

## USING SPECIALLY DESIGNED EXPONENTIAL FAMILIES FOR DENSITY ESTIMATION

BY BRADLEY EFRON<sup>1</sup> AND ROBERT TIBSHIRANI

*Stanford University*

We wish to estimate the probability density  $g(y)$  that produced an observed random sample of vectors  $y_1, y_2, \dots, y_n$ . Estimates of  $g(y)$  are traditionally constructed in two quite different ways: by maximum likelihood fitting within some parametric family such as the normal or by nonparametric methods such as kernel density estimation. These two methods can be combined by putting an exponential family “through” a kernel estimator. These are the specially designed exponential families mentioned in the title. Poisson regression methods play a major role in calculations concerning such families.

**1. Introduction.** Suppose that we wish to estimate the probability density  $g(y)$  that produced an observed random sample of vectors  $y_1, y_2, \dots, y_n$ ,

$$(1.1) \quad y_i \stackrel{\text{i.i.d.}}{\sim} g(y) \quad \text{for } i = 1, 2, \dots, n.$$

The vectors  $y_i$  take values in a sample space  $\mathcal{Y}$ . The numerical examples in this paper have  $\mathcal{Y}$  being portions of the real line or of the plane, but the methodology applies just as well to higher dimensionalities and to more complicated spaces.

Estimates of  $g(y)$  are traditionally constructed in two quite different ways: by maximum likelihood fitting within some parametric family such as the normal or by nonparametric methods such as kernel density estimation. These two methods can be combined by putting an exponential family “through” a nonparametric estimator. The resulting hybrid estimators are the specially designed exponential families of the title.

Figure 1 shows a simple example of this methodology. The  $y_i$  are pain scores for  $n = 67$  women, each obtained by averaging the results from a questionnaire administered after an operation. The scale runs from  $y = 0 =$  no pain to  $y = 4 =$  worst pain, so the sample space  $\mathcal{Y}$  is the interval  $[0, 4]$ . The 67 scores  $y_i$ , indicated by the histogram, run from 0.02 to 3.08. The dashed curve  $\hat{g}_0(y)$  is a normal kernel density estimator with window width  $\lambda = 1$ , described more carefully in Section 2. Also shown are two special exponential family estimates,  $\hat{g}_1(y)$  and  $\hat{g}_2(y)$ , described below.

---

Received January 1995; revised October 1995.

<sup>1</sup>Supported by NSF Grant DMS-95-04379 and Public Health Service Grant 5 ROI CA59039-20. AMS 1991 subject classifications. 62F05, 62G05.

*Key words and phrases.* Poisson regression, degrees of freedom, expected deviance, local and global smoothing, moment-matching.

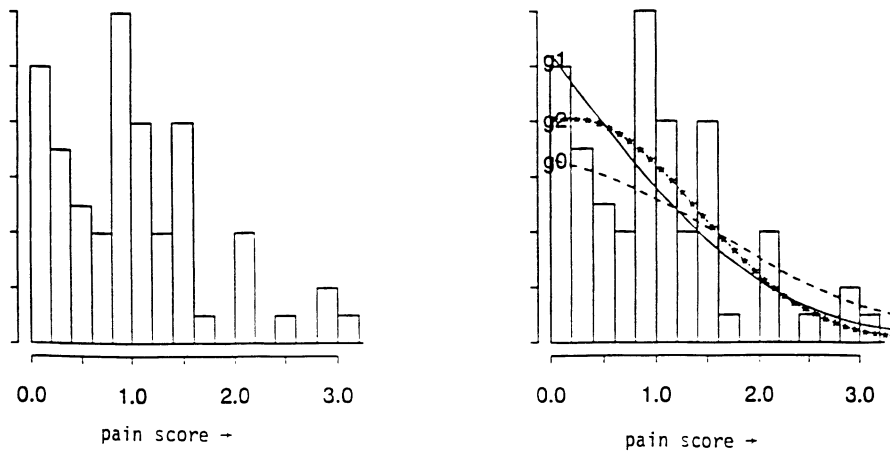


FIG. 1. Pain-score data. Left: Histogram of pain scores for 67 women following an operation: 0 = no pain; 4 = worst pain. Right:  $\hat{g}_0(y)$  is normal kernel density estimator, window width 1;  $\hat{g}_1(y)$  is the special exponential family through  $\hat{g}_0(y)$  with sufficient statistic  $t(y) = y$ ;  $\hat{g}_2(y)$  uses sufficient statistics  $t(y) = (y, y^2)$ . Density  $\hat{g}_2(y)$  matches the empirical mean and variance of the 67 data points.

An exponential family of densities on  $\mathcal{Y}$ ,  $\{g_\beta(y)\}$  is given by

$$(1.2) \quad g_\beta(y) = g_0(y) \exp(\beta_0 + t(y)\beta_1),$$

which is called the exponential family through  $g_0(y)$  with sufficient statistics  $t(y)$ . Here  $g_0(y)$  is a carrier density,  $t(y)$  is a  $1 \times p$  vector of sufficient statistics,  $\beta_1$  is a  $p \times 1$  parameter vector and  $\beta_0$  is a normalizing parameter that makes  $g_\beta(y)$  integrate to 1 over  $\mathcal{Y}$ . For example, the one-dimensional normal family, with all possible choices of expectation and variance, can be obtained using the standard normal carrier  $g_0(y) = \varphi(y) = \exp\{-0.5y^2\}/\sqrt{2\pi}$ , with the sufficient statistics  $t(y) = (y, y^2)$ . The densities in (1.2) are defined with respect to some background measure, which we will take to be Lebesgue measure. See Section 1.4 of Lehmann (1983).

The estimates  $\hat{g}_1(y)$  and  $\hat{g}_2(y)$  in Figure 1 are of the form

$$(1.3) \quad g_{\hat{\beta}}(y) = \hat{g}_0(y) \exp(\hat{\beta}_0 + t(y)\hat{\beta}_1),$$

where  $\hat{g}_0(y)$  is the kernel density estimate indicated by the dashed curve. Estimate  $\hat{g}_1(y)$  uses the single sufficient statistic  $t(y) = y$ , so  $p = 1$ , while  $\hat{g}_2(y)$  uses  $t(y) = (y, y^2)$ ,  $p = 2$ . The parameter values  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$  were chosen by maximum likelihood, that is, by maximizing  $\prod_{i=1}^n g_{\hat{\beta}}(y_i)$ , ignoring the fact that the carrier  $\hat{g}_0(y)$  is itself data-dependent. This choice of  $\hat{\beta}$  matches the  $t(y)$  moments of  $g_{\hat{\beta}}(y)$  to their empirical averages:

$$(1.4) \quad \int_{\mathcal{Y}} t(y) g_{\hat{\beta}}(y) dy = \frac{1}{n} \sum_{i=1}^n t(y_i).$$

Thus  $\hat{g}_2$  matches the first two empirical moments of the 67 pain scores or, equivalently, it matches the empirical mean and variance. Property (1.4) implies that the special exponential family estimate of  $g(y)$  is unbiased for the moments of  $t(y)$ . Linear transformations to “post-repair” the mean and variance of a kernel estimate are familiar in the density-estimation literature; see, for example, Jones (1993). The two-dimensional example of Section 4 shows a more ambitious example of moment-matching.

We can think of a special exponential family estimator in two complementary ways: (1) as being a standard exponential family estimator, except one that is preceded by an adaptive choice of the carrier, or (2) as being a standard nonparametric smoothing estimator, except one that is followed by a correction to match certain sample moments.

Either way, we will argue that special exponential families can be a favorable compromise between parametric and nonparametric density estimation. From the first point of view, we will be able to use exponential family theory more flexibly than when restricted to the usual small catalogue of normals, gammas, betas and so forth. An approach more in the spirit of a regression analysis than a density estimate is possible, including exploratory choices of the sufficient statistics  $t(y)$ .

From the second point of view, the moment-matching correction will usually reduce the bias of a nonparametric smoother, since the moments  $t(y)$  are estimated unbiasedly. This has an important practical consequence shown in the numerical examples: *it allows the nonparametric smoother to use a substantially greater window width without badly degrading the overall fit to the data.* The result is a substantially smoother estimate of the density  $g(y)$ . (This phenomenon is illustrated in the bivariate example of Figures 3 and 4.)

To put things another way, the special exponential family estimate (1.3) works at two different scales. The nonparametric smoother allows local adaptation to the data, while the exponential term matches some of the data’s global properties.

Sections 2–5 describe how to compute and interpret special exponential family (SEF) estimates such as those in Figure 1. Most of our computations are done using a Poisson regression model for density estimation (Section 2), originally introduced in Lindsey (1974a, b). Section 3 gives a delta-method formula for the covariance of  $\hat{\beta}_1$  in (1.3), which takes into account the data-based choice of the carrier  $\hat{g}_0(y)$ . We can use this formula in the usual way to select among possible choices of the sufficient statistic. Section 6 concerns formulas for choosing the window width of the smoother that produces  $\hat{g}_0(y)$ . This choice involves “degrees of freedom” calculations like those introduced by Hastie and Tibshirani (1990). Multisample SEF estimates are discussed in Section 7. Remarks appear in Section 9.

Special exponential families are an example of what Green and Silverman (1994) call *semiparametric methods*. Many other semiparametric methods have been proposed for density estimation. Hjort and Glad (1995) propose reversing the SEF order: first fit a parametric family to the data and then fit a nonparametric smoother to the residuals from the parametric estimator.



with sum  $\pi_+(\theta) = 1$ . Then  $\mathbf{s}$  has a multinomial distribution on  $K$  categories, with  $n$  draws and probability vector  $\boldsymbol{\pi}(\theta) = (\pi_1(\theta), \pi_2(\theta), \dots, \pi_K(\theta))$ ,

$$(2.4) \quad \mathbf{s} \sim \text{Mult}_K(n, \boldsymbol{\pi}(\theta)).$$

We could find the maximum likelihood estimate (MLE) of  $\theta$ , based on  $\mathbf{s}$ , by maximizing the multinomial probability of  $\mathbf{s}$ .

Instead we consider the  $s_k$  to be *independent* Poisson observations

$$(2.5) \quad s_k \overset{\text{ind}}{\sim} \text{Po}(\mu_k(\gamma, \theta)), \quad k = 1, 2, \dots, K,$$

with expectations

$$(2.6) \quad \mu_k(\gamma, \theta) = \gamma \pi_k(\theta).$$

Here  $\gamma$  is a free parameter, restricted only to be positive. Standard Poisson properties allow (2.5) and (2.6) to be expressed as

$$(2.7) \quad s_+ \sim \text{Po}(\gamma) \quad \text{and} \quad s|s_+ \sim \text{Mult}_K(s_+, \boldsymbol{\pi}(\theta)).$$

This means that the maximum likelihood estimates from (2.5) and (2.6) are

$$(2.8) \quad \hat{\gamma} = s_+ = n,$$

and  $\hat{\theta}$  is equal to the MLE for  $\theta$  in (2.4). Lindsey's method is to (approximately) maximize the original likelihood  $\prod_{i=1}^n g_\theta(y_i)$  by finding the Poisson MLE in (2.5) and (2.6). *Note:* The parameters  $\gamma$  and  $\theta$  are orthogonal in Cox and Reid's (1987) sense, so that the information for estimating  $\theta$  is the same in (2.4) and (2.5).

This method of finding the MLE is particularly convenient when the original densities are of the exponential family form (1.2). We will consider the density estimation problem (1.1) and (1.2) in the Poisson regression form (2.5) and (2.6):

$$(2.9a) \quad s_k \overset{\text{ind}}{\sim} \text{Po}(\mu_k(\boldsymbol{\beta})) \quad \text{for } k = 1, 2, \dots, K,$$

with

$$(2.9b) \quad \mu_k(\boldsymbol{\beta}) = \mu_k^o \exp(\beta_0 + t_k \beta_1).$$

Here  $\mu_k^o$  is proportional to  $\pi_k^o = \int_{\mathcal{Y}_k} g_0(y) dy$ , a discretized version of the carrier, and

$$(2.10) \quad t_k = t(y_{(k)}),$$

the sufficient vector  $t(y)$  evaluated at a convenient point  $y_{(k)}$  in  $\mathcal{Y}_k$ . The free parameter  $\beta_0$  corresponds to  $\gamma = e^{\beta_0}$  in (2.6).

Define  $X$  to be the  $K \times (p + 1)$  matrix whose  $k$ th row equals

$$(2.11) \quad x_k = (1, t_k).$$

The maximum likelihood equations for  $\boldsymbol{\beta} = (\beta_0, \beta_1)$  in the generalized linear model (2.9) are

$$(2.12) \quad X'[\mathbf{s} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})] = 0,$$

where  $\boldsymbol{\mu}(\hat{\beta})$  indicates the vector with  $k$ th component  $\mu_k^o \exp(x_k \hat{\beta})$ . Standard generalized linear model software easily solves for  $\hat{\beta}$  in (2.12), even for difficult nonstandard forms of the exponential family (1.2). This was Lindsey's principal point.

Here is how Figure 1 was constructed. The problem was discretized into  $K = 40$  cells as in Table 1. The carrier vector  $\boldsymbol{\mu}^o = (\mu_1^o, \mu_2^o, \dots, \mu_K^o)$  in (2.9b) was estimated using a  $K \times K$  smoothing matrix  $M(\lambda)$ ,

$$(2.13) \quad \hat{\boldsymbol{\mu}}^o = M(\lambda)\mathbf{s}.$$

Matrix  $M(\lambda)$  was taken to be a normal kernel smoother, having  $kj$ th element

$$(2.14) \quad M_{kj}(\lambda) = \frac{c_k}{\lambda} \varphi\left(\frac{y_{(k)} - y_{(j)}}{\lambda}\right),$$

with  $y_{(k)} = (k - 0.5)/10$ , the midpoint of cell  $\mathcal{Y}_k$ . The constants  $c_k$  were chosen to make  $M_{k+} = 1$ . The starred curve labelled  $\hat{g}_0$  in Figure 1 is actually  $\hat{\mu}_k^o$  plotted as a function of  $y_{(k)}$ , with the window width  $\lambda$  in (2.14) set equal to 1.

The curves labelled  $\hat{g}_1(y)$  and  $\hat{g}_2(y)$  are really the discrete analogs of the special exponential family (1.3), say  $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K)$ ,

$$(2.15) \quad \hat{\mu}_k = \hat{\mu}_k^o \exp(\hat{\beta}_0 + t_k \hat{\beta}_1),$$

plotted versus  $y_{(k)}$ :  $\hat{g}_1$  uses  $t_k = y_{(k)}$ , while  $\hat{g}_2$  is based on the quadratic vector  $t_k = (y_{(k)}, y_{(k)}^2)$ . The MLE estimates  $\hat{\beta}$  were obtained by iterative solution of (2.12), so

$$(2.16) \quad X'[\mathbf{s} - \hat{\boldsymbol{\mu}}] = 0,$$

where  $X$  is a  $40 \times 2$  matrix for  $\hat{g}_1$  and a  $40 \times 3$  matrix for  $\hat{g}_2$ . [The estimates were actually computed using a centered version of  $y$ ,  $\tilde{y}_{(k)} = (y_{(k)} - 2)/4$ .] To fit this model in the GLIM language or the `glm` function in SPlus, one simply includes  $\log \hat{\mu}^o$  as an offset in a Poisson generalized linear model.

Equation (2.16) shows that

$$(2.17a) \quad \hat{\mu}_+ = s_+ = n$$

and

$$(2.17b) \quad \sum_{k=1}^K \frac{\hat{\mu}_k}{n} t_k = \sum_{k=1}^K \frac{s_k}{n} t_k,$$

the discrete analog of the moment-matching property (1.4). Notice that because of (2.17a) and the fact that the cells are of length 0.1, the curves in Figure 1 integrate over  $\mathcal{Y}$  to 6.7 rather than to 1.

Changing  $K$  to 20 or to 80 made very little difference in Figure 1. The numerical calculations in this paper were insensitive to the form of discretization. In fact, discretization is not really necessary for any of our results, as discussed in Remark E of Section 9.

Nevertheless, it is conceptually easier to discuss special exponential family density estimation in terms of the discrete Poisson model (2.9). The problem becomes one of fitting a smooth regression curve to the independent observations  $s_k$ , and this lets us make use of the arsenal of regression tools. It also emphasizes the important point that density estimation is equivalent to Poisson regression, and not to ordinary least squares regression. The Poisson nature of the problem will be evident in the formulas developed below.

**3. Estimating the covariance of  $\hat{\beta}$ .** This section derives an approximate covariance matrix for  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ , the estimated parameters in a special exponential family such as (1.3) or (2.15). The formula for the covariance takes into account the data-based choice of the carrier. We will use the estimated standard errors of the components of  $\hat{\beta}$  for model building, checking the significance of the corresponding components of the sufficient statistic  $t(y)$  in the usual way.

We consider the Poisson form (2.9) of the SEF model,

$$(3.1) \quad \mathbf{s} \sim \text{Po}_K(\boldsymbol{\mu}(\boldsymbol{\beta})) \quad \text{with } \boldsymbol{\mu}(\boldsymbol{\beta}) = \boldsymbol{\mu}^\circ e^{X\boldsymbol{\beta}},$$

where  $\text{Po}_K(\boldsymbol{\mu})$  indicates a vector of  $K$  independent Poisson variates having expectations  $(\mu_1, \mu_2, \dots, \mu_K) = \boldsymbol{\mu}$ . The notation  $\boldsymbol{\mu}^\circ e^{X\boldsymbol{\beta}}$  indicates the vector with  $k$ th component  $\mu_k^\circ e^{x_k \beta}$ ,  $x_k = (1, t_k)$ , as in (2.11). Generalizing (2.13), we first estimate  $\boldsymbol{\mu}^\circ$  by some function of  $\mathbf{s}$ , say  $m(\mathbf{s})$ , and then solve for  $\hat{\beta}$  in the MLE equations (2.16):

$$(3.2a) \quad \hat{\boldsymbol{\mu}}^\circ = m(\mathbf{s}) \quad \text{and} \quad \hat{\beta}: X'[\mathbf{s} - \hat{\boldsymbol{\mu}}^\circ e^{X\hat{\beta}}] = 0.$$

The special exponential family estimate of  $\boldsymbol{\mu}$  is

$$(3.2b) \quad \hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}^\circ e^{X\hat{\beta}}.$$

LEMMA 1. Let  $\hat{D}$  be the  $K \times K$  diagonal matrix with  $k$ th diagonal element  $\hat{\mu}_k = \hat{\mu}_k^\circ e^{x_k \hat{\beta}}$  and let  $\hat{H}$  be the  $K \times K$  derivative matrix of  $\log(\hat{\boldsymbol{\mu}}^\circ) = (\log(\hat{\mu}_1^\circ), \log(\hat{\mu}_2^\circ), \dots, \log(\hat{\mu}_K^\circ))$  with respect to  $\mathbf{s}$ , with “ $j$ ” indexing columns,

$$(3.3) \quad \hat{H} = \frac{d \log \hat{\boldsymbol{\mu}}^\circ}{d\mathbf{s}} = \left( \frac{\partial \log(\hat{\mu}_k^\circ)}{\partial s_j} \right).$$

Then the  $(p + 1) \times K$  derivative matrix of  $\hat{\beta}$  with respect to  $\mathbf{s}$  is

$$(3.4a) \quad \frac{d\hat{\beta}}{d\mathbf{s}} = [X'\hat{D}X]^{-1}Z',$$

where

$$(3.4b) \quad Z' = X'(I - \hat{D}\hat{H}).$$

In case (2.13),  $\hat{\boldsymbol{\mu}}^\circ = M\mathbf{s}$ , this becomes

$$(3.4c) \quad Z' = X'(I - D(e^{X\hat{\beta}})M),$$

where  $D(e^{X\hat{\beta}})$  is the diagonal matrix with  $k$ th diagonal element  $e^{x_k \hat{\beta}}$ . (The proof appears below.)

By the usual delta-method argument, an approximate covariance matrix estimate for  $\hat{\beta}$  is given by

$$(3.5) \quad \left( \frac{d\hat{\beta}}{d\mathbf{s}} \right) \widehat{\text{Cov}}(\mathbf{s}) \left( \frac{d\hat{\beta}}{d\mathbf{s}} \right)'$$

COROLLARY 1. *The SEF vector  $\hat{\beta}$  obtained from (3.1) and (3.2) has approximate covariance matrix*

$$(3.6a) \quad \widehat{\text{Cov}}(\hat{\beta}) = [X' \hat{D}X]^{-1} [Z' \bar{D}Z] [X' \hat{D}X]^{-1},$$

where  $\bar{D}$  is the diagonal matrix with  $k$ th diagonal element  $s_k$ . An alternative estimate is

$$(3.6b) \quad \widehat{\text{Cov}}(\hat{\beta}) = [X' \hat{D}X]^{-1} [Z' \hat{D}Z] [X' \hat{D}X]^{-1}.$$

For any  $K$ -vector  $\mathbf{v}$ , we let  $D(\mathbf{v})$  be the  $K \times K$  diagonal matrix with  $k$ th diagonal element  $v_k$ . The true covariance of  $\mathbf{s} \sim \text{Po}_k(\boldsymbol{\mu})$  is  $D(\boldsymbol{\mu})$ . Approximation (3.6a) estimates  $\text{Cov}(\mathbf{s})$  by  $D(\mathbf{s}) = \bar{D}$  in (3.5), while (3.6b) uses  $D(\hat{\boldsymbol{\mu}}) = \hat{D}$ . The former may be preferred if model (3.1) is suspect. In our numerical examples the two formulas gave nearly the same results.

Table 2 applies the corollary to the pain-score data. The quadratic model, described in (2.13)–(2.16) has  $\hat{\beta}_1 = (-2.74, -3.80)$ , with standard errors (0.93, 2.45) according to (3.6a). [ $\beta_0$  serves only to normalize  $\hat{\boldsymbol{\mu}}$  to  $\hat{\mu}_+ = n$ , (2.17a), so its value is of no statistical interest.] The coefficient of  $\tilde{y} = (y - 2)/4$ , the centered version of  $y$ , is nearly 3 standard errors below zero. The coefficient of  $\tilde{y}^2$  is  $-1.55$  standard errors below zero, so it is not so clear that the quadratic term is significant. The cubic model, in which  $t(y) = (\tilde{y}, \tilde{y}^2, \tilde{y}^3)$ , has the cubic coefficient only 0.07 standard errors above zero. Either the linear or the quadratic model seem reasonable here; the cubic model is definitely excessive. Section 6 discusses model selection in more detail.

TABLE 2

Parameter estimates  $\hat{\beta}$  and estimated standard errors for the pain-score data of Table 1. The quadratic model is described in (2.13)–(2.16);  $\tilde{y} = (y - 2)/4$ , centered version of  $y$ ;  $\overline{\text{se}}$  square root of diagonal elements (3.6a);  $\widehat{\text{se}}$  from (3.6b); “jack” is jackknife standard error; “naive” is the usual exponential family sterr estimated ignoring data-based choice of carrier. The cubic model uses the same  $M$ , but  $t(y) = (\tilde{y}, \tilde{y}^2, \tilde{y}^3)$ ; the cubic term is not at all significant

|               | Quadratic model |                        |         |                       |      | Cubic model |                 |                        |
|---------------|-----------------|------------------------|---------|-----------------------|------|-------------|-----------------|------------------------|
|               | $\hat{\beta}_1$ | $\overline{\text{se}}$ | (ratio) | $\widehat{\text{se}}$ | jack | naive       | $\hat{\beta}_1$ | $\overline{\text{se}}$ |
| $\tilde{y}$   | -2.74           | 0.93                   | (-2.95) | 0.96                  | 1.07 | 1.18        | -2.78           | 1.24                   |
| $\tilde{y}^2$ | -3.80           | 2.45                   | (-1.55) | 2.50                  | 2.66 | 2.85        | -3.59           | 2.70                   |
| $\tilde{y}^3$ |                 |                        |         |                       |      |             | 0.59            | 8.30                   |



The jackknife standard errors for the components of  $\hat{\beta}_1$  [formula (11.5) of Efron and Tibshirani (1993)] were computed as a check on the corollary. They came out a little larger than the delta-method estimates from (3.6a) or (3.6b), as is often the case; see Section 2 of Efron (1992).

Suppose that we ignore the fact that the carrier  $\hat{\mu}^o$  in (3.2) is a function of the data  $\mathbf{s}$ . This amounts to taking  $\hat{H} = 0$  in (3.3), so  $Z' = X'$  in (3.4b). Then (3.6b) reduces to the usual covariance estimate for exponential families,

$$(3.7) \quad \widehat{\text{Cov}}(\hat{\beta}) = (X' \hat{D} X)^{-1},$$

while (3.6a) becomes the “sandwich” estimate  $(X' \hat{D} X)^{-1} (X' \bar{D} X) (X' \hat{D} X)^{-1}$ . The standard error estimates from (3.7), labelled “naive” in Table 2, are considerably larger than the standard errors that take into account the adaptive choice of the carrier.

The naive standard errors will usually exceed those from (3.6). The reason is that the adaptive choice of  $\hat{\mu}^o$  absorbs some of the variability in  $\hat{\beta}$ . Here is a simple normal-theory version of the same phenomenon: suppose we observe  $z \sim N(\mu^o + \beta, 1)$  and we wish to estimate  $\beta$ , the distance of the expectation  $\mu = \mu^o + \beta$  from some origin of measurement  $\mu^o$ . Then  $\hat{\beta} = z - \mu^o$  has standard error 1. However, if the origin is chosen adaptively, say by  $\hat{\mu}^o = 0.75 \cdot z$ , then  $\hat{\beta} = z - \hat{\mu}^o$  has standard error 0.25. Observing  $z = 4$ , for example, gives  $\hat{\beta} = 1$ , which in the adaptive case is 4 standard errors above 0. See Remark B in Section 9.

PROOF OF LEMMA 1. Using the  $D$  notation for diagonal matrices, (3.2) for  $\hat{\beta}$  can be expressed as

$$(3.8) \quad X' [\mathbf{s} - D(\hat{\mu}^o) e^{X \hat{\beta}}] = 0,$$

where  $e^{X \hat{\beta}}$  is the vector with components  $e^{x_k \hat{\beta}}$ . A small change  $d\mathbf{s}$  in  $\mathbf{s}$  produces change  $d\hat{\beta}$  in the MLE vector and change

$$(3.9) \quad d\hat{\mu}^o \doteq D(\hat{\mu}^o) \hat{H} d\mathbf{s}$$

in  $\hat{\mu}^o$ . Then (3.8) gives

$$(3.10) \quad \begin{aligned} \mathbf{0} &= X' [\mathbf{s} + d\mathbf{s} - D(\hat{\mu}^o + d\hat{\mu}^o) \exp(X(\hat{\beta} + d\hat{\beta}))] \\ &\doteq X' [\mathbf{s} + D(\hat{\mu}^o) \exp(X \hat{\beta}) + d\mathbf{s} - D(d\hat{\mu}^o) \exp(X \hat{\beta}) \\ &\quad - D(\hat{\mu}^o \exp(X \hat{\beta})) X d\hat{\beta}]. \end{aligned}$$

Using (3.8), (3.9) and the fact that  $D(d\hat{\mu}^o) e^{X \hat{\beta}} = D(e^{X \hat{\beta}}) d\hat{\mu}^o$ , (3.10) becomes

$$(3.11) \quad \begin{aligned} 0 &= X' d\mathbf{s} - X' D(\hat{\mu}^o e^{X \hat{\beta}}) \hat{H} d\mathbf{s} - X' D(\hat{\mu}^o e^{X \hat{\beta}}) X d\hat{\beta} \\ &= X' [I - \hat{D} \hat{H}] d\mathbf{s} - X' \hat{D} X d\hat{\beta}. \end{aligned}$$

This verifies (3.4).  $\square$

**4. A bivariate example.** Density estimation becomes more interesting, and more difficult, when the sample space  $\mathcal{Y}$  is of higher dimension. This section introduces an example when  $\mathcal{Y}$  is a portion of the plane. We will use this example to illustrate some of the advantages of special exponential family density estimation.

Figure 2 shows  $l = \log(\text{redshift})$  and  $m$  equal to the apparent magnitude for 486 galaxies taken from Loh and Spillar's (1988) redshift survey,

$$(4.1) \quad y_i = (l_i, m_i) \quad \text{for } i = 1, 2, 3, \dots, n = 486.$$

Hubble's law, that larger redshift implies greater distance from Earth, is apparent in the figure. The galaxies with larger values of  $l$  tend to appear dimmer, that is, to have larger apparent magnitudes, leaving the lower right corner nearly empty.

The data in Figure 2 are a truncated subsample of the 879 galaxies in the Loh–Spillar catalog. It is all of the catalog entries falling into the rectangle

$$(4.2) \quad \log(0.2) \leq l \leq \log(1.2) \quad \text{and} \quad 17.2 \leq m \leq 21.5.$$

We will take this rectangle to be the sample space  $\mathcal{Y}$ . Some of the scientific reasons for truncation are discussed in Efron and Petrosian (1992).

Figure 2 shows  $\mathcal{Y}$  partitioned into  $K = 285$  rectangular cells  $\mathcal{Y}_k$ , by dividing the  $l$  axis into 15 equal strips and dividing the  $m$  axis into 19 equal strips. The corresponding counts  $s_k$ , (3.2), are shown on the right side of the figure. It is convenient to index the cells by

$$(4.3) \quad k = (i, j), \quad i = 1, 2, \dots, 15, \quad j = 1, 2, \dots, 19.$$

The midpoint of rectangle  $\mathcal{Y}_k$  is  $y_{(k)} \equiv (l_{(i)}, m_{(j)})$ .

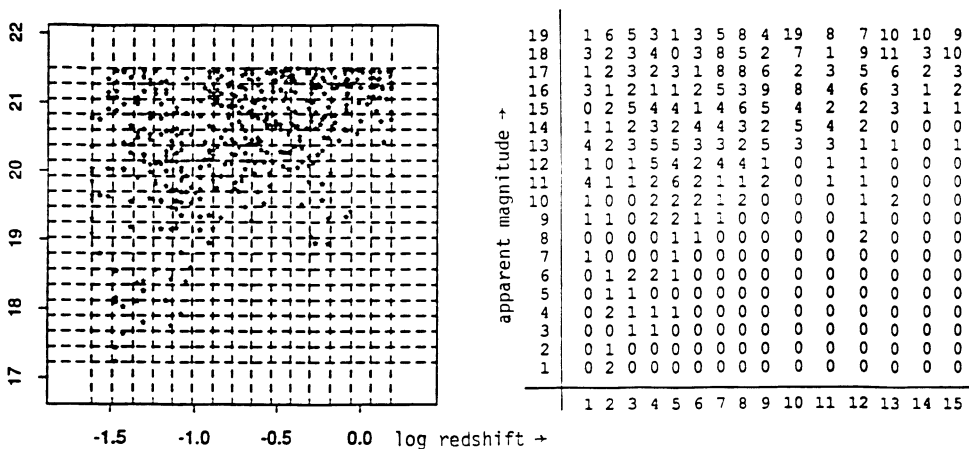


FIG. 2. The galaxy data: 486 galaxies from Loh and Spillar's (1988) redshift survey ( $\log$  redshift  $l$  and apparent magnitude  $m$ ) discretized into  $285 = 15 \times 19$  equal cells. The counts  $\mathbf{s}$  for the 285 cells are shown at right.

The right side of Figure 3 shows the results of applying a linear smoother  $\hat{\mu} = M(\lambda)\mathbf{s}$  to the 285-dimensional count vector  $\mathbf{s}$  given in Figure 2. The  $K \times K$  matrix  $M(\lambda)$  is a two-dimensional version of (2.14), with  $kk'$ th element

$$(4.4) \quad M_{kk'}(\lambda) = \frac{c_k}{\lambda^2} \exp\left(-\frac{1}{2\lambda^2}[(i-i')^2 + (j-j')^2]\right)$$

for  $k = (i, j)$  and  $k' = (i', j')$ . The  $c_k$  are chosen to make  $\sum_{k'} M_{kk'}(\lambda) = 1$ . The choice  $\lambda = 1.5$ , suggested by the expected deviance calculations of Section 6, gave the smoothed estimate  $\hat{\mu}^\circ = M(1.5) \cdot \mathbf{s}$  plotted versus  $y_{(k)}$  on the right side of Figure 3.

The left side of Figure 3 shows an SEF estimate  $\hat{\mu}$  of form(3.2), with

$$(4.5) \quad \hat{\mu}^\circ = M(2)\mathbf{s} \quad \text{and} \quad X = (\mathbf{1}, \mathbf{i}, \mathbf{j}, \mathbf{i}^2, \mathbf{j}^2, \mathbf{ij}).$$

Here  $\mathbf{i}$  is the  $K$  vector with  $k$ th element  $i - 8$ , and  $\mathbf{j}$  is the  $K$  vector with  $k$ th element  $j - 9$ . This choice amounts to using the usual sufficient statistics for a bivariate normal in (1.2),  $t(y) = (l, m, l^2, m^2, lm)$ . The fitted density matches the empirical means, variances and correlation of the galaxy data.

The variance calculations of Section 5 and the expected deviance calculations of Section 6 suggest that these two estimates are roughly equal in their overall ability to predict the true density. However, the SEF estimate is much smoother than the smoothing-only choice. This is obvious in Figure 4, which shows contour plots of the two density estimates.

Suppose that the carrier  $\hat{\mu}^\circ$  in (4.5) was taken to be the constant vector  $\hat{\mu}_k^\circ \equiv 1$  instead of  $M(2)\mathbf{s}$ . Then the SEF  $\hat{\mu}$  would be the (discretized) truncated bivariate normal MLE for the galaxy data, the truncation being to the rectangle (4.2). In fact, the SEF in Figure 3 looks like the lower corner of a bivariate normal, though there are some discrepancies due to the adaptation of  $\hat{\mu}^\circ$  to the galaxy data.

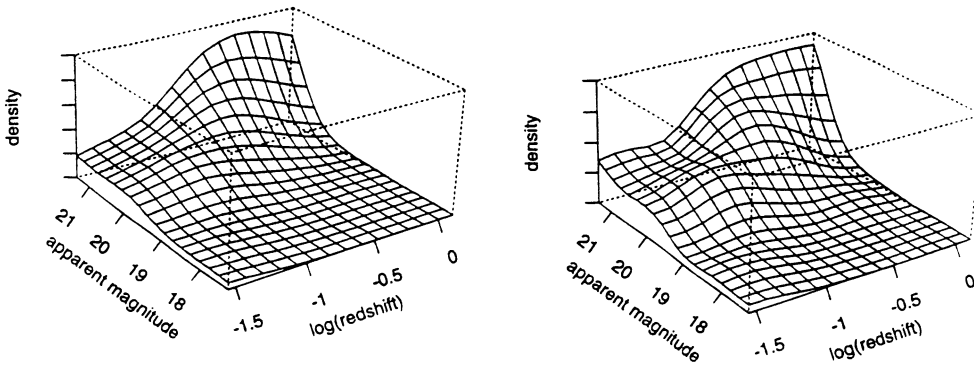


FIG. 3. Left panel: SEF estimate (3.2) for the galaxy data as discretized in Figure 2;  $\hat{\mu}^\circ = M(2) \cdot \mathbf{s}$ , quadratic matrix  $X$ , (4.5). Right panel: The smoothing-only estimate  $\hat{\mu}^\circ = M(1.5) \cdot \mathbf{s}$ . The calculations of Sections 5 and 6 suggest that the two estimates are roughly equal in accuracy. However, the SEF estimate is much smoother.

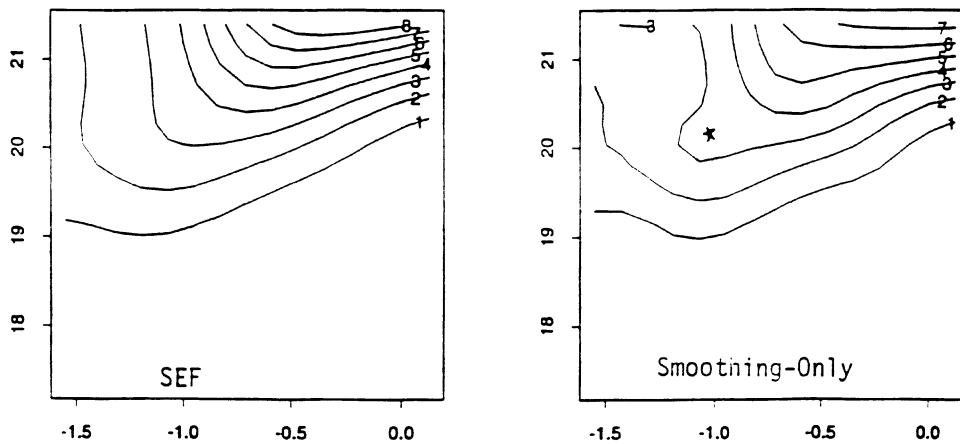


FIG. 4. Contour plots of the density estimates in Figure 3. The SEF contours (left panel) are much smoother than those from the smoothing-only estimate (right panel). The smoothing-only estimate has a weak second mode at the starred point.

Parametric models such as the bivariate normal are fierce data smoothers. This can be a big advantage if the statistician is interested in global properties of the density, like the general shape of its contours, and especially if the data-sampling process is suspect. In Figure 2 we can see that the galaxy data are quite patchy and clumpy, which is not surprising since the Loh-Spillar catalog was a census of the brighter galaxies in a few degrees of sky, and not a random sample of the full sky.

If those few degrees of sky are of great individual interest, then a narrow-band smoother like that in the right side of Figure 3 may be appropriate. [Notice that it estimates a weak second mode near  $(l, m) = (-1, 20)$ .] However if we really want to know the density for the whole sky, then it pays to *oversmooth* the estimate from a flawed sample like that in Figure 2. The SEF methodology allows us to oversmooth without losing much estimating efficiency compared to smoothing-only estimates and without making the drastic assumptions of the usual parametric models.

Various SEF models besides the quadratic choice of  $X$  in (4.5) were tried. One of these added a further cross-term to (4.5),

$$(4.6) \quad \hat{\mu}^o = M(2)\mathbf{s} \quad \text{and} \quad X = (\mathbf{1}, \mathbf{i}, \mathbf{j}, \mathbf{i}^2, \mathbf{j}^2, \mathbf{ij}, (\mathbf{i}^2\mathbf{j}^2)_{\text{orthog}}),$$

where  $(\mathbf{i}^2\mathbf{j}^2)_{\text{orthog}}$  is the component of  $\mathbf{i}^2\mathbf{j}^2$  orthogonal to  $(\mathbf{1}, \mathbf{i}, \mathbf{j}, \mathbf{i}^2, \mathbf{j}^2, \mathbf{ij})$ . Orthogonalization makes the first six components of  $\hat{\beta}$  have roughly the same values and standard errors as in (4.5). The term  $(\mathbf{i}^2\mathbf{j}^2)_{\text{orthog}}$  in (4.6) allows the regression surface to turn more quickly near the corners of the rectangle  $\mathcal{Z}$ .

The MLE vector  $\hat{\beta}_1$  for model (4.6) appears in Table 3, along with the standard error estimates  $\overline{\text{se}}$  from (3.6a) and the  $t$ -values  $\hat{\beta}/\overline{\text{se}}$ . All of the coefficients except the one for  $m^2$  are significantly different than zero [the

TABLE 3

Parameter estimates and standard errors for the components of  $\hat{\beta}_1$  in the SEF model (4.6). All of the components of  $\hat{\beta}_1$  are significantly nonzero, except for the coefficient of  $m^2$ .  $l$  indicates the coefficient corresponding to  $\mathbf{i}$ ,  $lm$  to  $\mathbf{ij}$  and so forth. Standard errors are from (3.6a). Model (4.5) gave similar estimates, standard errors and  $t$ -values for  $l, m, l^2, m^2, lm$

|                 | $l$    | $m$   | $l^2$   | $m^2$   | $lm$   | $(l^2m^2)_{\text{orthog}}$ |
|-----------------|--------|-------|---------|---------|--------|----------------------------|
| $\hat{\beta}_1$ | -0.105 | 0.084 | -0.0115 | -0.0018 | 0.0158 | 0.000267                   |
| se              | 0.019  | 0.015 | 0.0021  | 0.0014  | 0.0027 | 0.000071                   |
| $t$ -Value      | -5.6   | 5.8   | -5.5    | -1.4    | 5.8    | 3.75                       |

same is true for model (4.5)]. This includes the coefficient for the term  $(l^2m^2)_{\text{orthog}}$ , which takes us beyond the normal-theory SEF analogue (4.5). Compared to the left side of Figure 3, the SEF density estimate from (4.6) has its highest point at the upper right corner of the rectangle  $\mathcal{Z}$ .

The next two sections discuss several criteria for choosing among possible SEF models: observed deviance, expected deviance, degrees of freedom and so forth. However, none of these criteria is sharp enough to entirely free the statistician from model-choice quandries. A considerable amount of subjectivity must still go into the model-building process, just as in ordinary regression situations.

**5. Total relative variance.** A variety of diagnostic tools is available to assist model selection in standard regression situations. Similar tools are available for model selection in special exponential families. These ideas are developed in the next two sections, beginning here with the total relative variance, a simple measure of overall variability for an SEF density estimate. For example, looking ahead to Table 4 the reader can compare the total relative variances for the two galaxy-data SEF estimates in Figure 3: 6.8 for the quadratic model in the left panel versus 8.8 for the smoothing-only model in the right panel.

Let

$$(5.1) \quad \hat{\boldsymbol{\mu}} = \text{SEF}(\mathbf{s}; m, X)$$

indicate  $\hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}^o e^{X\hat{\boldsymbol{\beta}}}$ , the special exponential family estimate (3.2). First we will compute the  $K \times K$  derivative matrix of  $\hat{\boldsymbol{\mu}}$  with respect to  $\mathbf{s}$ , which leads immediately to a delta-method estimate of  $\text{Cov}(\hat{\boldsymbol{\mu}})$ . This derivative matrix involves the projection matrices

$$(5.2) \quad \hat{P} = X(X'DX)^{-1}X' \quad \text{and} \quad \hat{Q} = \hat{D}^{-1} - \hat{P},$$

where  $\hat{D} = D(\hat{\boldsymbol{\mu}})$  as before,  $\hat{P}$  is the symmetric projection matrix into the linear space spanned by the columns of  $X$ , in the inner product  $\hat{D}$ , the projection of vector  $\mathbf{v}$  being  $\hat{P}\hat{D}\mathbf{v}$ , and  $\hat{Q}$  is the projection orthogonal to  $X$ 's column space. Because they represent orthogonal projections, we have

$$(5.3) \quad \hat{P}\hat{D}\hat{P} = \hat{P}, \quad \hat{Q}\hat{D}\hat{Q} = \hat{Q} \quad \text{and} \quad \hat{P}\hat{D}\hat{Q} = 0.$$

LEMMA 2. *The derivative matrix of  $\hat{\boldsymbol{\mu}} = \text{SEF}(\mathbf{s}; m, X)$  with respect to  $\mathbf{s}$  is*

$$(5.4a) \quad \frac{d\hat{\boldsymbol{\mu}}}{d\mathbf{s}} = \hat{D}\hat{O}$$

where, in terms of  $\hat{H} = d \log \hat{\boldsymbol{\mu}}^\circ / d\mathbf{s}$ , (3.3),

$$(5.4b) \quad \hat{O} = \hat{P} + \hat{Q}\hat{D}\hat{H} = \hat{P} + \hat{H} - \hat{P}\hat{D}\hat{H}.$$

If  $\hat{\boldsymbol{\mu}}^\circ = M\mathbf{s}$ , then  $\hat{H} = D(1/\hat{\boldsymbol{\mu}}^\circ)M$  and

$$(5.4c) \quad \hat{O} = \hat{P} + \hat{Q}D(e^{X'\hat{\beta}})M.$$

(The proof appears below.)

The canonical parameter vector for the Poisson family (3.1) is

$$(5.5) \quad \boldsymbol{\eta} = \log(\boldsymbol{\mu}) = (\log(\mu_1), \log(\mu_2), \dots, \log(\mu_k)).$$

Likewise, define  $\hat{\boldsymbol{\eta}} = \log(\hat{\boldsymbol{\mu}})$  and  $\hat{\boldsymbol{\eta}}^\circ = \log(\hat{\boldsymbol{\mu}}^\circ)$ . Then we can write Lemma 2 as

$$(5.6) \quad \frac{d\hat{\boldsymbol{\eta}}}{d\mathbf{s}} = \hat{O} = \hat{P} + \hat{Q}\hat{D} \frac{d\hat{\boldsymbol{\eta}}^\circ}{d\mathbf{s}}.$$

This decomposes  $d\hat{\boldsymbol{\eta}}/d\mathbf{s}$  into a part  $\hat{P}$  coming from the exponential family factor  $e^{X\hat{\beta}}$  and an orthogonal part  $\hat{Q}\hat{D}(d\hat{\boldsymbol{\eta}}^\circ/d\mathbf{s})$  coming from the adaptive choice of the carrier.

Lemma 2 leads directly to delta-method estimates of the covariance matrix of  $\hat{\boldsymbol{\mu}}$ , as in (3.6),

$$(5.7) \quad \overline{\text{Cov}}(\hat{\boldsymbol{\mu}}) = \hat{D}\hat{O}\overline{D}\hat{O}'\hat{D} \quad \text{or} \quad \widehat{\text{Cov}}(\hat{\boldsymbol{\mu}}) = \hat{D}\hat{O}\hat{D}\hat{O}'\hat{D}$$

with  $\overline{D} = D(s)$ . The diagonal elements give variance estimates  $\overline{\text{var}}(\hat{\mu}_k)$  or  $\widehat{\text{var}}(\hat{\mu}_k)$  for the individual components. For Poisson variables it is natural to measure variance relative to the estimate  $\hat{\mu}_k$ . We define the *total relative variance* (TRV) estimate for  $\hat{\boldsymbol{\mu}}$  to be

$$(5.8) \quad \overline{\text{TRV}} = \sum_{k=1}^K \frac{\overline{\text{var}}(\hat{\mu}_k)}{\hat{\mu}_k} \quad \text{or} \quad \widehat{\text{TRV}} = \sum_{k=1}^K \frac{\widehat{\text{var}}(\hat{\mu}_k)}{\hat{\mu}_k}.$$

COROLLARY 2. *For  $\hat{\boldsymbol{\mu}} = \text{SEF}(\mathbf{s}; m, X)$  the total relative variance estimates are*

$$(5.9) \quad \overline{\text{TRV}} = \text{tr} \left[ \overline{D}\hat{P} + (\hat{H}\overline{D}\hat{H}')(\hat{D}\hat{Q}\hat{D}) \right] \quad \text{or} \\ \widehat{\text{TRV}} = (p+1) + \text{tr}(\hat{H}\hat{D}\hat{H}')(\hat{D}\hat{Q}\hat{D}),$$

where  $p+1$  is the number of columns of  $X$ ,  $\overline{D} = D(\mathbf{s})$ ,  $\hat{D} = D(\hat{\boldsymbol{\mu}})$  and  $\hat{P}$ ,  $\hat{Q}$  are as in (5.2).

The proof of Corollary 2 follows directly from (5.7) by writing  $\overline{\text{TRV}}$  in the trace form  $\text{tr } \hat{D}^{-1} \text{Cov}(\hat{\mu})$  and using the orthogonality relationship (5.3), and similarly for  $\overline{\text{TRV}}$ . Computationally more efficient expressions for  $\overline{\text{TRV}}$  and  $\overline{\text{TRV}}$  appear in Remark I, Section 9.

Table 4 shows  $\overline{\text{TRV}}$  for various SEF estimates for the galaxy data, as discretized in Figure 2. The three estimates correspond to three choices of  $X$  in (5.1): sef 2 is for  $X$  as in (4.5), sef 3 for  $X$  as in (4.6), and sef 0 for  $X = \mathbf{1}$ . This last choice is the smoothing-only estimate  $\hat{\mu}^o$  rescaled to  $(n/\hat{\mu}_+^o)\hat{\mu}^o$ , so that it sums to  $n = s_+$ . The carrier  $\hat{\mu}^o$  for the SEF is

$$(5.10) \quad \hat{\mu}^o = M(\lambda)\mathbf{s},$$

where  $M(\lambda)$  is the matrix (4.4).

All three TRV estimates decrease as the smoothing parameter  $\lambda$  increases because greater window width  $\lambda$  decreases the variability of  $\hat{\mu}^o$ . If the carrier  $\hat{\mu}^o$  were prechosen instead of adaptive, then sef2 would exceed sef 0 by about 5, this being the increased number of free parameters, and likewise sef 3 would exceed sef 2 by about 1. This is seen in (5.9) for the case  $\hat{H} = 0$ . In fact, this is nearly the case at the right side of Table 4, where  $\lambda$  is so large that  $\hat{\mu}^o$  has nearly constant entries. The difference between the three estimates decreases at smaller values of  $\lambda$  because the adaptability of the common carrier  $\hat{\mu}^o = M(\lambda) \cdot \mathbf{s}$  absorbs some of the difference in the exponential family fits  $e^{X\hat{\beta}}$ .

Small variability is a good property of course, but we also want  $\hat{\mu}$  to have small bias for estimating the true density vector  $\mu$ . The next section puts TRV into the context of bias–variance tradeoffs for Poisson regression estimates.

PROOF OF LEMMA 2. In the notation following (5.5),  $\hat{\eta} = \hat{\eta}^o + X\hat{\beta}$ . Differentiating this with respect to  $\mathbf{s}$  and using (3.4) gives

$$(5.11) \quad \frac{d\hat{\eta}}{d\mathbf{s}} = \hat{H} + X(X'\hat{D}X)^{-1}X'[I - \hat{D}\hat{H}] = \hat{P} + \hat{H} - \hat{P}\hat{D}\hat{H},$$

which is (5.6), the equivalent of Lemma 2.  $\square$

TABLE 4

Total relative variance  $\overline{\text{TRV}}$  for three SEF estimates, galaxy data, at increasing values of smoothing parameter  $\lambda$

|                        | $\lambda = 1.5$  | $\lambda = 2$    | $\lambda = 4$ | $\lambda = 6$ |
|------------------------|------------------|------------------|---------------|---------------|
| sef 0 (smoothing only) | 8.8 <sup>a</sup> | 5.2              | 1.2           | 0.4           |
| sef 2 (4.5)            | 9.6              | 6.8 <sup>b</sup> | 5.2           | 5.2           |
| sef 3 (4.6)            | 10.2             | 7.6              | 6.2           | 6.2           |

<sup>a</sup>Right panel of Figure 3.

<sup>b</sup>Left panel of Figure 3.

**6. Degrees of freedom and estimated deviance.** Selecting a good SEF estimate  $\hat{\mu}$  for a particular application involves making the usual tradeoffs between variance and bias. This section concerns estimating the total expected deviance of  $\hat{\mu}$  from the expectation vector  $\mu$ . This is a measure of accuracy for  $\hat{\mu}$  that involves both variance and bias. An important role is played by the *degrees of freedom* of the estimator  $\hat{\mu}$ , an idea related to the total relative variance of Section 5. The ideas here are an extension of those in Section 6.8 of Hastie and Tibshirani (1990).

Figure 5 relates to the *expected deviance* measure  $\widehat{\text{EDEV}}$  developed below, a diagnostic measure for comparing the goodness-of-fit of different SEF models. It shows  $\widehat{\text{EDEV}}$  for both the pain-score and galaxy examples.  $\widehat{\text{EDEV}}$  is plotted versus the smoothing parameter  $\lambda$  for the normal kernel estimates (2.13) and (4.4) used to obtain  $\hat{\mu}^o = M(\lambda) \cdot \mathbf{s}$ . The different curves correspond to different choices of  $X$  in the SEF formula (3.2). “Smoothing-only” refers to  $X = \mathbf{1}$ , which gives the estimate  $\hat{\mu}^o$  renormalized to sum to  $n$ . For the pain-score data, “linear” and “quadratic” are the cases referred to in (2.15). For the galaxy data, (4.5) and (4.6) are as in Table 4.

Notice that  $\widehat{\text{EDEV}}(\lambda)$  has a sharp minimum as a function of  $\lambda$  for both smoothing-only cases. (At  $\lambda = 0.63$  for the pain-score data and at  $\lambda = 1.5$  for the galaxy data.) This is not true for the genuine SEF estimates. They allow the smoothing parameter  $\lambda$  to be chosen much larger without incurring too much EDEV penalty. In other words, *the SEF methodology allows us to oversmooth the density estimate*, with the advantages seen in Figure 4.

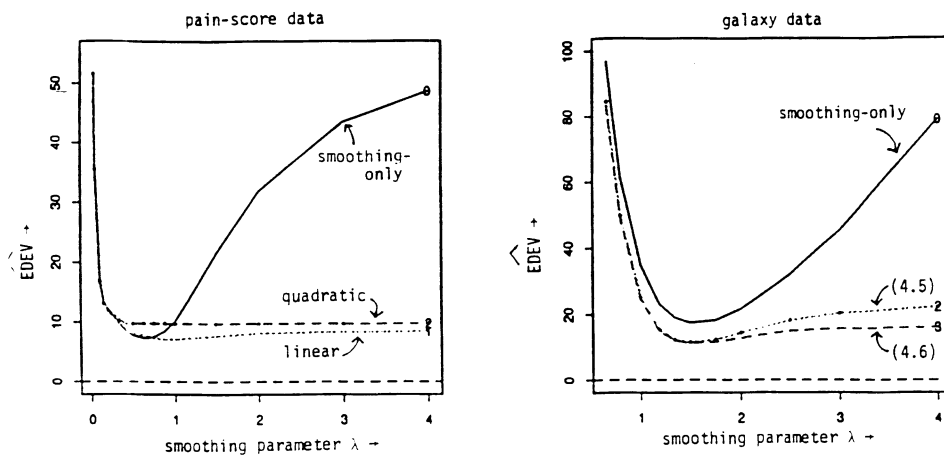


FIG. 5. Expected deviance estimates  $\widehat{\text{EDEV}}$ , (6.22). Left panel: SEF estimates for pain-score data (2.13)–(2.15). Right panel: SEF estimates for galaxy data as in Figure 5. In both cases the SEF estimates permit the use of larger smoothing parameters  $\lambda$ , compared to the smoothing-only estimates.



In order to motivate  $\widehat{\text{EDEV}}$ , it helps to begin the discussion with the normal case. Suppose that the statistician observes a  $K$ -dimensional normal vector

$$(6.1) \quad \mathbf{z} \sim N_K(\boldsymbol{\mu}, I)$$

and wishes to estimate  $\boldsymbol{\mu}$  using some linear estimator

$$(6.2) \quad \hat{\boldsymbol{\mu}} = S\mathbf{z},$$

$S$  being a  $K \times K$  matrix. Model selection amounts to making a good choice of  $S$ .

Let  $\text{err}$  be the observed total residual squared error of  $\hat{\boldsymbol{\mu}}$ :

$$(6.3) \quad \text{err} = \|\mathbf{z} - \hat{\boldsymbol{\mu}}\|^2.$$

Also define

$$(6.4) \quad \boldsymbol{\nu} = E\{\hat{\boldsymbol{\mu}}\}, \quad \text{DF} = \text{tr } S \quad \text{and} \quad \text{TV} = \text{tr } SS',$$

where DF stands for *degrees of freedom* and TV stands for *total variance*. Total variance TV equals  $\sum_k \text{var}(\hat{\mu}_k)$ , and if  $S$  is a projection matrix, then  $\text{tr } S$  is the usual degrees of freedom. It is easy to prove the following two relationships:

$$(6.5a) \quad E\{\text{err} - (K - 2\text{DF})\} = E\{\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2\}$$

and

$$(6.5b) \quad E\{\text{err} - (K - 2\text{DF} + \text{TV})\} = \|\boldsymbol{\nu} - \boldsymbol{\mu}\|^2.$$

Both (6.5a) and (6.5b) play an important role in normal-theory model selection. The first of these is essentially the  $C_p$  or AIC criterion. Its extension to nonlinear estimators is Stein's unbiased risk estimate. Extensions of formula (6.5b) are used in hypothesis testing. The quantity  $K - 2\text{DF} + \text{TV}$  equals the *residual degrees of freedom*:

$$(6.6) \quad \text{RDF} = \text{tr}(I - S)(I - S)'$$

If  $S$  is a projection matrix, then (6.5b) can be improved to

$$(6.7) \quad \text{err} \sim \chi_{\text{RDF}}^2(\|\boldsymbol{\nu} - \boldsymbol{\mu}\|^2),$$

where the notation indicates a noncentral chi-square variate with noncentrality parameter  $\|\boldsymbol{\nu} - \boldsymbol{\mu}\|^2$ . We test the adequacy of the estimate  $\hat{\boldsymbol{\mu}} = S\mathbf{z}$  by comparing  $\text{err}$  to a central chi-square distribution  $\chi_{\text{RDF}}^2$ .

We now return to the Poisson situation where we observe

$$(6.8) \quad \mathbf{s} \sim \text{Po}_K(\boldsymbol{\mu})$$

and estimate  $\boldsymbol{\mu}$  by some estimator  $\hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}(\mathbf{s})$ , not necessarily of the SEF form (3.2), having expectation

$$(6.9) \quad \boldsymbol{\nu} \equiv E\{\hat{\boldsymbol{\mu}}(\mathbf{s})\}.$$

The *observed error*  $\text{err} = \text{err}(\mathbf{s})$  in the Poisson context is

$$(6.10) \quad \text{err}(\mathbf{s}) = \text{Dev}(\mathbf{s}, \hat{\boldsymbol{\mu}}),$$

where Dev indicates the total Poisson deviance

$$\text{Dev}(\boldsymbol{\mu}, \boldsymbol{\nu}) = 2 \sum_k \{\log(\mu_k/\nu_k) - (\mu_k - \nu_k)\}.$$

The Poisson equivalents of  $K$ , DF and TV in (6.5) are

$$(6.11) \quad K(\boldsymbol{\mu}) = E\{\text{Dev}(\mathbf{s}, \boldsymbol{\mu})\}, \quad \text{DF}(\boldsymbol{\mu}) = E\left\{\sum_k (s_k - \mu_k) \log(\hat{\mu}_k)\right\},$$

$$\text{TRV}(\boldsymbol{\mu}) = 2 \sum_k E\{\mu_k \log(\nu_k/\hat{\mu}_k)\},$$

all expectations being with respect to  $\mathbf{s} \sim \text{Po}_K(\boldsymbol{\mu})$ .

LEMMA 3. *In the Poisson situation (6.8)–(6.11),*

$$(6.12a) \quad E\{\text{err}(\mathbf{s}) - [K(\boldsymbol{\mu}) - 2\text{DF}(\boldsymbol{\mu})]\} = E\{\text{Dev}(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}})\}$$

and

$$(6.12b) \quad E\{\text{err}(\mathbf{s}) - [K(\boldsymbol{\mu}) - 2\text{DF}(\boldsymbol{\mu}) + \text{TRV}(\boldsymbol{\mu})]\} = \text{Dev}(\boldsymbol{\mu}, \boldsymbol{\nu}).$$

(The proof is given below.) Formulas (6.12a, b) are the Poisson versions of (6.5a, b).

In order to use Lemma 3, we can approximate  $K(\boldsymbol{\mu})$ ,  $\text{DF}(\boldsymbol{\mu})$  and  $\text{TRV}(\boldsymbol{\mu})$  by their plug-in estimates  $K(\hat{\boldsymbol{\mu}})$ ,  $\text{DF}(\hat{\boldsymbol{\mu}})$  and  $\text{TRV}(\hat{\boldsymbol{\mu}})$ . Using bootstrap notation, let  $\mathbf{s}^*$  given  $\mathbf{s}$  have Poisson distribution

$$(6.13) \quad \mathbf{s}^* | \mathbf{s} \sim \text{Po}_K(\hat{\boldsymbol{\mu}}(\mathbf{s}))$$

and let  $E_*$  indicate expectations with respect to (6.13), with  $\mathbf{s}$  and  $\hat{\boldsymbol{\mu}}(\mathbf{s})$  fixed. Then

$$(6.14) \quad \begin{aligned} K(\hat{\boldsymbol{\mu}}) &= \sum_k E_* \{\text{Dev}(s_k^*, \hat{\mu}_k)\} \\ &= \sum_k E_* \{2[s_k^* \log(s_k^*/\hat{\mu}_k) - (s_k^* - \hat{\mu}_k)]\}. \end{aligned}$$

It is easy to evaluate (6.14) numerically since it is the sum of univariate Poisson deviances, the  $\hat{\mu}_k$  being just fixed constants in the  $E_*$  expectations.

Using this same notation we can write the degrees of freedom estimate  $\text{DF}(\hat{\boldsymbol{\mu}})$  as

$$(6.15) \quad \text{DF}(\hat{\boldsymbol{\mu}}) = E_* \left\{ \sum_k (s_k^* - \hat{\mu}_k) \log \hat{\mu}_k^* \right\} = E_* (\mathbf{s}^* - \hat{\boldsymbol{\mu}})' \hat{\boldsymbol{\eta}}^*,$$

where  $\hat{\boldsymbol{\eta}}^* = \log(\hat{\boldsymbol{\mu}}^*) = \log(\hat{\boldsymbol{\mu}}(\mathbf{s}^*))$ , the vector with  $k$ th component  $\log \hat{\mu}_k^*$ . Since by (5.6),  $d\hat{\boldsymbol{\eta}}/d\mathbf{s} = \hat{O}$ , (6.15) has the Taylor series approximation

$$(6.16) \quad \text{DF}(\hat{\boldsymbol{\mu}}) \doteq E_* (\mathbf{s}^* - \hat{\boldsymbol{\mu}})' \hat{O} (\mathbf{s}^* - \hat{\boldsymbol{\mu}}) = \text{tr } \hat{D}\hat{O} \equiv \overline{\text{DF}}.$$

See Remark K. An alternative estimate is  $\overline{\text{DF}} = \text{tr } \overline{D}\hat{O}$  as in (5.8).

The quadratic expansion  $\log(\hat{\mu}/\nu) \doteq (\hat{\mu}/\nu - 1) - \frac{1}{2}(\hat{\mu}/\nu - 1)^2$  gives

$$(6.17) \quad E\left\{ \mu_k \log\left(\frac{\nu_k}{\hat{\mu}_k}\right) \right\} \doteq \frac{\mu_k \operatorname{var}(\hat{\mu}_k)}{2 \nu_k^2} \doteq \frac{\operatorname{var}(\hat{\mu}_k)}{2 \mu_k},$$

permitting us to approximate (6.11) by

$$(6.18) \quad \operatorname{TRV}(\boldsymbol{\mu}) \doteq \sum_k \frac{\operatorname{var}(\hat{\mu}_k)}{\mu_k}.$$

The quantities  $\overline{\operatorname{TRV}}$  and  $\widehat{\operatorname{TRV}}$  in (5.8) are the obvious plug-in estimates for (6.18). The substitution of  $\mu_k$  for  $\nu_k$  in the denominator of (6.17) looks worrisome, but Remark K of Section 9 shows that the resulting error is asymptotically negligible.

The TRV estimates (5.8) can be written as

$$(6.19) \quad \overline{\operatorname{TRV}} = \operatorname{tr} \overline{D\hat{O}'\hat{D}\hat{O}} \quad \text{or} \quad \widehat{\operatorname{TRV}} = \operatorname{tr} \hat{D}\hat{O}'\hat{D}\hat{O},$$

using (5.7), compared to the DF estimates

$$(6.20) \quad \overline{\operatorname{DF}} = \operatorname{tr} \overline{D\hat{O}} \quad \text{or} \quad \widehat{\operatorname{DF}} = \operatorname{tr} \hat{D}\hat{O}.$$

Notice the similarity to the normal-theory definitions in (6.4). Comparing (6.5b) with (6.12b), we can also define residual degrees of freedom RDF for the Poisson situation, estimated by

$$(6.21) \quad \overline{\operatorname{RDF}} = K(\hat{\boldsymbol{\mu}}) - 2\overline{\operatorname{DF}} + \overline{\operatorname{TRV}} \quad \text{or} \quad \widehat{\operatorname{RDF}} = K(\hat{\boldsymbol{\mu}}) - 2\widehat{\operatorname{DF}} + \widehat{\operatorname{TRV}}.$$

See Remark L. The quantity graphed in Figure 5 was the expected deviance estimate obtained from (6.12a), (6.14) and (6.20):

$$(6.22) \quad \widehat{\operatorname{EDEV}} = \operatorname{err}(\mathbf{s}) - K(\hat{\boldsymbol{\mu}}) + 2\widehat{\operatorname{DF}}.$$

The vertical scale in Figure 5 is misleading. For the galaxy data, reducing  $\lambda$  from 4 to 1.5 reduces  $\widehat{\operatorname{EDEV}}$  for the quadratic SEF (4.5) from 22.2 to 11.5, a considerable amount. Is this significant? The observed deviance error (6.10),  $\operatorname{err}$ , decreases by 45.3, while the residual degrees of freedom  $\overline{\operatorname{RDF}}$ , (6.21), decreases by 29.3. This gives the naive chi-square significance value  $\operatorname{prob}\{\chi_{29.3}^2 > 45.3\} = 0.03$ . A more trustworthy significance level could be obtained by Monte Carlo methods, bootstrapping with  $\mathbf{s}^* \sim \operatorname{Prob}(\hat{\boldsymbol{\mu}})$ , with  $\hat{\boldsymbol{\mu}}$  the (1.5, 2) SEF vector.

The quantitative aspects of Figure 5 cannot be taken too literally. For instance, using  $\widehat{\operatorname{EDEV}}$  instead of  $\overline{\operatorname{EDEV}}$  moved the smoothing-only curve below the others for  $\lambda < 2$  in the galaxy data. In general the caretted (hat) formulas gave less erratic answers than the bar formulas, but there are really no strong reasons for preferring  $\widehat{\operatorname{EDEV}}$ . The fact is that it is usually difficult to estimate the performance of competing decision rules, and the SEF density estimates are no exception. Formulas like (5.9) and (6.12) help with model selection, but considerable subjectivity remains. The main point of Figure 5 is the qualitative one that the SEF estimates permit extensive oversmoothing.

**PROOF OF LEMMA 3.** Define  $\mathcal{E}(\boldsymbol{\mu}) = E\{\operatorname{Dev}(\mathbf{s}^\dagger, \hat{\boldsymbol{\mu}})\}$ , where  $\mathbf{s}^\dagger \sim \operatorname{Po}_K(\boldsymbol{\mu})$  independent of  $\mathbf{s}$ , so  $\mathcal{E}(\boldsymbol{\mu})$  is the expected prediction error of  $\hat{\boldsymbol{\mu}}(\mathbf{s})$ . Efron (1986)

shows that

$$(6.23a) \quad \mathcal{E}(\boldsymbol{\mu}) = E\{\text{err}(\mathbf{s}) + 2\text{DF}(\boldsymbol{\mu})\}.$$

We also have the identities

$$(6.23b) \quad \mathcal{E}(\boldsymbol{\mu}) = K(\boldsymbol{\mu}) + E\{\text{Dev}(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}})\}$$

and

$$(6.23c) \quad E\{\text{Dev}(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}})\} = \text{Dev}(\boldsymbol{\mu}, \boldsymbol{\nu}) + \text{TRV}(\boldsymbol{\mu}).$$

All three of these relationships are easy to prove directly from the definition of the total Poisson deviance. Substituting (6.23b) into (6.12a) gives (6.12a). Substituting (6.23c) on the right side of (6.12a) gives (6.12b).  $\square$

**7. Multisample problems.** So far we have only considered one-sample problems. SEF estimates are particularly useful for investigating density differences in multisample situations. We use the exponential family model (1.2) for the different densities, with a shared carrier  $g_0$  estimated nonparametrically, but with possibly different values of the exponential  $\beta$  parameters. An example will precede the theory.

Figure 6 concerns a two-sample application of SEF modelling. The data are the compliances of men in the Stanford arm of a randomized trial of the cholesterol-lowering drug Cholestyramine; see Efron and Feldman (1991). There were  $n_1 = 172$  men in the control group and  $n_2 = 165$  men in the treatment group. Compliance ran from 0 to 100%, so  $\mathcal{Y} = [0, 100]$ . A discretization  $\mathcal{Y} = \cup_k \mathcal{Y}_k$  partitioned  $\mathcal{Y}$  into  $K = 46$  intervals of equal length. The left panel of Figure 6 shows the counts in the two groups. Compliance is significantly worse in the treatment group, as shown by standard two-sample tests.

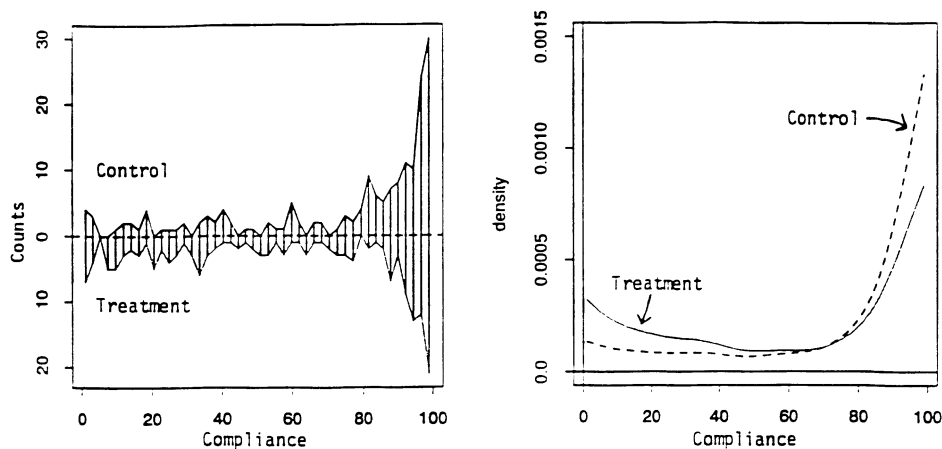


FIG. 6. An application of SEF modelling to the Cholestyramine trial compliance data of Efron and Feldman (1991). Left panel: The count vectors for the control group and the treatment group;  $K = 46$  equal divisions of  $\mathcal{Y} = [0, 100]$ . Right panel: SEF density estimates (7.4) for the two groups;  $m$  and  $X$  as in (7.11). The poorer compliance in the treatment group is graphically evident.

The right panel of Figure 6 is the output of a two-sample SEF analysis. It neatly displays the compliance differences between the two groups in terms of their estimated densities. The ratio  $\hat{g}_{\text{treat}}(y)/\hat{g}_{\text{cont}}(y)$  decreases almost linearly as  $y$  goes from 0 to 100%, but both densities are greatest at  $y = 100\%$ .

These densities were derived from a simple extension of the previous theory. In the multisample situation we observe independent random samples from  $L$  possibly different densities  $g_1, g_2, \dots, g_L$  on the same sample space  $\mathcal{Y}$ ,

$$(7.1) \quad y_i(l) \stackrel{\text{i.i.d.}}{\sim} g_l(y) \quad \text{for } i = 1, 2, \dots, n_l \text{ and } l = 1, 2, \dots, L.$$

We discretize the problem as in Section 2, obtaining count vector  $\mathbf{s}(l)$  for the  $l$ th sample,

$$(7.2) \quad s_k(l) = \#\{y_i(l) \in \mathcal{Y}_k\} \quad \text{for } k = 1, 2, \dots, K \text{ and } l = 1, 2, \dots, L.$$

The many-sample version of the Poisson regression model (3.1) is

$$(7.3) \quad \mathbf{s}(l) \stackrel{\text{ind.}}{\sim} \text{Po}_K(\boldsymbol{\mu}(l)) \quad \text{for } l = 1, 2, \dots, L, \text{ with } \boldsymbol{\mu}(l) = \boldsymbol{\mu}^\circ e^{X\beta(l)}.$$

The SEF estimates corresponding to (3.2) are

$$(7.4a) \quad \hat{\boldsymbol{\mu}}(l) = \hat{\boldsymbol{\mu}}^\circ e^{X\hat{\beta}(l)},$$

where

$$(7.4b) \quad \hat{\boldsymbol{\mu}}^\circ = m(\mathbf{s}(1), \mathbf{s}(2), \dots, \mathbf{s}(L)) \quad \text{and} \quad \hat{\beta}(l): X'[\mathbf{s}(l) - \hat{\boldsymbol{\mu}}^\circ e^{X\hat{\beta}(l)}] = 0.$$

In this model, a common carrier  $\hat{\boldsymbol{\mu}}^\circ$ , estimated from all of the data  $(\mathbf{s}(1), \mathbf{s}(2), \dots, \mathbf{s}(L))$  is modified by an exponential factor  $e^{X\hat{\beta}(l)}$  which can vary with  $l$ .

The many-sample version of Lemma 1 in Section 3 is the following lemma:

LEMMA 4. *Let*

$$(7.5) \quad \hat{H}(l) = \frac{\partial \log(\hat{\boldsymbol{\mu}}^\circ)}{\partial \mathbf{s}(l)}$$

*be the  $K \times K$  matrix with  $kj$ th entry  $\partial \log(\hat{\mu}_k^\circ)/\partial s_j(l)$ . Then the  $(p + 1) \times K$  derivative matrix of  $\hat{\beta}(l)$  with respect to  $\mathbf{s}(j)$  is*

$$(7.6a) \quad \frac{\partial \hat{\beta}(l)}{\partial \mathbf{s}(j)} = \hat{G}(l)^{-1} Z'_{lj},$$

where

$$(7.6b) \quad \hat{G}(l) = X' \hat{D}(l) X \quad \text{and} \quad Z'_{lj} = X'_{lj} [\delta_{lj} I - \hat{D}(l) \hat{H}(j)],$$

$\hat{D}(l) \equiv D(\hat{\boldsymbol{\mu}}(l))$ ,  $\delta_{ij}$  equalling 1 or 0 as  $l$  does or does not equal  $j$ . If

$$(7.7) \quad \hat{\boldsymbol{\mu}}^\circ = M \mathbf{s}_+ \quad \left( \mathbf{s}_+ \equiv \sum_l \mathbf{s}_l \right),$$

then

$$(7.8) \quad Z'_{lj} = X'_{lj} [\delta_{lj} I - D(e^{X\hat{\beta}(l)}) M].$$

(The proof is nearly the same as for Lemma 1 and will not be given here.)

Lemma 4 leads to delta-method estimates of the covariance matrix for  $\hat{\beta}(l)$ , just as in (3.6),

$$(7.9) \quad \begin{aligned} \overline{\text{Cov}}(\hat{\beta}(l)) &= \sum_{j=1}^L \hat{G}(l)^{-1} [Z'_{lj} \bar{D}(j) Z_{lj}] \hat{G}(l)^{-1} \quad \text{or} \\ \widehat{\text{Cov}}(\hat{\beta}(l)) &= \sum_{j=1}^L \hat{G}(l)^{-1} [Z'_{lj} \hat{D}(j) Z_{jl}] \hat{G}(l)^{-1}, \end{aligned}$$

$\bar{D}(j) \equiv D(\mathbf{s}_j)$ . We can also obtain covariance estimates for functions of the  $\hat{\beta}(l)$ . For example, if  $\hat{\gamma} = \hat{\beta}(2) - \hat{\beta}(1)$ , then

$$(7.10a) \quad \begin{aligned} \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(1)} &= \hat{G}(2)^{-1} Z'_{21} - \hat{G}(1)^{-1} Z'_{11}, \\ \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(2)} &= \hat{G}(2)^{-1} Z'_{22} - \hat{G}(1)^{-1} Z'_{12} \end{aligned}$$

and

$$(7.10b) \quad \begin{aligned} \overline{\text{Cov}}(\hat{\gamma}) &= \sum_{j=1}^2 \left( \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(j)} \right) \bar{D}(j) \left( \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(j)} \right)', \\ \widehat{\text{Cov}}(\hat{\gamma}) &= \sum_{j=1}^2 \left( \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(j)} \right) \hat{D}(j) \left( \frac{\partial \hat{\gamma}}{\partial \mathbf{s}(j)} \right)'. \end{aligned}$$

SEF model (7.4) was used to estimate the two compliance densities in Figure 6. The  $k$ th row of  $X$  was quadratic in compliance,

$$(7.11a) \quad x_k = (1, \tilde{y}_{(k)}, \tilde{y}_{(k)}^2),$$

where  $\tilde{y}_{(k)} = y_{(k)} - 50/100$ ,  $y_{(k)}$  being the midpoint of  $\mathcal{Z}_k$ ;  $\hat{\boldsymbol{\mu}}^o = M\mathbf{s}_1$  as in (7.7), with

$$(7.11b) \quad M = M(7),$$

as in (2.14). The estimated difference  $\hat{\gamma} = \hat{\beta}(\text{treatment}) - \hat{\beta}(\text{control})$  was

$$(7.12) \quad (0.26, -1.36, -0.44) \pm (0.22, 0.38, 1.51),$$

with the standard errors taken from  $\widehat{\text{Cov}}(\hat{\gamma})$ , (7.10b). We see that the linear coefficient is significantly negative,  $t$ -value =  $-3.6$ , but the quadratic coefficient is not.

**8. Other types of estimators.** In the SEF methodology the application of an initial nonparametric smoother is followed by the fitting of an exponential family parametric model. Hjort and Glad's (1995) semiparametric density

estimator reverses this order, with the parametric model coming before the smoother. *Backfitting*, as applied to generalized linear models like (3.1), repeatedly iterates between the parametric and nonparametric fitting methods. This section discusses the SEF methodology in the context of these other possibilities.

As in Section 6, we begin with the normal case where the statistician wishes to estimate  $\boldsymbol{\mu}$  having observed

$$(8.1) \quad \mathbf{z} \sim N_K(\boldsymbol{\mu}, I).$$

The linear model  $\boldsymbol{\mu} = \boldsymbol{\xi} + X\boldsymbol{\beta}$ , with  $\boldsymbol{\xi}$  a known *origin* vector and  $X$  a known  $K \times (p + 1)$  structure matrix, gives the estimate

$$(8.2) \quad \hat{\boldsymbol{\mu}} = \mathcal{P}(\mathbf{z}; \boldsymbol{\xi}) \equiv \boldsymbol{\xi} + P(\mathbf{z} - \boldsymbol{\xi}).$$

Here  $P = X(X'X)^{-1}X'$  is the  $K \times K$  projection matrix into  $\mathcal{L}(X)$ , the column space of  $X$ . In the normal case,  $\mathcal{P}(\mathbf{z}; \boldsymbol{\xi})$  plays the role of the parametric exponential family model estimator. The role of the nonparametric smoother is played by

$$(8.3) \quad \mathcal{M}(\mathbf{z}; \boldsymbol{\xi}) = \boldsymbol{\xi} + M(\mathbf{z} - \boldsymbol{\xi}),$$

where  $M$  is some fixed  $K \times K$  smoothing matrix such as  $M(\lambda)$  in (2.13). Once again  $\boldsymbol{\xi}$  represents a fixed and known origin. Often we take  $\boldsymbol{\xi} = \mathbf{0}$ , but it will be important here to consider more general choices of the origin.

The normal-theory analog of the SEF estimate (3.2) is

$$(8.4) \quad \begin{aligned} \hat{\boldsymbol{\mu}}_{\text{sef}} &= \mathcal{P}(\mathbf{z}; \mathcal{M}(\mathbf{z}; \mathbf{0})) = (P + M - PM)\mathbf{z} \\ &= (P + QM)\mathbf{z}. \end{aligned}$$

Here  $Q = I - P$  is the projection matrix into  $\mathcal{L}(X)^\perp$ , the orthocomplement to  $\mathcal{L}(X)$ . In other words, we begin with  $\mathbf{0}$  as the origin, apply  $\mathcal{M}$  to get an updated origin  $\hat{\boldsymbol{\mu}}^o = \mathcal{M}(\mathbf{z}, \mathbf{0}) = M\mathbf{z}$  and finally take the estimate of  $\boldsymbol{\mu}$  to be  $\hat{\boldsymbol{\mu}} = \mathcal{P}(\mathbf{z}; \hat{\boldsymbol{\mu}}^o) = \hat{\boldsymbol{\mu}}^o + P(\mathbf{z} - \hat{\boldsymbol{\mu}}^o)$ . Notice that

$$(8.5) \quad X'\mathbf{z} = X'\hat{\boldsymbol{\mu}}$$

according to (8.4), which says that  $\mathbf{z}$  and  $\hat{\boldsymbol{\mu}}$  have the same projection into  $\mathcal{L}(X)$ , namely,  $P\mathbf{z}$ . Equality (8.5) is the normal-theory analog of the moment-matching property (2.17).

Reversing the order of  $\mathcal{P}$  and  $\mathcal{M}$ , as Hjort and Glad do in the density estimation problem, gives the estimator

$$(8.6) \quad \hat{\boldsymbol{\mu}}_{\text{HG}} = \mathcal{M}(\mathbf{z}; \mathcal{P}(\mathbf{z}; \mathbf{0})) = (P + MQ)\mathbf{z}.$$

This no longer enjoys the moment-matching property (8.5), but we can restore it with one further application of  $\mathcal{P}$ . This defines the *symmetrized estimator*

$$(8.7) \quad \hat{\boldsymbol{\mu}}_{\text{sym}} = \mathcal{P}(\mathbf{z}; \hat{\boldsymbol{\mu}}_{\text{HG}}) = (P + QMQ)\mathbf{z}.$$

If  $M$  is a symmetric matrix with eigenvalues between 0 and 1, then so is  $P + QMQ$ . This makes  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  a formal Bayes estimator for  $\boldsymbol{\mu}$  in situation (8.1), versus a normal prior distribution on  $\boldsymbol{\mu}$  possibly having infinite variance in some directions.

The backfitting estimator  $\hat{\boldsymbol{\mu}}_{\text{back}}$  is defined as follows in Chapter 5 of Hastie and Tibshirani (1990): suppose we can find vectors  $\hat{\boldsymbol{\mu}}^o$  and  $\hat{\boldsymbol{\mu}}^1$  such that

$$(8.8a) \quad \hat{\boldsymbol{\mu}}^o = \mathcal{M}(\mathbf{z}; \hat{\boldsymbol{\mu}}^1) - \hat{\boldsymbol{\mu}}^1 \quad \text{and} \quad \hat{\boldsymbol{\mu}}^1 = \mathcal{P}(\mathbf{z}; \hat{\boldsymbol{\mu}}^o) - \hat{\boldsymbol{\mu}}^o.$$

Then

$$(8.8b) \quad \hat{\boldsymbol{\mu}}_{\text{back}} = \hat{\boldsymbol{\mu}}^o + \hat{\boldsymbol{\mu}}^1$$

[so  $\hat{\boldsymbol{\mu}}_{\text{back}} = \mathcal{M}(\mathbf{z}; \hat{\boldsymbol{\mu}}^1) = P(\mathbf{z}; \hat{\boldsymbol{\mu}}^o)$ .] Letting  $\mathbf{v} = \mathbf{z} - \hat{\boldsymbol{\mu}}^1$ , it is easy to show that

$$(8.9) \quad \mathbf{v} - M\mathbf{v} = (\mathbf{z} - \hat{\boldsymbol{\mu}}^1) - \hat{\boldsymbol{\mu}}^o \in \mathcal{L}^\perp(X),$$

which implies the moment-matching property (8.5). Hastie and Tibshirani show that  $\hat{\boldsymbol{\mu}}_{\text{back}} = S_{\text{back}}\mathbf{z}$  for a matrix  $S_{\text{back}}$  that is symmetric and with eigenvalues in  $[0, 1]$  if  $M$  has these same properties.

Further insight into the backfitting estimator can be gained by expressing it in an explicit form. The backfitting estimator satisfies

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (X'X)^{-1}X'(\mathbf{z} - \hat{\boldsymbol{\mu}}^o), \\ \hat{\boldsymbol{\mu}}^1 &= M(\mathbf{z} - X\hat{\boldsymbol{\beta}}). \end{aligned}$$

Solving these equations yields the explicit expression

$$\hat{\boldsymbol{\beta}} = [X'(I - M)X]^{-1}X'(I - M)\mathbf{z}.$$

This looks like a weighted least squares estimate with weight matrix  $I - M$  or, equivalently, variance matrix  $(I - M)^{-1}$ . It has the same form as Hjort and Glad's estimator except that they use an unweighted least squares estimator. The weights  $I - M$  can be justified from a mixed effects model in which  $\boldsymbol{\mu}^1$  is a random effect.

We have discussed four linear estimators  $\hat{\boldsymbol{\mu}} = S\mathbf{z}$ , three of which have the moment-matching property (8.5). These three can be described as follows: the  $\mathcal{L}(X)$  component of  $\hat{\boldsymbol{\mu}}^1$  matches the  $\mathcal{L}(X)$  component of  $\mathbf{z}$ ; the  $\mathcal{L}^\perp(X)$  component of  $\hat{\boldsymbol{\mu}}$  equals the  $\mathcal{L}^\perp(X)$  component of a point  $\mathbf{v}$  in the flat space  $\mathbf{z} \oplus \mathcal{L}(X)$ ;  $\mathbf{v} = \mathbf{z}$  for  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ ,  $\mathbf{v} = Q\mathbf{z}$  for  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  and  $\mathbf{v}$  equals the point or points satisfying the orthogonality condition (8.9) for  $\hat{\boldsymbol{\mu}}_{\text{back}}$ . *Note:* All four estimators are the same if  $M$  and  $P$  commute,  $MP = PM$ .

We calculated the "equivalent kernels" for each of the four estimators, as in Figure 2.5 of Hastie and Tibshirani (1990). The calculation was done for the situation that produced the quadratic estimator in Figure 1: 40 equally spaced  $x$  values,  $P$  based on a matrix  $X$  with rows  $(1, x, x^2)$  and  $M$  of the form (2.14). The smoothing parameter  $\lambda$  was chosen to give  $\text{tr} S = 5$  in each case. The SEF and backfitting kernels were remarkably similar, with both the HG and symmetrical kernels being slightly different.

We can define analogs to  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ ,  $\hat{\boldsymbol{\mu}}_{\text{HG}}$ ,  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  and  $\hat{\boldsymbol{\mu}}_{\text{back}}$  for the Poisson situation  $\mathbf{s} \sim \text{Po}_K(\boldsymbol{\mu})$ . Given a positive origin vector  $\boldsymbol{\xi}$  and structure matrix  $X$ , we let the analog of (8.2) be

$$(8.10) \quad \mathcal{P}(\mathbf{s}; \boldsymbol{\xi}) \equiv D(\boldsymbol{\xi})e^{X\hat{\boldsymbol{\beta}}(\boldsymbol{\xi})}, \quad \text{where } X'[\mathbf{s} - D(\boldsymbol{\xi})e^{X\hat{\boldsymbol{\beta}}(\boldsymbol{\xi})}] = 0;$$



$\hat{\boldsymbol{\mu}} = \mathcal{P}(\mathbf{s}; \boldsymbol{\xi})$  is the MLE estimate of  $\boldsymbol{\mu}$  in the “offset” generalized linear model  $\boldsymbol{\mu} = D(\boldsymbol{\xi})e^{X\boldsymbol{\beta}}$ . The analog of (8.3) is

$$(8.11) \quad \mathcal{M}(\mathbf{s}; \boldsymbol{\xi}) = D(\boldsymbol{\xi})MD(1/\boldsymbol{\xi})\mathbf{s} = D(\boldsymbol{\xi})M(\mathbf{s}/\boldsymbol{\xi});$$

$\hat{\boldsymbol{\mu}} = \mathcal{M}(\mathbf{s}; \boldsymbol{\xi})$  is a discretized version of what Hjort and Glad (1994) call a *nonparametric density estimate with a parametric start*, the “start” being the choice of  $\boldsymbol{\xi}$ .

The Poisson analogs of  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ ,  $\hat{\boldsymbol{\mu}}_{\text{HG}}$  and  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  [(8.4)–(8.7)] are

$$(8.12) \quad \hat{\boldsymbol{\mu}}_{\text{sef}} = \mathcal{P}(\mathbf{s}; \mathcal{M}(\mathbf{s}; \mathbf{1})), \quad \hat{\boldsymbol{\mu}}_{\text{HG}} = \mathcal{M}(\mathbf{s}; \mathcal{P}(\mathbf{s}; \mathbf{1})), \quad \hat{\boldsymbol{\mu}}_{\text{sym}} = \mathcal{P}(\mathbf{s}; \hat{\boldsymbol{\mu}}_{\text{HG}}).$$

The backfitting estimate (8.8) is now defined by  $\hat{\boldsymbol{\mu}}_{\text{back}} = \hat{\boldsymbol{\mu}}^o + \hat{\boldsymbol{\mu}}^1$ , where  $\hat{\boldsymbol{\mu}}^o$  and  $\hat{\boldsymbol{\mu}}^1$  satisfy the fixed-point relationships

$$(8.13) \quad \hat{\boldsymbol{\mu}}^o = D(1/\hat{\boldsymbol{\mu}}^1)\mathcal{M}(\mathbf{s}; \hat{\boldsymbol{\mu}}^1) \quad \text{and} \quad \hat{\boldsymbol{\mu}}^1 = D(1/\hat{\boldsymbol{\mu}}^o)\mathcal{P}(\mathbf{s}; \hat{\boldsymbol{\mu}}^o).$$

Chapter 6 of Hastie and Tibshirani (1990) discusses an iterative algorithm for computing  $\hat{\boldsymbol{\mu}}_{\text{back}}$ .

The moment-matching property (2.17) is satisfied by  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ ,  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  and  $\hat{\boldsymbol{\mu}}_{\text{back}}$ . Each of these estimators can be written in the form  $\hat{\boldsymbol{\mu}}^o e^{X\boldsymbol{\beta}}$ , (3.26), with

$$(8.14) \quad \hat{\boldsymbol{\mu}}_{\text{sef}}^o = M\mathbf{s}, \quad \hat{\boldsymbol{\mu}}_{\text{sym}}^o = \hat{\boldsymbol{\mu}}_{\text{HG}} \quad \text{and} \quad \hat{\boldsymbol{\mu}}_{\text{back}}^o = M(\mathbf{s}/\hat{\boldsymbol{\mu}}^1).$$

We have used  $\hat{\boldsymbol{\mu}}_{\text{sef}}^o$  in all of our examples because it makes the computation of  $\hat{H} = d \log(\hat{\boldsymbol{\mu}}^o)/d\mathbf{s}$ , a crucial part of the formulas in Sections 2–7, so simple. [An even simpler choice,  $\log(\hat{\boldsymbol{\mu}}^o) = H\mathbf{s}$  for some fixed matrix  $H$ , does not satisfy (5.14) and seems to have undesirable small-sample properties.]

In theory at least we can compute  $\hat{H}$  for any choice of  $\hat{\boldsymbol{\mu}}^o = m(\mathbf{s})$ . Here, without proof, is  $\hat{H}$  for  $\hat{\boldsymbol{\mu}}_{\text{sym}}^o$ :

LEMMA 5. For  $\hat{\boldsymbol{\mu}}^o = \hat{\boldsymbol{\mu}}_{\text{sym}}^o = \mathcal{M}(\mathbf{s}; \mathcal{P}(\mathbf{s}; \mathbf{1}))$ , the derivative matrix  $\hat{H} = d \log(\hat{\boldsymbol{\mu}}^o)/d\mathbf{s}$  is

$$(8.15a) \quad \hat{H} = D(\hat{\boldsymbol{\mu}}^{oo}/\hat{\boldsymbol{\mu}}^o)MD(1/\hat{\boldsymbol{\mu}}^{oo}) + [I - D(\hat{\boldsymbol{\mu}}^{oo}/\hat{\boldsymbol{\mu}}^o)MD(\mathbf{s}/\hat{\boldsymbol{\mu}}^{oo})]\hat{P}^{oo},$$

where

$$(8.15b) \quad \hat{\boldsymbol{\mu}}^{oo} = \mathcal{P}(\mathbf{s}; \mathbf{1}) \quad \text{and} \quad \hat{P}^{oo} = X[X'D(\hat{\boldsymbol{\mu}}^{oo})X]^{-1}X'.$$

The symmetrized version of the  $\hat{\boldsymbol{\mu}}_{\text{sef}}$  in Figure 3,  $\hat{\boldsymbol{\mu}}_{\text{sym}} = \mathcal{P}(\mathbf{s}; \mathcal{M}(\mathbf{s}; \hat{\boldsymbol{\mu}}_{\text{sef}}))$ , gave similar but somewhat rougher contours than those in the left panel of Figure 4. There is no compelling theoretical reason for preferring  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  or  $\hat{\boldsymbol{\mu}}_{\text{back}}$  to  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ , though they seem closer in structure to Bayes and maximum likelihood estimators. In practice, the specific choices of  $M$  and  $X$  seem more crucial to successful estimation than does the choice between  $\hat{\boldsymbol{\mu}}_{\text{sef}}$ ,  $\hat{\boldsymbol{\mu}}_{\text{sym}}$  or  $\hat{\boldsymbol{\mu}}_{\text{back}}$ .

**9. Remarks.** The following remarks apply to the indicated sections.

REMARK A (Section 2). There is an interesting connection between moment-matching and function-preserving properties of smoothers. For a smoother  $\hat{\boldsymbol{\mu}} = M\mathbf{s}$ , the condition  $\sum \hat{\boldsymbol{\mu}}_i = \sum s_i$  requires  $\mathbf{1}'M = \mathbf{1}'$  (where  $\mathbf{1}$  is a column of

ones), or equivalently  $M'\mathbf{1} = \mathbf{1}$ . Now most smoothers preserve constants so that  $M\mathbf{1} = \mathbf{1}$ . For symmetric  $M$  we see that matching the zeroth moment is the same as preserving the constant vector. However, most smoothers such as kernels are not symmetric, and hence will not match the zeroth moment exactly. Similarly, a smoother may preserve a vector  $\mathbf{t}$  (so  $M\mathbf{t} = \mathbf{t}$ ) without satisfying the moment-matching property  $M'\mathbf{t} = \mathbf{t}$ . In general, moment-matching is equivalent to function-preserving for the transpose of the smoother matrix.

REMARK B (Section 3). What is the true parameter  $\beta$  being estimated by  $\hat{\beta}$ ? Suppose that  $\mathbf{s} \sim \text{Po}_k(\boldsymbol{\mu})$ , but that  $\boldsymbol{\mu}$  is not necessarily of the form  $\boldsymbol{\mu}(\beta)$  in (3.1). From  $\boldsymbol{\mu}$  we determine  $\boldsymbol{\mu}^o = m(\boldsymbol{\mu})$  and then  $\beta$  the solution vector to the equations  $X'[\boldsymbol{\mu} - \boldsymbol{\mu}^o e^{X\beta}] = 0$ . Under reasonable conditions  $\hat{\beta}$  will be an asymptotically normal estimate for  $\beta$ , in the usual manner of an MLE. For instance if  $\boldsymbol{\mu} = n\boldsymbol{\pi}$ , with  $\boldsymbol{\pi}$  fixed and  $n \rightarrow \infty$ , and if the  $\boldsymbol{\mu}^o$  estimate is homogeneous,  $m(c\mathbf{s}) = cm(\mathbf{s})$ , then it is easy to show that

$$(9.1) \quad \sqrt{n}(\hat{\beta} - \beta) \rightarrow N_{p+1}(\mathbf{0}, A^{-1}BA^{-1}),$$

where, with  $Z' = X'[I - D(e^{X\beta})d\boldsymbol{\mu}^o/d\boldsymbol{\mu}]$ ,

$$(9.2) \quad A = \lim_{n \rightarrow \infty} X'D\left(\frac{\boldsymbol{\mu}^o e^{X\beta}}{n}\right)X \quad \text{and} \quad B = \lim_{n \rightarrow \infty} Z'D\left(\frac{\boldsymbol{\mu}}{n}\right)Z.$$

The bias of  $\hat{\beta}$  as an estimate of  $\beta$  is of order  $1/n$ .

REMARK C (Section 3). In case (2.13),  $\hat{\boldsymbol{\mu}}^o = M\mathbf{s}$ , the matrix  $Z$  in (3.4a) is a function of  $\hat{\beta}$  but not of  $\hat{\boldsymbol{\mu}}^o$ , say  $Z = Z(\hat{\beta})$ . If  $d\mathbf{s}$  is a  $K$  vector orthogonal to the columns of  $Z(\hat{\beta})$ , then

$$(9.3) \quad d\hat{\beta} = [X'\hat{D}X]^{-1}Z'(\hat{\beta})d\mathbf{s} = 0.$$

This implies that the level surfaces of constant  $\hat{\beta}$  value are  $K - (p + 1)$ -dimensional flat subspaces in the  $K$ -dimensional  $\mathbf{s}$  space. Flat level surfaces tend to make delta-method covariance estimates such as (3.6) more accurate.

REMARK D (Section 4). Here are the total Poisson deviances:

$$(9.4) \quad \text{Dev}(\mathbf{s}, \hat{\boldsymbol{\mu}}) = \sum_k 2 \left\{ s_k \log \left( \frac{s_k}{\hat{\mu}_k} \right) - (s_k - \hat{\mu}_k) \right\}$$

for four choices of  $\hat{\boldsymbol{\mu}}$ :  $\hat{\boldsymbol{\mu}}^o = M(2)\mathbf{s}$ ;  $\hat{\boldsymbol{\mu}}^1$  the SEF based on this  $\hat{\boldsymbol{\mu}}^o$  and  $X = (\mathbf{1}, \mathbf{i}, \mathbf{j})$ ;  $\hat{\boldsymbol{\mu}}^2$  the SEF (4.5);  $\hat{\boldsymbol{\mu}}^3$  the SEF (9.3):

$$(9.5) \quad \begin{array}{cccc} \hat{\boldsymbol{\mu}}^0 & \hat{\boldsymbol{\mu}}^1 & \hat{\boldsymbol{\mu}}^2 & \hat{\boldsymbol{\mu}}^3 \\ 247.1 & 240.3 & 226.1 & 220.3. \end{array}$$

The deviance decrease  $D(\mathbf{s}, \hat{\boldsymbol{\mu}}^2) - D(\mathbf{s}, \hat{\boldsymbol{\mu}}^3) = 5.77$  is much smaller than  $3.75^2$ , the square of the corresponding  $t$ -value in Table 3. The naive  $t$ -value, calculated using the standard error estimate from (3.7), is the right one for approximating the deviance decrease. In this case the naive  $t$ -value is  $0.000267/0.000111 = 2.74$ , predicting  $2.74^2 = 5.76$  for the deviance decrease.

REMARK E (Section 4). Results like (3.6) can be directly derived for continuous special exponential families (1.3) without going through the Poisson discretization argument. Suppose we estimate  $g_0(y)$  in (1.2) by a continuous version of (2.13),

$$(9.6) \quad \hat{g}_0(y) = \frac{1}{n} \sum_{i=1}^n M(y|y_i),$$

where, for any  $y_i$ ,  $M(y|y_i)$  is a distribution over  $\mathcal{Y}$  [not necessarily satisfying  $\int_{\mathcal{Y}} M(y|y_i) dy = 1$ ]. We define the SEF density estimate to be  $g_{\hat{\beta}}(y) = \hat{g}_0(y)\exp(\hat{\beta}_0 + t(y)\hat{\beta}_1)$ , where  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$  satisfies the “maximum likelihood” equations

$$(9.7) \quad \int_{\mathcal{Y}} t(y)g_{\hat{\beta}}(y) dy = \bar{t} = \frac{1}{n} \sum_{i=1}^n t(y_i),$$

as well as the constraint  $\int_{\mathcal{Y}} g_{\hat{\beta}}(y) dy = 1$ .

Define

$$(9.8) \quad z_i = (t(y_i) - \bar{t}) - \int_{\mathcal{Y}} (t(y) - \bar{t})\exp(\hat{\beta}_0 + (t(y) - \bar{t})\hat{\beta}_1)M(y|y_i) dy$$

and let  $\widehat{\text{Cov}}(t)$  indicate the covariance matrix of  $t(y)$  for the distribution on  $\mathcal{Y}$  corresponding to  $g_{\hat{\beta}}(y)$ . Then the continuous analog of (3.6a) is

$$(9.9) \quad \widehat{\text{Cov}}(\hat{\beta}_1) = \frac{1}{n} [\widehat{\text{Cov}}(t)]^{-1} \left[ \sum_{i=1}^n \frac{z'_i z_i}{n} \right] [\widehat{\text{Cov}}(t)]^{-1}.$$

(9.9) is the limit of (3.6a) as the discretization (2.1) becomes infinitely fine, after the superfluous parameter  $\hat{\beta}_0$  is removed.

In order to use (9.9), we need to evaluate the integrals over  $\mathcal{Y}$  involved in (9.7), (9.8) and  $\widehat{\text{Cov}}(t)$ . The discretization argument effectively does such integrals by summation over  $k = 1, 2, \dots, K$ . If  $\mathcal{Y}$  is high dimensional, we might prefer to work in the continuous mode, doing the integrals by some more efficient algorithm such as componentwise Simpson rules.

REMARK F (Section 5). We will often be interested in the probability vector

$$(9.10) \quad \hat{\pi} = \hat{\mu} / \hat{\mu}_+ = \hat{\mu} / n,$$

rather than in  $\hat{\mu}$  itself. Suppose that  $\hat{\mu}^o = m(\mathbf{s})$  is scale homogeneous,

$$(9.11) \quad m(c\mathbf{s}) = cm(\mathbf{s}) \quad (c > 0).$$

The  $\hat{\mu} = \text{SEF}(\mathbf{s}; m, X)$  is also scale homogeneous and  $\hat{\pi}(c\mathbf{s}) = \hat{\pi}(\mathbf{s})$ . Familiar homogeneity properties give

$$(9.12) \quad \mathbf{1}' \frac{d\hat{\mu}}{d\mathbf{s}} = \mathbf{1}', \quad \frac{d\hat{\mu}}{d\mathbf{s}} \mathbf{s} = \hat{\mu} \quad \text{and} \quad \frac{d\hat{\pi}}{d\mathbf{s}} = \frac{1}{\hat{\mu}_+} (I - \hat{\pi} \mathbf{1}') \frac{d\hat{\mu}}{d\mathbf{s}}.$$

Using (9.12) it is not difficult to show that

$$(9.13) \quad \overline{\text{Cov}}(\hat{\boldsymbol{\pi}}) \equiv \left( \frac{d\hat{\boldsymbol{\pi}}}{d\mathbf{s}} \right) \overline{D} \left( \frac{d\hat{\boldsymbol{\pi}}}{d\mathbf{s}} \right)' = \frac{1}{n^2} \left[ \overline{\text{Cov}}(\hat{\boldsymbol{\mu}}) - \frac{\hat{\boldsymbol{\mu}}\hat{\boldsymbol{\mu}}'}{n} \right]$$

and that the diagonal elements of  $\overline{\text{Cov}}(\hat{\boldsymbol{\pi}})$  satisfy

$$(9.14) \quad \overline{\text{trv}} \equiv n \sum_k \frac{\overline{\text{var}}(\hat{\pi}_i)}{\hat{\pi}_k} \equiv \overline{\text{TRV}} - 1.$$

Thus the comparisons in Table 4 remain valid for  $\overline{\text{trv}}$ . The results for the equivalent of  $\overline{\text{TRV}}$  are a little less neat.

REMARK G (Section 5). Let  $\mathbf{p} = \mathbf{s}/n$ , the vector of empirical probabilities, and define

$$(9.15) \quad \begin{aligned} \bar{d} &\equiv D(\mathbf{p}) = \overline{D}/n, & \hat{d} &\equiv D(\hat{\boldsymbol{\pi}}) = \hat{D}/n, \\ \hat{p} &\equiv X(X'\hat{d}X)^{-1}X' = n\hat{P}, & \hat{q} &\equiv \hat{d}^{-1} - \hat{p} = n\hat{Q}, & \hat{h} &\equiv n\hat{H}. \end{aligned}$$

[Under the homogeneity condition (5.14),  $\hat{h} = d \log(\hat{\boldsymbol{\pi}}^\circ)/d\mathbf{p}$ , for  $\hat{\boldsymbol{\pi}}^\circ \equiv \hat{\boldsymbol{\mu}}^\circ/n$ .] Then (5.9) can be written as

$$(9.16) \quad \overline{\text{TRV}} = \text{tr} \left[ \overline{D}\hat{p} + (\hat{h}\bar{d}\hat{h})'(\hat{d}\hat{q}\hat{d}) \right] \quad \text{or} \quad \widehat{\text{TRV}} = (p + 1) + \text{tr}(\hat{h}\bar{d}\hat{h})(\hat{d}\hat{q}\hat{d})$$

not depending on  $n$ . This shows that  $\overline{\text{TRV}}$  and  $\widehat{\text{TRV}}$  are  $O_p(1)$  as  $n \rightarrow \infty$ . See Remark K.

REMARK H (Section 5). Suppose that the discretization of  $\mathcal{Y}$  becomes infinitely fine, with  $K$  going to infinity and the  $k$ th cell  $\mathcal{Y}_k$  having volume  $\Delta_k$  going to zero. Under sufficient regularity conditions,  $\hat{\pi}(y_{(k)})/\Delta_k$  will approach  $\hat{g}(y_{(k)})$ , an estimate of the original continuous density with approximate variance

$$(9.17) \quad \overline{\text{var}}\{\hat{g}(y_{(k)})\} = \frac{\overline{\text{var}}\{\hat{\pi}(y_{(k)})\}}{\Delta_k^2}.$$

Then  $\overline{\text{trv}}$ , (9.14), will approach

$$(9.18) \quad n \int_{\mathcal{Y}} \frac{\overline{\text{var}}\{\hat{g}(y)\}}{\hat{g}(y)} dy = n \int_{\mathcal{Y}} \overline{\text{CV}}(y)^2 \hat{g}(y) dy,$$

with  $\overline{\text{CV}}(y) \equiv [\overline{\text{Var}}\{\hat{g}(y)\}]^{1/2}/\hat{g}(y)$ , a coefficient of variation measure for  $\hat{g}(y)$ . From Remark G we see that  $\overline{\text{CV}}(y)$  is typically  $O_p(1/\sqrt{n})$  as  $n \rightarrow \infty$ . The estimated average value for the coefficient of variation of SEF (4.5) was 0.15.

REMARK I (Section 5). A computationally more efficient expression than (5.9) is

$$(9.19) \quad \overline{\text{TRV}} = \text{tr} \overline{D} \left[ \hat{P} + \hat{H}'\hat{D}\hat{H} - \hat{H}'\hat{D}\hat{P}\hat{D}\hat{H} \right].$$

Each of the three terms in (9.19) can be evaluated using  $O(K^2)$  multiplications rather than  $O(K^3)$ . This makes values of  $K$  as large as 1000 practical. Letting  $\hat{G} = X'\hat{D}X$ , the  $O(K^2)$  computing formula is

$$(9.20) \quad \begin{aligned} \overline{\text{TRV}} &= \text{tr } \hat{G}^{-1}(X'\overline{DX}) + \|\overline{D}^{1/2}\hat{H}\hat{D}^{1/2}\|^2 \\ &\quad + \text{tr } \hat{G}^{-1}(X'\hat{D}\hat{H}\overline{D}^{1/2})(\overline{D}^{1/2}\hat{H}'\hat{D}X) \end{aligned}$$

( $\|a\|^2 \equiv \sum_i \sum_j a_{ij}^2$  for matrix  $a$ ), and similarly for  $\overline{\text{TRV}}$ , substituting  $\hat{D}$  for  $\overline{D}$ .

When  $\hat{\mu}^o = M\mathbf{s}$  the middle term in (9.20) equals  $\sum_k \hat{\mu}_k A_k$ , where

$$(9.21a) \quad A_k = \sum_j \hat{M}_{kj}^2 \tilde{s}_j / \left( \sum_j \tilde{M}_{kj} s_j \right)^2 \quad \left[ \tilde{M}_{kj} \equiv (\hat{\mu}_k / \hat{\mu}_k^o) M_{kj} \right],$$

and  $\tilde{s}_j$  equals  $s_j$  for  $\overline{\text{TRV}}$  or  $\hat{\mu}_j$  for  $\overline{\text{TRV}}$ . For the  $\overline{\text{TRV}}$  case  $A_k \leq 1$ , but  $A_k$  can blow up for  $\overline{\text{TRV}}$ . The calculations of Section 6 replace  $A_k$  with

$$(9.21b) \quad \min(A_k, 1)$$

in the  $\overline{\text{TRV}}$  case.

REMARK J (Section 6). The degrees of freedom formula (6.20) can be rewritten as

$$(9.22) \quad \overline{\text{DF}} = \text{tr } \overline{D} \left[ \hat{P} + \hat{H} - \hat{P}\hat{D}\hat{H} \right] \quad \text{or} \quad \overline{\text{DF}} = (p + 1) + \text{tr } \hat{D}(\hat{H} - \hat{P}\hat{D}\hat{H}).$$

The  $\hat{P}$  term corresponds to the  $e^{X\hat{\beta}}$  part of the SEF definition (3.2), while the  $\hat{H}$  term corresponds to the choice of the carrier  $\hat{\mu}^o$ . The  $\hat{P}\hat{D}\hat{H}$  subtraction term corrects for collinearity between the carrier and the exponential family terms. The degrees of freedom definition for  $\hat{\pi}$  rather than  $\hat{\mu}$ , as in Remark F, subtracts 1 from formula (6.20) or (9.22). It takes  $O(K^2)$  multiplications to compute  $\overline{\text{DF}}$  or  $\overline{\text{DF}}$  from either formula.

REMARK K (Section 6). Write  $\mathbf{s} = n\mathbf{p}$  as in Remark G and consider expression (6.15) for  $\text{DF}(\hat{\mu})$  as  $n \rightarrow \infty$  with  $\mathbf{p}$  fixed. Assuming  $m(\mathbf{s}) = cm(\mathbf{p})$ , (9.11), the approximation  $\overline{\text{DF}} = \hat{D}\hat{O}$  does not depend on  $n$  and so is  $O_p(1)$  just as in (9.16). Using higher-order expansions it is easy to show that

$$(9.23) \quad \text{DF}(\hat{\mu}) - \overline{\text{DF}} = O_p(1/n),$$

so at least in this sense  $\overline{\text{DF}}$  is a good approximation to  $\text{DF}(\hat{\mu})$ . The corresponding results hold for  $\overline{\text{DF}}$ ,  $\overline{\text{TRV}}$ , and  $\overline{\text{TRV}}$ .

REMARK L (Section 6). The Poisson residual degrees of freedom formula

$$(9.24) \quad \text{RDF}(\boldsymbol{\mu}) = K(\boldsymbol{\mu}) - 2\text{DF}(\boldsymbol{\mu}) + \text{TRV}(\boldsymbol{\mu})$$

is always nonnegative. To prove this we use (6.12b) to write

$$(9.25) \quad \text{RDF}(\boldsymbol{\mu}) = E\{\text{Dev}(\mathbf{s}, \hat{\boldsymbol{\mu}})\} - \text{Dev}(\boldsymbol{\mu}, \boldsymbol{\nu}),$$

and note that  $E\{(\mathbf{s}, \hat{\boldsymbol{\mu}})\} = (\boldsymbol{\mu}, \boldsymbol{\nu})$ . The proof is completed with the fact that the Poisson deviance  $\text{Dev}(\boldsymbol{\mu}, \boldsymbol{\nu})$  is a jointly convex function of  $(\boldsymbol{\mu}, \boldsymbol{\nu})$ .

The RDF estimates (6.21) are not necessarily nonnegative. They become so if we replace  $K(\boldsymbol{\mu})$  in (6.21) with the Taylor series estimates

$$(9.26) \quad \bar{K} = \text{tr } \bar{D}\hat{D}^{-1} \quad \text{or} \quad \hat{K} = \text{tr } \hat{D}\hat{D}^{-1} = (p + 1),$$

but these were poor approximations in our examples.

REMARK M (Section 8). We tried a “loess” smoother [Cleveland and Gross (1992)] for the galaxy data, in the form of a generalized additive Poisson model [Hastie and Tibshirani (1990)]. We used a degree 2 loess model, that is, one that fits second degree polynomials locally. This gave an expected deviance picture similar to Figure 5. This is not surprising given the similarity between the moment-matching and function-preserving properties mentioned in Remark A.

## REFERENCES

- AITKEN, M. (1993). Model choice in single samples from the exponential and double exponential families using the posterior Bayes factor. Technical report, Tel-Aviv Univ., Israel.
- CLEVELAND, W. C. and GROSS, E. (1992). Locally weighted regression. In *Statistical Models in S* (J. Chambers and T. Hastie, eds.). Wadsworth, Belmont, CA.
- COX, D. and REID, N. (1987). Parametric orthogonality and approximate conditional inference. *J. Roy. Statist. Soc. Ser. B* **49** 1–39.
- EFRON, B. (1986). How biased is the apparent error rate of a logistic regression? *J. Amer. Statist. Assoc.* **81** 461–470.
- EFRON, B. (1988). Logistic regression, survival analysis, and the Kaplan–Meier curve. *J. Amer. Statist. Assoc.* **83** 414–425.
- EFRON, B. (1992). Six questions raised by the bootstrap. In *Bootstrap Proceedings Volume* (R. LaPage and L. Billard, eds.). Wiley, New York.
- EFRON, B. and FELDMAN, D. (1991). Compliance as an explanatory variable in clinical trials (with comments and rejoinder). *J. Amer. Statist. Assoc.* **86** 9–25.
- EFRON, B. and PETROSIAN, V. (1992). A simple test of independence for truncated data with applications to redshift surveys. *Astrophysical Journal* **399** 345–352.
- EFRON, B. and TIBSHIRANI, R. J. (1993). *An Introduction to the Bootstrap*. Chapman and Hall, New York.
- GREEN, P. and SILVERMAN, B. (1994). *Nonparametric Regression and Generalized Linear Models*. Chapman and Hall, New York.
- HASTIE, T. and TIBSHIRANI, R. (1990). *Generalized Additive Models*. Chapman and Hall, New York.
- HJORT, N. and GLAD, I. (1995). Nonparametric density estimation with a parametric start. *Ann. Statist.* **23** 882–904.
- JONES, M. C. (1993). Kernel density estimation when the bandwidth is large. *Australian J. Statist.* **35** 319–326.
- KOOPERBERG, C. and STONE, C. (1991). A study of logspline density estimation. *Comput. Statist. Data Anal.* **12** 327–347.
- LEHMANN, E. (1983). *Theory of Point Estimation*. Wiley, New York.
- LINDSEY, J. (1974a). Comparison of probability distributions. *J. Roy. Statist. Soc. Ser. B* **36** 38–47.
- LINDSEY, J. (1974b). Construction and comparison of statistical models. *J. Roy. Statist. Soc. Ser. B* **36** 418–425.
- LINDSEY, J. and MERSCH, G. (1992). Fitting and comparing probability distributions with log linear models. *Comput. Statist. Data Anal.* **13** 373–384.

- LOH, E. and SPILLAR, E. (1986). Photometric redshift of galaxies. *Astrophysical Journal* **303** 154–161.
- MARDIA, K. V., KENT, J. T. AND BIBBY, J. M. (1979). *Multivariate Analysis*. Academic Press, New York.
- OLKIN, I. and SPIEGELMAN, C. (1987). A semiparametric approach to density estimation. *J. Amer. Statist. Assoc.* **82** 858–865.
- PIERCE, D. AND PETERS, D. (1992). Practical use of higher order asymptotics for multiparameter exponential families (with discussion). *J. Roy. Statist. Soc. Ser. B* **54** 701–737.
- STONE, C. (1994). The use of polynomial splines and their tensor products in multivariate function estimation. *Ann. Statist.* **22** 118–170.

DEPARTMENT OF STATISTICS  
STANFORD UNIVERSITY  
STANFORD, CALIFORNIA 94305

DEPARTMENT OF STATISTICS  
UNIVERSITY OF TORONTO  
TORONTO, ONTARIO  
CANADA