Biological Sequence Analysis Spring 2009          Dr. Richard Friedman
(212)851-4765 friedman@cancercenter.columbia.edu          824 ICRC

## Lesson 13: Functional Genomics II: Functional Genomics Databases and Overrepresentation Analysis and Clustering of Microarray Data

In this class we are going to learn the basics of using web-accessible function and pathway databases, and how to find the functions and pathways associated with a set of differentially expressed genes.

## Theory

1. Pevsner on functional databases:
>           A. OMIM p. 659-661.
>           B. GO 243-246.
>           C. KEGG 258-263.

2.  Ontools. Read sections on Onto-Express and Pathway-Express.

>      http://vortex.cs.wayne.edu/projects.htm

3. Pevsner on clustering or microarray data: p. 203-214.


Theory of Hierarchical  Clustering (Unsupervised Learning)

$d_{xy}$= distance in expression patterns between genes x and y.

$$d_{xy} = 1 - r_{xy} = 1 - \frac{1}{N} \sum_{i=1,N} \left( \frac{x_i}{\sqrt{\sum_{i=1,N} x_i^2}} \frac{y_i}{\sqrt{\sum_{i=1,N} y_i^2}} \right)$$

$r_{xy}$= modified Pearson correlation between expression of genes x and y with zero mean.
N= Number of samples.


Then a tree is inferred based on distances.


**Summary of Commands:**

>     Note: In this document different fonts have different meanings:

>     Times is used to explain commands.

```
Courier is used to indicate commands and command options.

Courier italics are used to indicate command parameters,
for   example, filenames.

Courier bold is used to indicate commands that are not
displayed.

Courier bold italics are used to indicate computer-generated
output.
```

Helvetica is used to indicate menu items.


| | |
|---|---|
| http://www.ncbi.nlm.nih.gov/entrez | Accesses the NCBI ENTREZ web page |
| Gene | Diverse information on genes including molecular interactions. |
| OMIM | Descriptive genotype-phenotype relationships in humans. |
| OMIA | Descriptive genotype-phenotype relationships in organisms other than humans. |
| http://www.geneontology.org/ | Gene Ontology Gives Biological Process, Molecular function, Cellular Component for proteins/ |
| Gene symbol/name, exact match | To search for a gene/protein This will give you all of the terms associated with a protein. |
| Terms | To search for a term. Following the links associated with a term will give you a definition of the term and the names of all of the proteins associated with that term. |
| http://bond.unleashedinformatics.com/Action? BOND: Biomolecular Object Network Database. | A database that gives interaction partners between biological molecules. |
| Search bind using a field specific query | Enables user to specify the kind of search. |
| Field: *Interaction description* contains *at least one of the following words* | Generates a list of interactions in which a specified molecule plays a role. |
| On the list page clicking on an: interaction identifier | gives a description of the interaction in question. |

On the interaction page, clicking on:

| | |
|---|---|
| expand all | gives a more detailed picture of the interaction in question. |

| | |
|---|---|
| http://www.genome.jp/kegg/ | KEGG: Kyoto Encylopedia of Genes and Genomes. |

Click on KEGG Table of Contents

| | |
|---|---|
| KEGG PATHWAY | Encyclopedia of Pathways. |
| KEGG GENES | Encyclopedia of Genes. |

SEARCH

| | |
|---|---|
| KEGG for *subject* | Searches Gene database for name of gene. |
| GENES for *genename* | Searches Gene database for name of gene. |
| PATHWAYS for *pathway* | Searches Pathway database for name of pathway. |
| SEARCH *organismcode* for *genename* | Searches for name of gene in specific organism. |

| | |
|---|---|
| Organism | Searches for organism code. |
| GENES for *Genename* | Searches Gene database for name of gene. |
| PATHWAYS for *Pathway* | Searches Pathway database for name of pathway. |

Analysis with Onto-tools

A. Preprocessing for ontools:
  1. Move output files to a directory that you can work in.
  2. Open output files with Microsoft Excel. Autoformat column length and save file as Excel Workbook with new name.
  3. Copy all genes with B >0 to new Excel Workbook and save as Workbook with new name.
  4. Delete all columns but the Affymetrix probeset identifier. Delete the header row. save as text file.
B. Analysis with Onto-Express and Pathway-Express.

| | |
|---|---|
| http://vortex.cs.wayne.edu/ontoexpress/ | Onto-Express web-site. |
| Login: | Login to web-site. |
| Onto-express | Select Onto-Express for finding overrepresented Gene Ontology Categories. |

Fill out the ontoexpress input window as below (settings for experiment- the parameters will of course vary with other files, organisms, and chip).

Onto-express output files:

*bioprocess*　　　　　　　　　　　Overrepresented biological process
　　　　　　　　　　　　　　　　　categories.

*molecular*　　　　　　　　　　　Overrepresented molecular function
　　　　　　　　　　　　　　　　　categories.

*cellular*　　　　　　　　　　　　Overrepresented cellular location
　　　　　　　　　　　　　　　　　categories.

*chromosome*　　　　　　　　　　Assignment of genes to chromosomes.


To process output files with Excel:

Open　　　　　　　　　　　　　　Opens file from inside Excel.

　　Files as type *textfiles*　　　Lists text files in the menu.

Delimited　　　　　　　　　　　　Reads file as entried separated by delimiters.

　　Tab　　　　　　　　　　　　　Recognizes tabs as delimiters.

　　Semi-colon　　　　　　　　　Recognizes semi-colons as delimiters.

Format->Column->Autofit　　　　　Fits selected columns.

| | |
|---|---|
| Data->Sort | Sorts by the selected column. |
|     Expand Selection | Sorts rest of row according to the selected column. |
| File->Save As | Save as file with name typed in box. |
|     Save As type *Microsoft Excel Workbook* | |
| | Save as file in Excel Format. |

1. Open each of these files with Excel (check delimited and tab and semicolon).Autofit selection for column, sort by corrected p-value, and save as a Workbook.



The above box will output the same 4 files as before only with the genes included.
In addition, it will output, "treeview", a file with all of the gene ontology categories detected and "treeview_input", a file with all of the gene-ontology categories detected along with input genes from the list. The above files cannot be sorted.by corrected p-values without getting mixed up, because there is more than one row per Gene Ontology category.

| | |
|---|---|
| Pathway-Express | Select Pathway-Express for finding overrepresented KEGG pathwayCategories. Use list of pathways with log fold changes for input. |

Advanced options
        Corrections = fdr                                          Corrects pathways for false discoveries with
                                                                   false discovery rate.

Pathway express output-windows and options
Bar Graph                                                          Display of overrepresented pathways..
Pathway Details                                                    Shows list of pathways in tables.
        (Right mouse button over pathway link)   Lists pathway options.
                Show Pathway Genes                                 Changes pathway in Pathway Gene Details.
                Show Pathway Details                               Shows KEGG diagram of pathways with
                                                                   genes on list indicated.

                Save Table                                         Saves Pathway Details file.

Pathway Gene Details                                               Shows list of genes in selected pathway and
                                                                   indicates whether or not they are part of the
                                                                   input list.

        (Right mouse button over genelink)       Lists gene options options, which are
                                                                   analogous to Pathway Details options.

Input Details                                                      Shows list of genes in input lists and
                                                                   indicates whether or not they are part of a
                                                                   pathway. Mouse options analogous to other
                                                                   windows.

Hierarchical clustering with Cluster 3.0 and JavaTreeview.

1. Start R and go to the estrogen working directory.



2. Load the Affy Program:

3. In R window:

```
estrogenEset<-ReadAffy()
```
Loads Estrogen cell files into Expression object.
```
estrogenmas5 <-mas5(estrogenEset)
```
Normalizes cel files by the MAS5 algorithm.
```
write.exprs(estrogenmas5,"estrogenmas5.txt")
```
Saves normalized arrays as text file.

4. Cut and paste file to directory which you are working in for lab.

5. Copy file with differential expression values from last week to file called id_gene.
    A, Open it with excel.
    B. Delete all columns other than Gene-and-symbol.
    C. Sort by ID.
    D. Rename Symbol Column "Gene".
    The spreadsheet should look like this:

6. Open Estrogenmas5 within Excel.
   A. Make 2 new columns.
   B. Paste in the contents of ID gene so that the id column should be aligned
   The file should look like this:



7. Delete the 2 id columns so that the spreadsheet looks like this:



8. Save file as text file.called "estrogenmas5gene"

9.  Open Cluster 3.0 on the PC and load mas5gene



10. Click adjust data. Check
    A.  Log Transform
    B.  Center Genes
    C.  Center Arrays



D. Click Apply

Click filter data

11. In filter data Check
- A.  SD (Gene Vector ) 2.0
- B.  At least one observation with abs(val) >= 2.0
- C.  MaxVal-MinVal>=2.0



D.  Click Apply filter. The GUI should say "70 passed out of 12625"



E.  Click Accept filter

12, Click on Hierarchical tab. Check
>    A. Cluster Genes.
>    B. Cluster Arrays



Then click "Average linkage"

Files entitled:

>    A.  estrogenmas_gene.atr      (ATR)
>    B.  estrogenmas_gene.          (CDT)
>    C.  estrogenmas_gene.          (ATR)

Should appear in your directory. You should now be ready to display the cluster diagram and heatmap.

13. Java Treeview
      A. Open Java Treeview.
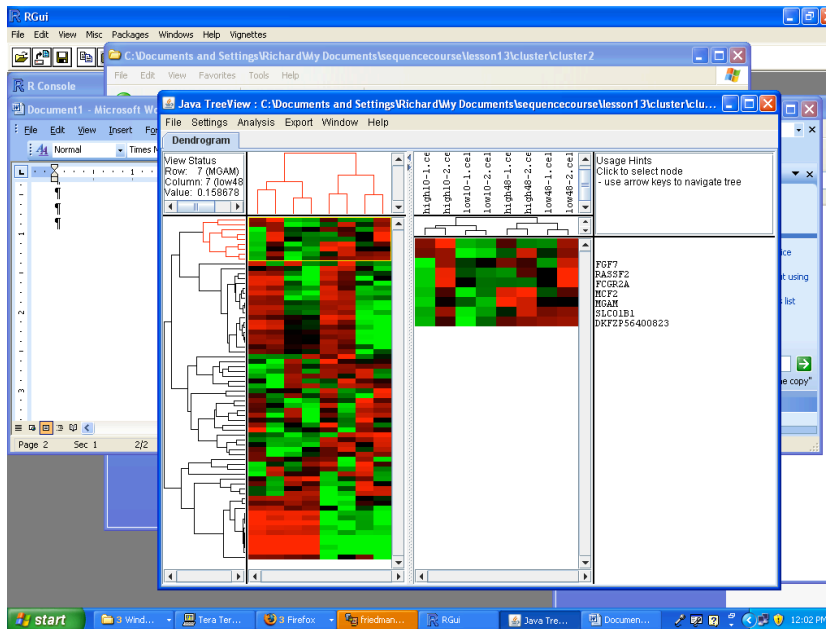      B. Select Open-> File
      C. Scroll to directory containing CDT file
      D. Open CDT file



14. You should get  a heat map that looks like this.

15. You can select and magnify clusters by clicking on nodes:



16. The resulting heat map has good expansion for exploring clusters.

## Lab

1. Characterize the function, interactions and pathways of the c-src (or the protein of your choice) using web-accessible databases. Save the necessary information.
2. What is the phenotype associated with the BRCA1 185DELAG mutation in humans?
3. From the list of differentially expressed genes obtained in the 10 hour estrogen experiment explained in part I, generate the following lists in Excel Workbook format:

   A. The Biological Function Gene-ontology values overrepresented in the Estrogen experiment sorted by corrected p-values.
   B. The Biological Function Gene-ontology values overrepresented in the Estrogen experiment with gene symbols and other identifiers (unsorted).
   C. The KEGG pathways overrepresented in the Estrogen experiment.

4. Produce a heatmap representing the gene expression patterns of the estrogen dataset from the last lab clustered by both gene and array.