

# Key Aspects of the Design & Analysis of DNA Microarray Studies for Diagnostic & Prognostic Prediction

Richard Simon, D.Sc.  
Chief, Biometric Research Branch  
National Cancer Institute  
<http://linus.nci.nih.gov/brb>

# <http://linus.nci.nih.gov/brb>

- Powerpoint presentation
  - Including slides not presented
- Bibliography
- Reprints & Technical Reports
  - Support for my dogmatic assertions
- BRB-ArrayTools software
  - Performs all analyses described

- *Design and Analysis of DNA Microarray Investigations*
  - R Simon, EL Korn, MD Radmacher, L McShane, G Wright, Y Zhao. Springer (2003)

# Microarray Expression Profiling

- Would like to know the concentration of each protein in a cell
  - Proteins do the work of cells
  - Proteins have many shapes and parallel assays for all proteins have not been developed

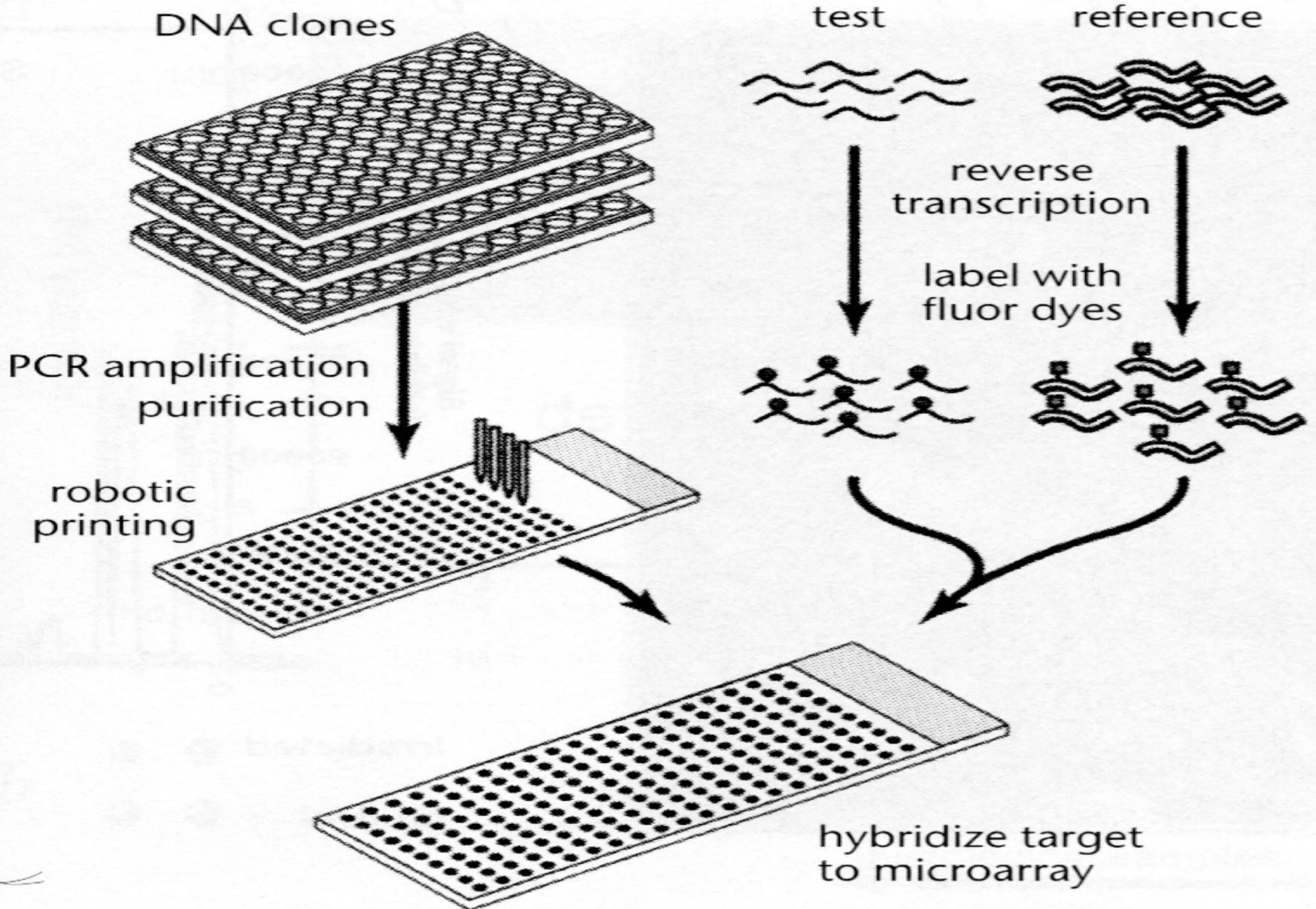
# Microarray Expression Profiling

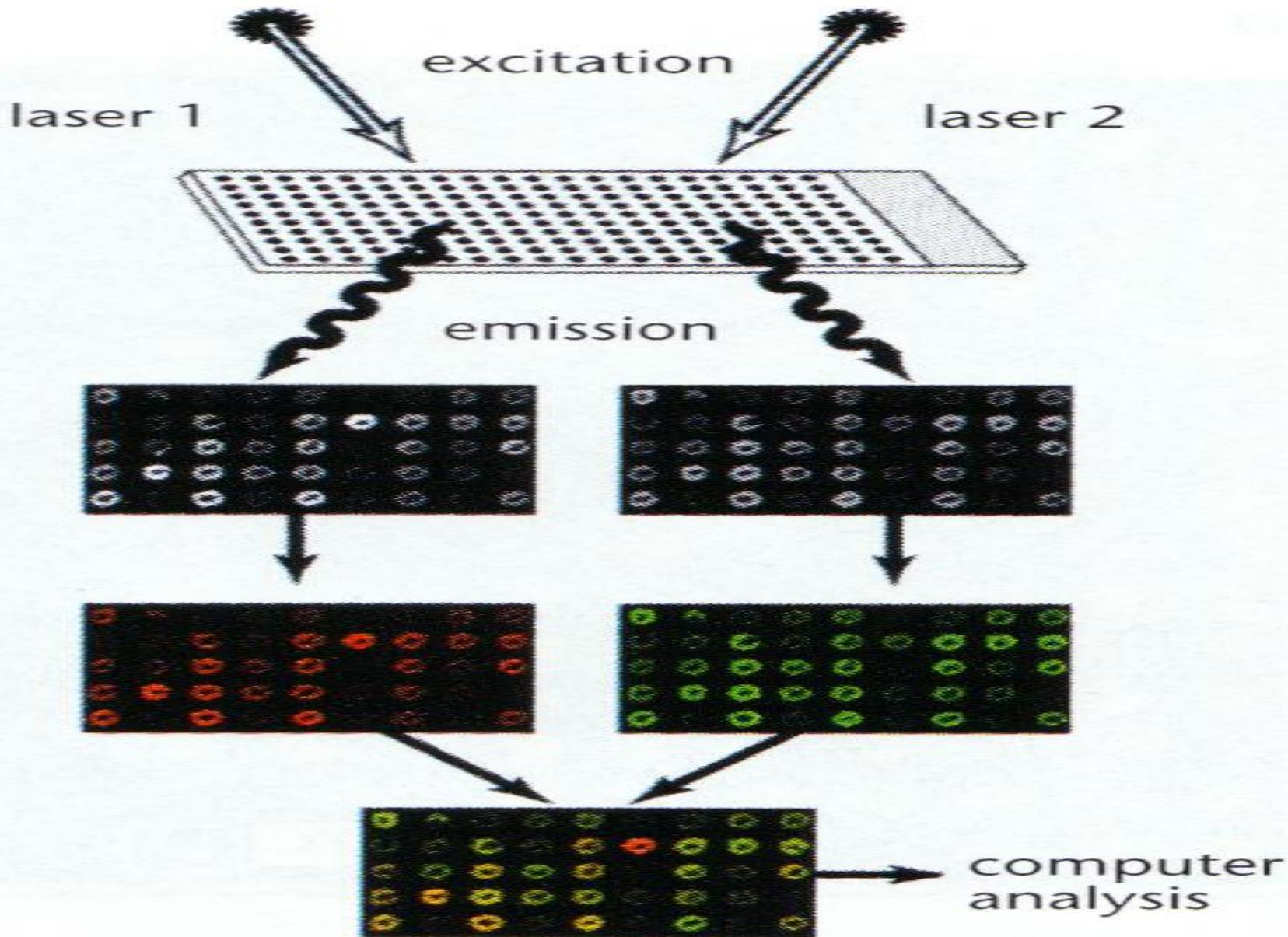
- One gene *transcription* produces one mRNA molecule produces one protein molecule
- # genes  $\cong$  # mRNA types
- mRNA molecules are similar in structure to DNA
  - Proteins are not
- mRNA molecule can be reverse transcribed into DNA and will bind only to the gene from which it was originally transcribed (to which it is *homologous*)

# Microarray Expression Profiling

- Estimates abundance of mRNA molecules of each type present in cells
  - Assay not sensitive enough to analyze single cells so estimate is for average of sample of cells
- Microarray contains a spot of DNA corresponding to each gene
  - Spots are in known fixed positions
  - Spots contain fewer nucleotides than the full gene

# cDNA Array





A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
NAME	TYPE	ACC	CLID	SPOT	svcc77 CH1D	svcc77 CH2D	svcc77 FLAG	svcc78 CH1D	svcc78 CH2D	svcc78 FLAG	svcc86 CH1D	svcc86 CH2D	svcc86 FLAG	svcc104 CH1D
zinc finger	cDNA	AA406467	753234	1	1846	2088	0	6650	10328	1	2404	2608	0	1635
small proli	cDNA	AA447835	813614	577	352	492	0	527	275	0	733	193	0	742
zinc finger	cDNA	T57959	71626	1153	2601	1883	0	2677	1685	0	4424	2797	0	4512
486544	cDNA	AA043334	486544	1729	3527	2745	0	2059	1431	0	3773	3329	0	3508
zinc finger	cDNA	H17047	50794	2	3444	3039	0	6830	9446	1	2822	2993	0	2287
small indu	cDNA	H62985	205633	578	842	1292	0	1330	1514	0	1224	1953	0	1082
Human PC	cDNA	AA425602	768644	1154	648	666	0	1275	860	0	1320	979	0	1766
small indu	cDNA	AA425102	768561	1730	4951	1797	0	2714	1055	0	5518	3272	0	5428
ESTs, Higl	cDNA	W16724	302190	3	2170	1897	0	6442	13371	0	1993	1945	0	4401
Sjogren sy	cDNA	H29484	49970	579	3710	2356	0	10049	2749	0	6335	4325	0	9894
zinc finger	cDNA	AA088564	511814	1155	3106	1890	0	3429	2074	0	3435	2085	0	5241
signal rec	cDNA	AA411407	754998	1731	3538	3169	0	3523	1513	0	4089	3632	0	4301
Down sync	cDNA	H19439	51408	4	1694	1256	0	6505	6343	1	1469	749	0	982
sex hormo	cDNA	T69346	82871	580	387	676	0	502	406	0	648	522	0	855
alpha thal	cDNA	AA410435	753430	1156	94	185	0	929	798	0	612	643	0	2351
selectin L	cDNA	H00756	149910	1732	1338	949	0	900	764	0	1721	1282	0	1189
wingless-t	cDNA	W49672	324901	5	6690	3050	0	8633	5541	1	6634	1342	0	4471
selectin E	cDNA	H39560	186132	581	207	517	0	254	954	0	689	1262	0	517
wingless-t	cDNA	T99653	122762	1157	272	2075	0	365	1404	0	362	1414	0	437
sarcoglyc	cDNA	AA234982	666829	1733	476	586	0	654	529	0	816	722	0	419
von Hippel	cDNA	H73054	234856	6	2517	2255	0	3252	2177	1	2737	2276	0	1974
steroid sul	cDNA	H15215	49591	582	5075	5328	0	7069	5692	0	8264	10019	0	9560
visinin-like	cDNA	H65066	210575	1158	1098	666	0	1046	448	0	959	637	0	1969
SRY (sex-	cDNA	AA400739	753184	1734	5026	7070	0	2180	3533	0	4460	8656	0	4563
ESTs, Moc	cDNA	H69834	213280	2305	112	374	0	420	274	0	169	215	0	161
phosphoin	cDNA	AA464765	810372	2881	879	409	1	867	87	1	927	290	0	1512
Homo sapi	cDNA	AA026831	469345	3457	338	1010	0	443	1263	0	566	987	0	403
peroxisom	cDNA	H25923	162211	4033	2574	2527	0	3506	2247	0	2116	2810	0	3123
keratin 13	cDNA	W23757	327676	2306	517	931	0	495	483	0	460	435	0	298
peroxisom	cDNA	H10964	47142	2882	1530	1466	1	1731	1657	1	2267	1805	0	3138
50182	cDNA	H17882	50182	3458	581	800	0	1241	980	0	765	1807	0	871
peripheral	cDNA	R26960	133273	4034	3795	9832	0	4501	4569	0	2551	6870	0	3832

# [Affymetrix] Hybridization Oligo Array



# Affymetrix GeneChips

- Contain multiple probes (spots) per gene
- Probes corresponding to the same gene must be processed to give a probe-set summary intensity for each gene
- Single label system
  - Higher reproducibility makes use of dual-labels unnecessary

# Important Topics I Will Not Discuss

- Image Analysis
- Background adjustment
- Probe-set intensity summaries for Affymetrix GeneChips
- Normalization
  - Red vs green bias on dual label arrays
  - Across arrays for single channel arrays

# Myths & Truths About Microarray Expression Profiling

<http://linus.nci.nih.gov/brb>

# Myth

- That microarray investigations should be unstructured data-mining adventures without clear objectives

- Good microarray studies have clear objectives, but not generally gene specific mechanistic hypotheses
- Design and analysis methods should be tailored to study objectives

# Common Types of Objectives

- **Class Comparison**
  - Identify genes differentially expressed among predefined classes such as diagnostic or prognostic groups.
- **Class Prediction**
  - Develop multi-gene predictor of class for a sample using its gene expression profile
- **Class Discovery**
  - Discover clusters among specimens or among genes

# Do Expression Profiles Differ for Two Defined Classes of Arrays?

- Not a clustering problem
  - Global similarity measures generally used for clustering arrays may not distinguish classes
  - Supervised methods
- Requires multiple biological samples from each class
  - Contrary to published statistical methods and widely used software

# Levels of Replication

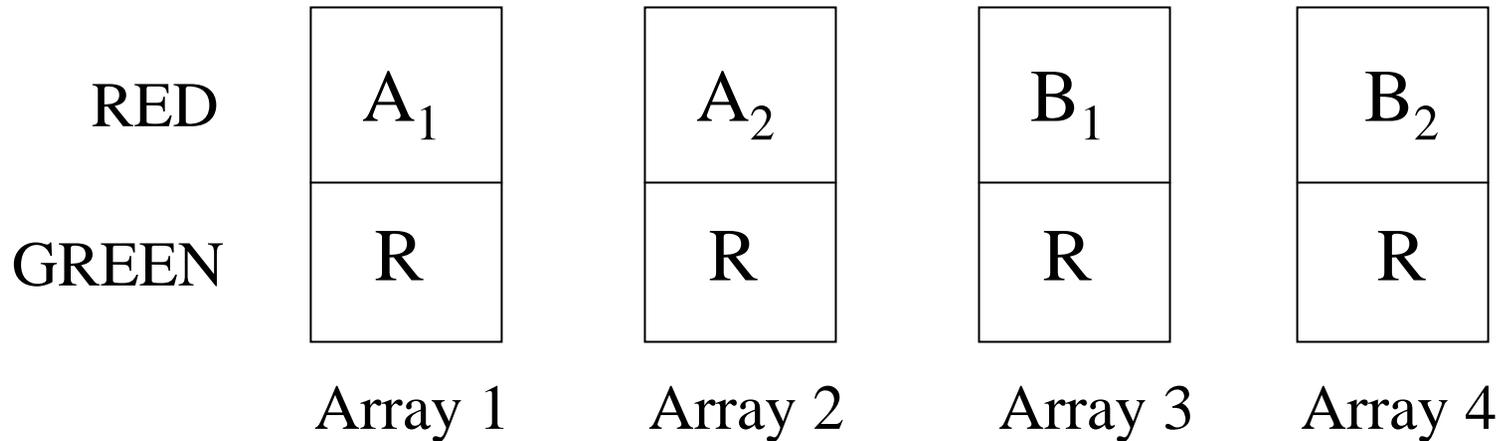
- Technical replicates
  - RNA sample divided into multiple aliquots and re-arrayed
- Biological replicates
  - Multiple subjects
  - Replication of the tissue culture experiment

- Biological conclusions generally require independent biological replicates. The power of statistical methods for microarray data depends on the number of biological replicates.
- Technical replicates are useful insurance to ensure that at least one good quality array of each specimen will be obtained.

# **Allocation of Specimens to Dual Label Arrays for Simple Class Comparison Problems**

- Common Reference Design
- Balanced Block Design

# Common Reference Design

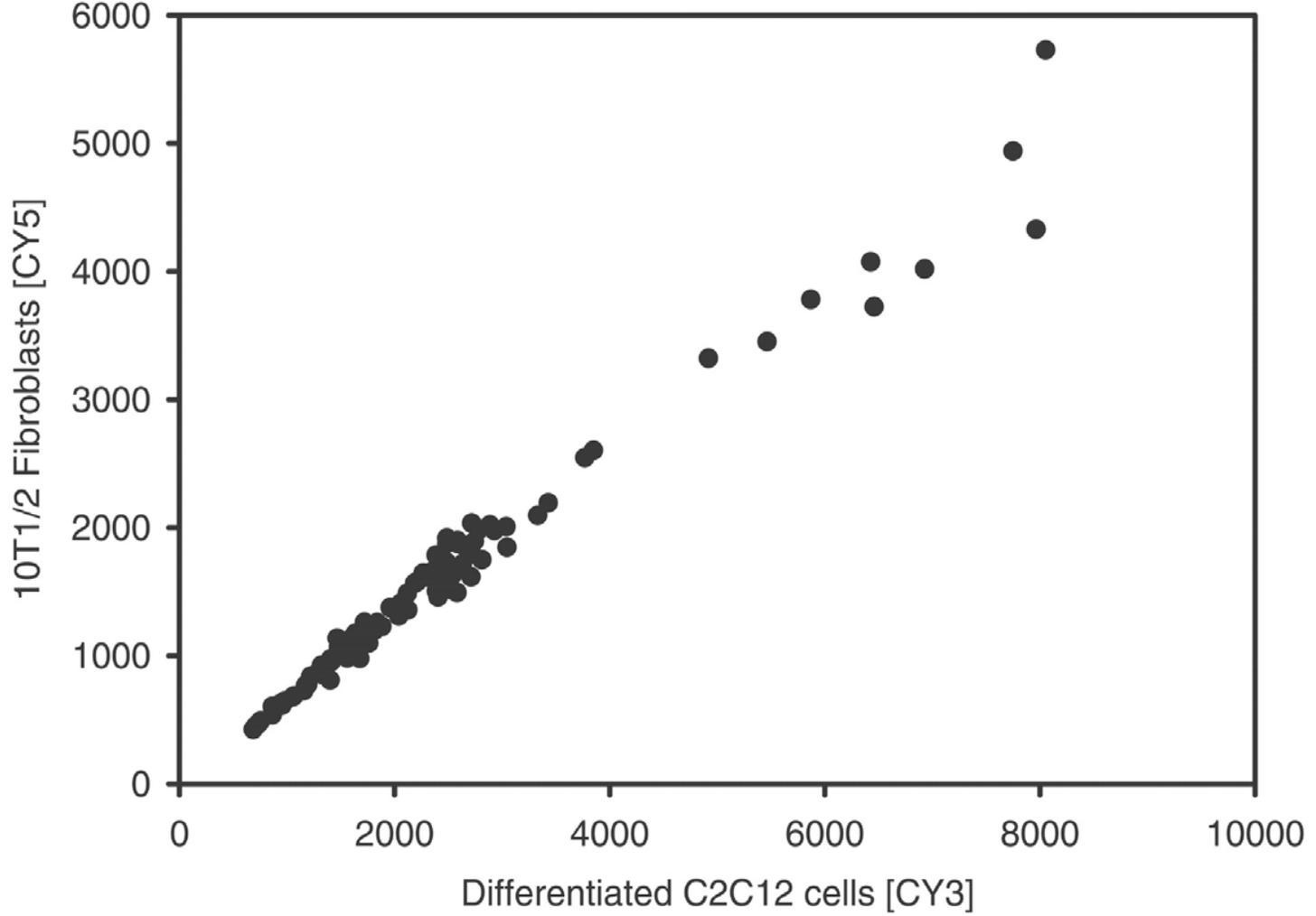


$A_i$  = *i*th specimen from class A

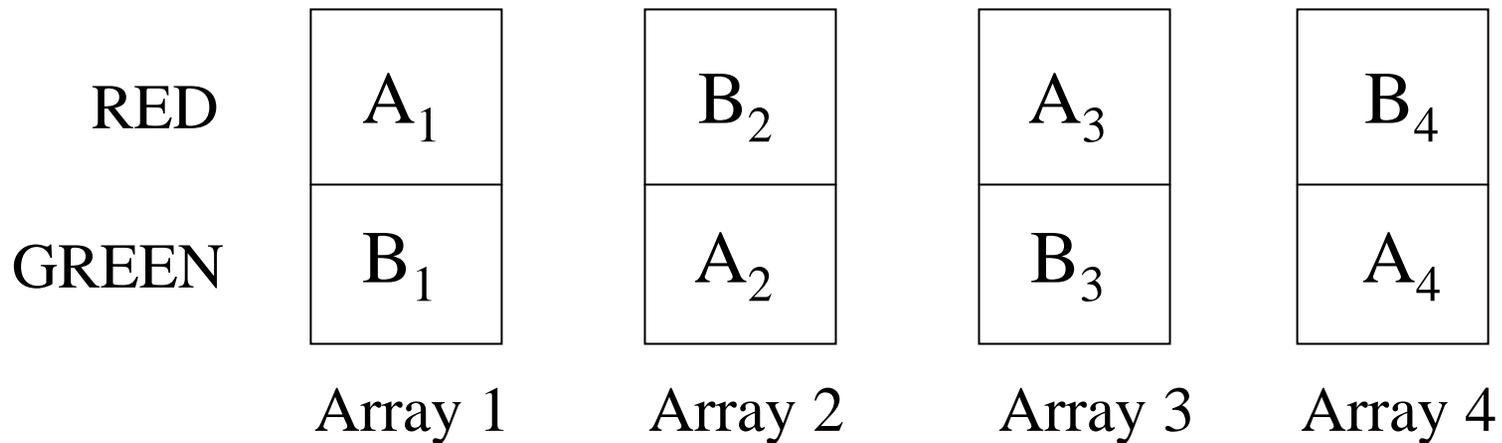
$B_i$  = *i*th specimen from class B

R = aliquot from reference pool

- The reference generally serves to control variation in the size of corresponding spots on different arrays and variation in sample distribution over the slide.
- The reference provides a relative measure of expression for a given gene in a given sample that is less variable than an absolute measure.
- The reference is not the object of comparison.
- The relative measure of expression will be compared among biologically independent samples from different classes.



# Balanced Block Design



$A_i$  =  $i$ th specimen from class A

$B_i$  =  $i$ th specimen from class B

- Detailed comparisons of the effectiveness of designs:
  - Dobbin K, Simon R. Comparison of microarray designs for class comparison and class discovery. *Bioinformatics* 18:1462-9, 2002
  - Dobbin K, Shih J, Simon R. Statistical design of reverse dye microarrays. *Bioinformatics* 19:803-10, 2003
  - Dobbin K, Simon R. Questions and answers on the design of dual-label microarrays for identifying differentially expressed genes, *JNCI* 95:1362-1369, 2003

- Common reference designs are very effective for many microarray studies. They are robust, permit comparisons among separate experiments, and permit many types of comparisons and analyses to be performed.
- For simple two class comparison problems, balanced block designs require many fewer arrays than common reference designs.
  - Efficiency decreases for more than two classes
  - Are more difficult to apply to more complicated class comparison problems.
  - They are not appropriate for class discovery or class prediction.
- Loop designs are less robust, and dominated by either common reference designs or balanced block designs, and are not suitable for class prediction or class discovery.

# Myth

- For two color microarrays, each sample of interest should be labeled once with Cy3 and once with Cy5 in dye-swap pairs of arrays.

# Dye Bias

- Average differences among dyes in label concentration, labeling efficiency, photon emission efficiency and photon detection are corrected by normalization procedures
- Gene specific dye bias may not be corrected by normalization

- Dye swap technical replicates of the same two rna samples are rarely necessary.
- Using a common reference design, dye swap arrays are not necessary for valid comparisons of classes since specimens labeled with different dyes are never compared.
- For two-label direct comparison designs for comparing two classes, it is more efficient to balance the dye-class assignments for independent biological specimens than to do dye swap technical replicates

# Can I reduce the number of arrays by pooling specimens?

- Pooling all specimens is inadvisable because conclusions are limited to the specific RNA pool, not to the populations since there is no estimate of variation among pools
- With multiple biologically independent pools, some reduction in number of arrays may be possible but the reduction is generally modest and may be accompanied with a large increase in the number of independent biological specimens needed
  - Dobbin & Simon, Biostatistics (In Press).

Number of samples pooled per array	Number of arrays required	Number of samples required
1	25	25
2	17	34
3	14	42
4	13	52

$\alpha=0.001, \beta=0.05, \delta=1, \tau^2+2\sigma^2=0.25, \tau^2/\sigma^2=4$

# Sample Size Planning

- GOAL: Identify genes differentially expressed in a comparison of two pre-defined classes of specimens on dual-label arrays using reference design or single label arrays
- Compare classes separately by gene with adjustment for multiple comparisons
- Approximate expression levels (log ratio or log signal) as normally distributed
- Determine number of samples  $n/2$  per class to give power  $1-\beta$  for detecting mean difference  $\delta$  at level  $\alpha$

# Comparing 2 equal size classes

$$n = 4\sigma^2(z_{\alpha/2} + z_{\beta})^2/\delta^2$$

where  $\delta$  = mean log-ratio difference between classes

$\sigma$  = standard deviation

$z_{\alpha/2}, z_{\beta}$  = standard normal percentiles

- Choose  $\alpha$  small, e.g.  $\alpha = .001$
- Use percentiles of t distribution for improved accuracy

# Total Number of Samples for Two Class Comparison

$\alpha$	$\beta$	$\delta$	$\sigma$	Samples Per Class
0.001	0.05	1 (2-fold)	0.5 human tissue	13
			0.25 transgenic mice	6 (t approximation)

# Sample Size Methods Also Developed for

- Balanced Block Designs
- For finding genes significantly associated with a survival outcome

# Class Comparison Paradigm

- Evaluate extent to which each gene is differentially expressed among classes
  - Univariate F-statistics, regularized F-statistics
- Select the most differentially expressed genes in a manner that limits the false discovery number or false discovery rate to a specified level

# t-test Comparisons of Gene Expression

- $x_j \sim N(\mu_{j1}, \sigma_j^2)$  for class 1
- $x_j \sim N(\mu_{j2}, \sigma_j^2)$  for class 2
- $H_{0j}: \mu_{j1} = \mu_{j2}$

# Estimation of Within-Class Variance

- Estimate separately for each gene
  - Limited degrees of freedom
  - Gene list dominated by genes with small fold changes and small variances
- Assume all genes have same variance
  - Poor assumption
- Random (hierarchical) variance model
  - Wright G.W. and Simon R. *Bioinformatics* 19:2448-2455, 2003
  - Inverse gamma distribution of residual variances
  - Results in exact F (or t) distribution of test statistics with increased degrees of freedom for error variance
  - For any normal linear model

# Simple Procedures for Controlling Multiple Comparisons

- Bonferroni method for controlling FEW
  - Probability of any false discoveries  $\leq 0.05$
- If each gene is tested for significance at level  $\alpha$  and there are  $G$  genes, then the expected number of false discoveries is  $G\alpha$ .
  - To control  $E(\text{FD}) \leq u$
  - Conduct each of  $G$  tests at level  $\alpha = u/G$
- Benjamini and Hochberg Method of Controlling the Expected False Discovery Rate

# Problems With Simple Procedures

- Bonferroni control of FWE is very conservative
- p values based on normal theory are not accurate at extremes quantiles
- Difficult to achieve extreme quantiles for permutation p values of individual genes
- Controlling *expected* number or proportion of false discoveries may not provide adequate control because distributions of FD and FDP may have large variances
- Methods do not take advantage of correlation among genes

# Multivariate Permutation Procedures

(Simon et al. 2003, Korn et al. 2004)

Allows statements like:

**FD Procedure:** We are 90% confident that the (actual) number of false discoveries is no greater than 5.

**FDP Procedure:** We are 90% confident that the (actual) proportion of false discoveries does not exceed .10.

# Control

$$\Pr\{\text{Number of FD} > k\} \leq \alpha$$

- Determine  $y = \alpha$  quantile of the distribution of the  $(k+1)$  st smallest p value under the multivariate permutation distribution.
- Include the genes corresponding to the  $k$  smallest p values in the gene list
- Include gene corresponding to  $p_{(i)}$  if  $p_{(i)} < y$

# Multivariate Permutation Procedures

- Permutation-based
  - Independent of distribution
  - even if they use t statistics
- Preserve/exploit correlation among tests by permuting each profile *as a unit*

# Multivariate Permutation Procedures

- More effective than univariate permutation tests especially with limited number of samples
  - Based on the  $\alpha$  percentile of the distribution of the  $(k+1)$ st smallest p value under multivariate permutation distribution; not on the  $\alpha/G$  percentile of the distribution of the univariate p value for a specific gene
- Stronger control than simple methods which control only expected number and proportion of false discoveries

# Control

$$\Pr\{\text{FDP} > \gamma\} \leq \alpha$$

- If you reject the null hypotheses for genes corresponding to  $p_{(1)}, \dots, p_{(i)}$  then the probability that the FDR is greater than  $\gamma$  equals the probability that there are more than  $\lfloor \gamma i \rfloor$  false discoveries in the list.
- This probability is  $\leq \alpha$  if you require  $p_{(i)} < y(\lfloor \gamma i \rfloor)$  where
- $y(u) = \alpha$  quantile of the distribution of the  $(u+1)$ st smallest p value under the multivariate permutation distribution.

# Control

$$\Pr\{\text{FDP} > \gamma\} \leq \alpha$$

- Determine  $y(u) = \alpha$  quantile of the distribution of the  $(u+1)$ st smallest  $p$  value under the multivariate permutation distribution.
  - For  $u = 1, 2, 3, \dots$
- Include in the list of differentially expressed genes the gene corresponding to the  $i$ 'th smallest  $p$  value as long as  $p_{(i)} < y(\lfloor \gamma i \rfloor)$ 
  - Sequentially for  $i = 1, 2, \dots$
  - $\lfloor \gamma i \rfloor =$  largest integer less than or equal to  $\gamma i$

# Class Prediction

- Most statistical methods were developed for inference, not prediction.
- Most statistical methods for were not developed for  $p \gg n$  settings

# Components of Class Prediction

- Feature (gene) selection
  - Which genes will be included in the model
- Select model type
  - E.g. DLDA, Nearest-Neighbor, ...
- Fitting parameters (regression coefficients) for model

# Feature Selection

- Genes that are univariately differentially expressed among the classes at a significance level  $\alpha$  (e.g. 0.01)
  - The  $\alpha$  level is selected to control the number of genes in the model, not to control the false discovery rate
    - Methods for class prediction are different than those for class comparison
  - The accuracy of the significance test used for feature selection is not of major importance as identifying differentially expressed genes is not the ultimate objective

# Feature Selection

- Small subset of genes which together give most accurate predictions
  - Combinatorial optimization algorithms
    - Genetic algorithms
- Little evidence that complex feature selection is useful in microarray problems
  - Failure to compare to simpler methods
  - Some published complex methods for selecting combinations of features do not appear to have been properly evaluated

# Linear Classifiers for Two Classes

$$l(\underline{x}) = \sum_{i \in F} w_i x_i$$

$\underline{x}$  = vector of log ratios or log signals

$F$  = features (genes) included in model

$w_i$  = weight for i'th feature

decision boundary  $l(\underline{x}) >$  or  $<$  d

# Linear Classifiers for Two Classes

- Fisher linear discriminant analysis

$$\underline{w} = \underline{d}' S^{-1}$$

- Diagonal linear discriminant analysis (DLDA) assumes features are uncorrelated
  - Naïve Bayes classifier
- Compound covariate predictor (Radmacher) and Golub's method are similar to DLDA in that they can be viewed as weighted voting of univariate classifiers

# Linear Classifiers for Two Classes

- Compound covariate predictor

$$w_i \propto \frac{\bar{x}_i^{(1)} - \bar{x}_i^{(2)}}{\hat{\sigma}_i}$$

Instead of for DLDA

$$w_i \propto \frac{\bar{x}_i^{(1)} - \bar{x}_i^{(2)}}{\hat{\sigma}_i^2}$$

# Linear Classifiers for Two Classes

- Support vector machines with inner product kernel are linear classifiers with weights determined to minimize errors
  - Can be written as finding hyperplane with separates the classes with a specified margin and minimizes length of weight vector
- Perceptrons are linear classifiers

# Support Vector Machine

$$\text{minimize } \sum_i w_i^2$$

$$\text{subject to } y_j (\underline{w}' \underline{x}^{(j)} + b) \geq 1$$

where  $y_j = \pm 1$  for class 1 or 2.

# When $p \gg n$

- For the linear model, an infinite number of weight vectors  $w$  can always be found that give zero classification errors for the training data.
  - $p \gg n$  problems are almost always linearly separable
- Why consider more complex models?
  - Fisher LDA is too complex

# Myth

- That complex classification algorithms such as neural networks perform better than simpler methods for class prediction.

- Artificial intelligence sells to journal reviewers and peers who cannot distinguish hype from substance when it comes to microarray data analysis.
- Most comparative studies indicate that simpler methods work as well or better for microarray problems

# Other Simple Methods

- Nearest neighbor classification
- Nearest centroid classification
- Shrunk centroid classification

- Fitting complex functions to training data results in unstable classifiers unless there is a huge training dataset
- Lack of stability is synonymous with *overfitting*
- For unstable classifiers, the test sample error rate is generally much less than the generalization error rate

# Model Stability Can Be Improved By

- Restriction to models with fewer parameters
  - Complexity depends on number of parameters per *candidate* feature, not per selected feature
- Reducing number of candidate features
  - Principal components of features
  - Centroids or pc's of clusters of features
- Not minimizing training error
  - Regularization; including penalty for complexity
- Aggregating models
  - Bagging
- Use fitting criterion incorporating robustness to changes in data

# Evaluating a Classifier

- Fit of a model to the same data used to develop it is no evidence of prediction accuracy for independent data.
- Demonstrating statistical significance of prognostic factors is not the same as demonstrating predictive accuracy.
- Demonstrating goodness of fit of a model to the data used to develop it is not a demonstration of predictive accuracy.
- Demonstrating stability of identification of gene predictors is not necessary for demonstrating predictive accuracy.

# Split-Sample Evaluation

- Training-set
  - Used to select features, select model type, determine parameters and cut-off thresholds
- Test-set
  - Withheld until a *single* model is *fully* specified using the training-set.
  - Fully specified model is applied to the expression profiles in the test-set to predict class labels.
  - Number of errors is counted
  - Ideally test set data is from different centers than the training data and assayed at a different time

# Leave-one-out Cross Validation

- Omit sample 1
  - Develop multivariate classifier from scratch on training set with sample 1 omitted
  - Predict class for sample 1 and record whether prediction is correct

# Leave-one-out Cross Validation

- Repeat analysis for training sets with each single sample omitted one at a time
- $e$  = number of misclassifications determined by cross-validation
- Subdivide  $e$  for estimation of sensitivity and specificity

# Myth

- Cross-validation of a model can occur after selecting the genes to be used in the model

- Cross validation is only valid if the test set is not used in any way in the development of the model. Using the complete set of samples to select genes violates this assumption and invalidates cross-validation.
- With proper cross-validation, the model must be developed *from scratch* for each leave-one-out training set. This means that feature selection must be repeated for each leave-one-out training set.
- The cross-validated estimate of misclassification error is an estimate of the prediction error for model fit using specified algorithm to full dataset
- If you use cross-validation estimates of prediction error for a set of algorithms indexed by a tuning parameter and select the algorithm with the smallest cv error estimate, you do not have a valid estimate of the prediction error for the selected model

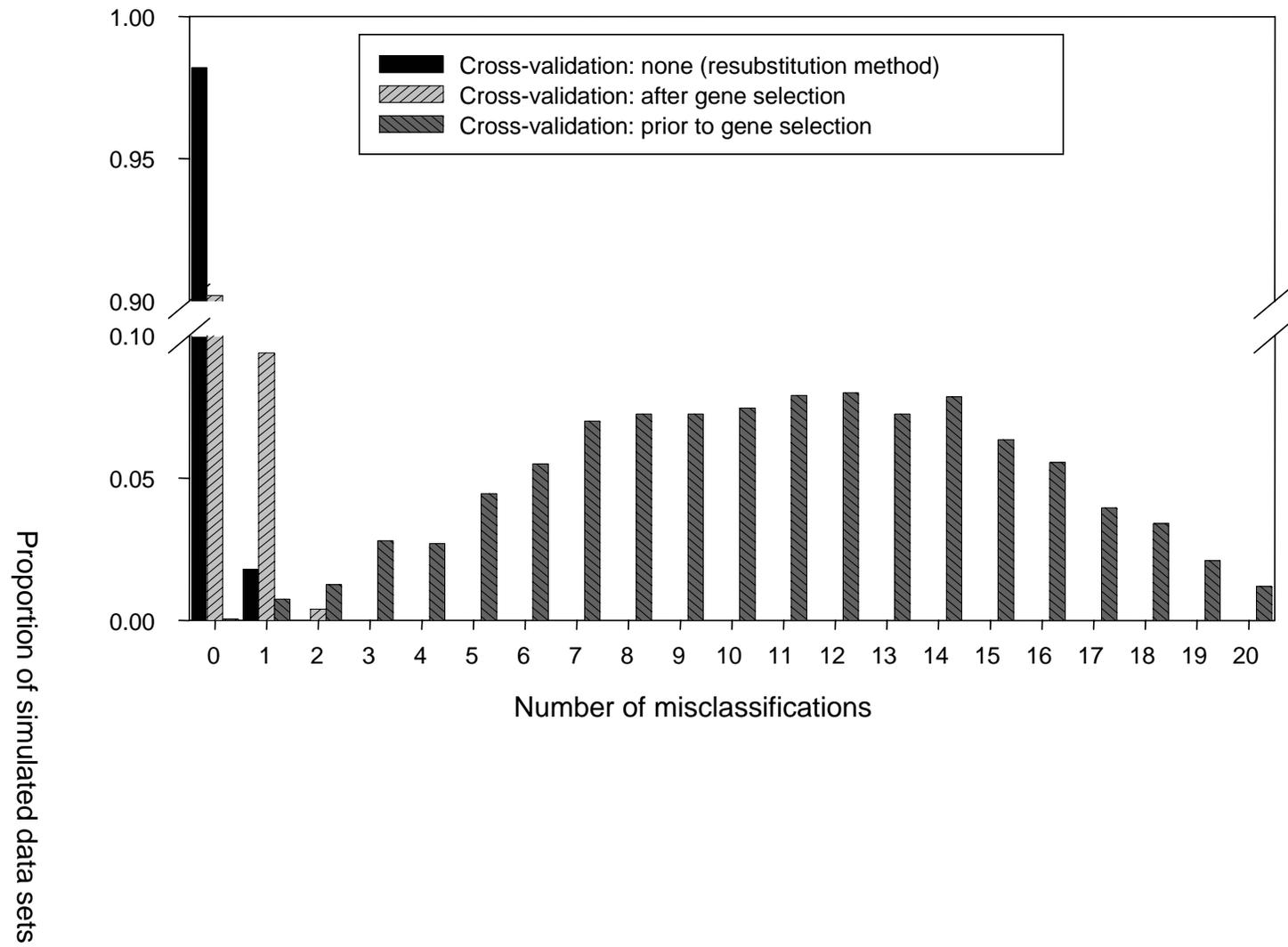
# Prediction on Simulated Null Data

## Generation of Gene Expression Profiles

- 14 specimens ( $P_i$  is the expression profile for specimen  $i$ )
- Log-ratio measurements on 6000 genes
- $P_i \sim \text{MVN}(\mathbf{0}, \mathbf{I}_{6000})$
- Can we distinguish between the first 7 specimens (Class 1) and the last 7 (Class 2)?

## Prediction Method

- Compound covariate prediction (*discussed later*)
- Compound covariate built from the log-ratios of the 10 most differentially expressed genes.



# Permutation Distribution of Cross-validated Misclassification Rate of a Multivariate Classifier

- Randomly permute class labels and repeat the entire cross-validation
- Re-do for all (or 1000) random permutations of class labels
- Permutation p value is fraction of random permutations that gave as few misclassifications as  $e$  in the real data

# Invalid Criticisms of Cross-Validation

- “You can always find a set of features that will provide perfect prediction for the training and test sets.”
  - For complex models, there may be many sets of features that provide zero training errors.
  - A modeling strategy that either selects among those sets or aggregates among those models, will have a generalization error which will be validly estimated by cross-validation.

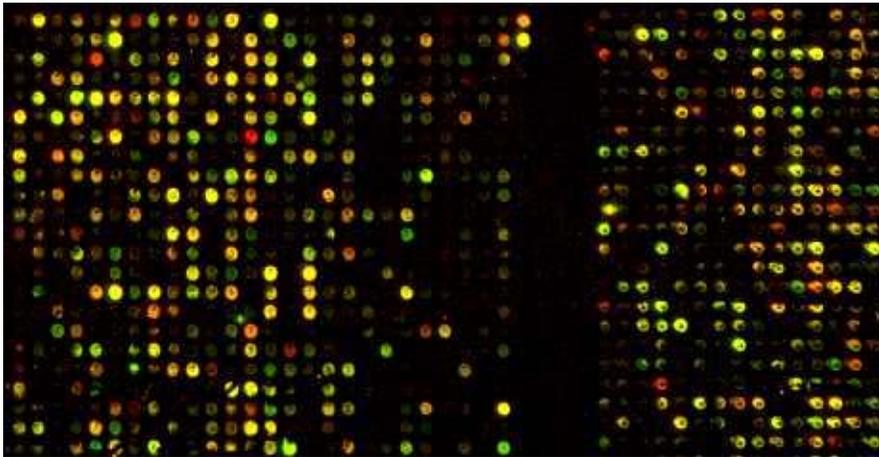
# Sources of Bias in Estimation of Error Rates

- Confounding by sample handling or assay effects
  - Cases collected and assayed at different times than controls
- Failure to incorporate important sources of future variability
  - Assay drift
- Change in distribution of unmodeled variables
  - In split sample validation, split samples by institution

# Gene-Expression Profiles in Hereditary Breast Cancer

## cDNA Microarrays

### *Parallel Gene Expression Analysis*



- Breast tumors studied:
  - 7 *BRCA1*+ tumors
  - 8 *BRCA2*+ tumors
  - 7 sporadic tumors
- Log-ratios measurements of 3226 genes for each tumor after initial data filtering

## RESEARCH QUESTION

Can we distinguish *BRCA1*+ from *BRCA1*- cancers and *BRCA2*+ from *BRCA2*- cancers based solely on their gene expression profiles?

# BRCA1

$\alpha_g$	# of significant genes	# of misclassified samples ( $m$ )	% of random permutations with $m$ or fewer misclassifications
$10^{-2}$	182	3	0.4
$10^{-3}$	53	2	1.0
$10^{-4}$	9	1	0.2

# BRCA2

$\alpha_g$	# of significant genes	$m = \#$ of misclassified elements (misclassified samples)	% of random permutations with $m$ or fewer misclassifications
$10^{-2}$	212	4 (s11900, s14486, s14572, s14324)	0.8
$10^{-3}$	49	3 (s11900, s14486, s14324)	2.2
$10^{-4}$	11	4 (s11900, s14486, s14616, s14324)	6.6

# Classification of BRCA2 Germline Mutations

Classification Method	LOOCV Prediction Error
Compound Covariate Predictor	14%
Fisher LDA	36%
Diagonal LDA	14%
1-Nearest Neighbor	9%
3-Nearest Neighbor	23%
Support Vector Machine (linear kernel)	18%
Classification Tree	45%

# Selected Features of BRB-ArrayTools

- Multivariate permutation tests for class comparison to control false discovery proportion with any specified confidence level
- Find Gene Ontology groups and signaling pathways that are differentially expressed
- Survival analysis
- Analysis of variance
- Class prediction models (6) with prediction error estimated by LOOCV, k-fold CV or .632 bootstrap, and permutation analysis of cross-validated error rate
  - DLDA, SVM, CCP, Nearest Neighbor, Nearest Centroid, Shrunken Centroids, Random Forests
- Clustering tools for class discovery with reproducibility statistics on clusters
- Visualization tools including rotating 3D principal components plot exportable to Powerpoint with rotation controls
- Extensible via R plug-in feature
- Links genes to annotations in genomic databases
- Tutorials and datasets

# Acknowledgements

- Kevin Dobbin
- Sudhir Varma
- Ed Korn
- Lisa McShane
- Michael Radmacher
- George Wright
- Yingdong Zhao
- Amy Peng Lam