

Generalized Additive Model

In the statistical analysis of **clinical trials** and **observational studies**, the identification and adjustment for **prognostic factors** is an important component. Valid comparisons of different treatments requires the appropriate adjustment for relevant prognostic factors. The failure to consider important prognostic variables, particularly in observational studies, can lead to errors in estimating treatment differences. In addition, incorrect modeling of prognostic factors can result in the failure to identify nonlinear trends or threshold effects on survival.

This article describes flexible statistical methods that may be used to identify and characterize the effect of potential prognostic factors on an outcome variable. These methods are called “generalized additive models”, and extend the traditional **general linear model**. They can be applied in any setting in which a linear or **generalized linear model** is typically used. These settings include standard continuous response **regression**, **categorical** or **ordered categorical** response data, count data, **survival** data and **time series**.

One of the most commonly used statistical models in medical research is the **logistic regression** model for **binary data**. We use it here as a specific illustration of a generalized additive mode. Logistic regression (and many other techniques) model the effects of prognostic factors x_j in terms of a linear predictor of the form $\sum x_j \beta_j$, where the β_j are parameters. The generalized additive model replaces $\sum x_j \beta_j$ with $\sum f_j(x_j)$, where f_j is a unspecified (“nonparametric”) function. This function is estimated in a flexible manner using a scatterplot smoother (see **Graphical Displays**). The estimated function $\hat{f}_j(x_j)$ can reveal possible nonlinearities in the effect of x_j .

We first give some background on the methodology, and then discuss the details of the logistic regression model and its generalization. Some related developments are discussed in the last section.

Smoothing Methods and Generalized Additive Models

The building block of the generalized additive model

algorithm is the scatterplot smoother. We will first describe scatterplot smoothing in a simple setting, and then indicate how it is used in generalized additive modeling.

Suppose that we have a scatterplot of points (x_i, y_i) such as that shown in Figure 1. Here y is a **response** or outcome variable, and x is a prognostic factor. We wish to fit a smooth curve $f(x)$ that summarizes the dependence of y on x . If we were to find the curve that simply minimizes $\sum [y_i - f(x_i)]^2$, the result would be an interpolating curve that would not be smooth at all.

The cubic **spline** smoother imposes smoothness on $f(x)$. We seek the function $f(x)$ that minimizes

$$\sum [y_i - f(x_i)]^2 + \lambda \int f''(x)^2 dx. \quad (1)$$

Notice that $\int f''(x)^2$ measures the “wiggleness” of the function f : linear f s have $\int f''(x)^2 = 0$, while nonlinear f s produce values greater than zero. λ is a nonnegative smoothing parameter that must be chosen by the data analyst. It governs the tradeoff between the **goodness of fit** to the data (as measured by $\sum [y_i - f(x_i)]^2$) and wiggleness of the function. Larger values of λ force f to be smoother.

For any value of λ , the solution to (1) is a cubic spline; that is, a piecewise cubic polynomial with pieces joined at the unique observed values of x in the dataset. Fast and stable numerical procedures are available for computation of the fitted curve. The right panel of Figure 1 shows a cubic spline fit to the data.

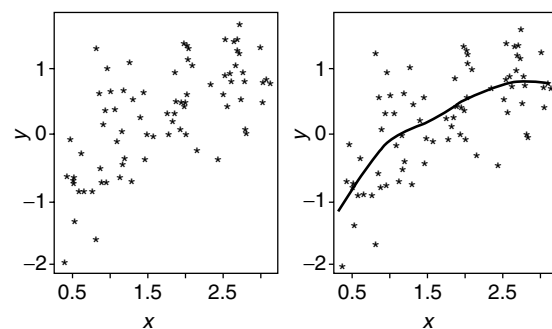


Figure 1 The left panel shows a fictitious scatterplot of an outcome measure y plotted against a prognostic factor x . In the right panel, a scatterplot smoother has been added to describe the trend of y on x

What value of λ did we use in Figure 1? In fact, it is not convenient to express the desired smoothness of f in terms of λ , as the meaning of λ depends on the units of the prognostic factor x . Instead, it is possible to define an “effective number of parameters” or “degrees of freedom” of a cubic spline smoother, and then use a numerical search to determine the value of λ to yield this number. In Figure 1 we chose the effective number of parameters to be 5. Roughly speaking, this means that the complexity of the curve is about the same as a **polynomial regression** of degree 4. However, the cubic spline smoother “spreads out” its parameters in a more even manner, and hence is much more flexible than a polynomial regression. Note that the degrees of freedom of a smoother need not be an integer.

The above discussion tells how to fit a curve to a single prognostic factor. With multiple prognostic factors, if x_{ij} denotes the value of the j th prognostic factor for the i th observation, we fit the additive model

$$\hat{y}_i \approx \sum_j f_j(x_{ij}). \quad (2)$$

A criterion such as (1) can be specified for this problem, and a simple iterative procedure exists for estimating the f_j s. We apply a cubic spline smoother to the outcome $y_i - \sum_{j \neq k} \hat{f}_j(x_{ij})$ as a function of x_{ik} , for each prognostic factor in turn. The process continues until the estimates \hat{f}_k stabilize. This procedure is known as “backfitting”, and the resulting fit is analogous to a multiple regression for linear models.

When generalized additive models are fit to binary response data (and in many other settings), the appropriate error criterion is a penalized log likelihood or a penalized log partial-likelihood (see **Penalized Maximum Likelihood**). To maximize it, the backfitting procedure is used in conjunction with a **maximum likelihood** or maximum **partial likelihood** algorithm. The usual Newton–Raphson routine (see **Optimization and Nonlinear Equations**) for maximizing log likelihoods in these models can be cast in an IRLS (iteratively reweighted least squares) form (see **Generalized Linear Model**). This involves a repeated weighted linear regression of a constructed response variable on the covariates: each regression yields a new value of the parameter estimates which give a new constructed variable, and the process is iterated. In the generalized additive model, the weighted

linear regression is simply replaced by a weighted backfitting algorithm. Details can be found in [7, Chapter 6].

The Generalized Additive Logistic Model

Generalized additive models can be used in virtually any setting in which linear models are used. The basic idea is to replace $\sum x_{ij}\beta_j$, the linear component of the model with an additive component $\sum f_j(x_{ij})$.

In the logistic regression model the outcome y_i is 0 or 1, with 1 indicating an event (such as death or relapse of a disease) and 0 indicating no event. We wish to model $p(y_i|x_{i1}, x_{i2}, \dots, x_{ip})$, the probability of an event given prognostic factors $x_{i1}, x_{i2}, \dots, x_{ip}$. The linear logistic model assumes that the log odds are linear:

$$\begin{aligned} \log \frac{p(y_i|x_{i1}, \dots, x_{ip})}{1 - p(y_i|x_{i1}, \dots, x_{ip})} \\ = \beta_0 + x_{i1}\beta_1 + \dots + x_{ip}\beta_p. \end{aligned} \quad (3)$$

The generalized additive logistic model assumes instead that

$$\begin{aligned} \log \frac{p(y_i|x_{i1}, \dots, x_{ip})}{1 - p(y_i|x_{i1}, \dots, x_{ip})} \\ = \beta_0 + f_1(x_{i1}) + \dots + f_p(x_{ip}). \end{aligned} \quad (4)$$

The functions f_1, f_2, \dots, f_p are estimated by an **algorithm** like the one described earlier.

To illustrate this, we describe a study on the survival of children after cardiac surgery for heart defects [13]. The data were collected during the period 1983–1988. A pre-operation warm-blood cardioplegia procedure, thought to improve chances for survival, was introduced in February 1988. This was not used on all of the children after February 1988, only on those for which it was thought appropriate and only by surgeons who chose to use the new procedure. The main question is whether the introduction of the warming procedure improved survival; the importance of risk factors age, weight, and diagnostic category is also of interest.

If the warming procedure was given in a randomized manner, we could simply focus on the post-February 1988 data and compare the survival of those who received the new procedure to those who did not. However, allocation was not random, so we can only

try to assess the effectiveness of the warming procedure as it was applied. For this analysis, we use all of the data (1983–1988). To adjust for changes that might have occurred over the five-year period, we include the data of the operation as a covariate. However, operation date is strongly confounded with the warming operation and thus a general nonparametric fit for date of operation might unduly remove some of the effect attributable to the warming procedure. To avoid this, we allow only a linear effect for operation date. Hence we must assume that any time trend is either a consistently increasing or decreasing trend.

We fit a generalized additive logistic model to the binary response death, with smooth terms for age and weight, a linear term for operation date, a categorical variable for diagnosis, and a binary variable for the warming operation. All the smooth terms are fitted with four degrees of freedom.

The resulting curves for age and weight are shown in Figure 2. As one would expect, the highest risk is for the lighter babies, with a decreasing risk over 3 kg. Somewhat surprisingly, there seems to be a low risk age around 200 days, with higher risk for younger and older children. Note that the numerical algorithm is not able to achieve exactly four degrees of freedom for the age and weight terms, but 3.80 and 3.86 degrees of freedom, respectively.

An analysis of deviance (*see Generalized Linear Model*) can be carried out for inference from a generalized additive model, analogous to that done for generalized linear models. The only new twist

is estimation of the degrees of freedom or effective number of parameters of the fitted model, which was discussed in the previous section. This analysis shows that the warming procedure is strongly beneficial to survival. There are strong differences in the diagnosis categories, while the estimated effect of operation date is not large.

Since a logistic regression is additive on the logit scale but not on the probability scale, a plot of the fitted probabilities is often informative. Figure 3 shows the fitted probabilities broken down by age and diagnosis, and is a concise summary of the findings of this study. The beneficial effect of the treatment at the lower weights is evident. As with all nonrandomized studies, the results here should be interpreted with caution. In particular, one must insure that the children were not chosen for the warming operation based on their prognosis. To investigate this, we perform a second analysis in which a **dummy variable** (say, period), corresponding to before vs. after February 1988, is inserted in place of the dummy variable for the warming operation. The purpose of this is to investigate whether the overall treatment strategy improved after February 1988. If this turns out not to be the case, it will imply that warming was used only for patients with a good prognosis, who would have survived anyway. A linear adjustment for operation date is included as before. The results are qualitatively very similar to the first analysis: age and weight are significant, with effects similar to those in Figure 2; diagnosis is significant, while operation date (linear effect) is not. Period is highly significant. Hence there seems to be a significant overall improvement in survival after February 1988. For more details, see [13].

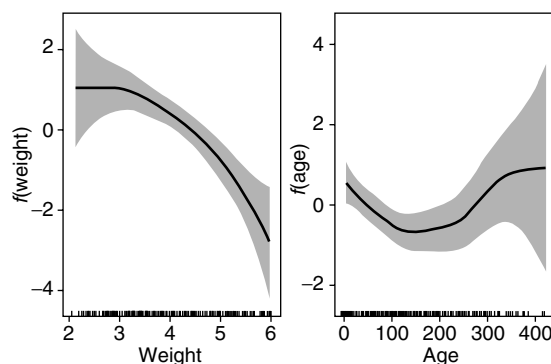


Figure 2 Estimated functions for weight and age for warm cardioplegia data. The shaded region represents twice the pointwise asymptotic standard errors of the estimated curve

Discussion

The nonlinear modeling procedures described here are useful for two reasons. First, they help to prevent model misspecification, which can lead to incorrect conclusions regarding treatment efficacy. Secondly, they provide information about the relationship between prognostic factors and disease risk that is not revealed by the use of standard modeling techniques. Linearity always remains a special case, and thus simple linear relationships can be easily confirmed with flexible modeling of covariate effects.

The most comprehensive source for generalized additive models is [7], from which the example was

4 Generalized Additive Model

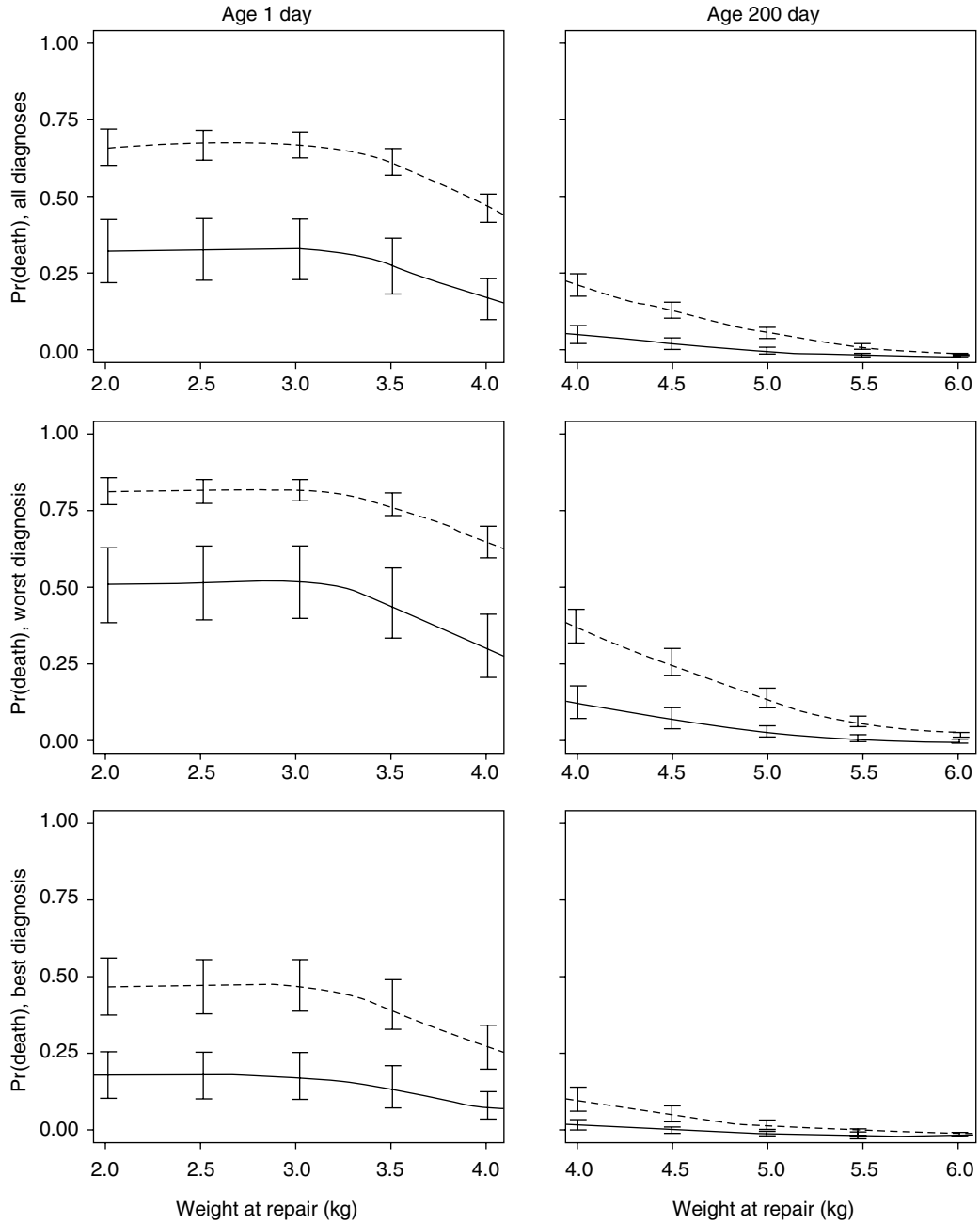


Figure 3 Estimated probabilities for warm cardioplegia data, conditioned on two ages (columns) and three diagnostic classes (rows). The broken line is standard treatment; the solid line is warm cardioplegia. Bars indicate ± 1

taken. A detailed example of the use of generalized additive models in the **proportional hazards** setting is given in [10]. Other medical applications are discussed in [6] and [9]. Penalization and spline models in a variety of settings are discussed in [5], and [12] is a good source for the mathematical background of spline models. See also [3] for an exposition of modern developments in statistics (including generalized additive models), for a nonmathematical audience.

There has been some recent related work in this area. A different method for flexible hazard modeling is described in [11] and a generalization of additive modeling that finds **interactions** among prognostic factors is proposed in [4]. Of particular interest in the **proportional hazards** setting is the *varying coefficient* model [8] (see **Semiparametric Regression**), in which the parameter effects can change with other factors such as time. The model has the form

$$h(t|x_{i1}, \dots, x_{ip}) = h_0(t) \exp \sum_{j=1}^p \beta_j(t)x_{ij}. \quad (5)$$

The parameter functions $\beta_j(t)$ are estimated by scatterplot smoothers in a similar fashion to the methods described earlier. This gives a useful way of modeling departures from the proportional hazards assumption by estimating the way in which the parameters β_j change with time.

Software for fitting generalized additive models is available in the **S/SPLUS** statistical environment [1, 2], in a FORTRAN program called gamfit available at statlib (in general/gamfit at the ftp site lib.stat.cmu.edu) and also in the GAIM package for MS-DOS computers, available from the authors.

Acknowledgment

Trevor Hastie was partially supported by grant DMS-9504495 from the National Science Foundation, and grant ROI-CA-72028-01 from the National Institutes of Health.

Robert Tibshirani was partially supported by a grant from the Natural Sciences and Engineering Research Council of Canada.

References

- [1] Becker, R., Chambers, J. & Wilks, A. (1988). *The New S Language*. Wadsworth International Group, Pacific Grove.
- [2] Chambers, J. & Hastie, T. (1991). *Statistical Models in S*. Wadsworth/Brooks Cole, Pacific Grove.
- [3] Efron, B. & Tibshirani, R. (1991). Statistical analysis in the computer age. *Science* **253**, 390–395.
- [4] Friedman, J. (1991). Multivariate adaptive regression splines (with discussion), *Annals of Statistics* **19**, 1–141.
- [5] Green, P. & Silverman, B. (1994). *Nonparametric Regression and Generalized Linear Models: a Roughness Penalty Approach*. Chapman & Hall, London.
- [6] Hastie, T. & Herman, A. (1990). An analysis of gestational age, neonatal size and neonatal death using nonparametric logistic regression, *Journal of Clinical Epidemiology* **43**, 1179–1190.
- [7] Hastie, T. & Tibshirani, R. (1990). *Generalized Additive Models*. Chapman & Hall, London.
- [8] Hastie, T. & Tibshirani, R. (1997). Discriminant analysis by Gaussian mixtures, to appear.
- [9] Hastie, T., Botha, J. & Schnitzler, C. (1989). Regression with an ordered categorical response, *Statistics in Medicine* **8**, 785–794.
- [10] Hastie, T., Sleeper, L. & Tibshirani, R. (1992). Flexible covariate effects in the proportional hazards model, *Breast Cancer Research and Treatment* **22**, 241–250.
- [11] Kooperberg, C., Stone, C. & Truong, Y. (1993). Hazard Regression, *Technical Report*. Department of Statistics, University of California, Berkeley.
- [12] Wahba, G. (1990). *Spline Models for Observational Data*. SIAM, Philadelphia.
- [13] Williams, W., Rebeyka, I., Tibshirani, R., Coles, J., Lightfoot, N., Freedom, R. & Trusler, G. (1990). Warm induction cardioplegia in the infant: a technique to avoid rapid cooling myocardial contracture, *Journal of Thoracic and Cardiovascular Surgery* **100**, 896–901.

(See also **Nonparametric Regression**)

TREVOR HASTIE & R. TIBSHIRANI