

Joint longitudinal and time-to-event models via Stan

Sam Brilleman, Michael Crowther, Margarita Moreno-Betancur,
Jacqueline Bueros Novik, Rory Wolfe*

Abstract

The joint modelling of longitudinal and time-to-event data has received much attention in the biostatistical literature in recent years. In this notebook, we describe the implementation of a shared parameter joint model for longitudinal and time-to-event data in Stan. The methods described in the notebook are a simplified version of those underpinning the `stan_jm` modelling function that has recently been contributed to the **rstanarm** R package. This notebook will proceed as follows. In Section 1 we provide an introduction to the joint modelling of longitudinal and time-to-event data, including briefly describing the potential motivations for using such an approach. In Section 2 we describe the formulation of a multivariate shared parameter joint model and introduce its log likelihood function. In Section 3 we describe some of the more important features of the Stan code required to implement the model. In Section 4 we present a short applied example to demonstrate estimation of the model and the type of inferences that can be obtained. In Section 5 we close with a discussion.

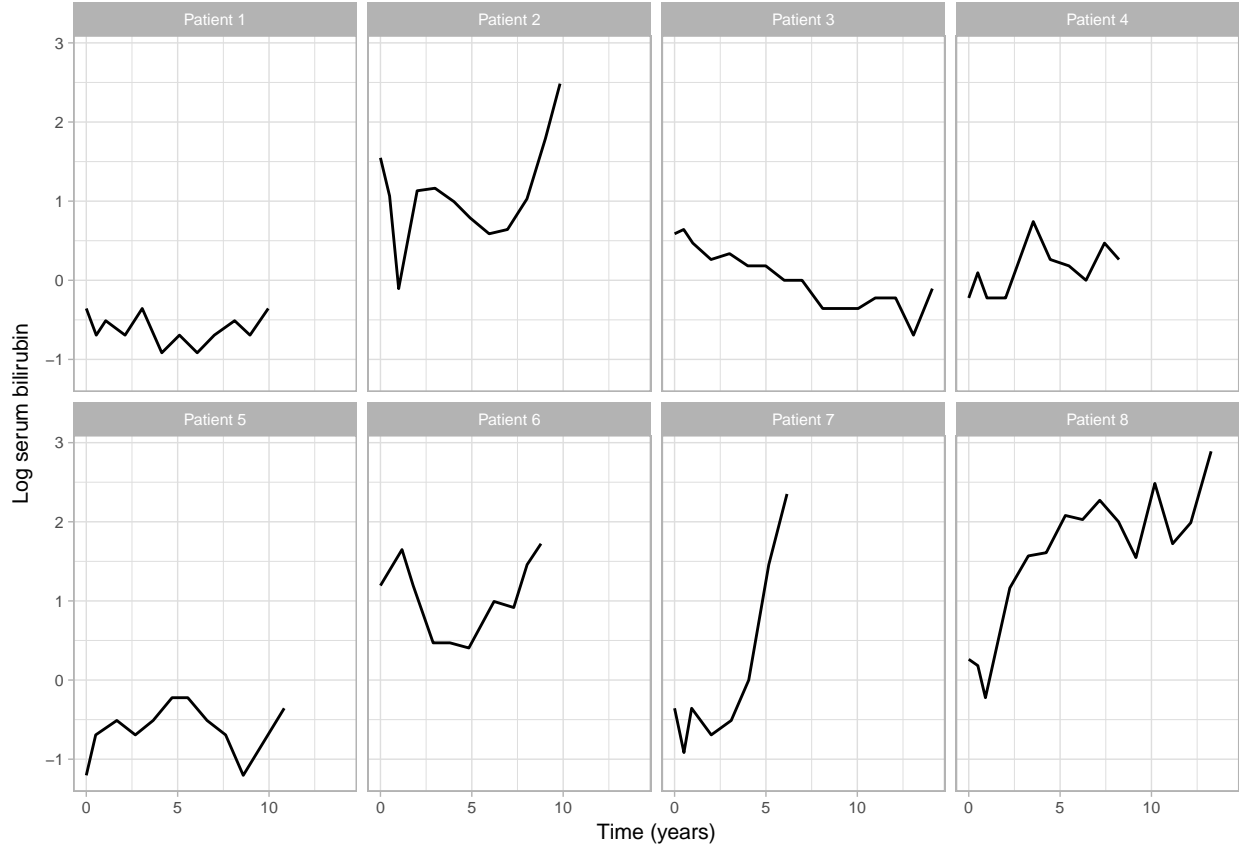
Date this notebook was compiled: 29 November 2017.

1 Introduction

Joint modelling can be broadly defined as the simultaneous estimation of two or more statistical models which traditionally would have been separately estimated. When we refer to a shared parameter joint model for longitudinal and time-to-event data, we generally mean the joint estimation of: 1) a longitudinal mixed effects model which analyses patterns of change in an outcome variable that has been measured repeatedly over time (for example, a clinical biomarker) and 2) a survival or time-to-event model which analyses the time until an event of interest occurs (for example, death or disease progression). Joint estimation of these so-called “submodels” is achieved by assuming they are correlated via individual-specific parameters (i.e. individual-level random effects).

Over the last two decades the joint modelling of longitudinal and time-to-event data has received a significant amount of attention [1-5]. Methodological developments in the area have been motivated by a growing awareness of the benefits that a joint modelling approach can provide. In clinical or epidemiological research it is common for a clinical biomarker to be repeatedly measured over time on a given patient. In addition, it is common for time-to-event data, such as the patient-specific time from a defined origin (e.g. time of diagnosis of a disease) until a terminating clinical event such as death or disease progression to also be collected. The figure below shows observed longitudinal measurements (i.e. observed “trajectories”) of log serum bilirubin for a small sample of patients with primary biliary cirrhosis. From the plots, we can observe between-patient variation in the longitudinal trajectories for log serum bilirubin, with some patients showing an increase in the biomarker over time, others decreasing, and some remaining stable. Moreover, there is variation between patients in terms of the frequency and timing of the longitudinal measurements.

*Corresponding author: sam.brilleman@monash.edu.



From the perspective of clinical risk prediction, we may be interested in asking whether between-patient variation in the log serum bilirubin trajectories provides meaningful prognostic information that can help us differentiate patients with regard to some clinical event of interest, such as death. Alternatively, from an epidemiological perspective we may wish to explore the potential for etiological associations between changes in log serum bilirubin and mortality. Joint modelling approaches provide us with a framework under which we can begin to answer these types of clinical and epidemiological questions.

More formally, the motivations for undertaking a joint modelling analysis of longitudinal and time-to-event data might include one or more of the following:

- One may be interested in how *underlying changes in the biomarker influence the occurrence of the event*. However, including the observed biomarker measurements directly into a time-to-event model as time-varying covariates poses several problems. For example, if the widely used Cox proportional hazards model is assumed for the time-to-event model then biomarker measurements need to be available for all patients at all failure times, which is unlikely to be the case [3]. If simple methods of imputation are used, such as the “last observation carried forward” method, then these are likely to induce bias [6]. Furthermore, the observed biomarker measurements may be subject to measurement error and therefore their inclusion as time-varying covariates may result in biased and inefficient estimates. In most cases, the measurement error will result in parameter estimates which are shrunk towards the null [7]. On the other hand, joint modelling approaches allow us to estimate the association between the biomarker (or some function of the biomarker trajectory, such as rate of change in the biomarker) and the risk of the event, whilst allowing for both the discrete time and measurement-error aspects of the observed biomarker.
- One may be interested primarily in the evolution of the clinical biomarker but *may wish to account for what is known as informative dropout*. If the value of future (unobserved) biomarker measurements are related to the occurrence of the terminating event, then those unobserved biomarker measurements will be “missing not at random” [8,9]. In other words, biomarker measurements for patients who have an

event will differ from those who do not have an event. Under these circumstances, inference based solely on observed measurements of the biomarker will be subject to bias. A joint modelling approach can help to adjust for informative dropout and has been shown to reduce bias in the estimated parameters associated with longitudinal changes in the biomarker [1,9,10].

- Joint models are naturally suited to the task of *dynamic risk prediction*. For example, joint modelling approaches have been used to develop prognostic models where predictions of event risk can be updated as new longitudinal biomarker measurements become available. Taylor et al. [11] jointly modelled longitudinal measurements of the prostate specific antigen (PSA) and time to clinical recurrence of prostate cancer. The joint model was then used to develop a web-based calculator which could provide real-time predictions of the probability of recurrence based on a patient's up to date PSA measurements.

In this notebook, we describe the implementation of a shared parameter joint model for longitudinal and time-to-event data in Stan. In Section 2 we describe the formulation for a multivariate joint model, that is, a joint model for multiple (i.e. more than one) longitudinal biomarkers and the time to a terminating event. In Section 3 we describe the important features of the Stan code required to fit the model. In Section 4 we present a brief applied example to demonstrate estimation of the model and the type of inferences that can be obtained. Note that the methods and code described in this paper are a simplified version of the `stan_jm` modelling function that is being contributed to the **rstanarm** R package [12,13], see <https://github.com/sambrilleman/rstanarm>.

2 Model formulation

A shared parameter joint model consists of related submodels which are specified separately for each of the longitudinal and time-to-event outcomes. These are therefore commonly referred to as the *longitudinal submodel(s)* and the *event submodel*. The longitudinal and event submodels are linked using shared individual-specific parameters, which can be parameterised in a number of ways. We describe each of these submodels below.

2.1 Longitudinal submodel(s)

We assume $y_{ijm}(t) = y_{im}(t_{ij})$ corresponds to the observed value of the m^{th} ($m = 1, \dots, M$) biomarker for individual i ($i = 1, \dots, N$) taken at time point t_{ij} , $j = 1, \dots, n_i$. We specify a (multivariate) generalised linear mixed model that assumes $y_{ijm}(t)$ follows a distribution in the exponential family with mean $\mu_{ijm}(t)$ and linear predictor

$$\eta_{ijm}(t) = g_m(\mu_{ijm}(t)) = \mathbf{x}_{ijm}^T(t)\boldsymbol{\beta}_m + \mathbf{z}_{ijm}^T(t)\mathbf{b}_{im} \quad (1)$$

where $\mathbf{x}_{ijm}^T(t)$ and $\mathbf{z}_{ijm}^T(t)$ are both row-vectors of covariates (which likely include some function of time, for example a linear slope, cubic splines, or polynomial terms) with associated vectors of fixed and individual-specific parameters $\boldsymbol{\beta}_m$ and \mathbf{b}_{im} , respectively, and g_m is some known link function.

The distribution and link function are allowed to differ over the M longitudinal submodels. We assume that the dependence across the different longitudinal submodel (i.e. the correlation between the different longitudinal biomarkers) is captured through a shared multivariate normal distribution for the individual-specific parameters; that is, we assume

$$\begin{pmatrix} \mathbf{b}_{i1} \\ \vdots \\ \mathbf{b}_{iM} \end{pmatrix} = \mathbf{b}_i \sim \text{Normal}(0, \boldsymbol{\Sigma}) \quad (2)$$

for some unstructured variance-covariance matrix Σ .

2.2 Event submodel

We assume that we also observe an event time $T_i = \min(T_i^*, C_i)$ where T_i^* denotes the so-called “true” event time for individual i (potentially unobserved) and C_i denotes the censoring time. We define an event indicator $d_i = I(T_i^* \leq C_i)$. We then model the hazard of the event using a parametric proportional hazards regression model of the form

$$h_i(t) = h_0(t) \exp \left(\mathbf{w}_i^T(t) \boldsymbol{\gamma} + \sum_{m=1}^M \sum_{q=1}^{Q_m} \alpha_{mq} f_{mq}(\boldsymbol{\beta}_m, \mathbf{b}_{im}; t) \right) \quad (3)$$

where $h_i(t)$ is the hazard of the event for individual i at time t , $h_0(t)$ is the baseline hazard at time t , $\mathbf{w}_i^T(t)$ is a row-vector of individual-specific covariates (possibly time-dependent) with an associated vector of regression coefficients $\boldsymbol{\gamma}$ (log hazard ratios), and the α_{mq} are also coefficients (log hazard ratios).

The longitudinal and event submodels are assumed to related via an “association structure” based on shared individual-specific parameters and captured via the $\sum_{m=1}^M \sum_{q=1}^{Q_m} \alpha_{mq} f_{mq}(\boldsymbol{\beta}_m, \mathbf{b}_{im}; t)$ term in the linear predictor of equation (3). The coefficients α_{mq} are referred to as the “association parameters” since they quantify the strength of the association between the longitudinal and event processes, while the $f_{mq}(\boldsymbol{\beta}_m, \mathbf{b}_{im}; t)$ (for some functions $f_{mq}(\cdot)$) can be referred to as the “association terms” and can be specified in a variety of ways which we describe in the next section.

We assume that the baseline hazard $h_0(t)$ is modelled parametrically. For simplicity, in the formulation of the joint model presented in this notebook we will restrict ourselves to modelling the log baseline hazard using B-splines. Note however that in the **rstanarm** package’s `stan_jm` modelling function the baseline hazard can be specified as either an approximation using B-splines (the default), a Weibull distribution, or a piecewise constant baseline hazard (sometimes referred to as piecewise exponential). In the case of the piecewise constant or B-splines baseline hazard, the user can control the flexibility by explicitly specifying the knot points or degrees of freedom.

2.3 Association structures

As mentioned in the previous section, the dependence between the longitudinal and event submodels is captured through the association structure, which can be specified in a number of ways. In this notebook, we focus on the simplest association structure

$$f_{mq}(\boldsymbol{\beta}_m, \mathbf{b}_{im}; t) = \eta_{im}(t) \quad (4)$$

This is often referred to as a *current value* association structure since it assumes that the log hazard of the event at time t is linearly associated with the value of the longitudinal submodel’s linear predictor also evaluated at time t . This is the most common association structure used in the joint modelling literature to date. In the situation where the longitudinal submodel is based on an identity link function and normal error distribution (i.e. a linear mixed model) the *current value* association structure can be viewed as a method for including the underlying “true” value of the biomarker as a time-varying covariate in the event submodel.¹

¹By “true” value of the biomarker, we mean the value of the biomarker which is not subject to measurement error or discrete time observation. Of course, for the expected value from the longitudinal submodel to be considered the so-called “true” underlying biomarker value, we would need to have specified the longitudinal submodel appropriately!

However, there are a variety of other association structures that could be specified. For example, we could assume the log hazard of the event is linearly associated with the *current slope* (i.e. rate of change) of the longitudinal submodel’s linear predictor, that is

$$f_{mq}(\boldsymbol{\beta}_m, \mathbf{b}_{im}; t) = \frac{d\eta_{im}(t)}{dt} \quad (5)$$

Moreover, there are a whole range of possible association structures, many of which have been discussed in the literature [14-16]. Also note that the full set of association structures that are accommodated in the **rstanarm** package’s **stan_jm** modelling function are not described here but are discussed in the documentation for the **stan_jm** function itself.

2.4 Conditional independence assumption

A key assumption of the multivariate shared parameter joint model is that the observed longitudinal measurements are independent of one another (both across the M biomarkers and across the n_i time points), as well as independent of the event time, conditional on the individual-specific parameters \mathbf{b}_i . That is, we assume

$$\text{Cov}(y_{im}(t), y_{im'}(t) | \mathbf{b}_i) = 0 \quad (6)$$

$$\text{Cov}(y_{im}(t), y_{im}(t') | \mathbf{b}_i) = 0 \quad (7)$$

$$\text{Cov}(y_{im}(t), T_i | \mathbf{b}_i) = 0 \quad (8)$$

for some $m \neq m'$ and $t \neq t'$.

Although this may be considered a strong assumption, it is useful in that it allows the full likelihood for joint model to be factorised into the likelihoods for each of the component parts (i.e. the likelihoods for each of the submodels and the likelihood for the distribution of the individual-specific parameters).

2.5 Log posterior distribution

Under the conditional independence assumption, the log posterior for the i^{th} individual can be specified as

$$p(\boldsymbol{\theta}, \mathbf{b}_i | \mathbf{y}_i, T_i, d_i) \propto \log \left[\left(\prod_{m=1}^M \prod_{j=1}^{n_i} p(y_{ijm} | \mathbf{b}_i, \boldsymbol{\theta}) \right) p(T_i, d_i | \mathbf{b}_i, \boldsymbol{\theta}) p(\mathbf{b}_i | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \right] \quad (9)$$

which we can rewrite as

$$p(\boldsymbol{\theta}, \mathbf{b}_i | \mathbf{y}_i, T_i, d_i) \propto \left(\sum_{m=1}^M \sum_{j=1}^{n_i} \log p(y_{ijm} | \mathbf{b}_i, \boldsymbol{\theta}) \right) + \log p(T_i, d_i | \mathbf{b}_i, \boldsymbol{\theta}) + \log p(\mathbf{b}_i | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (10)$$

where $\sum_{j=1}^{n_i} \log p(y_{ijm} | \mathbf{b}_i, \boldsymbol{\theta})$ is the log likelihood for the m^{th} longitudinal submodel, $\log p(T_i, d_i | \mathbf{b}_i, \boldsymbol{\theta})$ is the log likelihood for the event submodel, $\log p(\mathbf{b}_i | \boldsymbol{\theta})$ is the log likelihood for the distribution of the group-specific

parameters (i.e. random effects), and $\log p(\boldsymbol{\theta})$ represents the log likelihood for the joint prior distribution across all remaining unknown parameters.²

We can rewrite the log likelihood for the event submodel as

$$\log p(T_i, d_i | \mathbf{b}_i, \boldsymbol{\theta}) = d_i * \log h_i(T_i) - \int_0^{T_i} h_i(s) ds \quad (11)$$

and then use Gauss-Kronrod quadrature with Q nodes to approximate $\int_0^{T_i} h_i(s) ds$, such that

$$\int_0^{T_i} h_i(s) ds \approx \frac{T_i}{2} \sum_{q=1}^Q w_q h_i\left(\frac{T_i(1+s_q)}{2}\right) \quad (12)$$

where w_q and s_q , respectively, are the standardised weights and locations (“abscissa”) for quadrature node q ($q = 1, \dots, Q$) [17]. In this notebook we choose to use $Q = 15$ quadrature nodes.³

Therefore, once we have an individual’s event time T_i we can evaluate the design matrices for the event submodel and longitudinal submodels at the $Q + 1$ necessary time points (which are the event time T_i and the quadrature points $\frac{T_i(1+s_q)}{2}$ for $q = 1, \dots, Q$) and then pass these to Stan’s data block. We can then evaluate the log likelihood for the event submodel by simply calculating the hazard $h_i(t)$ at those $Q + 1$ time points and summing the quantities appropriately. This calculation will need to be performed each time we iterate through Stan’s model block. The Stan code required to evaluate this log posterior is described in the next section.

3 Stan code

Here we describe the most important features of the Stan code used to estimate the joint model. The full Stan code is provided in the separate `jm.stan` file supplied with this notebook.

3.1 Data and transformed data blocks

The data block includes dimensions of the model, outcome vectors (observed biomarker values, event times, event indicators), design matrices for the different submodels, and hyperparameters for the prior distributions. We do not discuss the data or transformed data blocks here in any detail.

3.2 Parameters block

Most of the parameters defined in the parameters block are “primitive” or “unscaled”, meaning that they will be given a prior distribution with mean 0 and scale 1 and then converted into the actual parameters used in the regression equation via the transformed parameters block. Our parameters block therefore includes:

²In this notebook we assume normal prior distributions for all unbounded parameters (e.g. regression coefficients) and half-normals for all lower-bounded parameters (e.g. standard deviations). In the **rstanarm** package there is a variety of prior distributions available. Moreover, the prior for the variance-covariance of individual-specific parameters is taken from the **stan_glm** modelling function in the **rstanarm** package; we refer the reader to the documentation of that package for those details.

³The `stan_jm` modelling function in the **rstanarm** package allows the user to choose between $Q = 15$ (the default), 11, or 7 quadrature nodes.

- **gamma**: the intercept for each of the longitudinal submodels, combined into a single vector. These intercept parameters are unbounded, given that each longitudinal submodel in our application consists of a linear mixed model (i.e. in our application we assume an identity link function and normal error distribution for each longitudinal biomarker).
- **y_z_beta**: the primitive coefficients for each of the longitudinal submodels, combined into a single vector.
- **e_z_beta**, **a_z_beta**: the primitive coefficients and primitive association parameters for the event submodel.
- **y_aux_unscaled**: the unscaled standard deviations (SD) of the residual errors for each of the longitudinal submodels, combined into a single vector.
- **e_aux_unscaled**: the unscaled coefficients for the B-spline terms used in the baseline hazard.

The parameters block also includes the group-specific parameters, defined in the same way as for the **stan_glmer** modelling function in the **rstanarm** package. Since the specification of the group-specific terms is not the primary focus of this notebook we do not include the details here. Further details on specification of the group-specific terms can be found in the full **jm.stan** file supplied with this notebook, or in the help documentation for the **rstanarm** package.

```
parameters {
  vector[M] gamma;          // intercepts in long. submodels
  vector[y_K] y_z_beta;    // primitive coefs in long. submodels
  vector[e_K] e_z_beta;    // primitive coefs in event submodel (log hazard ratios)
  vector[a_K] a_z_beta;    // primitive assoc params (log hazard ratios)
  vector<lower=0>[M] y_aux_unscaled; // unscaled residual error SDs
  vector[basehaz_df] e_aux_unscaled; // unscaled coefs for baseline hazard
  ...
}
```

3.3 Transformed parameters block

The transformed parameters block includes code to alter the location and scale of the “primitive” or “unscaled” parameters, in order to obtain the actual parameters used in the regression submodels:

```
transformed parameters {
  ...
  // coefficients and association parameters
  y_beta = y_z_beta .* y_prior_scale + y_prior_mean;
  e_beta = e_z_beta .* e_prior_scale + e_prior_mean;
  a_beta = a_z_beta .* a_prior_scale + a_prior_mean;

  // auxiliary parameters
  y_aux = y_aux_unscaled .* y_prior_scale_for_aux + y_prior_mean_for_aux;
  e_aux = e_aux_unscaled .* e_prior_scale_for_aux + e_prior_mean_for_aux;
  ...
}
```

3.4 Model block

The model block consists of several distinct parts. We describe each of these separately.

In the first part of the model block, we evaluate the combined linear predictor for all longitudinal submodels, excluding their intercept terms. We then loop over the $1 : M$ longitudinal submodels, and for each submodel

m we add the intercept term `gamma[m]` onto the linear predictor, and then increment the target with the resulting likelihood.

```
model {
  //----- Log-lik for longitudinal submodels
  {
    vector[N] y_eta; // linear predictor for long. submodels

    // evaluate linear predictor for all long. submodels
    y_eta = X * y_beta + csr_matrix_times_vector(N, q, w, v, u, b);

    // evaluate linear predictor for just submodel m and accumulate log-lik
    for (m in 1:M) {
      vector[NM[m]] y_eta_m;
      y_eta_m = y_eta[idx[m,1]:idx[m,2]] + gamma[m];
      target += normal_lpdf(y[idx[m,1]:idx[m,2]] | y_eta_m, y_aux[m]);
    }
  }
  ...
}
```

To evaluate the event submodel likelihood we must evaluate $h_i(T_i)$ (i.e. the hazard for individual i at their event or censoring time) as well as the cumulative hazard $\int_0^{T_i} h_i(s)ds$. Since we are going to evaluate the cumulative hazard using Gauss-Kronrod quadrature, this means calculating the hazard $h_i(t)$ at T_i and 15 quadrature points between 0 and T_i . To do this, we have constructed the design matrices in R evaluated at the necessary times; these are passed to Stan's data block (not shown here) as `e_Xq`, `y_Xq`, etc. In the code below there are several steps:

- In **Step 1** we use the event submodel design matrices to evaluate the $\mathbf{w}_i^T(t)\boldsymbol{\gamma}$ part of the event submodel's linear predictor at the event time T_i and the 15 quadrature points between 0 and T_i .
- The remainder of the event submodel's linear predictor consists of the term corresponding to the association structure: $\sum_{m=1}^M \alpha_m \eta_{im}(t)$. This involves the current value of the longitudinal submodel's linear predictor, so we must also evaluate the longitudinal submodel's linear predictor at the event time T_i and the 15 quadrature points between 0 and T_i . This is shown in **Step 2** of the code below.
- In **Step 3** we evaluate the log baseline hazard at the event time T_i and the 15 quadrature points between 0 and T_i .
- In **Step 4** we combine the log baseline hazard with the event submodel linear predictor, that is, we evaluate

$$\log h_i(t) = \log h_0(t) + \left(\mathbf{w}_i^T(t)\boldsymbol{\gamma} + \sum_{m=1}^M \alpha_m \eta_{im}(t) \right)$$

- In **Steps 5** and **6** we use the evaluated hazards to construct the likelihood for the event submodel: the hazard evaluated at T_i is separated out from the hazard evaluated at each of the quadrature points between 0 and T_i , and the latter are used to evaluate the approximate cumulative hazard at the event time via the Gauss-Kronrod quadrature rule described in equation (12).
- In **Step 7** we evaluate the log likelihood for the event submodel as

$$\log p(T_i, d_i | \mathbf{b}_i, \boldsymbol{\theta}) = d_i * \log h_i(T_i) - \int_0^{T_i} h_i(s)ds$$

and increment the target with this result.

```
//----- Log-lik for event submodel (Gauss-Kronrod quadrature)
{
  vector[sum(nrow_y_Xq)] y_eta_q;
  vector[nrow_e_Xq] e_eta_q;
```



```

vector[nrow_e_Xq] log_basehaz;
vector[nrow_e_Xq] ll_haz_q;
vector[Npat] ll_haz_eventtime;
vector[Npat_times_quadnodes] ll_surv_eventtime;
real ll_event;

// Step 1: event submodel linear predictor at event time and quadrature points
e_eta_q = e_Xq * e_beta;

// Step 2: long. submodel linear predictor at event time and quadrature points
y_eta_q = y_Xq * y_beta +
  csr_matrix_times_vector(sum(nrow_y_Xq), q, w_Zq, v_Zq, u_Zq, b);

// Step 2 (continued): add on contribution from association structure to
// the event submodel linear predictor at event time and quadrature points
for (m in 1:M) {
  vector[nrow_y_Xq[m]] y_eta_q_m;
  y_eta_q_m = y_eta_q[idx_q[m,1]:idx_q[m,2]] + gamma[m];
  e_eta_q = e_eta_q + a_beta[m] * y_eta_q_m;
}

// Step 3: log baseline hazard at event time and quadrature points
log_basehaz = basehaz_X * e_aux;

// Step 4: log hazard at event time and quadrature points
ll_haz_q = e_d .* (log_basehaz + e_eta_q);

// Step 5: log hazard contribution to the likelihood
ll_haz_eventtime = segment(ll_haz_q, 1, Npat);

// Step 6: log survival contribution to the likelihood, obtained by
// summing over the quadrature points to get the approximate integral
// (NB quadweight already incorporates the (b-a)/2 scaling such that the
// integral is evaluated over limits (a,b) rather than (-1,+1))
ll_surv_eventtime = quadweight .*
  exp(segment(ll_haz_q, Npat + 1, Npat_times_quadnodes));

// Step 7: log likelihood for event submodel
ll_event = sum(ll_haz_eventtime) - sum(ll_surv_eventtime);
target += ll_event;
}

```

We then increment the target with the log priors for each of the intercepts, coefficients, auxiliary parameters (including coefficients for the B-splines baseline hazard), and group-specific terms (i.e. individual-level random effects).

```

//----- Log-priors

// intercepts
for (m in 1:M)
  target += normal_lpdf(y_gamma[m] |
    y_prior_mean_for_intercept, y_prior_scale_for_intercept);

```

```

// coefficients
target += normal_lpdf(y_z_beta | 0, 1);
target += normal_lpdf(e_z_beta | 0, 1);
target += normal_lpdf(a_z_beta | 0, 1);

// auxiliary parameters
target += normal_lpdf(y_aux_unscaled | 0, 1);
target += normal_lpdf(e_aux_unscaled | 0, 1);

// group-specific parameters
decov_lp(z_b, z_T, rho, zeta, tau, regularization, delta, shape, t, p);
}

```

4 Application

4.1 Data

In order to make this notebook freely available we use a motivating example based on a publically accessible dataset. The Mayo Clinic’s widely used primary biliary cirrhosis (PBC) data contains 312 individuals with primary biliary cirrhosis, who participated in a randomised placebo controlled trial of D-penicillamine conducted at the Mayo Clinic between 1974 and 1984 [18]. In our secondary analysis of this trial data, our primary research is *not* concerned with the efficacy of the randomised treatment but rather understanding how the clinical biomarker histories for these patients are associated with their overall survival. Specifically, we focus on the associations between two repeatedly measured clinical biomarkers, log serum bilirubin and serum albumin, and the risk of death. Given that the joint modelling methods are computationally intensive we restrict our analyses to a small random subset of just 40 patients from the PBC dataset. This ensures that the computation time for the joint models described in later sections are kept to a minimum and therefore this notebook can be compiled in a relatively short time. However, this also means that the clinical findings from this analysis should not to be overinterpreted. Rather, this notebook aims to simply demonstrate the joint modelling framework and describe how these models can be estimated using Stan.

The PBC data are contained in two separate data frames, each saved as an RDS object. The first data frame (saved as “Data/pbcLong.rds”), contains multiple-row per patient longitudinal biomarker information, as shown in

```
head(pbcLong)
```

##	id	age	sex	trt	year	logBili	albumin	platelet
## 1	1	58.76523	f	1	0.0000000	2.67414865	2.60	190
## 2	1	58.76523	f	1	0.5256674	3.05870707	2.94	183
## 3	2	56.44627	f	1	0.0000000	0.09531018	4.14	221
## 4	2	56.44627	f	1	0.4982888	-0.22314355	3.60	188
## 5	2	56.44627	f	1	0.9993155	0.00000000	3.55	161
## 6	2	56.44627	f	1	2.1026694	0.64185389	3.92	122

while the second data frame (saved as “Data/pbcSurv.rds”), contains single-row per patient survival information, as shown in

```
head(pbcSurv)
```

##	id	age	sex	trt	futimeYears	status	death
## 1	1	58.76523	f	1	1.095140	2	1
## 3	2	56.44627	f	1	14.151951	0	0
## 12	3	70.07255	m	1	2.770705	2	1

```
## 16  4 54.74059   f   1    5.270363      2    1
## 23  5 38.10541   f   0    4.120465      1    0
## 29  6 66.25873   f   0    6.852841      2    1
```

The variables included across the two datasets can be defined as follows:

- `age` in years
- `albumin` serum albumin (g/dl)
- `logBili` logarithm of serum bilirubin
- `death` indicator of death at endpoint
- `futimeYears` time (in years) between baseline and the earliest of death, transplantation or censoring
- `id` numeric ID unique to each individual
- `platelet` platelet count
- `sex` gender (m = male, f = female)
- `status` status at endpoint (0 = censored, 1 = transplant, 2 = dead)
- `trt` binary treatment code (0 = placebo, 1 = D-penicillamine)
- `year` time (in years) of the longitudinal measurements, taken as time since baseline

4.2 Estimation using the simplified `jm.stan` file

We fit a multivariate joint model to the two longitudinal biomarkers, log serum bilirubin and serum albumin, and time-to-death. Note that patients are censored if they had a transplant prior to death (here we ignore the fact that this is likely to be informative censoring). We fit a linear mixed model (identity link, normal distribution) for each biomarker with a patient-specific intercept and linear slope but no other covariates. In the event submodel we include gender (`sex`) and treatment (`trt`) as baseline covariates. Each biomarker is assumed to be associated with the log hazard of death at time t via it's expected value at time t (i.e. a *current value* association structure).

To save needing to carry out any data manipulation steps we instead used the `stan_jm` modelling function in **rstanarm** to generate the R list for passing to **rstan**. This data is saved as an RDS object and supplied with the notebook ("Stan/standata.rds"). In addition, a function to generate a list of initial values has also been supplied as an RDS object with the notebook ("Stan/staninit.rds"). Of course, the stan file containing the model is also supplied ("Stan/jm.stan"). We can therefore estimate this model using the **rstan** package:

```
standata <- readRDS("Stan/standata.rds")
staninit <- readRDS("Stan/staninit.rds")
mod1 <- with_filecache(
  stan(
    file = "Stan/jm.stan",
    data = standata,
    init = function() staninit,
    control = list(adapt_delta = 0.95),
    chains = 3, cores = 3, iter = 2000,
    seed = 12345, refresh = 20),
  filename = "mod1.rds")
```

Since our primary interest is in the association between the current value of each of the biomarkers (log serum bilirubin and serum albumin) and the hazard of death, we focus on the estimated association parameters. The summary of the posterior distribution for each of the association parameters follows:

```
print(mod1, pars = "a_beta")
```

```
## Inference for Stan model: jm.
## 3 chains, each with iter=2000; warmup=1000; thin=1;
```

```
## post-warmup draws per chain=1000, total post-warmup draws=3000.
##
##           mean se_mean   sd  2.5%  25%   50%   75% 97.5% n_eff Rhat
## a_beta[1]  0.89     0.01 0.28  0.38  0.70  0.88  1.07  1.46 3000   1
## a_beta[2] -2.81     0.03 0.80 -4.45 -3.34 -2.75 -2.25 -1.43  715   1
##
## Samples were drawn using NUTS(diag_e) at Wed Nov 29 10:12:37 2017.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
```

We see that a one unit increase in log serum bilirubin is associated with a 0.89 (95% CrI: 0.38 to 1.46) unit increase in the log hazard of death, equivalent to a 2.44-fold (95% CrI: 1.46 to 4.31) *increase* in the hazard of death. Similarly, a one unit increase in serum albumin is associated with a 2.81 (95% CrI: 1.43 to 4.45) unit *decrease* in the log hazard of death. These estimates are broadly in line with what we would expect from a clinical perspective; that is, that higher serum bilirubin is associated with *worse* patient outcomes (i.e. higher risk of mortality), whilst higher serum albumin is associated with *better* patient outcomes (i.e. lower risk of mortality). However, recall that we have estimated this model with a very small dataset only used for demonstration purposes. Moreover, the number of mortality events ($N = 29$) is even less than the number of patients since some patients are censored. In fact, if we were to not suppress Stan's warning messages, we would see that there is a number of divergent transitions when fitting the model, potentially owing to the limited amount of data we are using to estimate this model.

4.3 Estimation using the joint modelling functionality in rstanarm

The `jm.stan` file provided with this notebook is a simplified version of the Stan code underlying the `stan_jm` modelling function in the `rstanarm` package. However, estimating the model using the `rstanarm` provides us with much nicer output (for example, meaningful variable names!) as well as a broad range of post-estimation functionality, including model diagnostics, posterior predictions, dynamic predictions and more.

To see this, we can use the development version of `rstanarm` with joint modelling functionality to refit our model, this time using `stan_jm` with the customary R formula syntax and data frames:

```
mod2 <- with_filecache(
  stan_jm(
    formulaLong = list(
      logBili ~ year + (year | id),
      albumin ~ year + (year | id)),
    formulaEvent = survival::Surv(futimeYears, death) ~ sex + trt,
    dataLong = pbcLong, dataEvent = pbcSurv,
    assoc = "etavalue", time_var = "year", basehaz = "bs",
    chains = 3, cores = 3, iter = 2000, seed = 12345, refresh = 20),
  filename = "mod2.rds")
```

We can now examine the output from the fitted model, for example

```
print(mod2)

## stan_jm
## formula (Long1): logBili ~ year + (year | id)
## family (Long1): gaussian [identity]
## formula (Long2): albumin ~ year + (year | id)
## family (Long2): gaussian [identity]
## formula (Event): survival::Surv(futimeYears, death) ~ sex + trt
```

```

## baseline hazard: bs
## assoc:          etavalue (Long1), etavalue (Long2)
## -----
##
## Longitudinal submodel 1: logBili
##           Median MAD_SD
## (Intercept) 0.672  0.186
## year        0.227  0.042
## sigma       0.354  0.017
##
## Longitudinal submodel 2: albumin
##           Median MAD_SD
## (Intercept)  3.520  0.085
## year        -0.160  0.025
## sigma       0.290  0.013
##
## Event submodel:
##           Median MAD_SD exp(Median)
## (Intercept)    6.736  2.854 841.796
## sexf          -0.141  0.673  0.868
## trt           -0.497  0.493  0.609
## Long1|etavalue  0.796  0.286  2.217
## Long2|etavalue -3.045  0.895  0.048
## b-splines-coef1 -0.859  1.064    NA
## b-splines-coef2  0.553  0.895    NA
## b-splines-coef3 -2.556  1.267    NA
## b-splines-coef4 -0.467  1.774    NA
## b-splines-coef5 -1.164  1.780    NA
## b-splines-coef6 -2.586  1.876    NA
##
## Group-level error terms:
## Groups Name          Std.Dev. Corr
## id      Long1|(Intercept) 1.2540
##          Long1|year      0.1948  0.52
##          Long2|(Intercept) 0.5068 -0.65 -0.52
##          Long2|year      0.1022 -0.60 -0.82  0.47
## Num. levels: id 40
##
## Sample avg. posterior predictive distribution
## of longitudinal outcomes:
##           Median MAD_SD
## Long1|mean_PPD 0.588  0.030
## Long2|mean_PPD 3.344  0.023
##
## -----
## For info on the priors used see help('prior_summary.stanreg').

```

or we can examine the summary output for the association parameters alone:

```
summary(mod2, pars = "assoc")
```

```

##
## Model Info:
##

```

```

## function:      stan_jm
## formula (Long1): logBili ~ year + (year | id)
## family (Long1): gaussian [identity]
## formula (Long2): albumin ~ year + (year | id)
## family (Long2): gaussian [identity]
## formula (Event): survival::Surv(futimeYears, death) ~ sex + trt
## baseline hazard: bs
## assoc:        etavalue (Long1), etavalue (Long2)
## algorithm:    sampling
## priors:       see help('prior_summary')
## sample:       3000 (posterior sample size)
## num obs:      304 (Long1), 304 (Long2)
## num subjects: 40
## num events:   29 (72.5%)
## groups:      id (40)
## runtime:     1.2 mins
##
## Estimates:
##              mean    sd    2.5%   25%    50%    75%    97.5%
## Assoc|Long1|etavalue  0.798  0.293  0.223  0.601  0.796  0.987  1.397
## Assoc|Long2|etavalue -3.103  0.891 -5.011 -3.686 -3.045 -2.471 -1.527
##
## Diagnostics:
##              mcse  Rhat  n_eff
## Assoc|Long1|etavalue 0.005 1.003 3000
## Assoc|Long2|etavalue 0.021 1.002 1828
##
## For each parameter, mcse is Monte Carlo standard error, n_eff is a crude measure of effective sample

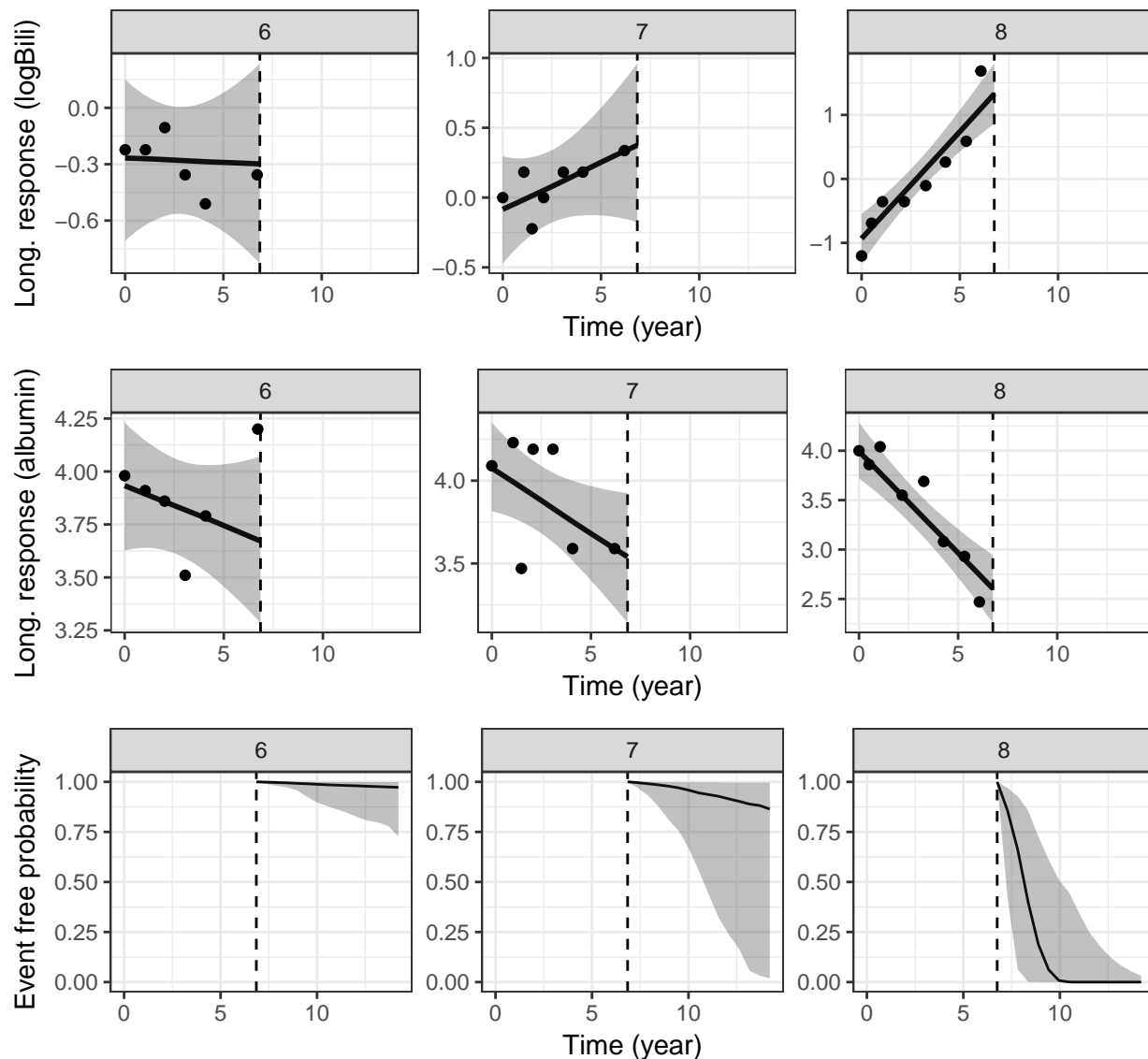
```

We can see that the estimated association parameters are similar to those obtained from the model in the previous section. However, we can now also access a range of post-estimation functions (described in the `stan_jm` and related help documentation; see for example `help(posterior_traj)` or `help(posterior_survfit)`). As an example, let's plot the predicted trajectories for each biomarker and the predicted survival function for three selected individuals in the dataset using `stan_jm` post-estimation functions:

```

p1 <- posterior_traj(mod2, m = 1, ids = 6:8)
p2 <- posterior_traj(mod2, m = 2, ids = 6:8)
p3 <- posterior_survfit(mod2, ids = 6:8, draws = 200)
pp1 <- plot(p1, plot_observed = TRUE, vline = TRUE)
pp2 <- plot(p2, plot_observed = TRUE, vline = TRUE)
plot_stack_jm(yplot = list(pp1, pp2), survplot = plot(p3))

```



Here we can see the strong relationship between the underlying values of the biomarkers and mortality. Patient 8 who, relative to patients 6 and 7, has a higher underlying value for log serum bilirubin and a lower underlying value for serum albumin at the end of their follow up has a far worse predicted probability of survival.

5 Discussion

In this notebook we have introduced the formulation of a shared parameter joint model for longitudinal and time-to-event data. The formulation of the joint model can allow for multiple longitudinal biomarkers along with a terminating event. The association between the longitudinal and event processes can be parameterised in a variety of ways, but here we have focussed on the so-called *current value* association structure which serves as the simplest and natural starting point.

The aim of this notebook was to introduce some of the ideas underpinning the estimation of these joint models in Stan. One key feature of the Stan code that we have tried to describe in detail is the implementation of the Gauss-Kronrod quadrature rule. The Gauss-Kronrod quadrature rule is required to approximate the

cumulative hazard in the likelihood of the event submodel. This aspect makes evaluating the log likelihood for the event submodel more computationally intensive than if there were a closed-form solution to the integral. In addition, the models are computationally intensive due to the relatively large number of group-specific parameters that often need to be estimated. Nonetheless, estimating joint models under a Bayesian framework can provide a number of benefits. The specification of complex association structures can be made much easier. Furthermore, a Bayesian approach can lead more naturally to dynamic predictions. For these, and other reasons, we believe it is of interest to try and optimise the estimation of these models in Stan. The hope is that by describing the Stan code in some detail as part of this notebook, those reading it will have the opportunity to provide guidance on how increases in speed, efficiency, or numerical stability might be achieved.

6 Acknowledgements

Much of the joint modelling functionality that has been contributed to the **rstanarm** package has been built upon code that was already included in that package, and that code was written primarily by Ben Goodrich and Jonah Gabry [12]. We are also grateful to them for their ongoing support in helping to get the joint modelling functionality up and running in **rstanarm**.

7 References

1. Henderson R, Diggle P, Dobson A. Joint modelling of longitudinal measurements and event time data. *Biostatistics* 2000;**1**(4):465-80.
2. Wulfsohn MS, Tsiatis AA. A joint model for survival and longitudinal data measured with error. *Biometrics* 1997;**53**(1):330-9.
3. Tsiatis AA, Davidian M. Joint modeling of longitudinal and time-to-event data: An overview. *Stat Sinica* 2004;**14**(3):809-34.
4. Gould AL, Boye ME, Crowther MJ, Ibrahim JG, Quartey G, Micallef S, et al. Joint modeling of survival and longitudinal non-survival data: current methods and issues. Report of the DIA Bayesian joint modeling working group. *Stat Med*. 2015;**34**(14):2181-95.
5. Rizopoulos D. *Joint Models for Longitudinal and Time-to-Event Data: With Applications in R* CRC Press; 2012.
6. Liu G, Gould AL. Comparison of alternative strategies for analysis of longitudinal trials with dropouts. *J Biopharm Stat* 2002;**12**(2):207-26.
7. Prentice RL. Covariate Measurement Errors and Parameter-Estimation in a Failure Time Regression-Model. *Biometrika* 1982;**69**(2):331-42.
8. Baraldi AN, Enders CK. An introduction to modern missing data analyses. *J Sch Psychol* 2010;**48**(1):5-37.
9. Philipson PM, Ho WK, Henderson R. Comparative review of methods for handling drop-out in longitudinal studies. *Stat Med* 2008;**27**(30):6276-98.
10. Pantazis N, Touloumi G. Bivariate modelling of longitudinal measurements of two human immunodeficiency type 1 disease progression markers in the presence of informative drop-outs. *Applied Statistics* 2005;**54**:405-23.
11. Taylor JM, Park Y, Ankerst DP, et al. Real-time individual predictions of prostate cancer recurrence using joint models. *Biometrics* 2013;**69**(1):206-13.
12. Stan Development Team. *rstanarm: Bayesian applied regression modeling via Stan*. R package version 2.14.1. <http://mc-stan.org/>. 2016.
13. R Core Team. *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing; 2015.
14. Crowther MJ, Lambert PC, Abrams KR. Adjusting for measurement error in baseline prognostic biomarkers included in a time-to-event analysis: a joint modelling approach. *BMC Med Res Methodol* 2013;**13**.

15. Hickey GL, Philipson P, Jorgensen A, Kolamunnage-Dona R. Joint modelling of time-to-event and multivariate longitudinal outcomes: recent developments and issues. *BMC Med Res Methodol* 2016;**16**(1):117.
16. Rizopoulos D, Ghosh P. A Bayesian semiparametric multivariate joint model for multiple longitudinal outcomes and a time-to-event. *Stat Med.* 2011;**30**(12):1366-80.
17. Laurie DP. Calculation of Gauss-Kronrod quadrature rules. *Math Comput* 1997;**66**(219):1133-45.
18. Therneau T, Grambsch P. *Modeling Survival Data: Extending the Cox Model* Springer-Verlag, New York; 2000. ISBN: 0-387-98784-3