# Gene Set Enrichment Analysis

Fernanda Díaz Espinosa

DNA MICROARRAY

RNA-SEQ

cDNA sample 1
cDNA sample 2
Fluorescent tag

cDNA sample 1
cDNA sample 2

Reference genome

Gene 1 | Gene 2 | Gene 3 | Gene 4

Sample 2

Sample 1

# What do we do with that MUCH information?

Fastq

```
@61DFRAAXX100204:1:100:10494:3070/1
AAACAACAGGGCACATTGTCACTCTTGTATTTGAAAAACACTTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCCC@@CACCCCCA
```
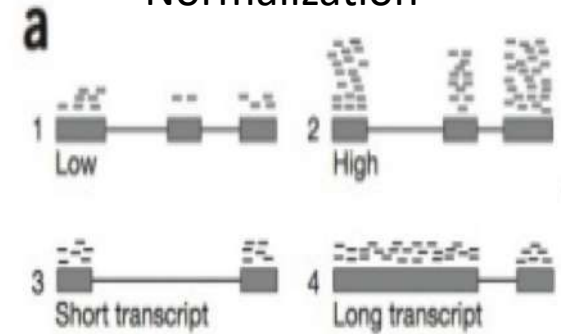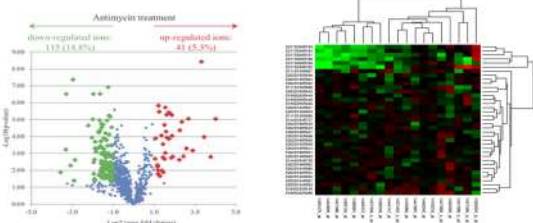
Trimming and QC

Reading Alignment

Normalization

SAM/BAM

Reference genome

Transcript assembly

Reference genome

## D.E table

| GENE | FC | p-value | FDR |
|------|-----|---------|---------|
| CTCF | -2 | 0.006 | 0.00013 |
| BORIS | +2 | 0.0009 | 0.0003 |

## Count

| GENE | A | B | C |
|------|-----|-----|-----|
| CTCF | 13 | 1 | 7 |
| BORIS | 0 | 18 | 5 |

Data display

Differential Expression

What do we do with that MUCH information?

# Pathway analysis

# Pathway analysis

And the Fold Change?

First described by Mootha et al., 2003.
- Analysis of second generation pathways.
- Significance analysis of function and expression (SAFE).

- Four step approach.

A) Calculate the local statistic (at the gene level).
B) Calculate the global statistic (at the level of pathway or set of genes).
C) Determine a significance value.
D) Adjust for multiple comparisons.

**A** Phenotype Classes — A, B

**B** Gene set S

Leading edge subset — Gene set S

Correlation with Phenotype

Random Walk

ES(S)

Maximum deviation from zero provides the enrichment score ES(S)

Gene List Rank

Ranked Gene List

# Files to start

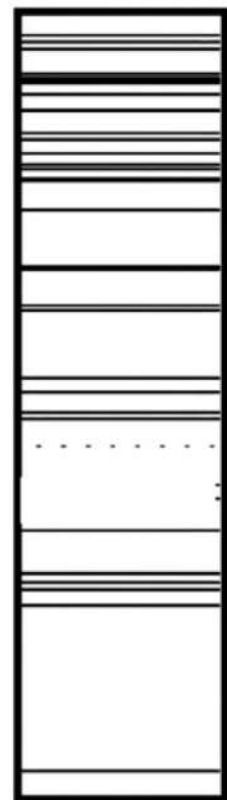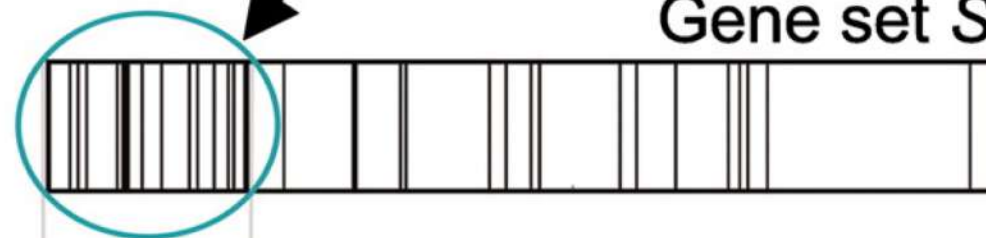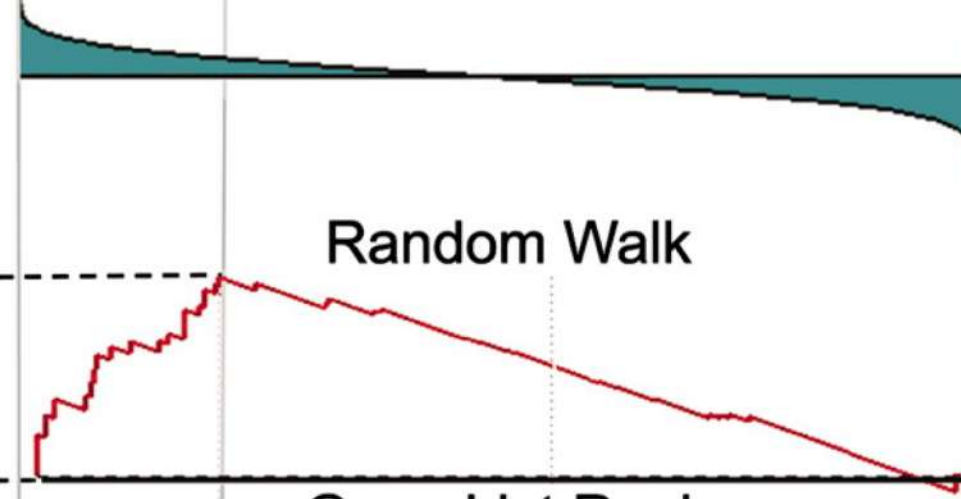| Data File | Content | Format | Source |
|---|---|---|---|
| Expression dataset | Contains features (genes or probes), samples, and an expression value for each feature in each sample. Expression data can come from any source (Affymetrix, Stanford cDNA, and so on). | res, gct, pcl, or txt | You create the file. |
| Phenotype labels | Contains phenotype labels and associates each sample with a phenotype. | cls | You create the file or have GSEA create it for you. |
| Gene sets | Contains one or more gene sets. For each gene set, gives the gene set name and list of features (genes or probes) in that gene set. | gmx or gmt | You use the files on the Broad ftp site, export gene sets from the Molecular Signature Database (MSigDb) or create your own gene sets file. |
| Chip annotations | Lists each probe on a DNA chip and its matching HUGO gene symbol. Optional for the gene set enrichment analysis. | Chip | You use the files on the Broad ftp site, download the files from the GSEA web site, or create your own chip file. |

# Choose your database

- Gene set Databases.

    a) Gene Ontology (GO).

    b) Molecular Signatures Database (MSigDB).


- Detailed biochemical pathway databases.

    a) Reactome.

    b) Panther.

    c) NetPath.

    d) HumanCyc.

    e) National Cancer Institute (NCI) Pathway Interaction Database (PID).

    f) KEGG.

# 1.Local statistic



The rank file is a list of detected genes and a rank metric score.

- At the top of the list are genes with the "strongest" up-regulation.

- At the bottom of the list are the genes with the "strongest" down-regulation.

- The genes not changing are in the middle.

**Metric**

- Signal-to-noise ratio (S2N; the default measure in GSEA)
- Absolute value of signal-to-noise ratio (|S2N|)
- Difference of expression means between classes (Ratio)
- $Log_2$ of Ratio ($log_2$(Ratio))
- T-test statistic (T-test).

https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1674-0

$$\frac{\mu_A - \mu_B}{\sigma_A + \sigma_B}$$

- **Signal2Noise** (default) uses the difference of means scaled by the standard deviation.

where μ is the mean and σ is the standard deviation; σ has a minimum value of .2 * absolute(μ), where μ=0 is adjusted to μ=1. The larger the signal-to-noise ratio, the larger the differences of the means (scaled by standard deviations);

$$\frac{\mu_A - \mu_B}{\sqrt{\dfrac{\sigma_A^2}{n_A} + \dfrac{\sigma_B^2}{n_B}}}$$

- **tTest** uses the difference of means scaled by the standard deviation and number of samples.

where μ is the mean, n is the number of samples, and σ is the standard deviation; σ has a minimum value of .2 * absolute(μ), where μ=0 is adjusted to μ=1. The larger the tTest ratio, the more distinct the gene expression is in each phenotype and the more the gene acts as a "class marker."

$$\frac{\mu_A}{\mu_B}$$

- **Ratio_of_Classes** (also referred to as fold change) uses the ratio of class means to calculate fold change for natural scale data:

where μ is the mean. The larger the fold change, the more distinct the gene expression is in each phenotype and the more the gene acts as a "class marker."

$$\mu_A - \mu_B$$

- **Diff_of_Classes** uses the difference of class means to calculate fold change for log scale data:
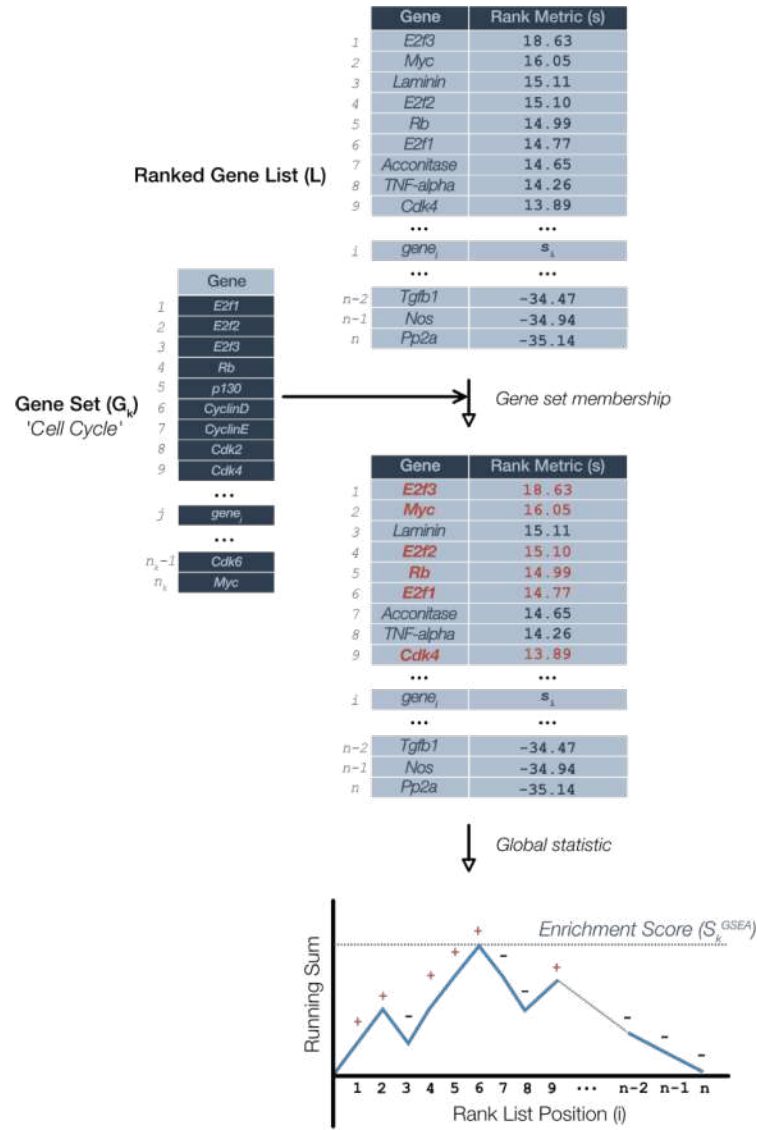
where μ is the mean. The larger the fold change, the more distinct the gene expression is in each phenotype and the more the gene acts as a "class marker."

$$\log 2\left(\frac{\mu_A}{\mu_B}\right)$$

- **log2_Ratio_of_Classes** uses the log2 ratio of class means to calculate fold change for natural scale data:

where μ is the mean. This is the recommended statistic for calculating fold change for natural scale data.

# 2. Global statistic



## *Enrichment Score (ES)*

Reflects the degree to which a gene set is overrepresented at the top or bottom of a ranked list of genes.

GSEA calculates the ES by walking down the ranked list of genes, increasing a running-sum statistic when a gene is in the gene set and decreasing it when it is not.
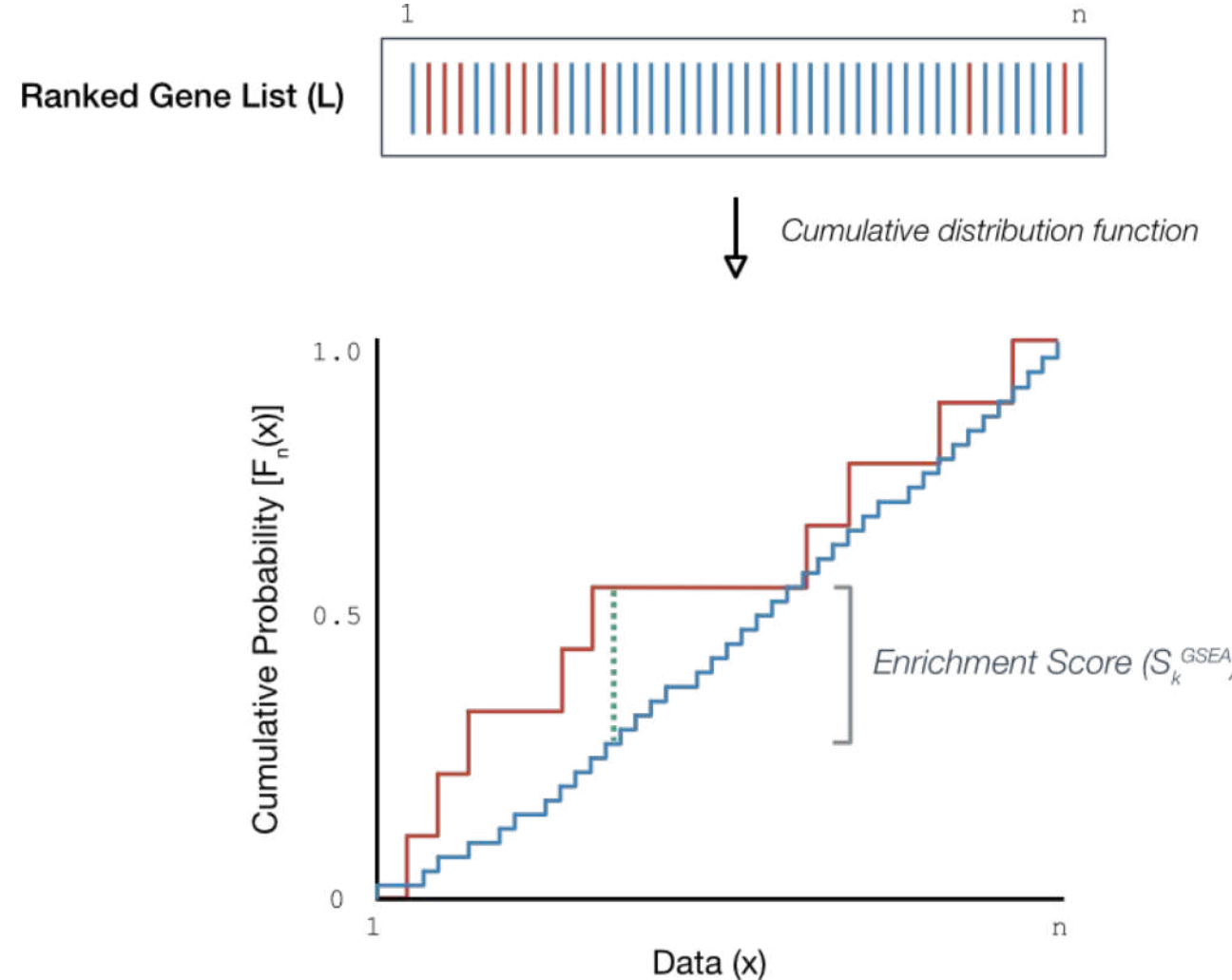
➢ The magnitude of the increment depends on the correlation of the gene with the phenotype.

➢ The ES is the maximum deviation from zero encountered in walking the list.

➢ A positive ES indicates gene set enrichment at the top of the ranked list; a negative

➢ ES indicates gene set enrichment at the bottom of the ranked list.

# Kolmogorov–Smirnov test

The Kolmogorov–Smirnov statistic quantifies
a distance between the empirical distribution function of
the sample and the cumulative distribution function of
the reference distribution.

 …Or between the empirical distribution functions of two
samples.

The null distribution of this statistic is calculated under
the null hypothesis that the sample is drawn from the
reference distribution.
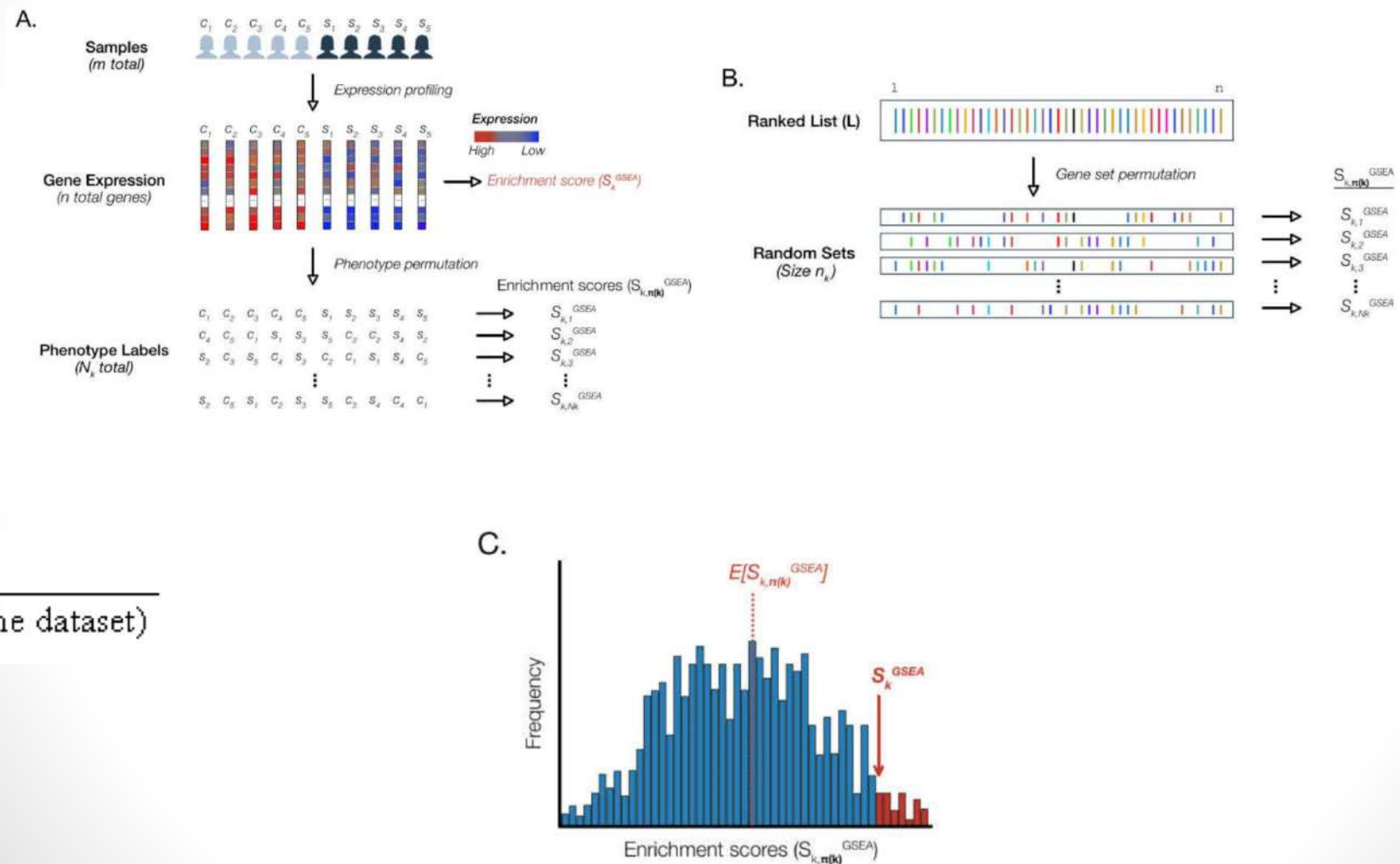
# 3. "Significance" value

## *Normalized Enrichment Score (NES)*

The normalized enrichment score (NES) is the primary statistic for examining gene set enrichment results.

By normalizing the enrichment score, GSEA accounts for differences in gene set size and in correlations between gene sets and the expression dataset.

Therefore, the normalized enrichment scores (NES) can be used to compare analysis results across gene sets.

$$NES = \frac{actual\ ES}{mean(ESs\ against\ all\ permutations\ of\ the\ dataset)}$$

# 4. Adjust for multiple comparisons.

## *False Discovery Rate (FDR)*

The false discovery rate (FDR) is the estimated probability that a gene set with a given NES represents a false positive finding.

For example, an FDR of 25% indicates that the result is likely to be valid 3 out of 4 times.

When we test a family of hypotheses, the chance of observing a statistic with a small p-value increases.
When smaller than the significance level, they can be erroneously classified as discoveries or Type I errors.
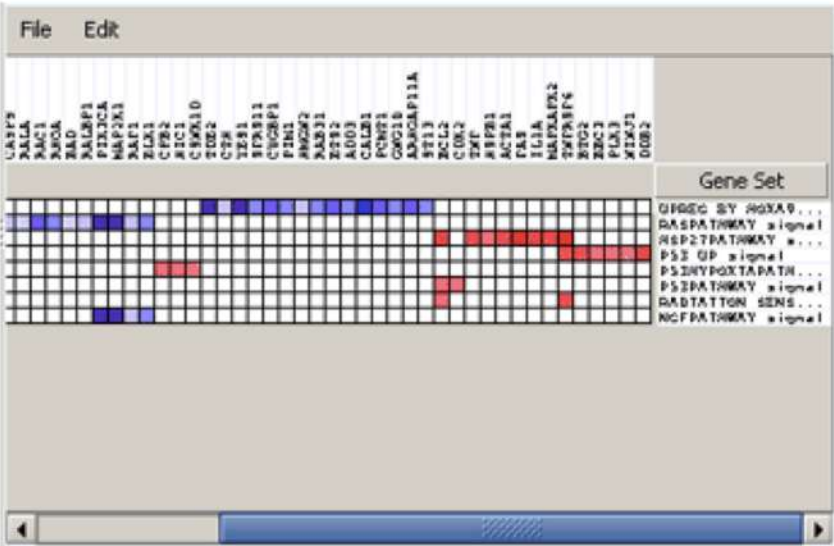Multiple testing procedures attempt to quantify and control for these.

$$\hat{FDR} = \frac{s_{null}}{s_{obs}} \cdot \frac{n_{obs}}{n_{null}}$$
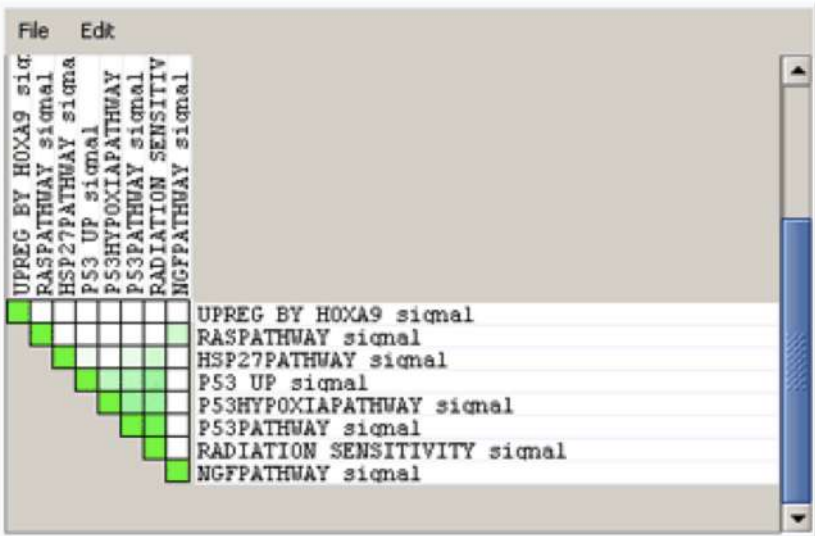
In GSEA, the collection of gene sets interrogated against the observed data are a family of hypotheses.

The recommended procedure for quantifying Type I errors is the false discovery rate (FDR). The FDR is defined as the expected value of the fraction of rejected null hypotheses that are in fact true. In practice, GSEA establishes this proportion empirically.
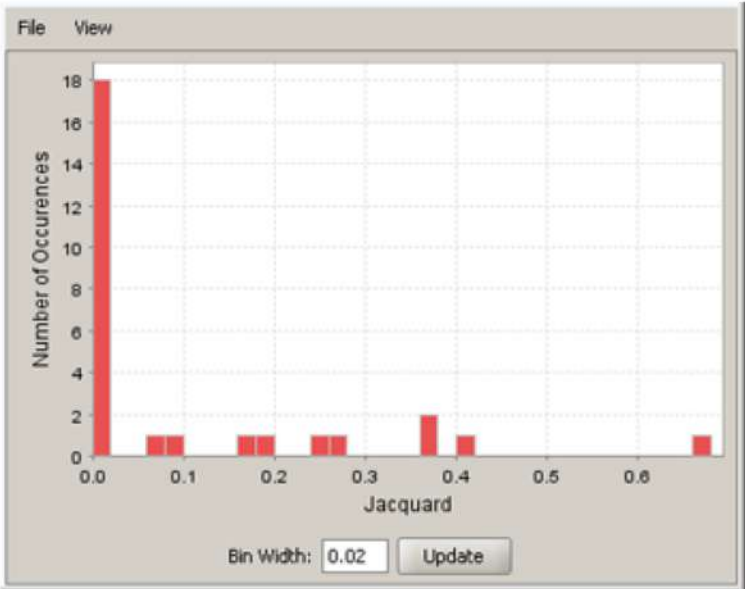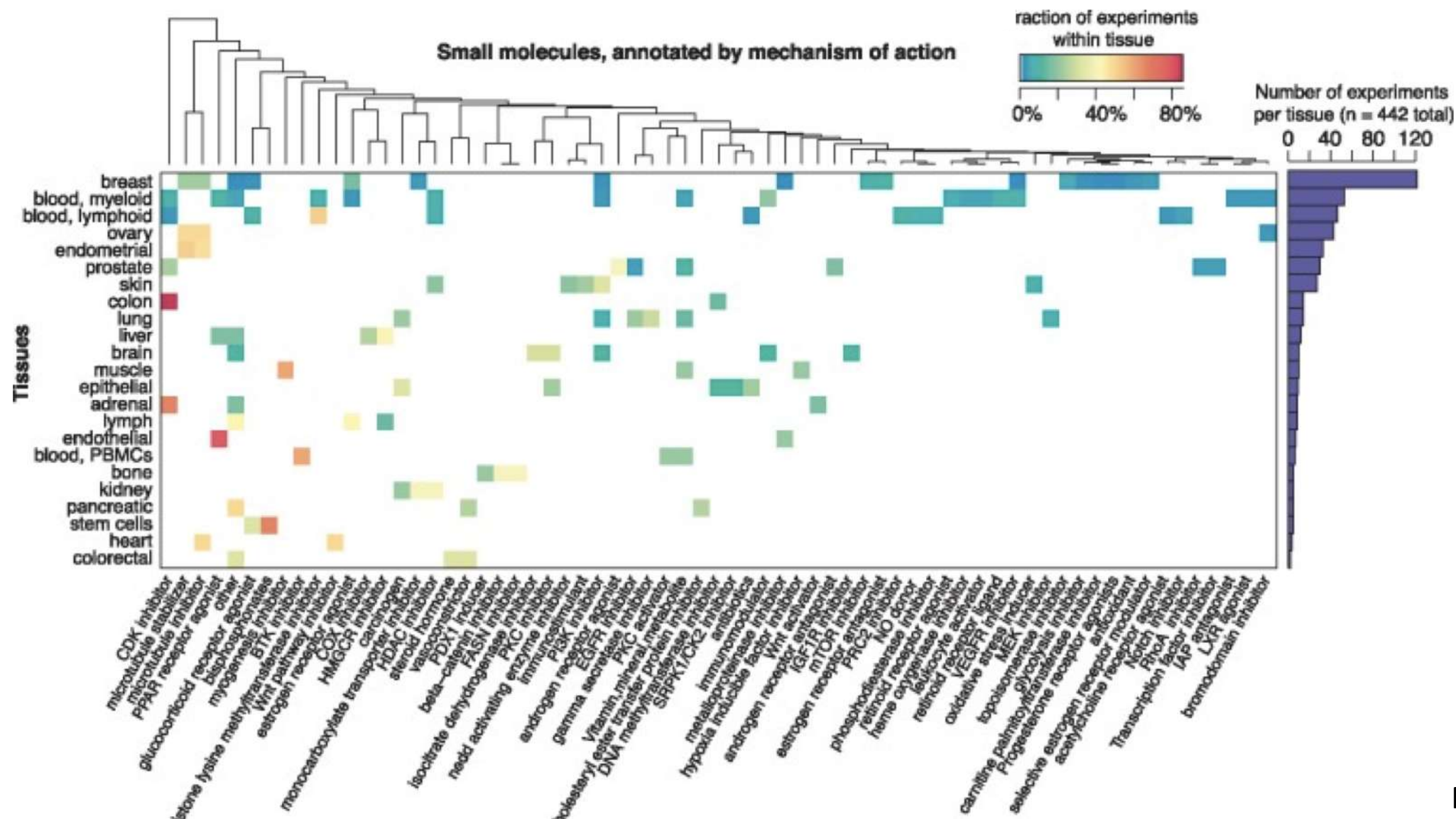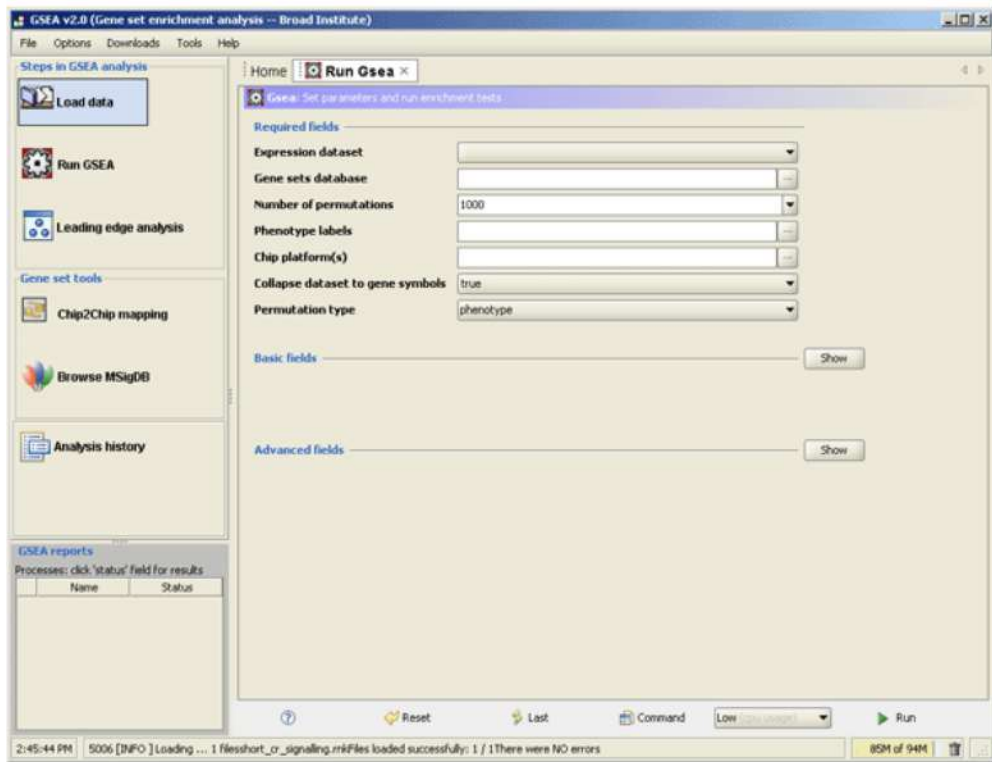
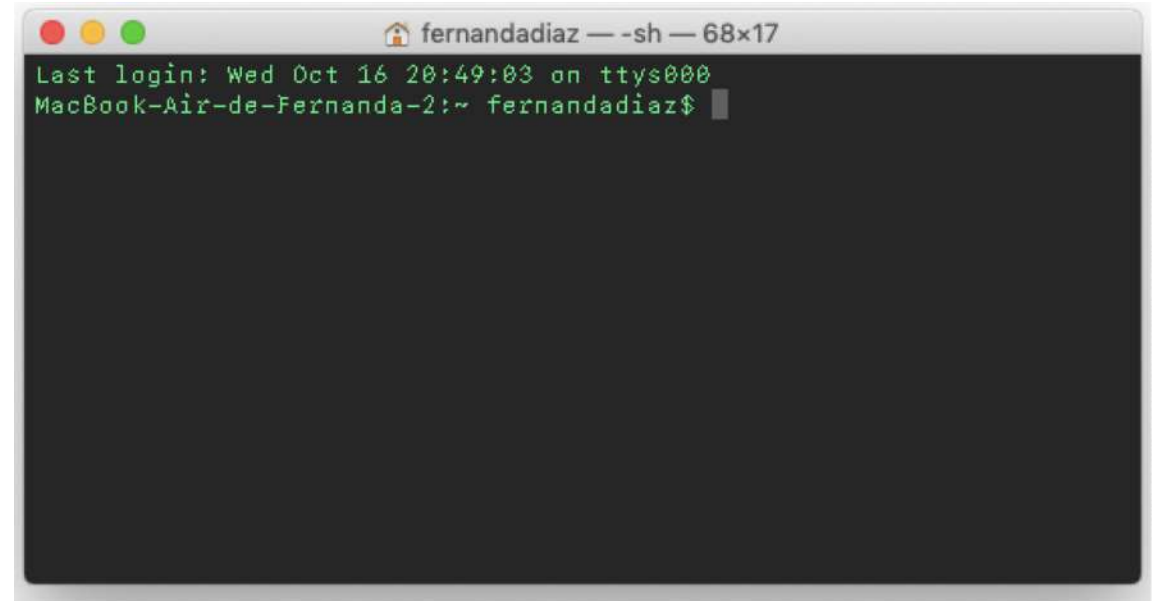# Graphs



HEATMAP                    SET TO SET                    Histograms

Small molecules, annotated by mechanism of action

Powers 2018

# Where to run



Java desktop app



Command line



R package

# Alternatives



WebGestalt — WEB-based GEne SeT AnaLysis Toolkit

*Translating gene lists into biological insights...*

Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., ... & Mesirov, J. P. (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*, *102*(43), 15545-15550.

https://www.pathwaycommons.org/guide/primers/data_analysis/gsea/