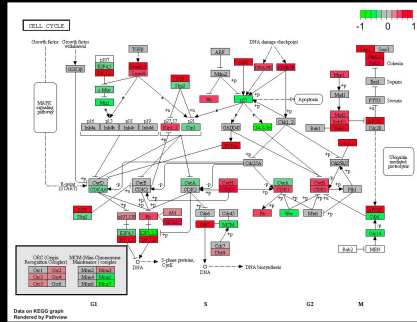


## What next ?

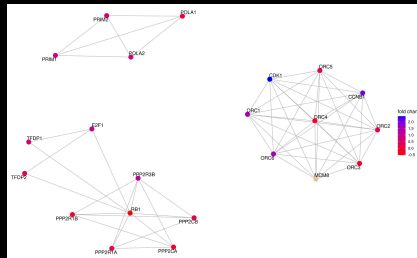


Fichier d'entrée :  
trop spécifique (**ENSEMBL**)

Nombre d'**organismes modèles** restreint ...  
et les autres ?



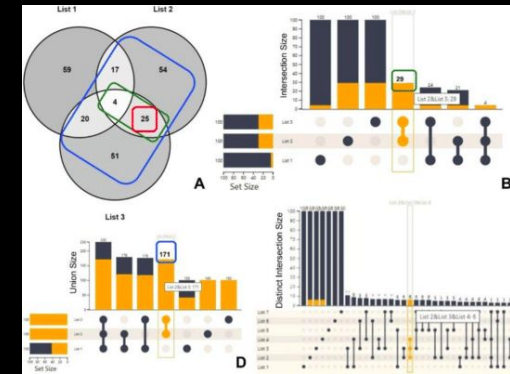
Pathview :  
optimiser l'affichage  
des **réseaux**



Reactome :  
Obtenir un  
**réseau lisible**  
(viewPathway())



Plus de **deux conditions**  
à comparer ?

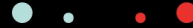


MERCI DE VOTRE ATTENTION



Annotation fonctionnelle pour :

- GO
- KEGG
- Domaines Protéiques



- Fichier d'entrée : trop spécifique (ENSEMBL)
- KEGG
  - Pathview : affichage des réseaux
  - Reactom : fonction `viewPathway()`

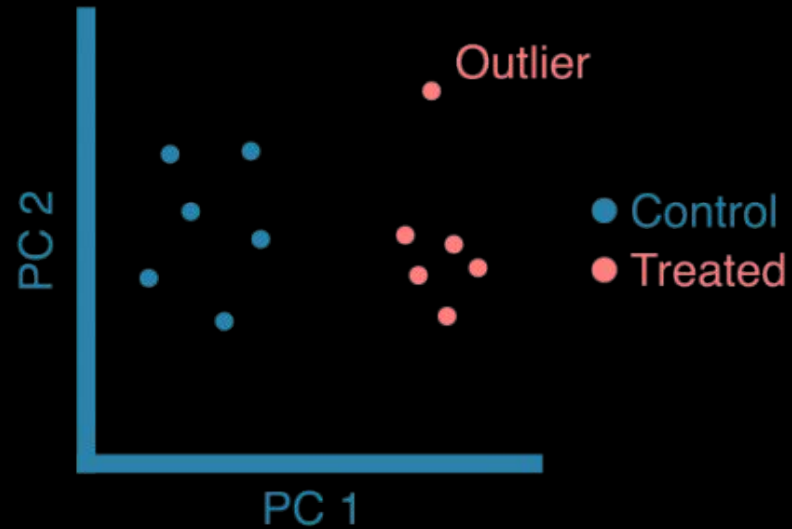


Organismes modèles... ok, et les autres ?

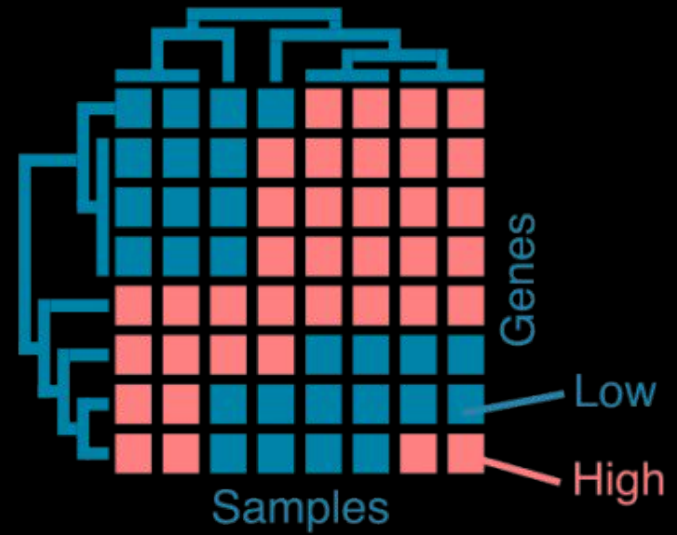
Restreint à un petit nombre d'organismes modèles

Enrichissement fonctionnel : souvent plus de deux conditions à comparer

## Principal component analysis



## Expression heatmap



## GO Term Enrichement : Plots

server.R

```

1  output$dotplot_gsea_go <-renderPlot({
2      gse<-goGse_annot()
3      require(DOSE)
4      dotplot(gse, showCategory = input$showCategory_dotplot, title = "gsea dotplot" , split=".sign") + facet_grid(.~.sign)
5  })
6
7  output$ridgeplot_go<-renderPlot({
8      gse<-goGse_annot()
9      require(DOSE)
10     ridgeplot(gse, showCategory =input$showCategory_ridgeplot)
11 })
12
13 output$gsea_plot_go <-renderPlot({
14     gse<-goGse_annot()
15     require(DOSE)
16     gseaplot2(gse, geneSetID = input$showCategory_gseaplot)
17 })

```

ui.R

```

1  conditionalPanel(
2
3      condition = "input.method_go == 2",
4
5      box(title = "Dot Plot GSEA", solidHeader = T, status = "success", width = 12, collapsible = T,id = "dotplot",
6          fluidRow(
7              column(3,
8                  wellPanel(
9                      numericInput("showCategory_dotplot", "number of categories to show", value = 5))),
10             column(9,
11                 wellPanel(
12                     shinycustomloader::withLoader(plotlyOutput("dotplot_gsea_go", height = "450px"), type = "image", loader = "wait.gif")

```



## GO Term Enrichment : Table

server.R

```

1 output$go_enrich_table_test <- DT::renderDataTable(DT::datatable({
2   if (input$method_go == 1){
3     data <- as.data.frame(goGse_enrich() )%%>%
4     mutate(interpro_link = paste0("<a href='https://amigo.geneontology.org/amigo/term/'", Description,"' target='_blank'>", Description,"</a>"))
5     col_a_afficher = c("URL",
6       "Description",
7       "GeneRatio",
8       "BgRatio",
9       "pvalue",
10      "p.adjust",
11      "qvalue")
12     data[col_a_afficher]
13   }
14   else if (input$method_go == 2){
15     data <- as.data.frame(goGse_annot() )%%>%
16     mutate(interpro_link = paste0("<a href='https://amigo.geneontology.org/amigo/term/'", Description,"' target='_blank'>", Description,"</a>"))
17     col_a_afficher = c("URL",
18       "Description",
19       "enrichmentScore",
20       "pvalue",
21       "p.adjust",
22       "rank")
23     data[col_a_afficher]
24   }
25   updateSelectInput(session,"paths")
26   DT::datatable(data)
27 },
28 escape = FALSE))
29

```

ui.R

```

1 box (
2   shinycustomloader::withLoader(dataTableOutput("go_enrich_table"), type = "image", loader = "wait.gif"),
3   width = 12) ,

```

```
re <- reactive({
  file <- input$file1
  ext <- tools::file_ext(file$datapath)
  req(file)
  validate(need(ext == "csv", "Invalid file. Please upload a .csv file"))
  data <- read.csv(file$datapath, header = TRUE, sep = ";")
  data <- na.omit(data)
  required_columns <- c("ID", "baseMean", "log2FC", "pval", "padj")
  column_names <- colnames(data)
  shiny::validate(need(all(required_columns %in% column_names), "Missing required columns"))
  data
})
```

Choose CSV File

Browse...

example.csv

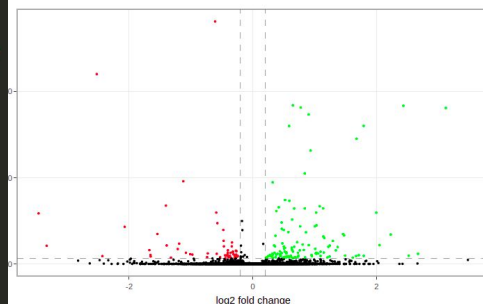
Upload complete

Input example



```
p <- ggplot(
  data = data_plot_plotly,
  mapping = aes(x = log2FC, y=-log10(padj), col=diffexpressed, key = key)) +
  geom_point(size = 0.45) +
  theme_bw() +
  scale_color_manual(values=c("red", "black", "green")) +
  xlab("log2 fold change") +
  ylab("-log10(p-value)")
```

Volcano Plot



# Enrichissement en domaines protéiques

```
interpro_id_raw <- getBM(
  attributes=c('interpro', 'interpro_description', 'ensembl_transcript_id', 'ensembl_gene_id'),
  filters = 'ensembl_gene_id',
  values = resOrdered$ID,
  uniqueRows = FALSE,
  mart = ensembl)
```

```
create_table_enrichment = function(GeneList, GeneRef){
```

```
  Gene.Bg.ratio = get_Gene_and_Bg_ratio(GeneList = GeneList, GeneRef = GeneRef)
  Bg.ratio = signif(100 * Gene.Bg.ratio$m/(Gene.Bg.ratio$m + Gene.Bg.ratio$n), 3)
  Gene.ratio = signif(100 * Gene.Bg.ratio$x / Gene.Bg.ratio$k, 3)
```

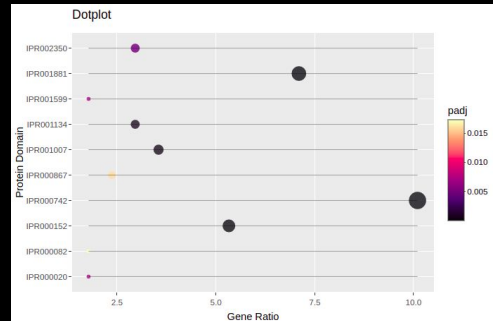
```
  test = hypergeom_test(x = Gene.Bg.ratio$x, k = Gene.Bg.ratio$k,
    m = Gene.Bg.ratio$m, n = Gene.Bg.ratio$n)
```

```
  table.enrich = data.frame(interpro_ID = Gene.Bg.ratio$Term,
    pvalue = signif(test$pvalue_fisher.test, 3),
    padj = signif(test$padj, 3),
    BgRatio = Bg.ratio,
    GeneRatio = Gene.ratio,
    count = Gene.Bg.ratio$x,
  )
```

```
  return (table.enrich[order(table.enrich$pval), ])
```

Result table

Interpro ID	Interpro Description	probe	padj	logratio	logratio_count	count	count_count	adj_pval	adj_pval_count
IPR000001	Anaphylatoxin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000002	SEA domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000003	EGF-like domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000004	EGF-like domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000005	Insulin-like growth factor-binding protein, IGFBP	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000006	Netrin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000007	Netrin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000008	Netrin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000009	Netrin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10
IPR000010	Netrin domain	0.000010	0.000000	0.0000	0.0000	10	10	0.0000	10



Alpha-2-macroglobulin, bait region domain  
 Alpha-2-macroglobulin  
 Collagen triple helix repeat  
 Insulin-like growth factor-binding protein, IGFBP  
 SEA domain  
 Alpha-2-macroglobulin, conserved site  
 Macroglobulin domain MG3  
 EGF-like domain  
 Alpha-macroglobulin, receptor-binding  
 Kazal domain superfamily  
 Netrin domain  
 Anaphylatoxin/fibulin

## Procédure [Irizarry *et al.*, 2009, Mootha *et al.*, 2003, Subramanian *et al.*, 2005]

- Soit un ensemble de gènes  $S$  pré-défini (ex : Terme GO)
- ① Calcul d'un score d'enrichissement ( $Y$  et  $N$  sont de même signe)

- Gènes sont ordonnés (par p-value, par LFC ou autre)
- Score calculé :

$$S_{c_i} = \begin{cases} +Y & \text{si le gène } i \in S \\ -N & \text{sinon.} \end{cases}$$

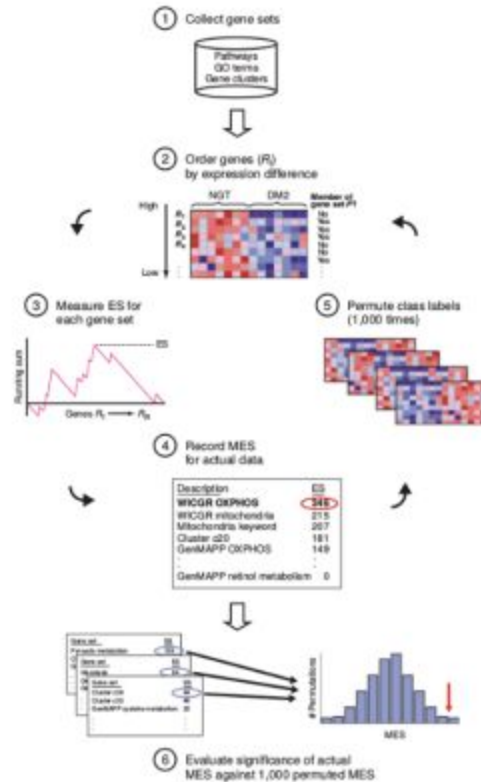
$$ES = \max_{1 \leq i \leq n} \sum_{j=1}^i S_{c_j} \text{ si } Y > 0 \text{ ou } ES = \min_{1 \leq i \leq n} \sum_{j=1}^i S_{c_j} \text{ si } Y < 0$$

- ② Estimation du niveau de significativité du score  $ES$ 
  - Permutation de l'ensemble  $S$  (1 000)
  - Calcul du score  $ES$  pour les ensembles permutés qu'on stocke
- ③ Comparaison du score  $ES$  calculé à celui des ensembles permutés
  - P-valeur calculée :

$$p = \frac{1}{1000} \sum_{j=1}^{1000} \mathbf{1}_{ES \geq ES_j} \text{ si } Y > 0 \text{ ou } p = \frac{1}{1000} \sum_{j=1}^{1000} \mathbf{1}_{ES \leq ES_j} \text{ si } Y < 0$$



# Étapes



## Résumé

- Test de randomisation/Kolmogorov-Smirnov
  - Détermine si un échantillon suit bien une loi donnée
  - Hypothèses sont :

$$H_0 : D_1 = D_2 \text{ vs } H_1 : D_1 \neq D_2$$

- Hypothèse nulle : si l'enrichissement est présent, cela est dû au hasard

### 3 App 1: Connecting **ui** and **server**

