

PAC2 Desenvolupament el treball - Fase 1

Vasyl Druchkiv

Estudiant de Màster de Bioestadística i Bioinformàtica

18 de Març 2019

Índice

1	Identificació del treball i data de l'informe	2
2	Descripció de l'avanç del projecte	2
3	L'anàlisi comú de GO, KEGG i Reactome	3
3.1	ORA	3
3.1.1	GO	3
3.1.2	KEGG	5
3.1.3	Reactome	7
3.2	GSEA	8
3.2.1	GO	8
3.2.2	KEGG	9
3.2.3	Reactome	9
3.3	Bar-Plots	10
3.4	Dot-Plots	11
3.5	Enrichment Plots	11
3.6	Category-Gene-Network Plot	12
3.7	GSEA Plot	12
4	L'anàlisi específic de GO, KEGG i Reactome	13
4.1	GO Plot	13
4.2	KEGG Pathway	14
4.3	Reactome Pathway	14

1 Identificació del treball i data de l'informe

2 Descripció de l'avanç del projecte

A la data d'avui he desenvolupat l'aplicació d'anàlisi de les rutes. L'aplicació es completament funcional localment i ofereix l'anàlisi a partir de les bases de dades GO, KEGG i Reactome. Com estava previst, l'usuari indica l'especie, puja l'arxiu amb els gens i els LogRatios provinents d'estudi de microarrays o NGS.

L'aplicació està dividida doncs en 4 parts substancials:

1. Entrada de les dades;
2. Anàlisi GO;
3. Anàlisi KEGG;
4. Anàlisi Reactome.

Figure 1: Pàgina d'entrada

L'aplicació ofereix dos mètodes d'anàlisi: d'una banda es pot fer ORA (Over-Representation Analysis) i d'altra banda l'anàlisi GSEA (Gene Set Enrichment Analysis). Recordem que l'ORA consisteix en seleccionar els gens diferencialment expressats i basant-se en GO, KEGG o Reactome comprobar si una de les agrupacions de gens suggerides per aquestes bases de dades està sobre o sotraexpressada en els gens seleccionats. Per dur a terme l'ORA l'usuari té opció de definir un *cut-off* de Log-Ratio per formal el conjunt dels gens que s'hi utilitzara (*gene set*). ORA és una bona eina per veure els efectes grans però els efectes petits li escapen. Els efectes petits derivats dels gens individuals poden acumular-se en un efect conjunt substancial el qual ORA no serà capaç de detectar. És aquí on GSEA mostra la seva utilitat.

Els apartats d'anàlisi (GO, KEGG i Reactome) ofereixen tan representacions comunes com representacions específiques.

Els anàlisis i representacions en comú són:

- Taula dels resultats ORA;
- Taula dels resultats GSEA;
- Gràfic de barres del resultat ORA;
- Gràfic de punts del resultat ORA;
- El mapa d'enriquement (Enrichment Map);
- La red dels gens en categories (Category-gene-network);
- El gràfic lde GSEA.

Els anàlisis específics són:

- GO → Gràfic GO
- KEGG → Rutes de la base de dades KEGG
- Reactome → Rutes de la base de dades Reactome

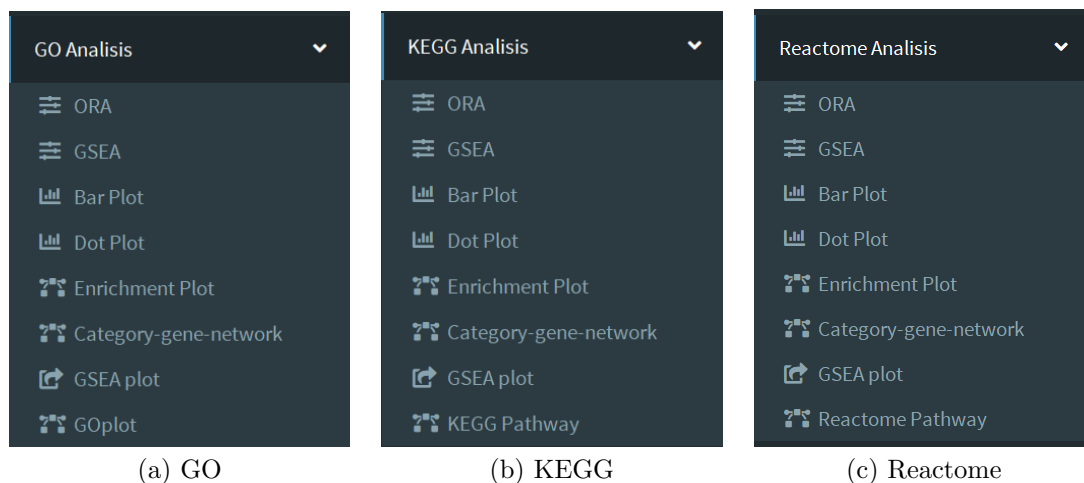


Figure 2: Els elements de les seccions d'anàlisi

3 L'anàlisi comú de GO, KEGG i Reactome

3.1 ORA

3.1.1 GO

Per realitzar l'anàlisi ORA per a termes GO s'utilitza la funció `enrichGO` del paquet `clusterProfiler`.

```
enrichGO(gene, OrgDb, keyType = "ENTREZID", ont = "MF", pvalueCutoff = 0.05,
pAdjustMethod = "BH", universe, qvalueCutoff = 0.2, minGSSize = 10, maxGSSize = 500,
readable = FALSE, pool = FALSE)
```

He implementat els valors per defecte amb la possibilitat per a usuari d'ellegir entre:

- Ontologies GO
 - Molecular function, Biological proces, Cellular Components;
- Nivell de significació basant-se en els valors de P ajustats
 - 0.1, 0.05, 0.01, 0.001;
- Mètode d'ajustament
 - Holm; Hochberg; Hommel; Bonferroni; BH; BY; FDR; None.

The screenshot shows a web application for pathway analysis. The sidebar on the left lists various analysis methods, with 'GO Analysis' and 'ORA' being prominent. The main content area is divided into three sections for configuration: 'Select ontology:' with radio buttons for CC, MF, and BP (BP is selected); 'Select adjustment method' with a dropdown menu showing 'BH'; and 'Select P-Value threshold:' with radio buttons for 0.1, 0.05, 0.01, and 0.001 (0.1 is selected). A 'Calculate Result' button is located at the bottom of the main area.

Figure 3: Especificació d'ORA dels termes GO

L'execució de la funció és un procés temporalment costós. Per aquest motiu he afegit el botó d'acció, en lloc de deixar la funció reactiu. D'aquesta manera l'usuari ha de fer una decisió consient de repetir l'anàlisi amb altres valors.

Apretant el botó apareix la taula i el botó nou mitjançant el qual l'usuari pot descarregar els resultats en format .csv. He formateat la taula amb els paquets `knitr`, `kableExtra`, `formattable` i `dplyr`. Amb els dos últims he afegit les barres de color per el nombre dels gens diferencialment expressats del terme específic de GO i el gradient de color del verd fins vermell pels valors de més petits fins els més grans.

Els camps més interessants de la taula són:

Calculate Result								
Download Results as .csv								
ID	Description	GeneRatio	BgRatio	pvalue	p.adjust	qvalue	Count	geneID
GO:0140014	mitotic nuclear division	33/193	232/11468	0.000	4.00e-18	0.000	33	CDCA8/CDC20/KIF23/CENPE/MYBL2/CCNB2/...
GO:0000280	nuclear division	35/193	316/11468	0.000	4.50e-16	0.000	35	CDCA8/CDC20/KIF23/CENPE/MYBL2/CCNB2/...

Figure 4: El resultat d'anàlisi ORA. GO.

- Description. El nom del terme GO;
- GeneRatio. El quotient: $\frac{\text{Nombre dels gens diferencialment expressats}}{\text{Nombre total dels gens en la mostra}}$;
- BgRatio. El quotient: $\frac{\text{Nombre dels gens de la ruta}}{\text{Nombre total dels gens en la base de dades GO}}$;
- p.adjust. El valor de P ajustat.

3.1.2 KEGG

Per l'ORA de base de dades KEGG he utilitzat la funció `enrichKEGG()` del paquet `clusterProfiler`.

```
enrichKEGG(gene, organism = "hsa", keyType = "kegg", pvalueCutoff = 0.05,
pAdjustMethod = "BH", universe, minGSSize = 10, maxGSSize = 500,
qvalueCutoff = 0.2, use_internal_data = FALSE)
```

Com en el cas de l'anàlisi dels termes GO també aquí l'usuari té la llibertat d'elegir l'organisme, el *cut-off* del valor de P i el mètode d'ajustament. Perquè la funció necessita el codi kegg', 'ncbi-geneid', 'ncbi-proteinid' o 'uniprot' l'usuari ha d'especificar altra vegada l'especie. La llista de les especies disponibles per a anàlisi de KEGG és molt llarga. Per aquest motiu he habilitat l'eina de cerca d'espècie per a usuari.

Pathway analysis

📁 Data input

GO Analysis

KEGG Analysis

📊 ORA

📊 GSEA

📊 Bar Plot

📊 Dot Plot

📊 Enrichment Plot

📊 Category-gene-network

📊 GSEA plot

📊 KEGG Pathway

Reactome Analysis

Enter Search Term for Specie

homo

Select KEGG Specie

Homo sapiens

Select P-Value threshold:

☒ 0.1 ☐ 0.05 ☐ 0.01 ☐ 0.001

Calculate Result

Figure 5: El resultat d'anàlisi ORA. KEGG.

6

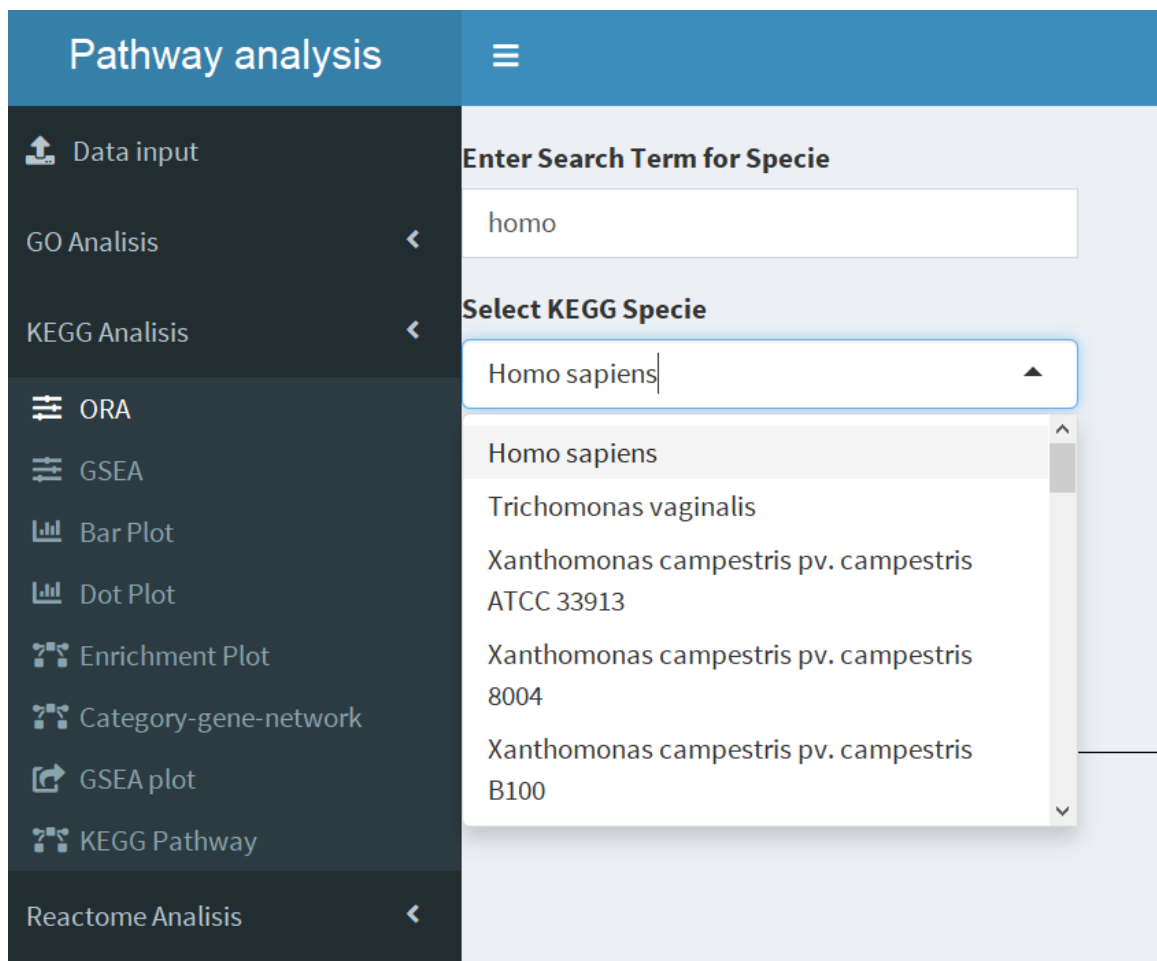


Figure 6: L'eina de cerca d'espècie. KEGG.

Una vegada introduïts els paràmetres i apretat el botó **Calculate** apareix el botó **Download .csv** i la taula previsualitzada. Els camps de la taula són els mateixos com d'anàlisi dels termes GO.

Calculate Result								
Download Results as .csv								
ID	Description	GeneRatio	BgRatio	pvalue	p.adjust	qvalue	Count	geneID
hsa04110	Cell cycle	11/92	124/7841	0.000	3.48e-05	0.000	11	8318/991/9133/890/983/4085/7272/1111/891/4174/9232
hsa04114	Oocyte meiosis	10/92	125/7841	0.000	1.70e-04	0.000	10	991/9133/983/4085/51806/6790/891/9232/3708/5241
hsa04218	Cellular senescence	10/92	160/7841	0.000	1.04e-03	0.001	10	2305/4605/9133/890/983/51806/1111/891/776/3708

Figure 7: El resultat d'anàlisi ORA. KEGG.

3.1.3 Reactome

Al cas de Reactome el procediment és similar. La funció usada és `enrichPathway()` del paquet `ReactomePA`:

```
enrichPathway(gene, organism = "human", pvalueCutoff = 0.05,
```

```
pAdjustMethod = "BH", qvalueCutoff = 0.2, universe, minGSSize = 10,
maxGSSize = 500, readable = FALSE)
```

Aquí l'usuari ha de seleccionar l'altra vegada l'organisme. Les opcions disponibles són: "human", "rat", "mouse", "celegans", "yeast", "zebrafish", "fly".

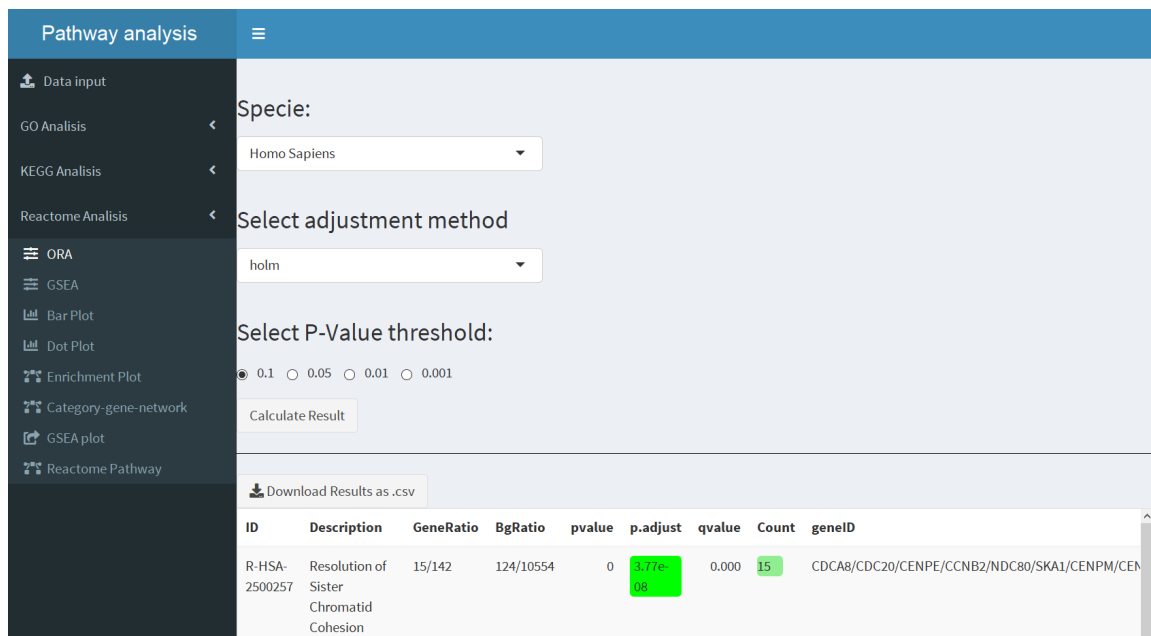


Figure 8: El resultat d'anàlisi ORA. Reactome.

3.2 GSEA

3.2.1 GO

El mètode GSEA per a termes GO es calcula amb la funció `gseGO()` del paquet `clusterProfiler`.

```
gseGO(geneList, ont = "BP", OrgDb, keyType = "ENTREZID",
      exponent = 1, nPerm = 1000, minGSSize = 10, maxGSSize = 500,
      pvalueCutoff = 0.05, pAdjustMethod = "BH", verbose = TRUE,
      seed = FALSE, by = "fgsea")
```

L'usuari pot elegir l'ontologia GO, el *cut-off* del valor P i el mètode d'ajustament.

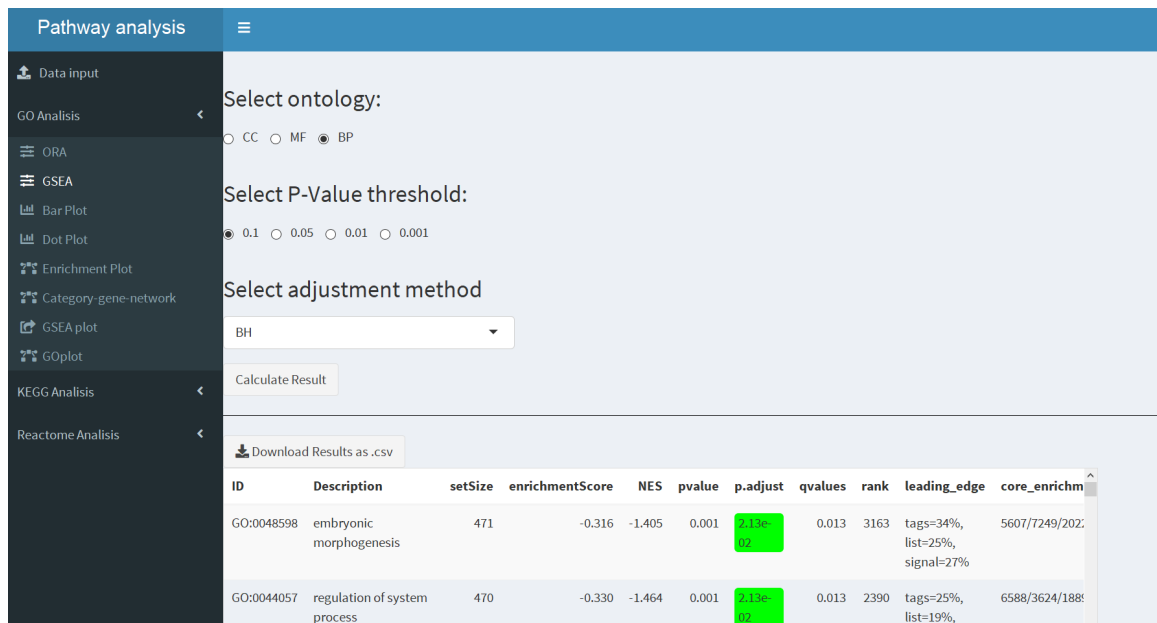


Figure 9: El resultat d'anàlisi GSEA. GO.

3.2.2 KEGG

De la mateixa manera es calcula GSEA amb la funció `gseKEGG()` del paquet `clusterProfiler`:

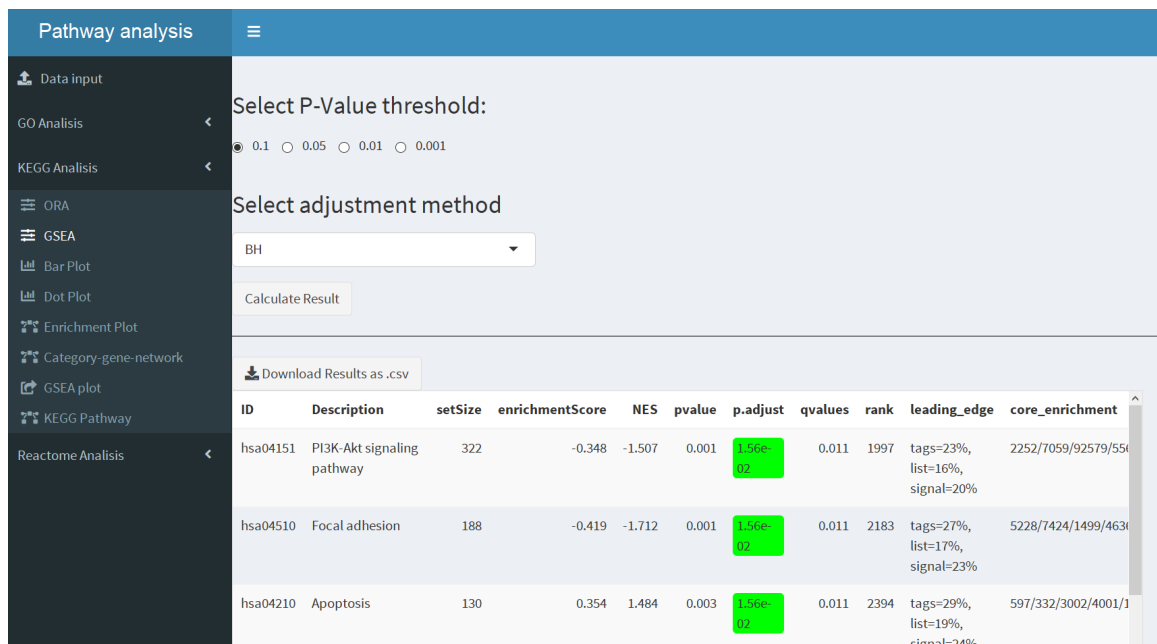


Figure 10: El resultat d'anàlisi GSEA. KEGG.

3.2.3 Reactome

Per completar l'anàlisi l'usuari pot calcular GSEA per a base de dades Reactome. Com als altres casos utilitzo el paquet `clusterProfiler` i específicament la funció `gsePathway()`

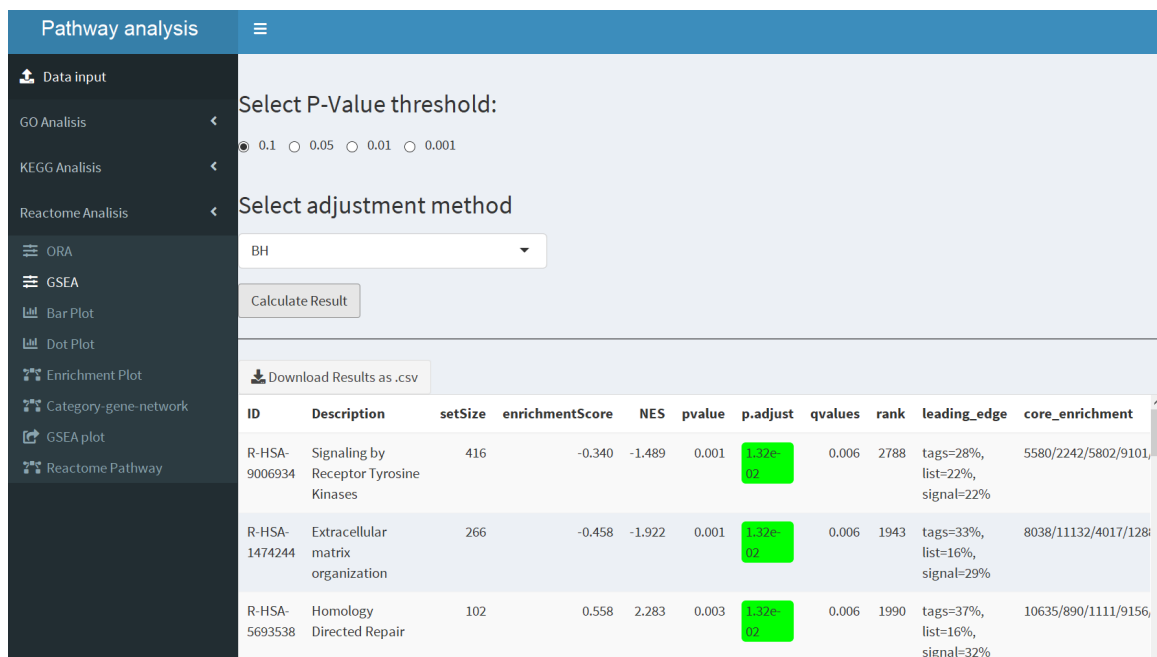


Figure 11: El resultat d'anàlisi GSEA. Reactome.

3.3 Bar-Plots

Els resultats de `enrichGO`, `enrichKEGG` i `enrichPathway` es pot visualitzar amb el gràfic de barres. L'usuari pot elegir el nombre de les categories visualitzades entre 2 i 30. Es dona l'opció per descarregar el gràfic en format .png.

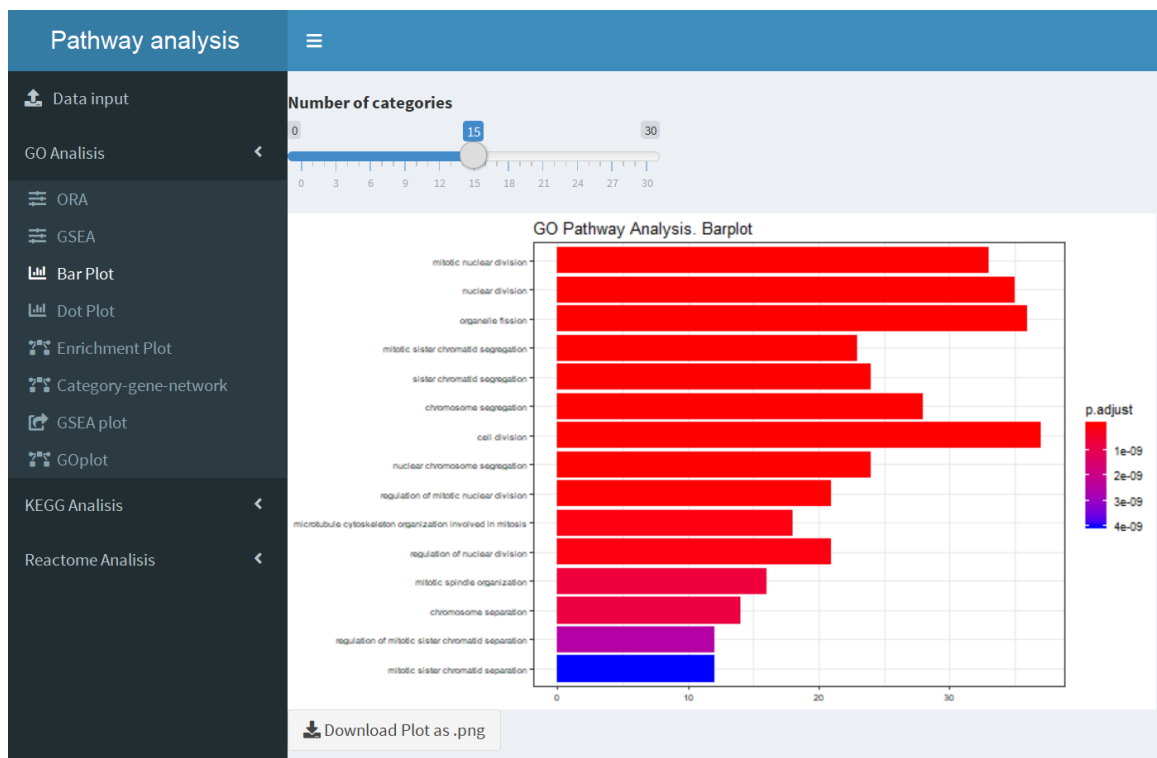


Figure 12: Bar-Plot. GO.

3.4 Dot-Plots

El *dot plot* visualitza addicionalment el *gen ratio*. També aquí l'usuari pot seleccionar el nombre de les categories.

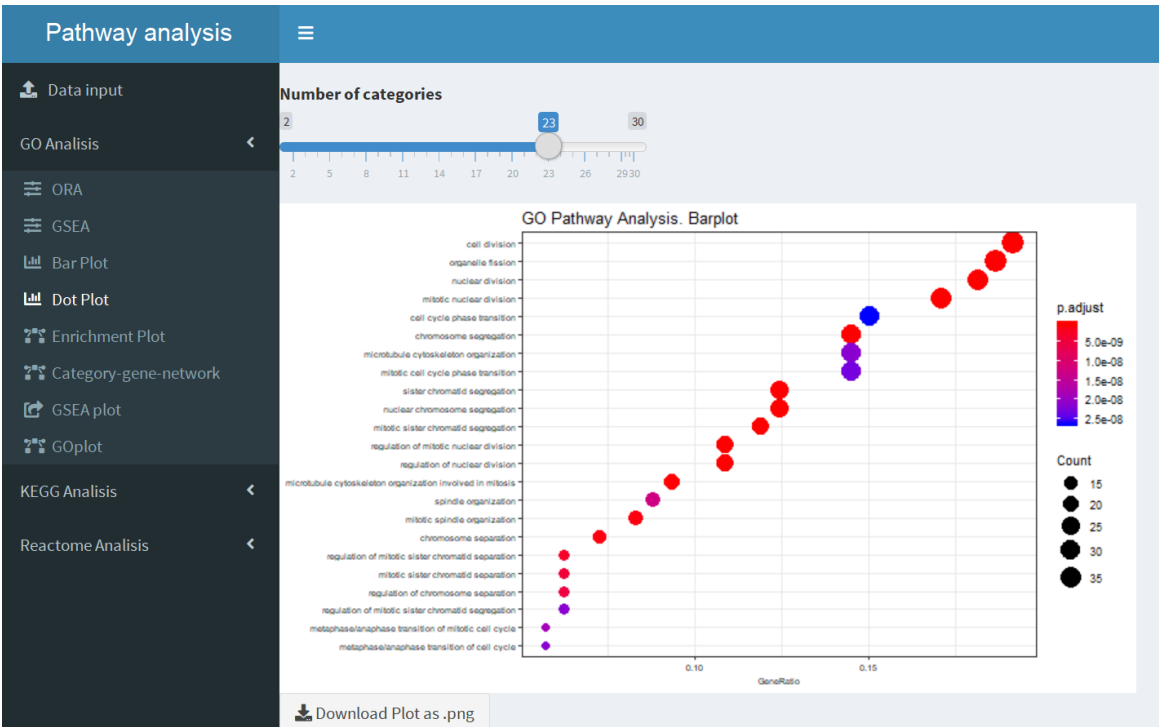


Figure 13: Bar-Plot. GO.

3.5 Enrichment Plots

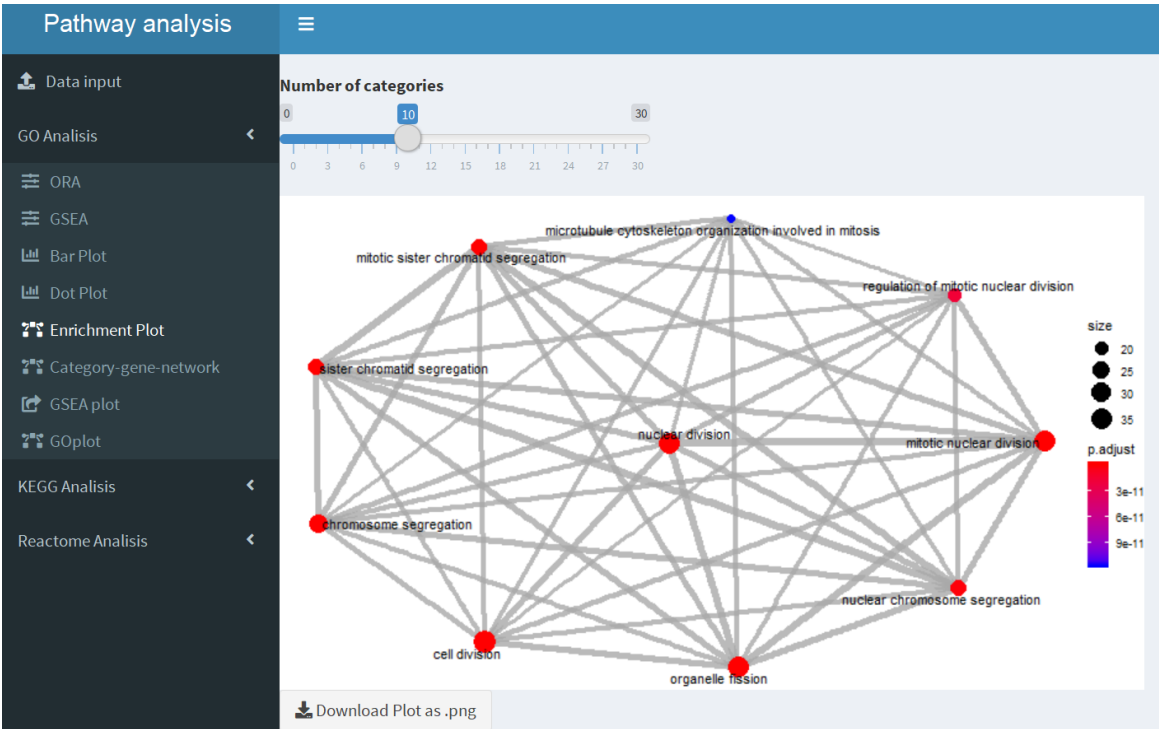


Figure 14: Bar-Plot. GO.

3.6 Category-Gene-Network Plot

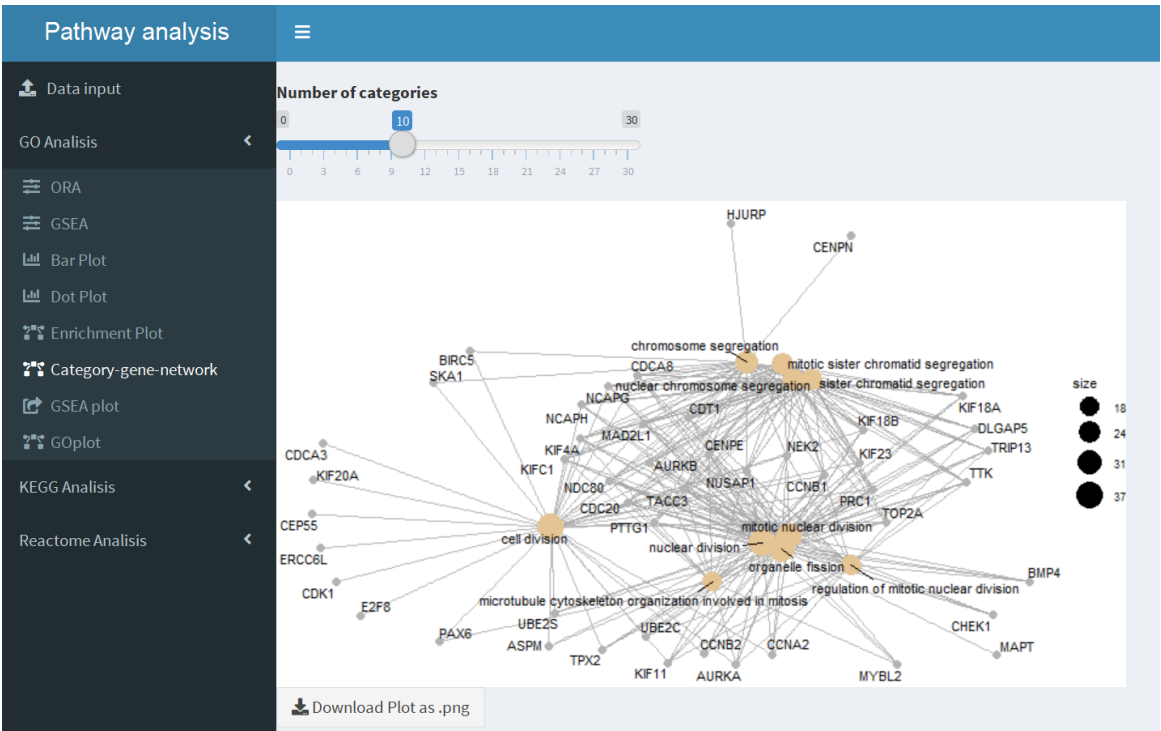


Figure 15: Category-Gene-Network Plot. GO.

3.7 GSEA Plot

L'usuari pot visualitzar una de les categories disponibles via *dropdown list*. El llistat inclou totes les rutes generades durant l'anàlisi GSEA en els apartats *Go Analysis*→*GSEA*; *KEGG*→*GSEA*

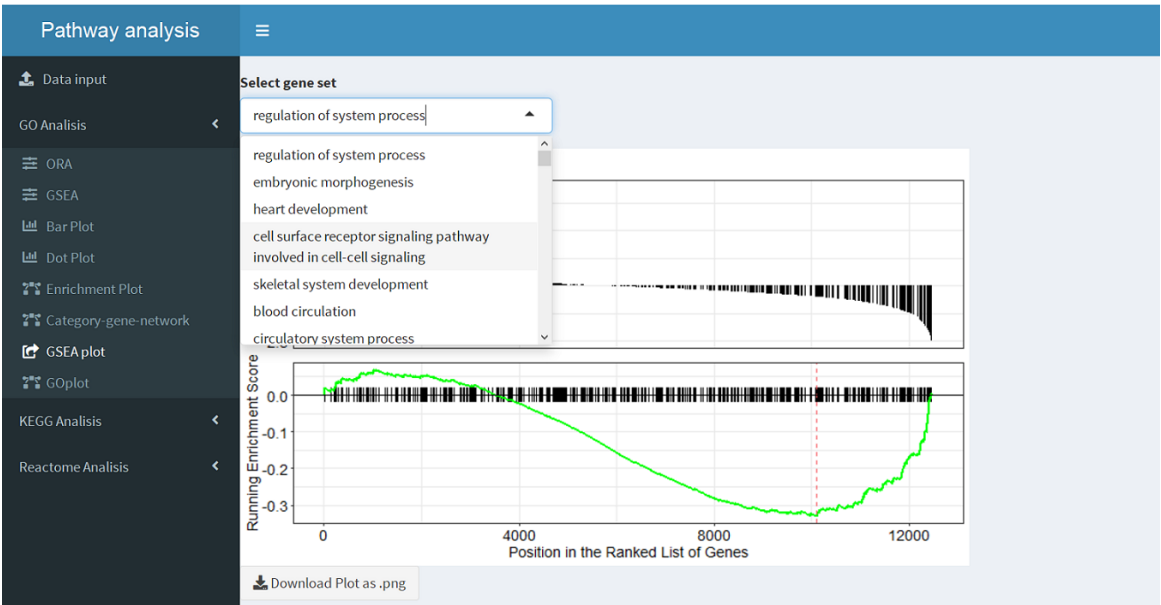


Figure 16: GSEA Plot. GO.

4.2 KEGG Pathway

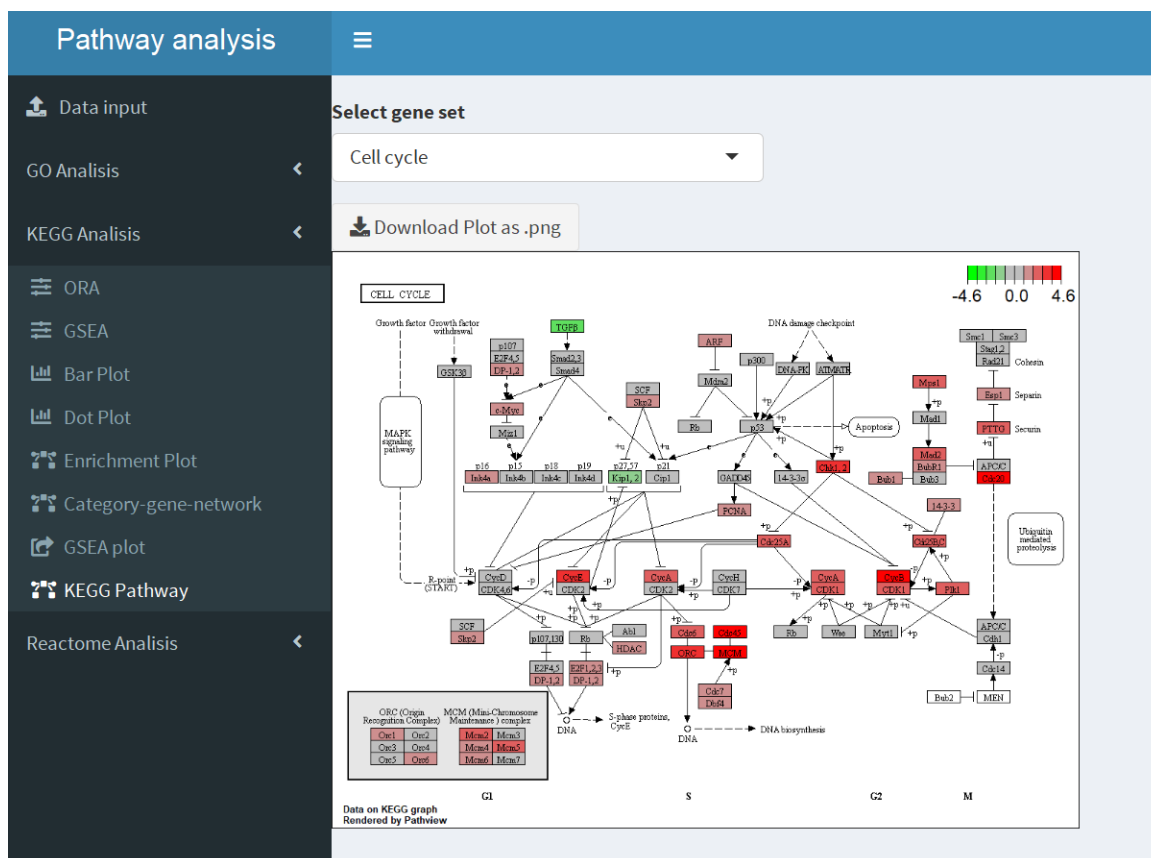


Figure 18: KEGG pathway

4.3 Reactome Pathway

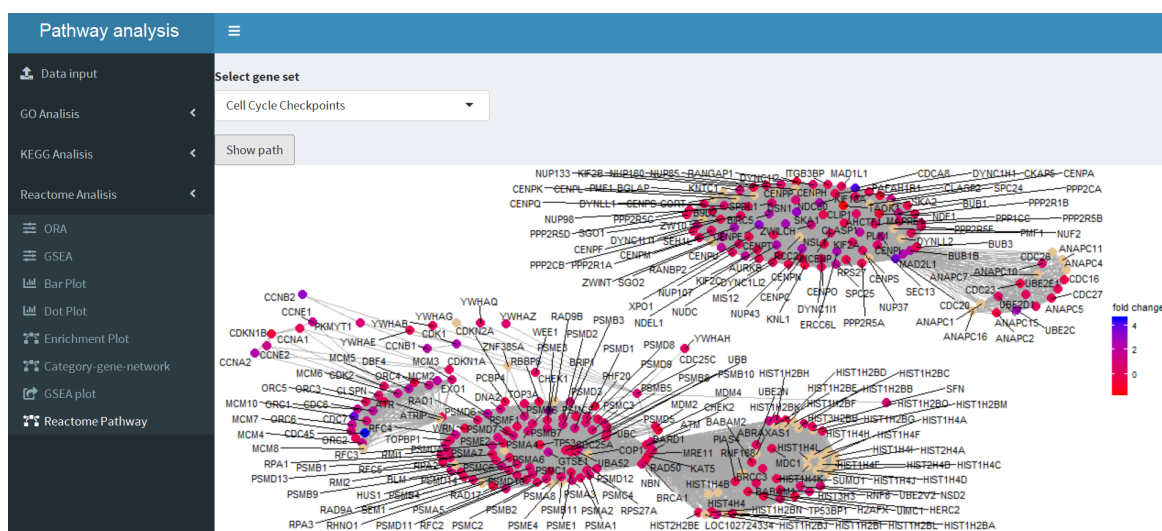


Figure 19: Reactome pathway