# Gene set enrichment analysis in R

Kim Dill-McFarland

2021.07.14

# Gene set

two or more genes with shared function, structure, localization, or any other defining similarity

# Gene set features

- Defined by individual study to general knowledge

- Vary in size

- Not exclusive (1 gene in many sets)

- Redundant and overlapping

# Broad Molecular Signatures Database

Hallmark

- summarize and represent specific well-defined biological processes
- generated by computational methodology based on identifying overlaps between gene sets in other MSigDB collections and retaining genes that display coordinate expression

- Total: 50
- Examples: glycolysis, inflammatory response, apoptosis

# Broad Molecular Signatures Database

Curated (C2) Canonical pathways

- from pathway databases including KEGG, REACTOME, etc
- Canonical representations of a biological process compiled by domain experts


- Total: 2922
- Examples: Caspase pathway, signaling by NOTCH1, DNA repair

# Broad Molecular Signatures Database

Gene ontology (C5) Biological process

- molecular-level activities performed by gene products


- Total: 7481

- Examples: viral life cycle, vitamin D biosynthetic process, mitochondrial calcium-ion transmembrane transport

# Broad Molecular Signatures Database

- too specific (individual studies)
  - C2 chemical and genetic perturbations
  - C7 immunologic signatures

- areas not relevant to our experimental design
  - C1 chromosome position
  - C3 gene regulation
  - C4/C6 cancer-oriented

# Hypergeometric enrichment
# (aka Fisher's Exact test)

Probability that the number of significant genes in a gene set occurred by chance
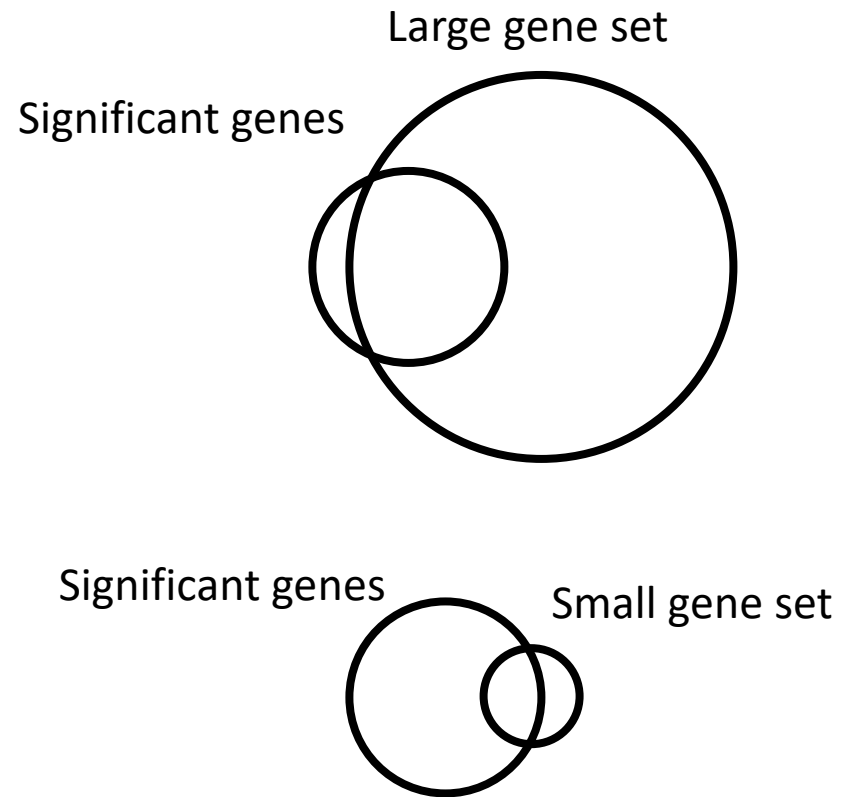
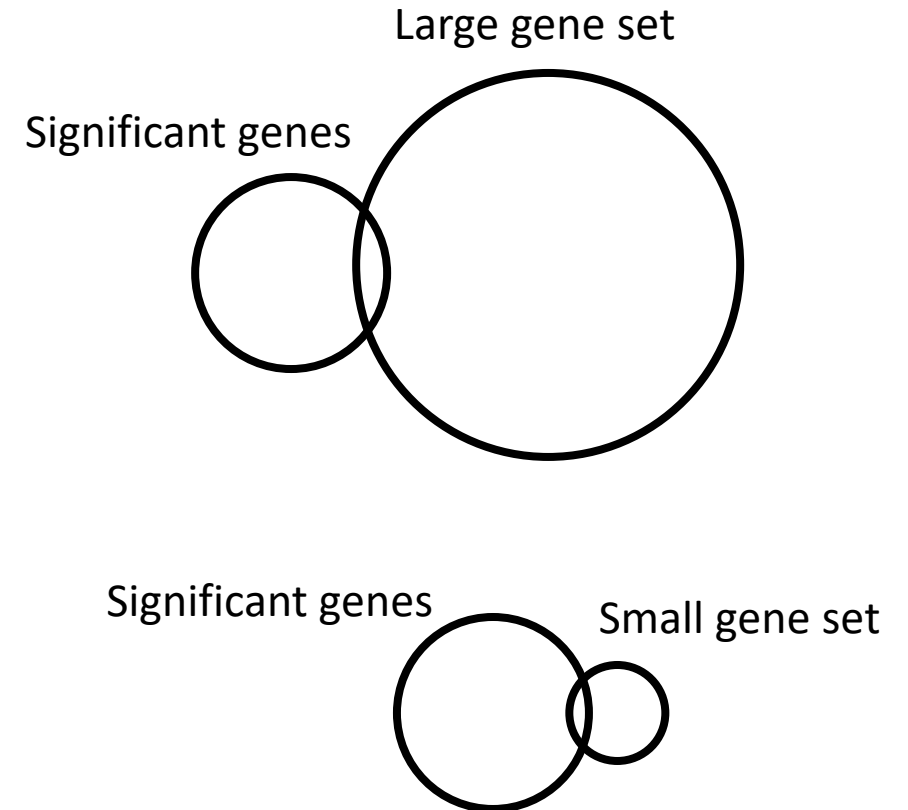# Hypergeometric enrichment

1. Define genes in a gene set


2. Define genes significant for your variable of interest (in our case, media vs Mtb-infected)


3. Calculate proportion of significant genes in gene set


4. Estimate probability and significance (P-value)

# Hypergeometric enrichment

**Likely enriched**

Large gene set

Significant genes

Significant genes

Small gene set

**Not likely enriched**

Large gene set

Significant genes

Significant genes

Small gene set

# Gene set enrichment analysis (GSEA)

Compares expression in two biological states and determines if gene sets show significant, concordant change
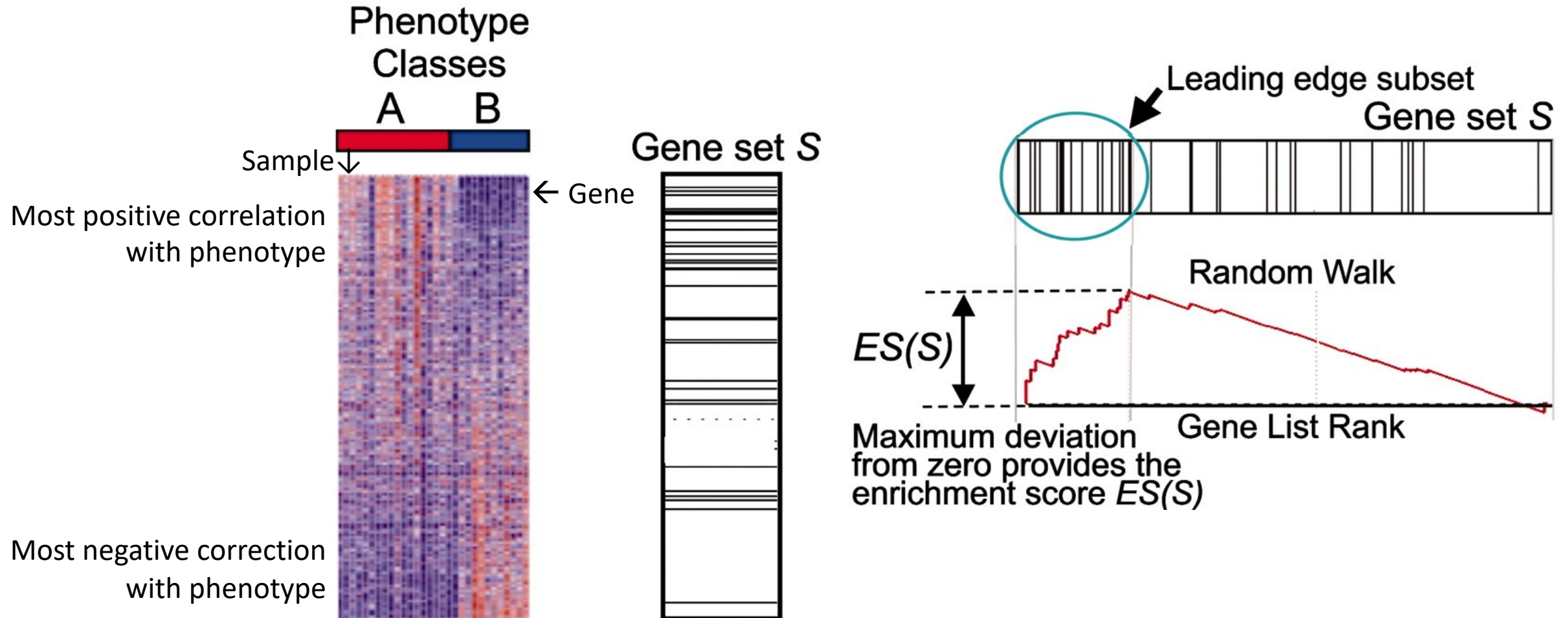
aka

Determines if genes in a set generally increase or decrease in expression

# Gene set enrichment analysis (GSEA)

1. Define genes in a gene set

2. Calculate fold change

3. Order genes by fold change and map to gene set

4. Estimate enrichment score and significance (P-value)

# Gene set enrichment analysis (GSEA)

## Hypergeometric enrichment

- Significant genes

- Binary significant vs not

- Depends on how you define significance

- Compare 2+ states

## Gene set enrichment analysis

- All genes

- Numeric fold change values

- Mean fold change per gene

- Compare only 2 states