

# **Penalized Survival Models and Frailty**

Terry M. Therneau, Patricia M. Grambsch, V. Shane Pankratz

Technical Report #66

June 2000

Copyright 2000 Mayo Foundation

# Penalized Survival Models and Frailty

Terry M Therneau, Patricia M Grambsch, and V. Shane Pankratz

June 5, 2000

## Abstract

Interest in the use of random effects in the survival analysis setting has been increasing. However, the computational complexity of such frailty models has limited their general use. While fitting frailty models has traditionally been a difficult task, standard algorithms for fitting Cox semi-parametric and parametric regression models can be readily extended to include penalized regression. We demonstrate that solutions for gamma shared frailty models can be obtained exactly via penalized estimation. Similarly, Gaussian frailty models are closely linked to penalized models. This makes it possible to apply penalized estimation to other frailty models using Laplace approximations. Fitting frailty models with penalized likelihoods can be made quite rapid by taking advantage of computational methods available for penalized models. We have implemented penalized regression for the `coxph` function of `S-plus` and illustrate the algorithms with examples using the Cox model.

KEY WORDS Cox model, penalized likelihood, proportional hazards, random effects

## 1 Introduction

In the last several years there has been significant and active research concerning the addition of random effects to survival models. In this setting, a random effect is a continuous variable which describes excess risk or *frailty* for distinct categories, such as individuals or families, over and above

---

any measured covariates. The idea is that individuals have different frailties, and that those who are most frail will die earlier than the others. Aalen [1] provides theoretical and practical motivation for frailty models by discussing the impact of heterogeneity on analyses, and by illustrating how random effects can deal with it.

Frailties are useful in modeling correlations in multivariate survival and event history data. Examples include recurrent events such as epileptic seizures or depressive episodes, where an individual's frailty influences the occurrence of events, and community trials, where the different events within each community share a common frailty. The simplest model, implicit in these examples, is the shared frailty model. In this model, all the units within each category share a common frailty, each unit belongs to precisely one category, and frailties of different categories are independent. More complex models are possible. Frailties can be nested; individuals within a family may share a common frailty, while families within communities share another common frailty. Frailties can also be correlated, as in studies of pedigrees. Due to its simplicity, we emphasize the shared frailty model here.

Frailties are usually viewed as unobserved covariates. This has led to the use of the EM algorithm as an estimation tool. However, the algorithm is slow, variance estimates require further computation, and no implementation has appeared in any of the more widely available packages.

Penalized models provide an alternate approach. The frailty terms are treated as additional regression coefficients which are constrained by a penalty function added to the log-likelihood. They are computationally similar to other shrinkage methods for penalized regression such as ridge regression, the lasso and smoothing splines. Standard algorithms for fitting Cox semi-parametric and parametric models can be simply extended to include penalty functions. These methods usually converge quickly and produce both point and variance estimates for model parameters.

We discuss below the link between penalized estimation and frailty models. In particular, we demonstrate that if the frailty has a gamma distribution, then the shared frailty model can be written exactly as a penalized likelihood. We also show that Gaussian frailty models are closely linked to penalized models. We then turn to computational issues in implementing penalized techniques for

fitting proportional hazard frailty models. We describe our S-plus implementation and illustrate the algorithms with several examples.

## 2 Frailty Models

Assume that the data for subject  $i$ , who is a member of the  $j$ th of  $q$  families, follows a proportional hazards shared frailty model. The hazard can be written as

$$\lambda_i(t) = \lambda_0(t)\varpi_{j(i)}e^{X_i\beta}, \quad (1)$$

where  $j(i)$  denotes that individual  $i$  belongs to family  $j$ ,  $\varpi_{j(i)} = \varpi_j$  is the frailty for family  $j$ ,  $X$  is the covariate matrix of dimension  $n$  by  $p$ , and  $\beta$  is a vector of regression coefficients. The  $\varpi$ 's are independent and identically distributed from some positive scale family with density function  $f(\varpi; \theta)$ , having mean 1 and variance  $\theta$  for identifiability

If the  $\varpi$ 's are known, the complete data log-likelihood is

$$\begin{aligned} \sum_{i=1}^n & \left[ \int_0^\infty Y_i(t) [\log(\lambda_0(t)) + \log(\varpi_{j(i)}) + X_i\beta] dN_i(t) \right. \\ & \left. - \int_0^\infty Y_i(t) \varpi_{j(i)} \exp(X_i\beta) \lambda_0(t) dt + \log f(\varpi_{j(i)}; \theta) \right]. \end{aligned}$$

If the  $\varpi$  are viewed as missing data, the problem can be approached using the EM algorithm. Parner [12] lays out a general framework. Let  $\phi(s) = \phi(s, \theta)$  be the Laplace transform of the distribution of  $\varpi$ , and let  $\phi^{(n)}(s)$  be its  $n$ th derivative with respect to  $s$ . Let  $A_j = A_j(\beta, \lambda_0) = \sum \int_0^\infty Y_i(s) \exp(X_i\beta) d\Lambda_0(s)$ , where the sum is over the members of family  $j$ , and let  $d_j$  be the number of events in the  $j$ th family. The log-likelihood of the observed data,

$$L_m(\beta, \lambda_0, \theta) = \sum_{i=1}^n \delta_i \log \left( \int_0^\infty Y_i(t) e^{X_i\beta} \lambda_0(t) \right) + \sum_{j=1}^q \log [(-1)^{d_j} \phi^{(d_j)}(A_j)], \quad (2)$$

is found by integrating over the distribution of  $\varpi$ . For any fixed value of  $\theta$ , Parner suggests maximizing this likelihood for  $\beta$  and  $\lambda_0$  by an EM algorithm, which alternates between the following steps.

- 1 M-step. Treat the current estimate of  $\varpi$  as a fixed value or *offset*, and update  $\beta$  and  $\lambda_0$  as in

usual Cox regression. Note that for given  $\beta$  and  $\varpi$ ,

$$d\hat{\Lambda}_0(t, \beta, \varpi) = \sum dN_i(t) / \sum Y_i(t) \varpi_{j(i)} \exp(X_i \beta). \quad (3)$$

2. E-step. Compute  $\varpi$  as the expected value given the current values  $\beta$  and  $\lambda_0$  and the data.

$$\varpi_j = -\frac{\phi^{(d_j+1)}(\hat{A}_j)}{\phi^{(d_j)}(\hat{A}_j)}, \quad (4)$$

where  $\hat{A}_j = A_j(\beta, \hat{\lambda}_0(\beta, \omega))$ .

Equations 2 and 4 require the shared frailty model and unfortunately do not hold for more complex models. Parner suggests that estimation of  $\theta$  be done by maximizing the profile log-likelihood

$$L_m(\theta) = L_m(\hat{\beta}(\theta), \hat{\lambda}_0(\theta), \theta). \quad (5)$$

Although  $\varpi$  is not an explicit parameter of the observed log-likelihood, the EM algorithm provides an estimate of this vector.

The penalized regression formulation for the shared frailty model is most easily developed from an alternative version of the hazard,

$$\lambda_i(t) = \lambda_0(t) e^{X_i \beta + Z_i \omega}, \quad (6)$$

which is equivalent to Equation 1. In this case,  $\varpi_j = \exp(\omega_j)$ ,  $Z$  is matrix of  $g$  indicator variables such that  $Z_{ij} = 1$  when subject  $i$  is a member of family  $j$  and 0 otherwise, and each individual belongs to only one family. Estimation under this model is done by maximizing a penalized partial log-likelihood

$$PPL = PL(\beta, \omega; \text{data}) - g(\omega; \theta)$$

over both  $\beta$  and  $\omega$ . Here PL is the log of the usual Cox partial likelihood,

$$PL(\beta, \omega) = \sum_{i=1}^n \int_0^{\infty} \left[ Y_i(t)(X_i \beta + Z_i \omega) - \log \left\{ \sum_k Y_k(t) \exp(X_k \beta + Z_k \omega) \right\} \right] dN_i(t) \quad (7)$$

and  $g$  is a penalty function chosen by the investigator to restrict the values of  $\omega$ . The parameter  $\theta$  is a tuning constant which may be pre-specified or adapted to the data. Typically, one would choose the penalty function to “shrink”  $\omega$  toward zero and use  $\theta$  to control the amount of shrinkage.

To estimate  $\beta$  and  $\omega$ , one solves the score equations. Because the penalty function does not involve  $\beta$ ,  $\partial PPL/\partial\beta = \partial PL/\partial\beta$ . Therefore, the score equations for  $\beta$  are identical to those for an ordinary Cox model treating  $Z\omega$  as an offset term. If we define

$$\bar{z}_j(t) = \bar{z}_j(\beta, \omega, t) = \frac{\sum Z_{ij} Y_i(s) \exp[X_i\beta + Z_i\omega]}{\sum Y_i(s) \exp[X_i\beta + Z_i\omega]}, \quad (8)$$

then

$$\frac{\partial PPL}{\partial\omega_j} = \sum_{i=1}^n \int_0^\infty (Z_{ij} - \bar{z}_j(t)) dN_i(t) - \frac{\partial g(\omega; \theta)}{\partial\omega_j}. \quad (9)$$

Recall that for given  $\beta$  and  $\omega$ , the Breslow estimator of the underlying hazard is

$$d\hat{\Lambda}_0(t, \beta, \omega) = \sum dN_i(t) / \sum Y_i(t) \exp(X_i\beta + Z_i\omega),$$

which is just Equation 3 in different notation. Let  $\hat{\lambda}_i = \hat{\lambda}_i(\beta, \omega) = \int_0^\infty Y_i(s) d\hat{\Lambda}_0(t; \beta, \omega)$ . Simple algebra shows that the score equation for  $\omega_j$  is

$$\frac{\partial PPL}{\partial\omega_j} = \sum_{i=1}^n \left[ Z_{ij} \delta_i - Z_{ij} \hat{\lambda}_i e^{X_i\beta + Z_i\omega} \right] - \frac{\partial g(\omega; \theta)}{\partial\omega_j} = 0. \quad (10)$$

Because of the structure of the matrix  $Z$ , this equation simplifies to

$$\frac{\partial PPL}{\partial\omega_j} = \left[ d_j - \hat{A}_j e^{\omega_j} \right] - \frac{\partial g(\omega; \theta)}{\partial\omega_j} = 0, \quad (11)$$

where  $d_j$  and  $\hat{A}_j$  are as defined above.

The penalized likelihood can be fit with the Newton-Raphson algorithm. In addition to the score vectors  $\partial PPL/\partial\beta$  and  $\partial PPL/\partial\omega$ , this requires the Hessian of the penalized partial log-likelihood:

$$H = H(\beta, \omega) = \mathcal{I} + \begin{pmatrix} 0 & 0 \\ 0 & g'' \end{pmatrix}, \quad (12)$$

where  $\mathcal{I} = \mathcal{I}(\beta, \omega)$  is the usual Cox model information matrix, or the second derivative matrix of  $PL$  with respect to  $\beta$  and  $\omega$ .

## 2.1 Gamma frailty

Details of the EM approach for the shared gamma frailty model can be found in Nielsen et al. [11] and Klein [6]. Equations 4 and 2 can be used to re-derive their results, and help make the connection

to penalized methods. Here we demonstrate that for any fixed  $\theta$ , the penalized log-likelihood with appropriate choice of penalty function and the observed-data log-likelihood in Equation 2 have the same solution.

Let the frailty have a gamma distribution with mean 1 and variance  $\theta = 1/\nu$ . The density of  $\varpi$  can be written as

$$\log[f(\varpi; \nu)] = (\nu - 1) \log(\varpi) - \nu\varpi + \nu \log(\nu) - \log \Gamma(\nu).$$

This has a Laplace transform of  $\phi(s) = (1 + s/\nu)^{-\nu}$ . The derivatives of  $\phi(s)$  are

$$\phi^{(d)}(s) = \left(-\frac{1}{\nu}\right)^d \left(1 + \frac{s}{\nu}\right)^{-(\nu+d)} \prod_{i=0}^{d-1} (\nu + i),$$

and Equation 4 reduces to

$$e^{\omega_j} = \frac{d_j + \nu}{\hat{A}_j + \nu}. \quad (13)$$

### Lemma

The solution to the penalized partial likelihood model, with penalty function

$$g(\omega; \theta) = -1/\theta \sum_{j=1}^q [\omega_j - \exp(\omega_j)],$$

coincides with the EM solution for any fixed value of  $\theta$ .

### Proof

For  $\beta$ , the EM and penalized methods have the same score equation, which includes  $Z\omega$  as a fixed offset. Thus if the solutions for  $\omega$  are the same, those for  $\beta$  will be also. Let  $(\hat{\beta}, \hat{\omega})$  be a solution to the EM process. Then  $\hat{\omega}$  must satisfy Equation 13 exactly, not just as an update step. Rearranging terms, we see that  $\hat{A}_j = \exp(-\hat{\omega}_j)(d_j + \nu) - \nu$ . Substituting this into the penalized score equation and simplifying with  $\nu = 1/\theta$  a fixed quantity, we see that

$$\begin{aligned} \frac{\partial PPL(\hat{\beta}, \hat{\omega})}{\partial \hat{\omega}_j} &= \left[ d_j - \hat{A}_j e^{\hat{\omega}_j} \right] - \frac{\partial g(\hat{\omega}; \theta)}{\partial \hat{\omega}_j} \\ &= \left[ d_j - e^{-\hat{\omega}_j} \left( d_j + \frac{1}{\theta} - \frac{1}{\theta} e^{\hat{\omega}_j} \right) e^{\hat{\omega}_j} \right] + \frac{1}{\theta} (1 - e^{\hat{\omega}_j}) \\ &= 0. \end{aligned}$$

This shows that the solution to the EM algorithm is also a solution to the penalized score equations.

Therefore, for any fixed  $\theta$ , the penalized log-likelihood and the observed-data log-likelihood in Equation 2 have the same solution, although these two equations are *not* equal to one another.

Furthermore, if we let  $PPL(\theta) = PPL(\hat{\beta}(\theta), \hat{\omega}(\theta), \theta)$ , then we can write Equation 5, the profile log-likelihood for  $\theta$ , as  $PPL(\theta)$  plus a correction that only involves  $\theta$  and the  $d_j$ s. Using the fact that each row of  $Z$  has exactly one 1 and  $q - 1$  0s, we see that the Cox PL for  $(\hat{\beta}, \hat{\omega})$  must be the same as that for  $(\hat{\beta}, \hat{\omega} + c)$  for any constant  $c$ . Simple algebra shows that the value of  $c$  which minimizes the penalty portion of the  $PPL$  is such that

$$\sum_{i=1}^q e^{\hat{\omega}_i} = q. \quad (14)$$

Using the identities in Equations 13 and 14, recalling that they hold only at the solution point, we show in the appendix that

$$L_m(\theta) = PPL(\theta) + \sum_{j=1}^q \nu - (\nu + d_j) \log(\nu + d_j) + \nu \log \nu + \log \left( \frac{\Gamma(\nu + d_j)}{\Gamma(\nu)} \right). \quad (15)$$

It is useful to consider  $L_m(\theta) + \sum_{j=1}^q d_j$ , rather than  $L_m(\theta)$ , because the profile log-likelihood converges to  $PL(\hat{\beta}) - \sum d_j$  as the variance of the random effect goes to zero. Adding  $\sum d_j$ , to  $L_m(\theta)$  makes the maximized marginal likelihood from a frailty model with small  $\theta$  comparable to the maximized likelihood from a non-frailty model.

The fitting program for a shared gamma frailty consists of an inner and outer loop. For any fixed  $\theta$ , Newton-Raphson iteration is used to solve the penalized model in a few (usually 3-5) steps, and return the corresponding value of the PPL. The outer loop chooses  $\theta$  to maximize the profile likelihood in Equation 15, which is easily done as it is a unimodal function of one parameter.

All of the results presented in this section were dependent on the correct choice of a penalty function. For gamma frailties, the penalty function that links the penalized and EM results is directly related to the density of the random effect; the log of the density for  $\omega$ , where  $\exp(\omega)$  has a gamma distribution, is equal to  $[\omega - \exp(\omega)]/\theta$  plus additional terms not involving  $\omega$ . Similarly, the penalty we use for a Gaussian frailty is related to a log-density, as discussed in the next section



## 2.2 Gaussian Frailty

McGilchrist and Aisbet [10, 9], suggest a Gaussian density for  $\omega$  in a shared frailty model. This leads to the penalized partial likelihood

$$PPL = PL - (1/2\theta) \sum_{j=1}^q \omega_j^2, \quad (16)$$

where  $\theta$  is the variance of the random effect.

The authors do not provide an exact connection to the marginal likelihood that can be used to choose the variance parameter  $\theta$ . Instead, they note the similarity of the Cox model's Newton-Raphson step to an iteratively re-weighted least-squares calculation. Using this observation, they propose using standard estimators from Gaussian problems. This leads to choosing  $\theta$  such that it satisfies

$$\theta = \frac{\sum_{j=1}^q \omega_j^2 + r}{q}. \quad (17)$$

The value of  $r$  varies depending on the estimation technique used. For BLUP,  $r = 0$ ; for MLE,  $r = \text{trace}[(H_{22})^{-1}]$ ; and for REML,  $r = \text{trace}[(H^{-1})_{22}]$ , where  $H$  is the Hessian of the penalized partial log likelihood in Equation 12 and  $H_{22}$  is the lower right  $q \times q$  submatrix corresponding to the random effects.

The Gaussian approach is justified and expanded in Ripatti and Palmgren [14]. Let the random effects have a positive definite covariance matrix  $D = D(\theta)$ . This provides a rich class of models for the random effects; for example, setting  $D = \theta I$  results in a shared frailty model. The marginal log-likelihood is

$$L_m(\beta, \theta) = -1/2 \log |D| + \log \left\{ \int \exp[PL(\beta, \omega) - 1/2\omega' D^{-1/2} \omega] d\omega \right\}.$$

Note that, unlike Parner's approach, this marginal log-likelihood does not involve  $\lambda_0(t)$ , that has already been partialled out to give the Cox partial log-likelihood. Following the methods of Breslow and Clayton [2], Ripatti and Palmgren use a Laplace approximation to the above integral to get an approximate marginal log-likelihood.

$$L_m(\beta, \theta) \approx PL(\beta, \tilde{\omega}) - 1/2 (\tilde{\omega}' D^{-1} \tilde{\omega} + \log |D| + \log |H(\beta, \tilde{\omega}_{22})|), \quad (18)$$

where  $\tilde{\omega} = \tilde{\omega}(\beta, \theta)$  solves

$$\sum_{i=1}^n \int_0^{\infty} (Z_{i,j} - \bar{z}_j(t)) dN_i(t) - D(\theta)^{-1} \tilde{\omega} = 0$$

which is comparable to Equation 9. As a result, the first two terms of the approximate marginal log-likelihood correspond to a penalized partial likelihood with  $g(\omega; \theta) = 1/2 \tilde{\omega}' D(\theta)^{-1} \tilde{\omega}$ . This reduces to Equation 16 in the case of a shared frailty model. We can ignore the third term of Equation 18 as  $D$  is constant for fixed  $\theta$ . Ignoring the fourth term can influence the estimates, but Ripatti and Palmgren suggest that the loss of information is slight.

As shown in Ripatti and Palmgren [14], the estimating equation for  $\theta_j$  is

$$\text{trace} \left[ D^{-1} \frac{\partial D}{\partial \theta_j} \right] + \text{trace} \left[ (H_{22})^{-1} \frac{\partial D^{-1}}{\partial \theta_j} \right] - \omega' D^{-1} \frac{\partial D}{\partial \theta_j} D^{-1} \omega = 0. \quad (19)$$

The Fisher information matrix, obtained by taking the expectation with respect to  $\omega$ , has a  $jk$  element of

$$\begin{aligned} (1/2) \text{ trace} \left[ D^{-1} \frac{\partial D}{\partial \theta_j} D^{-1} \frac{\partial D}{\partial \theta_k} + D^{-1} \frac{\partial^2 D}{\partial \theta_j \partial \theta_k} \right] + \\ (1/2) \text{ trace} \left[ (H_{22})^{-1} \frac{\partial D}{\partial \theta_j} (H_{22})^{-1} \frac{\partial D}{\partial \theta_k} - (H_{22})^{-1} \frac{\partial^2 D^{-1}}{\partial \theta_j \partial \theta_k} \right]. \end{aligned} \quad (20)$$

For the shared frailty model the estimating equation reduces to

$$\hat{\theta} = \frac{\omega' \omega + \text{trace}[(H_{22})^{-1}]}{q},$$

which is equivalent to the MLE formula of McGilchrist [9].

Yau and McGilchrist [16] display a similar formula for the ML estimate for an arbitrary correlation matrix  $D$ , and apply the results to the CGD data set using an AR(1) structure for the multiple infections within subject. (Unfortunately, differences in how ties are handled make it impossible to replicate their fits). In that paper, they also define an REML estimate, which is identical to Equations 19 and 20 above, but with  $(H^{-1})_{22}$  replacing  $(H_{22})^{-1}$ . Additionally, their simulations show Equation 20 to be an overestimate of the actual standard error.

### 3 Computational Issues

Thus far, we have discussed the relationship between frailty models and penalized likelihood estimation. In this section, we describe several issues important to the computational implementation of penalized likelihood methods for Cox models with random effects.

#### 3.1 Penalized Likelihood Inference

Consider a Cox model with both constrained and unconstrained effects, as shown in Equation 6. The model is fit by maximizing the penalized partial log-likelihood (PPL). We assume that  $\theta$  is fixed. Consider testing the set of hypotheses  $z = C(\beta', \omega')' = 0$ , where  $(\beta', \omega')'$  is the combined vector of  $p + q$  parameters, and  $C$  is a  $k \times p + q$  matrix of full row rank  $k$ ,  $k \leq p + q$ . Gray [3] suggests that

$$V = H^{-1} \mathcal{I} H^{-1} \tag{21}$$

be used as the covariance estimate of the parameter estimates. He recommends a Wald type test statistic,  $z'(CH^{-1}C')^{-1}z$ , with generalized degrees of freedom

$$df = \text{trace}[(CH^{-1}C')^{-1}(CVC')].$$

The total degrees of freedom for the model ( $C = I$ ) simplifies to

$$\begin{aligned} df &= \text{trace}[HV] \\ &= \text{trace}[H(H^{-1}(H - G)H^{-1})] \\ &= (p + q) - \text{trace}[GH^{-1}]. \end{aligned} \tag{22}$$

Under  $H_0$ , the distribution of the test statistic is asymptotically the same as  $\sum e_i X_i^2$ , where the  $e_i$  are the  $k$  eigenvalues of the matrix  $(CH^{-1}C')^{-1}(CVC')$  and the  $X_i$  are iid standard Gaussian random variables. In non-penalized models, the  $e_i$  are all either 0 or 1, and the test statistic has an asymptotic chi-square distribution on  $\sum e_i$  degrees of freedom. In penalized models, the test statistic has mean  $\sum e_i$  and variance  $2 \sum e_i^2 < 2 \sum e_i$ , because  $0 \leq e_i \leq 1$ . Using a reference chi-square distribution with  $df = \sum e_i$ , will tend to be conservative.

Verweij and Van Houwelingen [15] discuss penalized Cox models in the context of restricting the parameter estimates. They use  $H^{-1}$  as a “pseudo standard error”, and an “effective degrees of freedom” identical to Equation 22. With this variance matrix, the test statistic  $z'(CH^{-1}C')^{-1}z$  is a usual Wald test. To choose an optimal model they recommend either the Akaike Information Criterion (AIC) which uses the degrees of freedom described above or the cross-validated (partial) log-likelihood CVL, which uses a degrees of freedom estimate based on a robust variance estimator.

Our algorithm makes both  $H^{-1}$  and  $H^{-1}TH^{-1}$  available. Significance tests are based on  $H^{-1}$  as the more conservative choice. Simulation experiments for the related problem of penalized smoothing splines in Cox regression (not shown) suggest that this is the more reliable choice for tests, but we do not have more definitive results to support this.

In our implementation, the computation of the degrees of freedom and variance matrices are specialized to avoid any intermediate steps that would give a  $q$  by  $q$  result, where  $q$  is the number of constrained coefficients.

### 3.2 Sparse computation

When performing estimation with frailty models, memory and time considerations can become an issue. For instance, if there are 300 families, each with a frailty term, and 4 other variables, then the full information matrix has  $304^2 = 92416$  elements. The Cholesky decomposition must be applied to this matrix within each Newton-Raphson iteration. In our S-plus implementation, we have applied a technique that can provide significant savings in space and time.

If we partition the information matrix of a Cox shared frailty model according to the rows of  $X$  and  $Z$ , and arrange the matrix as

$$\mathcal{I} = \begin{pmatrix} \mathcal{I}_{ZZ} & \mathcal{I}_{ZX} \\ \mathcal{I}_{XZ} & \mathcal{I}_{XX} \end{pmatrix},$$

then the upper left corner will be a diagonally dominant matrix, having almost the form of the variance matrix for a multinomial distribution. Adding the penalty further increases the dominance

of the diagonal. Therefore, using a *sparse* computation option, where only the diagonal of  $\mathcal{I}_{ZZ}$  is retained, should not have a large impact on the estimation procedure.

Ignoring a piece of the full information matrix has a number of implications. First, the speed of the Cholesky factorization is increased dramatically. Second, the savings in space can be considerable. If we use the sparse option with the example above, the information matrix consists of only  $\mathcal{I}_{XZ}$  and  $\mathcal{I}_{XX}$ , with  $304 * 4 = 1216$  elements, along with the 300 element diagonal of  $\mathcal{I}_{ZZ}$ , a savings of over 95% in memory space. Third, because the score vector and likelihood are not changed, the solution point is identical to the one obtained in the non-sparse case, discounting trivial differences due to distinct iteration paths. Fourth, the Newton-Raphson iteration may undergo a slight loss of efficiency so that 1-2 more iterations are required. However, because each N-R iteration requires the Cholesky decomposition of the information matrix, the sparse problem is much faster per-iteration than the full matrix version. Finally, the full information matrix is a part of the formulas for the post-fit estimates of degrees of freedom and standard error. In a small number of simple examples, the effect of the sparse approximation on these estimates has been surprisingly small.

We have found two cases where our sparse method does not perform acceptably. The first is if the variance of the random effect is quite large ( $>5$ ). In this case, each N-R iteration may require a large number ( $>15$ ) iterations. The second is if one group contains a majority of the observations. The off diagonal terms are too important to ignore in this case, and the approximate N-R iteration does not converge.

## 4 Examples

We now present two examples where we use our S-plus functions to obtain estimates from frailty models. The first deals with the survival of kidney catheters. The second examines the effect of UDCA in patients with primary biliary cirrhosis.

## 4.1 Survival of kidney catheters

The following data set is presented in McGilchrist and Aisbett [10]. Each observation is the time to infection, at the point of insertion of the catheter, for kidney patients using portable dialysis equipment. Catheters may be removed for reasons other than infection, in which case the observation is censored. There are 38 patients, each with exactly 2 observations. Variables are the subject id, age, sex (1=male, 2=female), disease type (glomerulo nephritis, acute nephritis, polycystic kidney disease, and other), and the time to infection or censoring for each insertion. We first fit two ordinary Cox models, followed by a gamma frailty fit.

```
> kfit1 <- coxph(Surv(time, status) ~ age + sex, data=kidney)
> kfit2 <- coxph(Surv(time, status) ~ age + sex + disease,
                 data=kidney)
> kfit3 <- coxph(Surv(time, status) ~ age + sex + disease +
                 frailty(id), data=kidney)
> kfit3
```

	coef	se(coef)	se2	Chisq	DF	p
age	0.00318	0.0111	0.0111	0.08	1	7.8e-01
sex	-1.48314	0.3582	0.3582	17.14	1	3.5e-05
diseaseGN	0.08796	0.4064	0.4064	0.05	1	8.3e-01
diseaseAN	0.35079	0.3997	0.3997	0.77	1	3.8e-01
diseasePKD	-1.43111	0.6311	0.6311	5.14	1	2.3e-02
frailty(id)				0.00	0	9.5e-01

Iterations: 6 outer, 29 Newton-Raphson

Penalized terms:

Variance of random effect= 1.47e-07 M-likelihood = -179.1

Degrees of freedom for terms= 1 1 3 0

Likelihood ratio test=17.6 on 5 df, p=0.00342 n= 76

Many of the labels in this output are self-explanatory. Several may need some clarification. The `se(coef)` estimates are from the diagonal elements of  $H^{-1}$ , and `se2` uses the diagonal entries in  $H^{-1} \mathcal{I} H^{-1}$ . In this particular data set, they are identical, but that is not always the case. The M-likelihood is the marginal likelihood in Equation 15, evaluated at the MLE.

The partial log-likelihood values for first two models are -184.3 and -179.1, with 2 and 5 degrees of freedom respectively. Hence, the disease variable is a significant addition. In the third fit, the program provided an estimate of the MLE of  $\theta$ , the variance of the random effect, that was essentially 0.

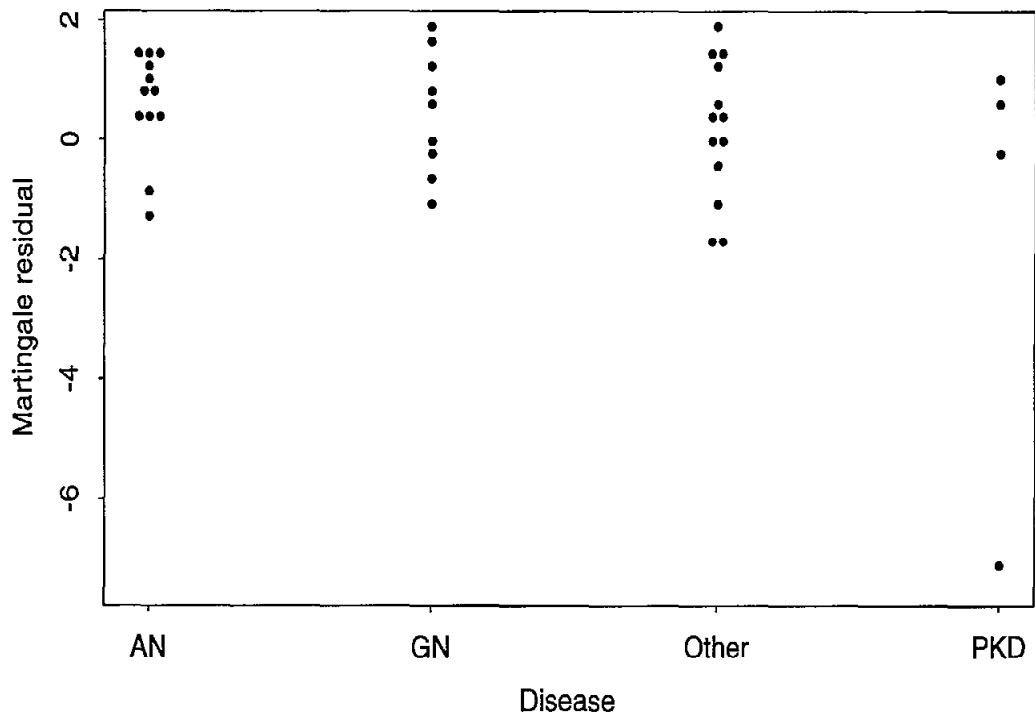


Figure 1. Residuals for the kidney data from model kfit1.

When the disease variable is left out of the random effects model, however, we get a quite different result.

```
> kfit4 <- coxph(Surv(time, status) ~ age + sex + frailty(id),
  data=kidney)
> kfit4
```

	coef	se(coef)	se2	Chisq	DF	p
age	0.00522	0.0119	0.0088	0.19	1.0	0.66000
sex	-1.58335	0.4594	0.3515	11.88	1.0	0.00057
frailty(id)				22.97	12.9	0.04100

```
Iterations: 7 outer, 49 Newton-Raphson
Variance of random effect= 0.408 M-likelihood = -181.6
Degrees of freedom for terms= 0.6 0.6 12.9
Likelihood ratio test=46.6 on 14.06 df, p=2.36e-05 n= 76
```

In this case, both the approximate Wald test and the likelihood ratio test indicate that the variance of the random effect is greater than zero. The Wald test shown in the printout is not as accurate as the the comparison of the marginal likelihood to that from `kfit1` (-184.3 vs -181.6), which gives a chi-square statistic of 5.4 on 1 degree of freedom for a p-value of 0.02. As discussed in Nielsen et al [11], this chi-square test for  $\theta$  is not affected by the boundary at zero.

Figure 1 shows the reason for the discrepancy of the results between the two models. The graph shows the martingale residuals for each subject (the sum of the residuals from the two observations), based on the simplest model, `kfit1`. Note the outlier in the lower right, corresponding to a 46 year old male whose age was quite close to the median for the study (45.5 years). There were 10 males and most had early failures: 2 observations were censored at 4 and 8 days, respectively, and the remaining 16 male kidneys had a median time to infection of 19 days. Subject 21, however, had failures at 154 and 562 days. With this subject removed, neither the disease ( $p=0.53$ ) nor the frailty ( $p>0.9$ ) are important. With this subject in the model, it is a toss-up whether the disease or the frailty term will be credited with 'significance'. Using a Gaussian frailty with REML gives partial importance to each.

```
> mfit1 <- coxph(Surv(time,status) ~ age + sex + disease +
  frailty(id, dist='gauss', sparse=F), data=kidney)
> mfit1
```

	coef	se(coef)	se2	Chisq	DF	p
age	0.00492	0.0149	0.0108	0.11	1.0	0.74000
sex	-1.70204	0.4631	0.3613	13.51	1.0	0.00024



```

diseaseAN 0.39442 0.5428 0.4052 0.53 1.0 0.47000
diseaseGN 0.18173 0.5413 0.4017 0.11 1.0 0.74000
diseasePKD -1.13160 0.8175 0.6298 1.92 1.0 0.17000
frailty(id, dist = "gauss" 18.13 12.3 0.12000

```

Iterations: 6 outer, 17 Newton-Raphson

Variance of random effect= 0.509 M-likelihood = -171.9

Degrees of freedom for terms= 0.5 0.6 1.7 12.3

Likelihood ratio test=118 on 15.14 df, p=0 n= 76

The sparse routines have some impact on the solution for a Gaussian model, since the REML estimate depends on the matrix  $H$ . Using the `sparse=T` option in the `frailty` function, the routine required 32 Newton-Raphson iterations and gave a solution of  $\theta = 0.493$ , but with about one third the total computing time.

The standard error estimates reported by a penalized coxph model in S-plus are computed under the assumption that  $\theta$  is fixed. For some models, such as a smoothing spline with user specified degrees of freedom, For the above frailty models it clearly is not and the standard errors are an underestimate. Using the bootstrap, we found the standard error to be *much* higher for this data set, which is not surprising given the inordinate influence of a single subject. More useful bootstrap results appear in the second example.

These answers differ slightly from the original authors' [9] results. Their paper presents formulas that are completely valid only for untied data, and this data set has 5 tied pairs and one quadruple. This is a small proportion of the data, and in a standard Cox model the ties would barely perturb the answers. Unfortunately, the REML solution for  $\theta$  can be very sensitive to small changes in the data.

## 4.2 UDCA in Patients With PBC

Primary biliary cirrhosis (PBC) is a chronic cholestatic liver disease characterized by progressive destruction of the bile ducts. PBC frequently progresses to cirrhosis, which may lead to death from liver failure unless liver transplant is offered – an extensive and costly procedure. Trials have been held for several promising agents, but an effective therapy remains elusive. Although progression of disease is inexorable the time course can be very long, many patients survive 10 or more years from

their initial diagnosis before requiring a transplant.

	UDCA	Placebo
Death	6	10
Transplant	6	6
Drug toxicity	0	0
Voluntary withdrawal	11	18
Histologic progression	8	12
Development of varices	8	17
Development of ascites	1	5
Development of encephalopathy	3	1
Doubling of bilirubin	2	15
Worsening of symptoms	7	9

Table 1: Total number of events in the UDCA trial

A randomized double-blind trial of a new agent, ursodeoxycholic acid (UDCA), was conducted at the Mayo Clinic from 1988 to 1992 and enrolled 180 patients. The data are reported in Lindor et al [8]; the analysis shown here has slightly longer follow-up. The endpoints of the study were pre-defined and are shown in Table 1. Although nearly all of the comparisons favored UDCA, none were significant individually. The primary report was based on an analysis of time to the first event; 58/84 placebo and 34/86 UDCA patients have at least one event. An analysis that used all of the events data would seem to be more complete, however, since it would be based on 93 placebo and 52 UDCA events, a gain in “information” of 57%.

The event endpoints are all unique, i.e., no single patient had more than one instance of death, transplant, doubling of bilirubin, etc. Three possible methods of analysis present themselves. The simplest is time to the first adverse event. In this case, each patient has a single observation and correlation is not an issue. The second is a marginal analysis in the manner of Lin [7]. The a third

involves the use of a frailty model. The data set for the latter two options is essentially a concatenation of the 9 individual data sets that would be created for an analysis of time to death (censoring all other causes), time to transplant, time to withdrawal, etc., with the event type as a stratification variable.

The covariates in each of these models are treatment and two of the stratification factors used in treatment assignment. The resulting parameter estimates and their standard errors are shown in Table 2. The robust standard error estimates for the frailty model were obtained from 1000 bootstrap realizations where  $\theta$  was fixed at the original model's estimate. They show the ordinary standard error to be quite reliable as an estimate when  $\theta$  is incorrectly fixed in advance. The bootstrap standard error estimates obtained when  $\theta$  was estimated were higher than those shown in Table 2, being 0.42, 0.51, and 0.50 for treatment, bilirubin, and stage respectively.

Upon examination of Table 2, two outcomes are immediately obvious. First, the naive variance is an underestimate in the multiple event model; accounting for the within-patient correlation is important. Second, the multiple-event robust variances and the frailty variance estimates are slightly larger than the variances for first events only. The use of multiple events added no information to the analysis!

A closer look at the data reveals the cause of the difficulty. Patients participating in the study returned for evaluation once a year, which is the point at which most of the outcomes were measured. For instance, one patient had 5 events, 4 of which were recorded on 20 July 1990. The fifth, death, occurred on 22 July. Similar outcomes are seen for many others. Figure 2 shows the event times for the 31 subjects with multiple adverse outcomes, with a circle marking each event. The data has been jittered slightly to avoid overlap. It appears that the use of multiple event types was useful in this study only to make the detection of "liver failure" more sensitive. Given that failure has occurred, the number of positive markers for failure was irrelevant.

In this situation we expect the frailty model to show significant within patient correlation, and indeed this is the case. The variance of the random effect is estimated as 1.47 in a shared gamma frailty model and is highly significant ( $\chi^2 = 31.6$  on 1 df). The estimated value of Kendall's  $\tau$  is

	$\beta$	$se(\beta)$	robust se
first event			
treatment	-0.94	0.22	0.22
bilirubin	0.74	0.24	0.23
stage	-0.02	0.25	0.25
marginal model			
treatment	-0.80	0.17	0.23
bilirubin	0.77	0.18	0.25
stage	0.05	0.21	0.28
frailty model			
treatment	-0.96	0.28	0.32
bilirubin	0.74	0.31	0.32
stage	0.31	0.32	0.35

Table 2: Results of 3 models for the UDCA data

$$\theta/(2 + \theta) = 0.42.$$

## 5 Concluding Remarks

Penalized estimation techniques are useful estimation tools. We have now shown that estimation using shared Gamma frailty models can be performed exactly with penalized likelihood methods. This is true for models with time-dependent covariates as well as for models with time-independent covariates, which we focused on in an attempt to keep the notation simple. We have yet to find such a correspondence for more general Gamma frailty models, such as the nested frailty model of Guo and

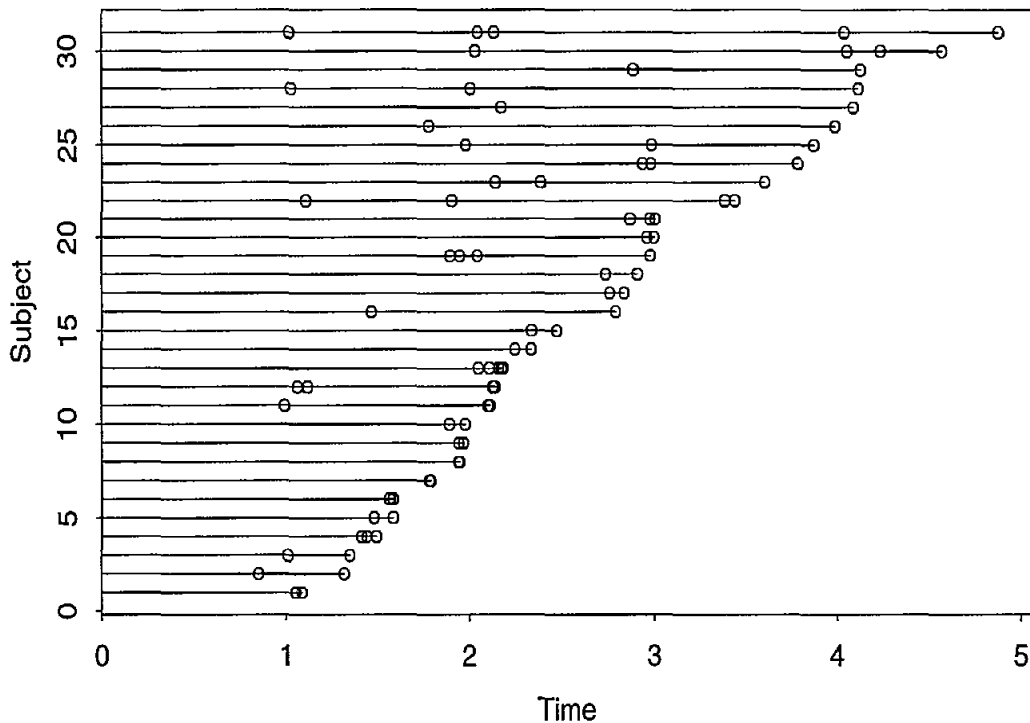


Figure 2: Multiple failure times for the UDCA data.

Rodríguez [4]. However, more general Gaussian frailty models can be approximately estimated using penalized estimation procedures.

Also, the programs support the use of AIC or corrected AIC [5] as a selection criteria. With this approach, models can be fit beyond those for which a formal ML-penalized correspondence has been worked out, such as models with multiple frailty terms or other frailty distributions. Using AIC as the optimization criteria for  $\theta$  and the log of a t-distribution density as the penalty term, for instance, appears to give similar results to more formal MCMC methods on two (small) local examples. Wider experience and/or formal results are needed to understand the relative merits of likelihood and degrees-of-freedom based approaches.

We outline several important issues regarding the variance of the random effect,  $\theta$ , below.

- The software does not print an estimate of the variance of  $\theta$ . However, a plot of the profile likelihood can easily be obtained by fitting a sequence of models with fixed  $\theta$ . This profile likelihood is often seriously asymmetric with a longer right tail, raising concerns about the utility of  $\text{se}(\theta)$  for either confidence intervals or tests. The current computer code does iteration on the  $\sqrt{\theta}$  scale. While this seems to speed convergence, other scales may be more appropriate.
- The estimate of the random effect is much less precise than the estimates of the coefficients  $\hat{\beta}$ . It is unclear how much data is needed for reliable estimation.
- The software prints out an approximate Wald test,  $\omega'(H^{-1})_{22}\omega$ , based on the fitted frailty coefficients. Since the number of frailty coefficients often grows with sample size, while the *effective* number might not, the statistical properties of the test are unknown. The printed test seems to be successful as a first “very significant/not at all significant” approximation, but final judgement should be based on the likelihood ratio test derived by comparing the printed *M-likelihood* value to the fit without the frailty term.
- The standard errors of the estimate are calculated as though  $\theta$  were fixed. This is true for some penalized problems, but false for the two examples given here. A bootstrap evaluation with  $\theta$  fixed at  $\hat{\theta}$  gives standard errors for the other parameters that agree with our asymptotic formula, but with  $\theta$  free the standard errors are larger by 30 to 60 percent.

Beyond its extendability, an important benefit of the penalized approach is speed. The computer code is fast enough that we can use it with computationally intensive secondary techniques such as the bootstrap. For instance, it took just over eleven minutes to perform 1000 bootstrap realizations of the kidney data holding  $\theta$  fixed, and 35 minutes when  $\theta$  was allowed to vary.

In summary, certain classes of frailty models can be formulated as penalized likelihoods. Because of its connection to other work in penalized regression, computational improvements are possible for selected models. For shared frailty models, use of a sparse Cholesky factorization provides significant computational advantages. Other, similar, gains can be made with other frailty models. As an

example, genetic family studies can be cast as a frailty model with one random effect per subject, and correlations among random effects that are block diagonal with one block per family. This can be efficiently handled using a more general sparse Cholesky algorithm.

## Appendix: Correspondence of Marginal Log-likelihoods at the Solution Point

Here, we obtain the realized value of the marginal log-likelihood at the solution point in terms of the penalized likelihood for the gamma shared frailty model. This justifies Equation 15.

Expanding Equation 2 gives

$$L_m(\beta, \lambda_0; \theta) = \sum_{i=1}^n \delta_i \log \left( \int Y_i(t) e^{X_i \beta} d\Lambda_0(t) \right) + \sum_{j=1}^q [-d_j \log \nu - (\nu + d_j) \log(1 + A_j/\nu) + \log\{\Gamma(\nu + d_j)/\Gamma(\nu)\}].$$

The log profile likelihood for  $\theta$  is just this function restricted to the one-dimensional curve defined by the maximizing values of  $\hat{\beta}(\theta)$ ,  $\hat{\omega}(\theta)$ ,  $\hat{\lambda}_0(\theta)$ . On that curve  $\hat{A}_j = (d_j + \nu - \nu e^{\hat{\omega}_j})/e^{\hat{\omega}_j}$  (see Equation 13).

With this substitution, after some rearrangement we get

$$L_m(\theta) = \sum_{i=1}^n \delta_i \log \left( \hat{\lambda}_i e^{X_i \hat{\beta} + Z_i \hat{\omega}} \right) + \sum_{j=1}^q [-(\nu + d_j) \log(\nu + d_j) + \nu \log(\nu e^{\hat{\omega}_j}) + \log \Gamma(\nu + d_j) - \log \Gamma(\nu)],$$

where  $\delta_i$  is a 0/1 indicator for an event for individual  $i$ .

Subtracting and adding the penalty function  $g(\omega; \theta) = -1/\theta \sum_{j=1}^q \omega_j - \exp(\omega_j)$ , evaluated at  $\hat{\omega}$  results in

$$\begin{aligned} L_m(\theta) &= \sum_{i=1}^n \delta_i \log(\hat{\lambda}_i e^{X_i \hat{\beta} + Z_i \hat{\omega}}) - g(\hat{\omega}; \theta) \\ &\quad + \sum_{j=1}^q [-\nu \hat{\omega}_j + \nu e^{\hat{\omega}_j} - (\nu + d_j) \log(\nu + d_j) + \nu \log(\nu e^{\hat{\omega}_j}) + \log \Gamma(\nu + d_j) - \log \Gamma(\nu)] \\ &= PPL(\theta) + \sum_{j=1}^q \left[ \nu - (\nu + d_j) \log(\nu + d_j) + \nu \log \nu + \log \left( \frac{\Gamma(\nu + d_j)}{\Gamma(\nu)} \right) \right], \end{aligned}$$

where the last step follows from Equation 14.

Note that, because considerable loss of accuracy can occur if one subtracts values of the log-gamma function, it is computationally advantageous to use

$$\log \left( \frac{\Gamma(\nu + d_j)}{\Gamma(\nu)} \right) = \sum_{i=0}^{d_j-1} \log \left( \frac{\nu + i}{\nu + d_j} \right)$$



rather than

$$\log\left(\frac{\Gamma(\nu + d_j)}{\Gamma(\nu)}\right) = \log(\Gamma(\nu + d_j)) - \log(\Gamma(\nu)).$$

## References

- [1] O O Aalen. Heterogeneity in survival analysis. *Statistics in Medicine*, 7:1121–1137, 1988.
- [2] N.E. Breslow and D.G. Clayton. Approximate inference in generalized linear mixed models. *J Amer Stat Assoc*, 88:9–25, 1993.
- [3] R J. Gray. Flexible methods for analyzing survival data using splines, with applications to breast cancer prognosis. *J Amer Stat Assoc*, 87:942–951, 1992.
- [4] G. Guo and G. Rodríguez. Estimating a multivariate proportional hazards model for clustered data using the EM algorithm, with an application to child survival in Guatemala. *J Amer Stat Assoc*, 87:969–976, 1992.
- [5] C. M. Hurvich, J S Simonoff, and C.-L. Tsai. Smoothing parameter selection in nonparametric regression using an improved akaike information criterion. *J. Royal Stat. Soc. B*, 60:271–293, 1998.
- [6] J P. Klein. Semiparametric estimation of random effects using the Cox model based on the EM algorithm. *Biometrics*, 48:795–806, 1992.
- [7] D.Y. Lin. Cox regression analysis of multivariate failure time data, the marginal approach. *Statistics in Medicine*, 13:2233–2247, 1994.
- [8] K.D. Lindor, E R. Dickson, W.P. Baldus, R.A. Jorgensen, J. Ludwig, P.A. Murtaugh, J M. Harrison, R.H. Wiesner, M.L Anderson, S.M. Lange, G. LeSage, S.S Rossi, and A.F. Hofman. Ursodeoxycholic acid in the treatment of primary biliary cirrhosis. *Gastroenterology*, 106:1284–90, 1994.

- [9] C.A. McGilchrist. REML estimation for survival models with frailty. *Biometrics*, 49:221–225, 1993.
- [10] C.A. McGilchrist and C W. Aisbett. Regression with frailty in survival analysis. *Biometrics*, 47:461–466, 1991.
- [11] G.G. Nielsen, R.D. Gill, P K. Andersen, and T I. Sørensen. A counting process approach to maximum likelihood estimation of frailty models. *Scandinavian J of Statistics*, 19:25–43, 1992.
- [12] E. Parner. Inference in semiparametric frailty models. Technical report, Ph.D. dissertation, University of Aarhus, Denmark, 1997.
- [13] R.L. Prentice and J. Cai. Marginal and conditional models for the analysis of multivariate failure time data. In J.P Klein and P.K. Goel, editors, *Survival Analysis, State of the Art*, pages 393–406. Kluwer Academic Publishers, Netherlands, 1991.
- [14] J Ripatti, S. and Palmgren. Estimation of multivariate frailty models using penalized partial likelihood. Technical report, Research Report 99/1, Department of Biostatistics, University of Copenhagen, 1999.
- [15] J.M. Verweij and H. C. Van Houwelingen. Penalized likelihood in cox regression *Statistics in Medicine*, 13:2427–2436, 1994
- [16] K.K.W. Yau and C.A. McGilchrist. ML and REML estimation in survival analysis with time dependent correlated frailty. *Statistics in Medicine*, 17:1201–1213, 1998.