

# MIXTURES OF PROPORTIONAL HAZARDS REGRESSION MODELS

ORI ROSEN<sup>1</sup> AND MARTIN TANNER<sup>2\*</sup>

<sup>1</sup>*Department of Statistics, Northwestern University, 2006 Sheridan Road, Evanston, IL 60208, U.S.A.*

<sup>2</sup>*Department of Statistics, University of Pittsburgh, 139 University Place, Pittsburgh, PA 15260, U.S.A.*

## SUMMARY

This paper presents a mixture model which combines features of the usual Cox proportional hazards model with those of a class of models, known as mixtures-of-experts. The resulting model is more flexible than the usual Cox model in the sense that the log hazard ratio is allowed to vary non-linearly as a function of the covariates. Thus it provides a flexible approach to both modelling survival data and model checking. The method is illustrated with simulated data, as well as with multiple myeloma data. Copyright © 1999 John Wiley & Sons, Ltd.

## 1. INTRODUCTION

This paper describes a generalization of the Cox proportional hazards model, utilizing ideas from a class of models, known as mixtures-of-experts (ME). ME, originally introduced by Jacobs *et al.*<sup>1</sup> assume that the process generating the data can be decomposed into a set of subprocesses defined on possibly overlapping regions of the covariate space. As an illustration, we consider the ME generalization of the usual linear regression model,  $y_s = \beta^T \mathbf{x}_s + \varepsilon_s$ , where  $\varepsilon_s \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$ ,  $s = 1, \dots, n$ . The mixtures-of-experts model assumes that  $y_s$  is generated as follows. For each covariate  $\mathbf{x}_s$

- (i) a label  $r$  is chosen from a multinomial distribution with probability  $g_r^{(s)} = P(r|\mathbf{x}_s, V)$ , where  $V = [\mathbf{v}_1, \dots, \mathbf{v}_{I-1}]$  is the  $d \times (I-1)$  matrix of parameters underlying the multinomial distribution and  $d$  is the dimension of  $\mathbf{x}_s$ . The multinomial probabilities depend on  $\mathbf{x}_s$ ,  $s = 1, \dots, n$ , according to the generalized linear model

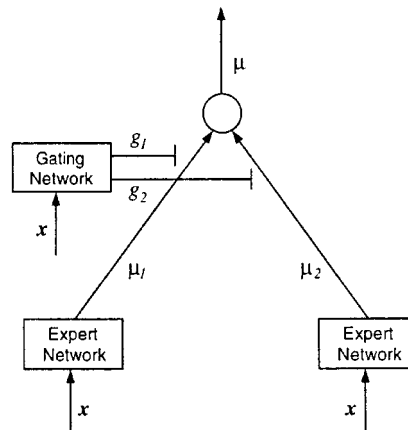
$$\log \frac{g_r^{(s)}}{g_I^{(s)}} = \mathbf{v}_r^T \mathbf{x}_s, \quad r = 1, \dots, I-1$$

where  $\sum_{r=1}^I g_r^{(s)} = 1$ ;

\* Correspondence to: Martin Tanner, Department of Statistics, Northwestern University, 2006 Sheridan Road, Evanston, IL 60208, U.S.A.

Contract/grant sponsor: NIH  
 Contract/grant number: CA35464

Contract/grant sponsor: Mellon Foundation

Figure 1. A mixtures-of-experts model (from Peng *et al.*<sup>2</sup>)

(ii) for the label,  $r$ , a response  $y^{(s)}$  is generated with probability

$$P(y_s | r, \mathbf{x}_s, \boldsymbol{\beta}_r, \sigma_r) = (2\pi\sigma_r^2)^{-1/2} \exp\left\{-\frac{(y_s - \boldsymbol{\beta}_r^T \mathbf{x}_s)^2}{2\sigma_r^2}\right\}.$$

In the above illustration, the expected conditional value of the response  $y_s$ , denoted by  $\mu_r^{(s)}$ , is given by  $\boldsymbol{\beta}_r^T \mathbf{x}_s$ . The total probability of generating  $y_s$  from  $\mathbf{x}_s$  is given by the mixture density

$$P(y_s | \mathbf{x}^{(s)}, \Theta) = \sum_{r=1}^I P(r | \mathbf{x}_s, V) P(y_s | r, \mathbf{x}_s, \boldsymbol{\beta}_r, \sigma_r)$$

where  $\Theta$  denotes all the parameters:  $\nu_1, \dots, \nu_{I-1}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_I, \sigma_1, \dots, \sigma_I$ . Assuming independently distributed data, the total probability of the data set  $\chi = \{(\mathbf{x}_s, y_s)\}_{s=1}^n$  is the product of  $n$  such densities with the likelihood given by

$$L(\Theta | \chi) = \prod_{s=1}^n \sum_{r=1}^I g_r^{(s)} \frac{1}{\sqrt{(2\pi\sigma_r^2)}} \exp\left\{-\frac{1}{2\sigma_r^2} (y_s - \boldsymbol{\beta}_r^T \mathbf{x}_s)^2\right\}.$$

Thus, the ME model can be viewed as a mixture model where the mixing weights and mixed distributions depend on the covariate  $\mathbf{x}$ .

ME can also be viewed as a tree-structured model. Figure 1 presents a graphical representation of ME in the case  $I = 2$ . The model consists of  $I$  modules, each referred to as an *expert*. These experts approximate the data within each region of the covariate space; expert  $r$  maps its input  $\mathbf{x}$  to an output  $\mu_r$ . It is assumed that different experts are appropriate in different regions of the covariate space. Thus, the model attempts to solve problems using a 'divide-and-conquer' strategy; complex tasks are decomposed into simpler subtasks. Consequently the model requires a module, referred to as a *gating network*, that identifies for an input  $\mathbf{x}_s$  the expert or blend of experts whose output is most likely to approximate the corresponding response  $y_s$ . The gating network outputs are the probabilities  $g_r$  that weight the contributions of the various experts. The total output of the model is  $\mu = \sum_{r=1}^I g_r \mu_r$ . In general, it is assumed that the conditional distribution of  $y_s$  given  $\mathbf{x}_s$  and the label  $r$ , is a member of the exponential family. Thus, in general,

ME provide a richer class of models than ordinary generalized linear models. In the linear regression case, in particular, they can be viewed as an alternative to non-parametric regression.

Jordan and Jacobs<sup>3</sup> present an extension of ME, called hierarchical mixtures-of-experts (HME). Peng *et al.*<sup>2</sup> discuss full Bayesian inference in ME and HME with an application to speech recognition. A detailed description of Bayesian inference in ME and HME with applications to regression and classification is given in Jacobs *et al.*<sup>4</sup> A Bayesian approach to model selection in ME and HME is proposed by Jacobs *et al.*<sup>5</sup>

The paper by Kooperberg *et al.*<sup>6</sup> presents an approach to hazard regression modelling using splines. Other spline approaches to the non-parametric modelling of survival time on covariates include the work by Hastie and Tibshirani,<sup>7</sup> O'Sullivan,<sup>8</sup> Sleeper and Harrington<sup>9</sup> and Gray.<sup>10</sup> LeBlanc and Crowley<sup>11</sup> develop an approach motivated by the CART methodology of Breiman *et al.*<sup>12</sup>

The outline of the paper is as follows. In Section 2 we present our method, preceding it with necessary background on the Cox model. The proposed method is based on a mixture embedded within the usual Cox partial likelihood, resulting in a more flexible model. A method for selecting the number of components in this mixture is also described. Section 3 illustrates the method with simulated, as well as real data. Section 4 concludes with a discussion.

## 2. A MIXTURE OF COX EXPERTS

### 2.1. The model

Prior to presenting our method, we describe the familiar set-up of the Cox model, which will be used in the sequel. The proportional hazards model<sup>13</sup> assumes that

$$\lambda(t; x_1, \dots, x_p) = \lambda_0(t) \exp \left( \sum_{i=1}^p \beta_i x_i \right) \quad (1)$$

where  $\lambda(t; x_1, \dots, x_p)$  is the hazard at time  $t$ , given covariate values  $\mathbf{x} = (x_1, \dots, x_p)^T$ , and  $\lambda_0$  is an unspecified baseline hazard function. The available data are of the form  $(y_1, \mathbf{x}_1, \delta_1), \dots, (y_n, \mathbf{x}_n, \delta_n)$ , where the survival time  $y_i, i = 1, \dots, n$ , is complete if  $\delta_i = 1$  and censored if  $\delta_i = 0$ . The distinct failure times are denoted by  $t_1 < t_2 < \dots < t_k$  with  $d_i$  failures at  $t_i, i = 1, \dots, k$ . The parameters  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  are estimated without specification of  $\lambda_0(t)$  by maximizing the partial likelihood

$$\prod_{i=1}^k \frac{\exp(\sum_{j \in \mathcal{D}(t_i)} \boldsymbol{\beta}^T \mathbf{x}_j)}{(\sum_{l \in \mathcal{R}(t_i)} e^{\boldsymbol{\beta}^T \mathbf{x}_l})^{d_i}}. \quad (2)$$

In (2),  $\mathcal{D}(t_i)$  is the set of indices of the failures at  $t_i$ , and  $\mathcal{R}(t_i)$  is the set of indices of the subjects at risk at time  $t_i^-$ . This approximation for ties is due to Peto<sup>14</sup> and Breslow.<sup>15</sup> For a fuller discussion of ties, see Kalbfleisch and Prentice.<sup>16</sup> The construction of the partial likelihood can be justified by arguments given in Kalbfleisch and Prentice.<sup>16</sup>

One of the goals of fitting a survival model is to estimate the survivor function of a subject with covariate values  $\mathbf{x}_0$ . Let  $\Lambda_0(t) = \int_0^t \lambda_0(u) du$  be the cumulative baseline hazard. Based on the

maximum partial likelihood estimate  $\hat{\beta}$ , Breslow<sup>17</sup> estimates  $\Lambda_0(t)$ , for  $t$  in  $[t_i, t_{i+1})$ ,  $i = 1, \dots, k-1$ , by

$$\hat{\Lambda}_0(t) = \sum_{j=1}^i \frac{d_j}{\sum_{l \in \mathcal{R}(t_j)} e^{\hat{\beta}^T \mathbf{x}_l}}. \quad (3)$$

Using (3), the survivor function  $S(t; \mathbf{x}_0)$  of a subject with covariate values  $\mathbf{x}_0$  is then estimated by

$$\hat{S}(t; \mathbf{x}_0) = \exp(-\hat{\Lambda}_0(t) e^{\hat{\beta}^T \mathbf{x}_0}). \quad (4)$$

Our approach to the construction of a partial likelihood for the ME generalization is to adapt a derivation of the original proportional hazards model, presented in Miller.<sup>18</sup> Given the  $I$  experts, the probability of a failure in  $[t_i, t_i + \Delta]$ , for small  $\Delta$ , given  $\mathcal{R}(t_i)$  is approximately

$$\sum_{l \in \mathcal{R}(t_i)} \sum_{r=1}^I g_r^{(l)} e^{\beta_r^T \mathbf{x}_l} \lambda_0(t_i) \Delta \quad (5)$$

where  $g_r^{(l)}$  depends on  $\mathbf{x}_l$  according to the generalized linear model

$$\log \frac{g_r^{(l)}}{g_I^{(l)}} = \mathbf{v}_r^T \mathbf{x}_l, \quad r = 1, \dots, I-1 \quad (6)$$

and  $\sum_{r=1}^I g_r^{(l)} = 1$ . It then follows that the probability of a failure of a subject with covariates  $\mathbf{x}_j$  at time  $t_i$ , given one failure in  $\mathcal{R}(t_i)$  at time  $t_i$  is

$$\frac{\sum_{r=1}^I g_r^{(j)} e^{\beta_r^T \mathbf{x}_j}}{\sum_{l \in \mathcal{R}(t_i)} \sum_{r=1}^I g_r^{(l)} e^{\beta_r^T \mathbf{x}_l}}.$$

Taking the product of these conditional probabilities (with an adjustment for ties) gives the partial likelihood for the mixture of experts

$$\text{PL}_g = \prod_{i=1}^k \prod_{j \in \mathcal{D}(t_i)} \frac{\sum_{r=1}^I g_r^{(j)} e^{\beta_r^T \mathbf{x}_j}}{\sum_{l \in \mathcal{R}(t_i)} \sum_{r=1}^I g_r^{(l)} e^{\beta_r^T \mathbf{x}_l}}. \quad (7)$$

As described in Section 1, the idea underlying ME is the ‘divide-and-conquer’ principle, according to which different experts are appropriate in different regions of the covariate space. In the application of ME to proportional hazards regression, this idea is reflected in equation (5). In this equation the probability of a failure is written as a mixture model, where the mixing weights  $g_r^{(l)}$  depend on the covariate vector  $\mathbf{x}_l$ . Analogous to the tree structure in Figure 1, each expert provides an assessment of the conditional probability of a failure in the infinitesimal interval  $[t_i, t_i + \Delta]$ . These assessments are then merged by the gating network via the weighting scheme given in (6). Note that when  $I = 1$ , (7) reduces to the ordinary partial likelihood (2).

## 2.2. Inferential issues

The parameters of the mixture model are estimated by maximizing  $l = \log \text{PL}_g$ . Approximate standard errors can be obtained utilizing the Hessian of  $l$ . The first and second derivatives of  $l$  are given in the Appendix.

Next, we use the Hessian to derive approximate pointwise confidence intervals for  $\eta^{(s)} = \sum_{r=1}^I g_r^{(s)} \beta_r^T \mathbf{x}_s$ . A graph of  $\hat{\eta}^{(s)} = \sum_{r=1}^I \hat{g}_r^{(s)} \hat{\beta}_r^T \mathbf{x}_s$  versus  $x_s^i$ , the  $i$ th component of  $\mathbf{x}_s$ , fixing all the other components of  $\mathbf{x}_s$ , provides a useful diagnostic for the proportional hazards model.

Non-linearity in this plot suggests a lack of fit of the proportional hazards model. The pointwise standard error of  $\hat{\eta}^{(s)}$  is obtained via the delta method. Specifically, the gradient of  $\eta^{(s)}$  is given by

$$\nabla^T \eta^{(s)} = \left( \frac{\partial \eta^{(s)}}{\partial \mathbf{v}_1^T}, \dots, \frac{\partial \eta^{(s)}}{\partial \mathbf{v}_{I-1}^T}, \frac{\partial \eta^{(s)}}{\partial \boldsymbol{\beta}_1^T}, \dots, \frac{\partial \eta^{(s)}}{\partial \boldsymbol{\beta}_I^T} \right)$$

where  $\partial \eta^{(s)} / \partial \mathbf{v}_l = g_l^{(s)} [(1 - g_l^{(s)}) \boldsymbol{\beta}_l^T \mathbf{x}_s - \sum_{r \neq l} g_r^{(s)} \boldsymbol{\beta}_r^T \mathbf{x}_s] \mathbf{x}_s'$ ,  $l = 1, \dots, I-1$ ,  $\partial \eta^{(s)} / \partial \boldsymbol{\beta}_m = g_m^{(s)} \mathbf{x}_s$ ,  $m = 1, \dots, I$ , and  $\mathbf{x}_s' = (1, \mathbf{x}_s^T)^T$ . It then follows that approximate  $(1 - \alpha)$  level pointwise confidence intervals for  $\eta^{(s)}$  are formed by

$$\hat{\eta}^{(s)} \pm z_{1-\alpha/2} (\nabla^T \hat{\eta}^{(s)} (-\hat{H}_l)^{-1} \nabla \hat{\eta}^{(s)})^{1/2} \quad (8)$$

where  $H_l$  is the Hessian matrix of  $l$ , and  $z_{1-\alpha/2}$  is the  $(1 - \alpha/2)$ th percentile of the standard normal distribution.

For estimating the baseline cumulative hazard and the survivor, we use weighted versions of formulae (3) and (4):

$$\begin{aligned} \hat{\Lambda}_0(t) &= \sum_{j=1}^i \frac{d_j}{\sum_{l \in \mathcal{R}(t_j)} \sum_{r=1}^I \hat{g}_r^{(l)} e^{\hat{\boldsymbol{\beta}}_r^T \mathbf{x}_i}} \\ \hat{S}(t; \mathbf{x}_0) &= \exp(-\hat{\Lambda}_0(t) \sum_{r=1}^I \hat{g}_r^{(0)} e^{\hat{\boldsymbol{\beta}}_r^T \mathbf{x}_0}). \end{aligned} \quad (9)$$

For the examples in Section 3, the function `nlimb` of S-PLUS was used to maximize  $l$ . In order to avoid local maxima for the examples in this paper, the maximization process was repeated 20 times with random starting values. Thus, the presented estimates represent the maximizer over the 20 maximizations. The use of multiple starting points is quite standard in such problems and not terribly onerous. For the examples in this paper, the algorithm converged fairly quickly, and, for the most part, the global maximum was not hard to find. For example, in the first example of Section 3, 14 out of the 20 maximizations resulted in the global maximum. However, as noted by the referee, multiple modes can be more of an issue as  $I$  is made larger. All computations were performed on a Sun Ultra 1 workstation.

### 2.3. Selecting the number of experts

An important issue in the specification of the mixture of proportional hazards models is the selection of the number of experts. In the context of mixtures of generalized linear models, this issue is addressed by Jacobs *et al.*<sup>5</sup> In the context of mixtures of marginal models, this issue is addressed by Rosen *et al.*<sup>19</sup> To develop an approach to selecting the number of experts in the mixture of proportional hazards models context, we again draw upon the motivating logic in Miller.<sup>18</sup>

According to expert  $r$ , the (unconditional) probability of a failure in the interval  $[t_i, t_i + \Delta]$  is given by  $g_r^{(l)} e^{\boldsymbol{\beta}_r^T \mathbf{x}_i} \lambda_0(t_i) \Delta$ . The total (unconditional) probability of a failure in this infinitesimal interval according to all  $I$  experts is given by  $\sum_{j=1}^I g_j^{(l)} e^{\boldsymbol{\beta}_j^T \mathbf{x}_i} \lambda_0(t_i) \Delta$ . Thus, to motivate a measure of the value of expert  $r$ , we consider the ratio

$$\frac{g_r^{(l)} e^{\boldsymbol{\beta}_r^T \mathbf{x}_i} \lambda_0(t_i) \Delta}{\sum_{j=1}^I g_j^{(l)} e^{\boldsymbol{\beta}_j^T \mathbf{x}_i} \lambda_0(t_i) \Delta} = \frac{g_r^{(l)} e^{\boldsymbol{\beta}_r^T \mathbf{x}_i}}{\sum_{j=1}^I g_j^{(l)} e^{\boldsymbol{\beta}_j^T \mathbf{x}_i}} = h_r^{(l)}.$$

Clearly, if expert  $r$  plays little or no role in modelling the response surface in the region of  $\mathbf{x}_l$ , then  $h_r^{(l)}$  will be negligible. However, if expert  $r$  takes the lead role in modelling the response surface in the region of  $\mathbf{x}_l$ , then  $h_r^{(l)}$  may approach unity. To examine the overall value or worth of expert  $r$ , we then average these  $h_r^{(l)}$ s over all the distinct failure times, that is, we consider the worth index ( $w_r$ ) of expert  $r$  as  $(1/k) \sum_{l=1}^k h_r^{(l)}$ . If the worth indices for each of the experts are all of similar magnitudes, then this suggests that the current model may be too small, and that a model with additional experts should be considered. Alternatively, if the worth index for an expert is small relative to that of other experts, then there is evidence to suggest that this expert can be pruned from the model. A useful more formal criterion suggested by Jacobs *et al.*<sup>5</sup> to determine the number of experts is the minimum number of experts with the largest worth indices for which the sum of their worth indices exceeds some critical value  $\kappa$ , that is

$$\min \left\{ J : \sum_{r=1}^J w_{(r)} \geq \kappa \right\} \quad (10)$$

where  $w_{(1)} \geq w_{(2)} \dots$ . All other experts are pruned from the model. Jacobs *et al.*<sup>5</sup> take  $\kappa = 0.8$ . Although the 0.8 cut-off criterion is arbitrary, based on empirical work, Jacobs *et al.*<sup>5</sup> have found this value to work well in practice. For example, suppose that the worth indices for the four experts that comprise a model have the values 0.4, 0.3, 0.2 and 0.1. According to (10), this suggests that three experts are required ( $0.4 + 0.3 + 0.2 = 0.9$ ), and one expert can be pruned. Care, however, has to be taken; if an expert is used for a specific type of data which occurs relatively infrequently in the data set, then dropping the expert may not be optimal. If the expert worth indices are similar in magnitude, for example, they are all within 10 per cent of  $1/I$ , where  $I$  is the number of experts, then one may wish to add experts to the model.

### 3. EXAMPLES

We illustrate the mixture model with a simulated data set, as well as using multiple myeloma data found in Krall *et al.*<sup>20</sup> and analysed by Collett.<sup>21</sup>

#### 3.1. Simulated data

The data considered here were simulated from the Weibull model

$$\log T_i = |x_i| + 0.2\varepsilon_i, \quad i = 1, \dots, 200 \quad (11)$$

where  $T_i$  is the survival time of the  $i$ th subject,  $x_i$  is the corresponding covariate, taking equally spaced values on a grid between  $-1$  and  $1$ , and the  $\varepsilon_i$ 's are i.i.d. Gumbel variates. The censoring times were simulated from an exponential random variable, resulting in 10 per cent non-informative censoring. Model (11) implies that the survival times follow a Weibull distribution with hazard function

$$\lambda(t; x) = 5t^4 e^{-5|x|}. \quad (12)$$

Note that this model can also be written as

$$\lambda(t; x) = 5t^4(g_1 e^{5x} + g_2 e^{-5x})$$

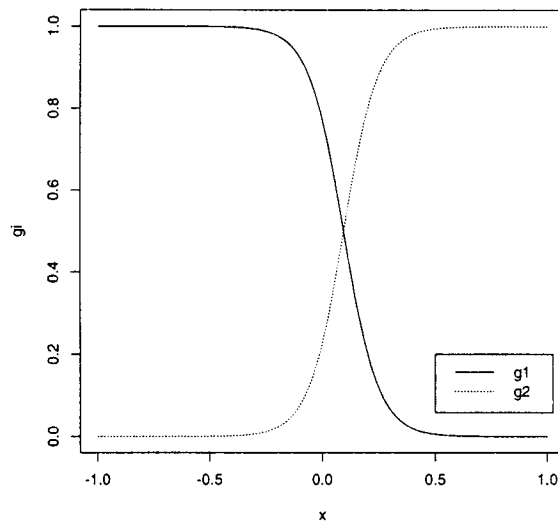


Figure 2. Mixing probabilities versus covariate values

where

$$g_1 = \begin{cases} 1 & x < 0 \\ 0 & x \geq 0 \end{cases}$$

and  $g_2 = 1 - g_1$ . Note also that the true log hazard ratio for this model, is  $-5|x|$ . Equations (1) and (12) imply that for a given  $x$ , the underlying true  $\beta$  satisfies  $\beta = -5 \operatorname{sgn}(x)$ .

Model (7) with two experts ( $I = 2$ ) was fitted to these data, resulting in the estimates  $\hat{\nu}_1^T = (\hat{\nu}_{10}, \hat{\nu}_{11}) = (1.22, -12.86)$ ,  $\hat{\beta}_1 = 4.78$  and  $\hat{\beta}_2 = -5.12$ , with standard errors  $\operatorname{SE}(\hat{\nu}_{10}) = 0.79$ ,  $\operatorname{SE}(\hat{\nu}_{11}) = 2.47$ ,  $\operatorname{SE}(\hat{\beta}_1) = 0.48$  and  $\operatorname{SE}(\hat{\beta}_2) = 0.69$ .

Figure 2 shows how the covariate space was segmented between the two experts. In particular, the horizontal axis gives the covariate values and the vertical axis gives the values of the mixing probabilities  $\hat{g}_1^{(s)} = \exp(1.22 - 12.86x_s)/(1 + \exp(1.22 - 12.86x_s))$  and  $\hat{g}_2^{(s)} = 1 - \hat{g}_1^{(s)}$ ,  $s = 1, \dots, 200$ . While  $\hat{g}_1^{(s)}$  is large for  $x_s < 0$ ,  $\hat{g}_2^{(s)}$  is large for  $x_s \geq 0$ . The mixing probabilities thus dictate that expert 2 assumes the main responsibility for modelling the response corresponding to  $x_s \geq 0$ , whereas for  $x_s < 0$  expert 1 is the primary expert.

To assess whether two experts are sufficient in this case, we used the technique described in Section 2.3. In particular, a model with three experts was fit to the data, resulting in the worth indices 0.20, 0.21 and 0.59. Criterion (10), thereby, suggests that a model with two experts is sufficient. In fact, the fitted log hazard ratio using three experts is quite similar to the fit using two experts, and appealing to Ockham's Razor we prefer the two-expert model.

Figure 3 shows: (a) the true log hazard ratio  $-5|x|$ ; (b)  $\hat{\eta} = \sum_{r=1}^2 \hat{g}_r \hat{\beta}_r x$ ; (c)  $\hat{\beta}^{\text{PH}} x$ , based on the proportional hazards model, where  $\hat{\beta}^{\text{PH}} = -0.0204$ ; (d) Hastie and Tibshirani's estimate (see details below); and (e) approximate 95 per cent pointwise confidence intervals for  $\eta$ , based on the mixture model. The mixture estimate  $\hat{\eta}$  is seen to recover the true log hazard ratio. The

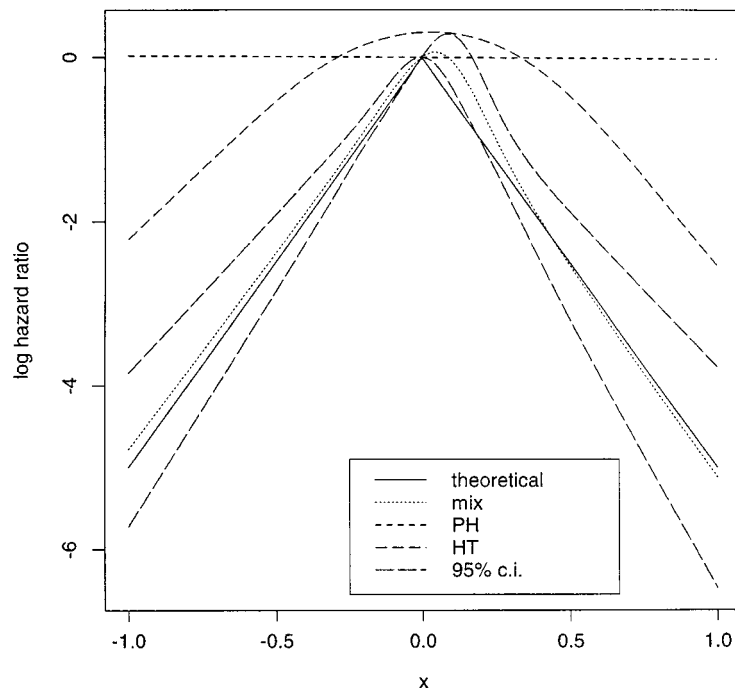


Figure 3. Theoretical log hazard ratio, mixture estimate of log hazard ratio, proportional hazards estimate of log hazard ratio, Hastie and Tibshirani's estimate and approximate 95 per cent pointwise confidence bands for the log hazard ratio

proportional hazards estimate clearly mismodels the data over the whole range of  $x$ . Of course, had the data analyst known that the log hazard ratio followed the absolute value of  $x$ , then it would have been appropriate to fit this form of the model. We assume here that the true structure is unknown to the data analyst. Moreover, while a careful residual analysis may uncover this structure, the point of this example is to illustrate that the mixture of Cox experts model can automatically detect this non-linear (and non-differentiable) structure. Hastie and Tibshirani's estimate was computed using statlib's S-PLUS function `coxgam()`, which implements the additive proportional hazards model described in Hastie and Tibshirani.<sup>7</sup> As can be seen, this estimate is also biased. Note that the true log hazard ratio is non-differentiable at  $x = 0$ , whereas the additive proportional hazards model assumes that the log hazard ratio is a sum of smooth functions, each differentiable and having an absolutely continuous first derivative.<sup>7</sup> Non-differentiable response surfaces arise in change-point models, which play an important role in cancer epidemiology.<sup>22</sup> The confidence bands based on (8) contain the log hazard ratio over most of the covariate range.

Figure 4 depicts the theoretical survivor function  $S(t; x) = \exp(-t^5 e^{-5|x|})$ , computed at  $x = 0$ , along with estimates based on the proportional hazards model and the mixture model ((4) and (9), respectively). It is evident that the survivor estimate based on the proportional hazards model is severely biased, whereas the bias of the mixture estimate is much smaller.



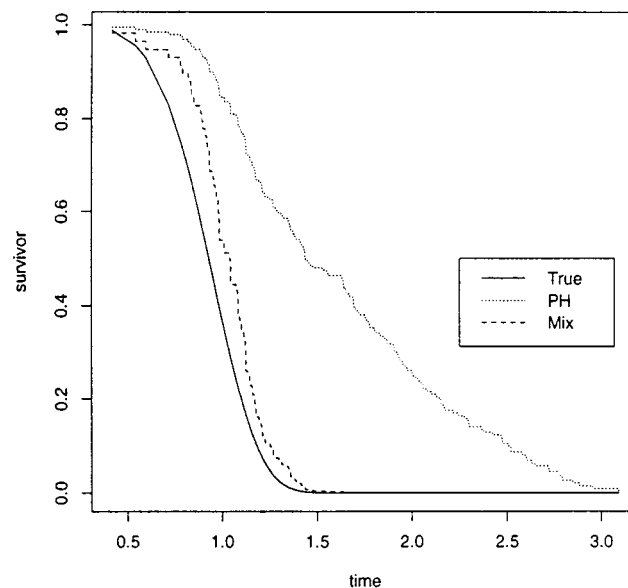


Figure 4. Theoretical survivor function and estimates

### 3.2. Multiple myeloma data

The multiple myeloma data relate to 48 patients at ages 50–80. Some of these patients were still alive when the study was completed, so their survival times are right-censored. A number of covariates were recorded for each patient. We fit a mixture of two experts ( $I = 2$ ) with the two covariates, serum haemoglobin ( $H$ ) and blood urea nitrogen ( $B$ ). Collett<sup>21</sup> found the proportional hazards model containing these two covariates to be the most satisfactory among several models.

Figure 5 shows, on the left-hand side, plots of  $\hat{\eta}_H$  against  $H$ , where  $B$  is fixed at its sample lower quartile, median and upper quartile, resulting in three of the displayed lines. Also shown is a graph of  $\hat{\beta}_H^{\text{PH}}H + \hat{\beta}_B^{\text{PH}}B$ , based on the proportional hazards model, where  $B$  is fixed at its sample upper quartile. On the right, the analogous quantities are plotted against  $B$ , where  $H$  is fixed at its sample quartiles. The point estimates in Figure 5 suggest that  $\eta$  may not be entirely additive in  $H$  and  $B$ , and that  $\eta$  may have non-linearities. A definitive conclusion, however, requires the use of confidence bands. Figure 6 shows for each of the two covariates: (a) one of the three lines from Figure 5, based on the mixture model, computed when the other covariate is fixed at its sample upper quartile; (b) approximate 95 per cent level confidence bands; (c) the estimate based on the standard proportional hazards model. The confidence bands on the left-hand side plot show that the bend in  $\hat{\eta}$  is probably due to chance variation (note the small number of observations to the left of  $H = 7.5$ ). The same is also true for the other two lines (not shown). Our analysis confirms that one expert, that is, the proportional hazards model, yields a satisfactory fit to the data. Thus, the ME approach not only provides a flexible scheme to model survival data, it also assists in model checking.

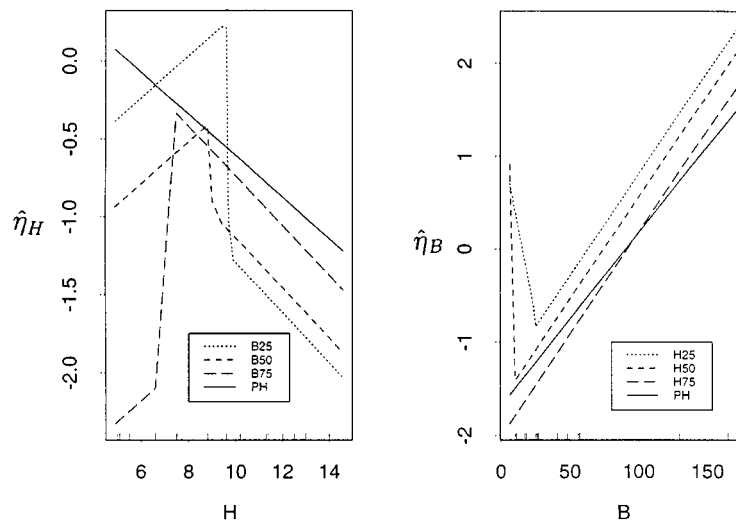


Figure 5.  $\hat{\eta}_H$  versus  $H$ , on the left-hand side, where  $B$  is fixed at its sample quartiles. On the right,  $\hat{\eta}_B$  versus  $B$ , where  $H$  is fixed at its sample quartiles. Also shown are  $\hat{\beta}_H^{\text{PH}}H + \hat{\beta}_B^{\text{PH}}B$ , where  $H$  and  $B$  are fixed in turn at their upper quartile

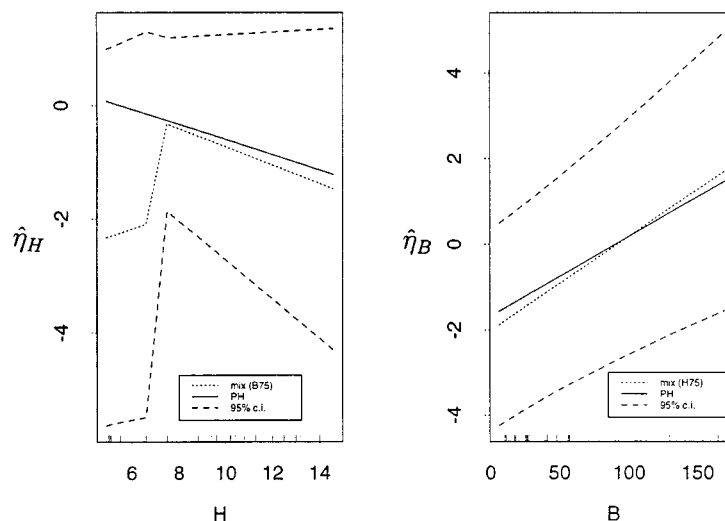


Figure 6. Approximate 95 per cent level confidence intervals, based on the mixture model, superimposed on  $\hat{\eta}_H$  and  $\hat{\beta}_H^{\text{PH}}H + \hat{\beta}_B^{\text{PH}}B$  versus  $H$  on the left-hand side, and the analogous quantities versus  $B$ , on the right. The covariates are fixed in turn at their upper quartile

Figure 7 shows estimates of the survivor function, based on the proportional hazards model and the mixture model, computed at the sample means  $\bar{H}$  and  $\bar{B}$ . There is a close agreement between the two estimates. The discrepancies between the two estimation methods, apparent in Figure 5 and Figure 6, have only a negligible effect on the corresponding survivor estimates.

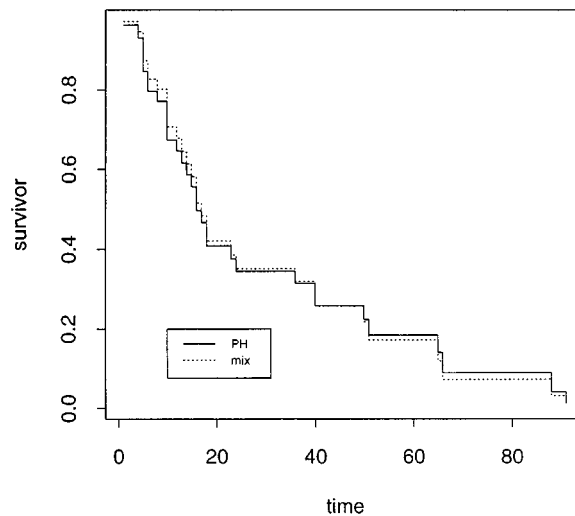


Figure 7. Survivor estimates

#### 4. CONCLUSION

Since any statistical model can be considered an ‘expert’, the mixtures-of-experts paradigm provides a generic approach for relaxing the parametric assumptions in the systematic component of *any* statistical model. Jacobs *et al.*<sup>4</sup> consider the ‘non-parameterization’ of the generalized linear model. Rosen *et al.*<sup>19</sup> consider the marginal model. This paper presents a mixture model which combines features of the usual proportional hazards model and those of mixtures-of-experts. In particular, the proportional-hazards assumption is retained, but the underlying log hazard-ratio can take a more general form than the linear combination of predictors inherent in the usual proportional hazards model. Thus, the mixtures-of-experts model provides a flexible approach to modelling survival data and to model checking.

Data analysts have at their disposal an extensive armamentarium for the analysis of censored data. The methodology of Hastie and Tibshirani<sup>7</sup> can address data with underlying smooth additive functions. The approach of LeBlanc and Crowley<sup>11</sup> can approximate piecewise constant functions that may not be additive. The methods of Kooperberg *et al.*<sup>6</sup> fall somewhere in between. In the examples considered in this paper, we see that the mixtures-of-experts approach can be used to model non-differentiable surfaces in a possibly non-additive manner. Jiang and Tanner<sup>23</sup> present formal regularity conditions for which a mixtures-of-experts model applied in the generalized linear model context is consistent, as well as present theoretical convergence rates. Under investigation are parallel results for mixtures of Cox experts.

#### APPENDIX: FIRST AND SECOND DERIVATIVES

We first define the following quantities ( $\mathbf{x}' = (1, \mathbf{x}^T)^T$ ):

$$c_j = \sum_{r=1}^I g_r^{(j)} \exp(\boldsymbol{\beta}_r^T \mathbf{x}_j)$$

$$\begin{aligned}
\mathbf{a}_{mj} &= \frac{\partial c_j}{\partial \boldsymbol{\beta}_m} = g_m^{(j)} \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) \mathbf{x}_j \\
A_{mj} &= \frac{\partial \mathbf{a}_{mj}}{\partial \boldsymbol{\beta}_m} = g_m^{(j)} \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) \mathbf{x}_j \mathbf{x}_j^T \\
\mathbf{b}_{mj} &= \frac{\partial c_j}{\partial \mathbf{v}_m} = g_m^{(j)} [(1 - g_m^{(j)}) \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) - \sum_{i \neq m} g_i^{(j)} \exp(\boldsymbol{\beta}_i^T \mathbf{x}_j)] \mathbf{x}'_j \\
B_{mj} &= \frac{\partial \mathbf{b}_{mj}}{\partial \mathbf{v}_m} = g_m^{(j)} (1 - 2g_m^{(j)}) [(1 - g_m^{(j)}) \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) - \sum_{i \neq m} g_i^{(j)} \exp(\boldsymbol{\beta}_i^T \mathbf{x}_j)] \mathbf{x}'_j \mathbf{x}'_j{}^T \\
D_{mnj} &= \frac{\partial \mathbf{b}_{mj}}{\partial \mathbf{v}_n} = g_m^{(j)} g_n^{(j)} \left[ (2g_m^{(j)} - 1) \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) - (1 - g_n^{(j)}) \exp(\boldsymbol{\beta}_n^T \mathbf{x}_j) \right. \\
&\quad \left. + \sum_{i \neq m} g_i^{(j)} \exp(\boldsymbol{\beta}_i^T \mathbf{x}_j) + \sum_{i \neq m, n} g_i^{(j)} \exp(\boldsymbol{\beta}_i^T \mathbf{x}_j) \right] \mathbf{x}'_j \mathbf{x}'_j{}^T \\
F_{mj} &= g_m^{(j)} (1 - g_m^{(j)}) \exp(\boldsymbol{\beta}_m^T \mathbf{x}_j) \mathbf{x}'_j \mathbf{x}_j^T \\
G_{mnj} &= g_n^{(j)} g_m^{(j)} \exp(\boldsymbol{\beta}_n^T \mathbf{x}_j) \mathbf{x}_j \mathbf{x}'_j{}^T.
\end{aligned}$$

The first and second derivatives of  $l = \log \text{PL}_g$  are given as functions of the above expressions:

$$\begin{aligned}
\frac{\partial l}{\partial \mathbf{v}_m} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{\mathbf{b}_{mj}}{c_j} - \frac{\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml}}{\sum_{l \in \mathcal{R}_i} c_l} \right\} \\
\frac{\partial l}{\partial \boldsymbol{\beta}_m} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{\mathbf{a}_{mj}}{c_j} - \frac{\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml}}{\sum_{l \in \mathcal{R}_i} c_l} \right\} \\
\frac{\partial^2 l}{\partial \boldsymbol{\beta}_m \partial \boldsymbol{\beta}_m^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{c_j A_{mj} - \mathbf{a}_{mj} \mathbf{a}_{mj}^T}{c_j^2} - \frac{(\sum_{l \in \mathcal{R}_i} c_l) (\sum_{l \in \mathcal{R}_i} A_{ml}) - (\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml}) (\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\} \\
\frac{\partial^2 l}{\partial \boldsymbol{\beta}_m \partial \boldsymbol{\beta}_n^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{-\mathbf{a}_{mj} \mathbf{a}_{nj}^T}{c_j^2} + \frac{(\sum_{l \in \mathcal{R}_i} \mathbf{a}_{nl}) (\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\}, \quad m \neq n \\
\frac{\partial^2 l}{\partial \mathbf{v}_m \partial \mathbf{v}_m^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{c_j B_{mj} - \mathbf{b}_{mj} \mathbf{b}_{mj}^T}{c_j^2} - \frac{(\sum_{l \in \mathcal{R}_i} c_l) \sum_{l \in \mathcal{R}_i} B_{ml} - (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml})^T (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\} \\
\frac{\partial^2 l}{\partial \mathbf{v}_m \partial \mathbf{v}_n^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{c_j D_{mnj} - \mathbf{b}_{mj} \mathbf{b}_{nj}^T}{c_j^2} - \frac{(\sum_{l \in \mathcal{R}_i} c_l) \sum_{l \in \mathcal{R}_i} D_{mnl} - (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml}) (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{nl})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\}, \quad m \neq n \\
\frac{\partial^2 l}{\partial \boldsymbol{\beta}_m \partial \mathbf{v}_m^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{c_j F_{mj} - \mathbf{a}_{mj} \mathbf{b}_{mj}^T}{c_j^2} - \frac{(\sum_{l \in \mathcal{R}_i} c_l) \sum_{l \in \mathcal{R}_i} F_{ml} - (\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml}) (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\} \\
\frac{\partial^2 l}{\partial \mathbf{v}_m \partial \boldsymbol{\beta}_n^T} &= \sum_{i=1}^k \sum_{j \in \mathcal{D}_i} \left\{ \frac{-c_j G_{mnj} - g_n^{(j)} \mathbf{a}_{mj} \mathbf{b}_{mj}^T}{c_j^2} - \frac{(\sum_{l \in \mathcal{R}_i} c_l) (-\sum_{l \in \mathcal{R}_i} G_{mnl}) - (\sum_{l \in \mathcal{R}_i} \mathbf{a}_{ml}) (\sum_{l \in \mathcal{R}_i} \mathbf{b}_{ml})^T}{(\sum_{l \in \mathcal{R}_i} c_l)^2} \right\}, \quad m \neq n.
\end{aligned}$$

## ACKNOWLEDGEMENTS

O. Rosen was supported by a postdoctoral fellowship grant from the Mellon Foundation. M. Tanner was supported by NIH grant CA35464. The authors wish to thank the editor and the referees for their comments and suggestions which greatly improved this paper.

## REFERENCES

1. Jacobs, R. A., Jordan, M. I., Nowlan, S. J. and Hinton, G. E. 'Adaptive mixtures of local experts', *Neural Computation*, **3**, 79–87 (1991).
2. Peng, F., Jacobs, R. A. and Tanner, M. A. 'Bayesian inference in mixtures-of-experts and hierarchical mixtures-of-experts models with an application to speech recognition', *Journal of the American Statistical Association*, **91**, 953–960 (1996).
3. Jordan, M. I. and Jacobs, R. A. 'Hierarchical mixtures of experts and the EM algorithm', *Neural Computation*, **6**, 181–214 (1994).
4. Jacobs, R. A., Tanner, M. A. and Peng, F. 'Bayesian inference for hierarchical mixtures-of experts with applications to regression and classification', *Statistical Methods in Medical Research*, **5**, 375–390 (1996).
5. Jacobs, R. A., Peng, F. and Tanner, M. A. 'A Bayesian approach to model selection in hierarchical mixture-of-experts architectures', *Neural Networks*, **10**, 231–241 (1997).
6. Kooperberg, C., Stone, C. and Truong, Y. K. 'Hazard regression', *Journal of the American Statistical Association*, **90**, 78–94 (1995).
7. Hastie, T. and Tibshirani, R. *Generalized Additive Models*, Chapman and Hall, London, 1990.
8. O'Sullivan, F. 'Nonparametric estimation of relative risk using splines and cross-validation', *SIAM Journal of Scientific and Statistical Computing*, **9**, 531–542 (1988).
9. Sleeper, L. and Harrington, D. 'Regression splines in the Cox model with application to covariate effects in liver disease', *Journal of the American Statistical Association*, **85**, 941–949 (1990).
10. Gray, R. 'Flexible methods for analyzing survival data using splines with application to breast cancer prognosis', *Journal of the American Statistical Association*, **87**, 942–951 (1992).
11. LeBlanc, M. and Crowley, J. 'Relative risk trees for censored data', *Biometrics*, **48**, 411–425 (1992).
12. Breiman, L., Friedman, J., Olshen, R. and Stone, C. *Classification and Regression Trees*, Wadsworth, Belmont, California, 1984.
13. Cox, D. R. 'Regression models and life-tables (with discussion)', *Journal of the Royal Statistical Society, Series B*, **34**, 187–220 (1972).
14. Peto, R. 'Contribution to the discussion of Cox', *Journal of the Royal Statistical Society, Series B*, **34**, 205–207 (1972).
15. Breslow, N. E. 'Covariance analysis of censored survival data', *Biometrics*, **30**, 89–99 (1974).
16. Kalbfleisch, J. D. and Prentice, R. L. *The Statistical Analysis of Failure Time Data*, Wiley, New York, 1980.
17. Breslow, N. E. 'Contribution to the discussion of Cox', *Journal of the Royal Statistical Society, Series B*, **34**, 216–217 (1972).
18. Miller, R. G. *Survival Analysis*, Wiley, New York, 1981.
19. Rosen, O., Jiang, W. and Tanner, M. 'Mixtures of marginal models', in press (1999).
20. Krall, J. M., Uthoff, V. A. and Harley, J. B. 'A step-up procedure for selecting variables associated with survival', *Biometrics*, **31**, 49–57 (1975).
21. Collett, D. *Modelling Survival Data in Medical Research*, Chapman and Hall, London, 1994.
22. Cronin, K., Slate, E. H., Turnbull, B. W. and Wells, M. T. 'Using the Gibbs sampler to detect changepoints: application to PSA as a longitudinal marker of prostate cancer', in Sall, J. and Lehman, A. (eds), *Computing Science and Statistics*, Interface Foundation of North America, 1994, pp. 314–318.
23. Jiang, W. and Tanner, M. A. 'Hierarchical mixtures-of-experts for exponential family regression models with generalized linear mean functions: A survey of approximation and consistency results', *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, pp. 296–303 (1998).