Transcriptome Lecture 3

Outline

- Multiple Testing Procedures
- Data Visualization, Distance Measures
- Clustering
- Gene Annotation and Enrichment Analysis

The problem

- After differential expression testing, we obtained a list of significantly differentially expressed probes, controlled for false discovery rate
- We want to understand the biological insight behind this list
 - 1. we need to map the gene annotation information to these probes or gene IDs
 - 2. we want to test/infer whether an annotation is significantly enriched in our list

Annotation mapping

- What annotation information can we map probes or gene IDs to?
 - Chromosome, genes, protein family, structure, sequence, variations...
 - Gene Ontology, KEGG Pathway,...
 - Published literatures...

Annotation mapping: example



Metabolic Pathways

- PMN: Plant Metabolic Network (<u>http://www.plantcyc.org/</u>)
- MetaCyc (<u>http://metacyc.org/</u>)
- KEGG: Kyoto Encyclopedia of Genes and Genomes (<u>http://www.genome.jp/kegg/kegg2.html</u>)
- Reactome (<u>http://www.reactome.org/</u>)
- PANTHER PATHWAYS (<u>http://www.pantherdb.org/pathway/</u>)
- Pathways Commons (<u>http://www.pathwaycommons.org/pc/home.do</u>)

KEGG Pathway

- KEGG Pathways:
 - Manually curated pathway maps representing the knowledge on the molecular interaction and reaction networks, for a large selection of organisms
 - The KEGG pathways include a collection of pathways important in:
 - Metabolism
 - Genetic Information Processing
 - Environmental Information Processing
 - Cellular Processes
 - Human Disease
 - ...

KEGG Pathway: An example



Annotation mapping: example



Gene Ontology (GO)

- Gene Ontology (GO) is a collection of controlled vocabularies describing the biology of a gene product in any organism
- <u>http://www.geneontology.org/</u>
- Very useful for interpreting biological insight of microarray data – and it is computable!

So what is ontology?

From a practical view, ontology is the representation of something we know about. "Ontologies" consist of a representation of things, that are detectable or directly observable, and the relationships between those things.







The GO has three Ontologies

Molecular Function

GO term: Malate dehydrogenase. GO id: GO:0030060 (S)-malate + <u>NAD(+)</u> = <u>oxaloacetate</u> + <u>NADH</u>.







Cellular Component

GO term: mitochondrion GO id: GO:0005739

Biological Process

GO term: tricarboxylic acid cycleSynonym:Krebs cycleSynonym:citric acid cycleGO id:GO:0006099



- Parent / child network organized • as a tree
- Terms get more detailed as you move down the network



Gene Ontology: Rule

- In GO, a gene can be •
 - present in any of the ontologies (MF / BP / CC)
 - a member of several GO terms
 - a gene must be a leaf in GO trees
- If a gene is a member of a term, it ۲ is also a member if the terms parents



Gene Ontology: files

- Ontology file: GO terms and relationships in a variety of formats. The ontology file is unique for all species.
- Annotation files: associations between gene products and GO terms submitted by members and associates of the GO consortium. Different species have different annotation files.
 - ✓ gene_association.tair
 - ✓ gene_association.goa_human

GO tools

- GO resources are freely available to anyone to use without restriction
 - Includes the ontologies, gene associations and tools developed by GO
- Other groups have used GO to create tools for many purposes:

http://amigo.geneontology.org/amigo

Gene Ontology: tools



Grouping by Biological process



Using GO in practice

statistical measure

how likely your differentially regulated genes fall

into that category by chance





The problem

- After differential expression testing, we obtained a list of significantly differentially expressed probesets, controlled for false discovery
- We want to understand the biological insight behind this list
 - 1. we need to map the gene annotation information to these probesets
 - 2. we need to test/infer whether an annotation is significantly enriched in our list

Annotation Testing (enrichment analysis)

- We want to ask:
 - Are there any GO terms overrepresented in the obtained gene list, compared with what would happen by chance?
 - Hypergeometric test
 - Fisher's exact test fisher.test()
 - Binomial test binom.test()
 - Chi-squared test chisq.test()
 - Kolmogorov-Smirnov test

Hypergeometric distribution

• The hypergeometric distribution arises from sampling from a fixed population.



Hypergeometric test



Hypergeometric test



- **TEST:** We want to calculate the probability for drawing 7 or more white balls out of 10 balls given the distribution of balls in the urn.
- The smaller the possibility is, the more significantly enriched.

Annotation Testing (Hypergeometric test)

- Example: we obtained a list of 80 significant genes from a gene expression experiment of yeast.
- Yeast has 6000 genes, and 100 of them can be mapped to a GO term called "Cell cycle". For the 80 significant differentially expressed genes from micrroarray/RNA-seq, 10 are mapped to this GO term.
 - Is this observation a significant event? Or, is the GO term "Cell cycle" significantly over-represented in our list of 80



GO enrichment analysis in R

- Gostat
 - <u>http://www.bioconductor.org/packages/2.3/</u>
 <u>bioc/html/GOstats.html</u>
- PGSEA
 - <u>http://www.bioconductor.org/packages/2.4/</u>
 <u>bioc/html/PGSEA.html</u>

Online tools

- DAVID http://david.abcc.ncifcrf.gov/
- GoMiner: <u>http://discover.nci.nih.gov/gominer</u>
- GOstat: http://gostat.wehi.edu.au
- GSEA: <u>http://www.broadinstitute.org/gsea/index.jsp</u>
- BIOBASE (Whitehead has license)
- BiNGO (uses Cytoscape)

DAVID: a function annotation tool

http://david.abcc.ncifcrf.gov/

