# Using Predictive Biomarkers in the Design of Pivotal Trials

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

#### BRB Website http://linus.nci.nih.gov/brb

- Powerpoint presentations and audio files
- Reprints & Technical Reports
- BRB-ArrayTools software
- BRB-ArrayTools Data Archive
- Sample Size Planning for Targeted Clinical Trials

## **Oncology Needs**

• Better treatments

Better targeting of treatments to the right patients

- Many cancer treatments benefit only a small proportion of the patients to which they are administered
- Targeting treatment to the right patients can greatly improve the therapeutic ratio of benefit to adverse effects
  - Smaller clinical trials needed
  - Treated patients benefit

### "Biomarkers"

- Surrogate endpoints
  - A measurement made before and after treatment to determine whether the treatment is working
  - Surrogate for clinical benefit
- Predictive classifiers
  - A measurement made before treatment to select good patient candidates for the treatment

# Surrogate Endpoints

- It is very difficult to properly validate a biomarker as a surrogate for clinical outcome. It requires a series of randomized trials with both the candidate biomarker and clinical outcome measured
  - Must demonstrate that treatment vs control differences for the candidate surrogate are concordant with the treatment vs control differences for clinical outcome
  - It is not sufficient to demonstrate that the biomarker responders survive longer than the biomarker nonresponders

### Cardiac Arrhythmia Supression Trial

- Ventricular premature beats was proposed as a surrogate for survival
- Antiarrythmic drugs supressed ventricular premature beats but killed patients at approximately 2.5 times that of placebo



- It is rare that we understand disease pathophysiology well enough to argue that a biomarker is self evidently a proper surrogate endpoint for clinical utility
- It is often more difficult and time consuming to properly "validate" an endpoint as a surrogate than to use the clinical endpoint in phase III trials
- The time frame for validating a surrogate is inconsistent with the time frame for initiating a pivotal study

- Biomarkers for use as endpoints in phase I or II studies need not be validated as surrogates for clinical outcome
- Unvalidated biomarkers can also be used for early "futility analyses" in phase III trials

## Validation=Fit for Purpose

- FDA terminology of "valid biomarker" and "probable valid biomarker" are inappropriate
- "Validation" has meaning only as fitness for purpose and the purpose of treatment selection classifiers are completely different than for surrogate endpoints
- Criteria for validation of surrogate endpoints should not be applied to biomarkers used for treatment selection

- The components of multi-gene expression based classifiers should not have to be "valid biomarkers"
- It is often much easier to develop an accurate predictive classifier than to elucidate the role of the component genes in disease biology

#### Oncology Needs Predictive Markers not Prognostic Factors

- Most prognostic factors are not used because they are not therapeutically relevant
- Most prognostic factor studies are poorly designed and not focused on a clear objective; they use a convenience sample of patients for whom tissue is available. Generally the patients are too heterogeneous to support therapeutically relevant conclusions
- Prognostic and predictive studies should be designed with as much care and statistical rigor as clinical trials

#### Pusztai et al. The Oncologist 8:252-8, 2003

- 939 articles on "prognostic markers" or "prognostic factors" in breast cancer in past 20 years
- ASCO guidelines only recommend routine testing for ER, PR and HER-2 in breast cancer
- "With the exception of ER or progesterone receptor expression and HER-2 gene amplification, there are no clinically useful molecular predictors of response to any form of anticancer therapy."

 Clinical trials of molecularly targeted drugs focused on patients whose tumors are expected to be susceptible to the drug can be much more efficient than traditional broad clinical trials

- In new drug development
  - The focus should be on evaluating the new drug in a population defined by a predictive classifier, not on "validating" the classifier
- In developing a predictive classifier for use in restricting a widely used treatment
  - The focus should be on evaluating the classifier; Is clinical outcome better if the classifier is used than if it is not used?

### New Drug Developmental Strategy (I)

- **Develop** a diagnostic classifier that identifies the patients likely to benefit from the new drug
- Develop a reproducible assay for the classifier
- Use the diagnostic to restrict eligibility to a prospectively planned evaluation of the new drug
- Demonstrate that the new drug is effective in the prospectively defined set of patients determined by the diagnostic

Develop Predictor of Response to New Drug



# Applicability of Design I

- Primarily for settings where the classifier is based on a single gene whose protein product is the target of the drug
- With substantial biological basis for the classifier, it will often be unacceptable ethically to expose classifier negative patients to the new drug

- Traditional parameters of sensitivity and specificity are not applicable to estimating relative efficacy of a new regimen versus a control with survival or progressionfree survival endpoint
  - The relevant parameters are treatment effect in classifier positive and classifier negative subsets
- "When your only tool is a hammer, everything looks like a nail"
- Forcing predictive medicine based drug development into square boxes developed for traditional medical devices creates a serious roadblock to the introduction of effective pharmacogenomic based therapeutics

#### Evaluating the Efficiency of Strategy (I)

- Simon R and Maitnourim A. Evaluating the efficiency of targeted designs for randomized clinical trials. Clinical Cancer Research 10:6759-63, 2004.
- Maitnourim A and Simon R. On the efficiency of targeted clinical trials. Statistics in Medicine 24:329-339, 2005.
- reprints and interactive sample size calculations at http://linus.nci.nih.gov/brb

Pharmacogenomic Model for Two Treatments With Binary Response

- •Molecularly targeted treatment E
- Control treatment C
- • $\gamma$  Proportion of patients that express target
- $\bullet p_c$  control response probability
- •response probability for E patients who express target is  $(p_c + \delta_1)$
- •Response probability for E patients who do not express target is ( $p_c + \delta_0$ )

### Approximations

Observed response rate ~ N(p,p(1-p)/n)

•  $p_e(1-p_e) \sim p_c(1-p_c)$ 

# **Two Clinical Trial Designs**

- Un-targeted design
  - Randomized comparison of E to C without screening for expression of molecular target
- Targeted design
  - Assay patients for expression of target
  - Randomize only patients expressing target

### Number of Randomized Patients Required

- Type I error  $\alpha$
- Power 1- $\beta$  for obtaining significance

$$n = 2(p_{c}q_{c} + p_{e}q_{e}) \left(\frac{k_{1-\alpha} + k_{1-\beta}}{p_{e} - p_{c}}\right)^{2}$$

For targeted design

$$-p_e = p_c + \delta_1$$
  
 $-p_e - p_c = \delta_1$ 

• For un-targeted design

$$- p_e = (1-\gamma)(p_c + \delta_0) + \gamma(p_c + \delta_1)$$
$$- p_e - p_c = \gamma \delta_1 + (1-\gamma) \delta_1$$

#### Randomized Ratio (normal approximation)

- RandRat =  $n_{untargeted}/n_{targeted}$ RandRat  $\approx \left(\frac{\delta_1}{\gamma \delta_1 + (1 - \gamma) \delta_0}\right)^2$
- $\delta_1 = rx$  effect in marker + patients
- $\delta_0 = rx$  effect in marker patients
- $\gamma$  =proportion of marker + patients
- If  $\delta_0$ =0, RandRat = 1/  $\gamma^2$
- If  $\delta_0 = \delta_1/2$ , RandRat = 4/( $\gamma$  +1)<sup>2</sup>

# Imperfect Assay Sensitivity & Specificity

•  $\lambda_{sens}$ =sensitivity

– Pr[assay+ | target expressed]

•  $\lambda_{spec}$ =specificity

- Pr[assay- | target not expressed]

### Proportion of Assay Positive Patients That Express Target

$$w_1 = \frac{\gamma \lambda_{sens}}{\gamma \lambda_{sens} + (1 - \gamma)(1 - \lambda_{spec})}$$

$\lambda_{ extsf{sens}}$	$\lambda_{ ext{spec}}$	γ	W <sub>1</sub>
0.9	0.9	0.75	0.96
0.9	0.9	0.5	0.9
0.9	0.9	0.25	0.75
0.9	0.9	0.10	0.50

### **Randomized Ratio**

• RandRat =  $n_{untargeted}/n_{targeted}$ 

$$RandRat = \left(\frac{w_1\delta_1 + (1 - w_1)\delta_0}{\gamma\delta_1 + (1 - \gamma)\delta_0}\right)^2$$

# Randomized Ratio sensitivity=specificity=0.9

γ Express target	δ <b>₀=0</b>	δ <sub>0</sub> = δ <sub>1</sub> /2
0.75	1.29	1.26
0.5	1.8	1.6
0.25	3.0	1.96
0.1	25.0	1.86

## Screened Ratio Imperfect Assay

$$N_{\text{targeted}} = \frac{n_{\text{targeted}}}{\gamma \lambda_{sens} + (1 - \gamma)(1 - \lambda_{spec})}$$
  
ScreenRat =  $[\gamma \lambda_{sens} + (1 - \gamma)(1 - \lambda_{spec})]Randra$ 

# Randomized Ratio sensitivity=specificity=0.9

γ Express target	δ <b>₀=0</b>	δ <sub>0</sub> = δ <sub>1</sub> /2
0.75	1.29	1.26
0.5	1.8	1.6
0.25	3.0	1.96

### Screened Ratio sensitivity=specificity=0.9

γ Express target	δ <b>₀=0</b>	δ <sub>0</sub> = δ <sub>1</sub> /2
0.75	0.9	0.88
0.5	0.9	0.80
0.25	0.9	0.59
0.1	4.5	0.33

#### Web Based Software for Comparing Sample Size Requirements

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

🚰 Biometric Research Branch home page - Microsoft Internet Explorer			×
File Edit View Favorites Tools Help			
🔇 Back + 🜍 + 💌 😰 🏠 🔎 Search 👷 Favorites 🤣 🎯 + 🌺 🔞 + 🛄 3%			
Address 🗃 http://linus.nci.nih.gov/~brb/	💌 🄁 Go	Links	>
Google + west hawaii cancer symposium 🕑 🖸 Search + 👰 🛷 🕥 + 🙀 - 😡 Options 🌽 🖏 west 👸 hawaii 👸 cancer 👸 symposium		1	
biomathematics, and computational biology, on topics ranging from methodology to facilitate understanding at the molecular level of the pathogenesis of cancer to methodology to enhar	nce the conduct of clinica	al	~



#### **Research** Areas

trials of new therapeutic and diagnostic approaches.

Clinical trials, Drug Discovery, Molecular Cancer Diagnosis, Biomedical Imaging, Computational and Systems Biology, and Biostatistical Research



#### **Technical Reports and Talks**

Download the PDF version of our most recent research papers and PowerPoint version of the talk slides.



**RRR Staff** Investigators and contact information



#### **BRB** ArrayTools Download the most advanced tools for microarray data

analysis



#### BRR Alumni



#### Sample Size Calculation



#### **BRB Annual Report 2005**



君

**Position Available** Post-doctoral fellow positions available



#### **Mathematics And Oncology**

- The Norton-Simon Hypothesis
- The Norton-Simon Hypothesis and Breast Cancer Mortality in National Randomized Trial



#### Software Download

- Accelerated Titration Design Software
- Optimal Two-Stage Phase II Design Software

19 Adobe Photoshop Ele. 💽 Inbox - Microsoft Out...

Biometric Research B...
File Edit View Favorites Tools Help	<u>ar</u>
🚱 Back 🔹 😰 🔹 🛃 🖍 🔎 Search 🧙 Favorites 🤣 🎯 + 🥁 🔞 🔹 🥨	
Address 🕘 http://linus.nci.nih.gov/~simonr/samplesize.html	Go Links 🎽
Google 🗸 west hawaii cancer symposium 🛛 🔽 🖸 Search 🔹 🕸 🕥 🗧 🖓 🔞 🔹 🚱 Options 🤌 👸 west 👸 hawaii 👸 cancer 👸 symposium	R •

Sample Size Calculation for Randomized Clinical Trials

• Optimal Two-Stage Phase II Design

- Biomarker Targeted Randomized Design\*
- 1. Binary Outcome Endpoint

S Adobe Photoshop El..

2. Survival and Time-to-Event Endpoint

\* Targeted design randomizes only marker positive patients to treatment or control arm. Untargeted design does not measure marker and randomizes all who otherwise are eligible.

© NIH, 2006

🔞 Inbox - Microsoft Ou ...

🛃 start

😂 🖸 🗿

😨 Connected - BlackBe..

Internet

Sample Size Calculation: Binary Outcome Endp	oint - Microso	oft Internet Exp	lorer			
File Edit View Favorites Tools Help						2
🜀 Back 🔹 🜍 🕤 🖹 📓 🚮 🔎 Search	📌 Favorites	🛛 🖉 - 🖉	🧯 🗹 🗧 🦓			
Address Address http://linus.nci.nih.gov/~simonr/boep.html						Go Links 2
Google - west hawaii cancer symposium	× C	Search 🔸 👯	🚿 💽 + 👔 🔹 🛃 Ostion	s 🤌 🔕 west 👸 hawaii 🛛	🖞 cancer 🛛 👸 symposium	
Evaluating the efficiency of targeted o	Samp	p <i>le Size</i> randomized c pc	Calculation: Bil linical trials and <u>Supple</u> 10:6759-6763,	mary Outcome ment by Richard Simo 2005)	Endpoint n and Aboubakar Maite	urnam. (Clinical Cancer Research
		gamma delta1				
		delta0 alpha	0.05			
		power	0.90 Submit			
	pc	= probabili	ty of "response" for con	trol arm		
	gamma	= proportio responsive ·	n of patients who are clo to new treatment	assifier negative (i.e. l	ess	
	delta1	= improvem positive pat	ent in response probabil tients	ity for new treatment	in classifier	
	delta0	= improvem negative pa	ent in response probabil tients	ity for new treatment	in classifier	
	alpha	= two-sided	significance level			
			© NIH, 200	06		
🗃 Done						🥩 Internet
W start Start Start	Ber 🚺 🙆 A	dobe Photoshop Ele	Microsoft Out	Romelo Siao Calculati	Decument - Microsof	

## Developmental Strategy (II)

Develop Predictor of Response to New Rx



#### Developmental Strategy (II)

- Do not use the diagnostic to restrict eligibility, but to structure a prospective analysis plan.
- Compare the new drug to the control overall for all patients ignoring the classifier.
  - If  $p_{overall} \leq 0.04\,$  claim effectiveness for the eligible population as a whole
- Otherwise perform a single subset analysis evaluating the new drug in the classifier + patients
  - If  $p_{\text{subset}} \leq 0.01$  claim effectiveness for the classifier + patients.

- This analysis strategy is designed to not penalize sponsors for having developed a classifier
- It provides sponsors with an incentive to develop genomic classifiers

## Key Features of Design (II)

 The purpose of the RCT is to evaluate treatment T vs C overall and for the predefined subset; not to re-evaluate the components of the classifier, or to modify or refine the classifier

## Sample Size Planning for Design II

- 1. Size for standard power (e.g. 0.9) for detecting usual treatment effect at significance level 0.04
- 2. Size for standard power (e.g. 0.9) for detecting larger treatment effect in positive subset
- 3. Size as in 1 but extend accrual of classifier positive patients if overall test is non-significant

#### Developmental Strategy (IIb)

- Do not use the diagnostic to restrict eligibility, but to structure a prospective analysis plan.
- Compare the new drug to the control for classifier positive patients
  - If  $p_+>0.05$  make no claim of effectiveness
  - If  $p_{\star}{\leq}~0.05\,$  claim effectiveness for the classifier positive patients and
    - Continue accrual of classifier negative patients and eventually test treatment effect at 0.05 level

#### Sample size Planning for IIb

- Accrue classifier positive and negative patients until there are sufficient classifier positive patients for standard power at significance level 0.05 for detecting large treatment effect D
- If treatment is found effective in classifier + patients, continue accrual of negative patients for standard power at significance level 0.05 for detecting usual size treatment effect d representing minimal useful clinical utility

Hazard ratio δ to be detected	Number of events required α=0.05	Number of events required α=0.04
1.2	632	669
1.3	305	323
1.4	186	196
1.5	128	135
1.6	95	101
1.7	75	79
1.8	61	64
1.9	51	54
2.0	44	46

Number of events required for detecting a proportional hazard treatment effect with 90% power

Hazard ratio $\delta$ to be detected	Number of events required
1.7	105
1.8	86
1.9	72
2.0	62
2.1	54
2.2	48
2.3	43

Number of events required for detecting a proportional hazard treatment effect with 1% two-sided significance level and 90% power

Hazard rate to be detected overall	Hazard rate to be detected in + subset	Proportion classifier +	Number events needed for overall analysis at .04 level	Number events needed for classifier + analysis at .01 level	Number of total events to accrue
1.3	2	0.33	323	62	323
1.5	2	0.33	135	62	186

The alternative design of separate testing of treatment effect in positive and negative subsets is problematic

- With classifier tightly linked to drug target, it may be ethically unacceptable to expose classifier negative patients
- With an empirically based classifier (C), if the treatment effect is not enhanced for C + patients, then 128 events are needed in both C+ and C- patients to detect a hazard ratio of 1.5 with significance level .05 and power .9 for each analysis
  - The chance of a false negative in at least one subset is 19%
  - the potential value of being able to do a subset analysis may not be worth the cost of having to demonstrate effectiveness in both subsets separately for broad labeling

# FDA Subset Catch 22

- Do not accept claims based on subset analysis
- Require sponsors to do subset analysis to establish that a claim based on overall treatment effect applies to all subsets

Event rate	Number of patients needed for trial	Separate analysis of classifier negative and classifier positive patients		
	ignoring classifier	Number of classifier negative patients needed	Number of classifier positive patients needed	
0.1	1280	1280	440	
0.2	640	640	220	
0.3	427	427	147	
0.4	320	320	110	
0.5	256	256	88	
0.6	214	214	74	
0.7	183	183	63	
0.8	160	160	55	

Two-sided significance level of 5% and power 90% for each comparison. Untargeted trial based on detecting hazard ratio of 1.5. Targeted trial based on detecting hazard ratio of 1.5 for classifier negative patients and 2.0 for classifier positive patients.

### Predictive Medicine not Correlative Science

- The purpose of the RCT is to evaluate treatment T vs C overall and for the pre-defined subset
- The purpose is not to re-evaluate the components of the classifier, or to modify or refine the classifier
- The purpose is not to demonstrate that repeating the classifier development process on independent data results in the same classifier

#### The Roadmap

- Develop a completely specified genomic classifier of the patients likely to benefit from a new drug
- 2. Establish reproducibility of measurement of the classifier
- 3. Use the completely specified classifier to design and analyze a new clinical trial to evaluate effectiveness of the new treatment with a pre-defined analysis plan.

# **Guiding Principle**

- The data used to develop the classifier must be distinct from the data used to test hypotheses about treatment effect in subsets determined by the classifier
  - Developmental studies are exploratory
  - Studies on which treatment effectiveness claims are to be based should be definitive studies that test a treatment hypothesis in a patient population completely pre-specified by the classifier

# Use of Archived Samples

- From a non-targeted "negative" clinical trial to develop a binary classifier of a subset thought to benefit from treatment
- Test that subset hypothesis in a separate clinical trial
  - Prospective targeted type (I) trial
  - Prospective type (II) trial
  - Using archived specimens from a second previously conducted clinical trial

### Development of Genomic Classifiers

- Single gene or protein based on knowledge of therapeutic target
- Empirically determined based on correlating gene expression to patient outcome after treatment

### Development of Genomic Classifiers

- During phase II development or
- After failed phase III trial using archived specimens.
- Adaptively during early portion of phase III trial.

## Development of Empirical Gene Expression Based Classifier

- 20-30 phase II responders are needed to compare to non-responders in order to develop signature for predicting response
  - Dobbin KK, Simon RM. Sample size planning for developing classifiers using high dimensional DNA microarray data, Biostatistics 8:101-117, 2007.

Adaptive Signature Design An adaptive design for generating and prospectively testing a gene expression signature for sensitive patients

#### Boris Freidlin and Richard Simon Clinical Cancer Research 11:7872-8, 2005

Adaptive Signature Design End of Trial Analysis

- Compare E to C for **all patients** at significance level 0.04
  - If overall  $H_0$  is rejected, then claim effectiveness of E for eligible patients
  - Otherwise

- Otherwise:
  - Using only the first half of patients accrued during the trial, develop a binary classifier that predicts the subset of patients most likely to benefit from the new treatment E compared to control C
  - Compare E to C for patients accrued in second stage who are predicted responsive to E based on classifier
    - Perform test at significance level 0.01
    - If H<sub>0</sub> is rejected, claim effectiveness of E for subset defined by classifier

Treatment effect restricted to subset. 10% of patients sensitive. Sensitivity genes are uncorrelated. 400 patients, 10,000 genes



Number of sensitivity genes

Treatment effect restricted to subset. 10% of patients sensitive.

Sensitivity genes are correlated, 400 patients, 10,000 genes.



#### Treatment effect restricted to subset. 10% of patients sensitive, 10 sensitivity genes, 10,000 genes, 400 patients.

Test	Power
Overall .05 level test	46.7
Overall .04 level test	43.1
Sensitive subset .01 level test (performed only when overall .04 level test is negative)	42.2
Overall adaptive signature design	85.3

#### Overall treatment effect, no subset effect. 10,000 genes, 400 patients.

Test	Power
Overall .05 level test	74.2
Overall .04 level test	70.9
Sensitive subset .01 level test	1.0
Overall adaptive signature design	70.9

Biomarker Adaptive Threshold Design W Jiang, B Freidlin, R Simon (submitted)

- Randomized pivotal trial comparing new treatment E to control C
- Quantitative biomarker B
- Survival or DFS endpoint

#### Biomarker Adaptive Threshold Design

- Compare E vs C overall using significance threshold of 0.04
  - If significant, claim broad effectiveness of E
  - If not significant, proceed as below

#### Biomarker Adaptive Threshold Design

- Test E vs C restricted to patients with biomarker B > b
  - Let T(b) be log likelihood ratio statistic
- Repeat for all values of b
- Let  $T = max{T(b)}$
- Compute null distribution of T by permuting treatment labels
- If the data value of T is significant at 0.01 level, then claim effectiveness of E for a patient subset
- Compute point and interval estimates of the optimal cut-point b

 $l(b) = \max \{ l(\mu, \eta, \gamma, b) \}$  $l(\mu, \eta, \gamma, b) = \log$  partial likelihood for model  $\log h(t) = \log h_0(t) + \mu \tau + \eta I(B > b) + \gamma \tau I(B > b)$  $\tau$  = binary treatment indicator  $\hat{b}$ =argmax{l(b)}  $\hat{\mathbf{b}}_* = \hat{b}$  value for bootstrap sample of cases  $\hat{F}_* =$  empirical distribution of  $\hat{b}_*$ *CI* for b based on percentiles of  $\hat{F}_*$  $\hat{F}_*(B)$  = probability patient with biomarker value B will benefit from treatment with E rather than C

Model	Hazard	Overall	Adaptive
	reduction	Power	Test
	for those		
	who benefit		
Everyone	33%	.775	.751
benefits			
50%	60%	.888	.932
benefit			
25%	60%	.429	.604
benefit			

#### Prostate Cancer DES (0.2 mg) vs Placebo

Covariate	#	Overall	Stage 2	Optimal
	Patients	Test	Test	Cut-Off
Acid	505	0.084	0.019	36
Phosphatate				

# Sample Size Planning (A)

- Standard broad eligibility trial were designed for 80% power to detect reduction in hazard D at significance level 5%
- Biomarker adaptive design is sized for 80% power to detect same reduction in hazard D at significance level 4% for overall analysis
Estimated Power of Broad Eligibility Design (n=386 events) vs Adaptive Design (n=412 events) 80% power for 30% hazard reduction

Model	Broad Eligibility Design	Biomarker Adaptive Design
40% reduction in 50% of patients (20% overall reduction)	.70	.78
60% reduction in 25% of patients (20% overall reduction)	.65	.91
79% reduction in 10% of patients (14% overall reduction)	.35	.93



ARTICLE

### Critical Review of Published Microarray Studies for Cancer Outcome and Guidelines on Statistical Analysis and Reporting

Alain Dupuy, Richard M. Simon

- Background Both the validity and the reproducibility of microarray-based clinical research have been challenged. There is a need for critical review of the statistical analysis and reporting in published microarray studies that focus on cancer-related clinical outcomes.
- Methods Studies published through 2004 in which microarray-based gene expression profiles were analyzed for their relation to a clinical cancer outcome were identified through a Medline search followed by hand screening of abstracts and full text articles. Studies that were eligible for cur analysis addressed one or more outcomes that were either an event occurring during follow-up, such as death or relapse, or a therapeutic response. We recorded descriptive characteristics for all the selected studies. A critical review of outcome-related statistical analyses was undertaken for the articles published in 2004.
- Results Ninety studies were identified, and their descriptive characteristics are presented. Sixty-eight (76%) were published in journals of impact factor greater than 6. A detailed account of the 42 studies (47%) published in 2004 is reported. Twenty-one (50%) of them contained at least one of the following three basic flaws: 1) in outcome-related gene finding, an unstated, unclear, or inadequate control for multiple testing; 2) in class discovery, a spurious claim of correlation between clusters and clinical outcome, made after clustering samples using a selection of outcome-related differentially expressed genes; or 3) in supervised prediction, a biased estimation of the prediction accuracy through an incorrect cross-validation procedure.
- Conclusions The most common and serious mistakes and misunderstandings recorded in published studies are described and illustrated. Based on this analysis, a proposal of guidelines for statistical analysis and reporting for clinical microarray studies, presented as a checklist of "Do's and Don'ts," is provided.

147 (1 of 11)

J Natl Cancer Inst 2007;99:147-57

DNA microarray technology has found many applications in biomedical research. In oncology, it is being used to better understand the biological mechanisms underlying oncogenesis, to discover new targets and new drugs, and to develop classifiers (predictors of good outcome versus poor outcome) for tailoring individualized treatments (1–4). Microarray-based clinical research is a recent and active area, with an exponentially growing number of publications. Both the reproducibility and validity of findings have been challenged, however (5,6). In our experience, microarray-based clinical investigations have generated both unrealistic hype and excessive skepticism. We reviewed published microarray studies in which gene expression data are analyzed for relationships with cancer outcomes, and we propose guidelines for statistical analysis and reporting, based on the most common and serious problems identified.

Medicine, followed by hand screening of abstracts and articles. The detailed process of selection is presented in Supplementary Note 1 (available online). The inclusion criteria were as follows: the work was an original clinical study on human cancer patients, published in English before December 31, 2004; it analyzed gene expression data of more than 1000 spots; and it presented statistical analyzes relating the gene expression profiling to a clinical outcome. Two types of outcome were considered: 1) A relapse or death occurring during the course of the disease. 2) A therapeutic response.

Correspondence to: Richard M. Simon, DSc, National Cancer Institute, 9000 Rockville Rike, MSC 7434, Bethesda, MD 20802 (e-mail: rsimon@nih.gov).

Affiliations of sechors: Biomotric Research Branch, Division of Cancer Treatment and Diagnosis, National Cancer Institute, National Institutes of Health, Bothesda, MD (AD, FMS); Univensita Paris VII Denis Diskord, Paris, France(AD); Assistance Publique-Hépitaux de Paris, Service de Dermatologie, Hépital Saint-Louis, Paris, France (AD).

## Major Flaws Found in 40 Studies Published in 2004

- Misleading use of cluster analysis
  - 13/28 studies invalidly claimed that expression clusters based on differentially expressed genes could help distinguish clinical outcomes
- Inadequate control of multiple comparisons in gene finding
  - 9/23 studies had unclear or inadequate methods to deal with false positives
    - 10,000 genes x .05 significance level = 500 false positives
- Misleading report of prediction accuracy
  - 12/28 reports based on incomplete cross-validation
- 50% of studies contained one or more major flaws

## t Tools Window Help

🕅 Sea	rch [ 🕎 ] 🕨 S	Select 👔	🗃 🔍 - 🚺 🕶 😁 100% - 📀 📑	• 🕜 Help → 🛛 🏹 🆓		
	9	Don't	Transform time-to-outcome data into a binary outcome	Use statistical methods suited for time-to-event data, unless		
			variable if the goal is to predict groups with different	you can ensure the absence of bias due to transformation.		
			survival probabilities.	See text and Supplementary Fig. 2 (available online).		
	Outcome-related gene finding†					
	10	Don't	Use only fold changes between groups to select the differentially expressed gapes	This does not take into account the variance of the genes' data values		
	11	Don't	Use a .05 P value threshold to select the differentially	A set of 10,000 genes will yield on average 500 false-positive		
		Don't	avorassad nanos	nones if this threshold is used		
	12	Do	Use a method for controlling the number of falsely	Lowering the P value threshold for selection (e.g. to .001) is		
		20	differentially expressed genes.	the simplest method. Others are available.		
	13	Do	Use a permutation test to assess the probability of finding	The result should be significant at $.05 P$ value level.		
		20	the same number of differentially expressed genes as	The result of our polymetric at the France foreit.		
			the one you found from your dataset.			
	Class dis	scovery				
	14	Don't	Use class discovery methods if you are interested in	Supervised prediction should be used for this purpose. It		
			classifying new samples in the future.	utilizes the outcome information to optimize predictive accuracy. See text.		
	15	Don't	Use a selection of outcome-related differentially expressed	Supervised clustering leads to a spurious correlation between		
			genes if you intend to correlate cluster-defined classes with the outcome.	cluster and outcome. See text and Fig. 1.		
	16	Don't	Select the clustering method that gives the best result.	Class discovery should not be result driven.		
	17	Do	Use methods for testing the reproducibility of cluster	Assessing the reproducibility of cluster finding without		
			finding.	using external information makes class discovery more convincing. See text.		
	18	Don't	Use conventional statistical tests for computing the	These tests assume independence between class definition		
			statistical significance of genes that are differentially	and expression profile data, which is not the case for		
			expressed between two clusters.	cluster-defined classes.		
	Supervised prediction					
	19	Do	Frame a therapeutically relevant question and select a	Classifiers developed outside a specific therapeutically		
			homogeneous set of patients accordingly.	relevant context are unlikely to be useful and utilized. See text.		
	20	Don't	Violate the fundamental principle of classifier validation,	Most of the "Don't" items on validation procedures are		
			i.e., no preliminary use of the tested samples.	illustrations of how this principle can be violated. See text		
	21	Den't	Attempt to prodict cluster defined classes	and Fig. 2 and supplementary Fig.1 (available online). Classos should be defined independently from the exercises		
	21	Don t	Attempt to predict cluster-defined classes.	classes should be defined independently from the expression profile data.		
	Evaluating the prediction on a separate test set			1		
	22	Don't	Use any information from the test set for developing the	The test set is to be used exclusively for evaluating the		
			classifier.	classifier performance. See text and Fig. 2.		
	23	Do	Access the test set only once and only for testing the	The test set must not be used to choose the best classifier.		
			samples with the fully specified classifier developed	See text and Fig. 2.		
			from the training set.			
	24	Do	Use the same outcome definition as the one used in the			
			training set			

hecklis	st		Comment	
Evalua	ting the p	rediction with a cross-validation procedure		
25	Don't	Use all the samples from the dataset to develop the classifier and test them.	The resubstitution estimate is not a cross-validation procedure. See text and Fig. 2.	
26	Don't	Use the same feature selection for all iterations.	This inflates the estimate of the prediction accuracy. See text and Fig. 2.	
27	Don't	Perform a cross-validation procedure on a selection of outcome-related differentially expressed genes.	Idem. Invalid although commonly done.	
28	Do	Report the estimates for all the classification algorithms if several have been tested, not just the most accurate.		
29	Don't	Consider that testing a few additional independent samples adds value to a correctly cross-validated estimate of the classifier prediction accuracy.	However, this may be valuable if the additional samples are in sufficient number and are representative of the samples in which the classifier might be used in the future. See text.	
30	Do	Report the fully specified classifier with its parameters.	So it can be used by others. Parameters are obtained from the whole training set in a separate test set procedure and from the whole dataset in a cross-validation procedure.	
31	Do	Report the correctly validated sensitivity and specificity or positive and negative apparent predictive values (for a binary outcome).	Receiver-operating characteristic curves may also be used. See text.	
32	Don't	Use an odds ratio to assess the performance of the prediction (for a binary outcome).	The odds ratio is a measure of association, not of prediction accuracy. See text and Supplementary Fig. 3 (available online).	
33	Do	Report the statistical significance of the prediction accuracy and, even better, of the sensitivity and specificity (for a binary outcome).	It states the probability of obtaining a prediction accuracy as high as actually observed if there was no relationship between the expression data and the outcome. See text.	
34	Don't	Use a Fisher's exact test or chi-square test to assess the statistical significance of the prediction accuracy for a binary outcome.	They do not test the statistical significance of the prediction. See text and Supplementary Fig. 3 (available online).	
35	Do	Pay attention to the imbalance between outcome categories when interpreting the prediction accuracy of a binary outcome.	90% prediction accuracy may be inadequate if outcome categories are highly imbalanced. See text and Supplementary Fig. 3 (available online).	
36	Don't	Use the log-rank test for testing the difference in survival between cross-validated groups.	The test is invalid because of a dependency among cases after cross-validation.	
37	Don't	Use standard regression models, e.g., logistic regression or proportional hazards model, with cross-validated predicted groups.	Idem.	
38	Don't	Assess the utility of the prediction based on the value of the regression coefficient or on its <i>P</i> value from multivariable regression models.	Regression coefficients are poor measures of prediction accuracy, and the test of statistical significance simply assesses if the coefficient is different from 0. See text.	
39	Do	Assess the added value of the classifier by examining its performance within the levels of the standard prognostic factors.	Other approaches can be used. See text.	
40	Do	Assess the utility of the classifier in a clinical context, for the therapeutically relevant question, and plan, if appropriate, further studies for external validation.		

## Good Microarray Studies Have Clear Objectives

- Class Comparison
  - Find genes whose expression differs among predetermined classes, e.g. tissue or experimental condition
- Class Prediction
  - Prediction of predetermined class (e.g. treatment outcome) using information from gene expression profile
- Class Discovery
  - Discover clusters of specimens having similar expression profiles
  - Discover clusters of genes having similar expression profiles

# Class Comparison and Class Prediction

- Not clustering problems
- Supervised methods

## **Class Prediction**

- A set of genes is not a classifier
- Testing whether analysis of independent data results in selection of the same set of genes is not an appropriate test of predictive accuracy of a classifier

#### ORIGINAL ARTICLE

## Concordance among Gene-Expression-Based Predictors for Breast Cancer

Cheng Fan, M.S., Daniel S. Oh, Ph.D., Lodewyk Wessels, Ph.D., Britta Weigelt, Ph.D., Dimitry S.A. Nuyten, M.D., Andrew B. Nobel, Ph.D., Laura J. van't Veer, Ph.D., and Charles M. Perou, Ph.D.

#### ABSTRACT

### RACKGROUND

From the Departments of Genetics (C.E., D.S.O., C.M.P.J. Stabilistics and Operations Research (A.B.N.J. and Pathology and Laboratory Medicine (C.M.P.J. University of North Carolina at Chapel Hill and Lineberger Comprehensive Cancer Center, Chapel Hills and the Divisions of Diagnostic Oncology (L.W., 8.W., L.J.N.) and Radiotherapy (D.S.A.N.), the Nethelands Cancer Institute, Amsterdam. Address reprint requests to Dr. Perou at Lineberger Comprehensive Cancer Center, University of North Carolina at Chapel Hill, Campus Box 7255, Chapel Hill, NC 27599, or at cperoughmed.unc.edu.

Drs. Fan and Oh contributed equally to this article.

N Engl.j Med 2006;355:560-9. Copyright © 2006 Massofusette Madical Soviety.

From the Departments of Genetics (C.F., Gene-expression-profiling studies of primary breast tumors performed by differtions Research (A.B.N., and Pathologrand Laboratories have resulted in the identification of a number of distinct prognostic profiles, or gene sets, with little overlap in terms of gene identity.

#### METHODS

To compare the predictions derived from these gene sets for individual samples, we obtained a single data set of 295 samples and applied five gene-expression-based models: intrinsic subtypes, 70-gene profile, wound response, recurrence score, and the two-gene ratio (for patients who had been treated with tamoxifen).

#### RESULTS

We found that most models had high rates of concordance in their outcome predictions for the individual samples. In particular, almost all tumors identified as having an intrinsic subtype of basal-like, HER2-positive and estrogen-receptor-negative, or luminal B (associated with a poor prognosis) were also classified as having a poor 70-gene profile, activated wound response, and high recurrence score. The 70-gene and recurrence-score models, which are beginning to be used in the clinical setting, showed 77 to 81 percent agreement in outcome classification.

#### CONCLUSIONS

Even though different gene sets were used for prognostication in patients with breast cancer, four of the five tested showed significant agreement in the outcome predictions for individual patients and are probably tracking a common set of biologic phenotypes.



N ENGL J MED 355:6 WWW.NEJM.ORG AUCUST 10, 2006

# Myth

 Complex classification algorithms such as neural networks perform better than simpler methods for class prediction.

- Artificial intelligence sells to nonspecialists who cannot distinguish hype from substance.
- Comparative studies generally indicate that simpler methods work as well or better for microarray problems because they avoid over-fitting the data.

## Linear Classifiers for Two Classes

- Fisher linear discriminant analysis
- Diagonal linear discriminant analysis (DLDA) assumes features are uncorrelated
- Compound covariate predictor (Radmacher et al )
- Weighted voting (Golub et al.)
- Support vector machines with inner product kernel
- Perceptron (Khan et al.)

# **Other Simple Methods**

- Nearest neighbor classification
- Nearest k-neighbors
- Nearest centroid classification
- Shrunken centroid classification

## 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

# **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

## **Non-Cross-Validated Prediction**

log-expression ratios



Prediction rule is built using full data set.
 Rule is applied to each specimen for class prediction.

## **Cross-Validated Prediction (Leave-One-Out Method)**



- 1. Full data set is divided into training and test sets (test set contains 1 specimen).
- 2. Prediction rule is built from scratch using the training set.
- 3. Rule is applied to the specimen in the test set for class prediction.
- 4. Process is repeated until each specimen has appeared once in the test set.

- 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.
  - Simon R, Radmacher MD, Dobbin K, McShane LM. Pitfalls in the analysis of DNA microarray data. Journal of the National Cancer Institute 95:14-18, 2003.
- The cross-validated estimate of misclassification error is an estimate of the prediction error for model fit using specified algorithm to full dataset

# Myth

• Split sample validation is superior to LOOCV for estimating prediction error

## Prediction Error Estimation: A Comparison of Resampling Methods

Annette M. Molinaro ab\* Richard Simon , Ruth M. Pfeiffer\*

<sup>a</sup>Biostatistics Branch, Division of Cancer Epidemiology and Genetics, NCI, NIH, Rockville, MD 20852, <sup>b</sup>Department of Epidemiology and Public Health, Yale University School of Medicine, New Haven, CT 06520, <sup>c</sup>Biometric Research Branch, Division of Cancer Treatment and Diagnostics, NCI, NIH, Rockville, MD 20852

### ABSTRACT

Motivation: In genomic studies, thousands of features are collected on relatively few samples. One of the goals of these studies is to build classifiers to predict the outcome of future observations. There are three inherent steps to this process: feature selection, model selection, and prediction assessment. With a focus on prediction assessment, we compare several methods for estimating the 'true' prediction error of a prediction model in the presence of feature selection.

Results: For small studies where features are selected from thousands of candidates, the resubstitution and simple splitsample estimates are seriously biased. In these small samples, leave-one-out (LOOCV), 10-fold cross-validation (CV), and the .632+ bootstrap have the smallest bias for diagonal discriminant analysis, nearest neighbor, and classification trees. LOOCV and 10-fold CV have the smallest bias for linear discriminant analysis. Additionally, LOOCV, 5- and 10-fold CV, and the .632+ bootstrap have the lowest mean square error. The .632+ bootstrap have the lowest mean square error. The .632+ bootstrap is quite biased in small sample sizes with strong signal to noise ratios. Differences in performance among resampling methods are reduced as the number of specimens available increase.

Availability: A complete compliation of results in tables and figures is available in Molinaro *et al.* (2005). R code for simulations and analyses is available from the authors. Contact: annette molinaro@vale.edu

#### **1 INTRODUCTION**

In genomic experiments one frequently encounters high dimensional data and small sample sizes. Microarrays simultaneously monitor expression levels for several thousands of genes. Proteomic profiling studies using SELDI-TOF (surface-enhanced laser desorption and ionization time-offlight) measure size and charge of proteins and protein fragments by mass spectroscopy, and result in up to 15,000 intensity levels at prespecified mass values for each spectrum. Sample sizes in such experiments are typically less than 100.

"to whom correspondence should be addressed

In many studies observations are known to belong to predetermined classes and the task is to build predictors or classifiers for new observations whose class is unknown. Deciding which genes or proteomic measurements to include in the prediction is called *feature selection* and is a crucial step in developing a class predictor. Including too many noisy variables reduces accuracy of the prediction and may lead to over-fitting of data, resulting in promising but often non-reproducible results (Ransohoff, 2004).

Another difficulty is model selection with numerous classification models available. An important step in reporting results is assessing the chosen model's error rate, or generalizability. In the absence of independent validation data, a common approach to estimating predictive accuracy is based on some form of resampling the original data, e.g., crossvalidation. These techniques divide the data into a learning set and a test set and range in complexity from the popular learning-test split to v-fold cross-validation, Monte-Carlo vfold cross-validation, and bootstrap resampling. Few comparisons of standard resampling methods have been performed to date, and all of them exhibit limitations that make their conclusions inapplicable to most genomic settings. Farly comparisons of resampling techniques in the literature are focussed on model selection as opposed to prediction error estimation (Breiman and Spector, 1992; Burman, 1989). In two recent assessments of resampling techniques for error estimation (Braga-Neto and Dougherty, 2004; Efron, 2004), feature selection was not included as part of the resampling procedures, causing the conclusions to be inappropriate for the high-dimensional setting.

We have performed an extensive comparison of resampling methods to estimate prediction error using simulated (large signal to noise ratio), microarray (intermediate signal to noise ratio) and proteomic data (low signal to noise ratio), encompassing increasing sample sizes with large numbers of features. The impact of feature selection on the performance of various cross validation methods is highlighted. The results clucidate the 'best' resampling techniques for

# **BRB-ArrayTools**

- Contains analysis tools that I have selected as valid and useful
- Analysis wizzard and multiple help screens for biomedical scientists
- Imports data from all platforms and major databases

## Predictive Classifiers in BRB-ArrayTools

- Classifiers
  - Diagonal linear discriminant
  - Compound covariate
  - Bayesian compound covariate
  - Support vector machine with inner product kernel
  - K-nearest neighbor
  - Nearest centroid
  - Shrunken centroid (PAM)
  - Random forrest
  - Tree of binary classifiers for kclasses
- Survival risk-group
  - Supervised pc's

- Feature selection options
  - Univariate t/F statistic
  - Hierarchical variance option
  - Restricted by fold effect
  - Univariate classification power
  - Recursive feature elimination
  - Top-scoring pairs
- Validation methods
  - Split-sample
  - LOOCV
  - Repeated k-fold CV
  - .632+ bootstrap

# **BRB-ArrayTools**

- Extensive built-in gene annotation and linkage to gene annotation websites
- Extensive gene-set enrichment tools for integrating gene expression with pathways and other biological information
- Publicly available for non-commercial use – <u>http://linus.nci.nih.gov/brb</u>

# **BRB-ArrayTools**

December 2006

- 6635 Registered users
- 1938 Distinct institutions
- 68 Countries
- 311 Citations

## Conclusions

 Developments in biotechnology and tumor biology make it increasingly feasible to identify which patients are most likely to benefit from a specified treatment

## Achieving the potential of new technology requires

- Paradigm changes in study design, moving from "correlative science" to predictive medicine
- New organizational structures and resource allocations to foster excellence in interdisciplinary research among biostatistical, laboratory and clinical scientists
  - Traditional core support structures are ineffective for high level collaboration
  - Major studies continue to be poorly designed and analyzed
  - Over-emphasis on software engineering at the expense of biostatistical collaboration
- FDA policies that encourage development of classifier targeted therapeutics

## Collaborators

- Alain Dupuy
- Boris Freidlin
- Wenyu Jiang
- Aboubakar Maitournam
- Yingdong Zhao

## Using Genomic Classifiers In Clinical Trials

•Dupue A and Simon R. Critical review of published microarray studies for clinical outcome and guidelines for statistical analysis and reporting, Journal of the National Cancer Institute 99:147-57, 2007

•Dobbin K and Simon R. Sample size planning for developing classifiers using high dimensional DNA microarray data. Biostatistics 8:101-117, 2007.

•Simon R. Development and validation of therapeutically relevant predictive classifiers using gene expression profiling, Journal of the National Cancer Institute 98:1169-71, 2006.

•Simon R. Validation of pharmacogenomic biomarker classifiers for treatment selection. Cancer Biomarkers 2:89-96, 2006.

•Simon R. Guidelines for the design of clinical studies for development and validation of therapeutically relevant biomarkers and biomarker classification systems. In Biomarkers in Breast Cancer, Hayes DF and Gasparini G, pp 3-15, Humana Press, 2006.

•Simon R. A checklist for evaluating reports of expression profiling for treatment selection. Clinical Advances in Hematology and Oncology 4:219-224, 2006.

•Simon R. Identification of pharmacogenomic biomarker classifiers in drug development. In Pharmacogenomics, Anti-cancer Drug Discovery and Response, F Innocenti (ed), Humana Press (In Press).

•Simon R, Lam A, Li MC, et al. Analysis of gene expression data using BRB-ArrayTools, Cancer Informatics 2:1-7, 2006

•Simon R. New challenges for 21st century clinical trials, Controlled Clinical Trials (In Press).

•Simon R. Development and validation of biomarker classifiers for treatment selection, Journal of Statistical Planning and Inference (In Press).

## Using Genomic Classifiers In Clinical Trials

•Simon R and Maitnourim A. Evaluating the efficiency of targeted designs for randomized clinical trials. Clinical Cancer Research 10:6759-63, 2004.

•Maitnourim A and Simon R. On the efficiency of targeted clinical trials. Statistics in Medicine 24:329-339, 2005.

•Simon R. When is a genomic classifier ready for prime time? Nature Clinical Practice – Oncology 1:4-5, 2004.

•Simon R. An agenda for Clinical Trials: clinical trials in the genomic era. Clinical Trials 1:468-470, 2004.

•Simon R. Development and validation of therapeutically relevant multi-gene biomarker classifiers. Journal of the National Cancer Institute 97:866-867, 2005.

•Simon R. A roadmap for developing and validating therapeutically relevant genomic classifiers. Journal of Clinical Oncology 23:7332-41,2005.

•Freidlin B and Simon R. Adaptive signature design. Clinical Cancer Research 11:7872-78, 2005.

•Simon R. and Wang SJ. Use of genomic signatures in therapeutics development in oncology and other diseases, The Pharmacogenomics Journal 6:166-73, 2006.

•Trepicchio WL, Essayan D, Hall ST, Schechter G, Tezak Z, Wang SJ, Weinreich D, Simon R. Designing prospective clinical pharmacogenomic trials- Effective use of genomic biomarkers for use in clinical decision-making. The Pharmacogenomics Journal 6:89-94,2006.