

Building a Predictive Biomarker

Noah Simon and Richard Simon

July 2016

Predictive Biomarker for a Continuous Measure

Randomize $n_T + n_C = n$ patients to **treatment/control**

On each patient measure

y_i - single continuous outcome
(eg. blood pressure, tumor growth)

\mathbf{x}_i - p -vector of features
(eg. SNPs, gene expression values)

Want to determine who will benefit from **treatment** over **control**.

Predictive Models from Prognostic Models

Given “true prognostic models”:

$$y_T = \beta_0 + \mathbf{x}^\top \beta + \epsilon$$

$$y_C = \alpha_0 + \mathbf{x}^\top \alpha + \epsilon$$

Would treat patients with

$$E[y_T|x] - E[y_C|x] > \delta$$

(with $\delta = 0$ if non-toxic)

Simplifies to patients with

$$\beta_0 - \alpha_0 + \mathbf{x}_i^\top (\beta - \alpha) > \delta$$

Estimating Predictive Models

Don't know “true models” so we estimate:

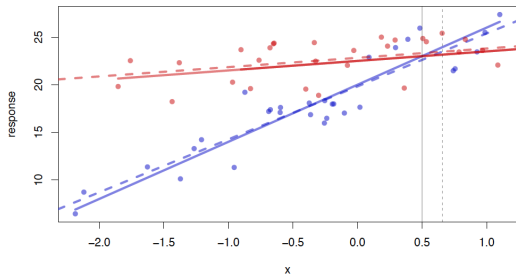
$$E[y_T] = \hat{\beta}_0 + \mathbf{x}^\top \hat{\beta}$$

$$E[y_C] = \hat{\alpha}_0 + \mathbf{x}^\top \hat{\alpha}$$

and classify new patients based on

$$\hat{\beta}_0 - \hat{\alpha}_0 + \mathbf{x}_i^\top (\hat{\beta} - \hat{\alpha})$$

Simple Example



Estimating Predictive Models

Can use complicated bells and whistles for prognostic models

$$E[y_T] = \hat{f}_T(x)$$

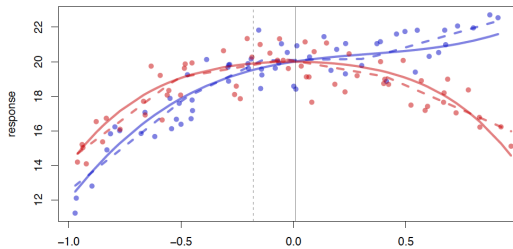
$$E[y_C] = \hat{f}_C(x)$$

and classify new patients based on

$$\hat{f}_T(x) - \hat{f}_C(x)$$

Using *Lasso* and *Basis expansions* keeps that difference simple (closed form).

Nonlinear example



Separate Models vs Interaction Model

In the linear (*non-lasso*) case, separate models is equivalent to:

$$y_i = \alpha_0 + \mathbf{x}_i^\top \alpha + \gamma_0 I_{\text{treatment}} + I_{\text{treatment}} \mathbf{x}_i^\top \gamma + \epsilon_i$$

Useful for testing individual regression terms

Vanilla Lasso will give different solutions to this vs separate fit.

We recommend thinking in terms of separate models (though one can use other tools, penalties, priors etc. . . to share strength)

Interaction Vs Qualitative Interaction

Qualitative Interaction:

An *interaction* (shape change) between $f_T(\cdot)$ and $f_C(\cdot)$ that indicates a crossing of the two curves.

This will inform treatment decisions.

A usual interaction (shape change) does not necessarily indicate crossing.

Without crossing, an interaction will not inform treatment decisions.

Predictive Biomarker for a Binary Measure

Randomize $n_T + n_C = n$ patients to **treatment/control**

On each patient measure

y_i - single binary outcome

(eg. Progression after a year, pCR)

\mathbf{x}_i - p -vector of features

(eg. SNPs, gene expression values)

Want to determine who will benefit from **treatment** over **control**.

Predictive Models from Prognostic Models

Given “true prognostic models”:

$$\text{logit}(P_T) = \beta_0 + \mathbf{x}^\top \beta$$

$$\text{logit}(P_C) = \alpha_0 + \mathbf{x}^\top \alpha$$

Would treat patients with

$$(P_T|x) - (P_C|x) > \delta$$

For $\delta = 0$ still simplifies to patients with

$$\beta_0 - \alpha_0 + \mathbf{x}^\top (\beta - \alpha) > 0$$

Predictive Biomarker for Survival Data

Randomize $n_T + n_C = n$ patients to **treatment/control**

On each patient measure

(t_i, s_i) - time, censoring-status

(eg. Disease free survival)

\mathbf{x}_i - p -vector of features

(eg. SNPs, gene expression values)

Predictive and Prognostic Models

A bit trickier; to make things comparable assume hazard is:

$$\lambda_T(t) = h(t)e^{\beta_0 + x^\top \beta}$$

$$\lambda_C(t) = h(t)e^{x^\top \alpha}$$

Need to fit a joint model due to baseline hazard.

Make decisions based on

$$\beta_0 + x^\top (\beta - \alpha)$$

Remember lower hazard is better