

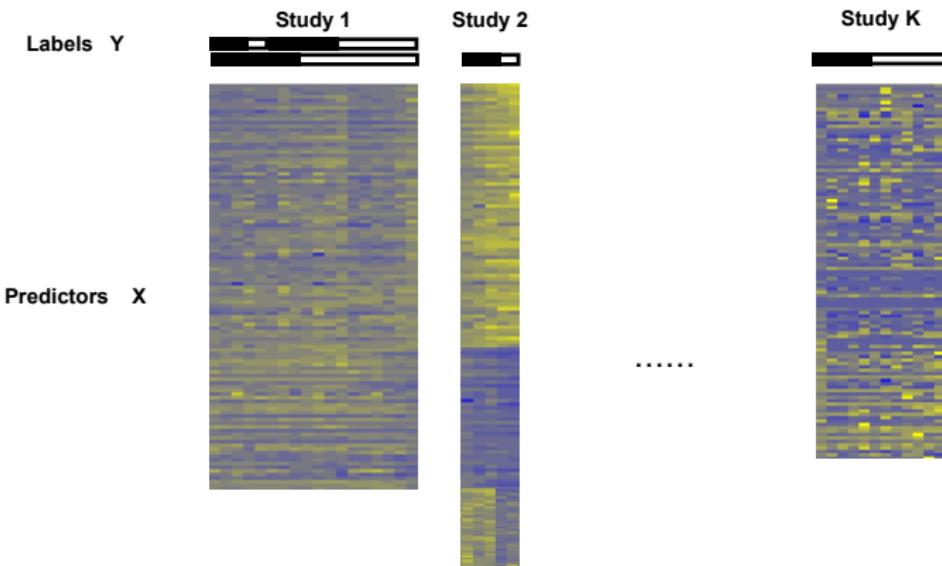
Replicability in Science: 2: The National Academies 2019 Report

giovanni_parmigiani@dfci.harvard.edu

Padova, September 19, 2023

- Left-overs
- Discussion 1: Who celebrates
- Discussion 2: Apples and CDs
- Extent of Non-Replicability in Research
- Sources of Non-Replicability in Research
- Improving Replicability in Research

data structure



False Discovery Rates

False Discovery Rates

A simple way to think about the discovery process is to focus on a list of promising candidates (or “discoveries”), as we did in our previous lecture. If we know the true data generating model for each gene, we could, for example, look at this table:

	No Discovery	Discovery	Total
Null hypotheses	U	R_0	G_0
Alternative hypotheses	T	R_A	$G - G_0$
Total	$G - R$	R	G

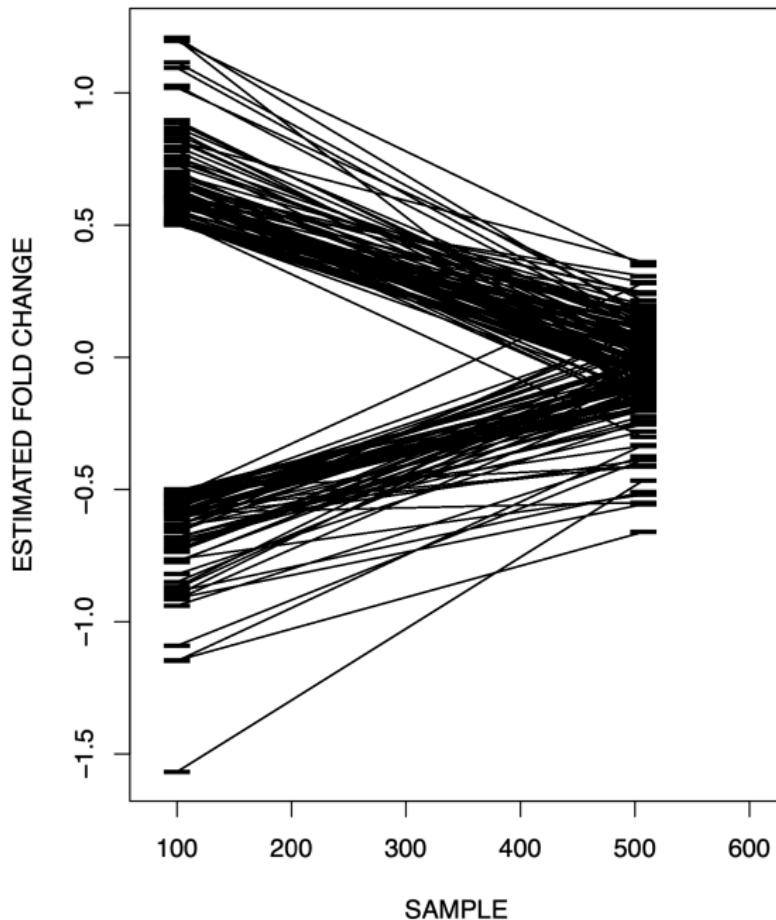
It is common to study the proportion R_0/R of discoveries for which the true data generating model is null.

Note

R_0/R is unknown. It depends on both the data and the parameter (the vector of indicators of whether each gene is null)

- When is R_0/R an appropriate quantity to focus on?

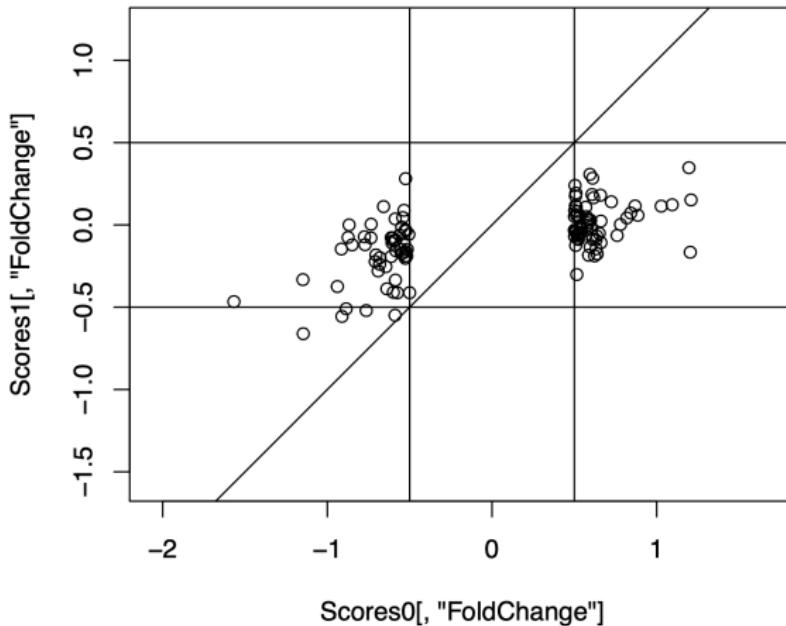
"half the effect size"



"half the effect size part 2: shrinkage"

Scatterplot version of the same data, to emphasize “shrinkage”.

```
Scores0 = ScoresSub[fcOnly,]
Scores1 = CompScores(XXX[fcOnly,],YYY)
plot(Scores0[,"FoldChange"],Scores1[,"FoldChange"],
     asp=1,ylim=range(Scores0[,"FoldChange"]))
abline(h=fcut); abline(h=-fcut); abline(v=fcut)
abline(v=-fcut); abline(0,1)
```



classifier cross-study validation matrix

Waldron et al 2014

Implemented Models

Validation Statistics for 14 Models in 10 Datasets

	Dataset Average	1.81	1.47	1.43	1.41	1.39	1.37	1.35	1.14	1.11	1.04	
TCGA11	2.05	2.28	1.58	1.85	1.36	1.64	1.97	1.94	1.07	1.07	1.53	TCGA11
Yoshihara12	2.44	9.65	1.21	1.8	1.02	1.77	1.97	1.21	1.35	1.35	1.42	Yoshihara1
Yoshihara10	2.69	1.38	1.15	1.93	1.45	1.57	1.47	1.33	0.7	0.7	7.27	Yoshihara1
Kernagis12	2.65	1.39	2.91	2.08	1.45	1.39	1.25	1.23	1.32	1.32	0.87	Kernagis12
Crijns09	1.22	1.92	1.49	1.51	1.2	1.44	1.28	1.1	3.04	3.04	1.21	Crijns09
Bentink12	1.94	1.01	1.89	1.44	1.2	1.14	1.62	1.16	1.26	1.26	1.45	Bentink12
Bonomo08_263genes	1.3	2.73	2	1.32	0.53	2.01	1.45	1.17	1.03	1.03	0.77	Bonomo08
Mok09	1.54	1.82	3.18	1.71	0.89	1.58	1.28	0.98	0.95	0.95	1.39	Mok09
Bonomo08_572genes	0.8	1.89	1.1	1.41	2.29	2.27	1.47	1.07	1.35	1.35	0.84	Bonomo08
Sabatier11	1.95	1.15	1.17	1.41	1.72	1.07	1.19	1.11	1.3	1.3	0.73	Sabatier11
Denkert09	2.6	0.76	1.33	1.31	2.25	1.04	1.29	1.08	1.15	1.15	0.79	Denkert09
Kang12	2.14	1.19	0.81	0.85	1.21	1.46	1.17	1.55	1.02	1.02	0.73	Kang12
Konstantinopoulos10	1.34	1.01	0.82	1.07	2.05	1.07	1	1.15	0.97	0.97	1.09	Konstantinopoulos10
Hernandez10	0.68	0.55	1.07	0.71	0.86	1.21	0.79	1.04	0.9	0.9	1.03	Hernandez10

Expression Datasets

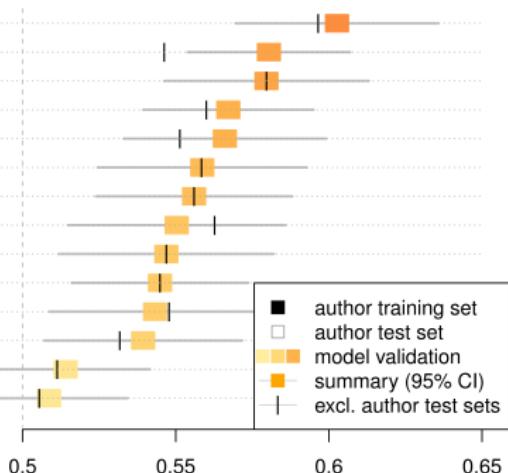
Dressman
Yoshihara 2012A
Mok
Tothill
Konstantinopoulos
Bonomo
Bentink
TCGA
Crijns
Yoshihara 2010

using c-stat instead

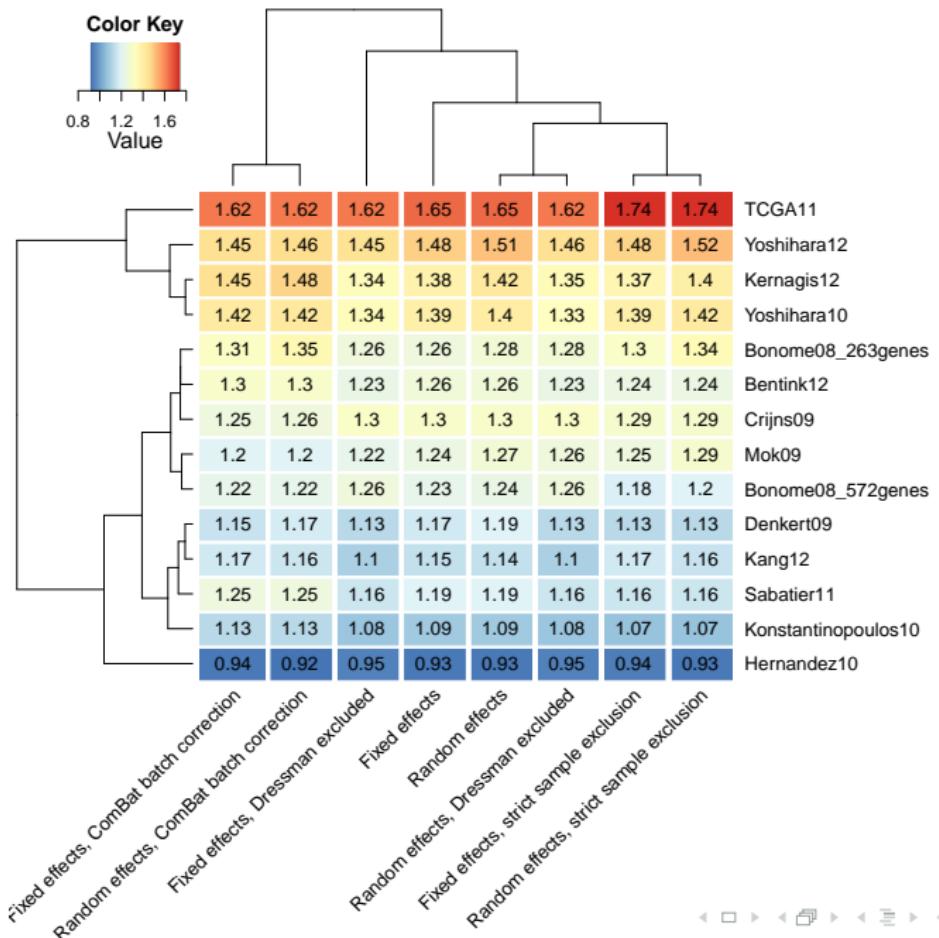
(A) Implemented Models Validation Statistics for 14 Models in 10 Datasets

<i>Dataset</i>	<i>Average</i>	0.61	0.58	0.57	0.56	0.56	0.55	0.55	0.54	0.54	0.53
TCGA11		0.62	0.69	0.6	0.63	0.61	0.47	0.57	0.6	0.64	0.55
Yoshihara12		0.63	0.81	0.64	0.6	0.62	0.51	0.5	0.58	0.57	0.55
Bonome08_263genes		0.57	0.68	0.58	0.6	0.62	0.53	0.6	0.54	0.56	0.52
Yoshihara10		0.7	0.55	0.62	0.53	0.55	0.53	0.54	0.8	0.56	0.52
Kernagis12		0.66	0.58	0.63	0.56	0.55	0.55	0.65	0.57	0.55	0.54
Sabatier11		0.64	0.54	0.56	0.57	0.54	0.62	0.55	0.57	0.56	0.52
Crijns09		0.5	0.6	0.59	0.55	0.58	0.55	0.56	0.47	0.54	0.67
Bentink12		0.65	0.56	0.55	0.61	0.55	0.57	0.57	0.53	0.53	0.52
Bonome08_572genes		0.57	0.6	0.54	0.55	0.64	0.63	0.55	0.5	0.53	0.54
		0.53	0.6	0.56	0.57	0.57	0.53	0.69	0.57	0.51	0.51
Kang12		0.63	0.54	0.52	0.54	0.57	0.54	0.49	0.54	0.58	0.52
Denkert09		0.67	0.52	0.54	0.53	0.53	0.58	0.53	0.51	0.52	0.55
Hernandez10		0.56	0.61	0.56	0.54	0.53	0.5	0.5	0.54	0.49	0.51
Konstantinopoulos10		0.57	0.5	0.52	0.48	0.49	0.6	0.5	0.51	0.53	0.5
<i>Expression Datasets</i>	Dressman	Yoshihara 2012A	Tothill	Bentink	Bonome	Konstantinopoulos	Mok	Yoshihara 2010	TCGA	Crijns	

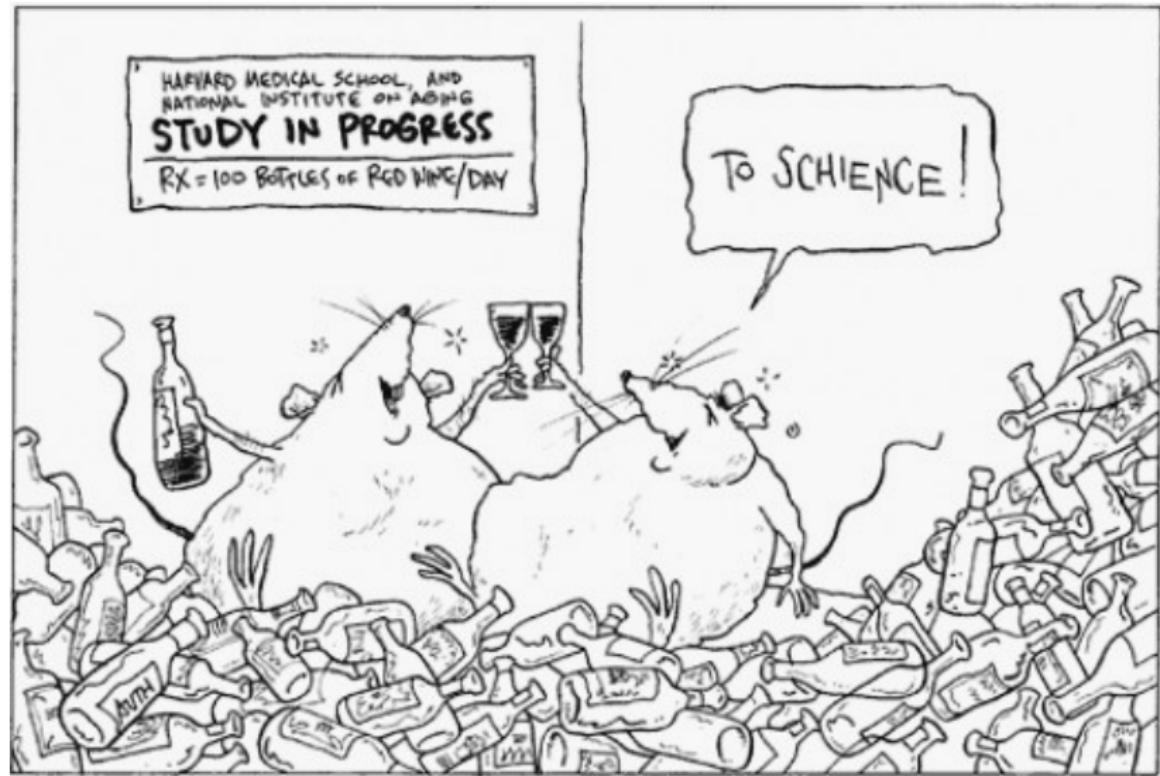
(B) Summary Validation Statistics



sensitivity analysis



DISCUSSION 1: WHO CELEBRATES?



no replicability. who wins?

Imagine a world-wide committee of replicability experts came to the conclusion that a specific field of science is not replicable.
Who would benefit?

replicability is built in. who wins?

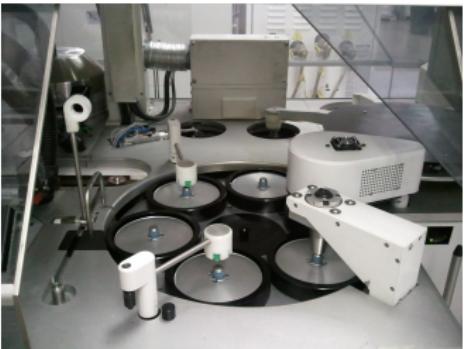
Imagine the same world-wide committee of replicability experts came to the conclusion that a specific field of science is intrinsically self correcting and thus that replicability is no concern.

Who would benefit?

DISCUSSION 2: APPLES AND CDS



...



replicability or variability?

- No two experiments will ever be the same.
Replicability is not a separate question from variability.
Tweet: Replicability does not exist".
- "No knowledge can be safely used in practice in future instances if it not replicable.
Tweet: Science does not exist without replicability."

EXTENT OF NON-REPLICABILITY



recommendation 5-1:

Researchers should, as applicable to the specific study, provide an accurate and appropriate characterization of relevant uncertainties when they report or publish their research.

Researchers should thoughtfully

communicate all recognized uncertainties and

estimate or acknowledge other potential sources of

uncertainty that bear on their results, including

stochastic uncertainties and

uncertainties in measurement, computation,

knowledge, modeling, and methods of analysis.

hues of replicability

NAS: Importantly, the assessment of replicability may not result in a binary pass/fail answer; rather, the answer may best be expressed as the degree to which one result replicates another.

restrictive and unreliable

NAS: it is restrictive and unreliable to accept replication only when the results in both studies have attained "statistical significance," that is, when the p-values in both studies have exceeded a selected threshold.

... rather

NAS: Rather, in determining replication, it is important to consider the distributions of observations and to examine how similar these distributions are. This examination would include summary measures, such as proportions, means, standard deviations (uncertainties), and additional metrics tailored to the subject matter.

cumulative evidence

NAS: A predominant focus on the replicability of individual studies is an inefficient way to assure the reliability of scientific knowledge. Rather, reviews of cumulative evidence on a subject, to assess both the overall effect size and generalizability, is often a more useful way to gain confidence in the state of scientific knowledge.

many replication studies are not reported

NAS: Because many scientists routinely conduct replication tests as part of a follow-on experiment and do not report replication results separately, the evidence base of non-replicability across all science and engineering research is incomplete.

SOURCES OF NON-REPLICABILITY



helpful and not

NAS: Non-replicability occurs for a number of reasons that do not necessarily reflect that something is wrong. Some occurrences of non-replicability may be helpful to science —discovering previously unknown effects or sources of variability— while others, ranging from simple mistakes to methodological errors to bias and fraud, are not helpful.

IMPROVING REPLICABILITY

CS452690



"OF COURSE YOU CAN'T REPLICATE MY...
EXPERIMENTS. THAT'S THE BEAUTY OF THEM."

recommendation 6-7

Journals and scientific societies requesting submissions for conferences should disclose their policies relevant to achieving reproducibility and replicability. The strength of the claims made in a journal article or conference submission should reflect the reproducibility and replicability standards to which an article is held, with **stronger claims reserved for higher expected levels of reproducibility and replicability.**

recommendation 6-7, continued

- ① set and implement desired standards of reproducibility and replicability and make this one of their priorities, such as deciding which level they wish to achieve for each Transparency and Openness Promotion guideline and working towards that goal;
- ② adopt policies to reduce the likelihood of non-replicability, such as considering incentives or requirements for research materials transparency, design, and analysis plan transparency, enhanced review of statistical methods, study or analysis plan preregistration, and replication studies; and
- ③ require as a review criterion that all research reports include a thoughtful discussion of the uncertainty in measurements and conclusions. inferences.

recommendation 6-8

Many considerations enter into decisions about what types of scientific studies to fund, including striking a balance between exploratory and confirmatory research.

recommendation 6-8, continued

If private or public funders choose to invest in initiatives on reproducibility and replication, two areas may benefit from additional funding:

- ① education and training initiatives to ensure that researchers have the knowledge, skills, and tools needed to conduct research in ways that adhere to the highest scientific standards; that describe methods clearly, specifically, and completely; and that express accurately and appropriately the uncertainty involved in the research; and
- ② reviews of published work, such as testing the reproducibility of published research, conducting rigorous replication studies, and publishing sound critical commentaries.

recommendation 6-9

Funders should require a thoughtful **discussion in grant applications of how uncertainties will be evaluated**, along with any relevant issues regarding replicability and computational reproducibility.

Funders should introduce **review of reproducibility and replicability guidelines and activities** into their merit-review criteria, as a low-cost way to enhance both.

recommendation 6-10

When funders, researchers, and other stakeholders are considering whether and where to direct resources for replication studies, they should consider the following criteria:

- ① The scientific results are important for individual decision-making or for policy decisions.
- ② The results have the potential to make a large contribution to basic scientific knowledge.
- ③ The original result is particularly surprising, that is, it is unexpected in light of previous evidence and knowledge.
- ④ There is controversy about the topic.

recommendation 6-10, continued

- ⑤ There was potential bias in the original investigation, due, for example, to the source of funding.
- ⑥ There was a weakness or flaw in the design, methods, or analysis of the original study.
- ⑦ The cost of a replication is offset by the potential value in reaffirming the original results.
- ⑧ Future expensive and important studies will build on the original scientific results.

too much replication?



DIGRESSIONS

the tale of the worms and the red wine



Phillips, P., Lithgow, G.J., and Driscoll, M. (2017).
A Long Journey to Reproducible Results. *Nature*, 548(7668), 387-388.

the tale of the worms and the red wine

Some research had found that resveratrol (found in red wine) could dramatically extend the life of worms in the lab, but other scientist had difficulty replicating the results.

Reasons for this lack of replicability:

Differences in lab protocol:

e.g. worms handled by gentle lab technicians lived a full day longer than others.

Difference in measurement:

e.g. one lab determined age on the basis of when an egg was laid; another began counting when it was hatched.

Once these were eliminated, researchers discovered inherent variability in the model system:

e.g. some cohorts of worms could partition into short-lived or long-lived modes of aging (previously unknown.)

Mechanisms of disease

Use of proteomic patterns in serum to identify ovarian cancer

Emanuel F Petricoin III, Ali M Ardekani, Ben A Hitt, Peter J Levine, Vincent A Fusaro, Seth M Steinberg, Gordon B Mills, Charles Simone, David A Fishman, Elise C Kohn, Lance A Liotta

Summary

Background New technologies for the detection of early-stage ovarian cancer are urgently needed. Pathological changes within an organ might be reflected in proteomic patterns in serum. We developed a bioinformatics tool and used it to identify proteomic patterns in serum that distinguish neoplastic from non-neoplastic disease within the ovary.

Methods Proteomic spectra were generated by mass spectroscopy (surface-enhanced laser desorption and ionisation). A preliminary "training" set of spectra derived from analysis of serum from 50 unaffected women and 50 patients with ovarian cancer were analysed by an iterative searching algorithm that identified a proteomic pattern that completely discriminated cancer from non-cancer. The discovered pattern was then used to classify an independent set of 116 masked serum samples: 50 from women with ovarian cancer, and 66 from unaffected women or those with non-malignant disorders.

Findings The algorithm identified a cluster pattern that, in the training set, completely segregated cancer from non-cancer. The discriminatory pattern correctly identified all 50 ovarian cancer cases in the masked set, including all 18 stage I cases. Of the 66 cases of non-malignant disease, 63 were recognised as not cancer. This result yielded a sensitivity of 100% (95% CI 93–100), specificity of 95% (87–99), and positive predictive value of 94% (84–99).

Introduction

Application of new technologies for detection of ovarian cancer could have an important effect on public health,¹ but to achieve this goal, specific and sensitive molecular markers are essential.^{1–5} This need is especially urgent in women who have a high risk of ovarian cancer due to family or personal history of cancer, and for women with a genetic predisposition to cancer due to abnormalities in predisposition genes such as *BRCA1* and *BRCA2*. There are no effective screening options for this population.

Ovarian cancer presents at a late clinical stage in more than 80% of patients,¹ and is associated with a 5-year survival of 35% in this population. By contrast, the 5-year survival for patients with stage I ovarian cancer exceeds 90%, and most patients are cured of their disease by surgery alone.^{1–6} Therefore, increasing the number of women diagnosed with stage I disease should have a direct effect on the mortality and economics of this cancer without the need to change surgical or chemotherapeutic approaches.

Cancer antigen 125 (CA125) is the most widely used biomarker for ovarian cancer.^{1–4} Although concentrations of CA125 are abnormal in about 80% of patients with advanced-stage disease, they are increased in only 50–60% of patients with stage I ovarian cancer.^{1–6} CA125 has a positive predictive value of less than 10% as a single marker, but the addition of ultrasound screening to CA125 measurement has improved the

- The researchers, from NIH & FDA, won widespread acclaim.
- Congressional resolution urged further funding for their research.
- The magazine "Health" named the test one of the top ten medical advances of the year.
- Commercial rights to develop the test were licensed from the US government to Correlogic Systems.
- Correlogic granted licenses to Quest Diagnostics and the Laboratory Corporation of America, hoping to market the test under the brand name OvaCheck.

Slide based on a presentation by S. Goodman at the National Cancer Policy Forum, 2018



Reproducibility of SELDI-TOF protein patterns in serum: comparing datasets from different experiments

Keith A. Baggerly*, Jeffrey S. Morris and Kevin R. Coombes

Department of Biostatistics, U.T. M.D. Anderson Cancer Center, 1515 Holcombe Blvd,
Box 447, Houston, TX 77030-4009, USA

Received on July 14, 2003; revised on October 14, 2003; accepted on October 16, 2003
Advance Access publication January 29, 2004

Compared SELDI proteomic spectra from serum from three experiments by the same group on separating ovarian cancer from normal tissue. These spectra are available on the web at <http://clinicalproteomics.steem.com>.

Results: In general, the results were not reproducible across experiments. Baseline correction prevents reproduction of the results for two of the experiments. In one experiment, there is evidence of a **major shift in protocol mid-experiment** which could bias the results. In another, **structure in the noise regions of the spectra allows us to distinguish normal from cancer**, suggesting that the normals and cancers were processed differently. Sets of **features found to discriminate well in one experiment do not generalize** to other experiments. Finally, the mass calibration in all three experiments appears suspect. Taken together, these and other concerns suggest that much of the structure uncovered in these experiments could be due to artifacts of sample processing, not to the underlying biology of cancer. We provide some guidelines for design and analysis in experiments like these to ensure better reproducible, biologically meaningful results.