

Gene expression

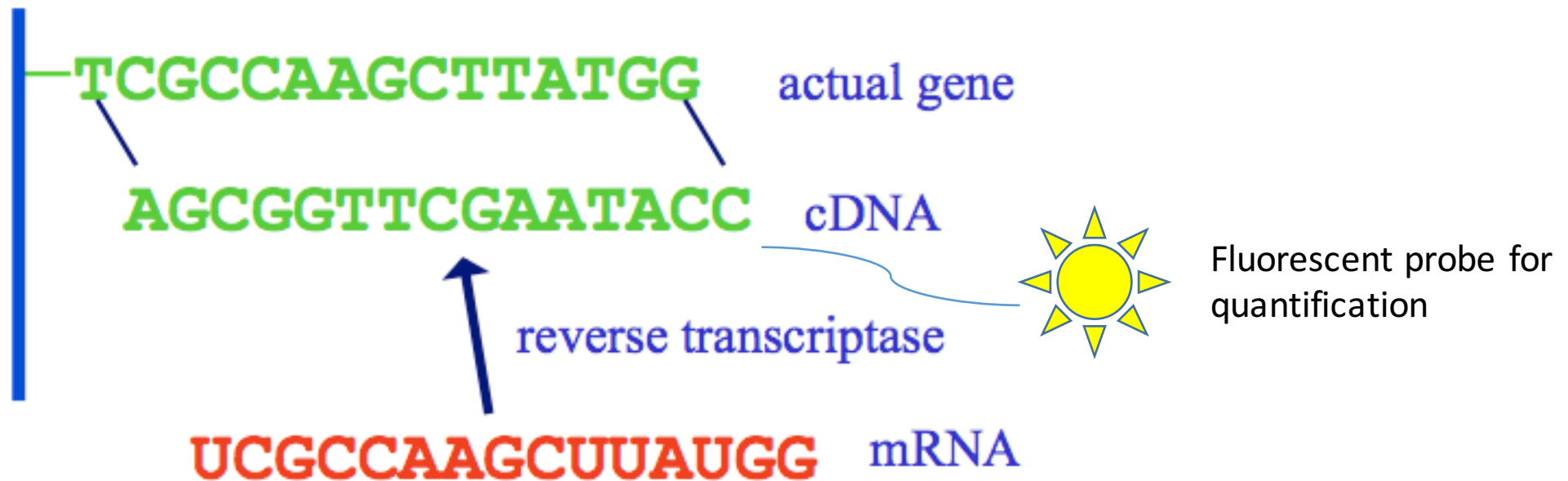
# Microarrays

- Each cell has the same genome but doesn't use it in the same way
- We can predict structural protein features from sequence and assign molecular function but some questions are difficult to answer looking at sequence
  - Sequence: this is a kinase with a an SH2 domain
  - Functional genomics: What tissue/organ/condition is this gene expressed in?
- Given the sequence of genes in the genome we can measure their simultaneous activity in a sample of interest
- Possible questions
  - What genes are different between cancer and normal tissue
  - What genes are required for response to ionizing radiation

# How it works

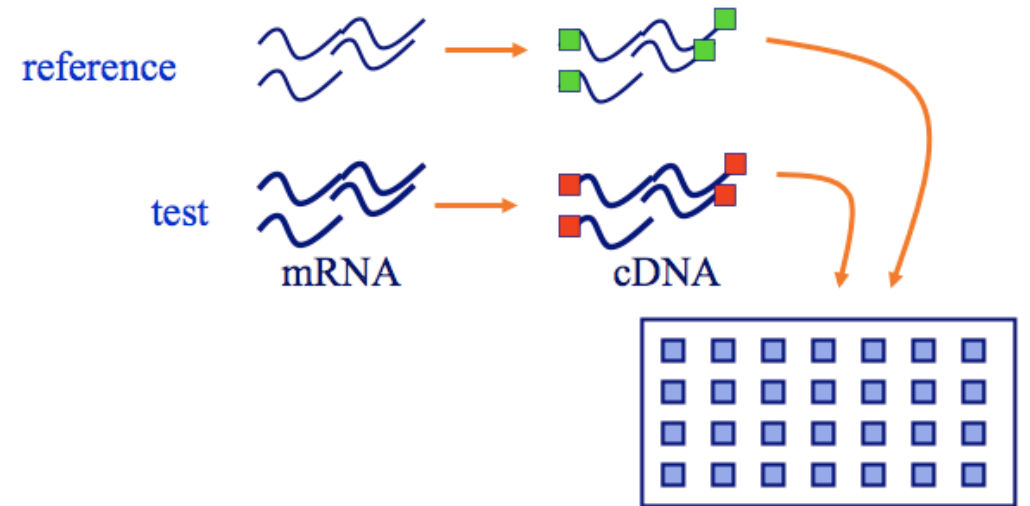
Complementary hybridization:

- Put a part of the gene sequence on the array
- convert mRNA to cDNA using reverse transcriptase



# Type of arrays

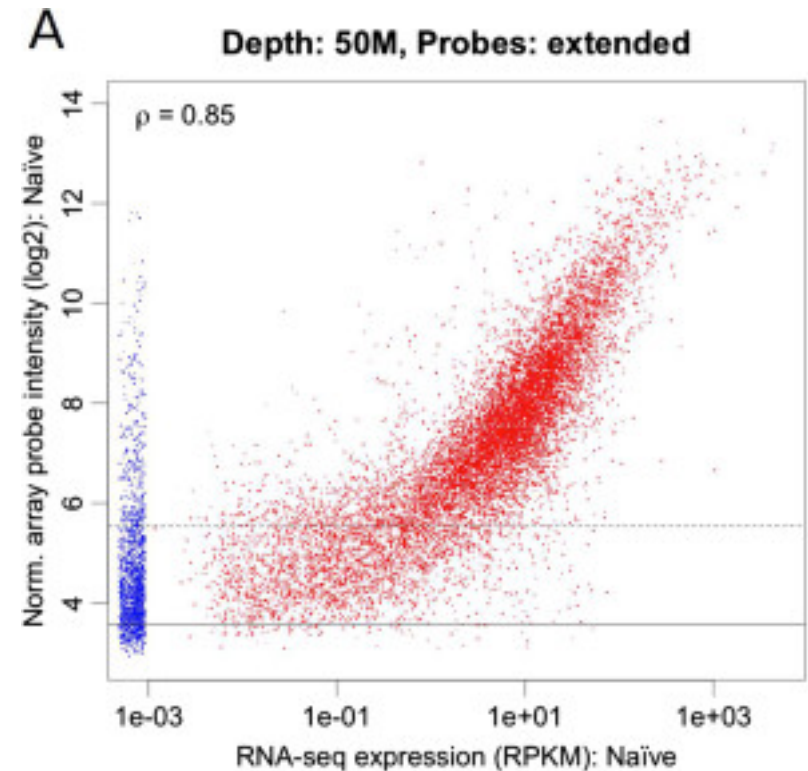
- Spotted (old) – probes are synthesized and then deposited
- Oligonucleotide – probes are synthesized in place
- 2-channel
  - Two mRNA samples (reference, test) are labeled with fluorescent dyes (Cy3, Cy5) and allowed to hybridize to array
  - No comparisons across probes
- Single channel
  - One sample is hybridized
  - Intensity is related to total abundance
- Most arrays today are single channel oligonucleotide





# RNAseq vs microarray

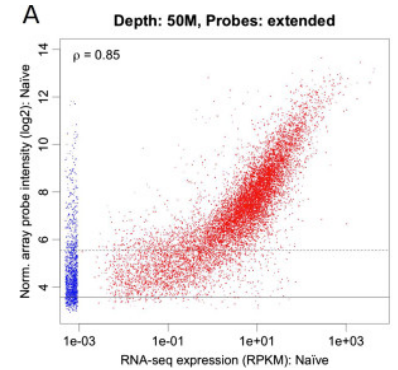
- RNAseq-direct sequencing of mRNA
  - Don't need to know what you are looking for
  - No probes
  - More certainty that you are detecting specific genes
  - Not based on fluorescent read out-better dynamic range
- Microarray vs RNAseq
  - Transcript misidentification
  - Saturation – low and high end



A comparison of RNA-seq and exon genome transcription profiling of the L5 spinal nerve transection model of neuropathic pain in the rat arrays for whole

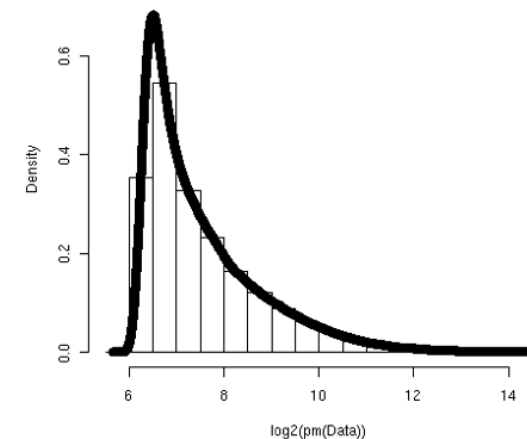
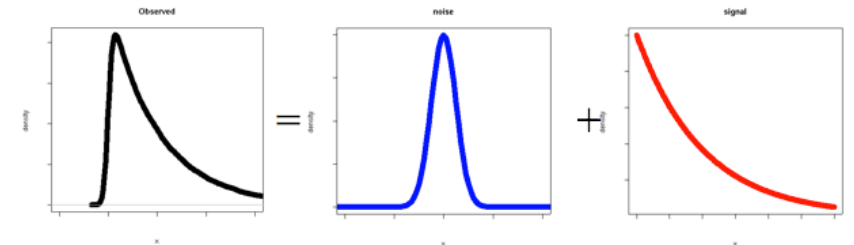
# Microarray background estimation

- Past some intensity point detection is not reliable
- General rule: for mammals about half of all possible genes are expressed in any given tissue/organ
- Use distribution characteristics to filter out unreliable measurements
- Signal intensity is modeled as a convolution of a normal and exponential distribution
- Illumina beadArray chips provide a detection p-value
  - These are often very close to a distribution based estimate



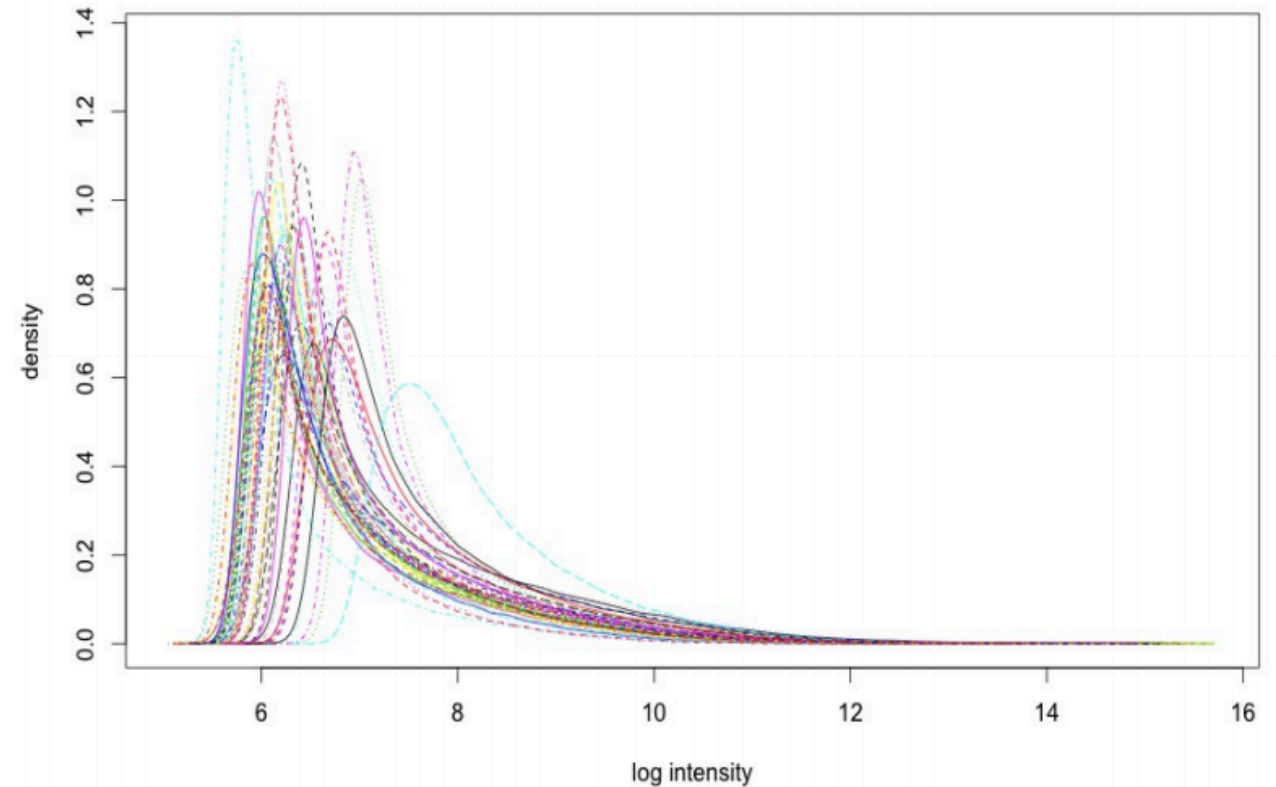
$$\text{Intensity} = \text{Background} + \text{Signal}$$

$N(\mu, \sigma^2)$        $\text{Exponential}(\alpha)$



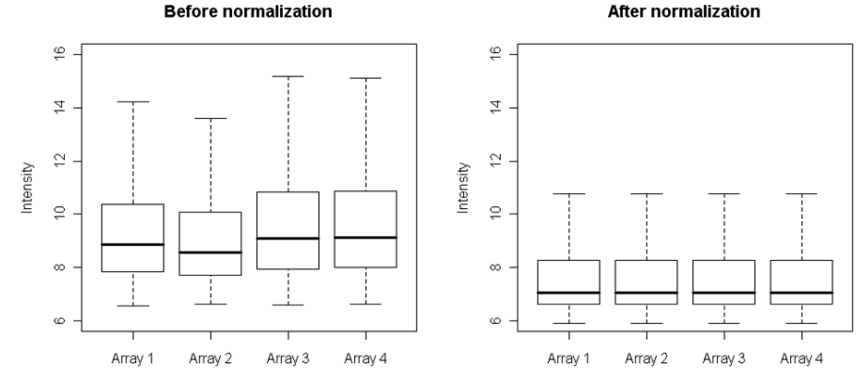
# Normalization

- We want to measure mRNA abundance but we measure fluorescence intensity
- Intensity is related through abundance through some arbitrary function
- This function depends on many experimental parameters and is different for different samples
- We have:  $A_i = F_i(I_i)$
- Ideally we would like:  $A_i = F(I_i)$ —all the functions are the same—though they still don't report true intensity

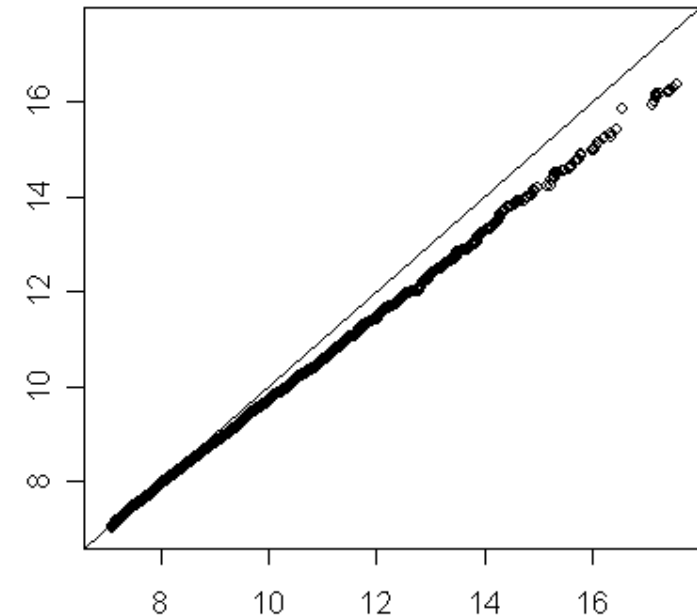
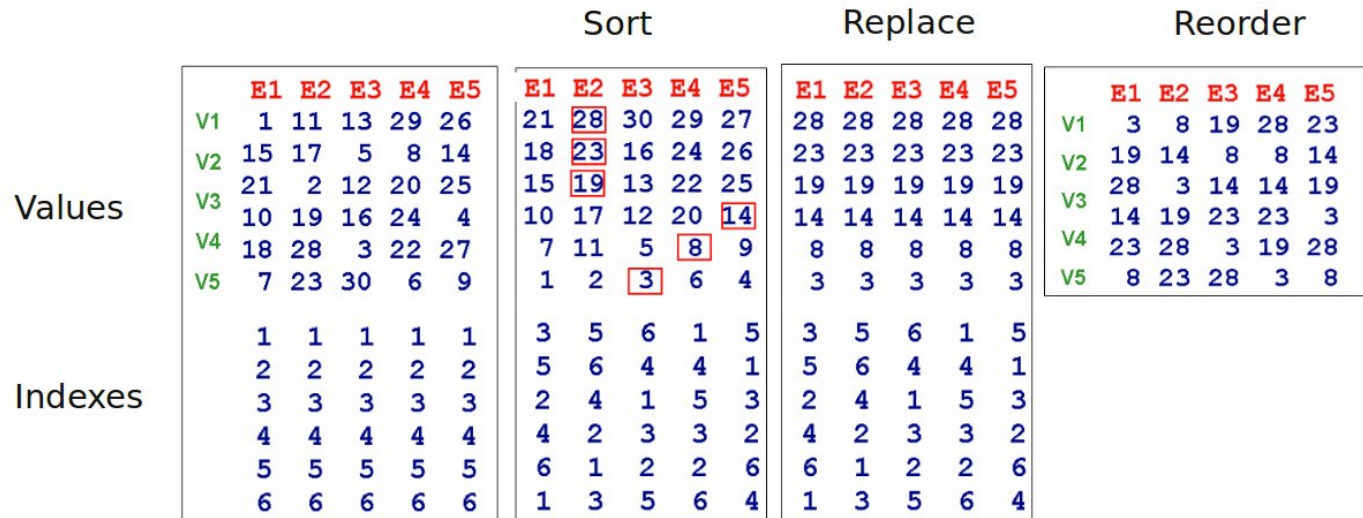


# Normalization

- Many methods have been proposed
- Most widely used is **quantile normalization**
- Force the distributions to be the same by assigning the gene in each rank to the median value in that rank (not necessarily the same gene) across samples



Quantile-quantile (QQ plot)



# Quantile normalization

- Assumption—abundance distributions are the same
  - May be very far from true for different tissues!
  - Not true in general
- Sophisticated methods can normalize just a subset truly equivalent genes
- Can be applied to other datatypes
- Works best when your set of measurements is large and complete—  
not preselected to test a specific hypothesis

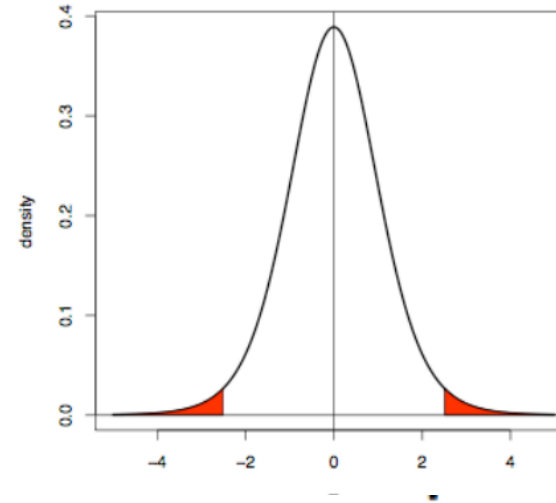
# Quantile normalization—biggest assumption

- Mapping of abundance to intensity is monotone! –intensity value depends only on the true abundance
- This is not true and is sample dependent
- One important factor: GC content of sequence which affects:
  - cDNA synthesis
  - Hybridization kinetics
- Non monotonicity factors must be known and modeled explicitly
- Many methods model GC content

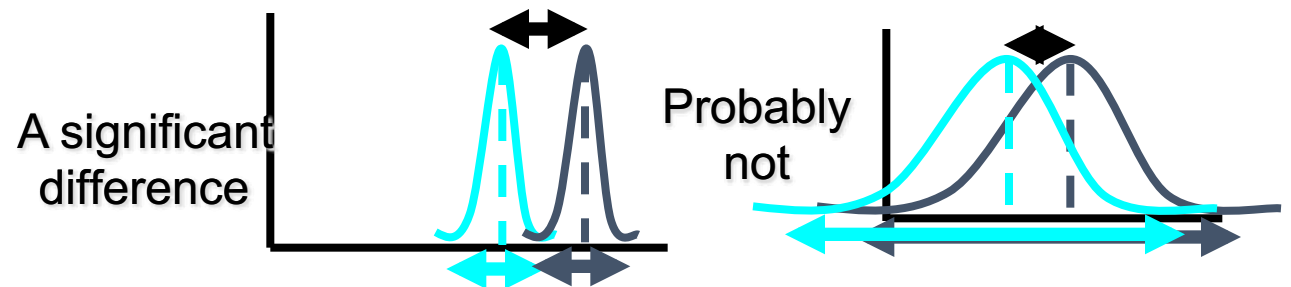
# Statistical inference

- Now that the data looks good what can we say about the biology
- Simplest experimental design 2 groups
- T-test—signal to noise ratio
- Has T distribution when the two means are actually equal
- Assign p-value –small p-values means the T statistic was very unlikely for equal means
- We did an experiment and found 150 genes are differentially expressed with a p-value<0.005
- Is this a good result?

We measured 30,000 genes



$$T_g = \frac{\bar{X}_{g1} - \bar{X}_{g2}}{S_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$



# Type I/ Type II Error

Your Statistical Decision	True state of null hypothesis	
	$H_0$ True No difference between means	$H_0$ False Difference between means
Reject $H_0$ (Conclude that samples are different)	<i>Type I error (<math>\alpha</math>) BIG MISTAKE</i>	<i>Correct</i>
Do not reject $H_0$ (ex: you conclude that there is insufficient evidence that the samples are different))	<i>Correct</i>	<i>Type II Error (<math>\beta</math>)</i>



# Testing many hypothesis at once: Error rates

- Per-family Error Rate

$$\text{PFER} = E(V)$$

- Per-comparison Error Rate

$$\text{PCER} = E(V)/m$$

- Family-wise Error Rate

$$\text{FWER} = p(V \geq 1)$$

- False Discovery Rate

$$\text{FDR} = E(Q), \text{ where}$$

$$Q = V/R \text{ if } R > 0;$$

$$Q = 0 \text{ if } R = 0$$

Decision \ Truth	# true H	# non-true H	totals
# rejected	V (Type I/big mistake)	S	R
# not rejected	U	T	m - R
totals	m <sub>0</sub>	m <sub>1</sub>	m

# Adjusted p-values

- If interest is in controlling, e.g., the FWER, the **adjusted p-value** for hypothesis  $H_j$  is:

$$p_j^* = \inf \{ \alpha : H_j \text{ is rejected at FWER } \alpha \}$$

- Hypothesis  $H_j$  is rejected at FWER  $\alpha$  if  $p_j^* \leq \alpha$

# Correction procedures

- $m$  is the total number of tests
- **Bonferroni single-step** adjusted p-values –controls FWER—probability of making at least one Type I error

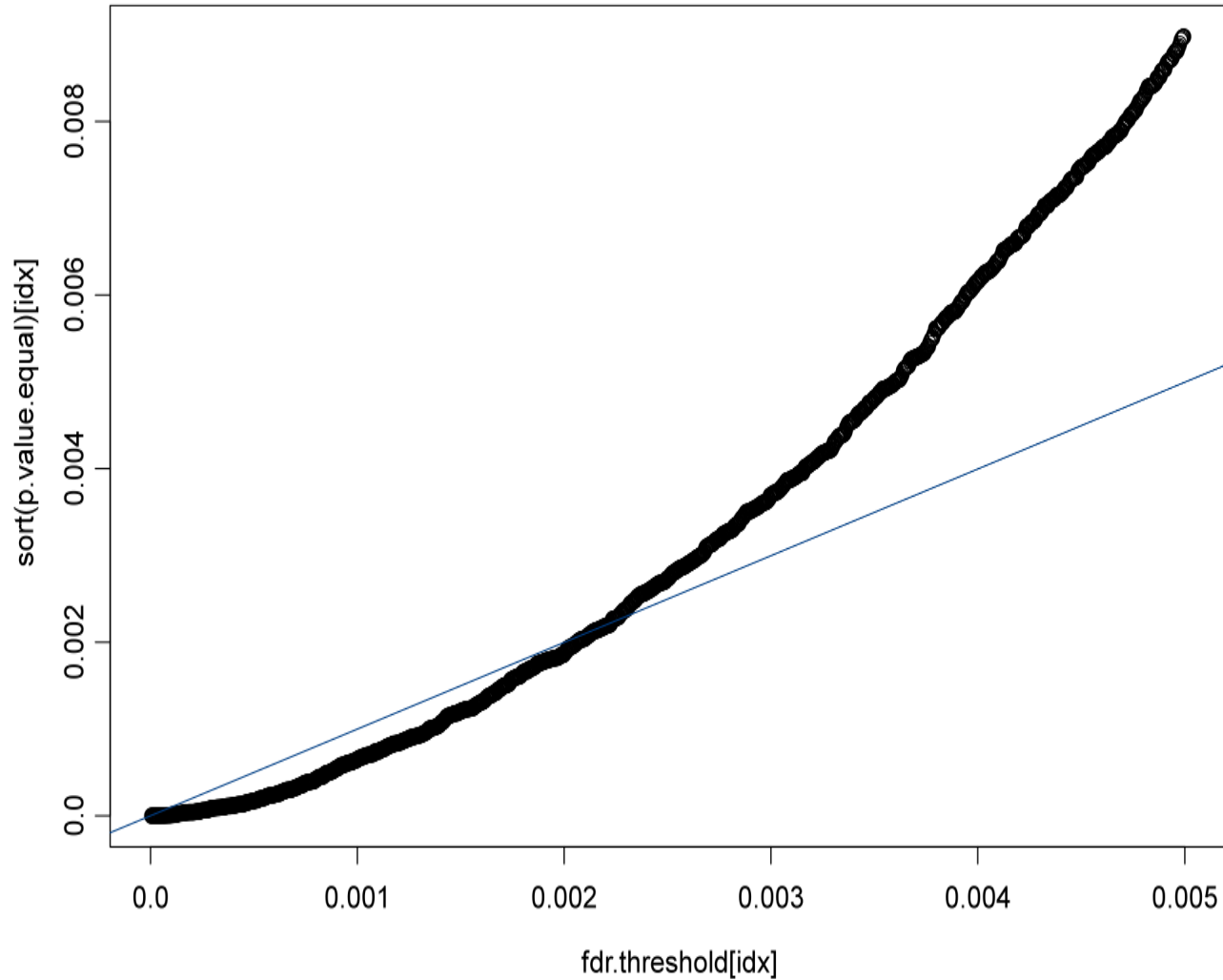
$$p_j^* = \min (mp_j, 1)$$

- **Benjamini & Hochberg (1995): step-up** procedure which controls the FDR under some dependency structures

$$p_{r_j}^* = \min_{k=j \dots m} \{ \min ([m/k] p_{r_k}, 1) \}$$

- In practice
- Sort p-values from smallest to largest
  - Multiply the first by  $m$ , the second by  $m/2$ , the third by  $m/3$  ....
  - Each  $p_j^*$  is at most the minim of the ones after it --monotonicity

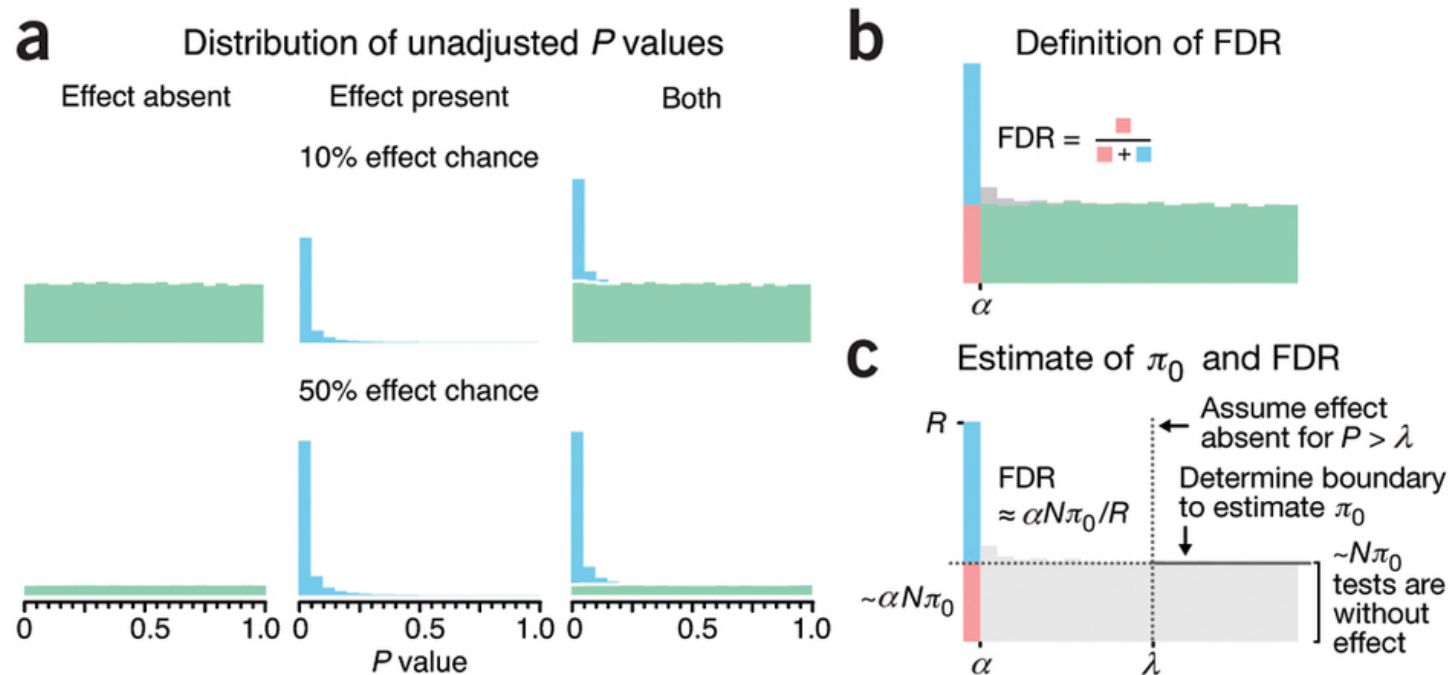
# Visual interpretation



$$mp/k < \alpha$$
$$p < k\alpha/m$$

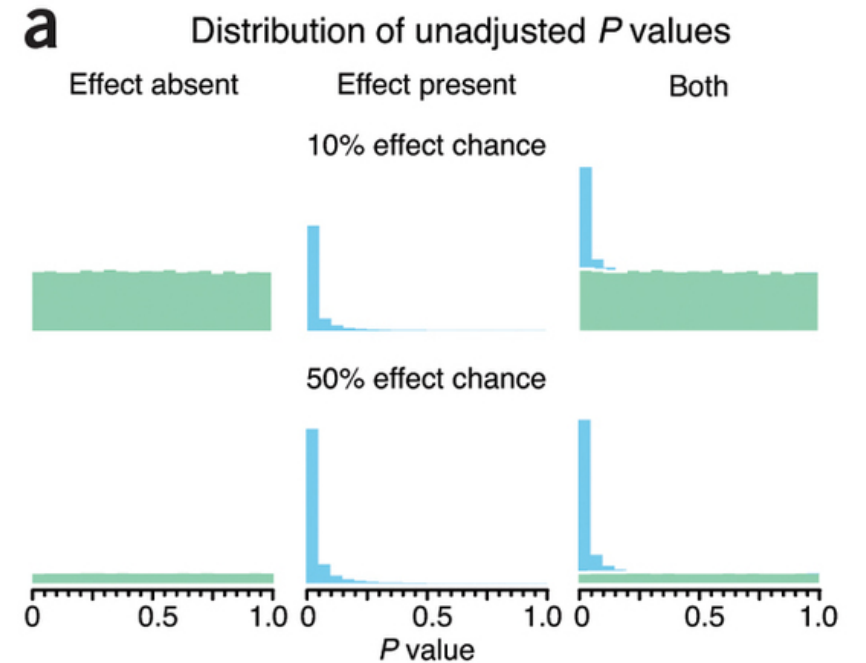
# Q-value

- Storey & Tibshirani, *PNAS*, 2003
- Empirically derived – uses the p-value distribution
- Storey's method first estimates the fraction of comparisons for which the null is true,  $\pi_0$ ,
- counting the number of  $P$  values larger than a cutoff  $\lambda$  (such as 0.5) relative to  $(1 - \lambda)N$  (such as  $N/2$ ), the count expected when the distribution is uniform
- Multiply the Benjamini & Hochberg FDR by  $\pi_0$ , strictly less conservative



# P-value summary

- P-value histogram can tell you there is an effect overall
  - Expect it to be uniform when there is no effect—even though individual test can return very small p-value
- $\pi_0 < 1$  can be used to argue that there is a difference even when no single gene is significant
  - Propose further testing such as aggregating across genes—pathway analysis (discussed later)



# Permutation test: simple example

We measured the expression of a single gene in each rat after treatment and obtain the following results:

	<u>Control</u>	<u>Drug</u>
<b>Expression</b>	9 12 14 17	18 21 23 26
<b>Average</b>	13	22

The difference in averages is  $22-13=9$ .

We wish to claim that this difference was caused by the drug.

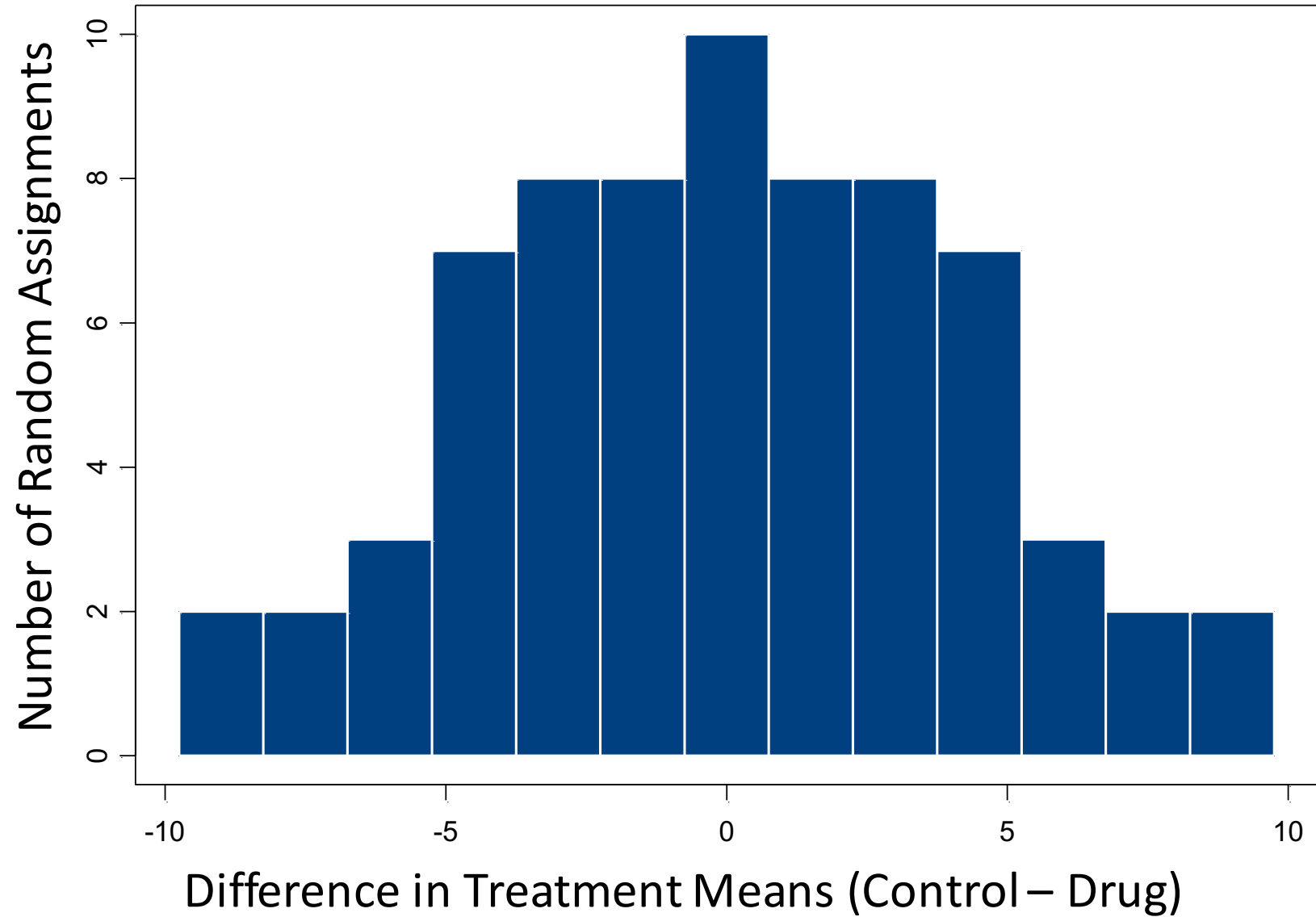
# Interpretation of the Results

- Clearly there is some natural variation in expression (not due to treatment) because the expression measures differ among rats within each treatment group.
- Maybe the observed difference ( $22-13=9$ ) showed up simply because I happened to choose the rats with larger expression for injection with the drug.
- What is the chance of seeing such a large difference in treatment means if the drug has no effect?



Random Assignment	Control				Drug				Difference in Averages
1	9	12	14	17	18	21	23	26	9.0
2	9	12	14	18	17	21	23	26	8.5
3	9	12	14	21	17	18	23	26	7.0
4	9	12	14	23	17	18	21	26	6.0
5	9	12	14	26	17	18	21	23	4.5
6	9	12	17	18	14	21	23	26	7.0
7	9	12	17	21	14	18	23	26	5.5
8	9	12	17	23	14	18	21	26	4.5
9	9	12	17	26	14	18	21	23	3.0
10	9	12	18	21	14	17	23	26	5.0
11	9	12	18	23	14	17	21	26	4.0
12	9	12	18	26	14	17	21	23	2.5
13	9	12	21	23	14	17	18	26	2.5
14	9	12	21	26	14	17	18	23	1.0
15	9	12	23	26	14	17	18	21	0.0
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
69	18	21	23	26	9	12	14	17	-8.5
70	18	21	23	26	9	12	14	17	-9.0

## Distribution of Difference between Treatment Means Assuming No Treatment Effect



# Conclusions

- Only 2 of the 70 possible random assignments would have led to a difference between treatment means as large as 9.
- Thus, under the assumption of no drug effect, the chance of seeing a difference as large as we observed was  $2/70 = 0.0286$ .
- 0.0286 is the empirical p-value for the gene—we generated an empirical null distribution
- Thousands of genes: we observe 100 with a mean  $>9$
- Generate 70 permutations: on average there were 7.7 genes with mean difference  $>9 \rightarrow \text{FDR} = 7.7/100$ 
  - All genes are used to calculate FDR

# Beyond T-test: Significance analysis of microarrays (SAM)

- Significance analysis of microarrays applied to the ionizing radiation response Virginia Goss Tusher, Robert Tibshirani, and Gilbert Chu
- 2001
- With small sample sizes low and high variance can occur by chance
- Variance depends on expression level
- Choose  $S_0$  so that variance is independent of expression level

Difference between the means of the two conditions

$$d(i) = \frac{\bar{x}_I(i) - \bar{x}_U(i)}{s(i) + S_0}$$

Estimate of the standard deviation of the numerator

Fudge Factor

# Assigning significance by permutation

- We have calculated a new statistics and we don't have a parametric description of the null distribution
- Solution: generate an empirical null distribution from a set of experiments where all hypotheses should be null
- Generate permutations of data labels so no difference is expected
- For each permutation  $p$ , calculate  $d_p(i)$ .
- Define FDR
  - Pick a threshold  $d_p$  for calling genes significant
  - Calculate the number of genes above the threshold  $X$
  - Calculate the number of expected falsely differentially expressed genes at that threshold  $Y$  from the permuted sample analysis
  - Compute  $Y/X$
  - 46 real DE genes , 8.4 average across permutation—FDR=.18 (8.4/46)

$$d_p(i) = \frac{\bar{x}_{G1}(i) - \bar{x}_{G2}(i)}{s(i) + s_0}$$

# More on permutations

- Very small experiment-random permutations may create unbalanced groups
  - Solution: restrict to balanced permutations-each permutation should split the real groups equally
- Can be applied to up/down regulated genes separately
- Permutation analysis can be applied to any complicated statistical procedure!

Balanced permutations

Number of red and cyan groups is equal

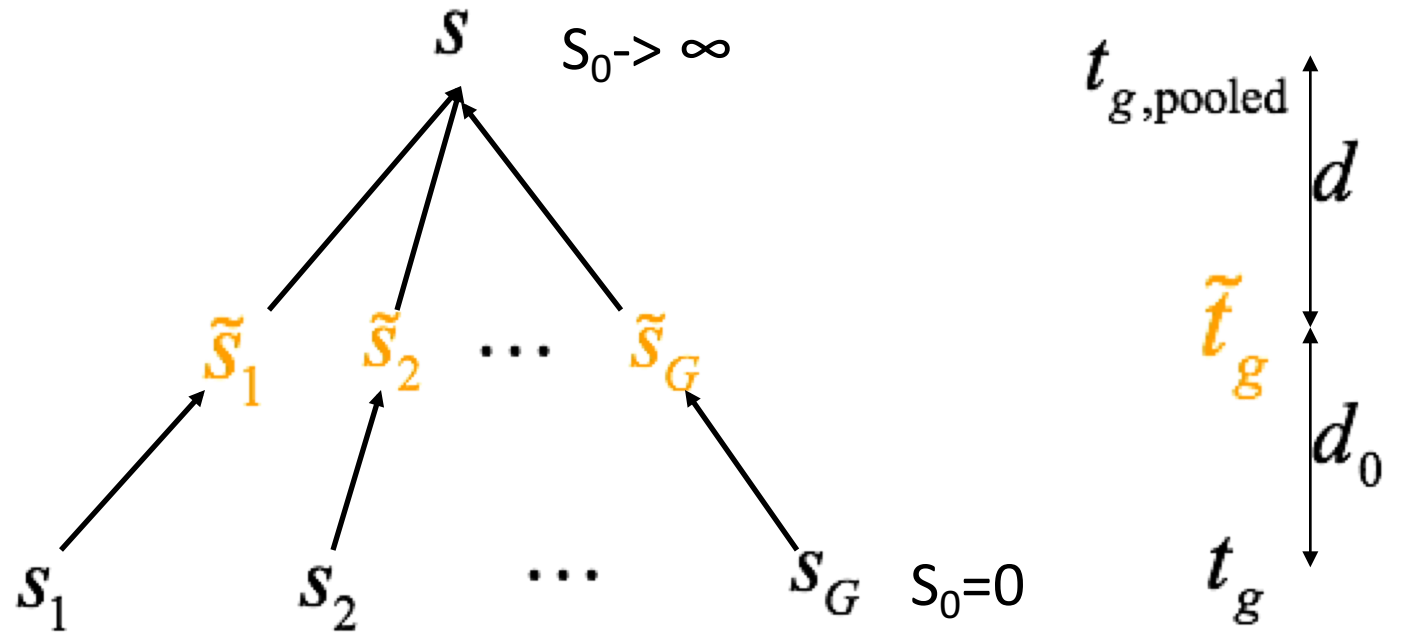
U1A	I1A
U1B	I1B
U2A	I2A
U2B	I2B

U1A	I1A
U1B	I1B
U2A	I2A
U2B	I2B

# Why does SAM work

- Sample variance is not an accurate assessment of the true variance
- What would the per gene variance be we had an infinite number of samples?
- SAM is an example of **moderated T statistic**
- Many current methods use a more principled **Bayesian method**

$$d_p(i) = \frac{\bar{x}_{G1}(i) - \bar{x}_{G2}(i)}{s(i) + s_0}$$



# Bayesian reasoning: short intro

- Synthesize prior knowledge and evidence
- Main theorem

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

- Simple derivation

$$P(A \text{ and } B) =$$

$$P(A|B)P(B) = P(B|A)P(A)$$



# Classical example

- Duchenne Muscular Dystrophy (DMD) can be regarded as a simple recessive sex-linked disease caused by a mutated X chromosome (X).
  - An XY male expresses the disease, whereas an XX female is a carrier but does not express the disease
- Suppose neither of a woman's parents expresses the disease, but her brother does. Then the woman's mother must be a carrier, and the woman herself therefore may be a carrier
- $P(C)=1/2$  ← prior
- What if she has a healthy son? ← observation

$$p(C|h.s.) = \frac{p(h.s.|C)p(C)}{p(h.s.)}$$

$$\frac{p(h.s.|C)p(C)}{p(h.s.|C)p(C) + p(h.s.|\bar{C})p(\bar{C})} =$$

$$\frac{(1/2) \cdot (1/2)}{(1/2) \cdot (1/2) + 1 \cdot (1/2)} = \frac{1}{3}$$

posterior

# Bayesian approach to statistics

- Last example: incorporate evidence into strong prior belief
- Statistics
  - Naïve approach: estimate the parameters from observation only
  - Bayesian approach: have some prior expectation
  - Prior expectation for gene expression:
    - Variance should be independent of mean expression--SAM
    - Gene-specific variance comes from an underlying variance distribution
- Bayesian statistical analyses:
  - begin with 'prior' distributions describing beliefs about the values of parameters in statistical models prior to analysis of the data at hand
  - requires specification of these parameters
  - 'Empirical Bayes' methods use the data at hand to guide prior parameter specification
  - Use all the data to define priors, compute posteriors of individual estimates

# A few more details

- True variance is unknown-only sample variance is known
- Want to estimate true variance from sample variance
- Naïve approach: sample variance == true variance
- Prior: gene specific variance is sampled from some distribution
- We need to make assumptions about the parameterization of the distribution
  - Assumption: gene variance comes from a scaled inverse chi-squared distribution

$$\sigma_1^2, \sigma_2^2, \dots, \sigma_J^2 \sim G(\theta)$$

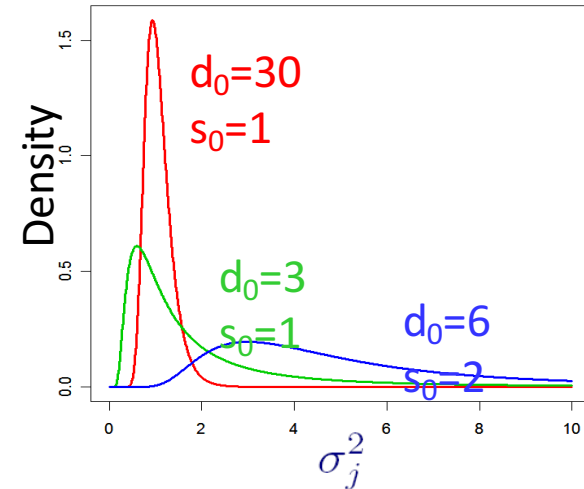
$$\frac{s_0^2}{\sigma_1^2}, \frac{s_0^2}{\sigma_2^2}, \dots, \frac{s_0^2}{\sigma_J^2} \sim \frac{\chi_{d_0}^2}{d_0}$$

Moderated estimate Sample estimate

$$\tilde{s}_j^2 = \frac{ds_j^2 + d_0 s_0^2}{d + d_0}$$

Distribution parameters

Degrees of freedom (n-2)



Software: Limma

Baldi & Long 2001, Wright & Simon 2003, Smyth 2004

# Summary

- Fudge factor is set in a principled way
- Resulting statistics have well understood theoretical behavior
  - We can use the T distribution to assess significance

# More complicated models

- So far we only consider 2 group experiments
- Many other possibilities
  - Factorial: two groups each has two treatments--Are treatment effects different across groups?
  - Continuous variables: dosage of a drug
  - Continuous discrete variables
    - 2 groups, 3 drug doses—do the drugs affect the groups differently?

# General framework for differential expression

- Linear models
- Model the expression of each gene as a linear function of explanatory variables
  - Groups
  - Treatments
  - Combinations of groups and treatments
  - Etc...

$$y = X\beta + \epsilon$$

The diagram shows the equation  $y = X\beta + \epsilon$  with three arrows pointing to the terms: a blue arrow from 'vector of observed data' to  $y$ , a red arrow from 'design matrix' to  $X$ , and a blue arrow from 'Vector of parameters to estimate' to  $\beta$ .

vector of observed data

design matrix

Vector of parameters to estimate

# Example of a design matrix

Normal sample x 2



Cancer Sample x 2



$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

$\beta_1$  = normal log-expression

$\beta_2$  = cancer - wt

$$E[y_1] = E[y_2] = \beta_1$$

$$E[y_3] = E[y_4] = \beta_1 + \beta_2$$

# More examples

- 6 samples
- 2 groups + drug treatment
- Group and treatment effect are additive

$$y = X\beta + \epsilon$$

Group1	Group 2- Group 1	Drug dose
1	0	0.25
1	0	1
1	0	4
1	1	0.25
1	1	1
1	1	4

3 coefficients to estimate



# More examples

$$y = X\beta + \epsilon$$

- 6 samples
- 2 groups + drug treatment
- Treatment affects groups differently

Group1	Group 2- Group 1	Drug dose	Drug dose + Group 2
1	0	0.25	0
1	0	1	0
1	0	4	0
1	1	0.25	0.25
1	1	1	1
1	1	4	4

4 coefficients to estimate

# Linear model parameter estimation

Model is specified –how do we find the coefficients

- Minimize squared error

$$y = X\beta + \epsilon$$

$$\epsilon'\epsilon = (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)$$

- Take derivative

$$\frac{d}{d\beta} ((\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)) = -2\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta)$$

- Set to 0

$$-2\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0}$$

- Solve

$$\mathbf{X}'\mathbf{Y} = (\mathbf{X}'\mathbf{X})\beta$$

- Significance of coefficients is tested with a T-test

$$\beta = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

# Hypothesis testing

$$y = X\beta + \epsilon$$

$\beta$  can be a vector. We can test the significance of any one coefficient  $\beta_i$  via a T test

$$t_{\text{score}} = \frac{\hat{\beta} - \beta_0}{SE_{\hat{\beta}}}$$

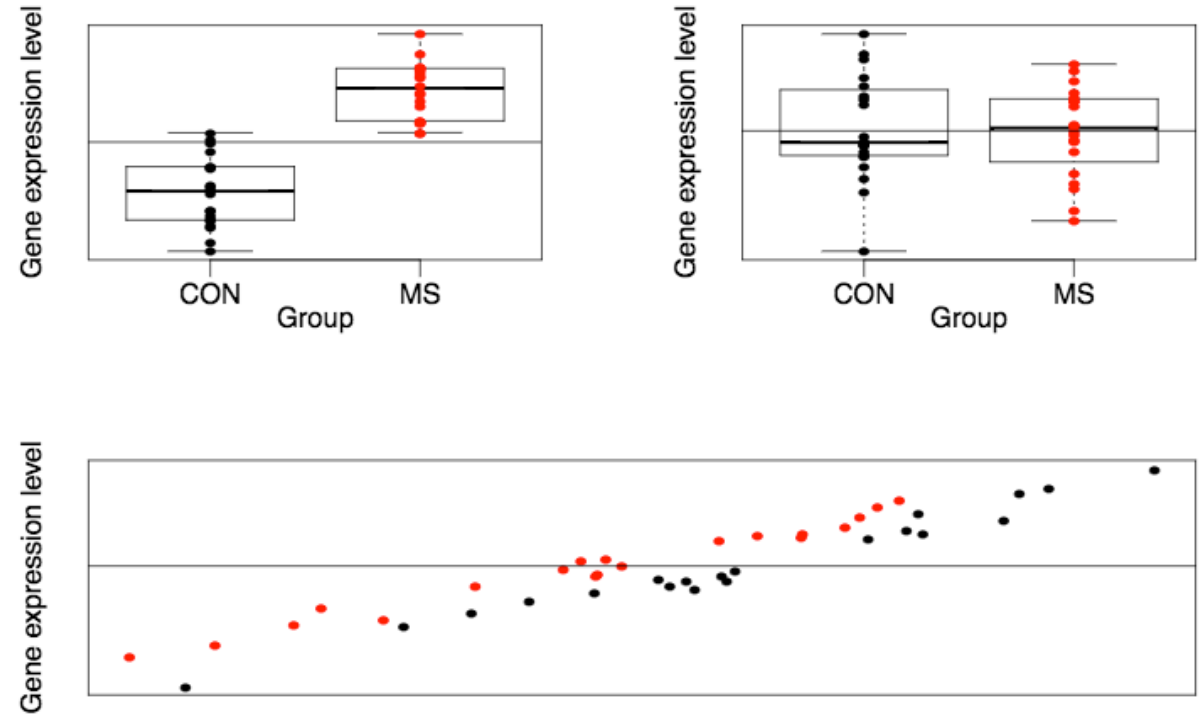
$$t_{\text{score}} = \frac{(\hat{\beta} - \beta_0)\sqrt{n-2}}{\sqrt{SSR / \sum_{i=1}^n (x_i - \bar{x})^2}}$$

SSR depends on the whole model

$$SSR = \sum_{i=1}^n \hat{\epsilon}_i^2 = \text{sum of squares of residuals.}$$

# Linear models for data clean up

- Linear models are useful for including nuisance variables-- Technical factors
- Variables that have an effect on measurements but are not themselves of interest
- 2 group design: Control vs MS (multiple sclerosis)
- Variable sample storage time
- Incorporating storage time gives us smaller residuals and thus larger T-stats for the disease coefficient

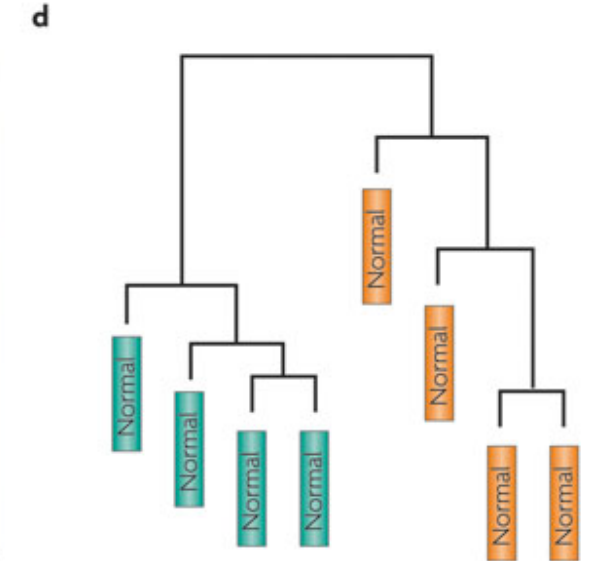
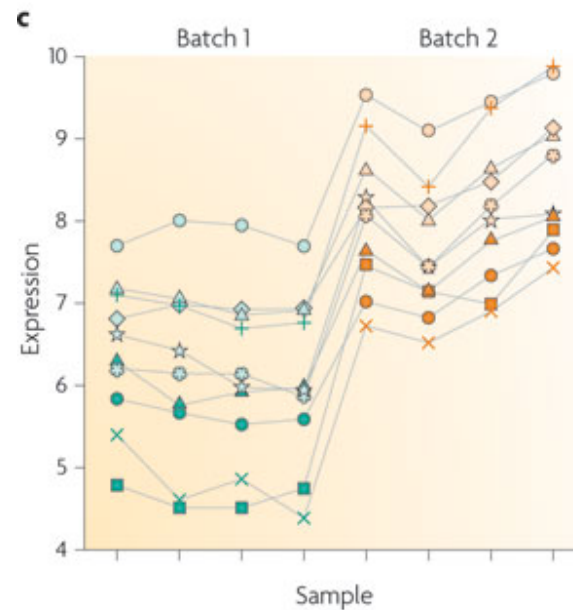
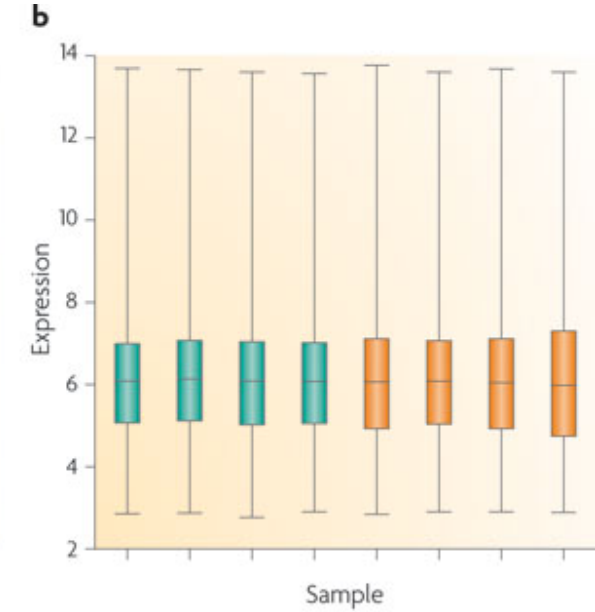
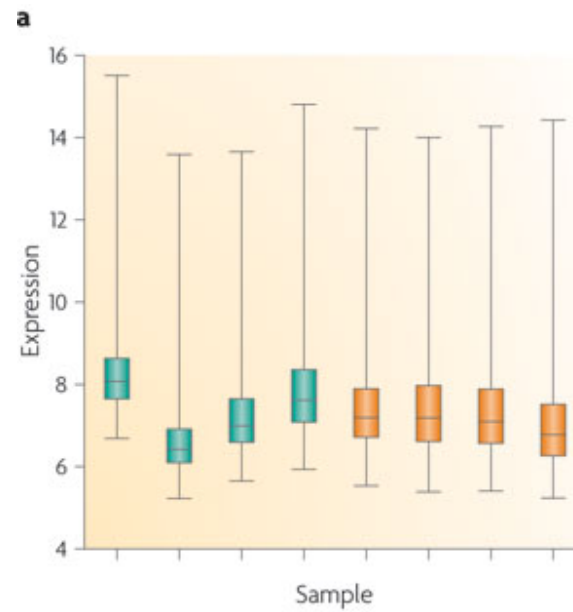


$X =$

sample storage time		
Control	MS-Control	Storage time
1	0	6
1	0	1
1	0	4
1	1	7
1	1	1
1	1	8

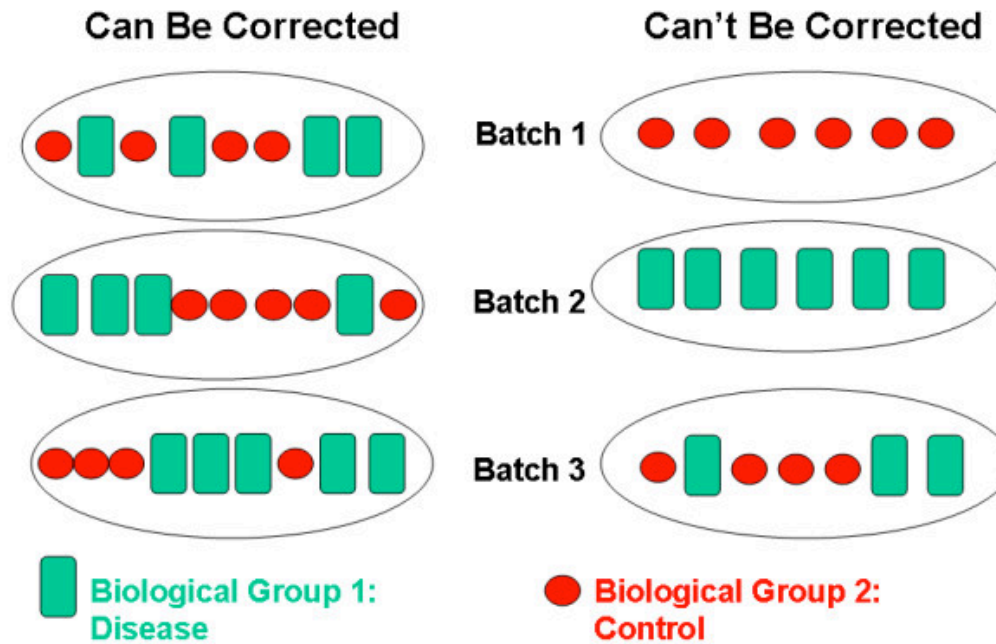
# Batch effects

- Batch effects: technical variables that effect gene expression
- Most obvious when samples were processed in “batches”
  - At different times
  - Different locations
  - Different technician
  - Different protocol
- Affects all high throughput measurement platforms



# Batch effects

- Batch variables are often discrete but can be continuous (such as: storage time)
  - With RNAseq we can model sample specific GC bias
- Batch effects can be corrected if they don't align (perfectly) with a variable of interest
- Experimental design is important!



We want to compare 3 groups in 3 replicates over 3 days  
How should we proceed?

Day 1	Day 2	Day 3
Group 1 replicate 1	Group 1 replicate 2	Group 1 replicate 3
Group 2 replicate 1	Group 3 replicate 2	Group 2 replicate 3
Group 3 replicate 1	Group 3 replicate 2	Group 3 replicate 3

Day 1	Day 2	Day 3
Group 1 replicate 1	Group 2 replicate 1	Group 3 replicate 1
Group 1 replicate 2	Group 3 replicate 2	Group 3 replicate 2
Group 1 replicate 3	Group 2 replicate 3	Group 3 replicate 3

# Pathway/geneset analysis

- All methods discussed so far apply to arbitrary high-dimensional data
  - All the gene labels can be hidden
- ...but for genes we know a lot about their identity
- We can assign genes to pathways and functional categories
- Examples
  - Genes in a signaling pathway—T-cell receptor signaling
  - Metabolic pathway—glycolysis enzymes
  - Genes known to be tissue cell/type specific—genes related to neuronal synapse formation
- We can use this to improve statistical interpretation of gene expression data

# Simple geneset analysis

- Are genes of functional category X overrepresented among the genes declared to be differentially expressed?
- Testing for proportion differences
- Fisher's exact test --- hypergeometric test

		Gene of Functional Category X?			
Declared to be Differentially Expressed?		yes	no	total	fraction
	yes	50	250	300	1 in 6
	no	50	19900	19700	
		100	19900	20000	1 in 200



# Problems Fisher's Exact Test for Detecting Overrepresentation

- The outcome of the overrepresentation test depends on the significance threshold used to declare genes differentially expressed.
- Functional categories in which many genes exhibit small changes may go undetected.
- Genes are not independent, so a key assumption of the Fisher's exact tests is violated.

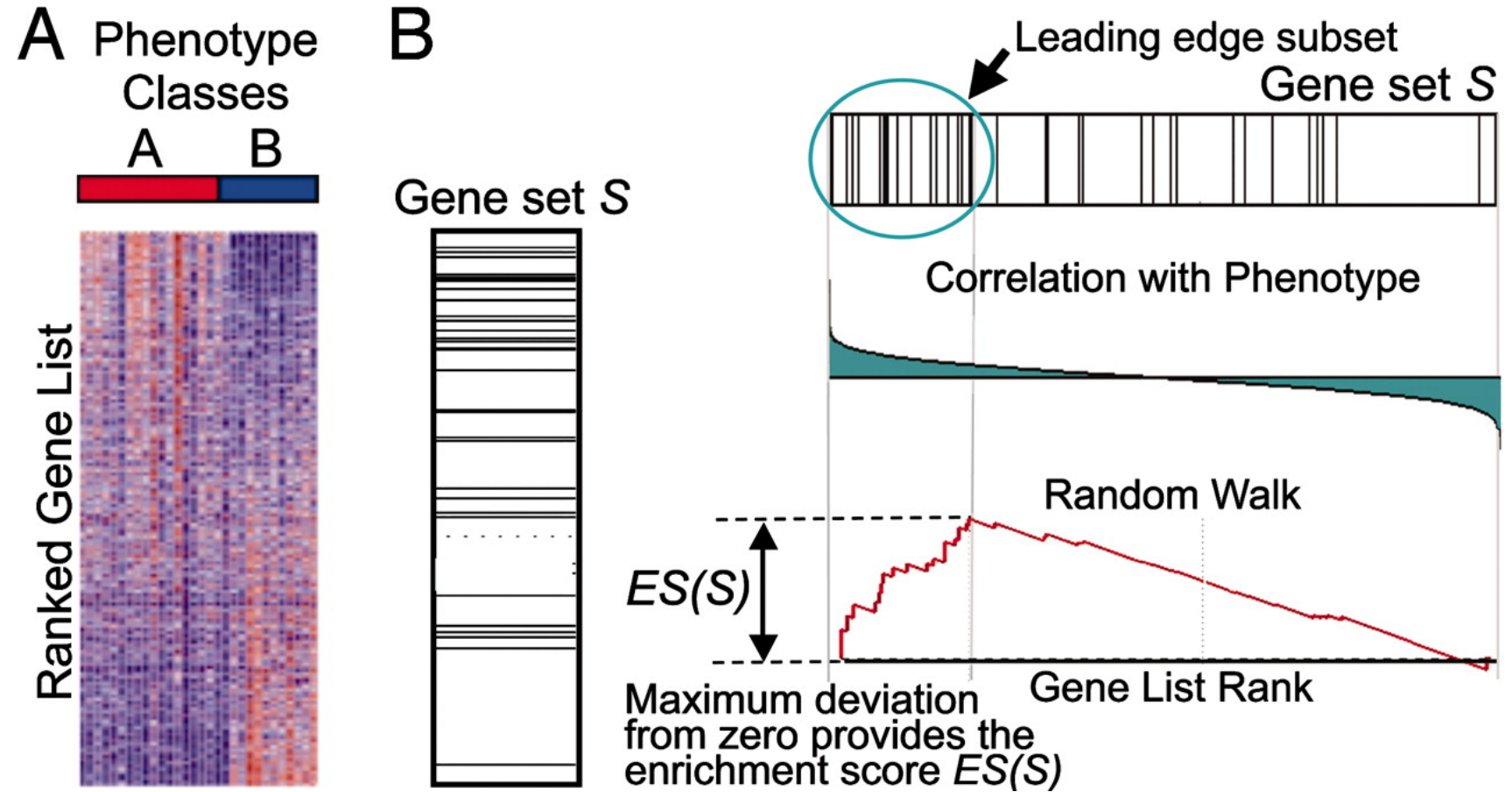
# Gene Set Enrichment Analysis (GSEA)

- Subramanian, et al. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*. **102**, 15545-15550.
- Compute a statistic (difference between 2 clinical groups) for each gene that measures the degree of differential expression between treatments.
- Create a list L of all genes ordered according to these statistics.
- Given a set of genes S we can see if these genes are non-randomly distributed in our list L
  - if the experiment produced random results we don't expect gene order to have biological coherence

# GSEA (continued)

- Calculate an enrichment score ( $ES$ ) that reflects the degree to which a set  $S$  is overrepresented at the extremes (top or bottom) of the entire ranked list  $L$ .
- The score is calculated by walking down the list  $L$  and ...
  - increase a running-sum statistic when we encounter a gene in  $S$
  - Decrease it when we encounter genes not in  $S$ .
- The magnitude of the increment depends on the correlation of the gene with the phenotype.
- The final enrichment score is the maximum deviation from zero encountered in the random walk
  - corresponds to a weighted Kolmogorov–Smirnov-like statistic

# GSEA



- Significance and FDR are calculated by permuting groups (phenotype classes)
- Positive results can be used to argue that the experiment worked (we detected molecular difference between clinical groups) even when no single gene meets an acceptable level of significance

# GSEA alternatives

- Other metrics to summarize unusual distribution of pathway specific genes are
  - Rank-sum (equivalent to AUC or wilcox test), [wilcoxGST in limma package](#)
  - Minmax, the maximum in absolute value of the positive and negative average T stat—can be show to have good statistical properties--Bradley Efron and Rob Tibshirani. Tech report. August 2006 -[Software: GSA](#)

# mSigDB-database of gene sets

<b>C2: curated gene sets</b> ( <a href="#">browse 4722 gene sets</a> )	Gene sets collected from various sources such as online pathway databases, publications in PubMed, and knowledge of domain experts. The gene set page for each gene set lists its source. <a href="#">details</a>	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>
CGP: chemical and genetic perturbations ( <a href="#">browse 3402 gene sets</a> )	Gene sets represent expression signatures of genetic and chemical perturbations. A number of these gene sets come in pairs: an xxx_UP (xxx_DN) gene set representing genes induced (repressed) by the perturbation. The gene set page for each gene set lists the PubMed citation on which it is based.	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>
CP: Canonical pathways ( <a href="#">browse 1320 gene sets</a> )	Gene sets from the pathway databases. Usually, these gene sets are canonical representations of a biological process compiled by domain experts. <a href="#">details</a>	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>
CP:BIOCARTA: BioCarta gene sets ( <a href="#">browse 217 gene sets</a> )	Gene sets derived from the BioCarta pathway database ( <a href="http://www.biocarta.com/genes/index.asp">http://www.biocarta.com/genes/index.asp</a> ).	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>
CP:KEGG: KEGG gene sets ( <a href="#">browse 186 gene sets</a> )	Gene sets derived from the KEGG pathway database ( <a href="http://www.genome.jp/kegg/pathway.html">http://www.genome.jp/kegg/pathway.html</a> ).	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>
CP:REACTOME: Reactome gene sets ( <a href="#">browse 674 gene sets</a> )	Gene sets derived from the Reactome pathway database ( <a href="http://www.reactome.org/">http://www.reactome.org/</a> ).	Download GMT Files <a href="#">original identifiers</a> <a href="#">gene symbols</a> <a href="#">entrez genes ids</a>

# General approaches to multi dimensional data

- Many more measurements than samples
- Use measurement distributions for normalization and filtering
- Borrow information across measurements for hypothesis testing

# Comparisons

- In general, for a given multiple testing procedure,

$$\text{PCER} \leq \text{FWER} \leq \text{PFER},$$

and

$$\text{FDR} \leq \text{FWER},$$

with  $\text{FDR} = \text{FWER}$  under the complete null



## Cluster analyses:

- 1) Usually outside the normal framework of statistical inference;
- 2) less appropriate when only a few genes are likely to change.
- 3) Needs lots of experiments

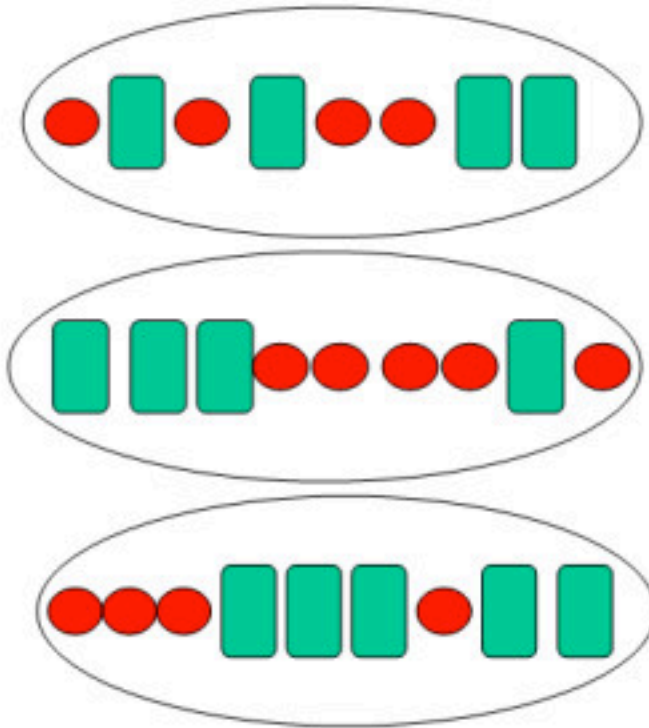
## Single gene tests:

- 1) may be too noisy in general to show much
- 2) may not reveal coordinated effects of positively correlated genes.
- 3) hard to relate to pathways.

# Example

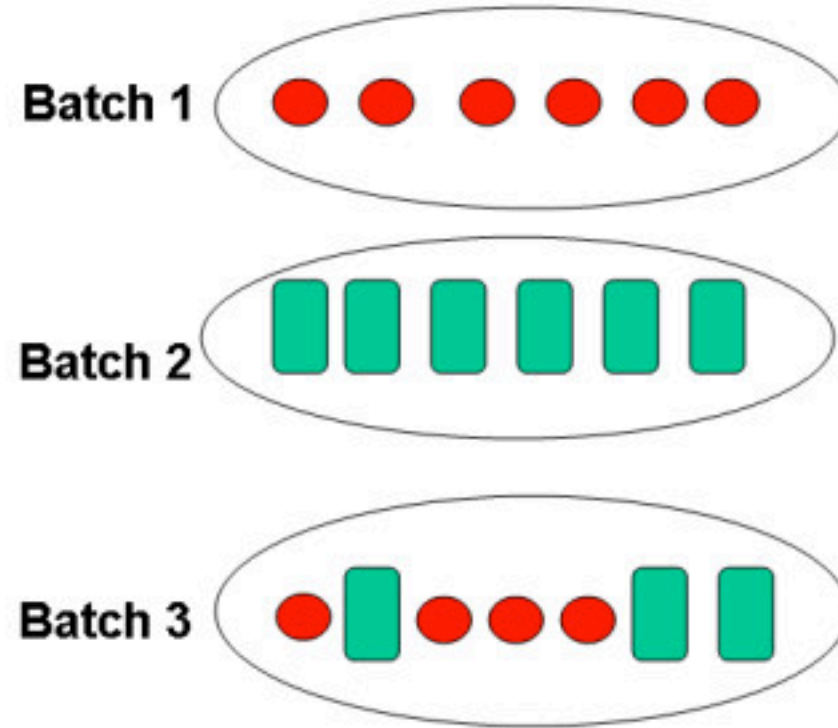
- 150 genes were difference with a T test p-value of  $< 0.05$
- Is this a good result?

## Can Be Corrected



 **Biological Group 1:  
Disease**

## Can't Be Corrected



 **Biological Group 2:  
Control**