# The Gene Set Enrichment Analysis (GSEA)
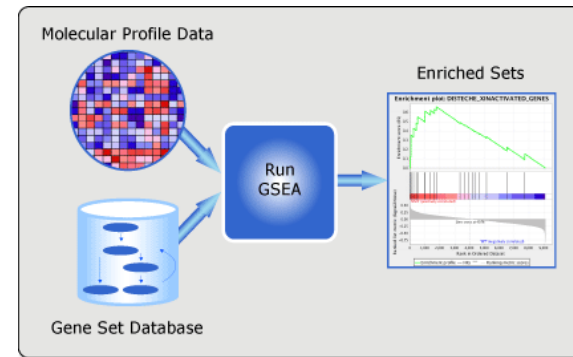
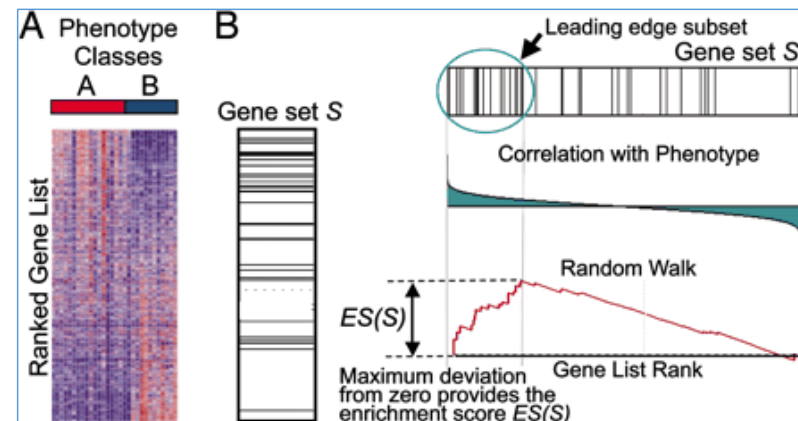**Rossella De Cegli, PhD**

(**TIGEM BAD days**)

Monday July 17th 2023

- GSEA is a computational method that determines whether an a priori defined set of genes shows **statistically significant, concordant differences** between two biological states (phenotypes, conditions, treatments).



- ✧ GSEA considers experiments with genome-wide expression profiles from samples belonging to two classes, (**TREATED vs UNTREATED, KO vs WT** etc..).

- ✧ Genes are ranked based on the correlation between their expression and the class distinction by using any suitable metric

https://www.gsea-msigdb.org/gsea/msigdb/index.jsp

# I Step:

**Generate your ranked list of interest (3'DGE…Bulk_RNA-seq…. proteomics)**

➢ The excel file will be converted in a **.rnk file**

➢ The complete gene list has to be in a format of two columns: one gene in one row

➢ To calculate the **rank**: use the formula **=-LOG10(FDR)\*logFC** and rank the signed ratio (array) or the LogFC (RNA-seq/quantseq) from the top up SIGNIFICANT (FDR<5%) to the top dw values

| ensembl_gene_id | hgnc_symbol | logFC_KO vs CTR | FDR_KO vs CTR | rank | | hgnc_symbol | rank |
|---|---|---|---|---|---|---|---|
| ENSG00000222328 | RNU2-2P | 3,137904757 | 1,43E-17 | 52,8545466 | | RNU2-2P | 52,8545466 |
| ENSG00000239002 | SCARNA10 | 3,227606082 | 9,07E-16 | 48,5516614 | | OASL | 51,474997 |
| ENSG00000212232 | SNORD17 | 2,709275191 | 6,46E-14 | 35,7340561 | | SCARNA10 | 48,5516614 |
| ENSG00000202538 | RNU4-2 | 3,227930663 | 1,98E-13 | 41,0066937 | | RSAD2 | 46,1273123 |
| ENSG00000200795 | RNU4-1 | 3,245419501 | 3,02E-13 | 40,6319126 | | RNU4-2 | 41,0066937 |
| ENSG00000274585 | RNU2-1 | 2,435178749 | 1,83E-12 | 28,5852458 | | RNU4-1 | 40,6319126 |
| ENSG00000252010 | SCARNA5 | 3,389028221 | 4,33E-12 | 38,5124226 | | SCARNA5 | 38,5124226 |
| ENSG00000263934 | SNORD3A | 2,443440214 | 4,93E-12 | 27,6285821 | | SNORD17 | 35,7340561 |
| ENSG00000251791 | SCARNA6 | 2,906943947 | 1,60E-10 | 28,4742979 | | CCL5 | 31,389108 |
| ENSG00000226869 | LHFPL3-AS1 | -1,903587798 | 2,71E-09 | -16,30861 | | SNORA53 | 31,3534592 |
| ENSG00000212443 | SNORA53 | 3,732852679 | 3,99E-09 | 31,3534592 | | RNU2-1 | 28,5852458 |
| ENSG00000135637 | CCDC142 | 2,449735305 | 9,81E-09 | 19,6185931 | | IFNL1 | 28,5796979 |
| ENSG00000200959 | SNORA74A | 3,164749749 | 1,53E-08 | 24,7337411 | | SCARNA6 | 28,4742979 |
| ENSG00000162892 | IL24 | 2,565797823 | 4,20E-08 | 18,9264374 | | SNORD3A | 27,6285821 |
| ENSG00000130766 | SESN2 | 2,821622207 | 1,12E-07 | 19,6086319 | | SNORA74A | 24,7337411 |
| ENSG00000141682 | PMAIP1 | 2,843487551 | 7,55E-07 | 17,4074253 | | IFIT2 | 24,2280865 |
| ENSG00000135114 | OASL | 8,427181242 | 7,79E-07 | 51,474997 | | TICAM2 | 21,7659126 |
| ENSG00000188467 | SLC24A5 | -1,951918932 | 8,16E-07 | -11,883513 | | IFNB1 | 21,185542 |
| ENSG00000277209 | RPPH1 | 2,17926035 | 1,14E-06 | 12,9502833 | | CCDC142 | 19,6185931 |
| ENSG00000124762 | CDKN1A | 2,024219404 | 1,20E-06 | 11,9882419 | | SESN2 | 19,6086319 |
| ENSG00000136826 | KLF4 | 3,150767798 | 2,18E-06 | 17,8383461 | | TCAF2C | 19,5672415 |
| ENSG00000134321 | RSAD2 | 8,673723027 | 4,81E-06 | 46,1273123 | | IFIT3 | 19,1834078 |
| ENSG00000160307 | S100B | -1,943647342 | 4,81E-06 | -10,336418 | | IL24 | 18,9264374 |
| ENSG00000080166 | DCT | -2,285294013 | 5,77E-06 | -11,972486 | | TRPV6 | 18,3373838 |

➢ Copy the two columns by using the application **TextWrangler (BBEdit)** to obtain the .rnk file.

➢ **What is a Custom list:** it is a list of interest as a list of targets of a specific TF; a list of lysosomal genes, or of immune response-related genes…)

➢ This list may be generated as an excel file but to load the data this excel must be converted in a **.gmt file** (please use the application **TextWrangler**).

➢ The list has to be in a format of one row (past special as TRANSPOSE)

➢ In the first column of the obtained one row-list please write the NAME of your study

## II StepB: Run a MsigDB GSEA

➢ The GSEA gene sets are divided into 8 major collections:

**The Frequently used:**

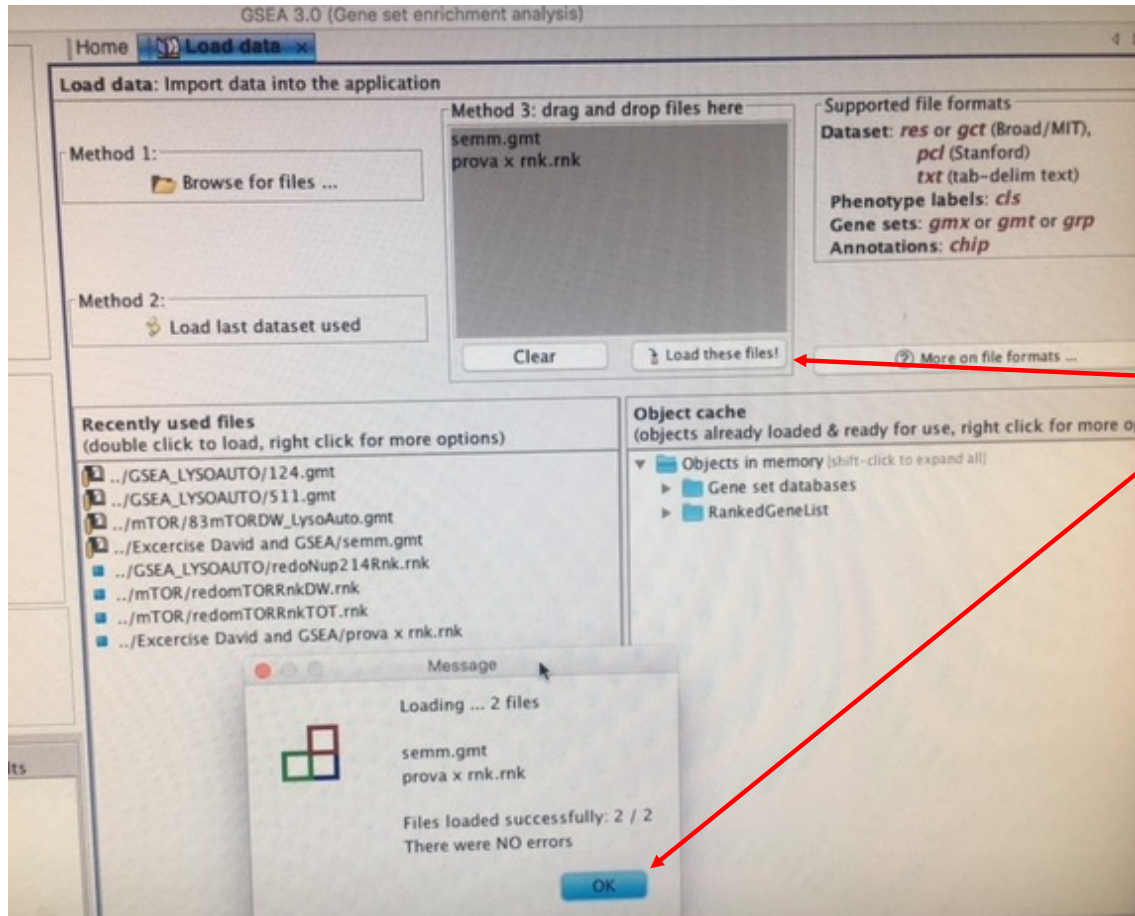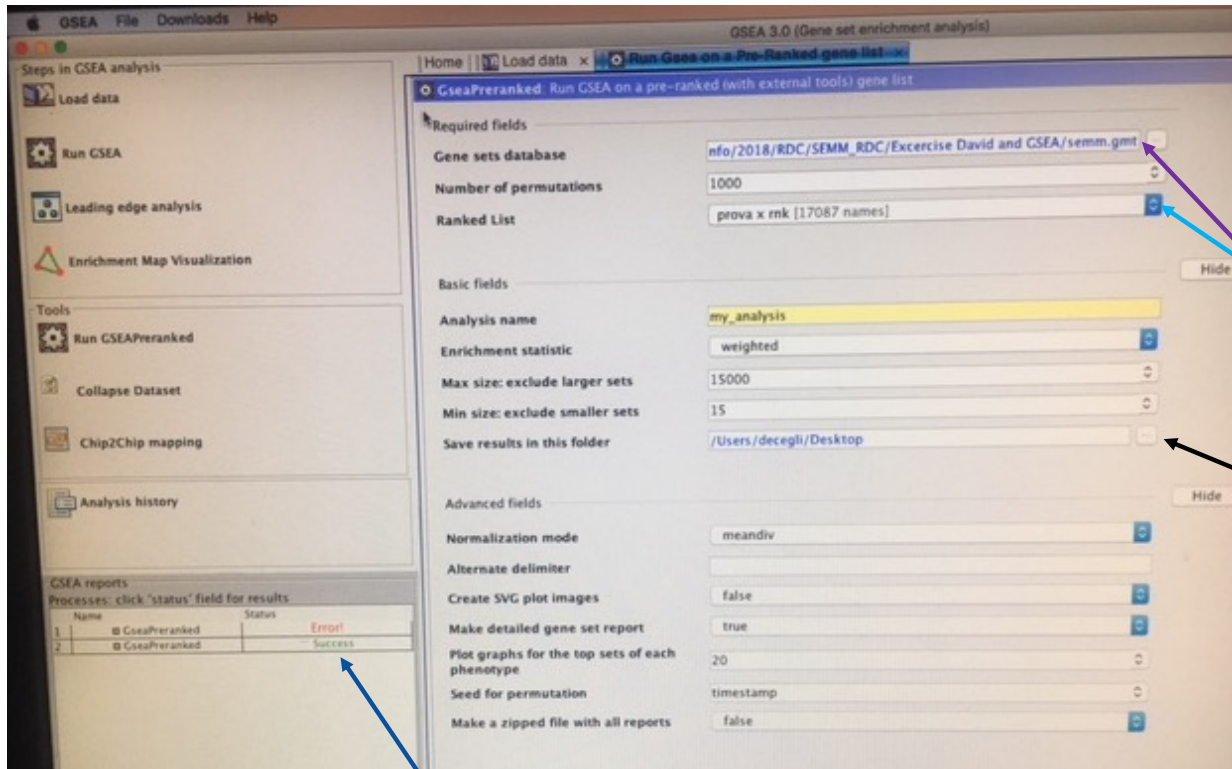| | |
|---|---|
| **H** | **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes. |
| **C1** | **positional gene sets** for each human chromosome and cytogenetic band. |
| **C2** | **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts. |
| **C3** | **regulatory target gene sets** based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites. |
| **C4** | **computational gene sets** defined by mining large collections of cancer-oriented microarray data. |
| **C5** | **GO gene sets** consist of genes annotated by the same GO terms. |
| **C6** | **oncogenic gene sets** defined directly from microarray gene expression data from cancer gene perturbations. |
| **C7** | **immunologic gene sets** defined directly from microarray gene expression data from immunologic studies. |

.rnk file
.gmt file

Where you wish to save the GSEA folder
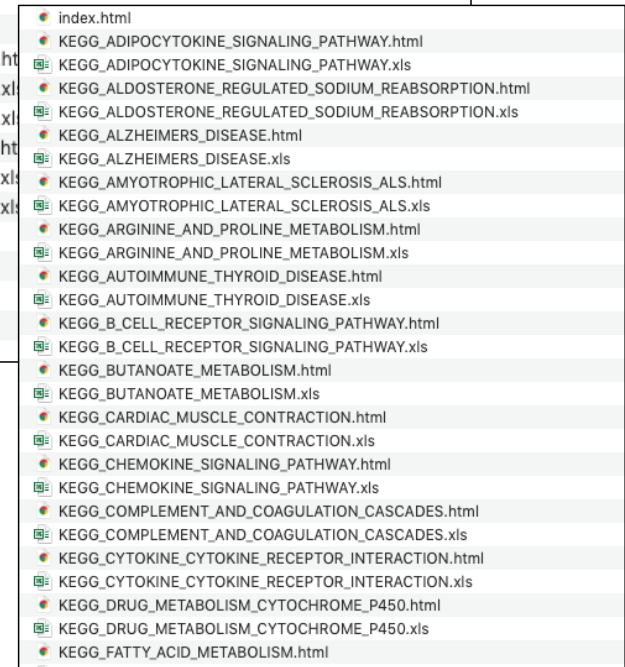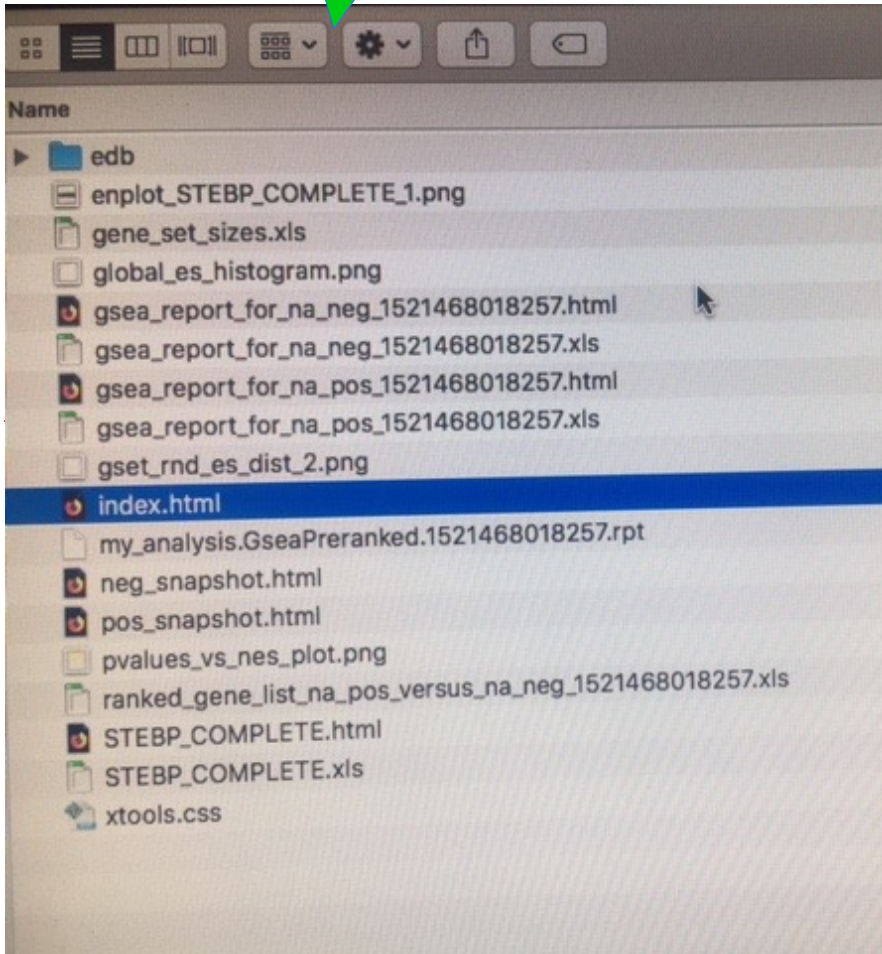
**GSEA**
**Success** or
**ERROR**

**Now ENJOY GSEA!!!!!**

RESULTS

GSEA FOLDER

# 2. GSEA

**Example of GOOD GSEA RESULT**

### Table: GSEA Results Summary

| | |
|---|---|
| Dataset | prova x rnk |
| Phenotype | NoPhenotypeAvailable |
| Upregulated in class | na_pos |
| GeneSet | SREBP_CHIP |
| Enrichment Score (ES) | 0.5765218 |
| Normalized Enrichment Score (NES) | 1.9980491 |
| Nominal p-value | 0.0 |
| FDR q-value | 0.0 |
| FWER p-Value | 0.0 |



**Fig 1: Enrichment plot: SREBP_CHIP**
*Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List*

| NAME | SIZE | ES | NES | NOM p-val | FDR q-val |
|---|---|---|---|---|---|
| GO_LYSOSOMAL_LUMEN | 69 | -0,8664851 | -1,9405853 | 0 | 0,00199854 |
| GO_PIGMENT_GRANULE | 96 | -0,8154515 | -1,8812673 | 0 | 0,00661115 |
| GO_DEVELOPMENTAL_PIGMENTATION | 39 | -0,9059212 | -1,8647331 | 0 | 0,00793784 |
| GO_VACUOLAR_LUMEN | 120 | -0,7977715 | -1,8842796 | 0 | 0,00844721 |
| GO_ANTIGEN_PROCESSING_AND_PRESENTATION_OF_PEPTIDE_ANTIGEN | 151 | -0,763359 | -1,8666888 | 0 | 0,00918049 |
| GO_ANTIGEN_PROCESSING_AND_PRESENTATION_OF_PEPTIDE_OR_POLYSACCHARIDE_ANTIGEN_VIA_MHC_CLASS_II | 81 | -0,8100114 | -1,8460969 | 0 | 0,01338742 |
| GO_VACUOLAR_PART | 430 | -0,6860185 | -1,8328513 | 0 | 0,01725247 |
| GO_AZUROPHIL_GRANULE | 112 | -0,749383 | -1,8148923 | 0 | 0,02436469 |
| GO_ANTIGEN_PROCESSING_AND_PRESENTATION | 172 | -0,741165 | -1,8196093 | 0 | 0,02458064 |
| GO_SPHINGOLIPID_METABOLIC_PROCESS | 57 | -0,8332759 | -1,7893295 | 0 | 0,03715184 |

| NAME | SIZE | ES | NES | NOM p-val | FDR q-val |
|---|---|---|---|---|---|
| GO_RESPONSE_TO_TYPE_I_INTERFERON | 68 | 0,9550428 | 1,7689795 | 0 | 0 |
| GO_REGULATION_OF_VIRAL_GENOME_REPLICATION | 74 | 0,94051254 | 1,7507216 | 0 | 0 |
| GO_VIRAL_GENOME_REPLICATION | 95 | 0,9257698 | 1,7638758 | 0 | 0 |
| GO_DEFENSE_RESPONSE_TO_VIRUS | 158 | 0,9159155 | 1,8068191 | 0 | 0 |
| GO_RESPONSE_TO_VIRUS | 216 | 0,9031988 | 1,8228018 | 0 | 0 |
| GO_DEFENSE_RESPONSE_TO_OTHER_ORGANISM | 238 | 0,88119084 | 1,7750893 | 0 | 0 |
| GO_SM_LIKE_PROTEIN_FAMILY_COMPLEX | 83 | 0,92373174 | 1,7440909 | 0 | 4,32E-04 |
| GO_REGULATION_OF_VIRAL_LIFE_CYCLE | 106 | 0,8948939 | 1,7309905 | 0 | 0,00202313 |
| GO_NEGATIVE_REGULATION_OF_VIRAL_GENOME_REPLICATION | 42 | 0,96347237 | 1,7121593 | 0 | 0,00567661 |
| GO_SPLICEOSOMAL_SNRNP_ASSEMBLY | 39 | 0,9559272 | 1,715771 | 0 | 0,0059322 |
| GO_NEGATIVE_REGULATION_OF_VIRAL_PROCESS | 67 | 0,9221198 | 1,7126523 | 0 | 0,0061437 |
| GO_REGULATION_OF_SYMBIOSIS_ENCOMPASSING_MUTUALISM_THROUGH_PARASITISM | 169 | 0,85546035 | 1,7057966 | 0 | 0,00872576 |
| GO_SNRNA_BINDING | 39 | 0,9428287 | 1,7008914 | 0 | 0,0106801 |
| GO_CYTOKINE_MEDIATED_SIGNALING_PATHWAY | 476 | 0,80943793 | 1,6872447 | 0 | 0,02351773 |
| GO_REGULATION_OF_T_CELL_MEDIATED_IMMUNITY | 34 | 0,95067674 | 1,6772937 | 0,00166667 | 0,0383134 |
| GO_RECEPTOR_SIGNALING_PATHWAY_VIA_STAT | 74 | 0,89211684 | 1,6725465 | 0 | 0,04628721 |
| GO_DOUBLE_STRANDED_RNA_BINDING | 66 | 0,89978564 | 1,653459 | 0 | 0,08593836 |
| GO_REGULATION_OF_RECEPTOR_SIGNALING_PATHWAY_VIA_STAT | 59 | 0,90091246 | 1,6586115 | 0,00154083 | 0,08643796 |

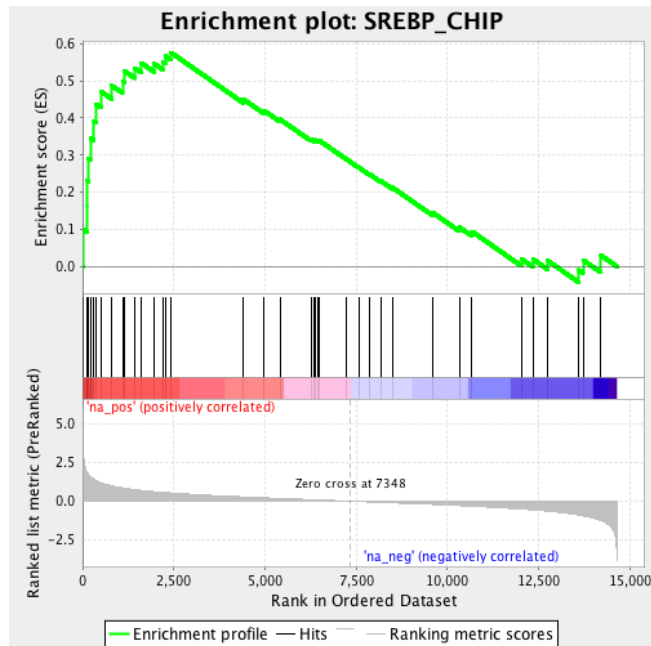**For Students: please refer to the file "GSEA_how create an input_example.xlsx", second sheet**

❖ The threshold for statistical significance of GSEA is **FDR<0,25** and **Enrichment Score >0.5** for induced gene sets and Enrichment Score **<-0.5** for inhibited GS
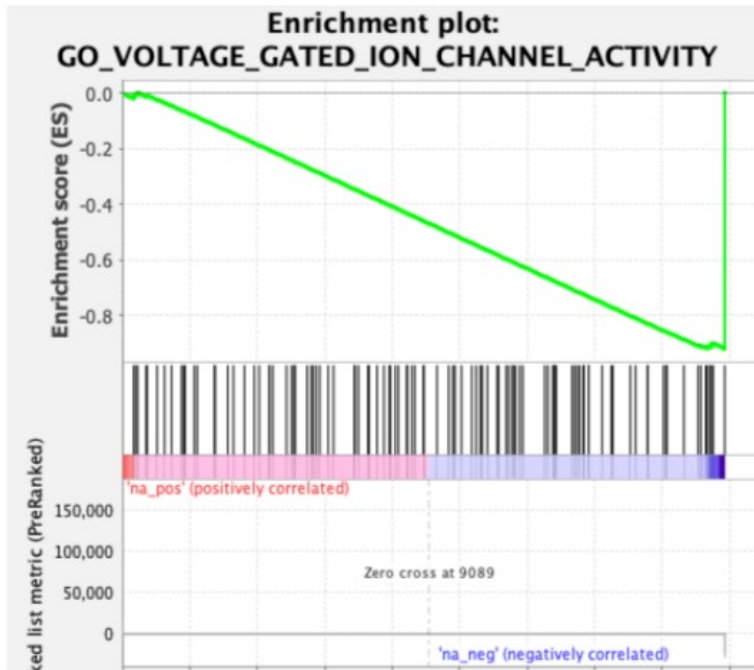
# 2. GSEA

**Table: GSEA Results Summary**

| Dataset | exp2Exp1 |
|---|---|
| Phenotype | NoPhenotypeAvailable |
| Upregulated in class | na_neg |
| GeneSet | GO_VOLTAGE_GATED_ION_CHANNEL_ACTIVITY |
| Enrichment Score (ES) | -0.92081845 |
| Normalized Enrichment Score (NES) | -1.4579594 |
| Nominal p-value | 0.1031941 |
| FDR q-value | 0.55201936 |
| FWER p-Value | 1.0 |

**Example of NOT significant RESULT**

**Enrichment plot:**
**GO_VOLTAGE_GATED_ION_CHANNEL_ACTIVITY**

Zero cross at 9089

'na_pos' (positively correlated)

'na_neg' (negatively correlated)

❖ The threshold for statistical significance of GSEA is **FDR<0,25** and **Enrichment Score >0.5** for induced gene sets and Enrichment Score **<-0.5** for inhibited GS