Two case studies in Genomics

WORKSHOP ON HUMAN GENOMIC VARIATION AND DISEASE: Norway, August 2008

> Robert Tibshirani Stanford University

The Tibshirani "lab"

- 4 grad students; work closely with Trevor Hastie and his students
- statistical and data mining techniques for genomic and proteomic data
- writing and supporting software packages for some of these tools Excel Add-ins and R language

Examples of our work

- SAM- Significance analysis of microarrays
- PAM- Prediction analysis of microarrays- for sample classification
- CGH-miner- for fiding "hot spots" in CGH data
- Superpc- Supervised principal components- survival analysis from genomic/proteomic data
- Currently- gene set analysis, complementary clustering



- two recent controversies in Cancer genomics: a gene expression study of lymphoma and a GWA of cancer
- some general suggestions for improving the state of statistical analyses in these areas



- Dave et al published a high-profile study in NEJM, reporting that they had found two sets of genes whose expression were highly predictive of survival in patients with Follicular Lyphoma.
- the paper got a lot of attention at the recent ASH meeting, because the genes in the clusters were largely expressed in non-tumor cells, suggesting that the host-response was the important factor
- One of my medical collaborators- Ron Levy, asked me to look over their paper- he wanted to apply their model to the Stanford FL patient population.

The NEW ENGLAND JOURNAL of MEDICINE

ESTABLISHED IN 1812 NOVEMBER 18, 2004 VOL.351 NO.21

Prediction of Survival in Follicular Lymphoma Based on Molecular Features of Tumor-Infiltrating Immune Cells

Sandeep S. Dave, M.D., George Wright, Ph.D., Bruce Tan, M.D., Andreas Rosenwald, M.D., Randy D. Gascoyne, M.D., Wing C. Chan, M.D., Richard I. Fisher, M.D., Rita M. Braziel, M.D., Lisa M. Rimsza, M.D., Thomas M. Grogan, M.D., Thomas P. Miller, M.D., Michael LeBlanc, Ph.D., Timothy C. Greiner, M.D., Dennis D. Weisenburger, M.D., James C. Lynch, Ph.D., Julie Vose, M.D., James O. Armitage, M.D., Erlend B. Smeland, M.D., Ph.D., Stein Kvaloy, M.D., Ph.D., Harald Holte, M.D., Ph.D. Jan Delabie, M.D., Ph.D., Joseph M. Connors, M.D., Peter M. Lansdorp, M.D., Ph.D., Qin Ouyang, Ph.D., T. Andrew Lister, M.D., Andrew J. Davies, M.D., Andrew J. Norton, M.D., H. Konrad Muller-Hermelink, M.D., German Ott, M.D., Elias Campo, M.D., Emilio Montserrat, M.D., Wyndham H. Wilson, M.D., Ph.D., Elaine S. Jaffe, M.D., Richard Simon, Ph.D., Liming Yang, Ph.D., John Powell, M.S., Hong Zhao, M.S., Neta Goldschmidt, M.D., Michael Chiorazzi, B.A., and Louis M. Staudt, M.D., Ph.D.

ABSTRACT

BACKGROUND

Patients with follicular lymphoma may survive for periods of less than 1 year to more From National Cancer Institute (S.S.D., than 20 years after diagnosis. We used gene-expression profiles of tumor-biopsy specimens obtained at diagnosis to develop a molecular predictor of the length of survival.

METHODS

with specimens in the test set to be divided into four quartiles with widely disparate mewith specimens in the existence of contract and the second second

CONCLUSIONS

CONCLUSIONS Staudt at the Nanosci Lacer manuer The length of survival among patients with follicular lymphoma correlates with the suplemilar fratures of nonmalignant immune cells present in the tumor at diagnosis. 2009,20 et Istaud@mal.nl.gov.

N Engl | Med 2004:351:2159-69. pyright © 2004 Massachusetts Medical Societ

N ENGL J MED 351;21 WWW.NEJM.ORG NOVEMBER 18, 2004

2159

Downloaded from www.neim.org at Stanford University on May 10, 2005 Copyright © 2004 Massachusetts Medical Society. All rights reserved

G.W., B.T., A.R., W.H.W., E.S.J., R.S., H.Z., N.G., M.C., L.M.S.); Center for Information Technology (L.Y., J.P.); and National Heart, Lung, and Blood Institute (S.S.D.) — all in

 MTHOOS
 Betlevela, ML: Similer Columba Cancer

 Gene-expression profiling was performed on 191 biopsy specimens obtained from particles with untreated follicular lymphoma. Supervised methods were used to discover
 PML, Q.O.; DKC, JMC, PML, Q.O.; DWC, C, TCG, DM, Q.O.; DMC, TC, TCG, DM, CC, TCG, DM, JCL, LY, JOA): Southwest Onemas. A molecular predictor of survival area constructed from these genes and validat

 remsn. A molecular predictor of survival area constructed from these genes and validat Intervisity of Red-Nationary (JAC), TCG, TCG, DM, JCL, JV, JOA): Southwest Onemas.

 Individual genes that predicted the length of survival were grouped into gene-express Intervisity of Red-Nationary of Automa (JAC), TPM, ML, DAD, Southwest Onemas, Construct a survival predictor. The two signatures allowed patients with widely disparate me

 with specimens in the test set to be divided into four quartiles with widely disparate me IND Stattare (MAL): Nonevergin Redum

Bethesda, Md.: British Columbia Cancer Cancer Research UK. St. Bartholomew's sity of Barcelona, Barcelona, Spain (E.C., E.M.). Address reprint requests to Dr. Staudt at the National Cancer Institute.

Summary of their findings

- They started with the expression of approximately 49,000 genes measured on 189 patient samples, derived from DNA microarrays. A survival time (possibly censored) was available for each patient
- they randomly split the data into a training set of 89 patients and a test set of 90 patients
- using a multi-step procedure (described below), they extracted two clusters of genes, called IR1 (immune response 1) and IR2 (immune response 2).
- They averaged the gene expression of the genes in each cluster, to create two "super-genes".

... continued

- They then fit these super-genes together in a Cox model for survival, and applied it to the training and test sets. The p-value in the training set was < 10e 7 and 0.003 in the test set. IR1 correlates with good prognosis; IR2 with poor prognosis
- In the remainder of the paper they interpret the genes in their model



- I downloaded the data
- Applied some familiar statistical tools- eg SAM (Significance analysis of microarrays), less familiar ones- supervised principal components, and also gave the data to Brad Efron. Our initial finding- no significant correlation between gene expression and survival.
- I spent 2-3 weeks emailing back and forth with their statistician (George Wright) and programming in R, to recreate their analysis
- Confession- it was fun being a "forensic statistician"



Details of their analysis

- Divide the data randomly into training and test sets of approximately equal numbers of patients. Apply the following recipe [steps 2–6] to the training set.
- 2) Choose all genes with univariate Cox score > 1.5 in absolute value. This reduced the number of genes from roughly 49,000 to roughly 3,000, with about a 50-50 split between good prognosis genes (negative scores) and poor prognosis genes (positive scores).
- 3) Do separate hierarchical clusterings (correlation metric, average linkage) of the good and poor prognosis genes.

Details continued...

- 4) Find all clusters in the dendrograms (clustering trees) containing between 25 and 50 genes, with internal correlation at least 0.5. Represent each cluster by the average expression of all genes in the cluster- a "supergene" Try every pair of supergenes as predictors in Cox models for predicting survival.
- 5) Choose the most significant pair from this process. The authors call the resulting pair of clusters IR1 (good prognosis) and IR2 (poor prognosis).
- 6) Finally use the model (IR1, IR2) in a Cox model to predict survival in the test set.





PREDICTION OF SURVIVAL IN FOLLICULAR LYMPHOMA

N ENGL J MED 351;21 WWW.NEJM.ORG NOVEMBER 18, 2004

2163

Downloaded from www.nejm.org at Stanford University on May 10, 2005 . Copyright © 2004 Massachusetts Medical Society. All rights reserved.



Panel A shows overall survival among the patients with biopsy specimens in the test set, according to the quartile of the survival-predictor score (SPS). Panel B shows overall survival according to the International Prognostic Index (IPI) risk group for all the patients for whom these data were available. Panel C shows overall survival among the patients with specimens in the test set for the indicated IPI risk group, stratified according to the quartile of the SPS.

Downloaded from www.nejm.org at Stanford University on May 10, 2005 Copyright © 2004 Massachusetts Medical Society. All rights reserved.

2165

PREDICTION OF SURVIVAL IN FOLLICULAR LYMPHOMA





The total number of points (cluster pairs) with test set p-values less than 0.05 (239) is far fewer than we'd expect to see by chance (735)

Univariate p-values







There are only 85 pairs out of 11628 that are significant in the test set at the 0.05 level, while we would expect 11628*.05=581 pairs just by chance.





Cluster size ranges (30,60) rather than (25,50)



The aftermath

- I published a short letter to NEJM in March 2005; full details of my re-analysis appear on my website
- The authors published a rebuttal in the same issue.
 "Nothing in Tibshirani's analysis calls into dispute the fact that we discovered and validated a strong association between gene expression in follicular lymphoma and overall survival." Their arguments:
 - (1) we followed standard statistical procedures, found a small p-value on the test set, therefore our finding is correct;
 - (2) our method found an interaction, which SAM can't find
 - (3) we get small p-values if we apply our model to random halves of the data (????!!!!!)

CORRESPONDENCE



Our predictor could not have been discovered fraction is due to this 12.6 percent contamination with the SAM method, which relies solely on uni- requires that the lymphoma cells in the CD19variate associations with survival. Rather, our pre- fraction have expression of the immune-response dictor derives its strength from the synergistic signatures that was more than eight times as high combination of two gene-expression signatures in as that in their counterparts in the CD19+ fraction. a multivariate model. Tibshirani confuses the abil- Second, Hong et al. incorrectly discount the imity of our method to discover a survival association mune-response 1 signature, which contributes sigwith the fact that we actually found one that validat- nificantly to the survival model in the test set ed the association. When he exchanged the training (P<0.001). Third, many of the immune-response set for the test set, he was unable to rediscover our signature genes are selectively expressed in T cells, gene-expression predictor because some genes in monocytes, or dendritic cells, or in more than one our predictor fell below the P value threshold for of these, but not in B cells, making the contention association with survival in the test set. This does of Hong et al. even more implausible. not negate the fact that our model is highly associ-Louis M. Staudt, M.D., Ph.D. ated with survival in the test set. Hong et al. have made three errors. First, in our Sandeep Dave, M.D.

George Wright, Ph.D.

sorted subpopulations, the CD19– fraction con-tained, on average, 12.6 percent contamination Bethesda MD 20892 with follicular-lymphoma cells, not 25 percent, as Istaudt@mail.nih.gov they claim. To believe that the higher expression of 1. Ransohoff DF. Rules of evidence for cancer molecular-marker

the immune-response signatures in the CD19- discovery and validation. Nat Rev Cancer 2004;4:309-14.

When Doctors Go to War

TO THE EDITOR: Like Bloche and Marks in their Per- challenging circumstances.² Physicians are despective article on doctors in combat (Jan. 6 is- fined by their common calling to prevent harm and sue),1 the American Medical Association (AMA) treat people who are ill or injured and by their uniapplauds the outstanding work of military physi- versal commitment to uphold recognized principles cians in treating wounded soldiers under extremely of medical ethics whenever patients rely on their

N ENGLIMED 352:14 WWW.NEIM.ORG APRIL 7, 2005

1497

Downloaded from www.nejm.org at Stanford University on May 10, 2005 Copyright © 2005 Massachusetts Medical Society. All rights reserved.



- Their finding is *fragile*. I don't believe that it is real or reproducible
- This experience uncovers a problem that is of general importance to our field:
 - with many predictors, it is too easy to overfit the data and find spurious results
 - we can inadvertently mislead the reader, and mislead ourselves. I have been guilty of this too

Some recommendations

- encourage authors to publish not only the raw data, but a script of their analysis
- encourage authors to use "canned" methods/packages, with built-in cross-validation to validate the model search process see "supervised principal components", Bair et al 2005