

# Biostatistics

Xing-Ming Zhao, School of Electronics and Information Engineering,  
Tongji University

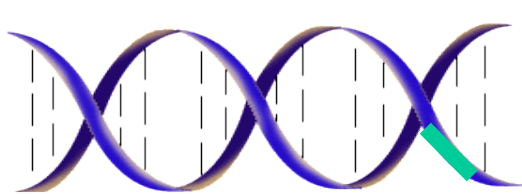
# Outline

1. Differential gene expression
2. Gene set analysis
3. Mining gene expression

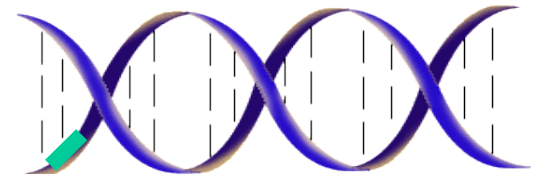
# Biological background

## Transcription

DNA



G T A A T C C T C  
| | | | | | | |  
C A T T A G G A G

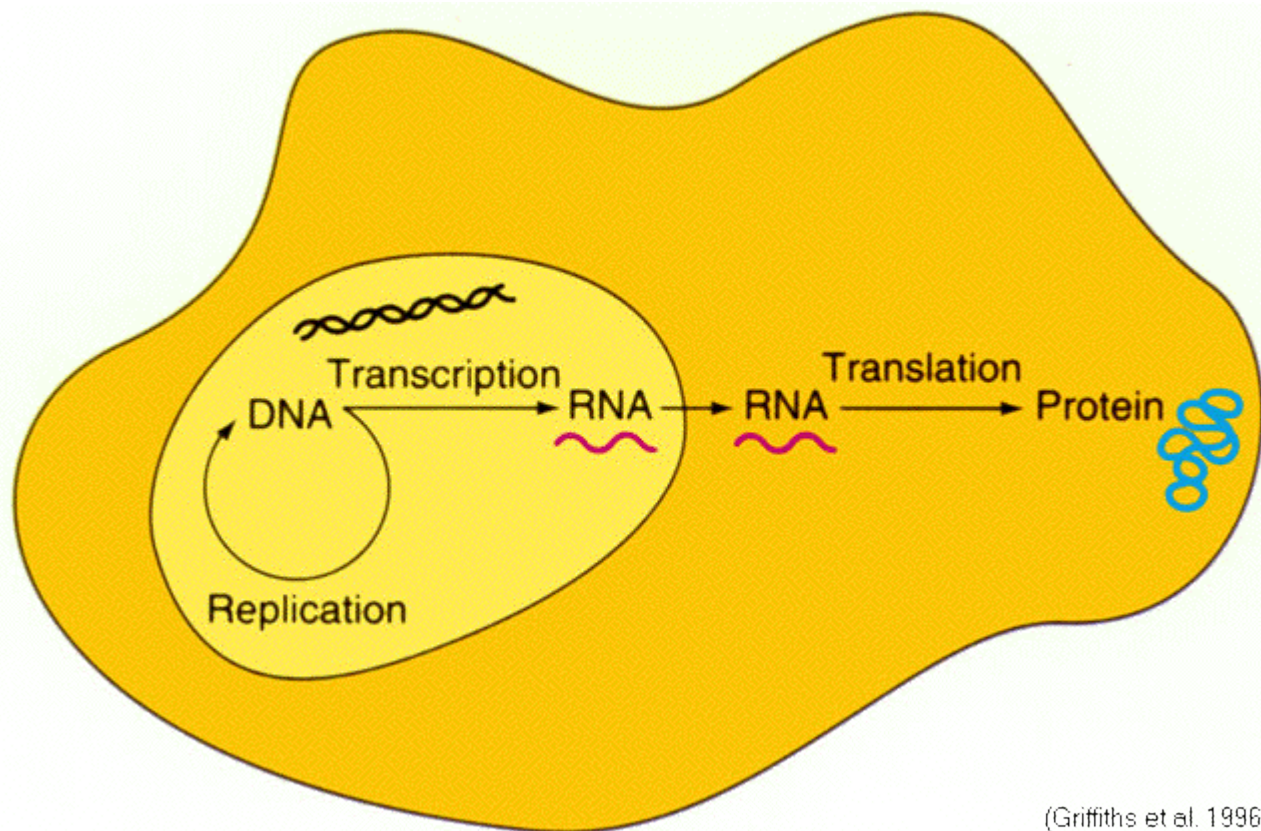


RNA  
polymerase

mRNA

G U A A U C C

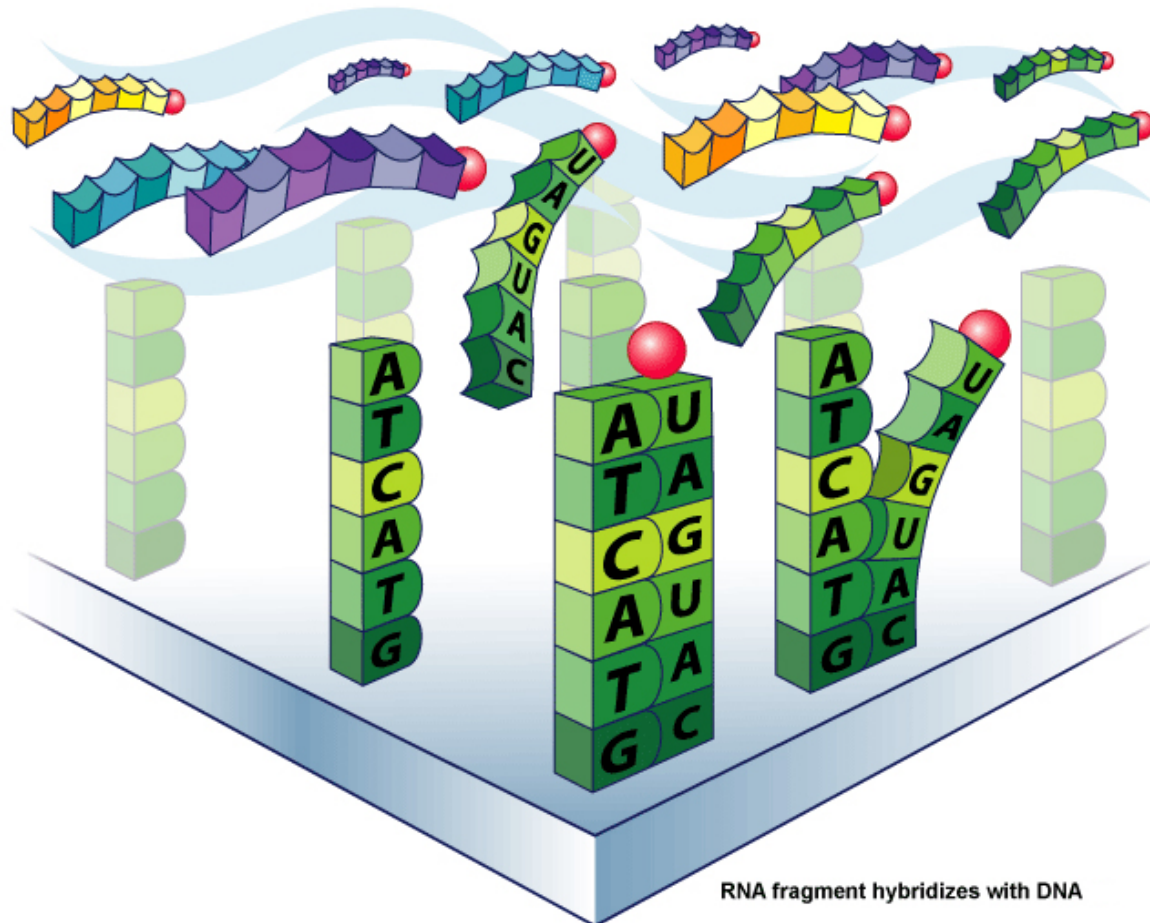
# Biological background



**the amount of mRNA tell which genes are being expressed in the cell**

# DNA microarrays

Microarrays are based on the ability of complementary strands of DNA (or DNA and RNA) to hybridize to one another in solution with high specificity.



# RNA-seq

RNA-seq is simple:

- isolate all mRNAs
- convert to cDNA using reverse transcriptase
- sequence the cDNA
- map sequences to the genome

The more times a given sequence is detected, the more abundantly transcribed it is. If enough sequences are generated, *a comprehensive and quantitative view* of the entire transcriptome of an organism or tissue can be obtained.

# What can we learn from gene expression?

mRNA levels compared in many different contexts

- Different tissues, same organism (brain vs liver)
- Same tissue, same organism (treatment vs control, tumor vs non-tumor)
- Same tissue, different organisms (evolution)
- Time course experiments (effect of treatment, development)

# What can we learn from gene expression?

- **Classifications:** for diagnosis, prediction...  
Cell-type, stage-specific, disease-related,  
treatment-related patterns of gene expression?
- **Gene Networks/Pathways:**  
Functional roles of genes in cellular processes?  
Gene regulation and gene interactions



# Public gene expression Repositories

- GeneX at US National Center for  
Genome Resources

<http://www.ncgr.org/research/genex/>

- ArrayExpress at European Bioinformatics  
Institute

<http://www.ebi.ac.uk/arrayexpress/>

# Public Repositories

- Stanford University Database

<http://genome-www4.stanford.edu/MicroArray/SMD/index.html>

- Gene Expression Omnibus at US National Center for Biotechnology Information

<http://www.ncbi.nlm.nih.gov/geo/>

# Differential Gene Expression

# Differential Gene Expression

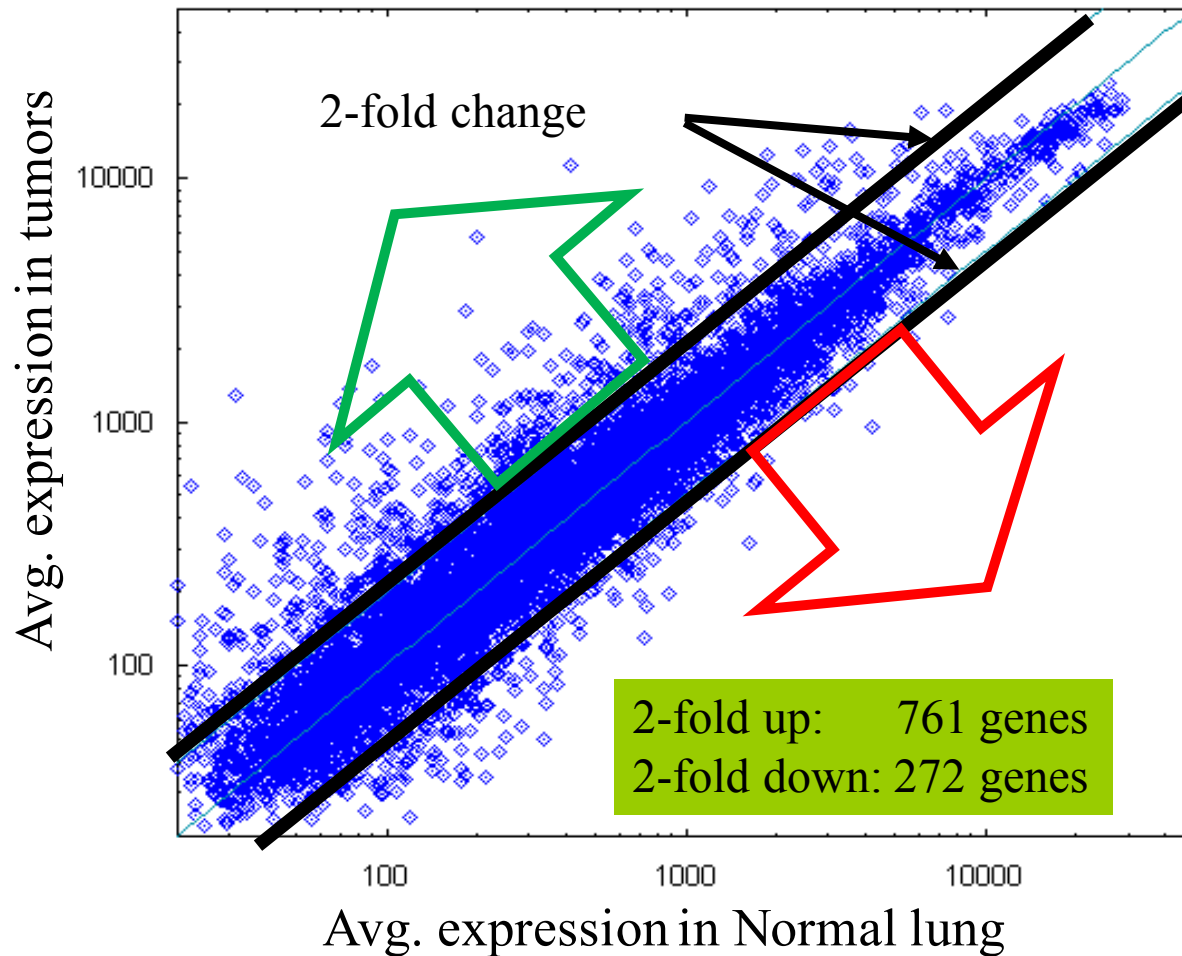
- Cell type specific
  - e.g. skin cell vs. brain cell
- Developmental stage
  - e.g. embryonic skin cell vs. adult skin cell
- Disease state
  - e.g. normal skin cell vs. skin tumor cell
- Environment-specific
  - e.g. skin cell untreated vs. treated  
drugs, toxins

# What We Need

- **Score** the genes, hopefully in a meaningful way..
- Attach a measure of **statistical significance** to the score so we can
  - Choose a subset of genes “wisely”
  - Have a measure of how strong our signal is

# Fold Change

-Simplest Score

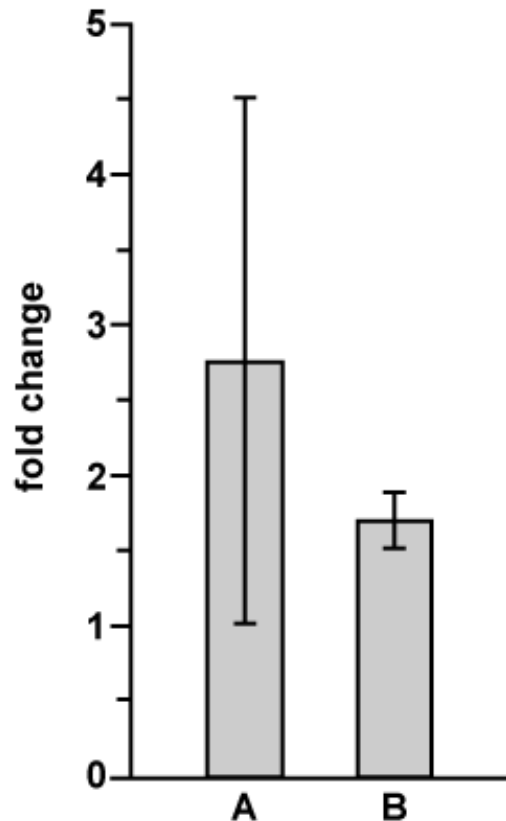


# Fold Change

- Arbitrary threshold is used for filtering (e.g. 2- or 3- fold)
  - Fold Change computation: Let  $E = \text{mean}(\text{group1})$ ,  $B = \text{mean}(\text{group2})$
  - $FC = (E - B) / \min(E, B)$
  - At low expression levels, a 3 fold change can be noise
  - At high expression levels a 1.1 fold change can be biologically important.

# Fold Change

- fold change cutoffs (e.g., 2x difference or above)
- not a very satisfying approach:
  - misses any small changes
  - doesn't take into account variance



Here, “A” has a fold change  $>2.5$ , but varies greatly between replicate experiments. “B” has a fold change of only 1.75, but changes reliably each time the experiment is performed.



# Statistical Tests

Using statistical tests is usually a better way of determining which genes are truly changed in your microarray experiment. The *p-value* tells you how much confidence to place in your result.

*p-value*: The chance of rejecting the null hypothesis by coincidence

For gene expression analysis we can say:  
the chance that a gene is categorized as differentially expressed by coincidence

# Statistical Tests

- Statistical tests seek to make conclusions about parameters (ex. Mean gene expression in groups 1 and 2) on the basis of data in a sample.
- Structure of statistical tests involve testing two hypotheses. The null hypothesis  $H_0$  and the alternative hypothesis  $H_A$ .
- The null hypothesis is usually a hypothesis of no difference.

# ***t***-test

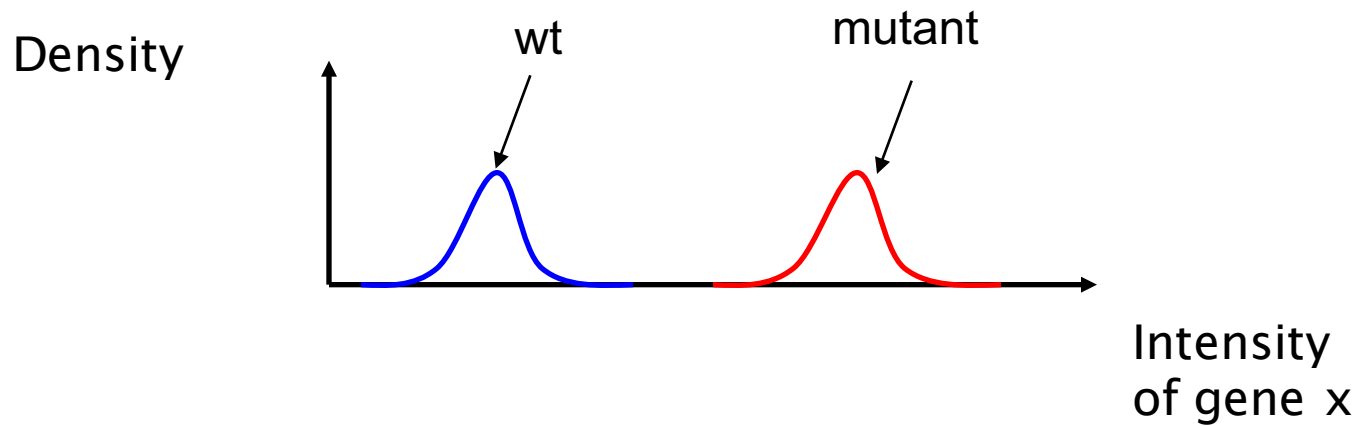
Calculate  $T$

$$T = \frac{\overline{X}_2 - \overline{X}_1}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Lookup  $T$  in a table

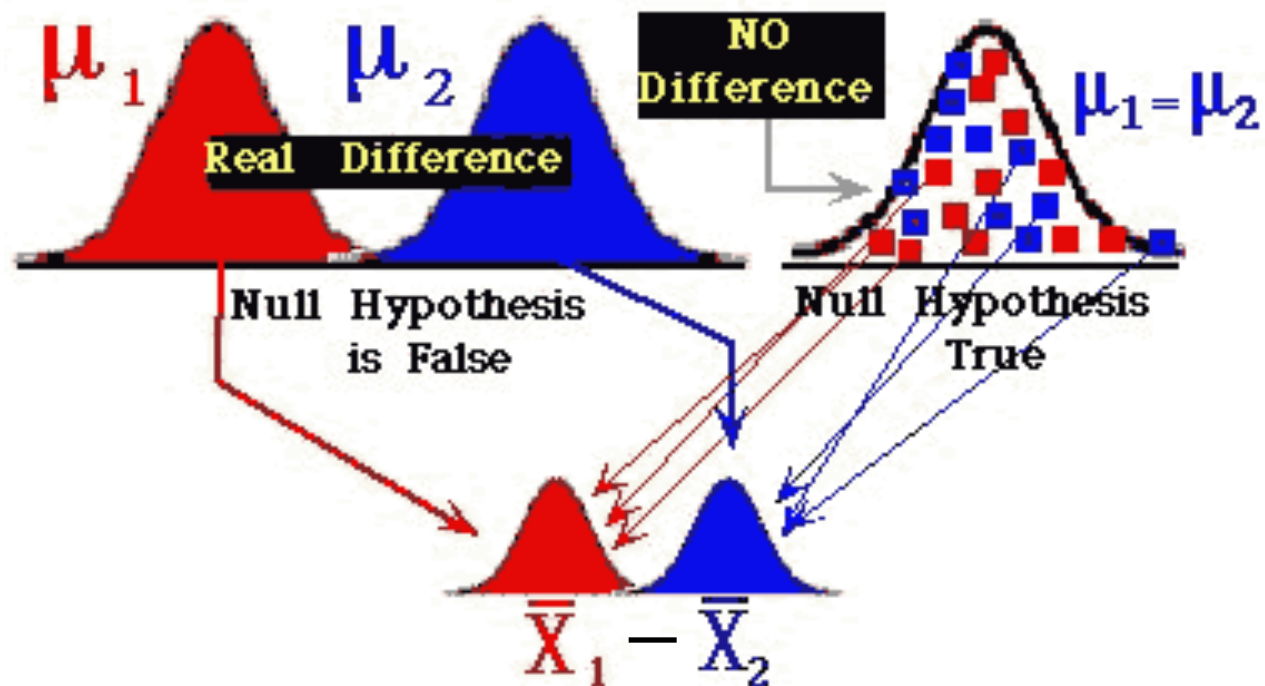
# $t$ -test

The  $t$ -test tests for difference in means



# *t*-test

The *t* statistic is based on the sample mean and variance



$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

# *t*-test: example

- Many genes
- 2 conditions
- A few replicates per condition

Gene	Condition 1			Condition 2		
	Rep 1	Rep 2	Rep 3	Rep 1	Rep 2	Rep 3
A	150	160	150	180	190	180
B	50	40	45	50	45	40
C	800	760	680	400	450	425
...						

# ***t*-test: example**

- Conditions can be control vs treated, different cell types, different time points, etc.
- Typical Question – Which genes' expression levels change due to condition?
  - *t*-test

# *t*-test

- Are we only considering up-regulated or down-regulated genes, or both?
  - If both, perform a 2-tailed test
- Can we assume that the variance of the gene is similar in both samples?
  - Yes => Homoscedastic (the standard t-test)
  - No => Heteroscedastic (Welch's test)



# ANOVA

- We want to measure how gene expression changes under different conditions.
  - Only two conditions and an *adequate* number of replicates → t-tests & extensions
  - More than two conditions / more than one factor

Analysis of Variance (ANOVA)

# ANOVA

- We want to determine when the variation due to gene expression is significant, but...
- There are multiple sources of variation in measurements *besides just gene expression*.
- We want to know when the variation in measurements is caused by
  - varying levels of gene expression
  - versus other factors.

# ANOVA

- Analysis of variance – like a multidimensional  $t$ -test
- Measure effects of multiple treatments and their interactions
- A thoughtful ANOVA design can help answer several questions with one analysis
- ANOVA can also analyze factors that should be controlled
  - just to confirm absence of confounding effects
- ANOVA generally identifies genes that are influenced by some factor – but then post-hoc tests must be run to identify the specific nature of the influence
  - Ex:  $t$ -test between all pairs of data

# ANOVA

- Some sources of variation in the measurements in microarray experiments are:
  - Array effects
  - Dye effects
  - Variety effects
  - Gene effects
  - Combinations

# One-way ANOVA

- Suppose you have a model for each measurement in your experiment:

$$y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

- $y_{ij}$  is  $j^{th}$  measurement for  $i^{th}$  group.
  - $\mu$  : overall mean effect (constant)
  - $\alpha_i$  :  $i^{th}$  group effect (constant)
  - $\varepsilon_{ij}$  : experimental error term  $\sim N(0, \sigma^2)$
- Therefore, observations from group  $i$  are distributed with mean  $\mu + \alpha_i$  and variance  $\sigma^2$ .

# One-way ANOVA

- Many genes
- Multiple conditions
- A few replicates per condition

# One way ANOVA

[illegible]

# One-way ANOVA

- Conditions can be treatments or chemicals, cell types, time points, etc.
- Question being asked is whether the expression level for each gene (taken one at a time) changes significantly as a function of dose.
- More specifically, it compares the variability within replicates for a given dose to the variability caused by changing the dose.
- If gene chip contains 1000 genes, then do 1000 ANOVAs.



# One-way ANOVA

gene	dose 0		dose 50 mM		dose 75 mM		ANOVA
	rep 1	rep 2	rep 1	rep 2	rep 1	rep 2	p value
AA848268_at	163	164	94	165	178	181	0.35
AA848421_at	-91	-13	-102	-18	-97	-73	0.79
AA848546_at	526	498	424	377	539	410	0.29
AA848563_s_at	283	343	3212	4392	1191	1911	0.02
AA859934_at	65	62	6	1	8	71	0.21
AA875509_at	442	415	259	355	265	280	0.07
AA933181_at	-17	38	4	30	29	27	0.78
AA956437_at	75	95	56	73	38	62	0.20
AA958273_at	42	34	-23	18	25	-24	0.35
AA958274_at	114	42	50	41	0	7	0.18

# Multiple Factor ANOVA

$$y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \varepsilon_{ijk}$$

- The model can be extended by adding
  - More factors ( $\alpha_i, \beta_j, \dots$ )
  - Interactions between them ( $\alpha\beta_{ij}, \dots$ )
  - Other ...
- This is used to model the different sources of variation appearing in microarray experiments

# Two-way ANOVA

- Many genes
- 2 Factors, multiple conditions per factor
  - For example, Factor 1 could be dose of a chemical, and Factor 2 could be time point after dosing
  - Can reveal significant effect of time, significant effect of dose, or a significant interaction between the two
- Multiple replicates per condition

# Two-way ANOVA

Gene	Dose 1				Dose 2				...
	Time 1		Time 2		Time 1		Time 2		...
	Rep 1	Rep 2	Rep 1	Rep 2	Rep 1	Rep 2	Rep 1	Rep 2	
A	150	160	150	180	190	180	150	155	
B	50	40	45	50	45	40	80	90	
C	800	760	680	400	450	425	200	220	
...									

# Two-way ANOVA

- Typical Question – Which genes' expression levels change due to time? Due to dose? Due to an interaction between the two?
  - 2-way ANOVA
- Or, eliminate one of the dimensions and ask the same questions as before – At time 1, which doses show similar expression levels among genes?

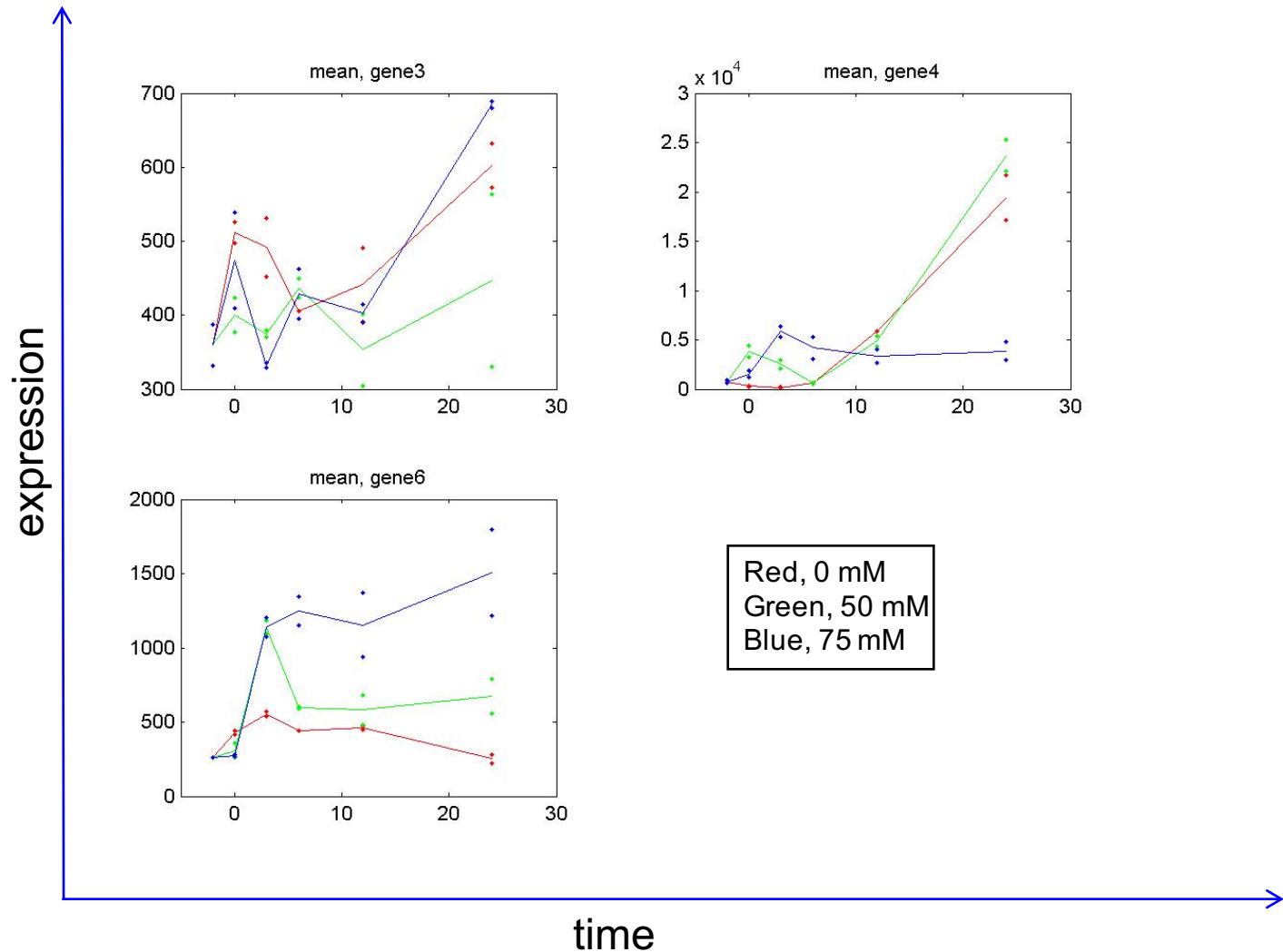
# Two-way ANOVA

Source	SS	df	MS	F	P
Time	28724	4	7181	9.20	7.3e-4
Dose	1143	2	572	0.73	0.498
Time*dose	22940	8	2868	3.67	0.016
Error	10930	14	781		
Total	64409	28			

# Two-way ANOVA

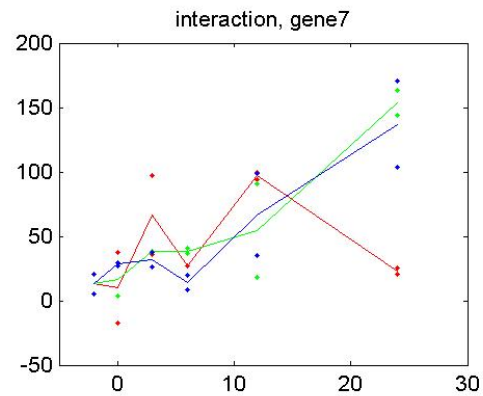
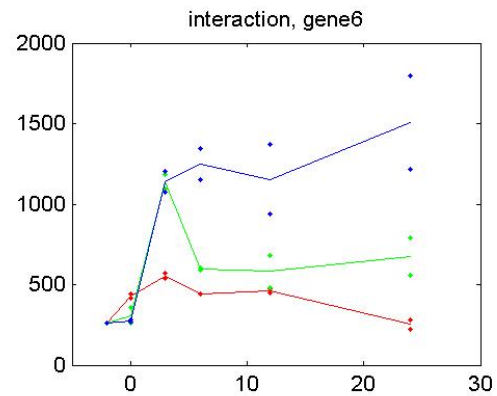
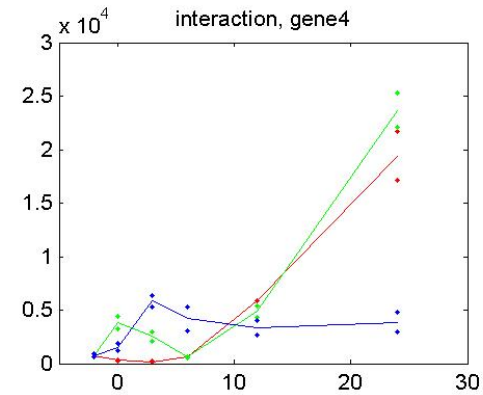
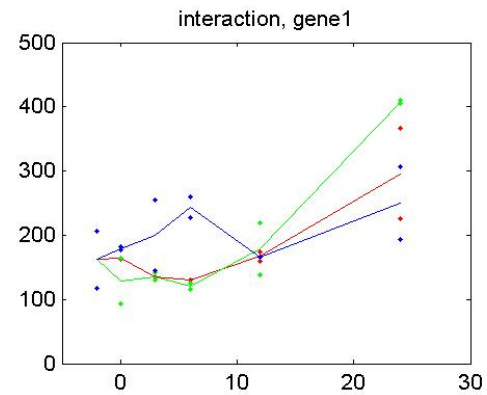
	time	dose	time*dose
AA848268_at	0.000	0.416	0.038
AA848421_at	0.000	0.609	0.913
AA848546_at	0.001	0.024	0.125
AA848563_s_at	0.000	0.000	0.000
AA859934_at	0.206	0.676	0.142
AA875509_at	0.000	0.000	0.001
AA933181_at	0.001	0.498	0.016
AA956437_at	0.434	0.607	0.693
AA958273_at	0.802	0.325	0.283
AA958274_at	0.344	0.242	0.127

# Two-way ANOVA





# Two-way ANOVA



# ANOVA

Let  $y_{ijk g}$  be the fluorescent intensity measured from Array  $i$ , Dye  $j$ , Variety  $k$ , and Gene  $g$ , on the appropriate scale (such as log).  
A typical analysis of variance (ANOVA) model is:

$$y_{ijk g} = \mu + A_i + D_j + V_k + G_g + (AG)_{ig} + (DG)_{jg} + (VG)_{kg} + e_{ijk g}$$

- $\mu, A, D, V$  are “normalization” terms
- $G$  are the overall gene effects
- $AG$ ’s are “spot” effects
- $DG$ ’s are gene-specific dye effects
- $VG$ ’s are the effects of interest
- $e$  is random error

# ANOVA

- ANOVA is a robust statistical procedure
- sources of variation, *e.g.* whether variation in gene expression is less in subset of data than in total data set
- Requires moderate levels of replication (4-10 replicates of each treatment)
- Expression judged according to statistical significance instead of by adopting arbitrary thresholds

# ANOVA

- Advantage: design does not need reference samples
- Concern: treatments should be randomised and all single differences between treatments should be covered

E.g., if male kidney and female liver are contrasted on one set, and female kidney and male liver on another, we cannot state whether gender or tissue type is responsible for expression differences observed

# Basic Statistical Comparisons between groups

- $t$ -Tests & ANOVA
  - Much better than fold change thresholds
  - Need to characterize (or stabilize) variances
  - Assumes normal distribution of expression values. Transformation helps.
- Non-Parametric Tests: Wilcoxon

# Wilcoxon Rank Test

- Another gene score:
  - Ignores absolute values
  - Takes into account only **order** of measurements

- Sort the expression values of both groups

+ + - - + + + - - + - - + + -

**a1 a2 a3 a4 a5 a6 a7 a8 a9 a10 a11 a12 a13 a14 a15**

- $W(g)$  = sum of ranks of the positive examples:

$$W(g) = 1 + 2 + 5 + 6 + 7 + 10 + 13 + 14 = 58$$

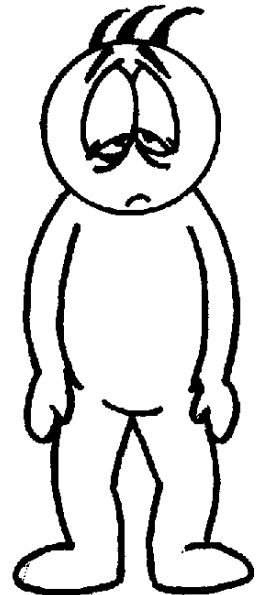
# Wilcoxon Rank Test

- A common test in statistics
- Again, we can compute p-values given the null hypothesis  $H_0$ 
  - $P(W(g) > s | n, k)$  = the probability of getting a score  $> s$  given a total of  $n$  samples, out of which  $k$  are labeled as (+).

# *P* value

**We can rank the genes according to the p-value**

**But, we can't really trust the p-value in a strict statistical way!**





# *P* value

## **Why not!**

For two reasons:

1. We are rarely fulfilling all the assumptions of the statistical test
2. We have to take multi-testing into account

# The *t*-test Assumptions

1. The observations in the two categories must be independent
2. The observations should be normally distributed
3. The sample size must be 'large' (>30 replicates)

# Multi-testing?

In a typical microarray analysis we test thousands of genes

If we use a significance level of 0.05 and we test 1000 genes. We expect 50 genes to be significant by chance

$$1000 \times 0.05 = 50$$

# Multi-testing correction

- ❑ Bonferroni correction
- ❑ False discovery rate

# Permutation testing

- Permutation testing
  - Scramble (randomly permute) the data to randomize relationship of interest, leaving other structure unchanged
  - Do (any) testing procedure on permuted data
  - Repeat  $K$  times.
  - Proportion of times randomly permuted data scores better than actual data is  $p$  value
- As  $K$  increases, becomes very accurate measure of  $p$  value. Makes no assumptions about test, so works with any test.
- Very time consuming (typical  $K = 1000$ )

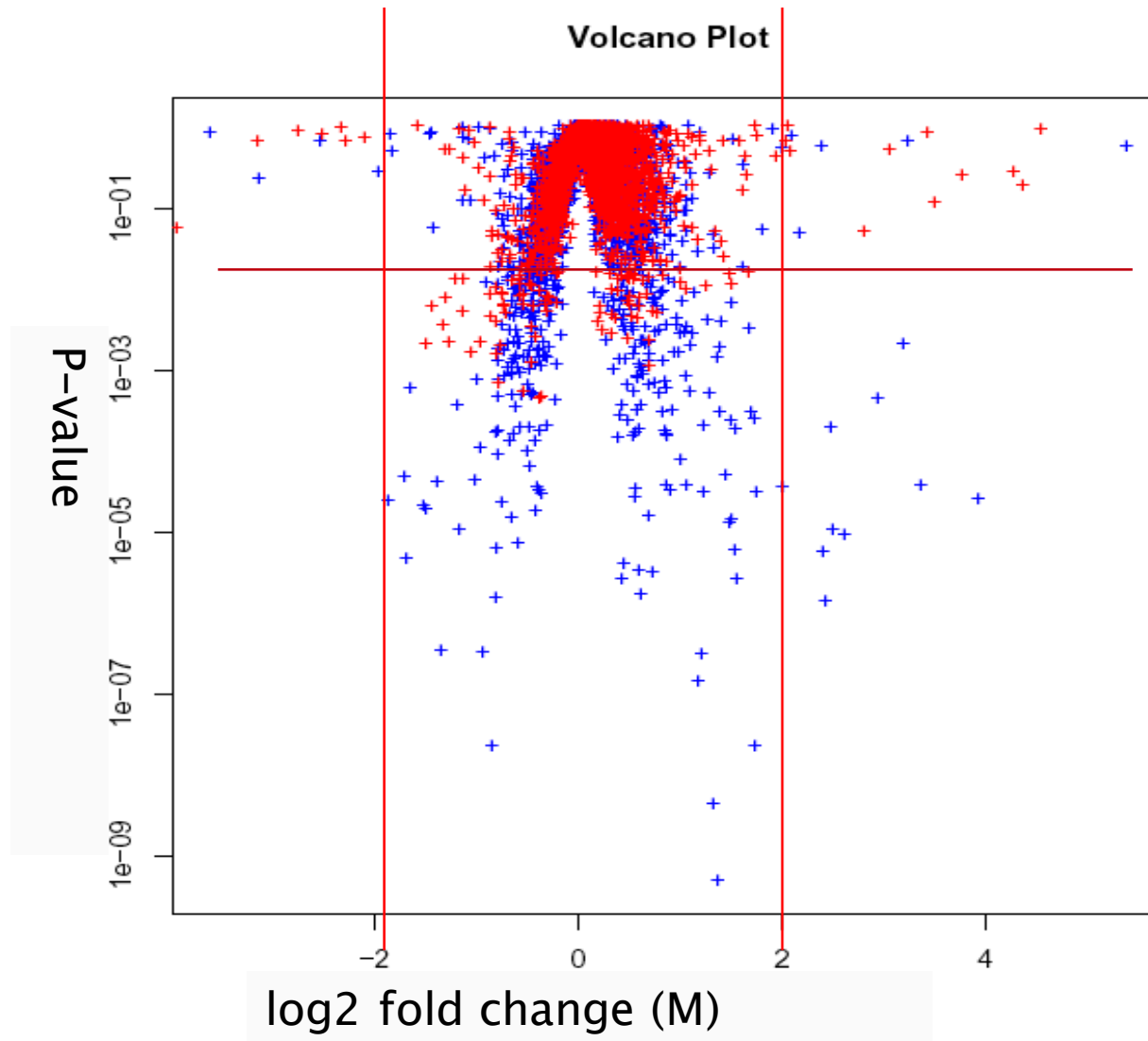
# Consequences of multiple testing corrections

- Multiple testing corrections cost statistical power (probability of correctly rejecting the null when there really is a difference)
  - The more genes tested, the greater the loss
- Therefore, eliminate as many genes as possible from consideration before testing for significant differences.
  - E.g. by eliminating genes that didn't change at all on any chips, regardless of experimental conditions.

# Combining $p$ -values and fold changes

- What's important biologically?
  - How significant is the difference?
  - How large is the difference?
- Both amounts can be used to identify genes.
- What cutoffs to use?
- How many genes should be selected?
- Where are your positive controls?

# Combining $p$ -values and fold changes





# Signal-to-Noise

Signal-to-Noise (S2N) filter used to select 50 most variable genes out of 19,892 genes on the array.

$$S2N(j) = \left| \frac{\mu_1(j) - \mu_2(j)}{\sigma_1(j) + \sigma_2(j)} \right|$$

where  $\mu$  and  $\sigma$  indicate means and std. dev. of expression levels of gene  $j$  for case and control groups.

# Gene set analysis

# Enrichment analysis

Long lists of differentially expressed genes

What happens next?

- Select some genes for validation?
- Do follow-up experiments on some genes?
- Publish a huge table with the results?
- Try to learn about all the genes on the list (read 100s of papers)?
- ....

Usually, some or all of the above will be done, and more.

Can we help further at this?

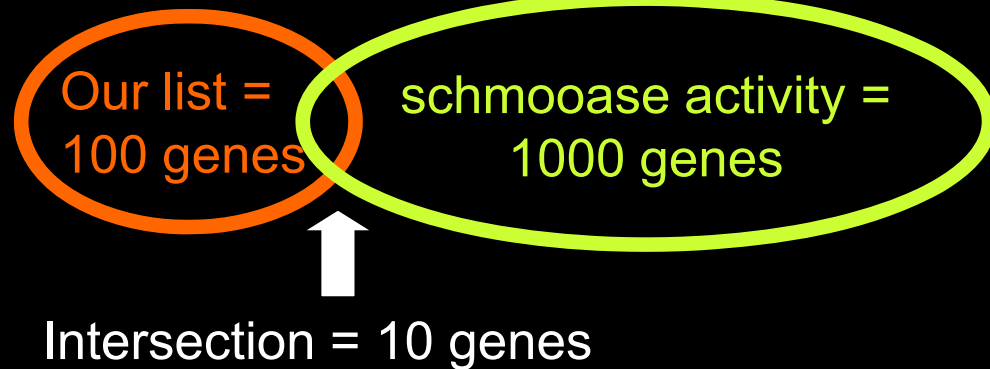
# Annotation sources

- Gene Ontology(most popular)
  - biological process, molecular function, cellular component
  - Terms may have  $>1$  “parent” (more general term)
  - GO Slim: includes only general categories
- KEGG; REACTOME pathways
- Genes sharing a motif regulated by the same protein
- Genes found on the same chromosome
- ...

# Why enrichment analysis?

- Too many genes to examine in detail
- Are we biased?
- How do we know that what we're seeing is by chance?

Genome =  
20,000 genes



# Main types of enrichment analysis

## List-based:

- A subset of all genes chosen by some relevant method
- A list of annotations, each linked to genes

## • Rank-based:

- A set of all genes ranked by some metric (ratio, fold change, etc.)
- A list of annotations, each linked to genes

## • List-based with relationships:

- A subset of all genes
- A list of annotations, each linked to genes, organized in some relationship (e.g., a hierarchy)

# Statistics to test for enrichment

Genome =  
20,000 genes

Our list =  
100 genes

schmooase activity =  
1000 genes

10%

5%

Intersection = 10 genes      $p=0.03$

Our list =  
1000 genes

stroumphase activity =  
20 genes

0.1%

0.2%

Intersection = 2 genes      $p=0.3$

# Statistics to test for enrichment

- What is the chance of observing enrichment due to chance?
- Different tests produce very different ranges of p-values
- All look for over-enrichment; some look for under-enrichment
- Recommendation: Use p-values as a tool to rank genes but don't take them literally
- Most methods correct for multiple testing (e.g., with FDR), which is necessary



# Statistics to test for enrichment

- Fisher's exact
- Hypergeometric
- Binomial
- Chi-squared
- Z
- Kolmogorov-Smirnov
- Permutation
- .....

# Fisher's test by hand in R

- `counts = (matrix(data = c(3, 297, 40, 19960), nrow = 2))`
- `counts`
- `fisher.test(counts)`
- `#` is better than
- `chisq.test(counts)`

|                   | Gene list | Genome |
|-------------------|-----------|--------|
| In anno group     | 3         | 40     |
| Not in anno group | 297       | 19960  |

## Fisher's Exact Test for Count Data

```
data: counts
p-value = 0.02552
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
 0.9918169 15.9604612
sample estimates:
odds ratio
5.039206
```

$(3/297) / (40/19960)$

# Hypergeometric test by hand in R

- `min(1 - cumsum(dhyper(0:(3-1), 40, 19960, 300) ))`
- 0.02193491

|                   | Gene list | Genome |
|-------------------|-----------|--------|
| In anno group     | 3         | 40     |
| Not in anno group | 297       | 19960  |

- Equation above tests only for over-enrichment

# The Gene Ontology Consortium

<http://www.geneontology.org>

The goal of the Gene Ontology™ (GO) Consortium is to produce a controlled vocabulary that can be applied to all organisms even as knowledge of gene and protein roles in cells is accumulating and changing. GO provides three structured networks of defined terms to describe gene product attributes

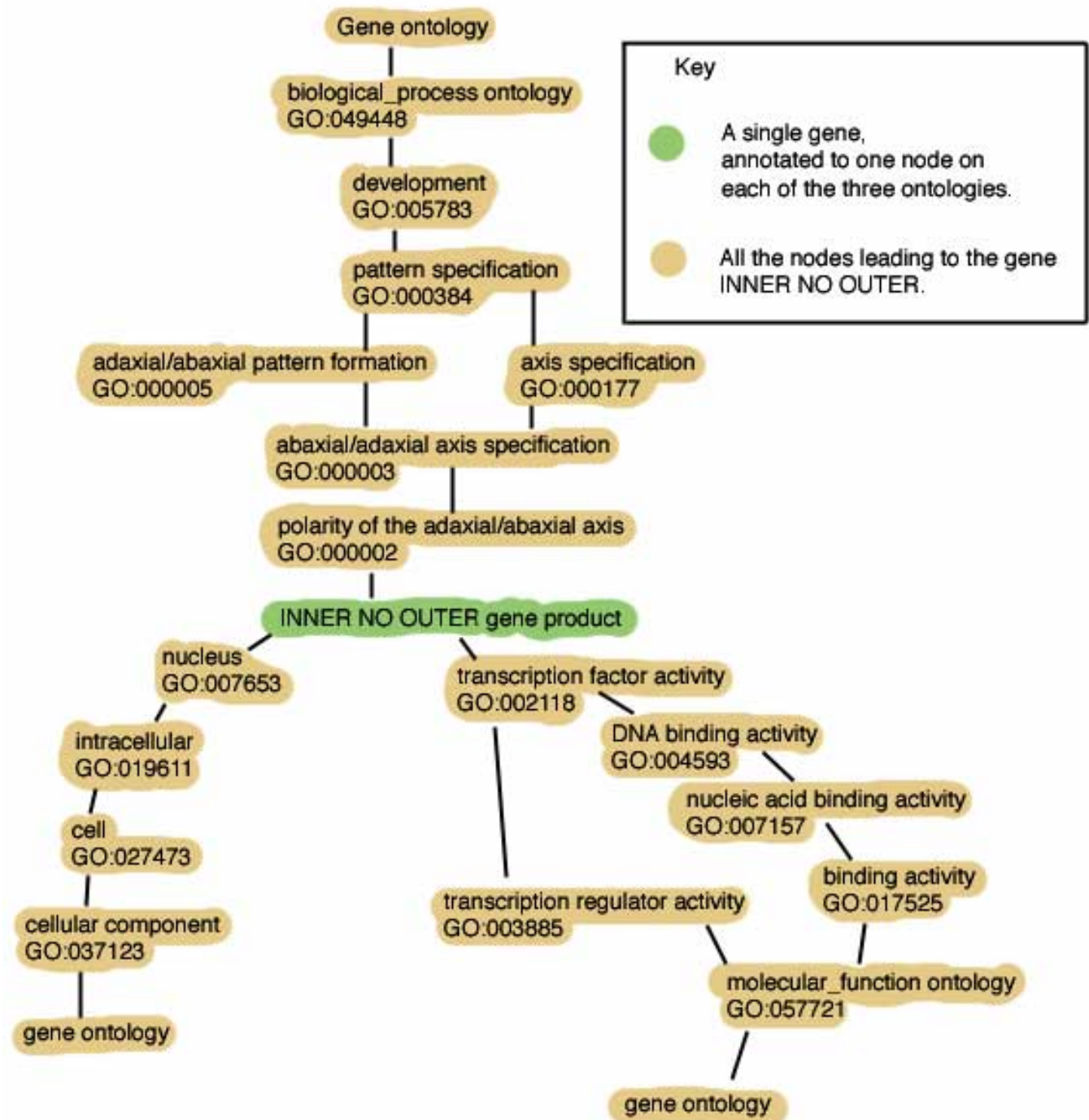
**Molecular Function Ontology:** the tasks performed by individual gene products; examples are *carbohydrate binding* and *ATPase activity*

**Biological Process Ontology:** broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

**Cellular Component Ontology:** subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *origin recognition complex*

From the GO web site. The *path* back to each ontology from a gene.

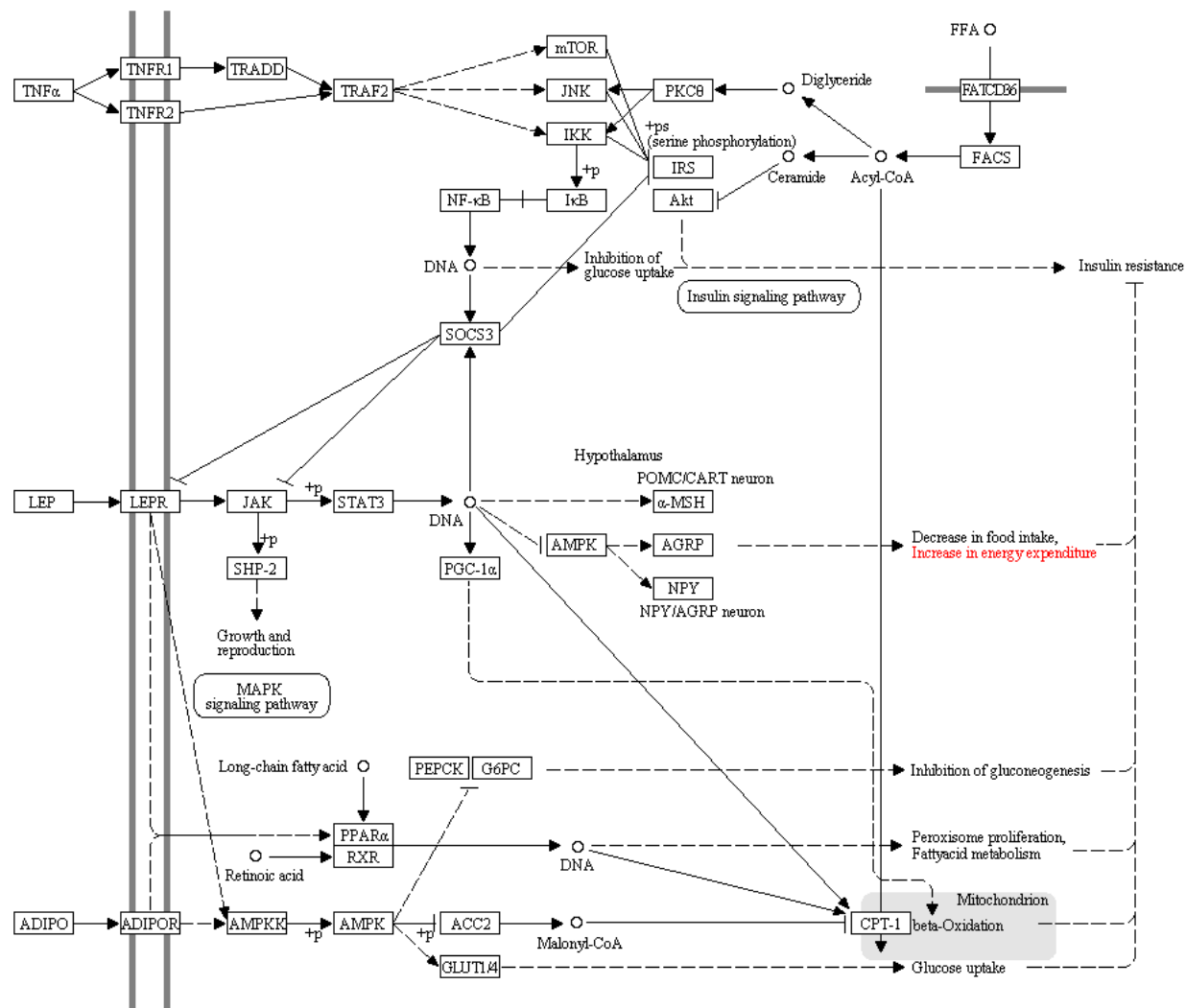
We will call each term in a path a *split*.



# KEGG pathway database

- KEGG = Kyoto Encyclopedia of Genes and Genomes
  - <http://www.genome.jp/kegg/pathway.html>
  - The pathway database gives far more detailed information than GO
    - Relationships between genes and gene products
  - But: this detailed information is only available for selected organisms and processes
  - Example: Adipocytokine signaling pathway

# ADIPOCYTOKINE SIGNALING PATHWAY



- MIT, Broad Institute
- V 2.0 available since Jan 2007



(Subramanian et al. PNAS. 2005.)

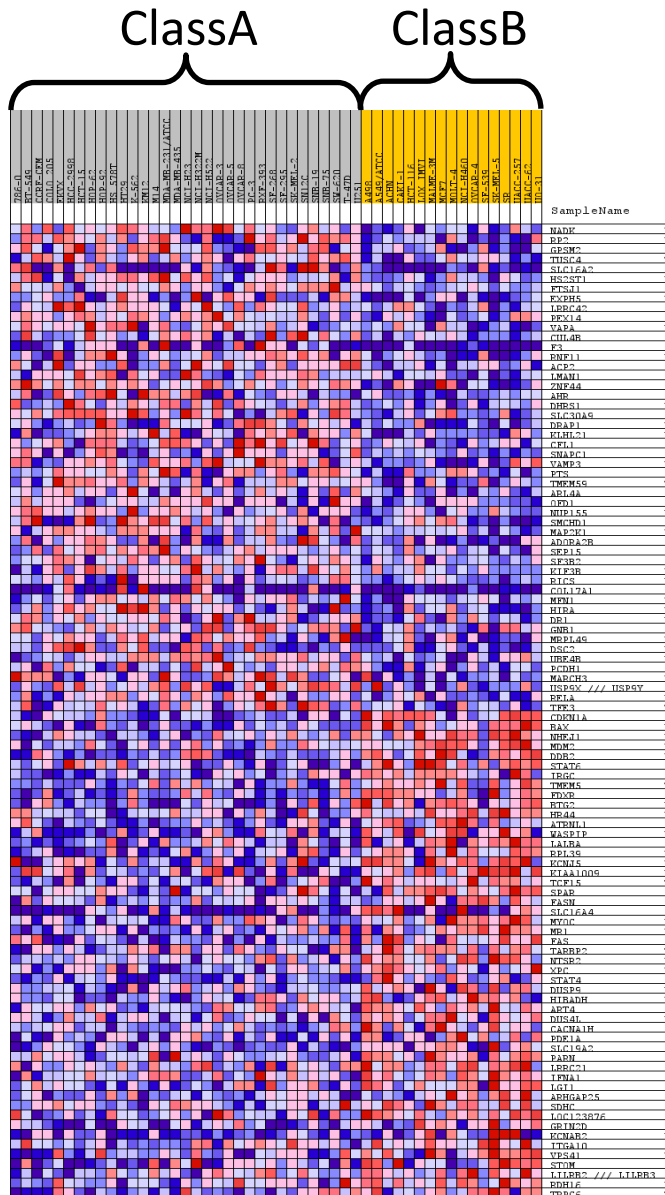


# GSEA key features

- Calculates a score for the enrichment of a **entire set of genes** rather than single genes!
- Does not require setting a cutoff!
- Identifies the set of relevant genes as part of the analysis!
- Provides a more robust statistical framework!

# Gene Set Enrichment Analysis

Genes ranked by expression correlation to Class A



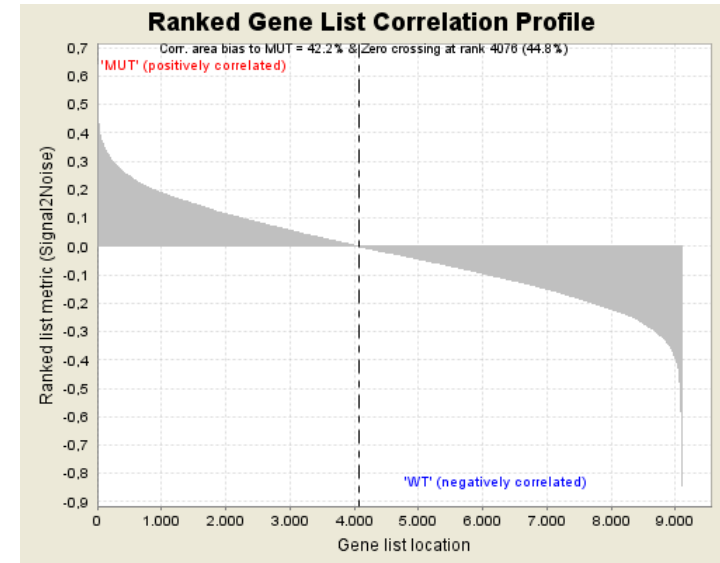
Function 1  
(e.g., metabolism)



Function 2  
(e.g., signaling)



Function 3  
(e.g., regulation)



**Running sum:**  
Increase when gene is in set  
Decrease otherwise

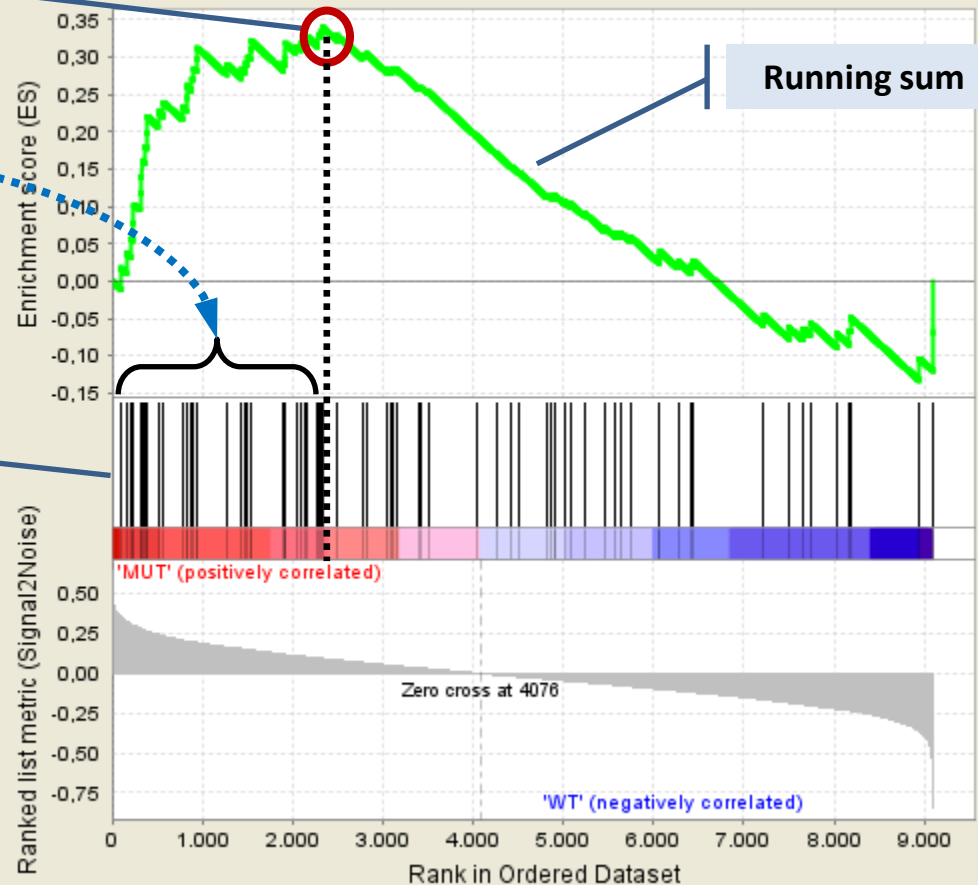
# Gene Set Enrichment Analysis

Enrichment score (ES) =  
max deviation from 0

Leading  
Edge genes

Genes within  
functional set  
(hits)

Enrichment plot: CELL\_CYCLE\_KEGG



Running sum

'MUT' (positively correlated)

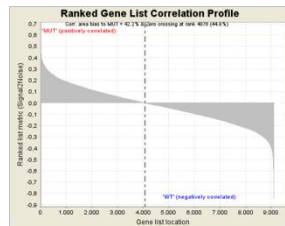
Zero cross at 4078

'WT' (negatively correlated)

Enrichment profile — Hits — Ranking metric scores

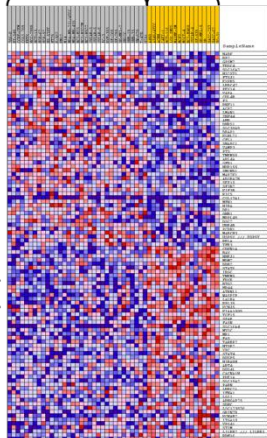
ClassA ClassB

Function 1  
Function 2  
Function 3



Running sum:  
Increase when gene is in set  
Decrease otherwise

Genes ranked by expression correlation to Class A



# Other statistical issues

- Goal: Identifying theme(s) of maximal biological significance
  - but this is not perfectly correlated with statistical significance
- What is your background gene set?
  - All genes that could appear in your list
- What about sparse annotation groups?
- Some annotation terms may be subsets of other terms.

# Some recommended tools

- DAVID
- GSEA
- BIOBASE (Whitehead has license)
- BiNGO (uses Cytoscape)
- GoMiner:  
<http://discover.nci.nih.gov/gominer>
- GOstat: <http://gostat.wehi.edu.au>

# **Mining gene expression**

# Mining gene expression

- After differential analysis, what is next?
  - Genes coregulated
  - Gene network
  - Genes involved in the same biological process

# Clustering

- Cluster analysis: dividing samples (genes) into homogeneous groups based on a set of features
- Steps: generate expression summaries  
measure pairwise distances  
cluster

## Distance (semi)metrics for pairwise measurements

Euclidean distance

$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

Manhattan distance

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

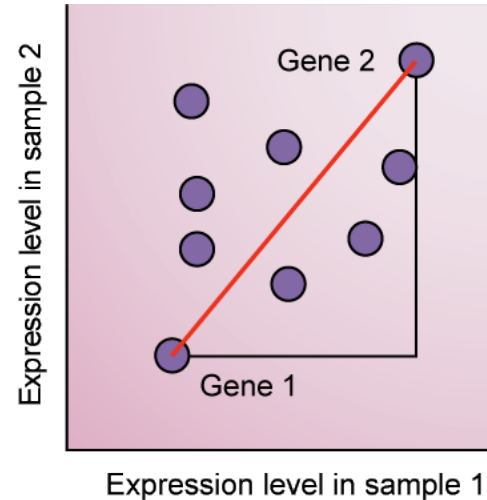
Correlation coefficient

$$r = \frac{\sum_{i=1}^p (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^p (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^p (y_i - \bar{y})^2}}$$

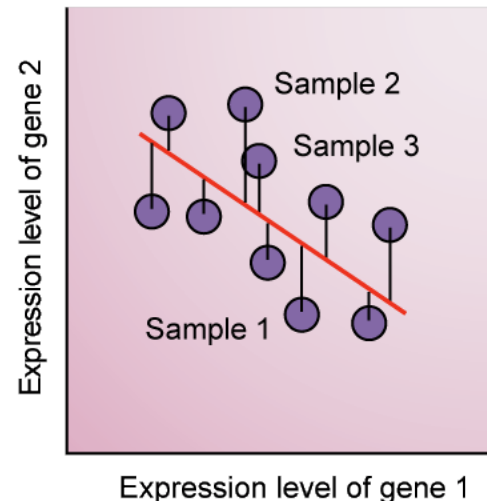


# Metrics for gene expression

- Need a method to measure how similar genes are based on expression
- Examples
  - Euclidean distance
  - Pearson correlation coefficient



**Euclidean distance**

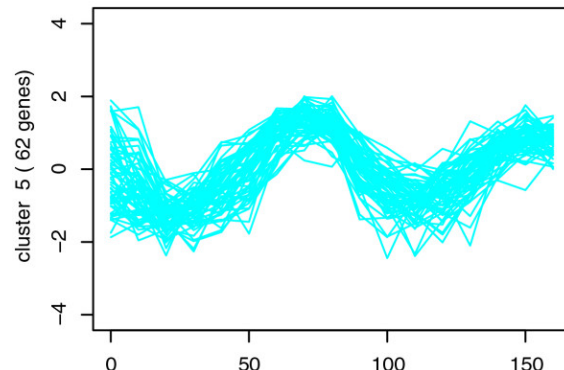
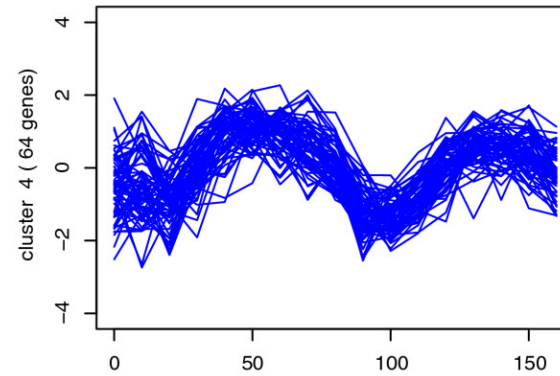
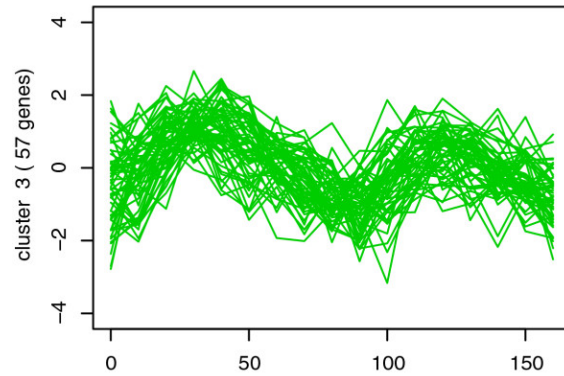
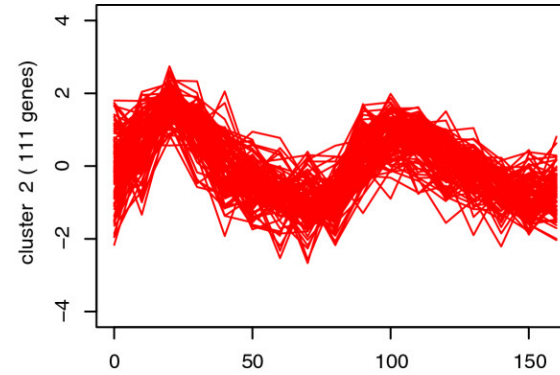
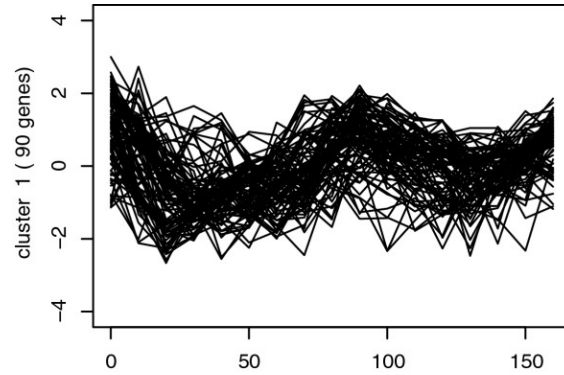


**Pearson correlation coefficient**

# Clustering

- make no assumptions about how the data should behave
- Typical approaches includes
  - hierarchical clustering
  - $k$ -means clustering
  - self-organizing maps (SOM)
  - principal component analysis (PCA)

# Clustering

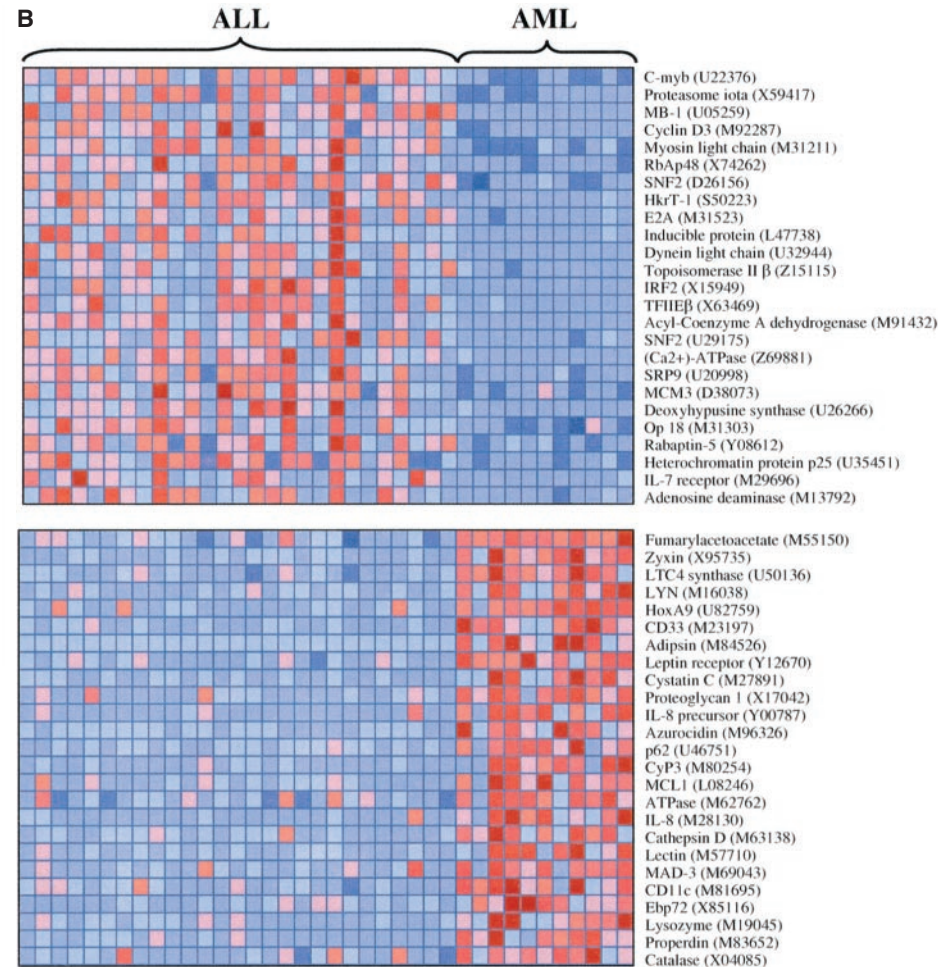
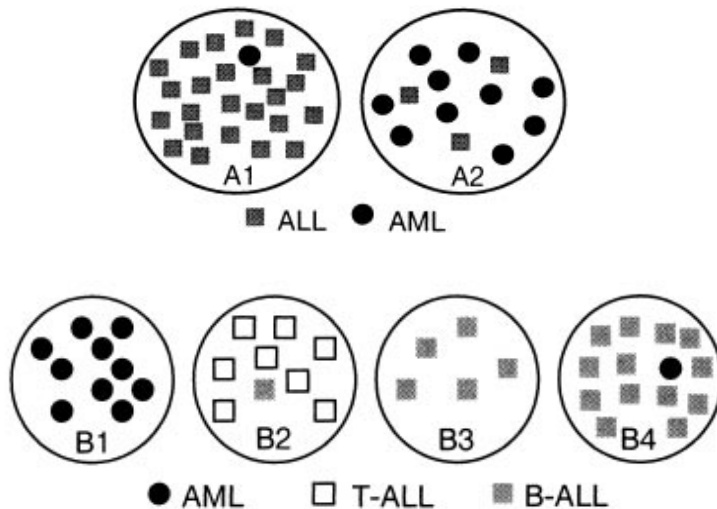


# Clustering

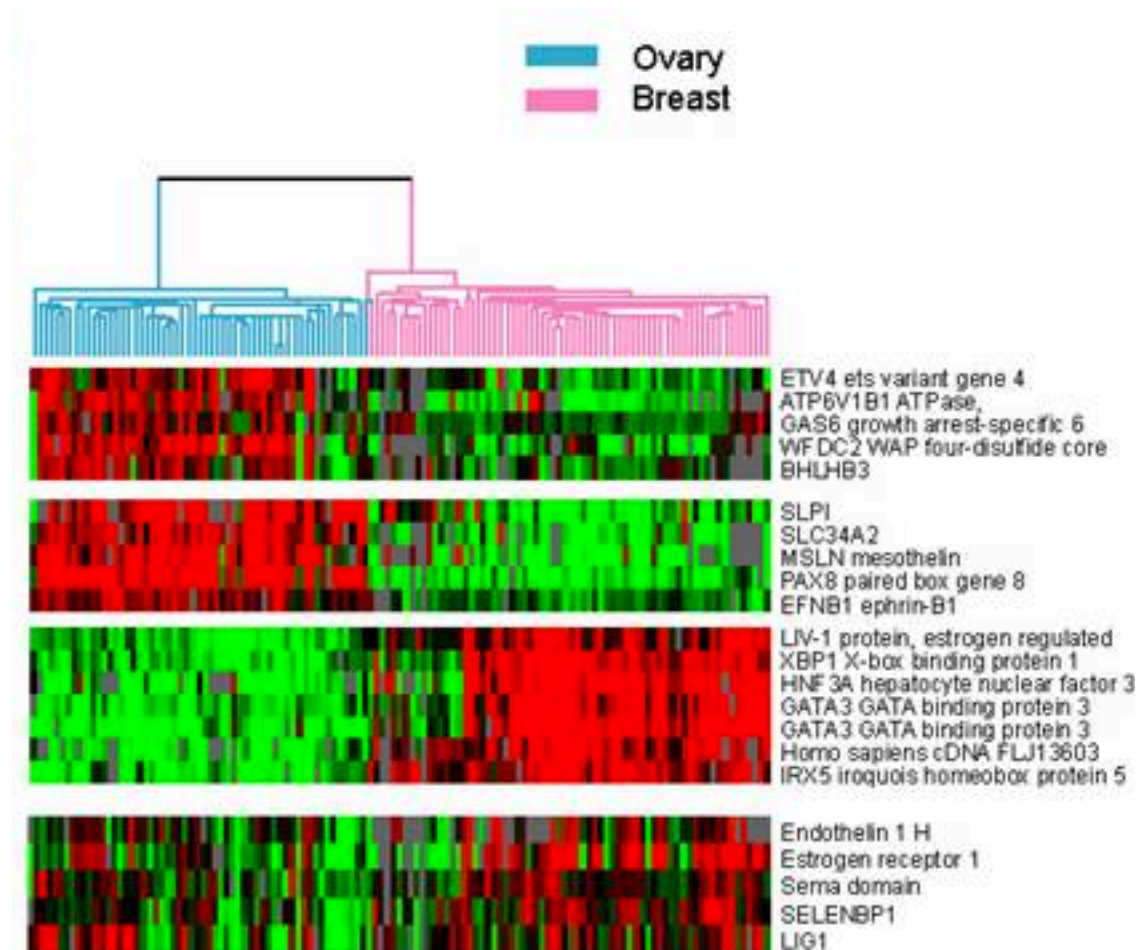
REPORTS

## Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub,<sup>1,2\*†</sup> D. K. Slonim,<sup>1†</sup> P. Tamayo,<sup>1</sup> C. Huard,<sup>1</sup>  
M. Gaasenbeek,<sup>1</sup> J. P. Mesirov,<sup>1</sup> H. Coller,<sup>1</sup> M. L. Loh,<sup>2</sup>  
J. R. Downing,<sup>3</sup> M. A. Caligiuri,<sup>4</sup> C. D. Bloomfield,<sup>4</sup>  
E. S. Lander<sup>1,5\*</sup>



# Hierarchical Clustering



# Clustering

- Different clustering methods will give you different views of the data
- There is no “correct” clustering method—clustering is just a guide that helps you to see the data in different ways
- On the other hand, there are “incorrect” methods—a certain amount of mathematical validation should be undertaken
- It’s often worth trying a variety of clustering methods to see what is the most useful for your purposes

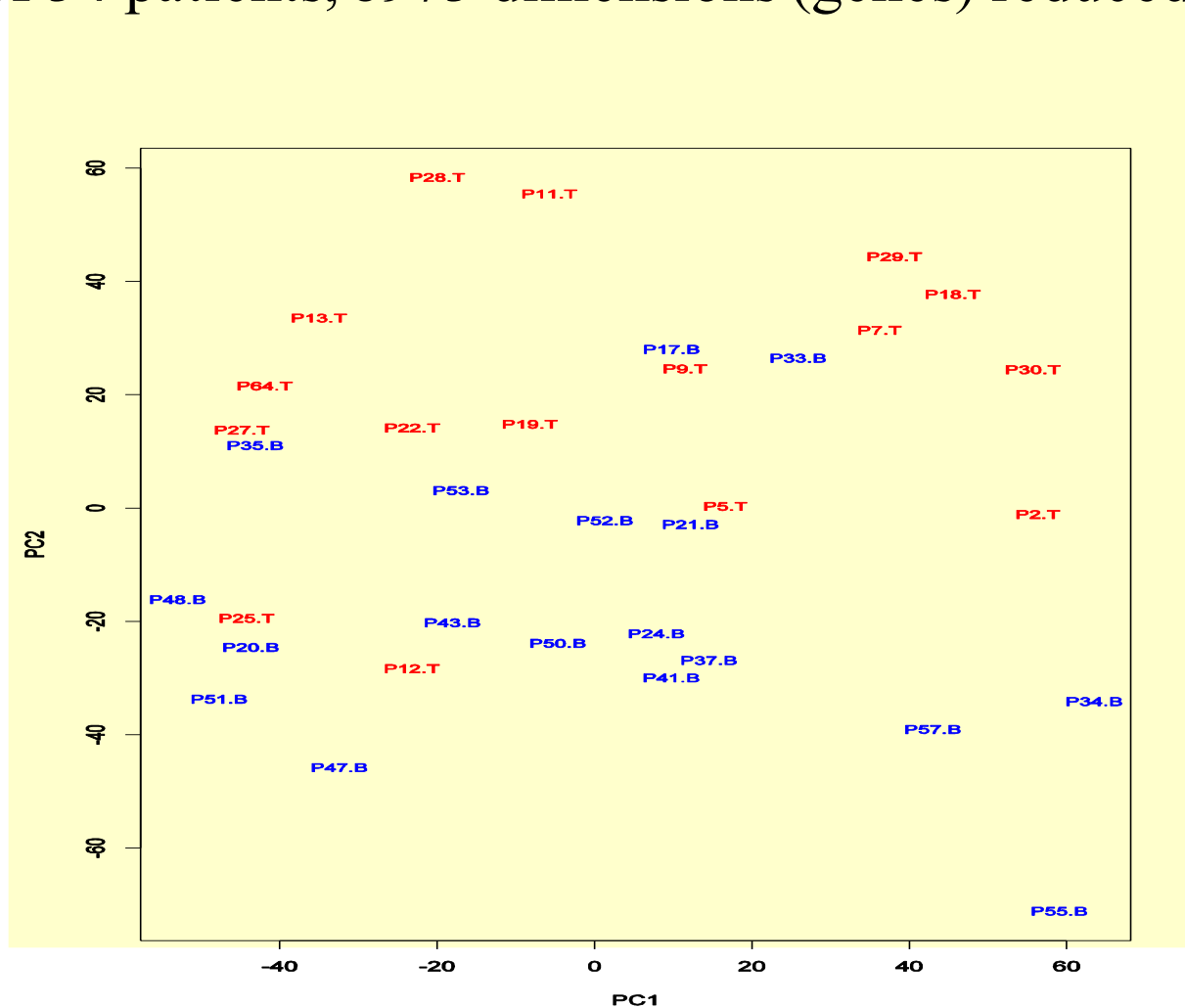
# Principal Component Analysis

- Reduces the dimensionality of the data set
  - Thousands of genes are combined in a few linear combinations to make 2 or 3 Principal Components (PC).  
Going from thousands of axes, with each axis representing the expression level for a gene, to 2 or 3 axes.
- These few PCs may capture most of the variability of the original data set
- Hope is that the first few PCs extract or expose the cluster structure of the original data set
  - i.e. Another clustering algorithm still needed after PCA

# PCA on leukemia data

## (precursor B and T)

- Plot of 34 patients, 8973 dimensions (genes) reduced to 2





# Student t-test can be used to filter out differential expressed genes

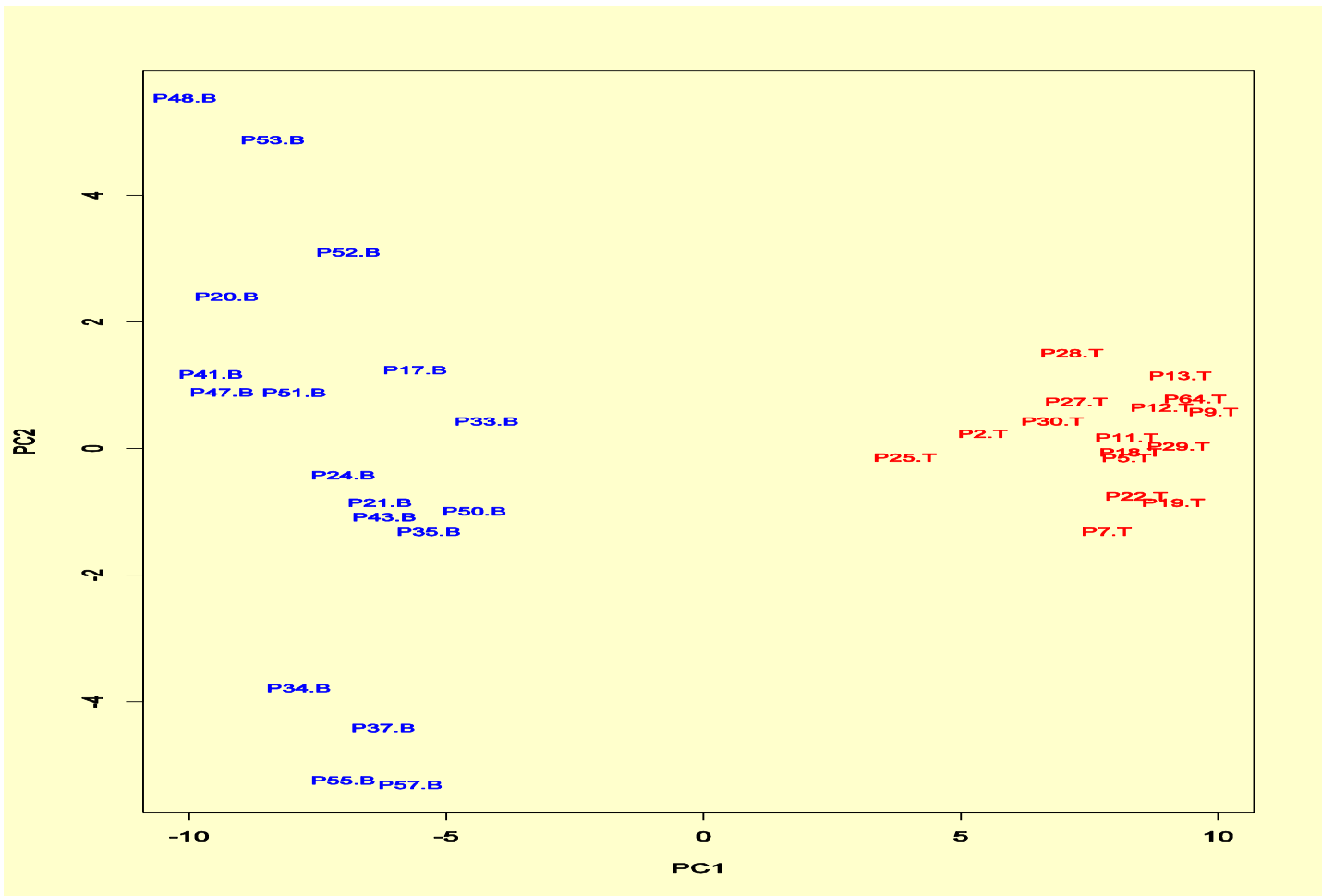
- Compares the means ( $\bar{x}_1$  &  $\bar{x}_2$ ) of two data sets
  - tells us if they can be assumed to be equal
- Can be used to identify significant genes
  - *i.e.* those that change their expression *a lot!*

$$t = \frac{\bar{x}_1 - \bar{x}_2}{SE}, \quad SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

# PCA on leukemia data

## (precursor B and T)

- Plot of 34 patients, 100 dimensions (genes) reduced to 2



# Clustering

Limitation of cluster analysis:

similarity in expression pattern suggests  
co-regulation but doesn't reveal cause-  
effect relationships

