

# survBayes: A introduction into the package

Volkmar Henschel, Christiane Hei, Ulrich Mansmann  
University of Munich  
Department of Medical Informatics, Biometry and Epidemiology  
Marchioninstr. 15, 81377 Munich, Germany

January 19, 2006

## Abstract

This software fits a multivariate proportional hazards model to interval censored event data by a Bayesian approach. Right and interval censored data and a lognormal or gamma frailty term can be fitted. An example is studied and the output analysed.

## 1 The basic model

The data, based on a sample of size  $n$ , consists of the triple  $(t_i, \delta_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$  where  $t_i$  is the time on study for subject  $i$ ,  $\delta_i$  is the event indicator for subject  $i$  ( $\delta_i = 1$  if event has occurred,  $\delta_i = 0$  if the observation is right censored),  $\mathbf{x}_i$  is the  $r$ -dimensional vector of covariate values for subject  $i$ .

The likelihood contribution of the  $i$ -th single observation is given by

$$\lambda_0(t_i|\mathbf{x}_i)^{\delta_i} S(t_i|\mathbf{x}_i) = \exp \left\{ \delta_i [h(t_i) + \beta' \mathbf{x}] - e^{\beta' \mathbf{x}} \int_0^{t_i} \exp[h(s)] ds \right\}$$

where  $h(s) = \ln[\lambda_0(s)]$ . The infinite dimensional problem gets to a finite dimensional one by partitioning the time axis  $[0, \infty[$  into disjoint intervals  $I_k = [\theta_{k-1}, \theta_k[$  for  $k = 1, \dots, K$  where  $\theta_k$  is the time of the  $k$ -th event and  $\theta_0 = 0$ . The largest event time observed is  $\theta_K$ . The function  $h$  is approximated by cubic B-splines

The priors for the components of the vector  $\beta$  will be multivariate normal distributed with mean 0 and a precision with a flat Wishart prior. The prior for the coefficients  $h_k$  of the function  $h$  will be a autoregressive process of order one with prior information on smoothness (Bayesian P-splines, see [4]). Writing  $h_k = h(\theta_k)$ ,  $k = 1, \dots, K$  the first order process is defined as  $h_k = h_{k-1} + \epsilon_k$  with  $\epsilon_k \sim N(0, \sigma_k^2)$  and  $h_0 \sim N(0, \sigma_0^2)$ , where  $h_0$  and  $\epsilon_k$ ,  $k = 1, \dots, K$  are pairwise independent. The variances are chosen as  $\sigma_k^2 = \Delta_k \sigma_1^2$  and  $\Delta_k$  depends on the interval lengths. The inverse of the covariance matrix,  $\Sigma^{-1}$ , is a band-matrix of bandwidth one. The parameters  $\frac{1}{\sigma_0^2} = \tau_0$  and  $\frac{1}{\sigma_1^2} = \tau_1$  are treated as hyperparameters with flat gamma priors setting both parameters equal to 0.001.

## 2 Sampling procedure

### Sampling for the parameter vector

Aitkin and Clayton [1] pointed out that the proportional hazards model can be interpreted as a generalized linear model.

Gamerman [2] describes how one can effectively sample the vector of covariates in generalized linear mixed models in a block updating step. This is a combination of the iterated least squares method (IWLS) as it is known in fitting such models with a Metropolis-Hastings sampling.

### Sampling for the baseline hazard

With the given structure of the log baseline hazard function one has to sample from a Gaussian Markov Random Field (GMRF), see Rue [7] and Knorr Held and Rue [3].

### Sampling for the dispersion parameters

For the dispersion parameters  $\sigma_0^2$  and  $\sigma_1^2$  a flat Gamma prior with rate  $\kappa$  and shape  $\nu$  is chosen. This leads to Gamma posteriors.

## 3 Extensions of the basic model

Data augmentation and a multiplicative frailty model is used to analyze clustered interval censored event data. Data augmentation is used to interfere unobserved event times. The potential clustering of event times within a statistical unit is modeled by introducing an unit specific random effect or frailty term into the proportional hazards model.

## 4 Example

Meisel et al. [6] present data on the shrinkage of aneurisms associated with cerebral arteriovenous malformations (cAVM) after embolization treatment. The time to a shrinkage of the aneurism to below 50% of the baseline volume was of interest. Several patients had multiple aneurisms. Each patient was inspected at a random inspection time *obs.t.* The censoring variable  $z$  was set to one, if at the inspection time sufficient shrinkage was observed, else the censoring indicator was set to zero.

Two covariates were considered: the degree of cAMV occlusion by embolization (dichotomized at 50%, variable *mo*) and the location of the aneurism, whether at the midline arteries or at other afferent cerebral arteries, variable *lok*.

The single aneurisms are not independent because aneurisms within a patient may shrink in the same way (because they share the same "environment"). Multiple aneurisms were observed per patient. This clustering of aneurisms is indicated by the grouping variable *gr*.

The data is loaded and inspected for the first eleven patients.

```
> library(survBayes)
```

```

Lade nötiges Paket: survival
Lade nötiges Paket: splines
Lade nötiges Paket: MCMCpack
Lade nötiges Paket: coda
Lade nötiges Paket: lattice
Lade nötiges Paket: MASS
##
## Markov Chain Monte Carlo Package (MCMCpack)
## Copyright (C) 2003-2006 Andrew D. Martin and Kevin M. Quinn
##
## Support provided by the U.S. National Science Foundation
## (Grants SES-0350646 and SES-0350613)
##

```

```

> data(AA.data)
> AA.data[1:11, ]

```

	z	mo	gr	lok	t.left	t.right
1	0	0	1	1	1.7698630	NA
2	0	1	2	1	0.9972603	NA
3	0	1	2	1	0.9972603	NA
4	0	1	2	1	0.9972603	NA
5	0	0	3	0	1.0712329	NA
6	0	0	3	1	1.0712329	NA
7	0	0	4	1	5.6547945	NA
8	0	0	5	1	1.5780822	NA
9	1	0	5	0	0.0000000	1.578082
10	1	0	5	0	0.0000000	1.578082
11	1	0	5	1	0.0000000	1.578082

The data is analyzed by applying the `survBayes` algorithm. The fit with `survBayes` gives an object which stores all sampled values in the required number after the burn in. The `str` function gives a survey over the output. The low number for the sample is only due to fast checking of the package in the CRAN. Please choose at least 5000.

```

> AA.res <- survBayes(Surv(t.left, t.right, z * 3, type = "interval") ~
+   mo + lok + frailty(gr, dist = "gauss"), data = AA.data, burn.in = 0,
+   number.sample = 10)
> str(AA.res)

```

```

List of 8
 $ t.where      : num [1:51] 0.0000 0.0275 0.0598 0.1299 0.1559 ...
 $ beta         : mcmc [1:10, 1:2] -0.424 -0.697 -0.750 -0.803 -0.478 ...
 ..- attr(*, "dimnames")=List of 2
 .. ..$ : NULL
 .. ..$ : chr [1:2] "mo" "lok"
 ..- attr(*, "mcpair")= num [1:3] 1 10 1
 ..- attr(*, "class")= chr "mcmc"
 $ cov.beta     : mcmc [1:10, 1:3] 0.329 0.111 2.449 8.149 4.632 ...
 ..- attr(*, "mcpair")= num [1:3] 1 10 1

```

```

..- attr(*, "class")= chr "mcmc"
$ lbh.coef      : mcmc [1:10, 1:53]  0.00000  0.00332 -0.01868 -0.32831 -0.38706 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:10] "lbh.coef" "lbh.coef" "lbh.coef" "lbh.coef" ...
.. ..$ : NULL
..- attr(*, "mcpair")= num [1:3] 1 10 1
..- attr(*, "class")= chr "mcmc"
$ sigma.lbh     : mcmc [1:10, 1:2]  0.00067  0.00610  0.09084  0.09739  1.02245 ...
..- attr(*, "mcpair")= num [1:3] 1 10 1
..- attr(*, "class")= chr "mcmc"
$ alpha.cluster : mcmc [1:10, 1:83]  0.000000  0.000743  0.002884 -0.002074 -0.004040 ..
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:10] "alpha.cluster" "alpha.cluster" "alpha.cluster" "alpha.cluster" ...
.. ..$ : NULL
..- attr(*, "mcpair")= num [1:3] 1 10 1
..- attr(*, "class")= chr "mcmc"
$ sigma.cluster :Class 'mcmc' atomic [1:10] 2.46e-06 5.09e-06 7.00e-06 8.76e-06 9.43e-06
.. ..- attr(*, "mcpair")= num [1:3] 1 10 1
$ m.h.performance: num [1:3] 7 7 9

```

The components are, if appropriate,

t.where: the time points which were chosen; the range of the Kaplan Meier estimate is divided by the number of grid points and transformed back to the time axis;

beta: samples of the vector of covariates;

cov.beta: samples of the covariance matrix of beta;

lbh.coef: samples of the log baseline hazard coefficients at the grid points;

sigma.lbh: samples of sigma.lbh.0 and sigma.lbh.1;

alpha.cluster: samples of the frailty values;

sigma.cluster: samples of frailty variance;

z.cluster: samples of the frailty values;

mu.cluster: samples of the rate and shape of the gamma prior;

m.h.performance: number of the successful performances of the Metropolis-Hastings step for beta, lbh and alpha.cluster or mu.cluster

The convergence is diagnosed by mean of CODA. The Raftery-Lewis diagnostic gives a good description of the convergence, see [5].

```
> raftery.diag(AA.res$beta)
```

```
Quantile (q) = 0.025
```

```
Accuracy (r) = +/- 0.005
```

```
Probability (s) = 0.95
```

You need a sample size of at least 3746 with these values of q, r and s

```
> raftery.diag(AA.res$cov.beta)
```

```
Quantile (q) = 0.025  
Accuracy (r) = +/- 0.005  
Probability (s) = 0.95
```

You need a sample size of at least 3746 with these values of q, r and s

```
> raftery.diag(AA.res$sigma.lbh)
```

```
Quantile (q) = 0.025  
Accuracy (r) = +/- 0.005  
Probability (s) = 0.95
```

You need a sample size of at least 3746 with these values of q, r and s

```
> raftery.diag(AA.res$sigma.cluster)
```

```
Quantile (q) = 0.025  
Accuracy (r) = +/- 0.005  
Probability (s) = 0.95
```

You need a sample size of at least 3746 with these values of q, r and s

```
> raftery.diag(AA.res$alpha.cluster)
```

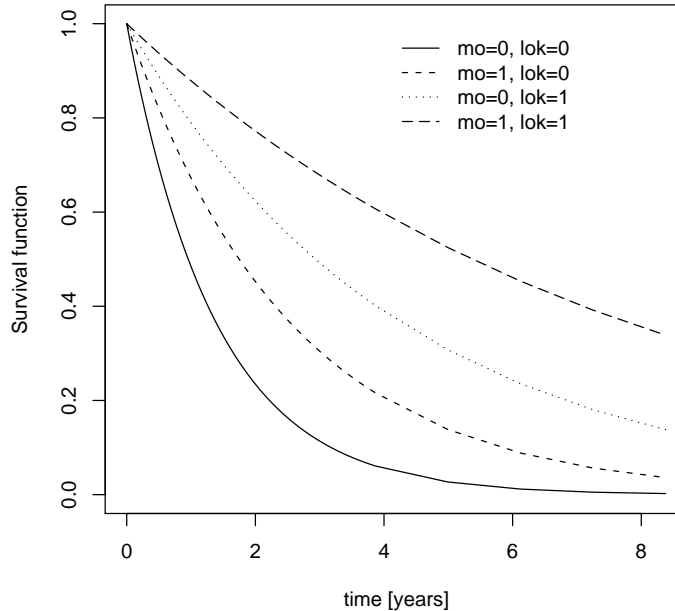
```
Quantile (q) = 0.025  
Accuracy (r) = +/- 0.005  
Probability (s) = 0.95
```

You need a sample size of at least 3746 with these values of q, r and s

This indicates that the sample size should be increased to at least 80000 samples.

The estimated coefficients and cumulative baseline hazard can be used to estimate and plot group specific survival curves.

```
> beta.est <- apply(AA.res$beta, 2, mean)  
> baseline.hazard <- survBayes.baseline.hazard(AA.res, start = 1,  
+       type = "cum")  
> time <- baseline.hazard$time  
> Lambda0 <- baseline.hazard$cum.base.haz  
> surv.base <- exp(-Lambda0)  
> plot(time, surv.base, type = "l", xlab = "time [years]", ylab = "Survival function",  
+       lty = 1, ylim = c(0, 1))  
> lines(time, surv.base*exp(beta.est["mo"]), type = "l", lty = 2)  
> lines(time, surv.base*exp(beta.est["lok"]), type = "l", lty = 3)  
> lines(time, surv.base*exp(sum(beta.est[c("mo", "lok")])), type = "l",  
+       lty = 5)  
> leg.names <- c("mo=0, lok=0", "mo=1, lok=0", "mo=0, lok=1", "mo=1, lok=1")  
> legend(4, 1, leg.names, lty = c(1, 2, 3, 5), bty = "n")
```



This work was supported by DFG grant MA 1723/2-1.

## References

- [1] M. Aitkin and D. Clayton. The fitting of exponential, Weibull and extreme value distributions to complex censored survival data using GLIM. *Applied Statistics*, 29:156–163, 1980.
- [2] D. Gamerman. Sampling from the posterior distribution in generalized linear mixed models. *Statistics and Computing*, 7:57–68, 1997.
- [3] L. Knorr-Held and H. Rue. On block updating in markov random field models for disease mapping. *Scandinavian Journal of Statistics*, 29:597–614, 2002.
- [4] S. Lang and A. Brezger. Bayesian p-splines. *Journal of Computational and Graphical Statistics*, 13:183–212, 2004.
- [5] U. Mansmann. Convergence Diagnosis for Gibbs Sampling Output. *Medical Infobahn for Europe*, A. Hasman et al. (Eds.), IOS Press, 83–87, 2000.
- [6] H. J. Meisel, U. Mansmann, H. Alvarez, G. Rodesch, M. Brock, and P. Lasjaunias. Cerebral arteriovenous malformations and associated aneurysms: Analysis of 305 cases from a series of 662 patients. *Neurosurgery*, 46:793–802, 2000.
- [7] H. Rue. Fast sampling of gaussian markov random fields. *Journal of the Royal Statistical Society B*, 63:325–338, 2001.

- [8] Terry M. Therneau and Patricia M. Grambsch. *Modeling Survival Data: extending the Cox model*. Springer, New York, 2000.
- [9] B. Turnbull. Nonparametric estimation of a survivorship function with doubly censored data. *Journal of the American Statistical Association*, 69:169–173, 1974.