

Swiss Institute of Bioinformatics

Introduction to RNA-Seq: Enrichment analysis

Wandrille Duchemin





What to do with your list of differentially expressed gene?



What to do with your list of differentially expressed gene?





What to do with your list of differentially expressed gene?





What to do with your list of differentially expressed gene?





Regrouping genes together in meaningful sets

- Same pathway
- Same location in the cell
- Same molecular function

• • • •



- Gene Ontologies
- Reactome
- KEGG
- MSigDB
- Custom set
- •



Gene Ontologies <u>geneontology.org</u> 3 domains of nested terms:

- Molecular Function
- Cellular Component
- Biological Process







MSigDB : <u>http://www.gsea-msigdb.org/gsea/msigdb/index.jsp</u>

Human, mouse, and rat only



hallmark gene sets are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

1 positional gene sets for each human chromosome and cytogenetic band.





regulatory target gene sets based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites.

computational gene sets defined by mining large collections of cancer-oriented microarray data.





oncogenic signature gene sets defined directly from microarray gene expression data from cancer gene perturbations.

7 **immunologic signature gene sets** represent cell states and perturbations within the immune system.



cell type signature gene sets curated from cluster markers identified in single-cell sequencing studies of human tissue.



KEGG : <u>http://www.genome.jp/kegg</u>KEGG PATHWAY



Projection of DE data onto a KEGG pathway map With R package pathview

https://bioconductor.org/packages/release/bioc/ html/pathview.html



Custom gene sets

- Derived from specialized litterature
- Tentative annotation

Very important in non-model organisms



2 main approaches

- Over Representation Analysis (ORA)
- Gene Set Enrichment Analysis (GSEA)



Over Representation Analysis (ORA)

• Basically a Fisher's exact test with p-value correction

	DE	Not DE	N = A+B+C+D #	N = A+B+C+D # total genes	
			M = A+B #	genes in set	
in gene set	А	В		9	
not in gene set	С	D	n = A+C #	DE genes	
			K = A #	JE genes in set	

$$p=1-\sum_{i=0}^{k-1}rac{\binom{M}{i}\binom{N-M}{n-i}}{\binom{N}{n}}$$



Gene Set Enrichment Analysis (GSEA)

- Does not rely on o/1 DE
- Use ranking along a continuous measure (eg. log2FC)
- Compute Enrichment Score
- Estimate significance with a permutation scheme





Enrichment Score visualized using functions from the R package enrichplot









+ SIB









Many more methods or implementations

- Signaling Pathway Impact Analysis
 <u>https://bioconductor.org/packages/release/bioc/html/SPIA.html</u>
- ISMARA (TF-based) ismara.unibas.ch/mara



Practical





Thank you



