

Statistical Analysis of Gene Expression Microarray Data

Lisa M. McShane, Ph.D.

Biometric Research Branch
National Cancer Institute

Class Overview

- Day 1: Discussion of statistical analysis of microarray data – Lisa M. McShane
- Day 2: Hands-on BRB ArrayTools workshop – Supriya Menezes

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

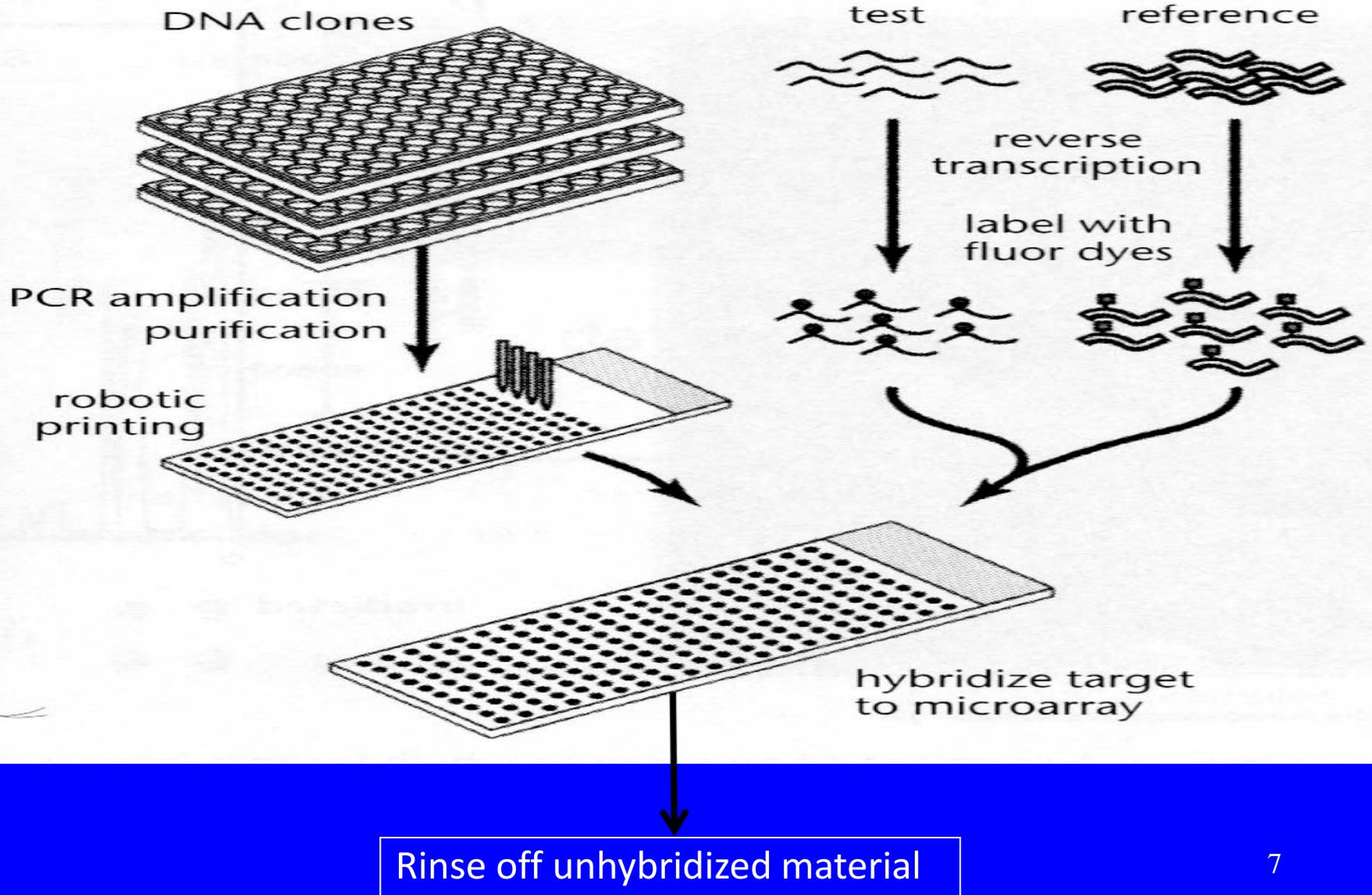
Gene Expression Microarrays

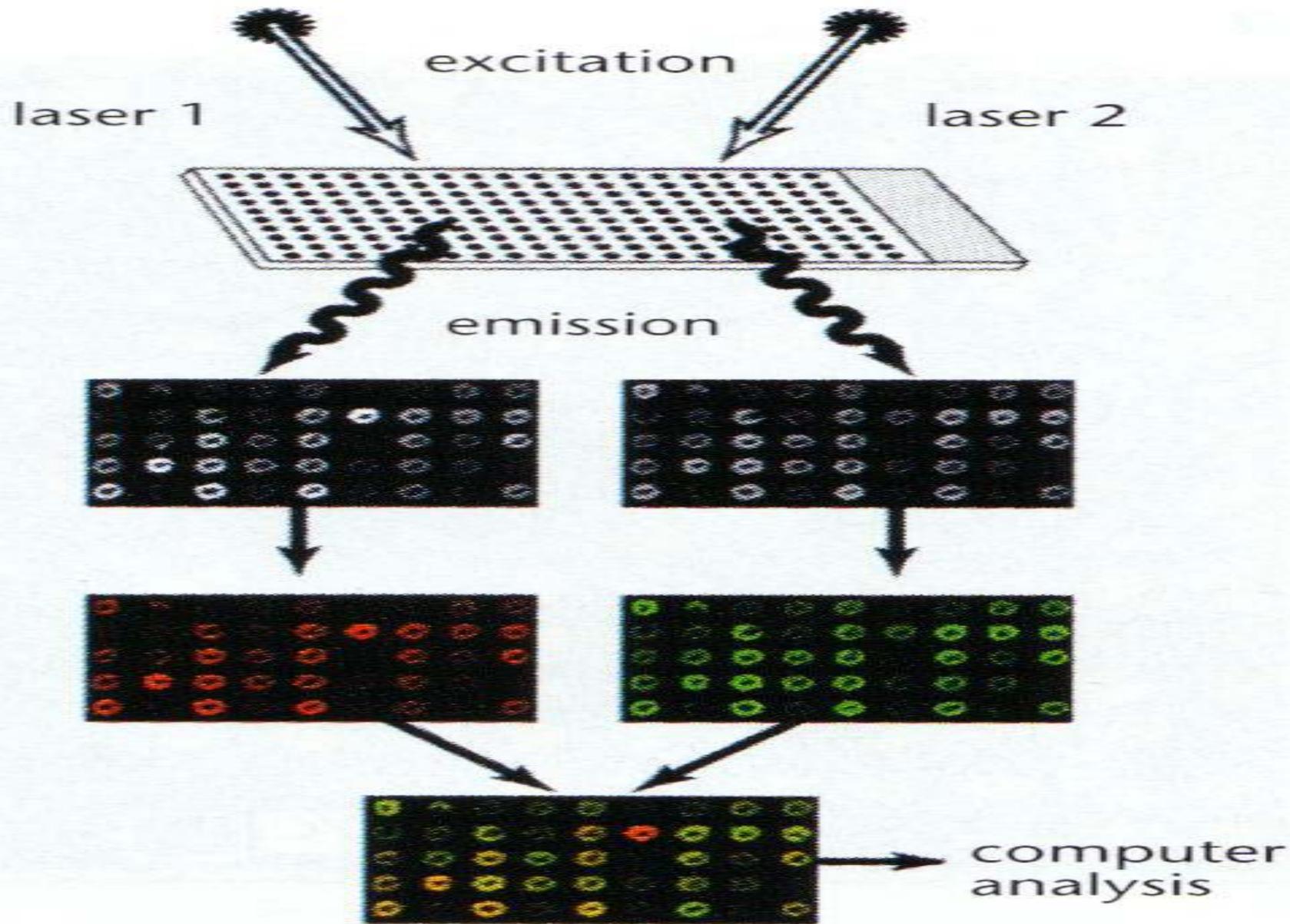
- Permit simultaneous evaluation of expression levels of thousands of genes
- Main platforms
 - Spotted cDNA arrays (2-color)
 - Affymetrix GeneChip (1-color)
 - Spotted oligo arrays (2-color or 1-color)
 - Bead arrays (e.g., Illumina-DASL)
 - Nylon filter arrays

Spotted cDNA Arrays (and other 2-color spotted arrays)

- cDNA arrays: Schena *et al.*, *Science*, 1995
- Each spot corresponds to a gene (sometimes multiple spots per gene)
- Two-color (two-channel) system
 - Two colors represent the two samples competitively hybridized
 - Each spot has “red” and “green” measurements associated with it

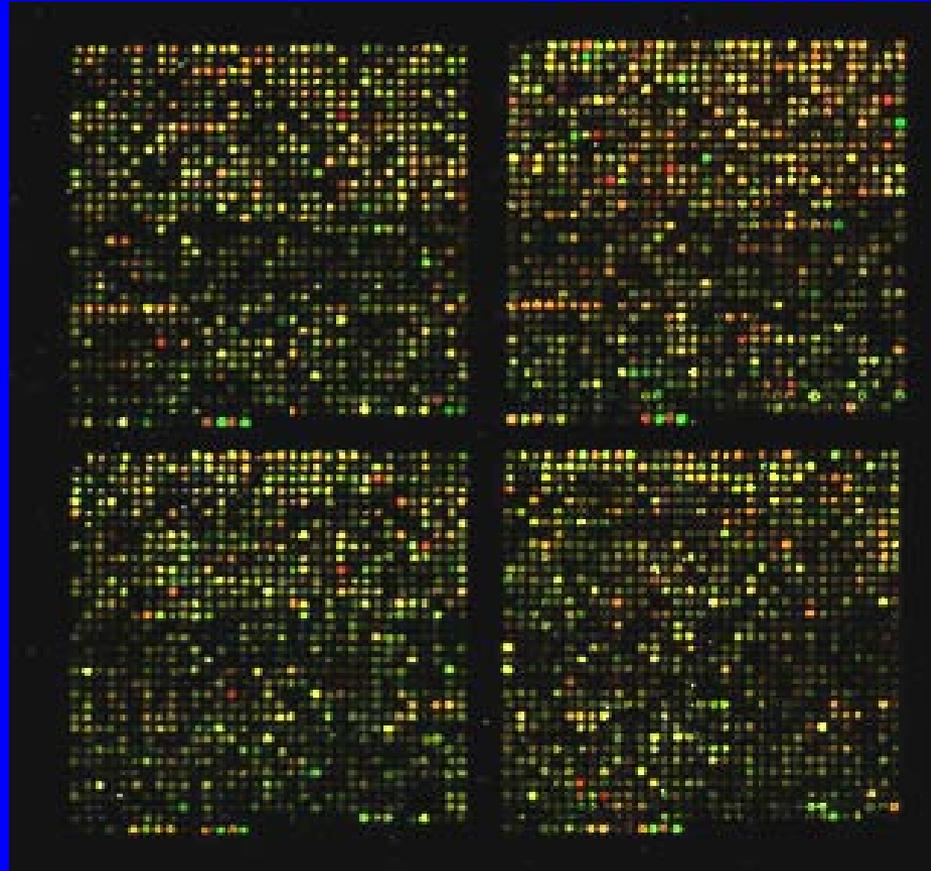
cDNA Array





cDNA Microarray Image

(overlaid “red” and “green” images)



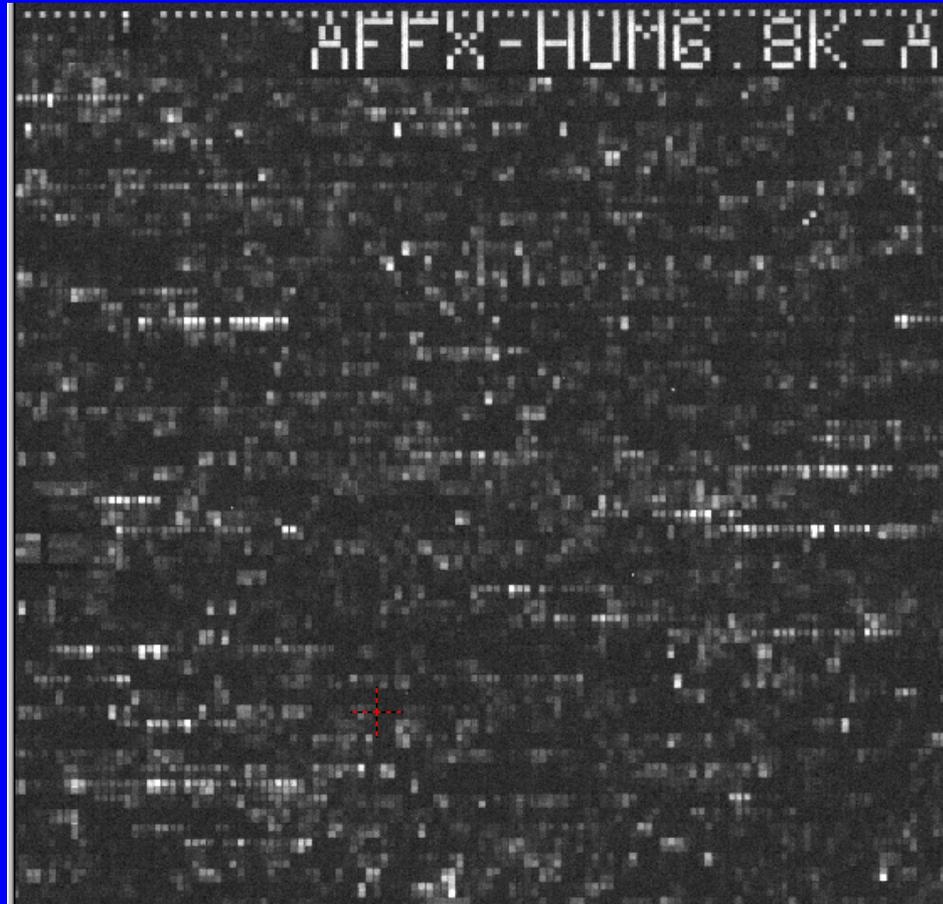
Affymetrix GeneChip Arrays

- Lockhart *et al.*, *Nature Biotechnology*, 1996
- Affymetrix: <http://www.affymetrix.com>
- Glass wafer (“chip”) – photolithography, oligonucleotides synthesized on chip
- Single sample hybridized to each array
- Each gene represented by one or more “probe sets”
 - One probe type per array “cell”
 - Typical oligo probe is 25 nucleotides in length
 - 11-20 PM:MM pairs per probe set (PM = perfect match, MM = mismatch)

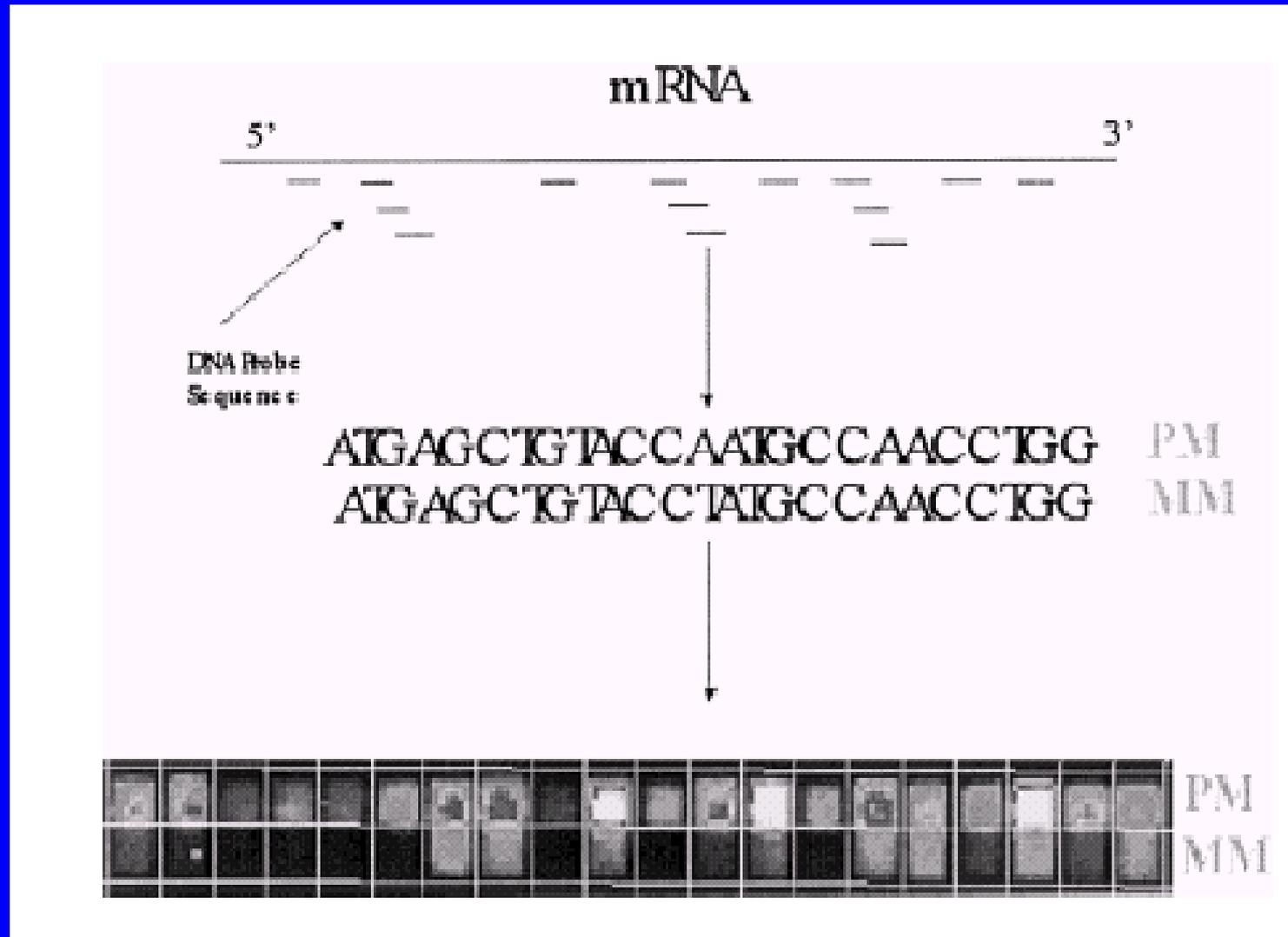
[Affymetrix] Hybridization Oligo “GeneChip” Array



Image of a Scanned Affymetrix GeneChip



Perfect Match - Mismatch Probe Pairs

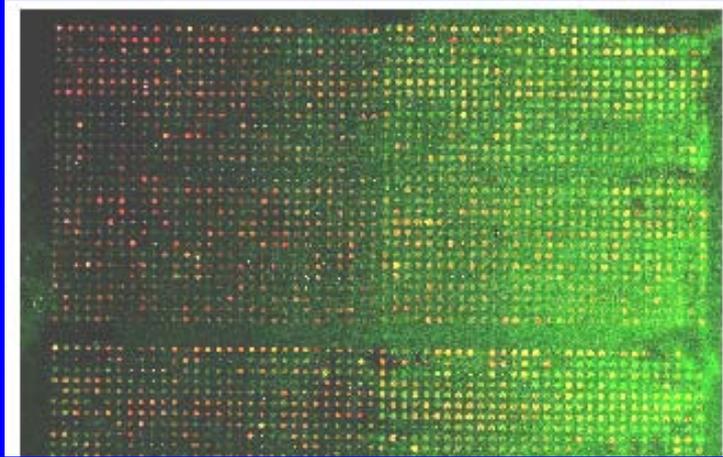


(Figure 2 from Schadt *et al.*, *Journal of Cellular Biochemistry*, 2001) 13

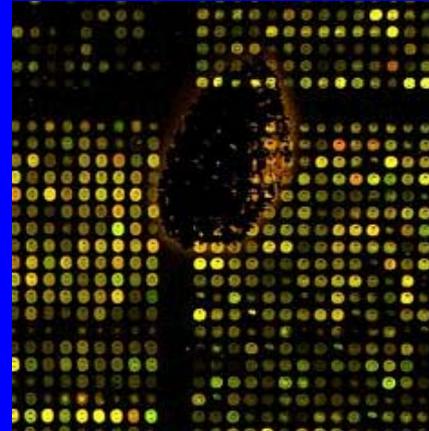
Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

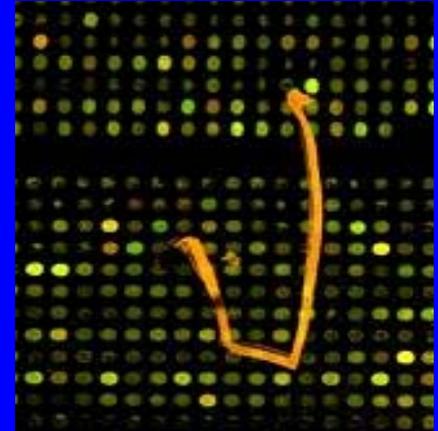
cDNA/Spotted Arrays: QC Issues



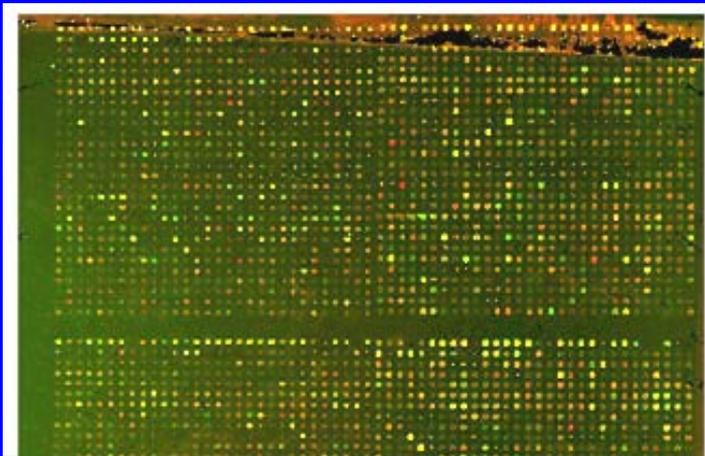
Background haze



Bubble



Scratch
or fiber?

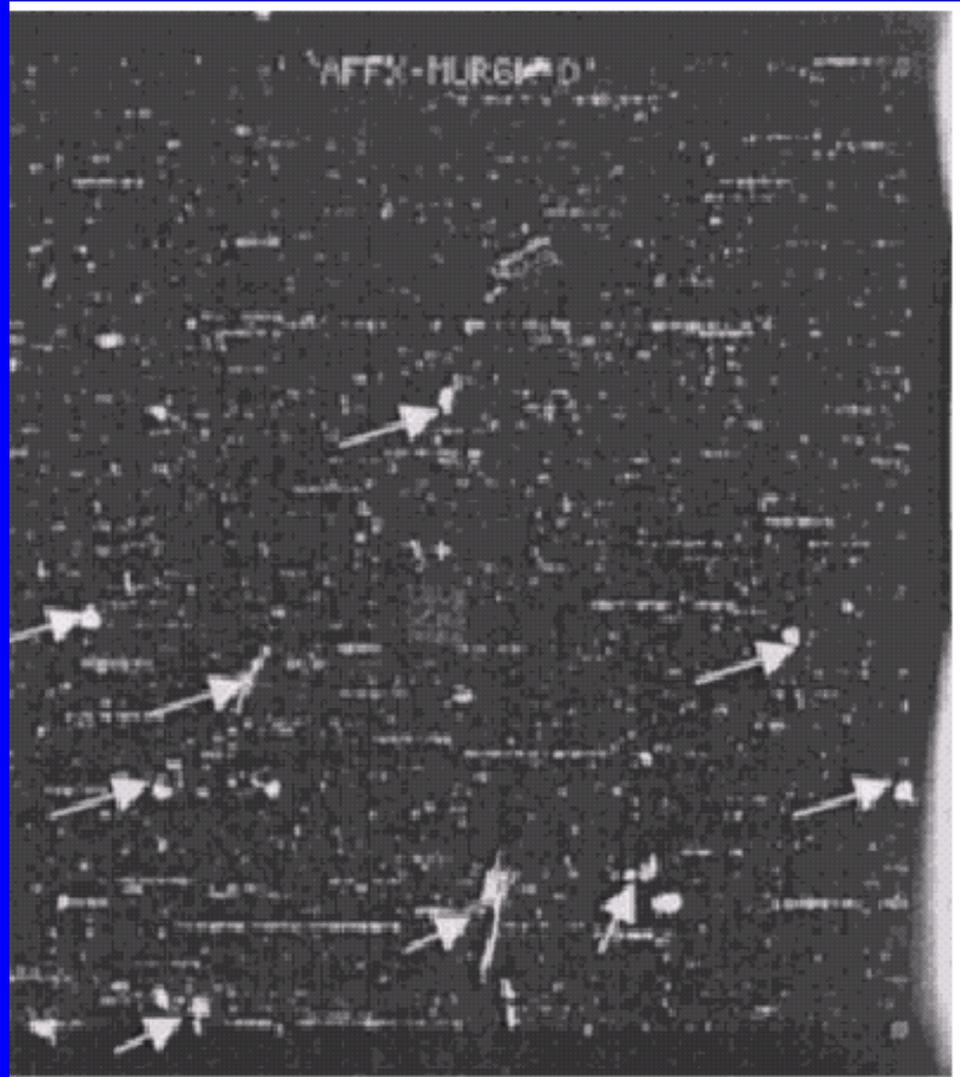


Edge effect

- Visual inspection of arrays advisable
- Danger: Garbage in/Garbage out

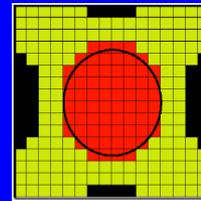
Affymetrix Arrays: Quality Problems

(Figure 1 from Schadt *et al.*, *Journal of Cellular Biochemistry*, 2001)



cDNA/2-color Spotted Arrays: Segmentation

- Segmentation - separation of feature (F) from background (B) for each spot.



(See software documentation)

- Summary measures computed for F (for each channel/color)
 - Intensity: mean or median over pixels
 - Additional measures: SD, # pixels (size), etc.

cDNA/2-color Spotted Arrays: Background Correction & Signal Calculation (for each channel/color)

- No background correction
Signal = F intensity
- Local background correction
Signal = F intensity - B_{local}
- Regional background correction
Signal = F intensity - B_{regional}

cDNA/2-color Spotted Arrays: Flagging Spots/arrays Exclusion

Exclude spots if
“signal” or “signal-to-
noise” measure(s)
poor (low):

- F
- F-B
- $(F-B)/SD(B)$
- Spot Size

Exclude whole arrays
or regions of arrays:

- Too many spots
flagged
- Narrow range of
intensities
- Uniformly low signals

cDNA/2-color Spotted Arrays: Gene-level Summaries

- Model-based methods
 - Work directly on signals from two channels
 - Color effects and interactions between color and experimental factors incorporated into statistical models
- Ratio methods*
 - Red signal/Green signal
 - “Green” sample serves as internal reference

* Today's lecture will focus on ratio methods

Affymetrix Arrays: Image Processing

- DAT image files → CEL files
- Each probe cell: 10x10 pixels
- Grid alignment to probe cells
- Signals:
 - Remove outer 36 pixels → 8x8 pixels
 - The probe cell signal, PM or MM, is the 75th percentile of the 8x8 pixel values
- Background correction: Average of the lowest 2% probe cell values in zone is taken as the background value and subtracted
- Summarize over probe pairs to get gene expression indices
 - Detection calls - present/absent

See Affymetrix documentation:

- Affymetrix website (<http://www.affymetrix.com>)
- *Affymetrix Microarray Suite User Guide*
- *Affymetrix Statistical Algorithms Description Document*

Affymetrix Arrays: Probe Set (Gene) Summaries

- $AvDiff_i = \Sigma(PM_{ij} - MM_{ij}) / n_i$ for each probe set i
(original Affymetrix algorithm)
- Revised Affymetrix algorithm to address negative signals (MAS 5.x series)
 - anti-log of a robust average (Tukey biweight) of the $\log(PM_{ij} - IM_{ij})$, where
 $IM = MM$, if $MM < PM$
= adjusted to be less than PM , if $MM \geq PM$

Affymetrix Arrays:

Model-based Probe Set (Gene) Summaries

- Li and Wong (*PNAS*, 2001; *Genome Biology*, 2001; incorporated into dChip)
 - $MBEI_i = \theta_i$ estimated from $PM_{ij} - MM_{ij} = \theta_i \phi_j + \varepsilon_{ij}$
=> weighted average difference
(ϕ_j = probe j sensitivity index, ε_{ij} = random error)
 - $MBEI_i^* = \theta_i^*$ estimated from $PM_{ij} = v_i + \theta_i^* \phi_j$:
probe set summaries are based on PM signals only.
(v_i = baseline response of probe *pair*, ϕ_j = probe j sensitivity index)

Affymetrix Arrays:

Model-based Probe Set (Gene) Summaries

(continued)

- Irizarry *et al.* (*Nucleic Acids Research*, 2003; *Biostatistics*, 2003)
 - $RMA_i = e_i$ estimated from $T(PM_{ij}) = e_i + a_j + \varepsilon_{ij}$, where $T(PM)$ represents the PM intensities which have been cross-hybridization corrected, (quantile-) normalized and log-transformed
- Wu, Irizarry, Gentleman, Murillo, Spencer (*J. Amer. Stat. Assoc.*, 2004)
 - Apply cross-hybridization correction that depends on G-C content of probe.

Affymetrix Arrays: Comparison of Cross-hybridization Corrections

- Affymetrix MAS 5.x series: Estimate cross-hybridization (CH) using MM probes for the gene
- RMA: Some target hybridizes to the MM probe; for high expressed genes MM is brighter than true cross-hybridization; use smaller CH estimate (approx. mode of the MM probes across all MM probe sets)

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

cDNA/2-color Spotted Arrays: Need for Normalization

- Unequal incorporation of labels
 - green brighter than red
- Unequal amounts of sample
- Unequal PMT voltage
- Autofluorescence greater at shorter scanning wavelength

Normalization Methods for cDNA/2-Color Spotted Arrays

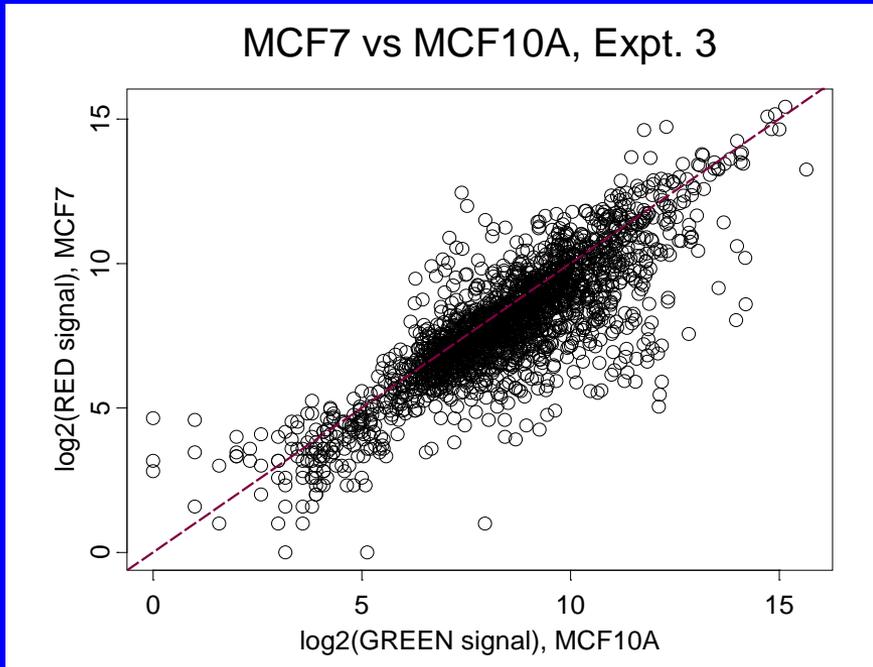
- Model-based methods
 - Normalization incorporated into model
- Ratio-based methods
 - Median (or Mean) Centering Method
 - Lowess Method
 - Multitude of other methods

Chen et al., Journal of Biomedical Optics, 1997

Yang et al., Nucleic Acids Research, 2002

- Scaling factors, separately by printer pin, etc.

Median (or Mean) Centering

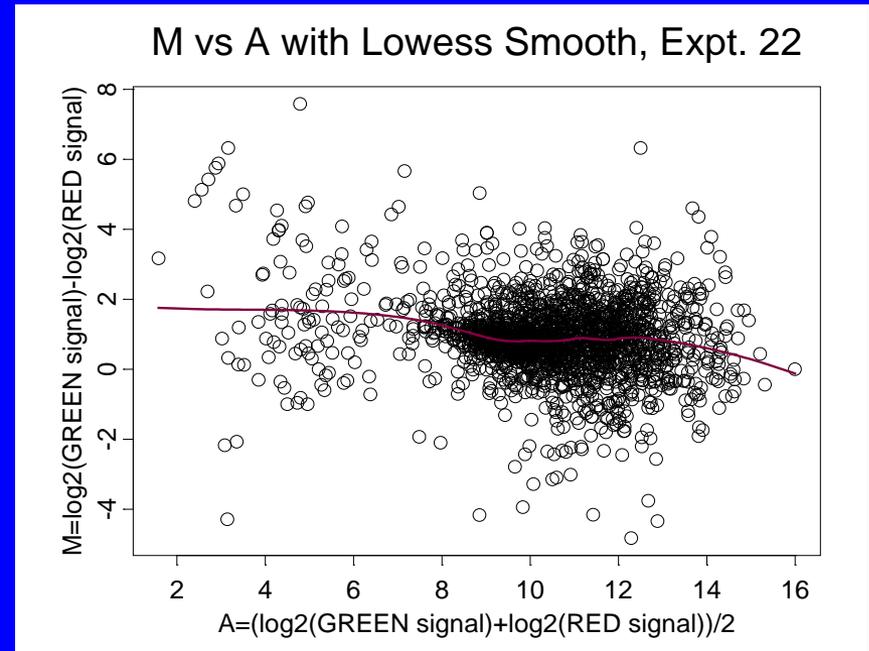
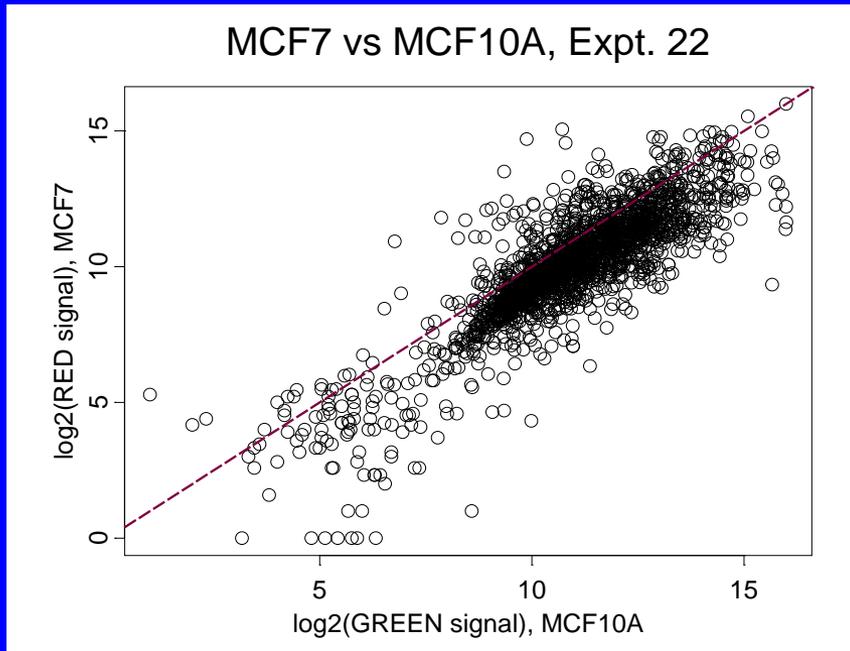


In plot of $\log(\text{red signal})$ versus $\log(\text{green signal})$, if point scatter is parallel to 45° line, adjust intercept to 0.

Subtract median or mean log-ratio (computed over all genes on the slide or only over housekeeping genes) from each log-ratio.

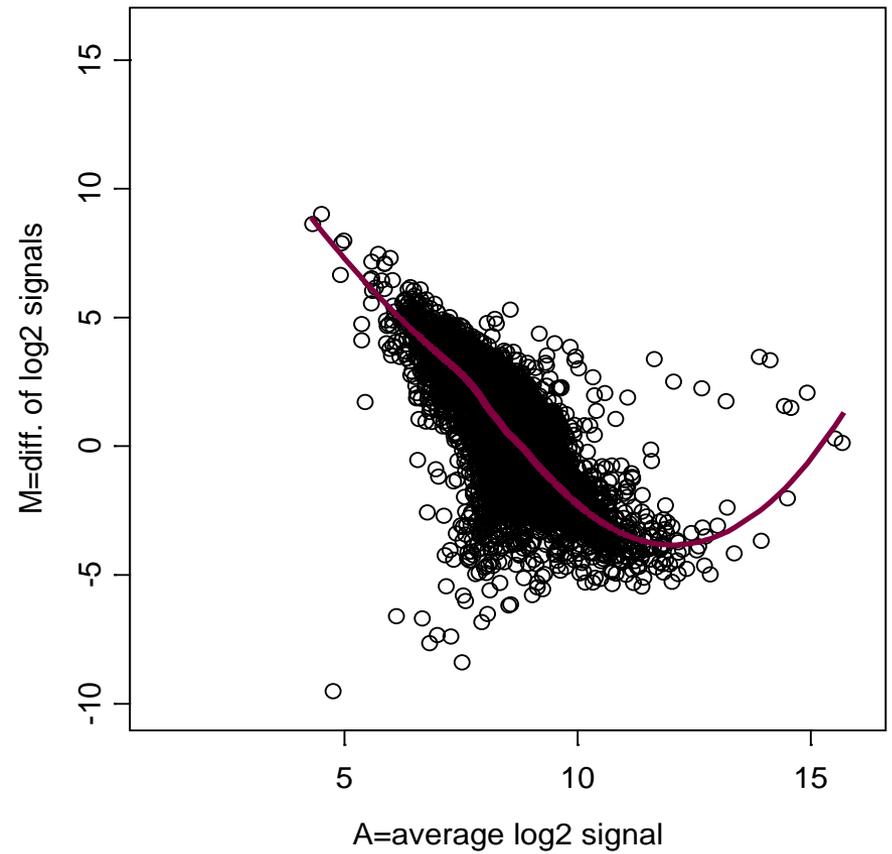
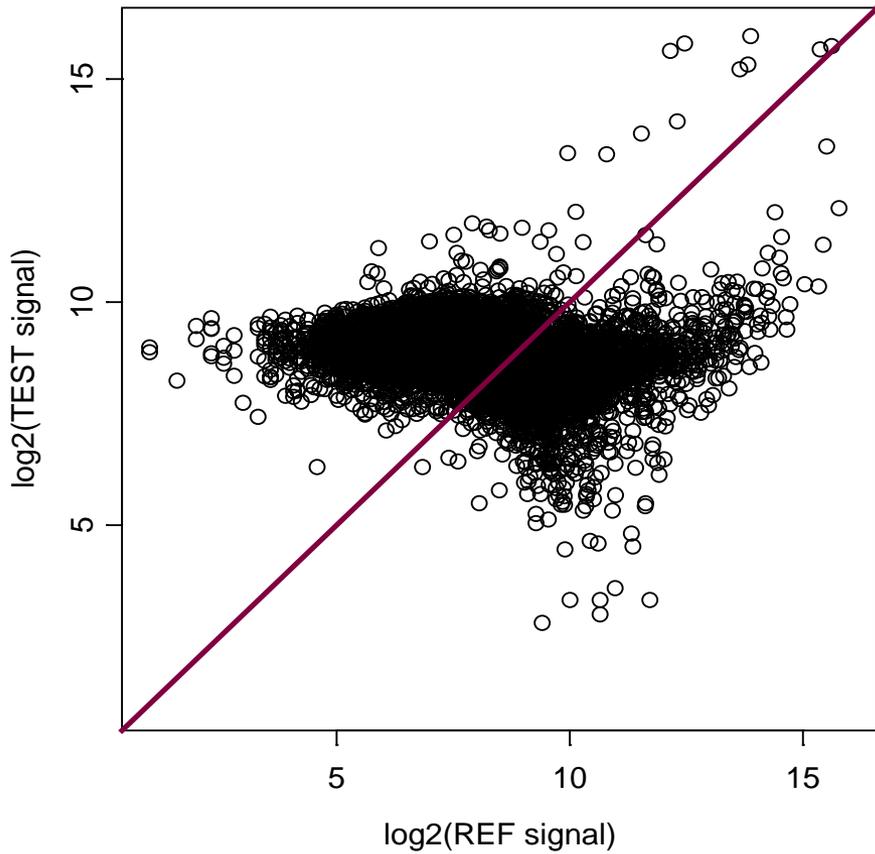
Lowess Normalization: M vs A Plots

Yang *et al.*, *Nucleic Acids Research*, 2002



$$M = \log_2(\text{GREEN signal}) - \log_2(\text{RED signal})$$
$$A = (\log_2(\text{GREEN signal}) + \log_2(\text{RED signal})) / 2$$

Bad Array Example



Normalization: Affymetrix Arrays

- Variations due to sample, chip, hybridization, scanning
- Probe set-level vs probe-level
- Quantile normalization, intensity-dependent, etc.
- Normalize across all arrays or pairwise
- PM-MM vs PM only
- Built in to dChip, RMA, and MAS 5.x series algorithms
 - Li and Wong (*PNAS*, 2001; *Genome Biology*, 2001)
 - Irizarry *et al.* (*Nucleic Acids Research*, 2003; *Biostatistics*, 2003)
 - Bolstad *et al.* (*Bioinformatics*, 2003)

Filtering Genes

- “Bad” or missing values on too many arrays
- Not differentially expressed across arrays (non-informative)
 - Variance (assumes approx. normality)

s^2_i = sample variance of gene i (log) measurements across n arrays.

Exclude gene i if (gene has smaller var than median)

$(n-1) s^2_i < \chi^2(\alpha, n-1) \times \text{median}(s^2_1, s^2_2, \dots, s^2_K)$, K = number of genes.
 - Fold difference

Max/Min < 3 or 4, (95th percentile/5th percentile) < 2 or 3

Filter if $k\%$ of genes have FC < 2 or 3 relative to median

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

Design and Analysis Methods Should Be Tailored to Study Objectives

- Class Comparison (supervised)
 - For predetermined classes, establish whether gene expression profiles differ, and identify genes responsible for differences
- Class Discovery (unsupervised)
 - Discover clusters among specimens or among genes
- Class Prediction (supervised)
 - Prediction of phenotype using information from gene expression profile

Class Comparison Examples

- Establish that expression profiles differ between two histologic types of cancer.
- Identify genes whose expression level is altered by exposure of cells to an experimental drug.

Class Discovery Examples

- Discover previously unrecognized subtypes of lymphoma.
- Cluster temporal gene expression patterns to get insight into genetic regulation in response to a drug or toxin.

Class Prediction Examples

- Predict from expression profiles which patients are likely to experience severe toxicity from a new drug versus who will tolerate it well.
- Predict which breast cancer patients will relapse within two years of diagnosis versus who will remain disease free.

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Expression Measures, Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

Analysis Strategies for Class Comparisons

- Global tests
 - Compare whole profiles
 - Permutation tests
- Gene-level analyses
 - Model-based methods (e.g., multi-parameter)
 - Non-model-based methods (e.g., t-tests, F-tests, nonparametric tests)
 - Hybrid variance methods

Global Tests for Differences in Profiles Between Classes

- Choice of summary measure of difference
 - Examples:
 - Sum of squared univariate t-statistics
 - Number of genes univariately significant at 0.001 level
- Statistical testing by permutation test
- BRB-ArrayTools uses the number of univariately significant genes as a summary measure for the global test for differences between profiles.

Summary of Results:

Number of genes significant at 0.001 level of the univariate test: 52

Probability of getting at least 52 genes significant by chance (at the 0.001 level) if there are no real differences between the classes: 0.001

Gene-Level Analyses

- Model-based methods*
 - Multi-parameter modeling of channel-level data (e.g., Gaussian mixed or ANOVA models), hierarchical Bayesian models, etc.
 - May borrow information across genes
 - May use multiple comparison adjustments
- Non-model-based methods
 - Log ratios or signal (e.g., Affymetrix)
 - T-test, F-test, or nonparametric counterparts (e.g., Wilcoxon)
 - Multiple comparison adjustment commonly used
- Random variance methods
 - Variance estimates borrow across genes

* Not discussed in today's lecture

Random Variance Methods for Small Sample Gene-Level Analyses

- Bayesian
 - Baldi and Long, *Bioinformatics*, 2001
- Frequentist
 - Wright and Simon, *Bioinformatics*, 2003
 - Available as the 'Random variance' option in BRB-ArrayTools for Class Comparison and Class Prediction analyses

Variance model: _____



Use randomized variance model for univariate tests.

Multiple Testing Procedures for Gene-Level Analyses

Identification of differentially expressed genes while controlling for false discoveries (genes declared to be differentially expressed that in truth are not)

- *Actual Number* of False Discoveries: FD
- *Expected Number* of False Discoveries: $E(\text{FD})$
- *Actual Proportion* of False Discoveries: FDP
- *Expected Proportion* of False Discoveries:
 $E(\text{FDP}) = \text{False Discovery Rate (FDR)}$

Simple Procedures

- Control expected number of false discoveries
 - $E(\text{FD}) \leq u$
 - Conduct each of k tests at level u/k
- Bonferroni control of familywise error (FWE) rate at level α
 - Conduct each of k tests at level α/k
 - At least $(1-\alpha)100\%$ confident that $\text{FD} = 0$

Problems With Simple Procedures

- Bonferroni control of FWE is very conservative
- Controlling *expected* number or proportion of false discoveries may not provide adequate control on *actual* number or proportion

Additional Procedures

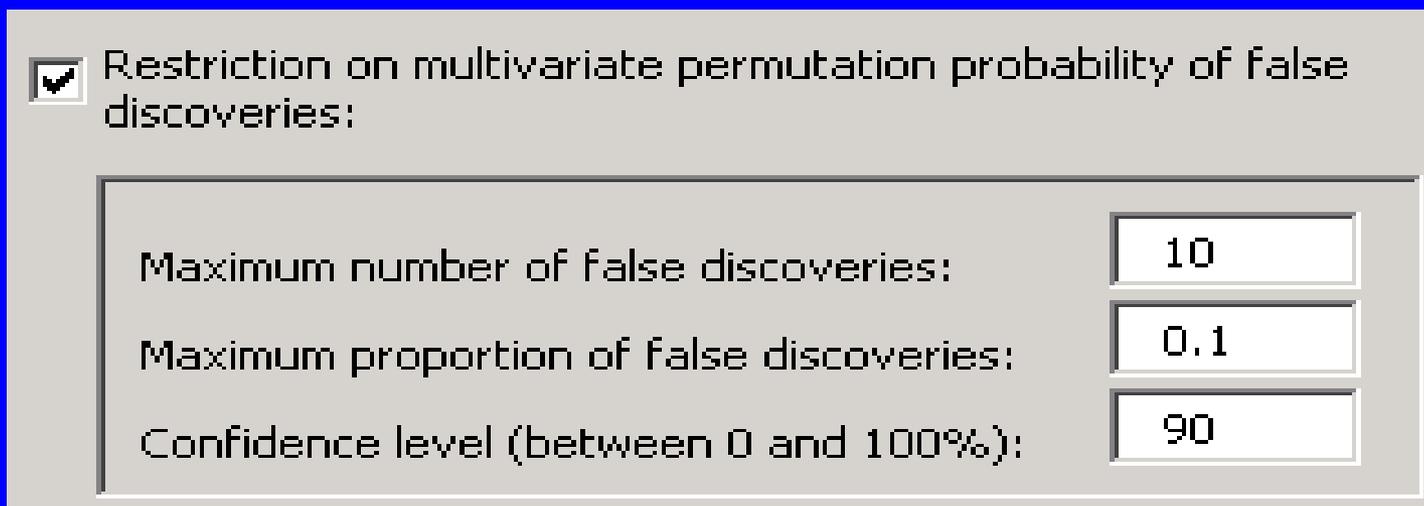
- Review by Dudoit *et al.* (*Statistical Science*, 2003)
- “SAM” - Significance Analysis of Microarrays
 - Tusher *et al.*, *PNAS*, 2001 and relatives
 - Estimate quantities similar to FDR (old SAM) or control FDP (newer versions of SAM)
- Bayesian
 - Efron *et al.*, *JASA*, 2001; *Stanford Tech Rep*, 2001
 - Manduchi *et al.*, *Bioinformatics* 2000
 - Newton *et al.*, *J Comp Biology* 2001
- Step-down permutation procedures
 - Westfall and Young, 1993 Wiley (FWE)
 - Korn *et al.*, *JSPI*, 2004 (FD and FDP control)

Examples of Types of Control

- Korn *et al.* FD procedure: “We are 95% confident that the (actual) number of false discoveries is no greater than 2.”
- Korn *et al.* FDP procedure: “We are 95% confident that the (actual) proportion of false discoveries does not exceed approximately 0.10.”
- Tusher *et al.* SAM: “On *average*, the false discovery proportion will be controlled at approximately 10%.”
- Current SAM – more similar to Korn FDP procedure
- Bayesian methods: “High posterior probability of differential expression”

Multiple Testing Procedures Available in BRB-ArrayTools

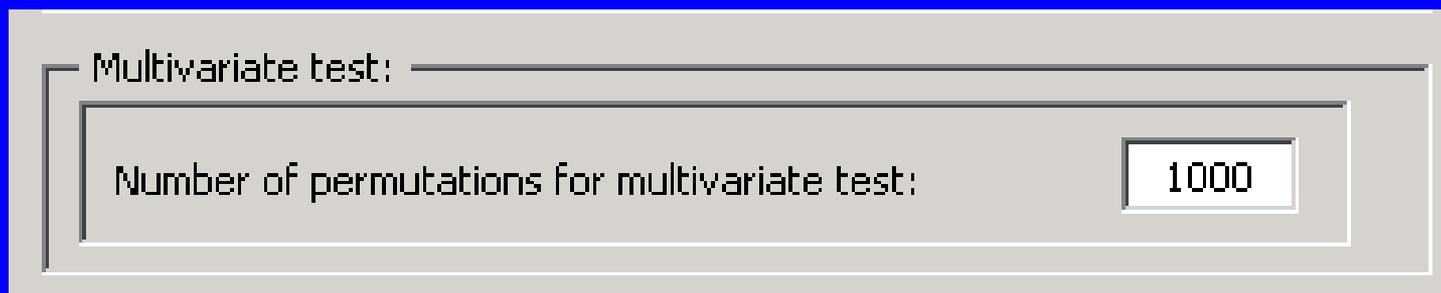
- The step-down permutation procedure for FD and FDP control (Korn, *et al.*) is available in BRB-ArrayTools for Class Comparison, Survival Analysis, and Quantitative Traits Analysis (finding genes significantly correlated with a quantitative variable). The following screenshot from the analysis dialog shows the default options:



Restriction on multivariate permutation probability of false discoveries:

Maximum number of false discoveries:	10
Maximum proportion of false discoveries:	0.1
Confidence level (between 0 and 100%):	90

- The number of permutations may be specified on the 'Options' page:



Multivariate test: _____

Number of permutations for multivariate test:	1000
---	------

Class Discovery

- Cluster analysis algorithms
 - Hierarchical
 - K-means
 - Self-Organizing Maps
 - Maximum likelihood/mixture models
 - Multitude of others
- Graphical displays
 - Hierarchical clustering
 - Dendrogram
 - “Ordered” color image plot (heatmap)
 - Multidimensional scaling plot

Hierarchical Agglomerative Clustering Algorithm

- Cluster genes with respect to expression across specimens
- Cluster specimens with respect to gene expression profiles
 - Filter genes that show little variation across specimens
 - Median or mean center genes

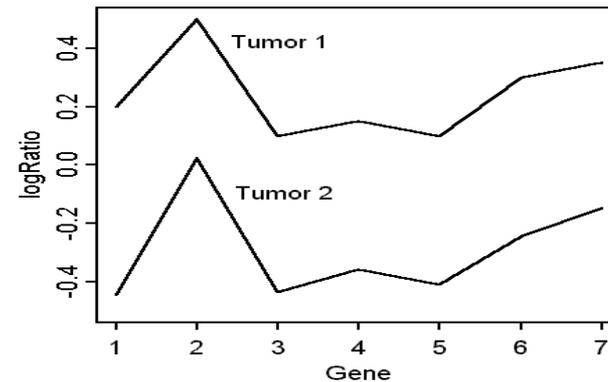
Hierarchical Agglomerative Clustering Algorithm

- Merge two closest observations into a cluster.
 - How is distance between individual observations measured?
- Continue merging closest clusters/observations.
 - How is distance between clusters measured?
 - Average linkage
 - Complete linkage
 - Single linkage

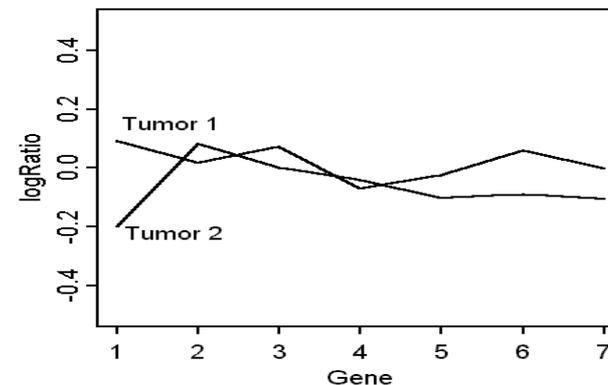
Common Distance Metrics for Hierarchical Clustering

- Euclidean distance
 - Measures absolute distance (square root of sum of squared differences)
- 1-Correlation
 - Large values reflect lack of linear association (pattern dissimilarity)

Euclidean distance large, 1-Correlation small

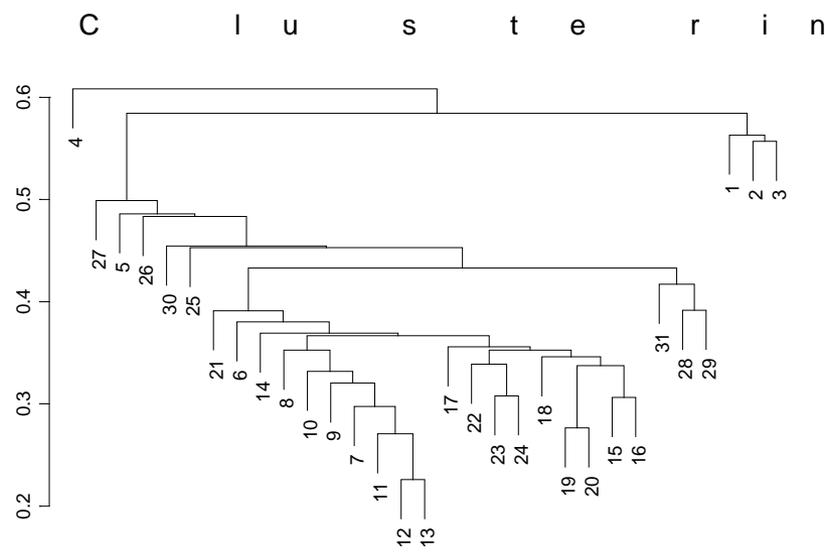
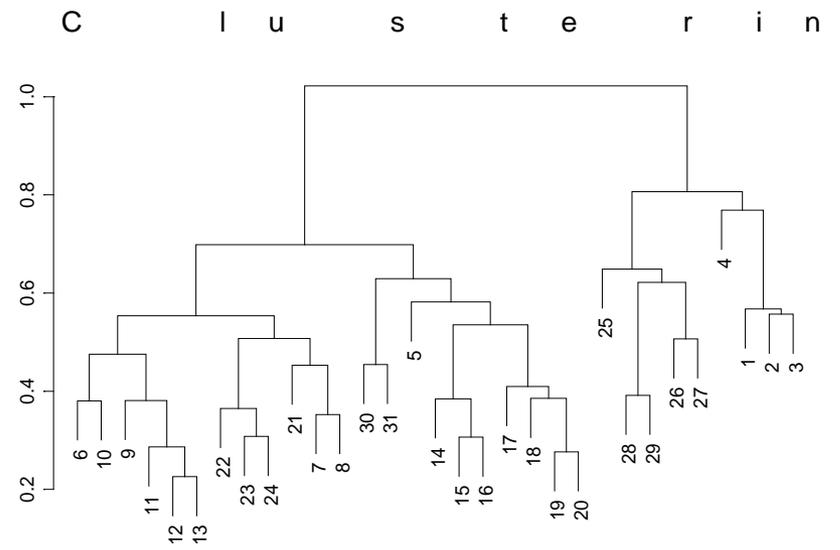
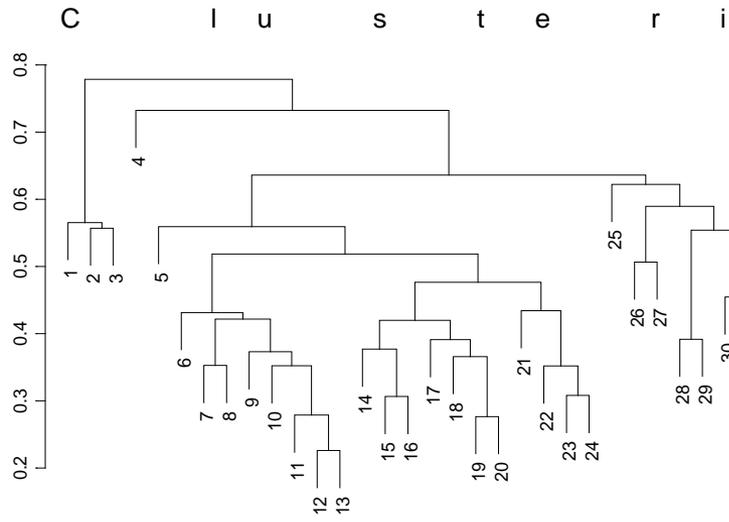


Euclidean distance small, 1-Correlation large



Linkage Methods

- Average Linkage
 - Merge clusters whose average distance between all pairs of items (one item from each cluster) is minimized
 - Particularly sensitive to distance metric
- Complete Linkage
 - Merge clusters to minimize the maximum distance within any resulting cluster
 - Tends to produce compact clusters
- Single Linkage
 - Merge clusters at minimum distance from one another
 - Prone to “chaining” and sensitive to noise

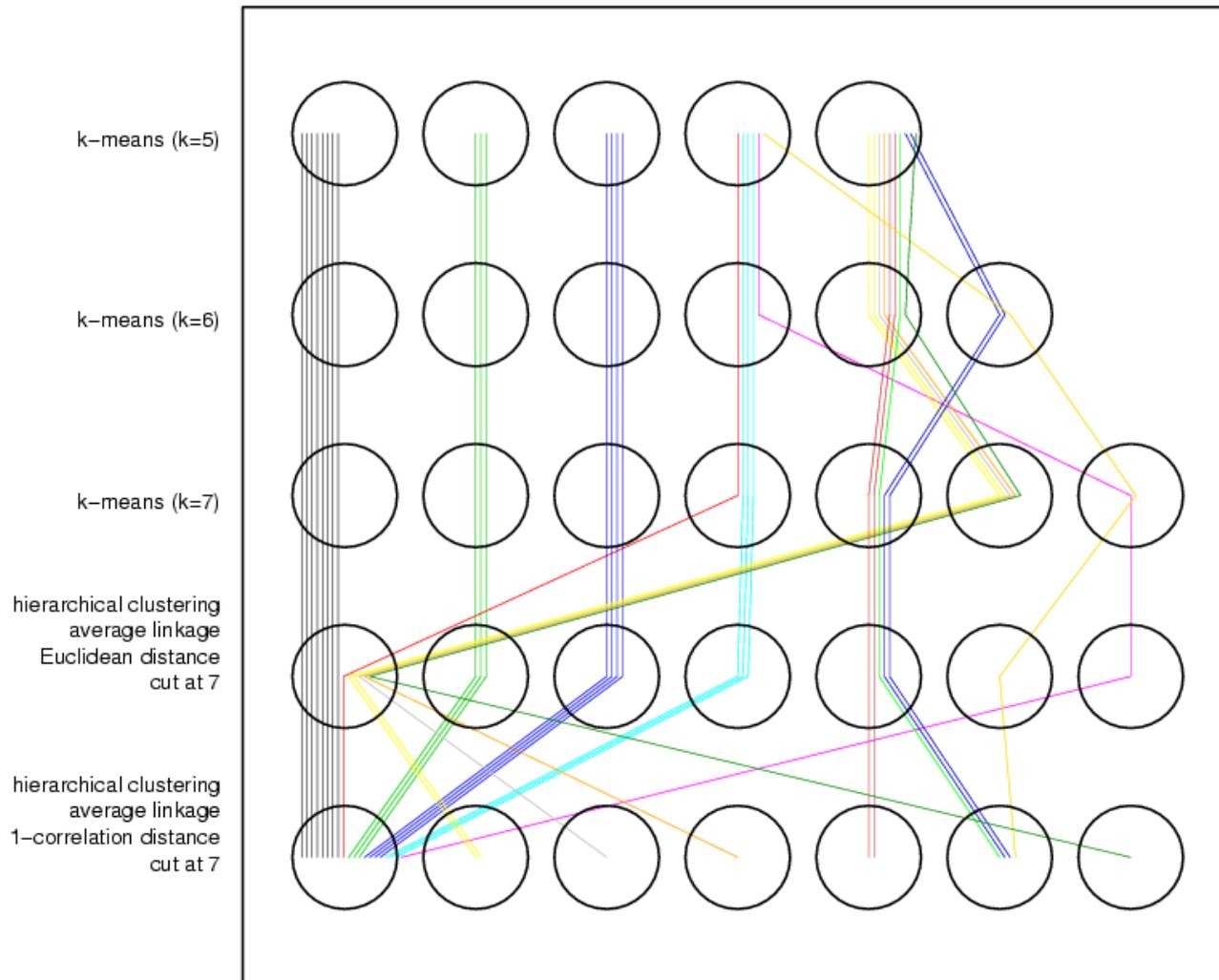


Dendrograms using 3 different linkage methods, distance = 1-correlation

(Data from Bittner *et al.*, *Nature*, 2000)

Does clustering method matter?

One set of specimens clustered by different methods



Interpretation of Cluster Analysis Results

- Cluster analyses always produce cluster structure
 - Where to “cut” the dendrogram?
 - Which clusters do we believe?
- Circular reasoning
 - Clustering using only genes found significantly different between two classes
 - “Validating” clusters by testing for differences between subgroups observed to segregate in cluster analysis
- Different clustering algorithms may find different structure using the same data

Assessing Clustering Results in BRB-ArrayTools

- Global test of clustering
 - Based on inter-sample distances in transformed dimension-reduced space
 - Available as an option in BRB-ArrayTools for the multidimensional scaling of samples.
- Assessment of reproducibility of individual clusters at selected cuts of the dendrogram in hierarchical clustering (could be generalized to other clustering methods)

(McShane *et al.*, *Bioinformatics*, 2002)

Assessing Cluster Reproducibility: Data Perturbation Methods

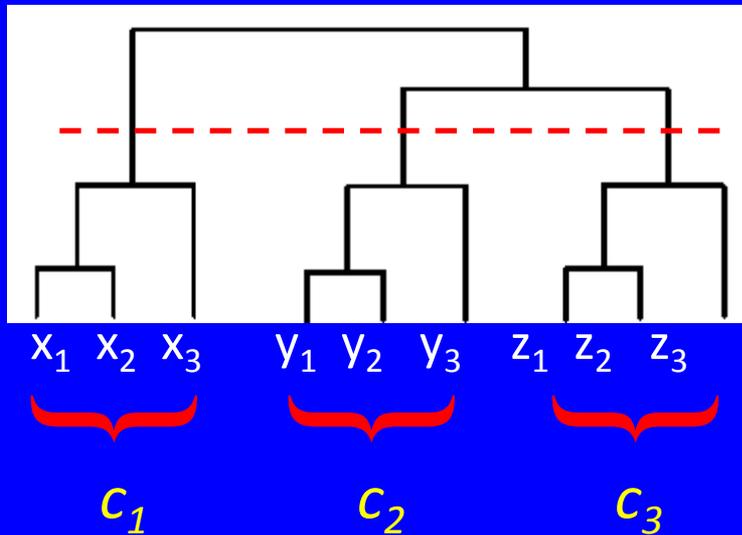
- Most believable clusters are those that persist given small perturbations of the data.
 - Perturbations represent an anticipated level of noise in gene expression measurements.
 - Perturbed data sets are generated by adding random errors to each original data point.
 - McShane *et al.*, *Bioinformatics*, 2002 – Gaussian errors
 - Kerr and Churchill, *PNAS*, 2001 – Bootstrap residual errors

Assessing Cluster Reproducibility: Data Perturbation Method in BRB-ArrayTools

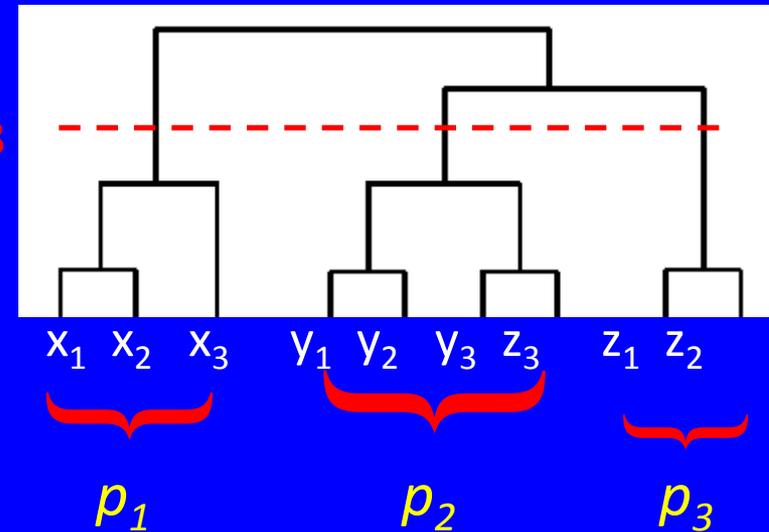
- Perturb the log-gene measurements by adding Gaussian noise and then re-cluster.
- For each original cluster:
 - Compute proportion of elements that occur together in the original cluster and remain together in perturbed data clustering when cutting dendrogram at the same level k .
 - Average the cluster-specific proportions over many perturbed data sets to get an *R-index* for each cluster.
 - The *R-index* may be obtained in BRB-ArrayTools for the hierarchical clustering of samples by selecting the ‘Compute cluster reproducibility measures’ option. The *R-index* option is not implemented for the hierarchical clustering of genes.
 - Hope for *R-index* ≥ 0.75

R-index Example

Original Data



Perturbed Data

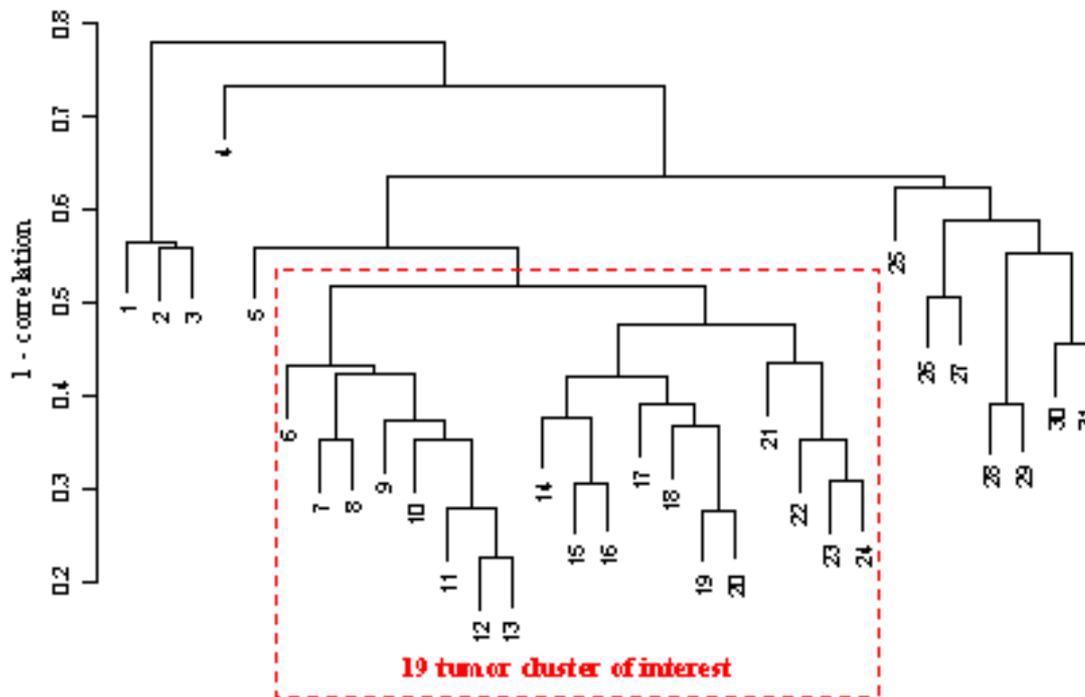


- 3 out of 3 pairs in c_1 remain together in perturbed clustering.
- 3 out of 3 in c_2 remain together.
- 1 out of 3 in c_3 remain together.
- $R\text{-index} = (3 + 3 + 1)/(3 + 3 + 3) = 0.78$

Cluster Reproducibility: Melanoma

(Bittner *et al.*, *Nature*, 2000)

- Expression profiles of 31 melanomas were examined with a variety of class discovery methods.
- A group of 19 melanomas consistently clustered together.



For hierarchical clustering, the cluster of interest had *R-index* = 1.0.

⇒ highly reproducible

Melanomas in the 19 element cluster tended to have:

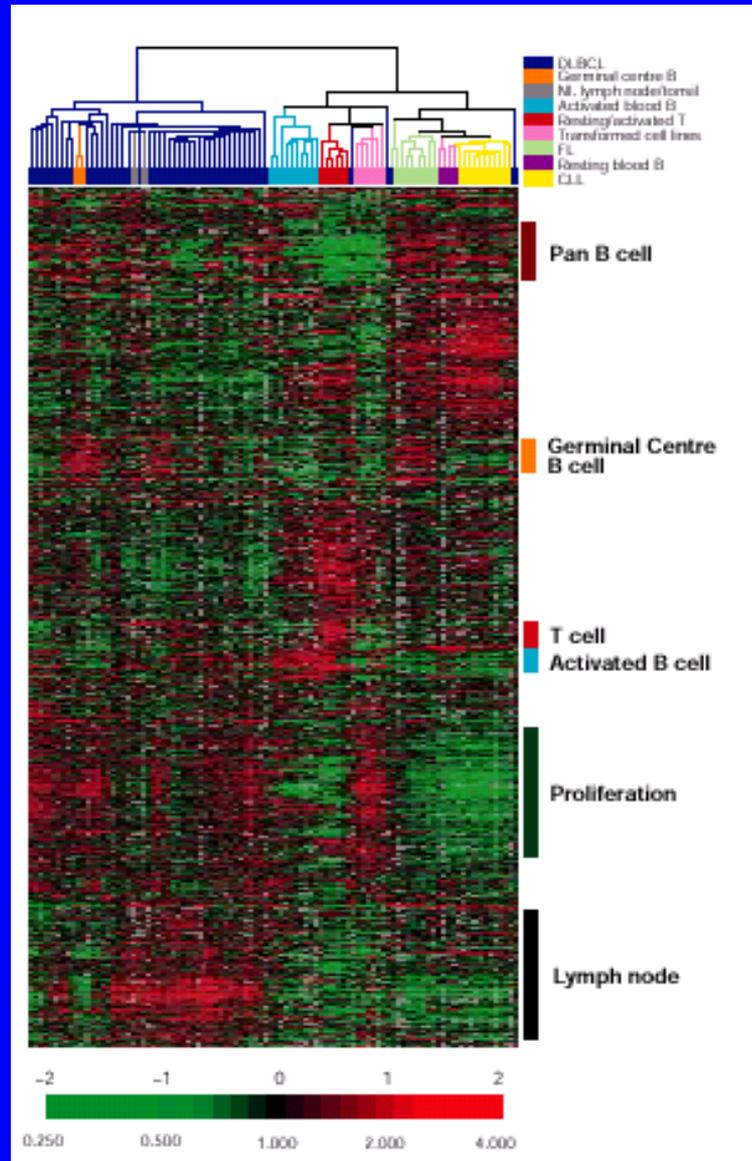
- reduced invasiveness
- reduced motility

Further References:

Estimating Number of Clusters

- GAP statistic (Tibshirani *et al.*, *JRSS B*, 2002) – detects too many false clusters (not recommended)
- Yueng *et al.* (*Bioinformatics*, 2001) – jackknife method, estimate # of *gene* clusters
- Dudoit *et al.* (*Genome Biology*, 2002) – prediction-based resampling
- Comparisons of methods for estimating number of clusters (Milligan and Cooper, *Psychometrika*, 1985) (uncertain performance in high dimensions)

Graphical Displays: Heat Map



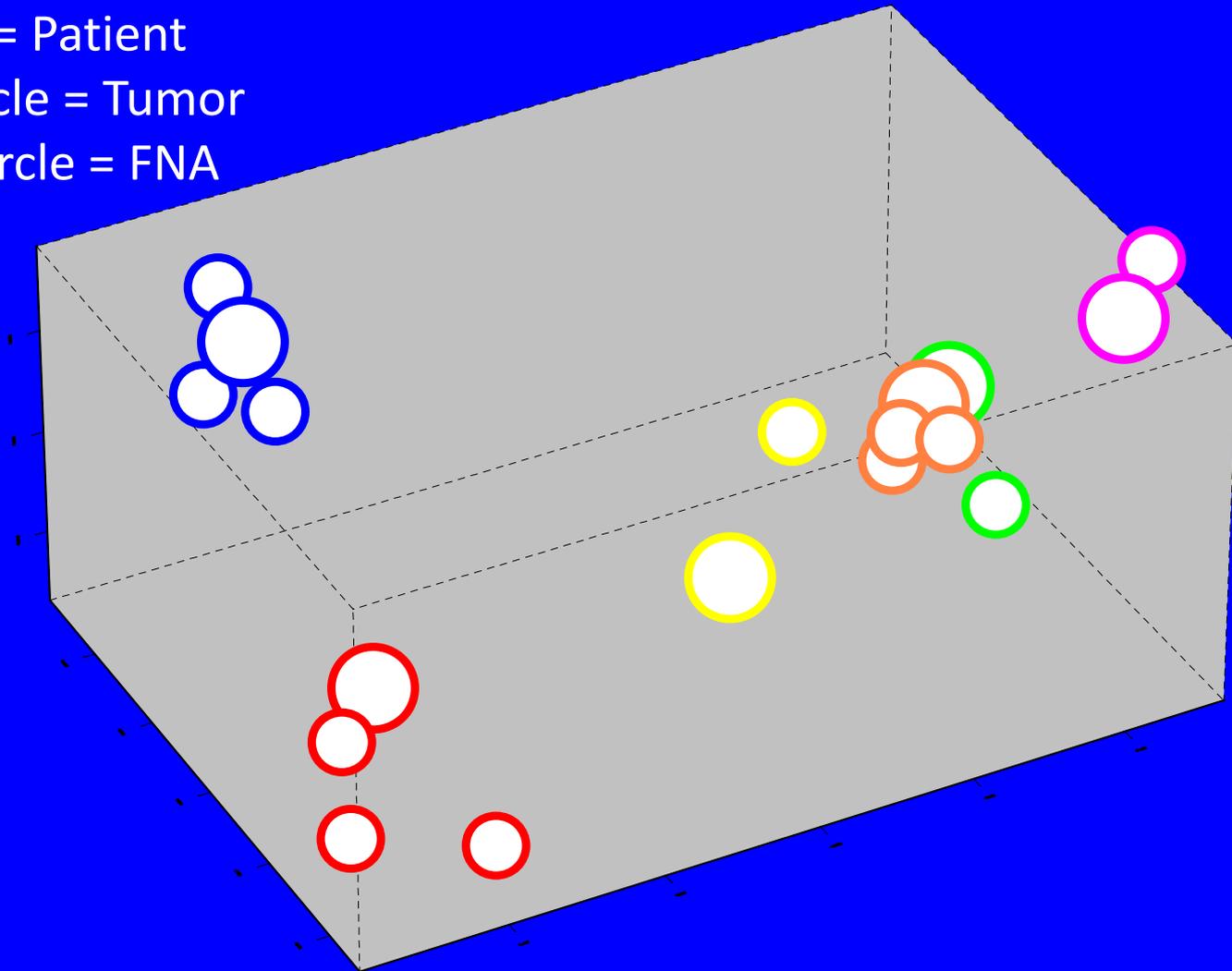
Hierarchical Clustering of Lymphoma Data (Alizadeh *et al.*, *Nature*, 2000) ⁶⁴

Graphical Displays: Multidimensional Scaling (MDS)

- High-dimensional (e.g. 5000-D) data points are represented in a lower-dimensional space (e.g. 3-D)
 - Principal components or optimization methods
 - Depends only on pairwise distances (Euclidean, 1-correlation, . . .) between points
 - “Relationships” need not be well-separated clusters

MDS: Breast Tumor and FNA Samples

Color = Patient
Large circle = Tumor
Small circle = FNA



(Assersohn *et al.*, *Clinical Cancer Research*, 2002)

Class Prediction Methods

Comparison of linear discriminant analysis, NN classifiers, classification trees, bagging, and boosting: tumor classification based on gene expression data (Dudoit, *et al.*, *JASA*, 2002)

Weighted voting method: distinguished between subtypes of human acute leukemia (Golub *et al.*, *Science*, 1999)

Compound covariate prediction: distinguished between mutation positive and negative breast cancers (Hedenfalk *et al.*, *NEJM*, 2001; Radmacher *et al.*, *J. Comp. Biology*, 2002)

Support vector machines: classified ovarian tissue as normal or cancerous (Furey *et al.*, *Bioinformatics*, 2000)

Neural Networks: distinguished among diagnostic subcategories of small, round, blue cell tumors in children (Khan *et al.*, *Nature Medicine*, 2001)

Compound Covariate Predictor (CCP)

(Tukey, *Controlled Clinical Trials*, 1993)

- Select “differentially expressed” genes by two-sample t -test with small α .

$$CCP_i = t_1 x_{i1} + t_2 x_{i2} + \dots + t_d x_{id}$$

t_j is the two-sample t -statistic for gene j .

x_{ij} is the log expression measure for gene j in sample i .

Sum is over all d differentially expressed genes.

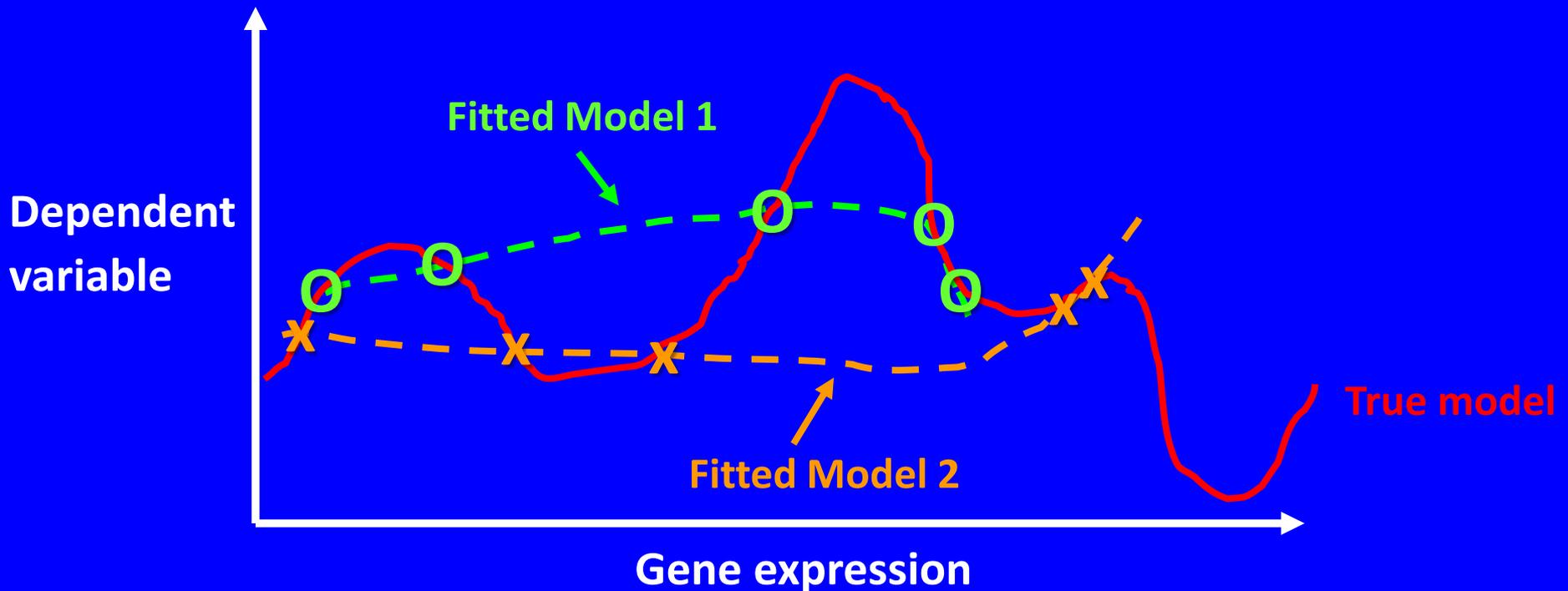
- Threshold of classification: midpoint of the CCP means for the two classes.

Classification: Avoiding Pitfalls

(Radmacher *et al.*, *J Comp Biology*, 2002;
Simon *et al.*, *JNCI*, 2003)

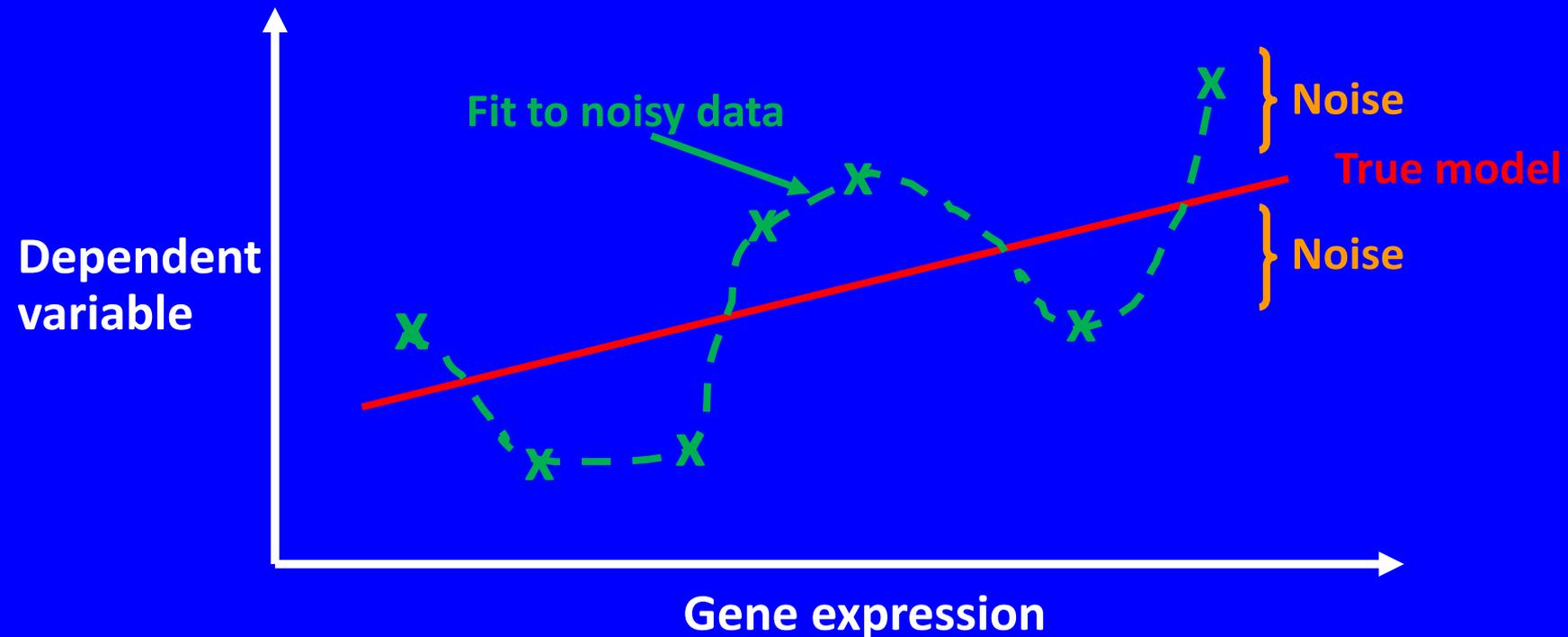
- When number of potential features is much larger than the number of cases, one can always fit a predictor to have 100% prediction accuracy on the data set used to build it
- If applied naively, more complex modeling methods (e.g., neural networks) are more prone to overfitting
- Estimating accuracy by “plugging in” data used to build a predictor results in highly biased estimates of prediction accuracy (re-substitution estimate)
- Internal and external validation of predictor are essential

Complexity of True Model



- Models in high dimension are usually complex
- Sample sizes are virtually always too small for precise estimation of the true model
- Look for simpler models that provide reasonable approximations

Overfitting

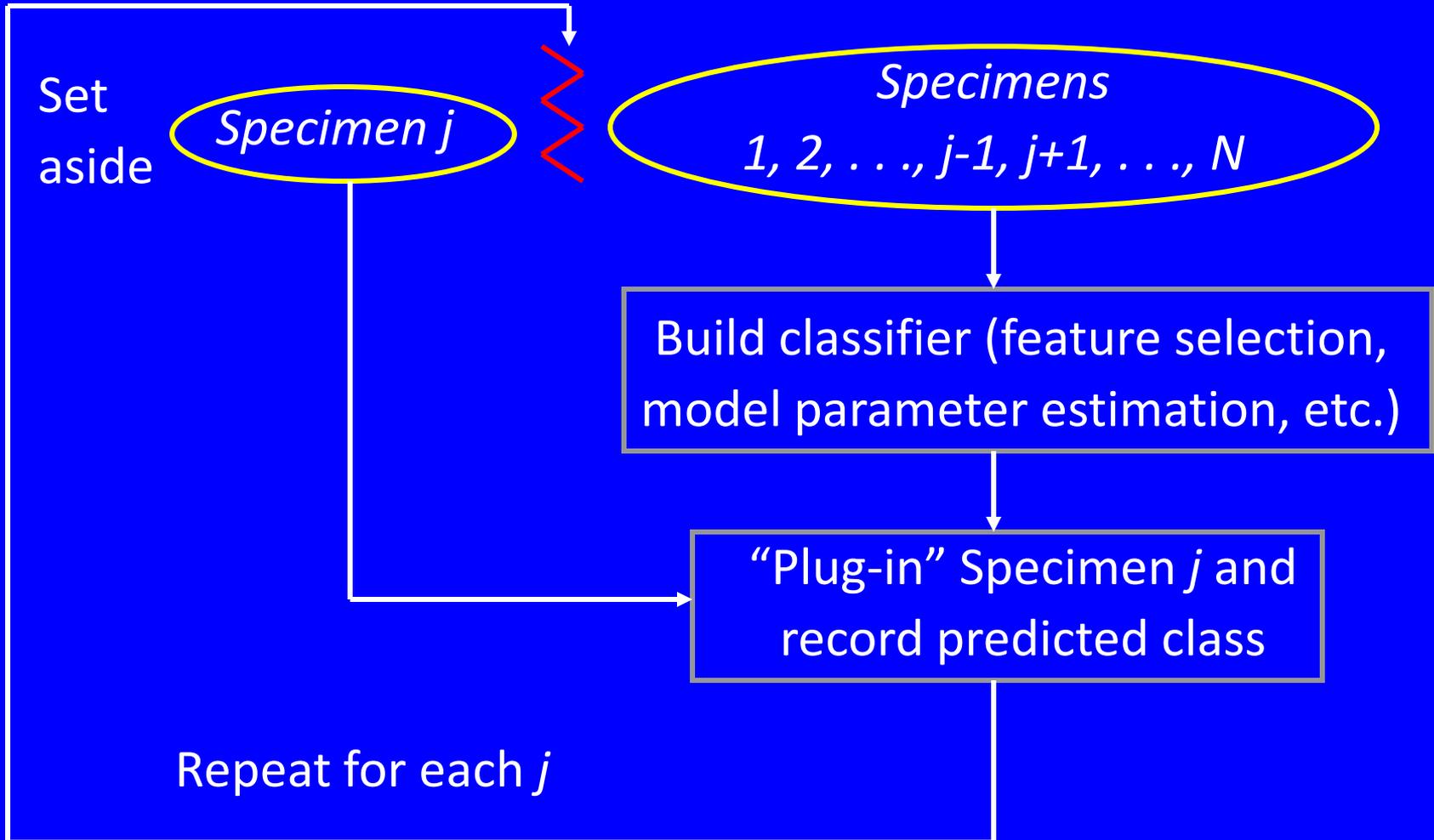


Attempting to fit complex models when sample size is small usually results in fitting to noise and producing models with no predictive value on independent data sets.

Validation Approaches

- Internal: within-sample validation
 - Cross-validation
(leave-one-out, split-sample, k-fold, etc.)
 - Bootstrap and other resampling methods
 - See Molinaro et al (*Bioinformatics* 2005) for comparison of methods
- External: independent-sample validation

Leave-One-Out Cross-Validation (LOOCV)



ALL steps, including feature selection, must be included in the cross-validation loop

Limitations of Within-Sample Validation

- Frequently performed incorrectly
 - Improper cross-validation (e.g., not including feature selection)
 - Special statistical inference procedures required (Lusa et al, *Statistics in Medicine* 2007; Jiang et al, *Stat Appl Genetics and Mol Biol* 2008)
- Large variance in estimated accuracy and effect sizes
- Doesn't protect against biases due to selective inclusion/exclusion of samples
- Built-in biases? (e.g., lab batch, specimen handling)

Prediction on Simulated Null Data

Generation of Gene Expression Profiles

- 20 specimens (P_i is the expression profile for specimen i)
- Log-ratio measurements on 6000 genes
- $P_i \sim \text{MVN}(\mathbf{0}, \mathbf{I}_{6000})$
- 10000 simulation repetitions
- Can we distinguish between the first 10 specimens (Class 1) and the last 10 (Class 2)? (class distinction is totally artificial since all 20 profiles were generated from the same distribution)

Prediction Method

- Compound covariate prediction
- Compound covariate built from the log-ratios of the 10 most differentially expressed genes.

Resubstitution Method

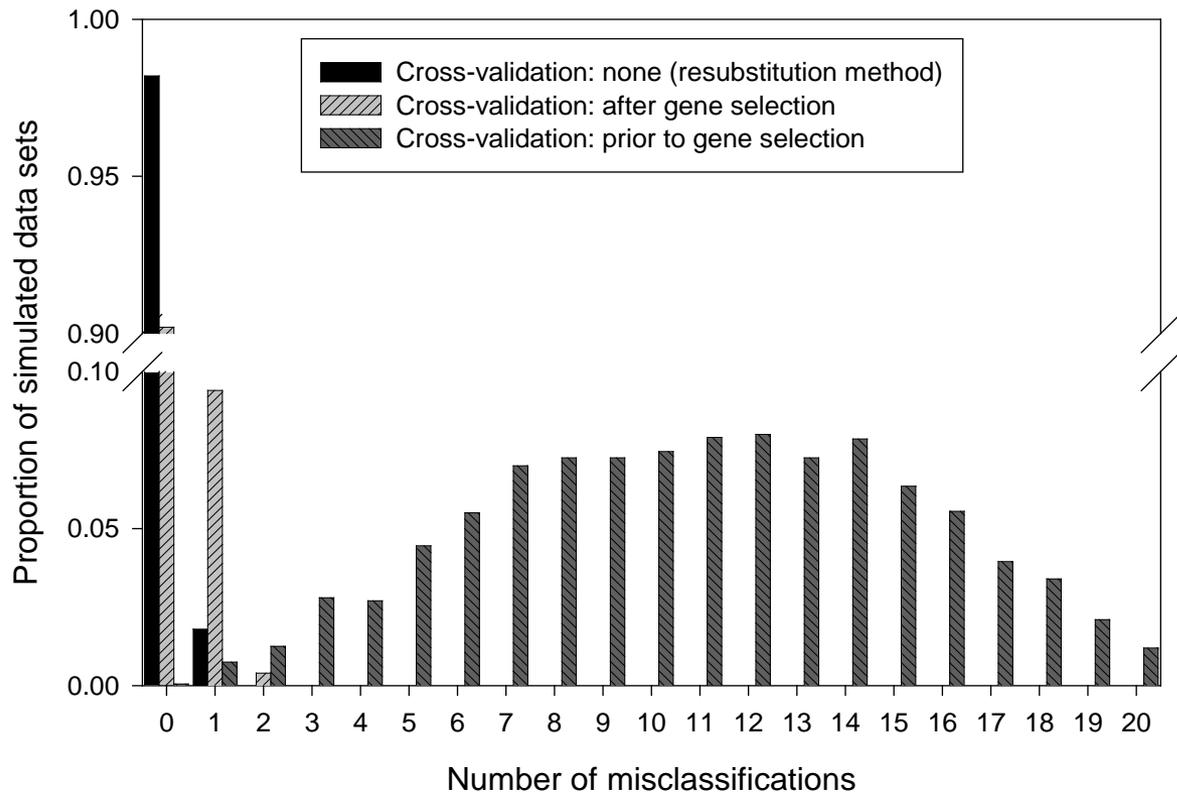
1. Build CCP from all data.
2. For $i=1\dots 20$, apply CCP to sample i .
3. Compare predicted class to actual class.

LOOCV Without Gene Selection

1. Select top 10 genes for CCP.
2. For $i=1, \dots, 20$:
 1. Leave out sample i .
 2. Build CCP(i) on other 19 samples.
 3. Apply CCP(i) to sample i .
3. Compare predicted class to actual class.

Full LOOCV

1. For $i=1\dots 20$
 1. Leave out sample i .
 2. Select top 10 genes and construct $CCP(i)$
 3. Apply $CCP(i)$ to sample i .
2. Compare predicted class to actual class.

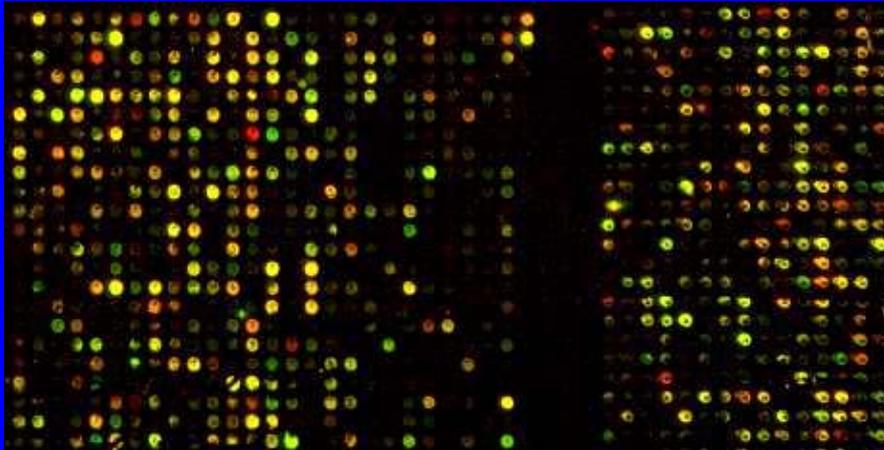


Gene-Expression Profiles in Hereditary Breast Cancer

(Hedenfalk *et al.*, *NEJM*, 2001)

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?

Classification of Hereditary Breast Cancers with Compound Covariate Predictor

Class labels	Number of differentially expressed genes (full data set, $\alpha = 0.0001$)	m = number of misclassifications using leave-one-out cross-validation	Proportion of random permutations with m or fewer misclassifications
BRCA1+ vs BRCA1-	9	1 (0 BRCA1+, 1 BRCA1-)	0.004
BRCA2+ vs BRCA2-	11	4 (3 BRCA2+, 1 BRCA2-)	0.043

Class Prediction in BRB-ArrayTools

- Variety of prediction methods available.
- Predictors are automatically cross-validated, and a significance test may be performed on the cross-validated mis-classification rate.
- Independent test samples may also be classified using the predictors formed on the training set.

Variance model: _____

Use randomized variance model for univariate tests.

Prediction methods: _____

- Compound covariate predictor
- K-nearest neighbors (for K=1 and 3)
- Nearest centroid
- Support vector machines
- Diagonal linear discriminant analysis

Predictors should only include genes: _____

- Significant univariately at level:
- With univariate misclassification rate below:
- With fold-ratio of geometric means between two classes exceeding:

Multivariate permutation test: _____

Do statistical significance test of cross-validated misclassification rate.

Outline

- 1) Introduction: Technology
- 2) Data Quality & Image Processing
- 3) Normalization & Filtering
- 4) Study Objectives
- 5) Analysis Strategies Based on Study Objectives
- 6) Design Considerations

Design Considerations

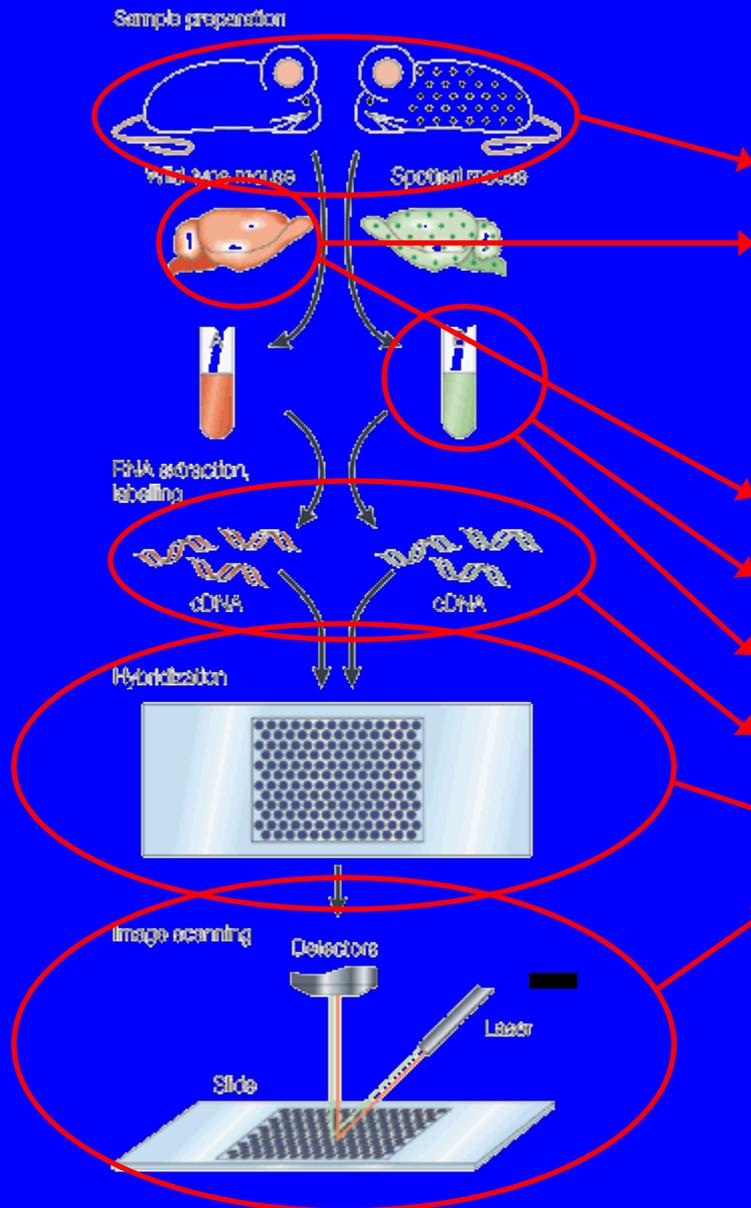
- Sample selection, including reference sample
- Sources of variability/levels of replication
- Pooling
- Sample size planning
- Controls
- For cDNA/2-color spotted arrays:
 - Reverse fluor experiments
 - Dobbin, Shih and Simon, *Bioinformatics*, 2003
 - Allocation of samples to (cDNA) array experiments
 - Kerr and Churchill, *Biostatistics*, 2001
 - Dobbin and Simon, *Bioinformatics*, 2002

Sample Selection

- Experimental Samples
 - A random sample from the population under investigation?
 - Broad versus narrow inclusion criteria?
- Reference Sample (cDNA array experiments using reference design)
 - In most cases, does not have to be biologically relevant.
 - Expression of most genes, but not too high.
 - Same for every array
 - Other situations exist (e.g., matched normal & cancer)

Sources of Variability

(cDNA Array Example)



- Biological Heterogeneity in Population
- Specimen Collection/ Handling Effects
 - Tumor: surgical bx, FNA
 - Cell Line: culture condition, confluence level
- Biological Heterogeneity in Specimen
- RNA extraction
- RNA amplification
- Fluor labeling
- Hybridization
- Scanning
 - PMT voltage
 - laser power

Levels of Replication

- Technical replicates
 - RNA sample divided into multiple aliquots and re-arrayed.
- Biological replicates
 - Use a different human/animal for each array.
 - In cell culture experiments, re-grow the cells under the same condition for each array (independent replication).

Summary: Replication Levels

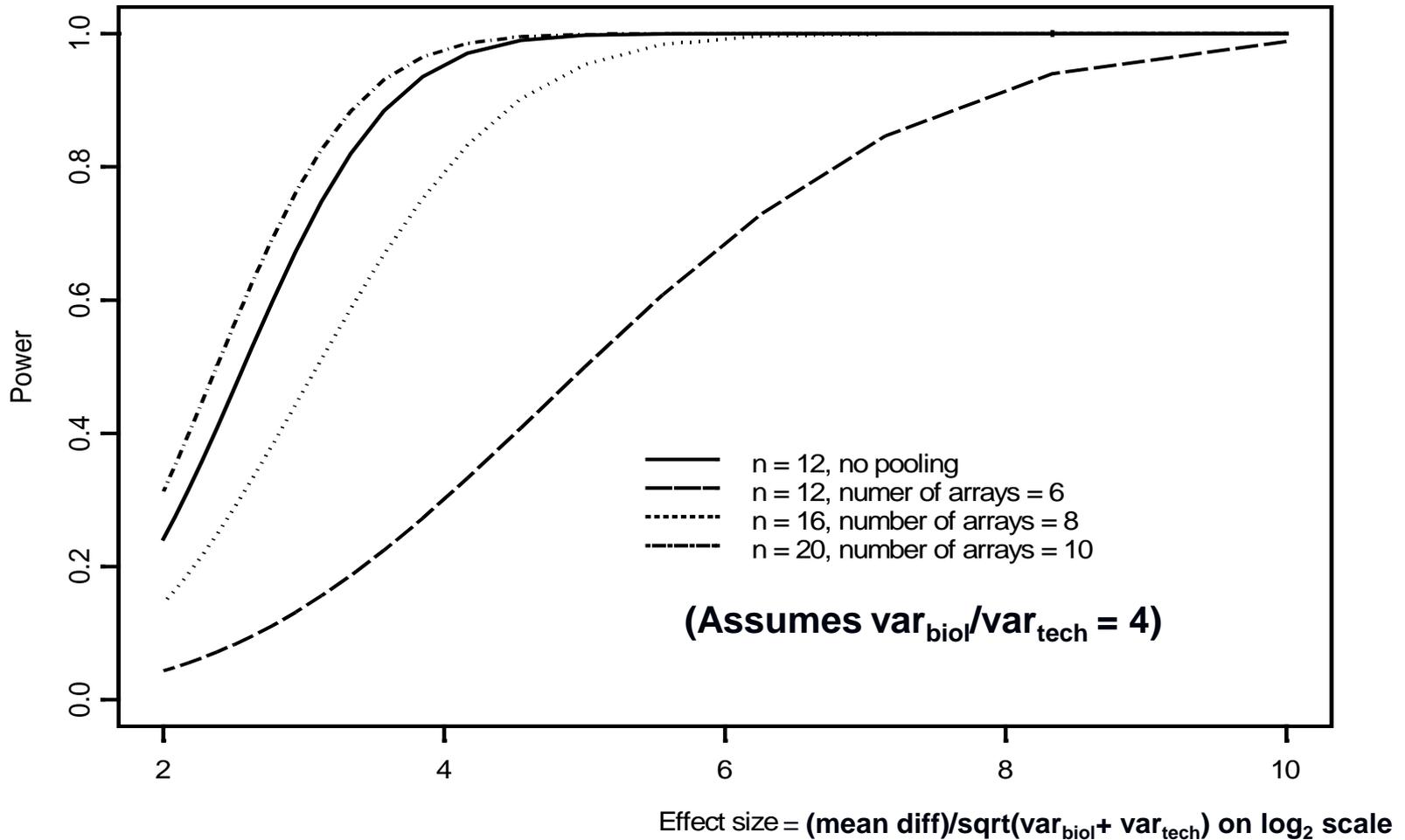
- Independent biological replicates are required for valid statistical inference.
- Maximizing biological replicates usually results in the best power for class comparisons.
- Technical replicates can be informative, e.g., for QC issues.
- But, systematic technical replication usually results in a less efficient experiment.

Is Pooling Advantageous?

- If RNA samples are tiny, pooling is an alternative to amplification.
- If RNA samples big enough, then there is not usually an advantage unless arrays are very expensive and samples very cheap.
- **NO FREE LUNCH:** Pooling samples for each array can reduce the number of arrays needed to achieve desired precision and power, but this will come at the **COST** of requiring that a larger number of biologically distinct samples be used.
- Single pool with many aliquots hybridized to arrays is **NOT** smart! Inference requires independent replication.

Power to detect specified expression difference between 2 groups

(McShane et al., JMGBN, 2003)



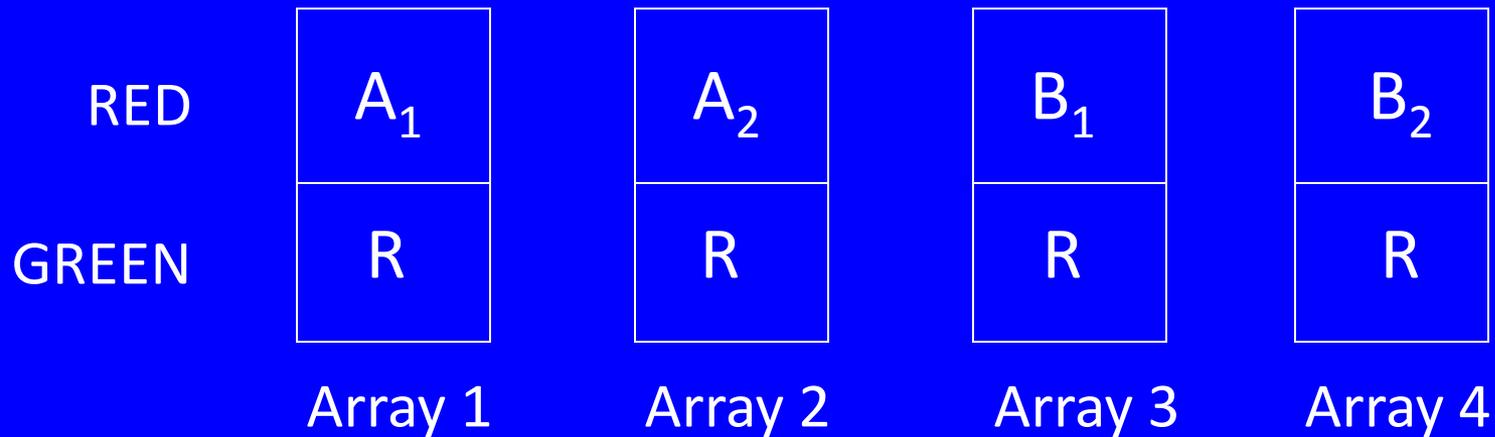
Kendziorski et al., *Biostatistics*, 2003

Shih et al., *Bioinformatics*, 2004

Class Comparison: Allocation of Specimens in cDNA Array Experiments

- Reference Design (traditional)
- Balanced Block Design
- Others
 - All pairs design
 - Loop Design (Kerr and Churchill, *Biostatistics*, 2001)
 - Variations on loop designs

Reference Design



A_i = *i*th specimen from class A

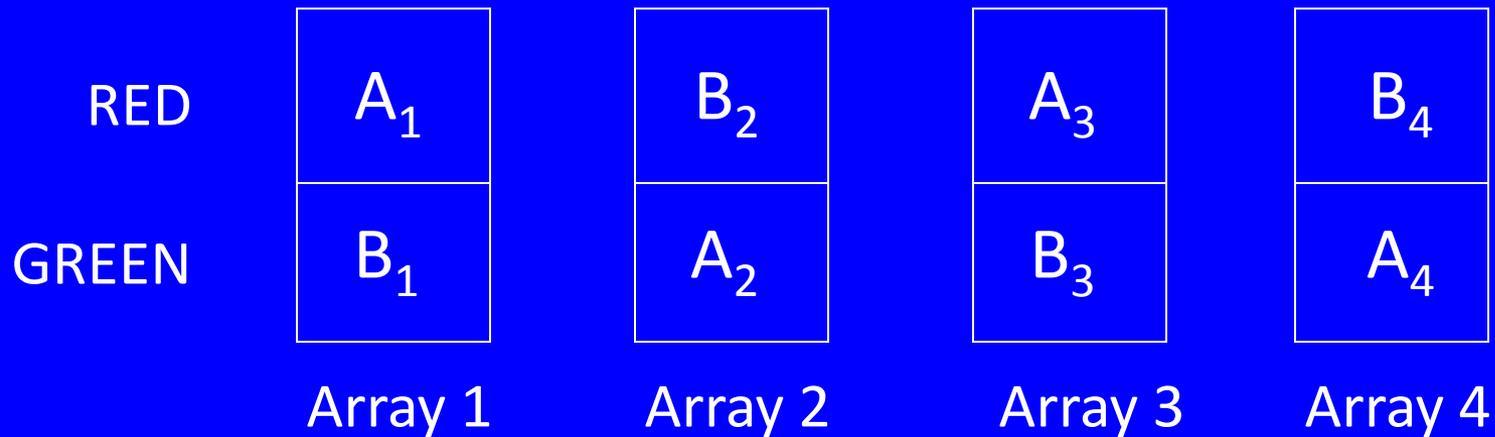
B_i = *i*th specimen from class B

R = aliquot from reference pool

Reference Design

- If the reference sample is not biologically relevant to the test samples, the class comparison is done between groups of arrays.
- If the comparison between the reference sample and test samples is biologically meaningful (e.g. reference sample is a mixture of normal samples, test samples are types of tumor samples), the class comparison is done between green and red channels – some reverse fluor experiments are required to adjust for potential dye bias.

Balanced Block Design



A_i = i th specimen from class A

B_i = i th specimen from class B

Summary Recommendations for Sample Allocation Schemes

- For 2-group comparison, block design is most efficient but precludes clustering.
- For cluster analysis or comparison of many groups, reference design is preferable.
- Reference design permits easiest analysis, allows greatest flexibility in making comparisons within and between experiments (using same reference), and is most robust to technical difficulties.
- The BRB-ArrayTools software performs class comparison between “groups of arrays” (e.g. reference designs) or between “red and green channels” (e.g. block designs).

Sample Size Planning

for 2-group comparisons with cDNA arrays using common reference design or with Affymetrix arrays

- No comprehensive method for planning sample size exists for gene expression profiling studies.
- In lieu of such a method...
 - Plan sample size based on comparisons of two classes involving a single gene.
 - Make adjustments for the number of genes that are examined.

Sample Size Planning

- Approx. total sample size required to compare two equal sized, independent groups:

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

where δ = mean difference between classes

σ = standard deviation

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

(δ and σ on log scale)

- More accurate iterative formulas recommended if n is approximately 60 or less

Sample Size Planning

Choosing α and β

Let K = # of genes on array

M = # of genes truly differentially expressed at
a fold difference of $\theta = 2^\delta$

Expected number of false positives:

$$\text{EFP} \leq (K-M) \times \alpha \quad (\alpha = \text{significance level})$$

Expected number of false negatives for θ -fold genes:

$$\text{EFN}_\theta = M \times \beta \quad (1-\beta = \text{power})$$

Popular choices for α and β :

$$\alpha = 0.001 \quad \beta = 0.05 \text{ or } 0.10$$

Sample Size Planning: Effect of α and β on FDR

- False Discovery Rate (FDR) is the expected proportion of false-positive genes on the gene list.

$$\text{FDR} = \frac{\alpha(1 - \pi)}{\alpha(1 - \pi) + (1 - \beta)\pi}$$

where

π = proportion of differentially expressed genes

π	α	$1 - \beta$	FDR
.005	.01	.95	68%
.005	.01	.80	71%
.005	.001	.95	17%
.005	.001	.80	20%
.05	.001	.95	2%

Sample Size Planning

Choosing σ and δ

- Value of σ will be determined by biology and experimental variation

Within a *single class*, what SD is expected for expression measure?

For \log_2 ratios, σ in range 0.25 – 1.0
(typically smallest for animal model
and cell line experiments)

- Value of δ is the size of mean difference (\log_2 scale) you want to be able to detect

$$\text{2-fold: } \delta = \log_2 (2) = 1$$

$$\text{3-fold: } \delta = \log_2 (3) = 1.59$$

Example Sample Size Calculation

$K = 10,000$ genes on array

$M = 100$ genes differentially expressed 2-fold

Specify $\alpha = 0.001, \beta = 0.05$
($z_{\alpha/2} = 3.291, z_{\beta} = 1.645$)
 $\sigma = 0.75$
 $\delta = 1$ (2-fold)

NEED $n = 55$ (approximately 28 per group)

Expect ≤ 10 false positives

Expect to miss approximately 5/100 2-fold genes

Sample Size Examples

($\alpha = .001$)

σ	δ	Fold-difference (2^δ)	n per group	Power(%)
.25	1	2	6	95
.5	1	2	14	95
.25	1	2	5	82
.5	1	2	5	14
.25	1.20	2.29	5	95
.5	2.39	5.24	5	95

Sample Size for Class Prediction

- Raises unique issues
 - The classes may mostly overlap, even in the high dimensional space.
 - There may be NO GOOD CLASSIFIER.
 - There will be an upper limit optimal performance that no classifier can exceed.
- Solution: Determine sample size big enough to get “close to optimal” performance
 - Dobbin and Simon, Biostatistics, 2007; Dobbin, Zhao and Simon, Clin Cancer Res, 2008.
 - Online interactive program website:
 - <http://linus.nci.nih.gov/brb/>

3 Essential Inputs for Sample Size for Class Prediction With Two Classes

- Number of genes on the array.
 - Example: ~22,000 features on an Affymetrix U113A array, ~54,000 on Affy + 2 arrays.
- The prevalence in each class.
 - Example: If 20% of patients respond to a drug– then the prevalence is 20% vs. 80%.
- The fold-change for informative genes:
 - Difference in class means divided by within class SD, on log base 2 scale

Sample Size Planning for Developing Classifiers Using High Dimensional Data

(Kevin Dobbin and Richard Simon, *Biostatistics* 8:101-17, 2007.)

Enter standardized fold change [> 0.2]

Enter number of genes on array [> 50]

Enter population prevalence in largest group (2 groups only) [between 0.5 and 0.85]

Note: This program provides estimates of the sample size required for a training set in order to ensure the resulting classifier has an expected accuracy within a tolerance of the optimal accuracy. Classifier performance should also be assessed. This can be done by cross-validation (resampling) or by applying the classifier to an independent validation set. The sample sizes given by this program are not appropriate for assessing classifier performance using cross-validation nor for determining the sample size required for a validation set. The program only considers sample sizes below 300. If >300 samples are needed, then an error message is returned indicating this.

Definition of standardized fold change: The standardized fold change is the difference between the class means divided by the within-class standard deviation, on the base 2 log scale. For example, if the raw fold change between the classes is 2, with $\log_2(2)=1$, and the within-class standard deviation for a typical gene is 0.71, then the standardized fold change is $1/0.71 \approx 1.4$. The 0.71 here is a typical median variance observed on

Sample Size Planning for Developing Classifiers Using High Dimensional Data Output Page - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Refresh Mail Print Send To Settings

Address <http://linus.nci.nih.gov/cgi-bin/simonr/R.cgi/ssc4pred.R> Go Links >>

Google Go Bookmarks 1665 blocked Check AutoLink AutoFill Send to Settings

Your Input:

Standardized fold change = 1.4

Number of genes on array = 22000

Population prevalence in largest group = 0.5

Result:

Training set sample size for Tolerance=0.05 is 63 , with 32 in class 1 and 31 in class 2

Training set sample size for Tolerance=0.10 is 50 , with 25 in class 1 and 25 in class 2

Output produced at Wed Sep 12 14:45:28 2007

Done Internet

Further Sample Size References

- Technical replicates for comparing 2 samples
 - Lee et al., PNAS, 2000
 - Black and Doerge, Bioinformatics, 2002
- Sample sizes for pooled RNA designs
 - Shih et al., Bioinformatics, 2004
- Sample sizes for balanced block designs, paired data, dye swaps, technical replicates, etc.
 - Dobbin et al., Bioinformatics, 2003
 - Dobbin and Simon, Biostatistics, 2005

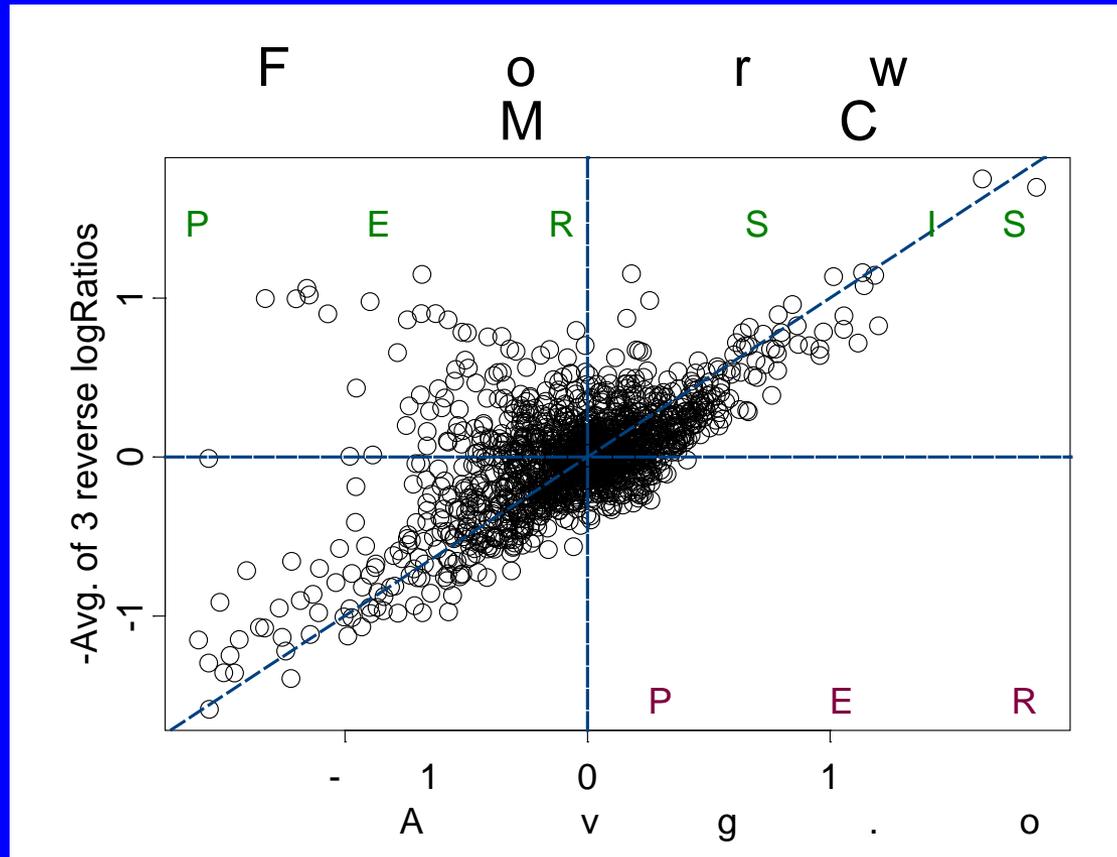
How Best to Allocate Effort?

- Microarrays can serve as a good high-throughput screening tool to identify potentially interesting genes.
- Verification of results via a different, more accurate, assay often preferable to running many arrays or technical replicates.
- Gene IDs associated with sequences can change over time, so periodic verification is advisable.

Controls

- Internal controls: Multiple clones (cDNA arrays) or probe sets (oligo arrays) for same gene spotted on array
- External controls: Spiked controls (e.g. yeast or *E. coli*)

cDNA/2-Color Spotted Arrays: Reverse Fluor Experiments



cDNA/2-color Spotted Arrays With Common Reference Design

Should reverse fluor “replicates” be performed for every array?

Usually NO!

See Dobbin, Shih and Simon, *Bioinformatics*, 2003 for a comprehensive discussion of reverse fluor replication

Reverse Fluors: cDNA/2-color Spotted Arrays With Common Reference Design

- When interested in interpreting individual ratios . . .
 - If gene-specific dye bias depends on gene sequence and not sample characteristics, dye bias can be adjusted for by performing *some* reverse fluor experiments.
 - If dye bias depends on both the gene and the sample, dye swaps won't help (Dobbin, Shih and Simon, 2005)!
- In BRB-ArrayTools reverse fluor arrays must be specified during the data importation (collation) step.

Reverse Fluors: cDNA/2-color Spotted Arrays With Common Reference Design

- When interested in class comparisons and using common reference design. . .
 - When comparing classes of non-reference samples tagged with the same dye, the dye bias should cancel out.
 - Reverse fluors are not required.

Reverse Fluors: cDNA/2-color Spotted Arrays With Balanced Block Design

- For each class, half the samples should be tagged with Cy3 and half with Cy5.
- When comparing different classes, dye bias will cancel out of the class comparisons.
- No reverse fluors are required.

Reverse Fluors: cDNA/2-color Spotted Arrays With Common Reference Design

- When interested in class discovery . . .
 - Usefulness of reverse fluor experiments and replicates will depend on nature and magnitude of both dye bias and experimental variability relative to between subject variability.
 - For some clustering methods (Euclidean distance), constant dye biases should “wash out”.
 - Some reverse fluors and replicates may be useful as informal quality checks.

Reverse Fluors: cDNA/2-color Spotted Arrays With Common Reference Design

- When interested in class prediction . . .
 - Dye bias may wash out for some predictors (e.g., nearest neighbor)
 - Dye bias may be incorporated in to some predictors (e.g., CCP)

Summary Remarks

- Data quality assessment and pre-processing are important.
- Different study objectives will require different statistical analysis approaches.
- Different analysis methods may produce different results. Thoughtful application of *multiple* analysis methods may be required.
- Chances for spurious findings are enormous, and validation of any findings on larger independent collections of specimens will be essential.
- Analysis tools can't compensate for poorly designed experiments.
- Fancy analysis tools don't necessarily outperform simple ones.
- Even the best analysis tools, if applied inappropriately, can produce incorrect or misleading results.

Helpful Websites

- NCI: <http://linus.nci.nih.gov/~brb>
 - Tech reports, talk slides, reference to book written by BRB members
 - BRB-ArrayTools software
 - .pdf of these talk slides: ftp://linus.nci.nih.gov/pub/techreport/CIT_course.pdf
- Berkeley: <http://www.stat.berkeley.edu/users/terry/Group/index.html>
- Harvard: <http://www.dchip.org>
- Hopkins: <http://biosun01.biostat.jhsph.edu/~ririzarr/Raffy/>
- Jackson Labs: <http://www.jax.org/staff/churchill/labsite/>
- Stanford:
 - <http://genome-www5.stanford.edu/MicroArray/SMD/restech.html>
 - <http://www-stat.stanford.edu/~tibs/> (R. Tibshirani)
- Bioconductor: <http://www.bioconductor.org/>
 - R-based, open source pre-processing and analysis tools
- Whitehead Institute (Cancer Genomics Group):
<http://www.broad.mit.edu/cancer/index.html>

Acknowledgements

- Richard Simon
- Joanna Shih
- Kevin Dobbin
- Michael Radmacher
- Other Members of the NCI Biometric Research Branch
- My NCI collaborators and students in the NIH microarray classes
- BRB ArrayTools Development Team