- 1. Find a research article for your presentation and final exam
- 2. Prepare your presentation
- 3.No class for the week of Midterm exam
- 4.http://sysbio.unl.edu/Teaching/ BIOS497897_2014/



Transcriptome

Lecture 1

Outline

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

The problem

- After differential expression testing (from RNAseq or Microarray assay), a list of P-value is obtained, one for each gene.
- Most investigators want to
 - Identify the genes that are differentially expressed
 - Estimate the proportion of errors in the list of selected "differentially expressed genes"

A single gene example (low throughput case)

- Suppose you are only interested in a single gene.
- You want to compare the expression level (the level of transcription) of this gene between two conditions (control and treatment).
- For each treatment, there are three biological samples.
- Experiments are performed on each sample to measure gene expression levels (e.g., quantitative PCR, gel blot).
- There are three independent treatment per condition.
- A t-test is performed and a p-value is obtained.
- Declare there is differential expression if p-value is below some threshold (e.g., 0.05).

Extreme parallel hypothesis testing

- With high throughput technology, we can and often perform the same hypothesis test on each and every gene.
- Thus tens of thousands of hypotheses are tested in parallel.

A naïve solution

- Since genes with small p-values are likely to be differentially expressed, why don't we just use the traditional (pre-specified) alpha = 0.05 to decide?
 - □Yes?
 - ⊠No? But why?

The Multiple Testing Problem

- Suppose one test of interest has been conducted for each of *m* genes in a microarray experiment. For example, through two-sample t-test.
- Let $p_1, p_2, ..., p_m$ denote the *p*-values corresponding to the *m* tests.
- Let H_{01} , H_{02} , ..., H_{0m} denote the null hypotheses corresponding to the *m* tests.

The Multiple Testing Problem

- Usually,
- *H*_{0i}: no differential expression for gene i
- *H*_{1i}: differential expression for gene i
- *i*=1,2,....*m*

The Multiple Testing Problem (continued)

- Suppose m_0 of the null hypotheses are true and m_1 of the null hypotheses are false.
- Let c denote a value between 0 and 1 that will serve as a cutoff for significance:
 - Reject H_{0i} if $p_i \le c$ (declare significant)

- Fail to reject (or accept) H_{0i} if $p_i > c$ (declare non-significant)

What is P-value?

- P-value is the probability of obtaining a test result as extreme as the one you are getting under the null hypothesis (i.e., area in both tails of the distribution).
 - Null hypothesis: The difference in average expression between the two groups is zero.
- The *lower* the p-value, the *less* probable the result is. (assuming the null hypothesis is true).
- Interpretation: if you repeat the same experiment many times (i.e., computing a T-statistics for each gene on a microarray), the p-value represents the proportion of times that you would expect to see a Tstatistic this extreme.

What is P-value? A more rigorous interpretation



- P-value = Prob(Type I Error) <- describe the false positive rate
- New Interpretation: if you repeat the same experiment many times (i.e., computing a t-statistic for each gene on a microarray), the p-value represents the proportion of times that you would commit a type I error (i.e., false positive call).

What does this mean to microarray data?

• The result is that we obtain one p-value for each gene

	Т1	Т 2	Т 3	N 1	N 2	N 3	T-statistics	P-value
G 1							T1	P1
G 2							Т2	p2
•••							•••	
G 20000							T20000	P20000

• 20,000 p-values...



Histogram of pvalues1

What does this mean to microarray data?

• The result is that we obtain one p-value for each gene



- If we use alpha=0.05 to decide differentially expressed genes, 5% of the 20,000 genes would then be selected by chance
- That means 1000 genes would be false positives...

A naïve solution

- Since genes with small p-values are likely to be differentially expressed, why don't we just use the traditional (pre-specified) alpha = 0.05 to decide?
 - **Yes**?
 - **No!** 20,000x0.05 = 1000 false positives!
 - If the investigator is interested in selecting 100 genes for downstream analysis, they could all be false positives by chance!
 - Other solutions?

The solutions

- To select differentially expressed genes, we need to do multiple testing (multiplicity) corrections
 - Familywise Error Rate (FWER), such as Bonferroni correction and Holm's method: adjust the p-value threshold from alpha to alpha/(number of genes)
 - Control False Discovery Rate: algorithm proposed by Benjamini & Hochberg
 - Re-sampling techniques (i.e., Permutation P-values)