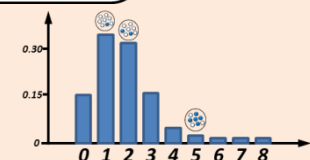
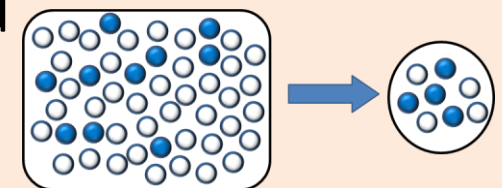
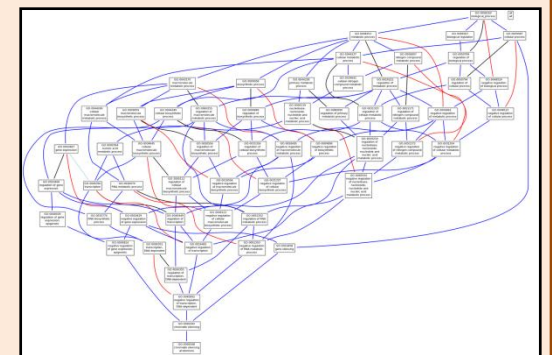


# Gene Set Enrichment Analysis

Genome 373  
Genomic Informatics  
Elhanan Borenstein

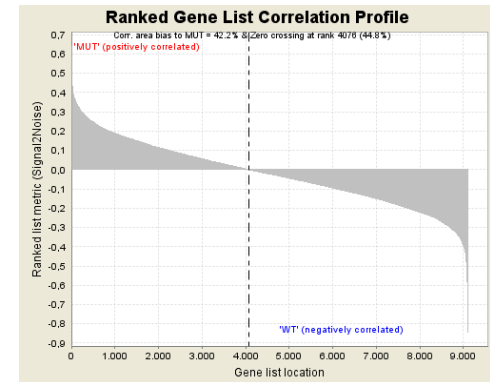
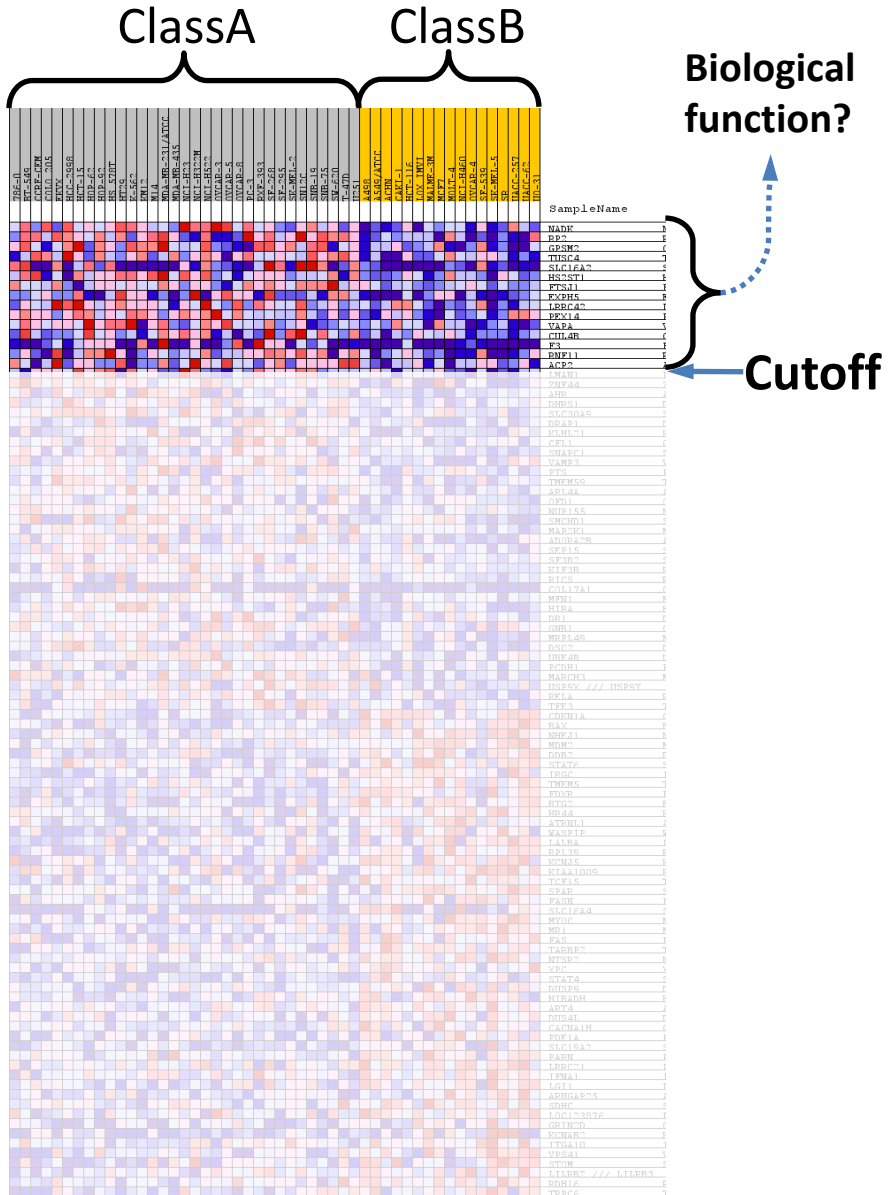
# A quick review

- Gene expression profiling
  - Which molecular processes/functions are involved in a certain phenotype (e.g., disease, stress response, etc.)
- The Gene Ontology (GO) Project
  - Provides shared vocabulary/annotation
  - Terms are linked in a complex structure
- Enrichment analysis:
  - Find the “most” differentially expressed genes
  - Identify **over-represented** annotations
  - Modified Fisher's exact test



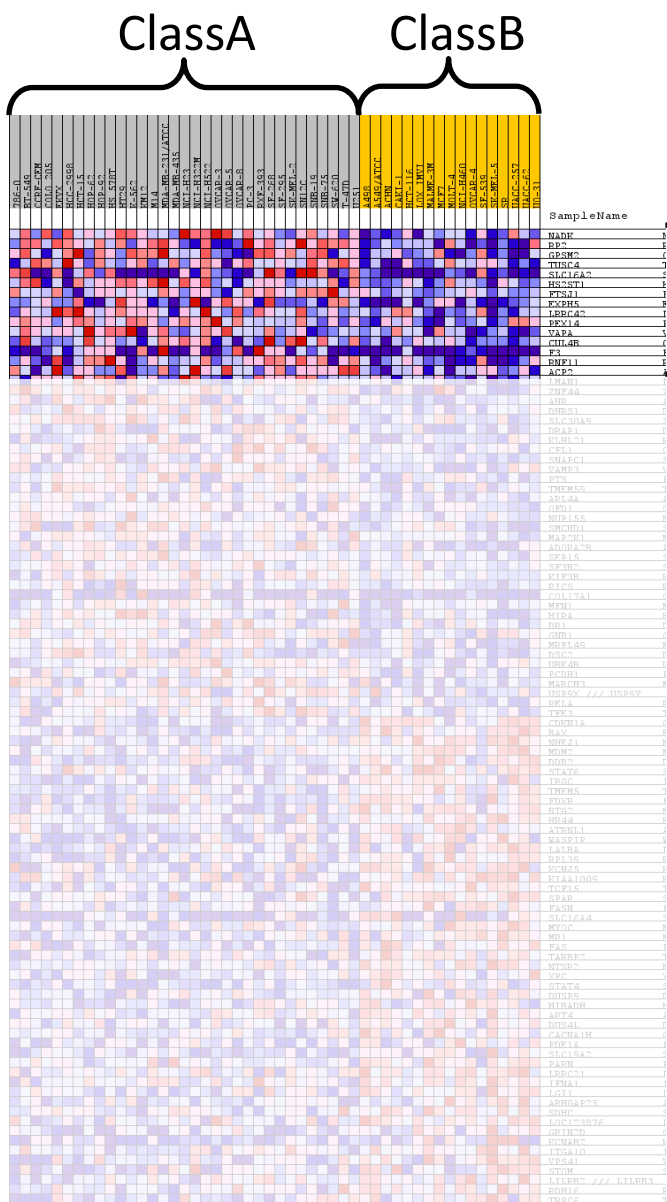
# Enrichment Analysis

Genes ranked by expression correlation to Class A

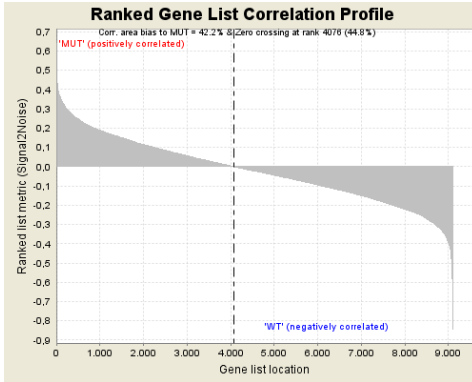
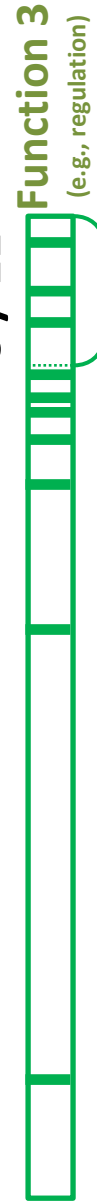
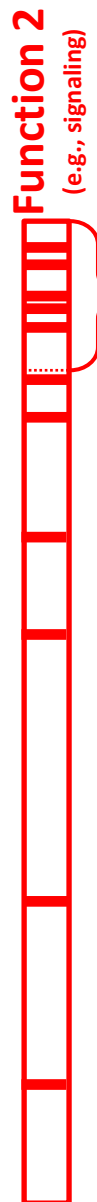


# Enrichment Analysis

Genes ranked by expression correlation to Class A

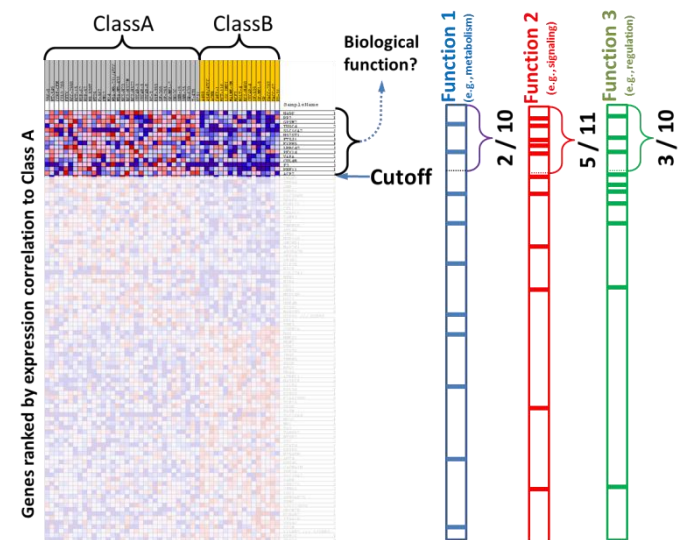


Biological function?  
Cutoff



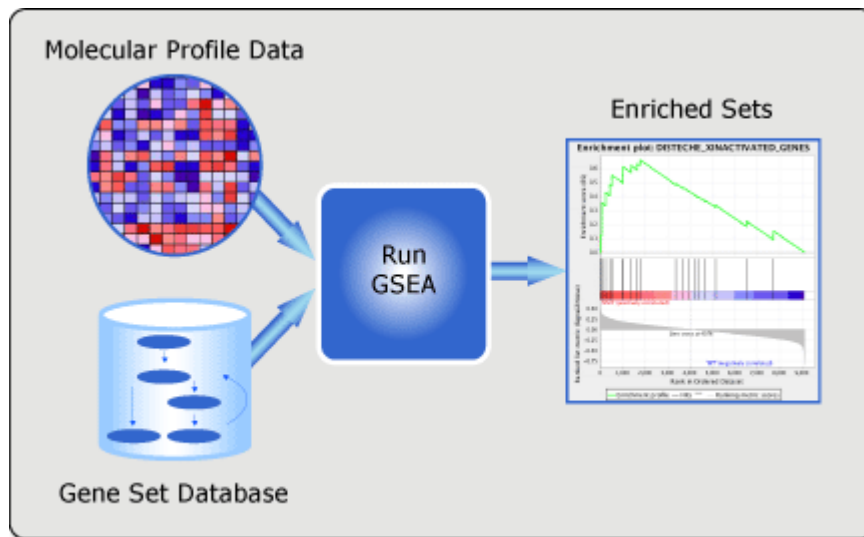
# Problems with cutoff-based analysis

- After correcting for multiple hypotheses testing, no individual gene may meet the threshold due to noise.
- Alternatively, one may be left with a long list of significant genes without any unifying biological theme.
- The cutoff value is often arbitrary!
- **We are really examining only a handful of genes, totally ignoring much of the data**



# Gene Set Enrichment Analysis

- MIT, Broad Institute
- V 2.0 available since Jan 2007



## Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles

Aravind Subramanian<sup>1,2</sup>, Pablo Tamayo<sup>1,3</sup>, Vamsi K. Mootha<sup>1,4</sup>, Sayan Mukherjee<sup>1</sup>, Benjamin L. Ebert<sup>1</sup>, Michael A. Gillette<sup>1</sup>, Amanda Paulovich<sup>1</sup>, Scott L. Pomeroy<sup>1</sup>, Todd R. Golub<sup>1,5</sup>, Eric S. Lander<sup>1,6,7,8,9,10</sup>, and Jill P. Mesirov<sup>1,4</sup>

<sup>1</sup>Broad Institute of Massachusetts Institute of Technology and Harvard, 320 Charles Street, Cambridge, MA 02141; <sup>2</sup>Department of Systems Biology, Albert Einstein College of Medicine, 1400 Pelham Avenue, Bronx, NY 10461; <sup>3</sup>Institute for Genome Sciences and Policy, Center for Interdisciplinary Engineering, Medicine, and Applied Sciences, Duke University, 151 Science Drive, Durham, NC 27708; <sup>4</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, 44 Binney Street, Boston, MA 02115; <sup>5</sup>Division of Pulmonary and Critical Care Medicine, Massachusetts General Hospital, 55 Fruit Street, Boston, MA 02114; <sup>6</sup>Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, C2-023, P.O. Box 19024, Seattle, WA 98109-1024; <sup>7</sup>Department of Neurology, Emory University School of Medicine, 1365 Clifton Road, NE, Atlanta, GA 30309; <sup>8</sup>Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139; <sup>9</sup>Department of Genetics, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, MA 02115; <sup>10</sup>Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139; and <sup>11</sup>Whitehead Institute for Biomedical Research, Massachusetts Institute of Technology, Cambridge, MA 02142

Contributed by Eric S. Lander, August 2, 2005

Although genome-wide RNA expression analysis has become a routine tool in biomedical research, extracting biological insight from such information remains a major challenge. Here, we describe a powerful analytical method called Gene Set Enrichment Analysis (GSEA) for interpreting gene expression data. The method derives its power by focusing on gene sets, that is, groups of genes that share common biological function, chromosomal location, or regulation. We demonstrate how GSEA yields insights into several cancer-related data sets, including leukemia and lung cancer. Notably, where single-gene analysis finds little similarity between two independent studies of patient survival in lung cancer, GSEA reveals many biological pathways in common. The GSEA method is embodied in a freely available software package, together with an initial database of 1,325 biologically defined gene sets.

microarray

Genome-wide expression analysis with DNA microarrays has become a mainstay of genomics research (1, 2). The challenge no longer lies in obtaining gene expression profiles, but rather in interpreting the results to gain insights into biological mechanisms.

In a typical experiment, mRNA expression profiles are generated for thousands of genes from a collection of samples belonging to one of two classes, for example, tumors that are sensitive vs. resistant to a drug. The genes can be ordered in a ranked list  $L$ , according to their differential expression between the classes. The challenge is to extract meaning from this list.

A common approach involves focusing on a handful of genes at the top and bottom of  $L$  (i.e., those showing the largest difference) to discern telltale biological clues. This approach has a few major limitations.

(i) After correcting for multiple hypotheses testing, no individual gene may meet the threshold for statistical significance, because the relevant biological differences are modest relative to the noise inherent to the microarray technology.

(ii) Alternatively, one may be left with a long list of statistically significant genes without any unifying biological theme. Interpretation can be daunting and ad hoc, being dependent on a biologist's area of expertise.

(iii) Single-gene analysis may miss important effects on pathways. Cellular processes often affect sets of genes acting in concert. An increase of 20% in all genes encoding members of a metabolic pathway may dramatically alter the flux through the pathway and may be more important than a 20-fold increase in a single gene.

(iv) When different groups study the same biological system, the list of statistically significant genes from the two studies may show distressingly little overlap (3).

To overcome these analytical challenges, we recently developed a method called Gene Set Enrichment Analysis (GSEA) that

evaluates microarray data at the level of gene sets. The gene sets are defined based on prior biological knowledge, e.g., published information about biochemical pathways or coclustering in previous experiments. The goal of GSEA is to determine whether members of a gene set  $S$  tend to occur toward the top (or bottom) of the list  $L$ , in which case the gene set is correlated with the phenotypic class distinction.

We used a preliminary version of GSEA to analyze data from muscle biopsies from diabetics vs. healthy controls (4). The method revealed that genes involved in oxidative phosphorylation show reduced expression in diabetics, although the average decrease per gene is only 20%. The results from this study have been independently validated by other microarray studies (5) and by *in vivo* functional studies (6).

Given this success, we have developed GSEA into a robust technique for analyzing molecular profiling data. We studied its characteristics and performance and substantially revised and generalized the original method for broader applicability.

In this paper, we provide a full mathematical description of the GSEA methodology and illustrate its utility by applying it to several diverse biological problems. We have also created a software package, called GSEA-P and an initial inventory of gene sets (Molecular Signature Database, MSigDB), both of which are freely available.

### Methods

**Overview of GSEA.** GSEA considers experiments with genome-wide expression profiles from samples belonging to two classes, labeled 1 or 2. Genes are ranked based on the correlation between their expression and the class distinction by using any suitable metric (Fig. 1A).

Given an *a priori* defined set of genes  $S$  (e.g., genes encoding products in a metabolic pathway, located in the same cytogenetic band, or sharing the same GO category), the goal of GSEA is to determine whether the members of  $S$  are randomly distributed throughout  $L$  or primarily found at the top or bottom. We expect

Freely available online through the PNAS open access option.

Abbreviations: ALL, acute lymphoid leukemia; AML, acute myeloid leukemia; ES, enrichment score; FDR, false discovery rate; GSEA, Gene Set Enrichment Analysis; MAPK, mitogen-activated protein kinase; MSigDB, Molecular Signature Database; NES, normalized enrichment score.

See Commentary on page 15278.

\*A.S. and P.T. contributed equally to this work.

<sup>†</sup>to whom correspondence may be addressed. E-mail: lander@broad.mit.edu or mesirov@broad.mit.edu.

© 2005 by The National Academy of Sciences of the USA.

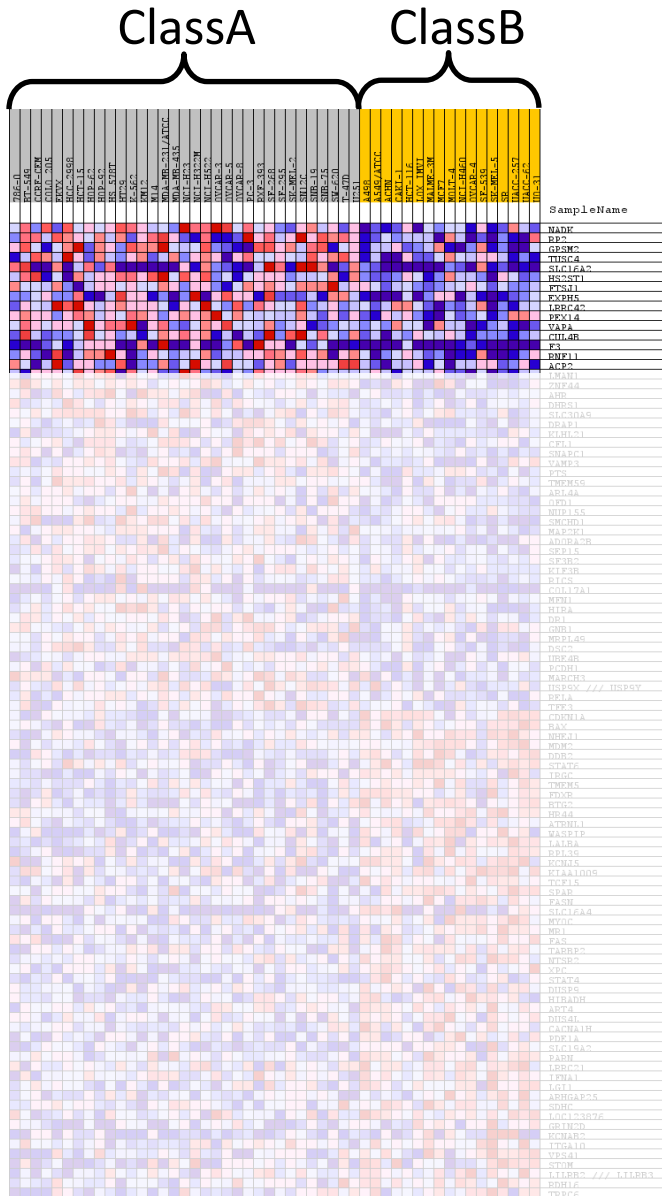
(Subramanian et al. PNAS. 2005.)

# GSEA key features

- Does not require setting a cutoff!
- Identifies the set of relevant genes as part of the analysis!
- Calculates a score for the enrichment of a **entire set of genes** rather than single genes!
- Provides a more robust statistical framework!

# Gene Set Enrichment Analysis

Genes ranked by expression correlation to Class A



Biological function?

Cutoff

Function 1  
(e.g., metabolism)

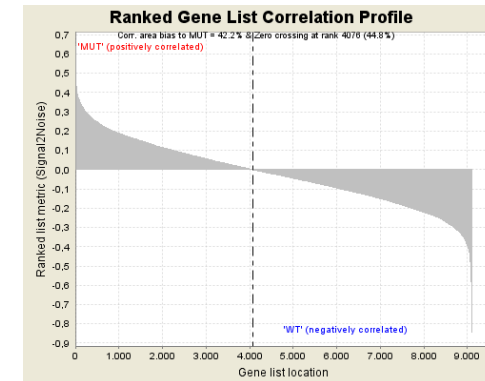
2 / 10

Function 2  
(e.g., signaling)

5 / 11

Function 3  
(e.g., regulation)

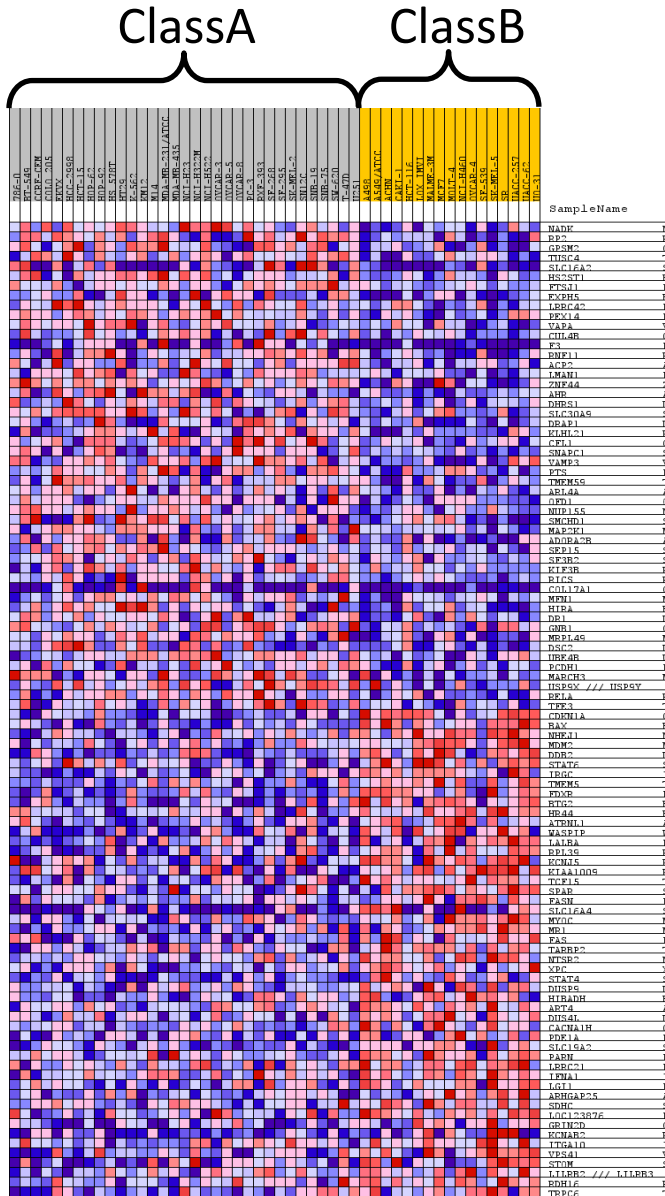
3 / 10





# Gene Set Enrichment Analysis

Genes ranked by expression correlation to Class A



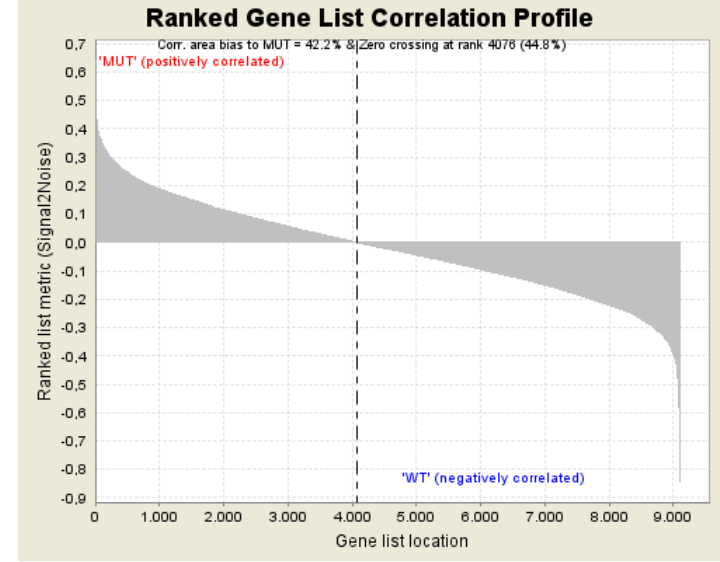
Function 1  
(e.g., metabolism)



Function 2  
(e.g., signaling)

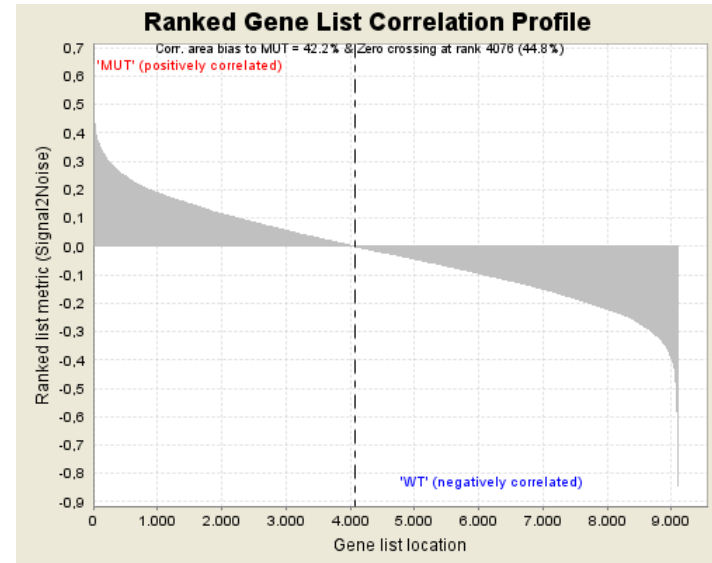
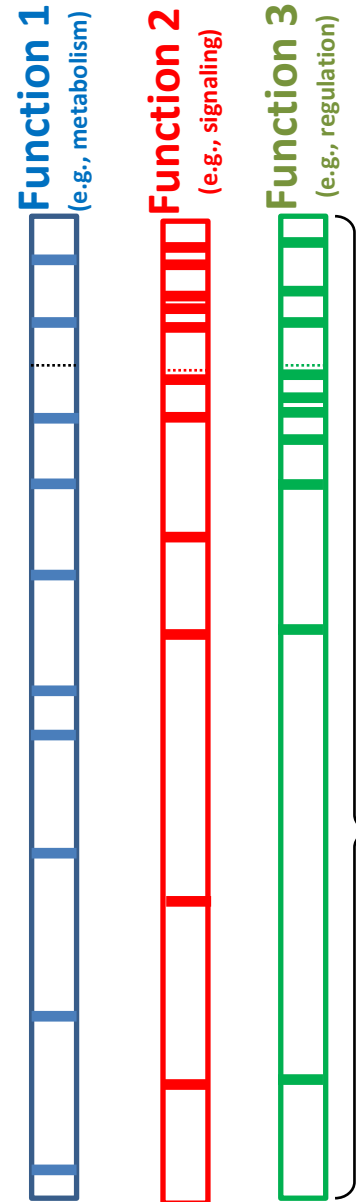
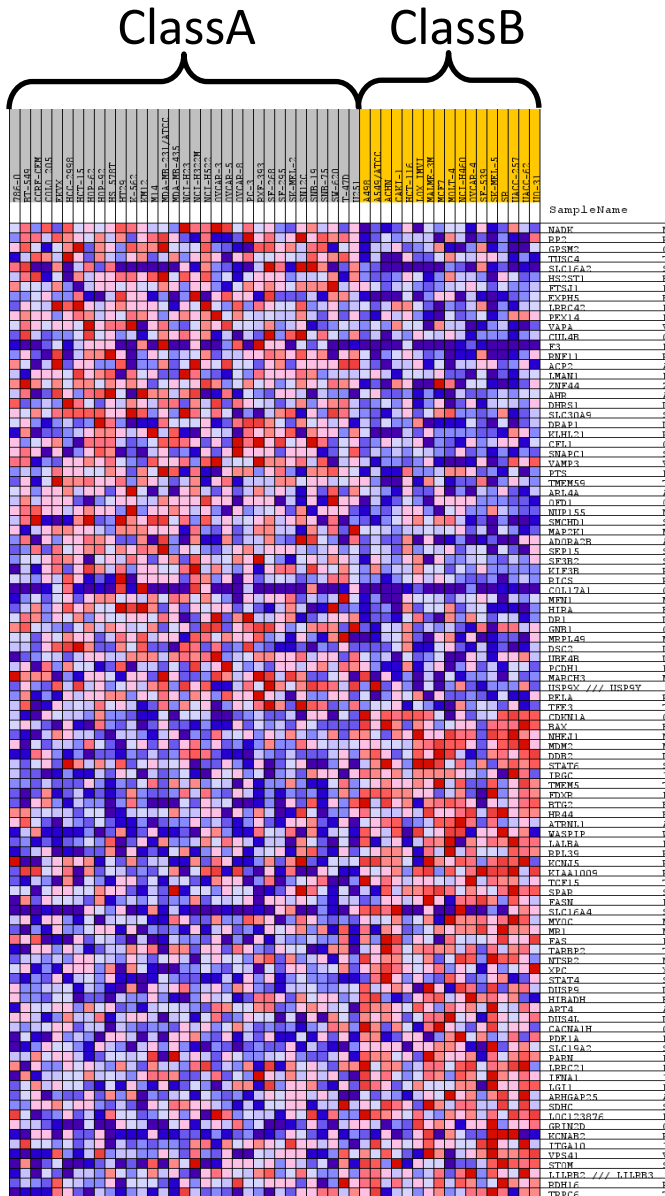


Function 3  
(e.g., regulation)



# Gene Set Enrichment Analysis

Genes ranked by expression correlation to Class A



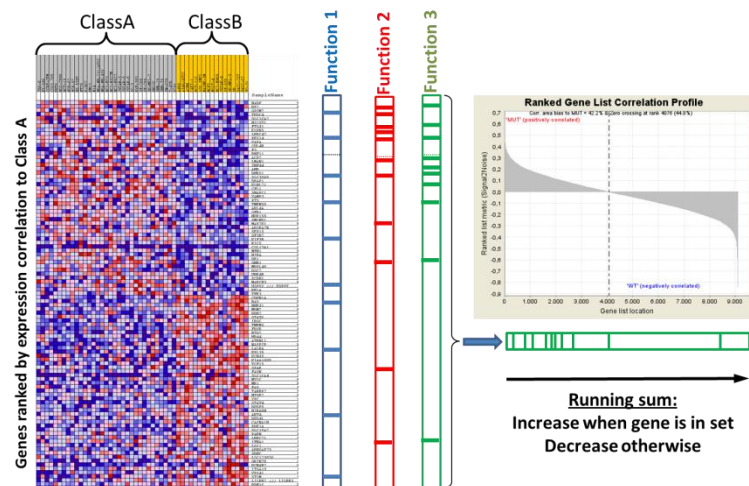
**Running sum:**  
 Increase when gene annotated with the function under study  
 Decrease otherwise

# Gene Set Enrichment Analysis

What would you expect if ALL genes annotated with this function cluster at the top of the list?

What would you expect if genes annotated with this function are randomly distributed?

What would you expect if most of the genes annotated with this function cluster at the top of the list?

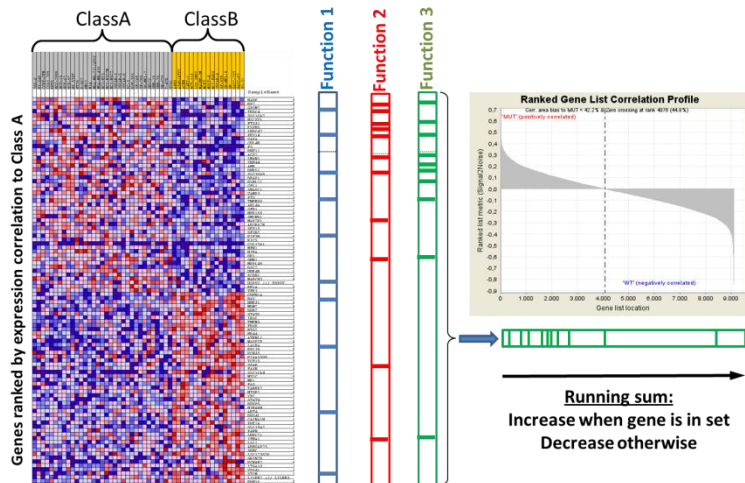
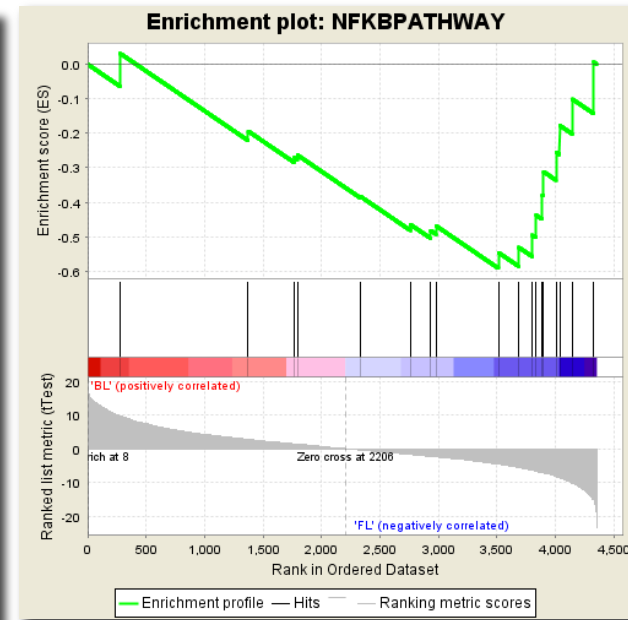
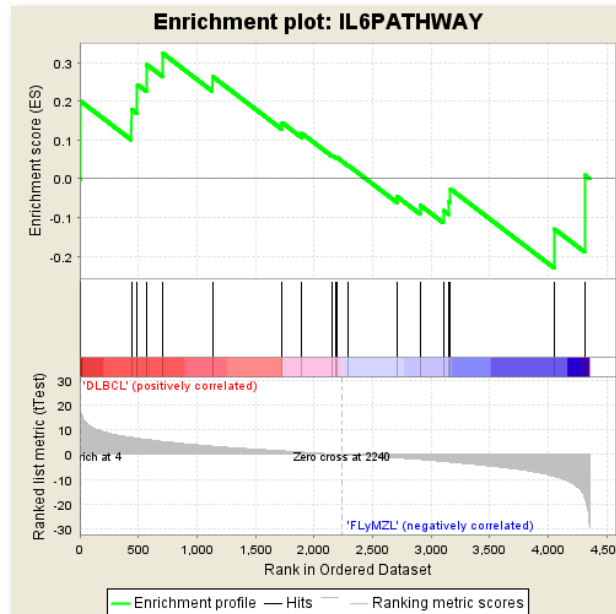
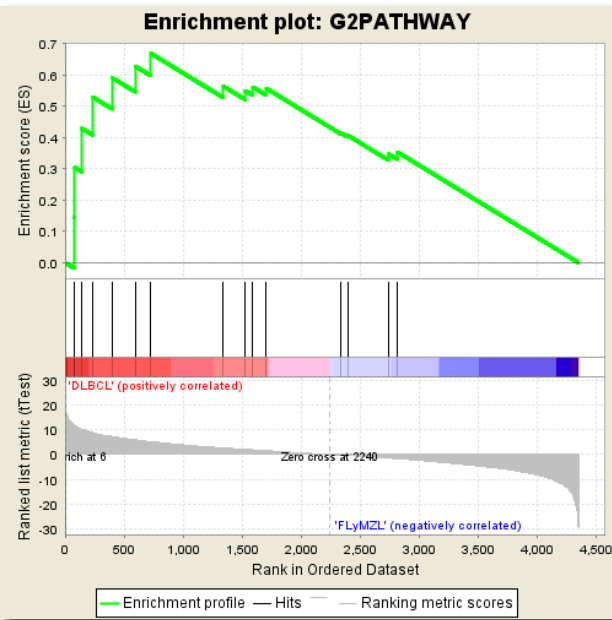


# Gene Set Enrichment Analysis

ES = 0.69

Low ES (evenly distributed)

ES = -0.59

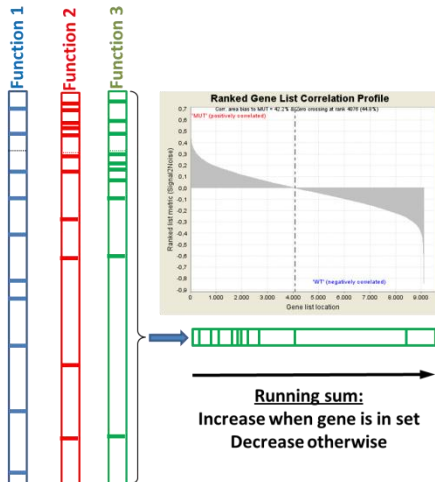
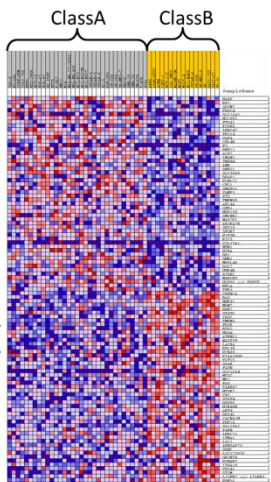
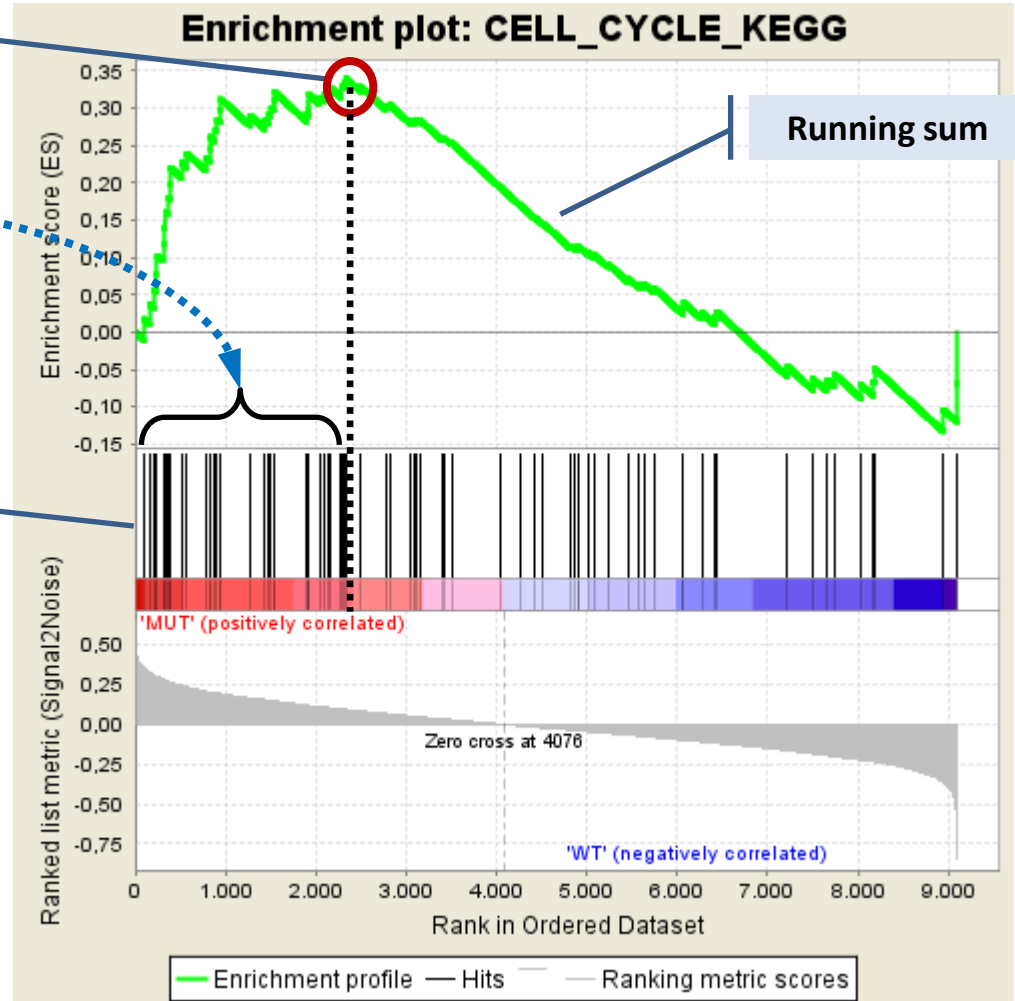


# Gene Set Enrichment Analysis

Enrichment score (ES) =  
max deviation from 0

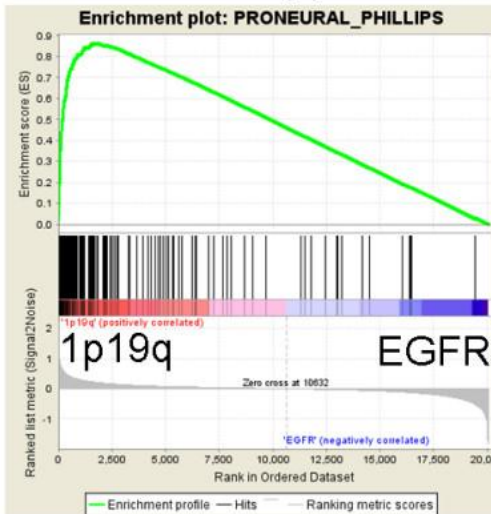
Leading  
Edge genes

Genes within  
functional set  
(hits)



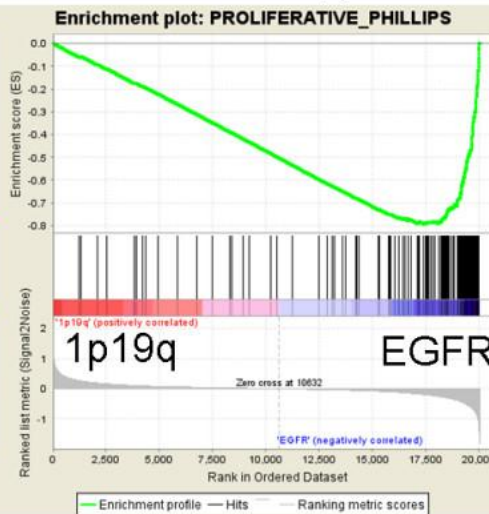
# Gene Set Enrichment Analysis

**A** ES=0.86,  $p<0.001$



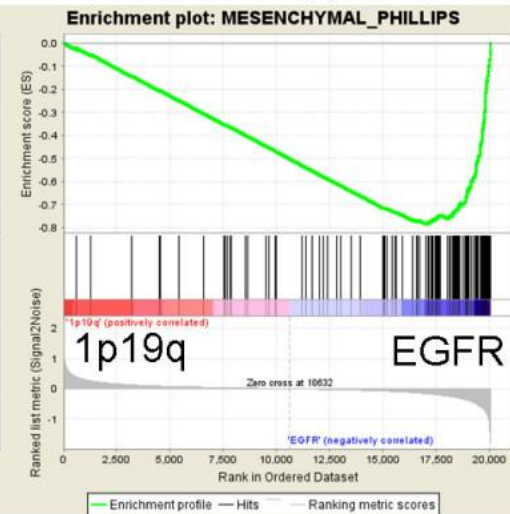
SNR

**B** ES= -0.79,  $p<0.001$



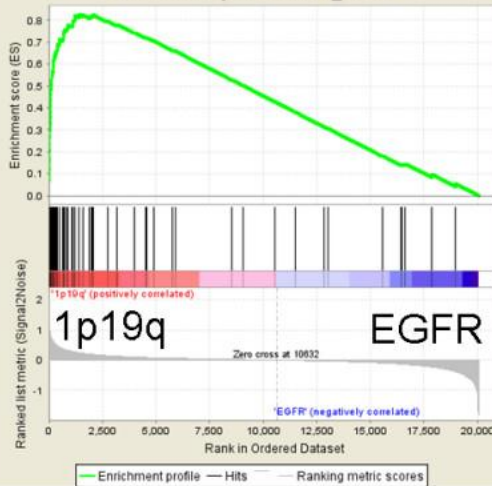
SNR

**C** ES= -0.78,  $p<0.001$



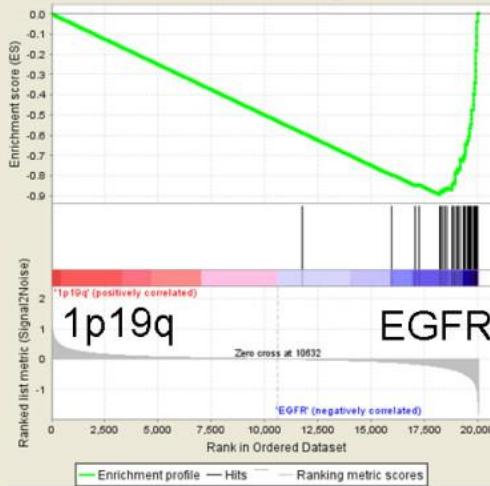
SNR

Enrichment plot: FREIJE\_HC1A



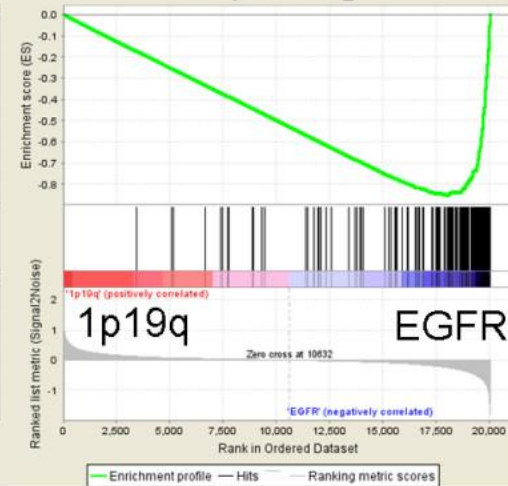
**D** ES=0.82,  $p<0.001$

Enrichment plot: FREIJE\_HC2A



**E** ES= -0.89,  $p<0.001$

Enrichment plot: FREIJE\_HC2B

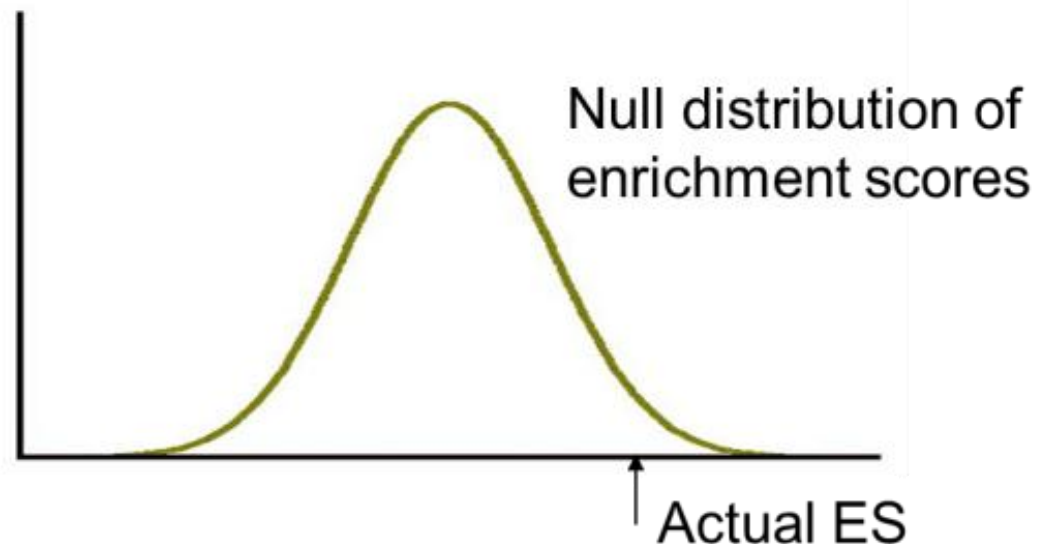


**F** ES= -0.85,  $p<0.001$

# Estimating Significance of ES

# Estimating Significance of ES

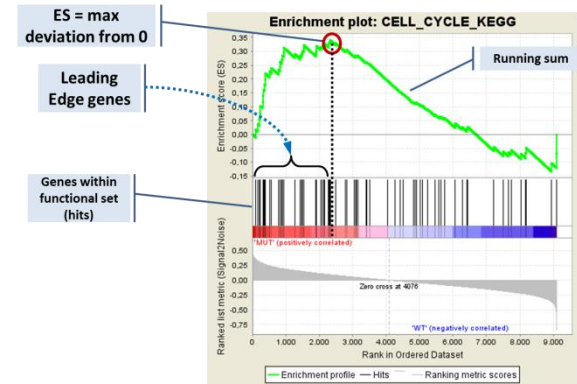
- An empirical permutation test
- Phenotype labels are shuffled and the ES for this functional set is recomputed. Repeat 1000 times.
- Generating a null distribution





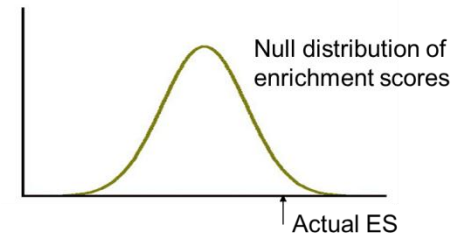
# GSEA Steps

1. Calculation of an enrichment score (ES) for each functional category



2. Estimation of significance level of the ES

- Shuffling-based null distribution



3. Adjustment for multiple hypotheses testing

- Necessary if comparing multiple gene sets (i.e., functions)
- Computes FDR (false discovery rate)

