

Efficiently analyzing large patient registries with Bayesian joint models for longitudinal and time-to-event data

Pedro Miranda Afonso^{1,2,*}, Dimitris Rizopoulos^{1,2}, Anushka K. Palipana³,
Grace C. Zhou^{3,4}, Cole Brokamp^{3,5}, Rhonda D. Szczesniak^{3,5,6} and
Eleni-Rosalina Andrinopoulou^{1,2}

¹Department of Biostatistics, Erasmus University Medical Center, the Netherlands

²Department of Epidemiology, Erasmus University Medical Center, the Netherlands

³Division of Biostatistics and Epidemiology, Cincinnati Children’s Hospital Medical Center, USA

⁴Division of Statistics and Data Science, University of Cincinnati, USA

⁵Department of Pediatrics, University of Cincinnati, USA

⁶Division of Pulmonary Medicine, Cincinnati Children’s Hospital Medical Center, USA

Abstract

The joint modeling of longitudinal and time-to-event outcomes has become a popular tool in follow-up studies. However, fitting Bayesian joint models to large datasets, such as patient registries, can require extended computing times. To speed up sampling, we divided a patient registry dataset into subsamples, analyzed them in parallel, and combined the resulting Markov chain Monte Carlo draws into a consensus distribution. We used a simulation study to investigate how different consensus strategies perform with joint models. In particular, we compared grouping all draws together with using equal- and precision-weighted averages. We considered scenarios reflecting different sample sizes, numbers of data splits, and processor characteristics. Parallelization of the sampling process substantially decreased the time required to run the model. We found that the weighted-average consensus distributions for large sample sizes were nearly identical to the target posterior distribution. The proposed algorithm has been made available in an R package for joint models, `JMbayes2`. This work was motivated by the clinical interest in investigating the association between ppFEV₁, a commonly measured marker of lung function, and the risk of lung transplant or death, using data from the US Cystic Fibrosis Foundation Patient Registry (35,153 individuals with 372,366 years of cumulative follow-up). Splitting the registry into five subsamples resulted in an 85% decrease in computing time, from 9.22 to 1.39 hours. Splitting the data and finding a consensus distribution by precision-weighted averaging proved to be a computationally efficient and robust approach to handling large datasets under the joint modeling framework.

Keywords: big data, consensus Monte Carlo, distributed inference, joint model.

*Correspondence at: Department of Biostatistics, Erasmus University Medical Center, PO Box 2040, 3000 CA Rotterdam, the Netherlands. E-mail address: p.mirandaafonso@erasmusmc.nl.

1 Introduction

The joint modeling of longitudinal and time-to-event data has become a popular tool in follow-up studies to explore the association between such outcomes (Rizopoulos 2012). The availability of large-scale datasets, such as patient registries, is a valuable resource for research to enhance our understanding of disease outcomes and management. However, applying sophisticated statistical models, such as joint models, to large datasets may result in prohibitive computing times. The long running times are particularly cumbersome during the model-building phase, which often requires fitting and comparing multiple models. Different approaches have been described in the body of Bayesian literature to overcome the time-consuming nature of the posterior samplers, such as streamlining the estimation process via analytical approximations, or reducing the amount of data or outcomes analyzed simultaneously. In particular, as an alternative to Markov chain Monte Carlo (MCMC) sampling methods, the model posterior distribution can be approximated with a simpler distribution determined asymptotically. This simplifies the parameter estimation and could therefore allow the analysis of more complex model structures. Examples of such approaches are the Laplace integral approximation (Rue et al. 2009) and variational Bayes (Jaakkola and Jordan 2000) methods. Rustand et al. (2022) recently presented a joint model for multivariate longitudinal markers and competing risks based on the integrated nested Laplace approximation. Mauff et al. (2020) describe a corrected two-stage method for fitting a joint model with multiple longitudinal outcomes and a survival outcome. While the two-stage approach substantially decreases the time required to run the model, the time gains are limited by the computation necessary to fit a multivariate model in the first stage. For multivariate joint models with a prohibitive number of outcomes, in subsequent work, Mauff et al. (2021) presented a Bayesian adaptation of the pairwise approach introduced by Fieuws and Verbeke (Fieuws and Verbeke 2006). The strategy proved less effective under the joint model framework; nevertheless, the results are promising and warrant further investigation.

Our work is motivated by the analysis of the US Cystic Fibrosis Foundation Patient Registry (CFFPR) (Knapp et al. 2016). Cystic fibrosis (CF) is a severe genetic disease that affects the entire body, but its main symptoms primarily affect the lungs. There have been major improvements in CF outcomes over the past few decades, but it remains a life-limiting condition (Farrell et al. 2008). The CFFPR contains detailed health-related longitudinal data on US individuals living with CF. It is a large, comprehensive dataset containing annual and encounter-based data on demographics and CF outcomes such as lung function. A commonly measured marker of lung function in CF individuals above six years of age is the percentage of predicted forced expiratory volume in one second (ppFEV_1). Respiratory failure is the primary cause of death for people with CF, and some patients undergo lung transplantation when the disease progresses quickly. Therefore, there is much clinical interest in investigating the association between CF ppFEV_1 decline and the risk of lung transplant or death using joint models. Previous work on the CFFPR has used smaller samples of the registry data to overcome computational burdens (Andrinopoulou et al. 2020). Many applications with CF data have utilized single-center cohorts (Schluchter et al. 2002; Piccorelli and Schluchter 2012; Schluchter and Piccorelli 2019; Su et al. 2021). Other applications of CFFPR and UK CF data have considered all available data but relied on simpler joint models (Li et al. 2017; Taylor-Robinson et al. 2020; Andrinopoulou et al. 2020; Barrett et al. 2015).

In the present study, we explore reducing the computational time required to fit a joint model by tackling the amount of data analyzed simultaneously using consensus Monte Carlo methods. In particular, we split the dataset into independent subsamples, and distribute the posterior sampling between different MCMC samplers that each use one of these subsamples. The multiple sets of posterior samples are combined into a consensus distribution approximating the full posterior distribution. The algorithm is “embarrassingly parallel” (Herlihy and Shavit 2012) because the MCMC samplers can be executed concurrently without communicating. Thus, when paired with the multi-core processors available in

today’s computers, which allow the execution of multiple processes in parallel, it can substantially reduce the computing time required to run the model. The advantages of this approach are that it is simple, robust, and invariant to the dimension of the parameter space. The consensus distribution is obtained by averaging the sampled chains across the subposterior distributions. The weights reflect the information in each subsample, estimated from the within-sample Monte Carlo variance. This algorithm relies heavily on the assumption of normality. However, according to the Bernstein–von Mises theorem, the posteriors are approximately Gaussian for large sample sizes. The consensus through averaging has shown promising results when applied to simple regression models (Scott et al. 2016)—in which the normality assumption is met for most posteriors—but its performance under more complex settings, such as joint models, has not been evaluated to date. Motivated by the CF application, we examine the association between ppFEV₁ and the risk of lung transplant or death using all available CFFPR data by adapting different consensus strategies in the context of complex methods such as the joint model framework. To make the joint analysis of longitudinal and time-to-event outcomes using large datasets easily applicable by others, we implement the reviewed consensus strategies in an R package for joint models, `JMbayes2` (Rizopoulos et al. 2022), which is publicly available.

The remainder of this article is organized into four main sections. In Section 2, we present a theoretical introduction to joint models and Bayesian inference and modeling. We describe the consensus Monte Carlo algorithm and three methods to combine MCMC posterior samples, and we discuss their implementation in the R package `JMbayes2` (Rizopoulos et al. 2022). In Section 3, we explore through a simulation study how different sample sizes, numbers of data splits, and computer processor characteristics affect the performance of these consensus methods under the joint model framework. Section 4 presents the CFFPR case study that motivates this work, our modeling approach, and the results. In Section 5, we provide concluding remarks and suggest future research directions.

2 Statistical methods

2.1 Joint modeling framework

To model the longitudinal ppFEV₁ with time to transplantation or death, we rely on the joint modeling of longitudinal and survival data framework (Rizopoulos 2012). Let T_i^* denote the true failure time in the event-process for the i th individual, $i = 1, \dots, n$, and C_i the corresponding independent censoring time. The observed failure time is then T_i , with $T_i = \min\{T_i^*, C_i\}$. The event indicator δ_i is equal to 1 if $T_i < C_i$ and 0 otherwise. We denote the longitudinal marker measured at time t by $y_i(t)$. Joint models assume a full joint distribution of the longitudinal and time-to-event processes $[y_i(t), T_i]$.

Different factorizations of the joint distribution have been proposed in the literature (Sousa 2011). In this work, we focus on the shared-parameter joint models. We assume that the time-to-event and longitudinal processes depend on an unobserved process, defined by random effects \mathbf{b}_i . The observed processes are assumed independent conditional on the random effects, that is, $[y_i(t), T_i | \mathbf{b}_i] = [y_i(t) | \mathbf{b}_i][T_i | \mathbf{b}_i]$. This is mathematically convenient but computationally intensive.

The proposed model takes the form

$$\begin{cases} \mathcal{G}\{\mu_i(t)\} = \eta_i(t) = \mathbf{x}_i^\top(t)\boldsymbol{\beta} + \mathbf{z}_i^\top(t)\mathbf{b}_i & \text{longitudinal marker} \\ h_i(t) = h_0(t) \exp[\mathbf{w}_i^\top(t)\boldsymbol{\gamma} + \mathcal{F}\{\eta_i(t), \mathbf{b}_i\}\alpha] & \text{terminal event} \end{cases},$$

where $i = 1, \dots, n$ represent individuals, $\mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{-1})$, and $\boldsymbol{\Sigma}$ is the covariance matrix. To describe the individual-specific time evolution of the longitudinal marker, we specify a generalized linear mixed model (Pinheiro and Bates 2006), with $\mathbf{x}_i(t)$ and $\mathbf{z}_i(t)$ denoting the design vectors for the fixed effects $\boldsymbol{\beta}$ and the random effects \mathbf{b}_i , respectively. The design vectors may incorporate baseline or time-varying exogenous covariates. The fixed effects describe the average evolution of the marker over time. The random effects account for the

individual-specific evolution and the within-individual correlation over time. The expected value of $y_i(t)$ is $\mu_i(t)$, $\mathcal{G}(\cdot)$ is the link function, and $\eta_i(t)$ is the linear predictor.

For the terminal event process, we rely on a proportional hazard risk model to describe the time to the terminal event (Cox and Oakes 2018). The design vector $\mathbf{w}_i(t)$ is the parameter vector of covariates with the corresponding vector of regression coefficients $\boldsymbol{\gamma}$; it can include either baseline or exogenous time-varying covariates. The term $\mathcal{F}\{\eta_i(t), \mathbf{b}_i\}$ describes the functional form that links the longitudinal and terminal event processes. Different functional forms have been described in the literature, such as underlying value, slope, and cumulative effect (Rizopoulos 2012; Mauff et al. 2017; Andrinopoulou and Rizopoulos 2016). The magnitude of the association between the two processes is quantified by α .

2.2 Inference and consensus methods

We used Bayesian inference to estimate the joint model parameters from the CFFPR dataset. The posterior distribution for the joint model parameter θ is defined as

$$p(\theta \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}) = \frac{p(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta} \mid \theta)p(\theta)}{a} \propto p(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta} \mid \theta)p(\theta), \quad (1)$$

where $p(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta} \mid \theta)$ is the likelihood of the sample $(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta})$, $p(\theta)$ is the prior distribution, and a is a normalizing constant to make the posterior density integrate to one. The posterior estimation is based on Markov chain Monte Carlo (MCMC) sampling methods, such as the the Metropolis–Hastings algorithm. The sampling of high-dimensional individual-specific random effects in shared-parameter joint models makes the sampling process time consuming. More details about Bayesian joint model inference are given by Brown et al. (2005) and Ibrahim et al. (2001). Some methods for Bayesian inference and modeling, and MCMC methods for posterior sampling, have been presented by Gelman et al. (2013).

MCMC samplers have revolutionized statistical inference, but the analysis of large datasets with high-dimensional parameter spaces remains a major challenge. As datasets grow, the

computational cost and time required to run MCMC algorithms increase dramatically, making them impractical for many applications. Motivated by the computational burden of fitting joint models in the CFFPR dataset, we describe three computationally inexpensive consensus strategies under the joint modeling framework, in particular, combining all draws together and averaging the individual draws using two different weighting schemes.

We randomly partition the dataset $(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta})$ into S disjoint subsamples $(\mathbf{y}_1, \mathbf{t}_1, \boldsymbol{\delta}_1), \dots, (\mathbf{y}_S, \mathbf{t}_S, \boldsymbol{\delta}_S)$. We ensure they are independent by putting together all measurements from the same hierarchy level. For example, for a longitudinal dataset with only two hierarchy levels, measurements from the same individual are placed in a single subsample, one subsample can include multiple individuals, and the subsamples do not overlap. Assuming conditional independence across the subsamples given the parameter of interest θ , the posterior distribution in (1) can be factorized as

$$p(\theta \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}) \propto p(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta} \mid \theta) p(\theta) = \left(\prod_{s=1}^S p(\mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s \mid \theta) \right) p(\theta),$$

where $p(\mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s \mid \theta)$ is the likelihood of the subsample $(\mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s)$. When appropriate, the prior information $p(\theta)$ should be scaled to $p(\theta)^{1/S}$ to adjust the prior knowledge to the subsample at hand:

$$p(\theta \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}) \propto \prod_{s=1}^S p(\mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s \mid \theta) p(\theta)^{1/S} = \prod_{s=1}^S p(\theta \mid \mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s).$$

As a result, the full posterior distribution can be seen as the product of the S independent subposterior distributions. Each subsample goes through an independent MCMC simulation (in parallel), each of which generates D draws per chain from the subposterior distribution $p(\theta \mid \mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s)$,

$$\boldsymbol{\theta}_s^k = (\theta_{s,1}^k, \dots, \theta_{s,D}^k) \sim p(\theta \mid \mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s),$$

with $s = 1, \dots, S$, for the Monte Carlo Markov chain $k \in \{1, \dots, K\}$. Here, $\theta_{s,1}^k$ is the first

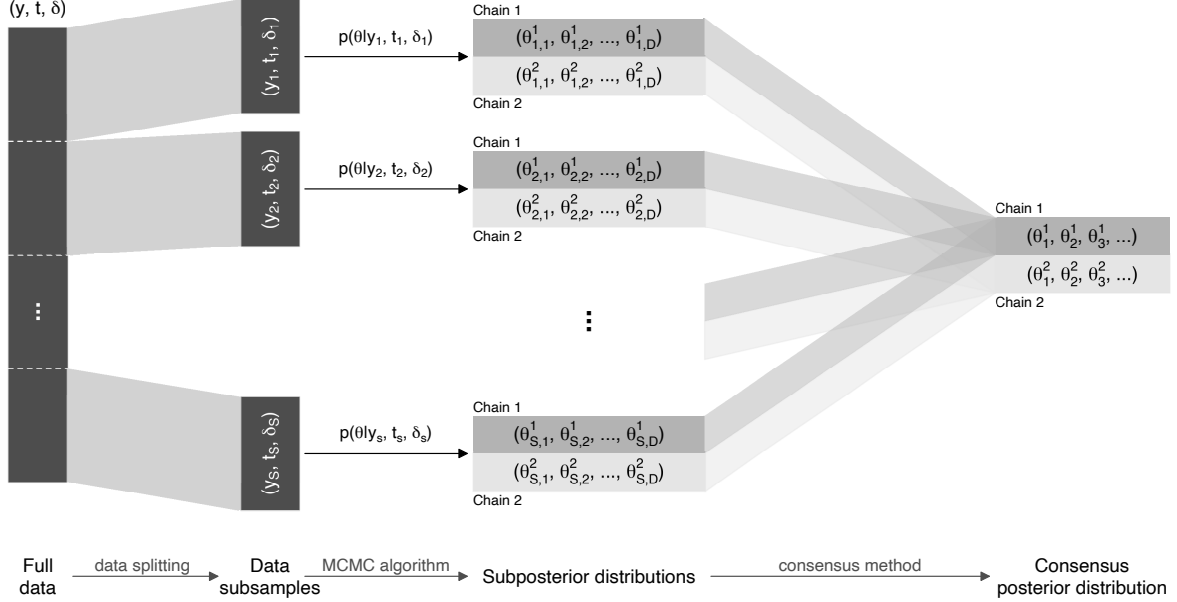


Figure 1: Illustration of a consensus Monte Carlo algorithm, considering S data splits with two Markov chains of D draws each.

MCMC draw in chain k obtained from the subsample s for the parameter θ . We obtain a consensus posterior distribution by combining the MCMC draws from all subsamples.

Figure 1 illustrates the consensus Monte Carlo process, considering two Markov chains. The draws from the k th chain of each subsample are combined into the k th posterior consensus chain. Before the consensus step, the algorithm is embarrassingly parallel, since the MCMC samplers run independently without communicating. This allows the sampling process to be carried out across independent machines to speed up the sampling process. The algorithm is transferable to any MCMC sampler, and the subposterior sampling is carried out in the same manner as sampling from the full-data posterior.

An intuitive approach to producing a consensus posterior distribution is creating a set containing all posterior draws for each chain,

$$\boldsymbol{\theta}_{\text{union}}^k = \left(\boldsymbol{\theta}_1^{k\top}, \dots, \boldsymbol{\theta}_S^{k\top} \right) \sim p(\theta \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}). \quad (2)$$

$(D \times S) \times 1$ $1 \times D$ $1 \times D$

The number of posterior consensus draws per chain is the product of the combined chains and the number of draws per chain. In the remainder of this article, we refer to this consensus approach as the union algorithm. Even though this approach is straightforward to implement, and the estimates are close to those obtained from the full dataset, it has the disadvantage that it increases the spread of the posterior distribution. To overcome this problem, we can average the individual posterior draws across the S subsamples and combine the D averaged draws. In this case, we have

$$\boldsymbol{\theta}_{D \times 1}^k = (\overline{\theta}_1^k, \dots, \overline{\theta}_D^k) \sim p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}), \quad (3)$$

where

$$\overline{\theta}_d^k = \sum_{s=1}^S w_s^k \theta_{s,d}^k, \quad (4)$$

for posterior draw $d \in \{1, \dots, D\}$. The weight applied to the MCMC draws in the k th chain from the subsample s is w_s^k . The final number of posterior draws equals the number of draws sampled from each subposterior distribution. Compared to the union algorithm, averaging the draws reduces the spread of the consensus posterior distribution around the mean. Table S1 illustrates how the individual posterior draws from the different subsamples are combined under the union and weighted-average consensus methods. For illustrative purposes, we consider D posterior draws for each of two Markov chains, and a dataset split into S subsamples.

One may argue that, for large subsamples, given that the dataset was split randomly, each subposterior sample has roughly the same information. Thus, we can give all subsamples the same relative importance by assigning them equal weights. In this case, each average draw is a simple average or arithmetic mean. Then w_s^k in (4) becomes

$$w_s^k = \frac{1}{S}, \quad (5)$$

for each subsample s . However, despite the random splitting, the shape of the subposterior densities might still vary between subsamples by chance. We can account for such differences by using information-based weighting to estimate the parameters unbiasedly. [Scott et al. \(2016\)](#) proposed a weighting scheme in which the weights reflect the precision in each subposterior sample. Recalling that precision is the inverse of variance, w_s^k in (4) becomes

$$w_s^k = \frac{w_s^{k'}}{a^k}, \quad (6)$$

where

$$w_s^{k'} = \text{Var}^{-1} \boldsymbol{\theta}_s^k = \text{Var}^{-1} (\theta_{s,1}^k, \dots, \theta_{s,D}^k), \quad (7)$$

and a^k is the normalizing constant for the k th chain by which all weights sum to one, $a^k = \sum_{s=1}^S w_s^{k'}$. Greater precision leads to a higher weight. Previous research on simpler regression models has demonstrated that the posterior achieved is almost identical to that obtained by fitting all data together ([Scott et al. 2016](#)). If the subposterior distributions are normal, that is, $\theta_{s,d}^k \sim p(\theta \mid \mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s) = \mathcal{N}(\mu_s, \tau_s)$, then by considering the weight $w_s^{k'}$ as the distribution precision τ_s , the averaged draws are exact random samples from the full-data posterior

$$\overline{\theta}_d^k \sim \mathcal{N} \left(\left(\sum_s \tau_s \right)^{-1} \sum_s \tau_s \theta_{s,d}^k, \sum_s \tau_s \right) = p(\theta \mid \mathbf{y}, \mathbf{t}, \boldsymbol{\delta}),$$

where

$$\overline{\theta}_d^k = \frac{1}{\sum_{s=1}^S \tau_s} \sum_{s=1}^S \tau_s \theta_{s,d}^k.$$

Recall that $p_1 p_2 \sim \mathcal{N}((\tau_1 + \tau_2)^{-1}(\tau_1 \mu_1 + \tau_2 \mu_2), \tau_1 + \tau_2)$ when $x_1 \sim p_1 = \mathcal{N}(\mu_1, \tau_1)$ and $x_2 \sim p_2 = (\mu_2, \tau_2)$. The precision τ_s is unknown, but the sample precision from the Monte Carlo draws, as defined in (7), can be used as it is the best estimate $\hat{\tau}_s = \text{Var}^{-1} \boldsymbol{\theta}_s^k$. Normality is a sufficient condition here but not a necessary one. The Bernstein–von Mises theorem—the Bayesian analogue to the central limit theorem—shows that for large sample sizes the

posterior distribution is approximately Gaussian (Van der Vaart 2000). The algorithm’s ability to capture characteristics of the joint model posterior distributions has yet to be evaluated.

2.3 Software and implementation

In recent years, joint models have seen considerable software development across primary statistical software tools, such as R, Stata, and SAS. These solutions cover frequentist and Bayesian models and offer a range of model customization options. Furgal et al. (2019) reviewed some of these software implementations. In this work, we implemented the consensus methods described in Subsection 2.2 in R package `JMbayes2`, available from the Comprehensive R Archive Network. `JMbayes2` is a user-friendly and versatile package that fits Bayesian joint models for longitudinal and time-to-event data (Rizopoulos et al. 2022).

Today, most computers have multi-core central processing units that enable more efficient processing of multiple tasks simultaneously. However, this does not reduce the computing time by a factor of the number of cores available because there is an overhead due to the indirect computing time used by the operating system to conduct the parallelization process. For fast computation processes, the cost of parallelization may outweigh the benefits of having more processing power, potentially making them take longer. Standard R is single-threaded, but there are packages such as `parallel` that allow a workload to be split across multiple cores. In our implementation, we concurrently run numerous joint models and various Markov chains to speed up sampling. The package runs parallel independent MCMC simulations for each subsample and returns a consensus distribution, according to the number of splits and the consensus method specified by the user. We run the joint models in parallel, and within each joint model we run its Markov chains in parallel. One core is required to run each Monte Carlo Markov chain. To maximize the algorithm’s efficiency, and thus minimize the computing time, the number of available cores should ideally match the product of the

number of splits and the number of Markov chains in each model. However, this is not strictly necessary to apply the techniques and benefit from a faster computation, as further results will show. In the supplementary Section D, we present an example of the use of consensus methods with `JMbayes2`.

3 Simulation study

We investigated how the union, equal-weighted, and precision-weighted consensus methods perform under the joint model framework through a simulation study. We explored a range of scenarios reflecting different sample sizes, numbers of data splits, and computer processor characteristics. In particular, we considered n individuals, $n \in \{500, 1000, 2500, 5000\}$, with repeated measurements, and for each dataset we used s data splits, $s \in \{1, 2, 5, 10\}$. The purpose of the single-split simulations ($s = 1$) was to replicate the gold standard approach, in which all data are used together. Furthermore, we investigated two processor architecture scenarios: i) unlimited availability of cores, and ii) 7 cores (to mimic a common 8-core architecture in which one core is set free for other processes outside R). We replicated each scenario 200 times. Table 1 outlines the simulation study.

3.1 Data generation

We generated data according to the following joint model:

$$\left\{ \begin{array}{l} y_i(t) = (\beta_0 + b_{0,i}) + (\beta_{\text{ns}_1} + b_{\text{ns}_1,i}) \text{ns}_1(t) + (\beta_2 + b_{2,i}) \text{ns}_2(t) + (\beta_3 + b_{3,i}) \text{ns}_3(t) \\ \quad + \varepsilon_i(t) \\ = \eta_i(t) + \varepsilon_i(t), \\ h_i(t) = h_0(t) \exp \{ \gamma_{\text{sex}} \text{sex}_{\text{male},i} + \gamma_{\text{age}} \text{age}_i + \eta_i(t) \alpha \}, \end{array} \right. \quad (8)$$

Table 1: Outline of the simulation study.

Step	Description
1:	Specify the number of individuals n , for $n \in \{500, 1,000, 2,500, 5,000\}$.
2:	Repeat 200 times:
3:	Simulate dataset from joint model (8). (Detailed description in Table S3.)
4:	Randomly split the $(\mathbf{y}, \mathbf{t}, \boldsymbol{\delta})$ into S subsamples $\{(\mathbf{y}_1, \mathbf{t}_1, \boldsymbol{\delta}_1), \dots, (\mathbf{y}_S, \mathbf{t}_S, \boldsymbol{\delta}_S)\}$, for $S \in \{1^\dagger, 2, 5, 10\}$.
5:	Specify the number of cores to use $\{7, \infty\}$.
6:	Run S separate MCMC algorithms to sample $\theta_{s,1}, \dots, \theta_{s,D} \sim p(\theta \mathbf{y}_s, \mathbf{t}_s, \boldsymbol{\delta}_s)$.
7:	Combine the resulting MCMC draws according to the consensus methods: union (2), equal- (3, 5), and precision-weighted average (3, 6).
<hr/> Datasets count: $9,600 = \underbrace{\left(\underbrace{4}_{\text{Step 1}} \times \underbrace{4}_{\text{Step 4}} \times \underbrace{3}_{\text{Step 7}} \right)}_{\text{No. data scenarios}} \times \underbrace{200}_{\text{Step 3}}_{\text{No. replicas}}$ <hr/>	

MCMC: Markov chain Monte Carlo.

\dagger Full data (gold standard).

$t > 0$, where $i = 1, \dots, n$ represent individuals, $(b_{0,i}, b_{\text{ns}_1,i}, b_{\text{ns}_2,i}, b_{\text{ns}_3,i})^\top \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{-1})$, $\varepsilon_i \sim \mathcal{N}(0, \sigma_y^2)$, and $(b_{0,i}, b_{\text{ns}_1,i}, b_{\text{ns}_2,i}, b_{\text{ns}_3,i}) \perp\!\!\!\perp \varepsilon_i$. The longitudinal marker is described by a linear mixed-effects model. It includes natural cubic splines with three degrees of freedom to accommodate the nonlinear progression of time for both fixed and random effects, represented by $\text{ns}_j(t)$, $j = 1, 2, 3$ in (8). The true level of the longitudinal outcome at time t , that is, the observed $y_i(t)$ without the measurement error $\varepsilon_i(t)$, is denoted by $\eta_i(t)$. We assume a Weibull baseline hazard $h_0(t)$ in the proportional hazard risk model, $h_0(t) = \phi\sigma_0 t^{\sigma_0-1}$. The relative risk of the terminal event is influenced by continuous and binary time-invariant variables: baseline age and sex, respectively. The parameter values used are in Table S2; the characteristics of the simulated datasets are in Table S4. We randomly split the individuals from each dataset according to the different strategies. By doing so, we ensure that all subsamples are independent of each other.

3.2 Model fitting

The joint models were run using `JMbayes2` (v0.4-5) on a 64-thread machine. We restricted the number of available cores to 7 and 63. For the data scenarios considered, the maximum number of cores required at the same time was 30 (3 chains \times 10 subsamples). Thus, 63 cores were enough to reproduce the unlimited number of cores scenario. We considered the true model structure, except for the risk model baseline hazard: we replaced the Weibull baseline risk with penalized B-spline functions (Eilers and Marx 1996) with 15 knots. For each model, we used the package’s default prior distributions and 3 Markov chains with 3500 iterations per chain, discarding the first 500 iterations as a warm-up. The convergence of the chains was assessed by the convergence diagnostic \hat{R} (Gelman and Rubin 1992) and by visual inspection of posterior traceplots of randomly chosen datasets within each scenario.

3.3 Results

There is a general agreement between the full-data (gold standard) estimates and the consensus algorithms’ estimates. The relative bias distribution for the association parameter α obtained across the 200 replications is shown in the left panel in Figure 2, which shows a box plot of relative bias against the consensus algorithm, the sample size, and the number of data splits. Section B in the supplementary material contains the same plot for the remaining parameters. In general, the consensus methods slightly increase the relative bias compared to the gold standard. A larger sample size reduces the estimation bias. We further notice that when bias is already present in the gold standard approach, minor differences emerge between the different consensus methods, with the precision-weighted consensus yielding a lower value (Figure S11, left panel, bottom row, and right column).

The three consensus methods accurately capture the location and spread of the posterior distribution of the model parameters. However, the union consensus leads to an

AMD Ryzen Threadripper PRO 3975WX 32-core 64-thread processor running at 3.49 GHz, using 256 GB of RAM, running Windows 11 Pro (v21H2).

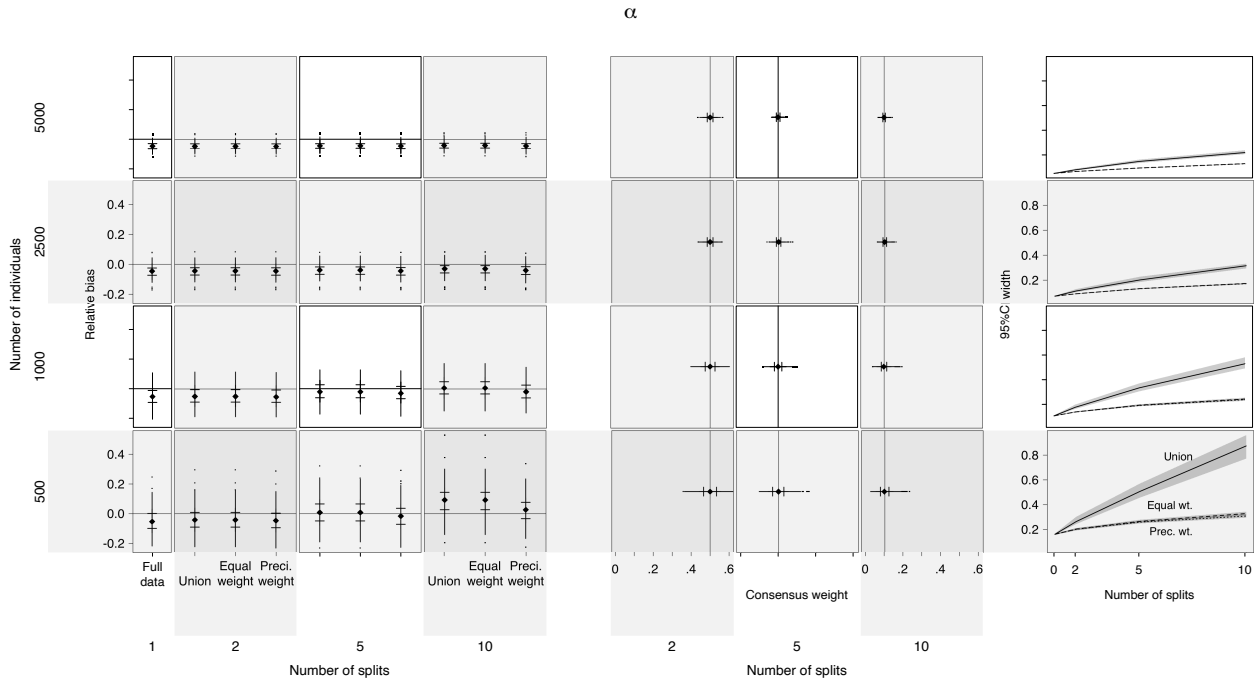


Figure 2: Left: Box plot of the relative bias for the α estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the α estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weights. Right: Median width of the 95% credible interval, with the associated interquartile range (IQR), of the α estimate—1st, 2nd, and 3rd quartiles—against the number of subsamples for different sample sizes.

over-dispersed distribution. The right panels in Figures 2 and S10–15 illustrate how the 95% credible interval (CI) width changes as the sample size increases from 50 to 2,500 individuals per subsample. The plot displays the 1st, 2nd, and 3rd quartiles of the 95% credible interval width against the sample size and the number of data splits. A higher number of data splits, and so less information per subsample, leads to wider credible intervals. However, this effect becomes negligible when the subsamples are sufficiently large. Nonetheless, caution should be taken when estimating effects that are only weakly supported by the data, which may be artifacts introduced by the consensus algorithm. The equal-weighted and precision-weighted algorithms produce similar estimates for the posterior mean and lower and upper percentiles. Given that we split the simulated data randomly, the information asymmetry across subsamples is low. This characteristic makes the precision weights similar across subsamples, approaching the values of the equal weights. The center panels in Figures 2 and S8–S13 present box plots of the standardized consensus precision weights for each model parameter across the 200 replicas against the sample size and the number of data splits. The vertical line represents the equal weight corresponding to the number of data splits. As the size of each subsample increases, the information asymmetry between the subposterior samples decreases and consequently the precision weights approach the equal weights.

In the left panel in Figure 3, we display the computing time required to fit the joint model (8) to datasets of different sizes for the scenario with an unlimited number of cores. As the concave upward curve reveals, the time does not increase linearly with the sample size but instead increases faster for larger sample sizes. Parallel computing comes with overhead costs, which reduce the efficiency of the consensus methods. According to Amdahl’s law, the speedup of a process is limited by the fraction of the process that cannot be parallelized, known as the serial fraction (Amdahl 1967). In the case of consensus methods, the overhead increases with the number of data splits since more processes need to be created by the operating system, leading to a higher serial fraction. Therefore, the performance of the techniques does not scale linearly as a function of the number of data splits. For example, if

the sample size is small, then the overhead may outweigh any advantages of the consensus methods and potentially make the computation take longer. Figure 3 (right panel) shows the computing time against the sample size, for the three numbers of data splits considered when all required cores were made available. The computing times are presented as proportional increases relative to the median time required to fit all data together. When the relative times are below one, the method is faster than fitting all the data together. The two-split strategy achieves faster computing times than the gold standard across all sample sizes considered. However, the five- and ten-split strategies only become more efficient for the larger sample sizes. Two data splits yield the lowest computing time when $n = 500$, while five data splits produce the lowest computing time when $n = 1,000$, $n = 2,500$, and $n = 5,000$. Although five splits result in the lowest median computing time for $n = 5,000$, the time is similar to that required for ten data splits. The ten-split approach requires higher sample sizes to outweigh its overhead and overcome the five-split strategy. The different splitting strategies allowed us to fit joint models in a timely fashion. For the sample size $n = 5,000$, all splitting strategies reduced the computing time by at least half. Splitting the data into two subsamples reduced the median running time by 56.17%. The five- and ten-split approaches show a median time reduction of about 0.22 and 0.26, respectively. The absolute computing times are available in Table S5. The time performance of a given splitting strategy is affected by the number of available cores. If the number of available cores is less than the number of processes to be parallelized, then some processes are placed into a waiting queue, increasing the computing time. The computing times obtained when restricting the number of available cores to seven are available in Figure S7 and Table S6. In our study, each subsample required three cores available to run its three Markov chains. Thus, the five- and ten-split strategies required more cores simultaneously than the seven made available—15 cores and 30 cores, respectively. The waiting time necessitated by this constraint downgraded the performance of both strategies.

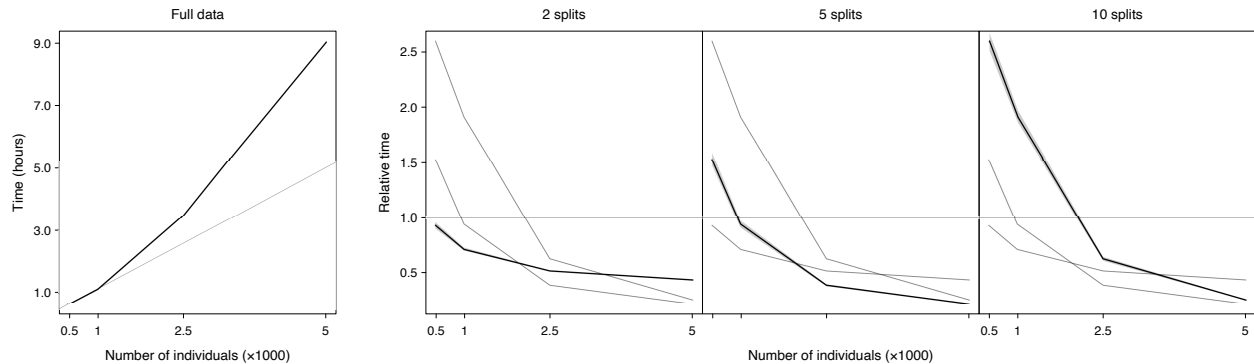


Figure 3: Left: Median computing time, with associated IQR, against the number of individuals, when using full data. The gray diagonal line shows a linear evolution. Right: Median and IQR computing time, relative to the time to fit all data together, against the number of individuals, in the scenario with an unlimited number of cores. The gray lines show the median time from the remaining panels.

4 Application

4.1 The CFFPR dataset

Our research goals require studying the association between ppFEV_1 and the risk of death or transplantation in CF individuals using the available CFFPR data. This dataset contains health-related information for 35,153 individuals aged over six years, who collectively contributed 1,523,406 observations and 372,366 years of follow-up. The demographic, social, and clinical characteristics of these individuals are displayed in Table S8. Females account for 48.34% of the group, and the median age at baseline is approximately 8.92 years (IQR 6.23–18.56). The median follow-up duration is 10.28 years (IQR 4.59–16.78). The dataset encompasses encounters between January 1, 1997, and December 31, 2017, with 50% of encounters between 2005 and 2014. This study focuses on a composite endpoint of death or lung transplantation. During the follow-up period, 23.47% of the individuals experienced one of the two. The median age to die from respiratory failure or to receive a double lung transplant is 27.12 years (IQR 21.36–36.00). The median age of the individuals at which their follow-up was censored is 21.33 years (IQR 14.12–30.94). Figure 4 (center panel) shows the Kaplan–Meier curve for the composite endpoint. Over the total follow-up, the median

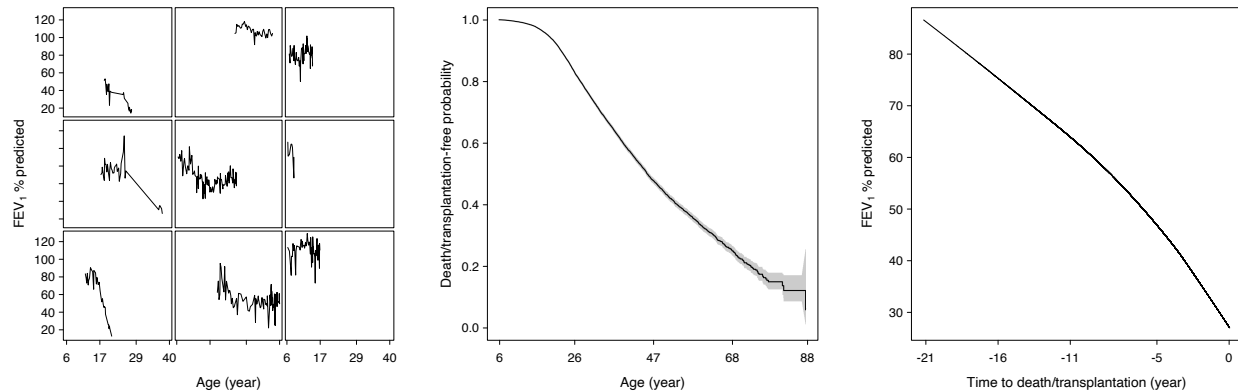


Figure 4: Left: ppFEV₁ measurements against age for nine randomly selected individuals. Center: Kaplan–Meier estimate of the death/transplantation-free probability, with associated 95% confidence interval. Right: ppFEV₁ measurements against time to transplantation or death, with loess smooth curve.

ppFEV₁ across all individuals is 73.60 (IQR 50.30–92.60). Figure 4 (left panel) displays the ppFEV₁ evolution experienced by nine randomly selected individuals over time. Figure 4 (right panel) suggests a negative association between ppFEV₁ and the risk of experiencing death/transplantation; it shows the mean ppFEV₁ evolution against the time to death or transplantation, for those who experienced one of the two events. It has been well established that lung function loss in individuals with CF is associated with a worse prognosis (Liou et al. 2001).

4.2 Analysis

We assumed a Bayesian joint model of longitudinal and survival data to study the association between ppFEV₁ and the risk of death/transplantation. To reduce the computing time required to evaluate the entire dataset, we split it into subsamples and analyzed them in parallel using the R statistical package `JMbayes2` (v0.4-5). We then used the union, equal-weighted, and precision-weighted strategies to obtain a consensus posterior. Motivated by our simulation results, we split the sample into five subsamples. By doing so, we expected to substantially reduce the computing time while keeping each subsample sufficiently large to

ensure they contained similar amounts of information. We fitted a joint model using all data together to assess the quality of the consensus estimates and the time savings. We term this model fit the "gold standard". This was only possible, however, using a non-conventional computer with sufficient memory. As shown in Table S8, the five subsamples are similar to each other and are each representative of the full sample. We fitted the following joint model:

$$\left\{ \begin{array}{l} \text{ppFEV}_{1i}(t) = (\beta_0 + b_{0,i}) + (\beta_{\text{ns}_1} + b_{\text{ns}_1,i}) \text{ns}_1(t) + (\beta_{\text{ns}_2} + b_{\text{ns}_2,i}) \text{ns}_2(t) + \beta_{\text{pa}} \text{Pa}_i(t) \\ \quad + \beta_{\text{ins}} \text{insurance}_{\text{yes},i}(t) + \beta_{\text{dpx}} \text{dep-idx}_i(t) + \beta_{[88,93]} \text{dob}_{[1988,1993],i} \\ \quad + \beta_{[93,98]} \text{dob}_{[1993,1998],i} + \beta_{[98,11]} \text{dob}_{[1998,2011],i} + \beta_{\text{sex}} \text{sex}_{\text{male},i} \\ \quad + \beta_{\text{htz}} \text{F508del}_{\text{htz},i} + \beta_{\text{oth}} \text{F508del}_{\text{other},i} + \varepsilon_i(t) \\ = \eta_i(t) + \varepsilon_i(t), \\ h_i(t) = h_0(t) \exp \{ \gamma_{\text{sex}} \text{sex}_{\text{male},i} + \eta_i(t) \alpha \}, \end{array} \right. \quad (9)$$

where $(b_{0,i}, b_{\text{ns}_1,i}, b_{\text{ns}_2,i})^\top \sim \mathcal{N}(\mathbf{0}, \Sigma^{-1})$, $\varepsilon_i \sim \mathcal{N}(0, \sigma_y^2)$, and $(b_{0,i}, b_{\text{ns}_1,i}, b_{\text{ns}_2,i}) \perp\!\!\!\perp \varepsilon_i$. We used penalized B-splines functions (Eilers and Marx 1996) $bs(t)$ to define the baseline hazard $h_0(t) = \exp\{\gamma_{h_0 0} + \sum_{q=1}^{15} \gamma_{h_0 q} bs_q(t)\}$. We assumed a non-linear evolution over time for the linear mixed-effects model, using natural cubic splines with two degrees of freedom, represented by $\text{ns}_{1,i}(t)$ and $\text{ns}_{2,i}(t)$ in (9). The initial time is at age six; for example, at $t = 2$ an individual is eight years old. We adjusted the average progression of ppFEV₁ for the following individual baseline characteristics: sex, birth cohort, and genotype. We furthermore assumed time-varying characteristics, such as the presence of infection by *Pseudomonas aeruginosa*, the possession of Medicaid insurance, and the deprivation index, as developed in previous work (Brokamp et al. 2019). For the random effects structure, we assumed a random intercept and the same non-linear effect of time considered for the fixed effects. Concerning the risk model, we assumed that the risk of death/transplantation is affected by

the individual’s sex and the underlying value of ppFEV_1 .

We fitted the models assuming 64 cores. We considered three Markov chains with 3,500 iterations, of which 500 were discarded for warm-up. The convergence of the chains was assessed by the convergence diagnostic \hat{R} (Gelman and Rubin 1992) and by visual inspection of the Markov chains’ traceplots.

4.3 Results

Splitting the full sample, with 35,153 individuals and 1,523,406 samples, into five subsamples produced an 84.89% decrease in computing time, from 9.22 to 1.39 hours. The results agree with the simulation study findings (Section 3.3). Figure 5 shows the estimated posterior means, with associated 95% credible interval, obtained from the full data (gold standard) and the different consensus methods. The union method consistently led to wider credible intervals. In particular, this method wrongly suggests that the deprivation index and possession of Medicaid insurance have a weak effect on ppFEV_1 , contradicting the findings from the gold standard method. The equal- and precision-weighted approaches performed equally well, yielding posterior estimates close to those obtained when using all data together. Each filled circle in Figure 6 (left panel) shows the estimated weights for the five different subsamples for that model parameter. The horizontal line represents equal weights of 0.2 each. The precision weights are close to the equal weights, because the random splitting process led to subsamples with similar information, as shown in Table S8.

The results suggest that ppFEV_1 is negatively associated with the risk of experiencing death or transplantation. The model obtained from the precision-weighted consensus method estimated an association of -0.119 (95% CI $-0.123, -0.113$), while the model leveraging the entire dataset yielded an estimate of -0.120 (95% CI $-0.123, -0.114$). A 10-unit decrease in the ppFEV_1 increases the hazard by approximately three times ($\text{HR} = 3.32$). Table 2 shows

AMD Ryzen Threadripper PRO 3975WX 32-core 64-thread processor running at 3.49 GHz, using 256 GB of RAM, running Windows 11 Pro (v21H2).

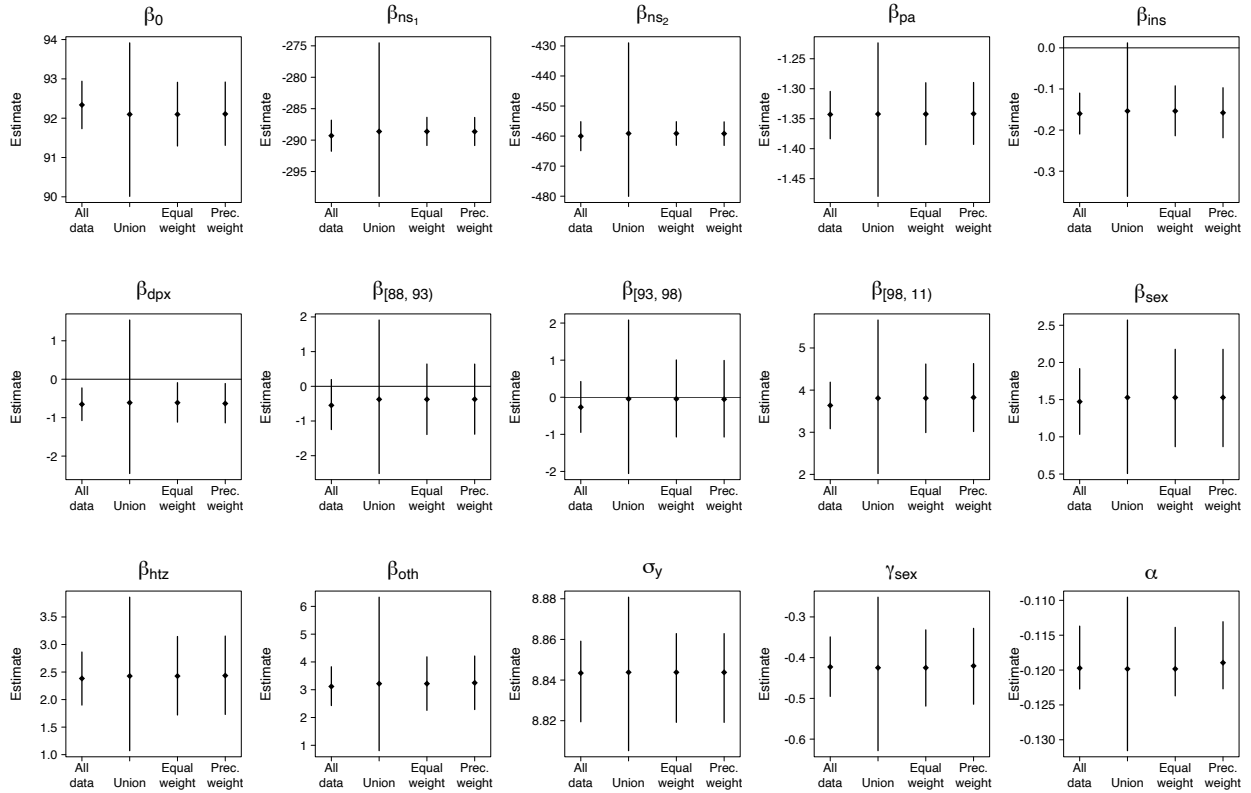


Figure 5: Estimated posterior means and 95% credible interval for joint model coefficients obtained from the full data (gold standard) and different consensus algorithms.

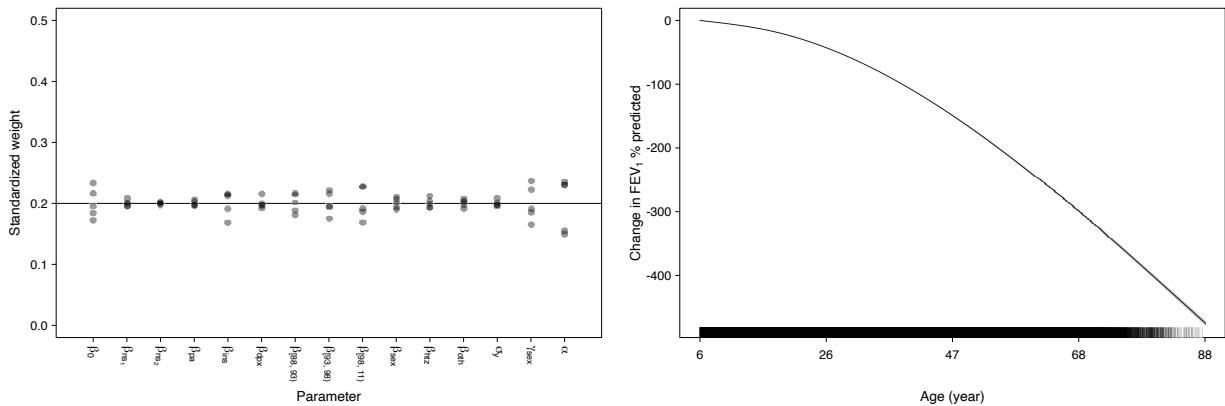


Figure 6: Left: Standardized weights used in the precision consensus method for each model parameter; the equal weights are represented by the horizontal line. Right: Main effect of time on the ppFEV₁ progression relative to its initial value at 6 years old.

the estimates for the remaining parameters obtained using the three consensus methods. For example, the detection of infection by *Pseudomonas aeruginosa* reduces the ppFEV₁ on average by approximately 1.34. Females present, on average, a lower ppFEV₁ of 1.5 and a 1.5 times higher hazard of death or transplantation (HR = 1.521). Figure 6 (right panel) displays a main effect plot for time: the average change in ppFEV₁ relative to the ppFEV₁ at baseline against age, considering the gold standard estimates. The concave downward curve indicates that the rate of ppFEV₁ decline increases as the patient gets older. That is, the deterioration of the individual’s lung function accelerates with aging.

5 Discussion

The collection of uniform observational data in patient registries constitutes an important epidemiological research tool that can be used to improve disease outcomes. However, processing such large amounts of data with complex statistical models can be computationally demanding. There is great clinical interest in determining the association between ppFEV₁, a commonly measured marker of lung function in CF patients, and their risk of death or lung transplantation, using all available CFFPR data. The CFFPR contains health-related information for around 35,000 individuals aged over six, who collectively contributed around 1,500,000 observations. Joint models provide a means to quantify the association between endogenous time-varying covariates and the relative risk of an event of interest (Rizopoulos 2012). These models are typically complex as they include multiple submodels with shared random effects. As a consequence, applying them to large datasets can be impractical or even impossible, due to the long running times or memory requirements. A range of approaches have been proposed in the literature to tackle this problem.

In this work, we studied the embarrassingly parallel consensus Monte Carlo algorithm applied to the joint modeling of longitudinal and survival data framework, with the goal of speeding up the posterior sampling. We randomly divided the dataset into non-overlapping

Table 2: Estimated posterior means and 95% credible interval for the joint model coefficients obtained from the full data (gold standard) and different consensus methods. We present the consensus estimates in terms of differences relative to full-data estimates; the original estimates can be found in Table S7.

Param.	Full data		Consensus methods						
	Mean	95% CI	Union		Equal weight		Precision weight		
			Δ Mean [†]	Δ 95% CI [‡]	Δ Mean [†]	Δ 95% CI [‡]	Δ Mean [†]	Δ 95% CI [‡]	
LME									
β_0	92.337	(91.734, 92.939)	0.240	(1.722, -0.976)	-0.240	(-0.441, -0.026)	-0.230	(-0.425, -0.020)	
β_{ns1}	-289.296	(-291.766, -286.828)	-0.667	(7.182, -12.297)	0.667	(0.913, 0.451)	0.652	(0.908, 0.432)	
β_{ns2}	-460.002	(-464.860, -455.208)	-0.873	(15.193, -26.215)	0.873	(1.798, -0.005)	0.821	(1.752, -0.061)	
β_{pa}	-1.343	(-1.383, -1.304)	-0.001	(0.096, -0.080)	0.001	(-0.010, 0.014)	0.001	(-0.010, 0.014)	
β_{ins}	-0.160	(-0.210, -0.110)	-0.006	(0.151, -0.122)	0.006	(-0.004, 0.018)	0.002	(-0.009, 0.013)	
β_{dpx}	-0.652	(-1.076, -0.227)	-0.041	(1.377, -1.767)	0.041	(-0.040, 0.137)	0.020	(-0.060, 0.115)	
$\beta_{[88,93]}$	-0.548	(-1.247, 0.196)	-0.172	(1.266, -1.714)	0.172	(-0.142, 0.446)	0.175	(-0.133, 0.444)	
$\beta_{[93,98]}$	-0.268	(-0.947, 0.424)	-0.223	(1.109, -1.657)	0.223	(-0.124, 0.581)	0.214	(-0.125, 0.567)	
$\beta_{[98,11]}$	3.636	(3.084, 4.188)	-0.173	(1.063, -1.475)	0.173	(-0.091, 0.431)	0.190	(-0.067, 0.441)	
β_{sex}	1.472	(1.034, 1.917)	-0.057	(0.527, -0.653)	0.057	(-0.166, 0.260)	0.057	(-0.166, 0.259)	
β_{htz}	2.383	(1.900, 2.863)	-0.044	(0.827, -0.996)	0.044	(-0.180, 0.282)	0.052	(-0.169, 0.290)	
β_{oth}	3.117	(2.431, 3.826)	-0.101	(1.624, -2.504)	0.101	(-0.173, 0.357)	0.131	(-0.144, 0.387)	
σ_y	8.843	(8.819, 8.859)	-0.001	(0.014, -0.022)	0.001	(0.000, 0.004)	0.001	(0.000, 0.004)	
PH									
γ_{sex}	-0.423	(-0.495, -0.349)	0.002	(0.133, -0.097)	-0.002	(-0.024, 0.017)	0.003	(-0.019, 0.021)	
α	-0.120	(-0.123, -0.114)	0.000	(0.009, -0.004)	0.000	(-0.001, 0.000)	0.001	(0.000, 0.001)	

CI: credible interval; LME: linear mixed-effects model; PH: proportional hazards model.

[†] Δ Mean = $\bar{\theta}_{\text{All data}} - \bar{\theta}_{\text{Consensus method}}$.

[‡] Δ 95% CI = $(P_{0.025} \{\theta_{\text{All data}}\} - P_{0.025} \{\theta_{\text{Consensus method}}\}, P_{0.975} \{\theta_{\text{All data}}\} - P_{0.975} \{\theta_{\text{Consensus method}}\})$.

and independent subsamples, used a multi-core processor to analyze them in parallel, and combined the resulting MCMC draws into a consensus distribution. During the sampling process, the MCMC samplers run independently on each subsample without communicating. This characteristic allows the sampling process to be carried out independently in multiple processor cores. This approach can be applied to any MCMC method.

We explored three consensus strategies—union, equal-weighted, and precision-weighted—that have previously been applied to simpler Bayesian regression models with satisfactory results. These techniques differ in how the posterior draws obtained from the data splits are combined into a single consensus posterior distribution. The union algorithm combines all the MCMC samples. The equal-weighted and precision-weighted algorithms compute a weighted average for each individual posterior draw across the subposterior samples. Through a simulation study, we investigated how these consensus methods perform under the joint model framework for a range of scenarios with different sample sizes, numbers of data splits, and computer processor characteristics. We applied the reviewed methods to the CFFPR case study that motivated this work. We split the CFFPR dataset into five subsamples and compared the computing time and the quality of the estimates against a model that leveraged the entire dataset (the gold standard model). To assist with future longitudinal and time-to-event analyses of large datasets, we implemented the consensus methods in an R package for joint models, `JMbayes2` (Rizopoulos et al. 2022).

Our simulation results reveal that the computing time required to fit the joint model increases rapidly with the sample size. Furthermore, they show that the consensus Monte Carlo methods can substantially reduce the computing time by parallelizing the sampling process. However, a larger number of data splits does not necessarily translate into reduced computing time. Parallel computing comes with overhead costs, which reduce the efficiency of the consensus methods. For relatively small sample sizes, the overhead associated with a larger number of data splits may outweigh any advantages of the techniques and thus increase the computing time. For example, dividing the data into two splits consistently led to a

reduction in computing time for samples of 500, 1,000, 2,500, and 5,000 individuals. However, using five or ten splits only achieved meaningful time savings for sample sizes over 1,000 individuals. These results are sensitive to the complexity of the joint model at hand, such as the number of independent variables and outcomes. More complex models require longer running times; thus, the overhead costs become less significant as the model complexity grows. The efficiency of the methods is also dependent on the processor characteristics. If the number of processes to be parallelized—the number of subsamples multiplied by the number of Markov chains—exceeds the number of available cores, then some of the processes must be placed into a queue, limiting the potential time gains. The number of data splits should therefore be chosen to match the available hardware resources. In our application, splitting the CFFPR into five subsamples yielded an 84.9% decrease in computing time, from 9.22 to 1.39 hours.

We found an agreement between the posterior mean estimates from the three consensus algorithms and the gold standard model in both the simulation and application studies. The differences lie in the widths of the estimated credibility intervals. The union consensus algorithm is simple to use but generates over-dispersed consensus distributions. This can be problematic in the presence of marginally significant effects. For example, in our application, the union method widened the parameters' credible intervals, leading to misinterpretation of the relevance of the effect of two of the independent variables in the model. Averaging the draws reduces the spread of the consensus posterior distribution around the mean. The equal-weighted and precision-weighted techniques fared similarly well, yielding posterior samples that were close to those obtained when using the full dataset together. We believe this is because the information asymmetry across the data splits was negligible due to the random splitting of the dataset and the sizes of the subsamples. For large, random subsamples, we expect that the precision-weighted and equal-weighted algorithms should produce similar results. However, in practice, it is difficult to evaluate in advance whether all subsamples are equally informative. Consequently, precision weighting should typically be preferred.

Our application findings suggest that ppFEV_1 is negatively associated with the risk of dying by respiratory failure or requiring a double lung transplant. That is, the lower the ppFEV_1 value, the higher the risk. While the gold standard approach produced an estimate for this association of -0.120 (95% CI $-0.123, -0.114$), the precision-weighted consensus algorithm estimated an association of -0.119 (95% CI $-0.123, -0.113$). The risk is increased by around three times ($\text{HR} = 3.32$) for each 10-unit drop in ppFEV_1 . In short, our findings show that our approach accelerates the sampling process while producing correct samples from the target distribution.

This work considered a simple joint model with one longitudinal and one terminal outcome. However, it would be straightforward to extend the reviewed consensus methods to more complex settings, such as joint models accommodating multiple longitudinal outcomes, competing risks, or intermediate, multi-state, or recurrent events. Our goal was to reduce the computing time required to fit a joint model to a large dataset using a single computer; thus, the algorithms presented were described in the context of parallel computation in a single multi-core machine. However, they could be adapted to run across multiple machines in multiple centers to overcome data confidentiality concerns—it would no longer be necessary to share data between centers, but only the model MCMC draws—or sequentially in the same machine to alleviate memory bottlenecks. In this work, we ignored the presence of different demographic groups within the population, that is, we assumed a homogeneous population. When there are large differences between groups, a stratified partition of the dataset might be required. Each subsample must contain enough information to produce unbiased estimates of the parameters: small-sample biases might be introduced if the data are split into many small subsamples. When this is a concern, the bias could be mitigated using, for example, jackknife bias correction ([Scott et al. 2016](#)). Further investigation could focus on the impact of the random splitting strategy on the consensus estimates. More specifically, the degree to which the estimated value varies with regard to different groups of random subsamples of the same dataset, to assess the reliability of the various consensus methods.

While we expect larger subsample sizes to equalize the variability between the weighted techniques, we anticipate the precision-weighted algorithm to yield a lower variability than the equal-weighted algorithm, since it is robust to asymmetrical subposteriors.

The question of which is the best consensus algorithm to apply to consensus Monte Carlo remains open. In this work, we explored three consensus algorithms. The union algorithm has no theoretical grounding and has major flaws in practice. The precision-weighted algorithm is rooted in Gaussian theory, but given that it does not rely on an explicit Gaussian approximation, it can also capture non-Gaussian features (e.g., skewness, fat tails) of the posterior distribution. Its main weakness is that it is limited to continuous parameter spaces. The equal-weighted consensus can be seen as a relaxation of the precision-weighted algorithm in which, to simplify computation, one assumes that for large sample sizes, the subposterior distributions encapsulate roughly the same information, so their samples can be combined with equal weights. Strategies other than averaging can also be applied to obtain consensus samples. Alternative algorithms have been described in the literature, such as importance resampling (Huang and Gelman 2005), and modeling the subposterior distributions parametrically (Huang and Gelman 2005), with kernels (Neiswanger et al. 2013; Wang and Dunson 2013; Wang et al. 2015; Srivastava et al. 2018), or with Gaussian process approximations (Nemeth and Sherlock 2018). Some of these techniques show promise, but they also add non-trivial complexity. Averaging has powerful advantages: it is simple, stable, and computationally inexpensive. More complex algorithms were not addressed in this work because the simpler techniques produced satisfactory results and the more complex algorithms would increase the computing time. Nonetheless, this should not discourage researchers from exploring other consensus methods under the joint model framework in future work.

Joint models provide a means to estimate individual-specific CF ppFEV₁ decline and its association with mortality or transplantation. The consensus Monte Carlo algorithm is an efficient solution for handling large datasets like the CFFPR without compromising the

amount of information taken into account or sacrificing model adequacy, thereby enhancing our understanding of CF ppFEV₁ decline. It could bring new insights into the progression of the disease and could contribute to better monitoring and treatment strategies. The availability of an easy-to-use statistical tool such as `JMbayes2` is likely to help applied researchers, in the era of big data, perform longitudinal and time-to-event data analyses in their everyday practice.

Acknowledgments

The authors thank the Cystic Fibrosis Foundation for providing the patient registry (CFFPR) data used to conduct this study, and the patients, care providers, and clinic coordinators at US CF centers for their contributions to the CFFPR.

Funding

This work was supported by grants from the National Institutes of Health (R01 HL141286) and the Cystic Fibrosis Foundation (SZCZES18AB0).

Data Availability Statement

The data that support the findings of this study are available from the Cystic Fibrosis Foundation. Restrictions apply to the availability of these data, which were used under license for this study. Requests for data may be sent to datarequests@cff.org.

References

- Amdahl, G. M. (1967). Validity of the single processor approach to achieving large scale computing capabilities. In *Proceedings of the April 18-20, 1967, spring joint computer conference*, pages 483–485.
- Andrinopoulou, E.-R., Clancy, J. P., and Szczesniak, R. (2020). Multivariate joint modeling to identify markers of growth and lung function decline that predict cystic fibrosis pulmonary exacerbation onset. *BMC Pulmonary Medicine* **20**, 1–11.
- Andrinopoulou, E.-R., Nasserinejad, K., Szczesniak, R., and Rizopoulos, D. (2020). Integrating latent classes in the bayesian shared parameter joint model of longitudinal and survival outcomes. *Statistical Methods in Medical Research* **29**, 3294–3307.
- Andrinopoulou, E.-R. and Rizopoulos, D. (2016). Bayesian shrinkage approach for a joint model of longitudinal and survival outcomes assuming different association structures. *Statistics in Medicine* **35**, 4813–4823.
- Barrett, J., Diggle, P., Henderson, R., and Taylor-Robinson, D. (2015). Joint modelling of repeated measurements and time-to-event outcomes: flexible model specification and exact likelihood inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **77**, 131–148.
- Bender, R., Augustin, T., and Blettner, M. (2005). Generating survival times to simulate cox proportional hazards models. *Statistics in Medicine* **24**, 1713–1723.
- Brokamp, C., Beck, A. F., Goyal, N. K., Ryan, P., Greenberg, J. M., and Hall, E. S. (2019). Material community deprivation and hospital utilization during the first year of life: an urban population-based cohort study. *Annals of Epidemiology* **30**, 37–43.
- Brown, E. R., Ibrahim, J. G., and DeGruttola, V. (2005). A flexible b-spline model for multiple longitudinal biomarkers and survival. *Biometrics* **61**, 64–73.

- Cox, D. R. and Oakes, D. (2018). *Analysis of survival data*. Chapman and Hall/CRC.
- Eilers, P. H. and Marx, B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science* **11**, 89–121.
- Farrell, P. M., Rosenstein, B. J., White, T. B., Accurso, F. J., Castellani, C., Cutting, G. R., Durie, P. R., LeGrys, V. A., Massie, J., Parad, R. B., et al. (2008). Guidelines for diagnosis of cystic fibrosis in newborns through older adults: Cystic fibrosis foundation consensus report. *The Journal of pediatrics* **153**, S4–S14.
- Fieuws, S. and Verbeke, G. (2006). Pairwise fitting of mixed models for the joint modeling of multivariate longitudinal profiles. *Biometrics* **62**, 424–431.
- Furgal, A. K., Sen, A., and Taylor, J. M. (2019). Review and comparison of computational approaches for joint longitudinal and time-to-event models. *International Statistical Review* **87**, 393–418.
- Gelman, A., Carlin, J., Stern, H., Dunson, D., Vehtari, A., and Rubin, D. (2013). *Bayesian Data Analysis, Third Edition*. Chapman & Hall/CRC Texts in Statistical Science. Taylor & Francis.
- Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical Science* pages 457–472.
- Herlihy, M. and Shavit, N. (2012). The art of multiprocessor programming, revised reprint. *ISBN-13* pages 978–0123973375.
- Huang, Z. and Gelman, A. (2005). Sampling for bayesian computation with large datasets. *Available at SSRN 1010107*.
- Ibrahim, J. G., Chen, M.-H., Sinha, D., Ibrahim, J., and Chen, M. (2001). *Bayesian survival analysis*, volume 2. Springer.

- Jaakkola, T. S. and Jordan, M. I. (2000). Bayesian parameter estimation via variational methods. *Statistics and Computing* **10**, 25–37.
- Knapp, E. A., Fink, A. K., Goss, C. H., Sewall, A., Ostrenga, J., Dowd, C., Elbert, A., Petren, K. M., and Marshall, B. C. (2016). The cystic fibrosis foundation patient registry. design and methods of a national observational disease registry. *Annals of the American Thoracic Society* **13**, 1173–1179.
- Li, D., Keogh, R., Clancy, J. P., and Szczesniak, R. D. (2017). Flexible semiparametric joint modeling: an application to estimate individual lung function decline and risk of pulmonary exacerbations in cystic fibrosis. *Emerging Themes in Epidemiology* **14**, 1–13.
- Liou, T. G., Adler, F. R., FitzSimmons, S. C., Cahill, B. C., Hibbs, J. R., and Marshall, B. C. (2001). Predictive 5-year survivorship model of cystic fibrosis. *American Journal of Epidemiology* **153**, 345–352.
- Mauff, K., Erler, N. S., Kardys, I., and Rizopoulos, D. (2021). Pairwise estimation of multivariate longitudinal outcomes in a bayesian setting with extensions to the joint model. *Statistical Modelling* **21**, 115–136.
- Mauff, K., Steyerberg, E., Kardys, I., Boersma, E., and Rizopoulos, D. (2020). Joint models with multiple longitudinal outcomes and a time-to-event outcome: a corrected two-stage approach. *Statistics and Computing* **30**, 999–1014.
- Mauff, K., Steyerberg, E. W., Nijpels, G., van der Heijden, A. A., and Rizopoulos, D. (2017). Extension of the association structure in joint models to include weighted cumulative effects. *Statistics in Medicine* **36**, 3746–3759.
- Neiswanger, W., Wang, C., and Xing, E. (2013). Asymptotically exact, embarrassingly parallel mcmc. *arXiv preprint arXiv:1311.4780* .

- Nemeth, C. and Sherlock, C. (2018). Merging mcmc subposteriors through gaussian-process approximations.
- Piccorelli, A. V. and Schluchter, M. D. (2012). Jointly modeling the relationship between longitudinal and survival data subject to left truncation with applications to cystic fibrosis. *Statistics in Medicine* **31**, 3931–3945.
- Pinheiro, J. and Bates, D. (2006). *Mixed-effects models in S and S-PLUS*. Springer Science & Business Media.
- Rizopoulos, D. (2012). *Joint models for longitudinal and time-to-event data: With applications in R*. CRC press.
- Rizopoulos, D., Papageorgiou, G., and Miranda Afonso, P. (2022). *JM-bayes2: Extended Joint Models for Longitudinal and Time-to-Event Data*. <https://drizopoulos.github.io/JMbayes2/>, <https://github.com/drizopoulos/JMbayes2>.
- Rue, H., Martino, S., and Chopin, N. (2009). Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **71**, 319–392.
- Rustand, D., van Niekerk, J., Krainski, E. T., Rue, H., and Proust-Lima, C. (2022). Fast and flexible inference approach for joint models of multivariate longitudinal and survival data using integrated nested laplace approximations. *arXiv preprint arXiv:2203.06256* .
- Schluchter, M. D., Konstan, M. W., and Davis, P. B. (2002). Jointly modelling the relationship between survival and pulmonary function in cystic fibrosis patients. *Statistics in Medicine* **21**, 1271–1287.
- Schluchter, M. D. and Piccorelli, A. V. (2019). Shared parameter models for joint analysis of longitudinal and survival data with left truncation due to delayed entry—applications to cystic fibrosis. *Statistical Methods in Medical Research* **28**, 1489–1507.

- Scott, S. L., Blocker, A. W., Bonassi, F. V., Chipman, H. A., George, E. I., and McCulloch, R. E. (2016). Bayes and big data: The consensus monte carlo algorithm. *International Journal of Management Science and Engineering Management* **11**, 78–88.
- Sousa, I. (2011). A review on joint modelling of longitudinal measurements and time-to-event. *Revstat Stat J* **9**, 57–81.
- Srivastava, S., Li, C., and Dunson, D. B. (2018). Scalable bayes via barycenter in wasserstein space. *The Journal of Machine Learning Research* **19**, 312–346.
- Su, W., Wang, X., and Szczesniak, R. D. (2021). Flexible link functions in a joint hierarchical gaussian process model. *Biometrics* **77**, 754–764.
- Taylor-Robinson, D., Schlüter, D. K., Diggle, P. J., and Barrett, J. K. (2020). Explaining the sex effect on survival in cystic fibrosis: a joint modeling study of uk registry data. *Epidemiology (Cambridge, Mass.)* **31**, 872.
- Van der Vaart, A. W. (2000). *Asymptotic statistics*, volume 3. Cambridge university press.
- Wang, X. and Dunson, D. B. (2013). Parallel mcmc via weierstrass sampler. *arXiv preprint arXiv:1312.4605* **24**,
- Wang, X., Guo, F., Heller, K. A., and Dunson, D. B. (2015). Parallelizing mcmc with random partition trees. *Advances in Neural Information Processing Systems* **28**,

Supplementary material

A Consensus methods

Table S1: Illustration of the union and weighted-average consensus methods, considering S data subsamples with two Markov chains of D draws each.

	Subsamples							Consensus methods			
	$s = 1$		$s = s$		\cdots	$s = S$		Union		Weighted average	
	Chain 1	Chain 2	Chain 1	Chain 2		Chain 1	Chain 2	Chain 1	Chain 2	Chain 1	Chain 2
Draw 1 $d = 1$	$\theta_{1,1}^1$	$\theta_{1,1}^2$	$\theta_{2,1}^1$	$\theta_{2,1}^2$	\cdots	$\theta_{S,1}^1$	$\theta_{S,1}^2$	$(\theta_{1,1}^1, \dots, \theta_{S,1}^1)^\top$	$(\theta_{1,1}^2, \dots, \theta_{S,1}^2)^\top$	$\bar{\theta}_1^1 = \sum_{s=1}^S w_s^1 \theta_{s,1}^1$	$\bar{\theta}_1^2 = \sum_{s=1}^S w_s^2 \theta_{s,1}^2$
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
Draw 3 $d = 3$	$\theta_{1,3}^1$	$\theta_{1,3}^2$	$\theta_{2,3}^1$	$\theta_{2,3}^2$	\cdots	$\theta_{S,3}^1$	$\theta_{S,3}^2$	$(\theta_{1,3}^1, \dots, \theta_{S,3}^1)^\top$	$(\theta_{1,3}^2, \dots, \theta_{S,3}^2)^\top$	$\bar{\theta}_3^1 = \sum_{s=1}^S w_s^1 \theta_{s,3}^1$	$\bar{\theta}_3^2 = \sum_{s=1}^S w_s^2 \theta_{s,3}^2$
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
Draw D $d = D$	$\theta_{1,D}^1$ $= \underbrace{\theta_{1,D}^1}_{D \times 1}$	$\theta_{1,D}^2$	$\theta_{2,D}^1$	$\theta_{2,D}^2$	\cdots	$\theta_{S,D}^1$	$\theta_{S,D}^2$	$(\theta_{1,D}^1, \dots, \theta_{S,D}^1)^\top$ $= \underbrace{\theta_{(D \times S) \times 1}^1}_{(D \times S) \times 1}$	$(\theta_{1,D}^2, \dots, \theta_{S,D}^2)^\top$	$\bar{\theta}_D^1 = \underbrace{\sum_{s=1}^S w_s^1 \theta_{s,D}^1}_{= \theta_{D \times 1}^1 \text{ w-avg}}$	$\bar{\theta}_D^2 = \sum_{s=1}^S w_s^2 \theta_{s,D}^2$

B Simulation study

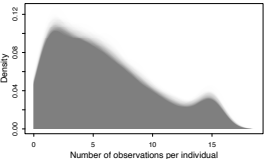
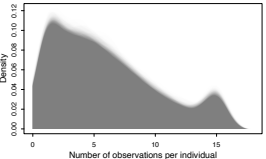
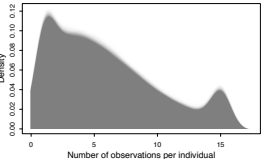
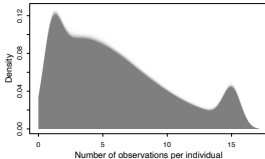
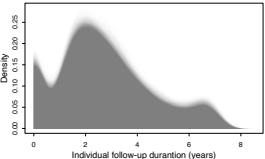
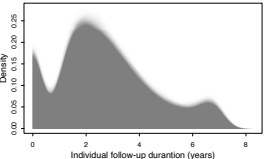
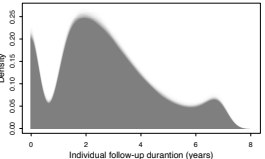
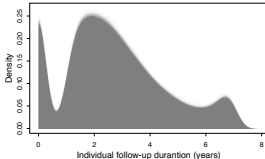
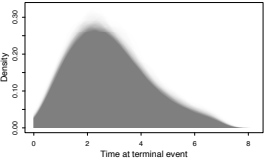
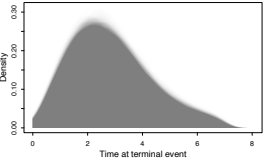
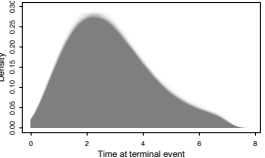
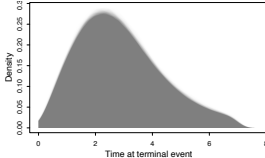
Table S2: Parameter values used in the joint model for data generation in the simulation study.

Model	Parameter	Value
Linear mixed-effects	$(\beta_0, \beta_{\text{ns}_1}, \beta_{\text{ns}_2}, \beta_{\text{ns}_3})$	(6.94, 1.30, 1.84, 1.82)
	σ_y^2	0.6^2
	Σ	$\begin{bmatrix} 0.71 & 0.33 & 0.07 & 1.26 \\ & 2.68 & 3.81 & 4.35 \\ & & 7.62 & 5.40 \\ & & & 8.00 \end{bmatrix}$
Proportional hazards	ϕ_0	$\exp\{-9\}$
	σ_0	2
	$(\gamma_{\text{sex}}, \gamma_{\text{age}})$	(0.5, 0.05)
	α	0.5

Table S3: Outline of the joint model data simulation process.

Longitudinal outcome (1/2):	
1:	Generate n random samples from $\mathcal{N}(\mathbf{0}, \Sigma^{-1})$ for the individual-specific random effects, \mathbf{b}_i : \mathbf{b} . $n \times 3$
2:	Generate $(n \times (n_i - 1))$ random samples from $\mathcal{U}(0, 7)$ for the visiting times t_{ij} , and add the time 0 for all individuals: \mathbf{t} . $(n \cdot n_i) \times 1$
3:	Generate the B-spline basis vectors for a natural cubic spline with 3 degrees of freedom for the visiting times \mathbf{t} : $[\mathbf{t}_{\text{ns}_1} \quad \mathbf{t}_{\text{ns}_2} \quad \mathbf{t}_{\text{ns}_3}]$. $(n \times n_i) \times 4$
4:	Generate the n vectors of n_i individual underlying longitudinal responses with $\mathbf{\eta}_i = [\mathbf{1} \quad \mathbf{t}_{\text{ns}_1, i} \quad \mathbf{t}_{\text{ns}_2, i} \quad \mathbf{t}_{\text{ns}_3, i}] \boldsymbol{\beta} + [\mathbf{1} \quad \mathbf{t}_{\text{ns}_1, i} \quad \mathbf{t}_{\text{ns}_2, i} \quad \mathbf{t}_{\text{ns}_3, i}] \mathbf{b}_i$. $n_i \times 1$ $n_i \times 4$ 4×1 $n_i \times 4$ 4×1
5:	Generate $(n \times n_i)$ random samples from $\mathcal{N}(0, \sigma_y^2)$ for the observation measurement error, ε_i : ε . $(n \cdot n_i) \times 1$
6:	Obtain the observed longitudinal response by summing the vectors $\boldsymbol{\eta} + \boldsymbol{\varepsilon}$: \mathbf{y} . $(n \cdot n_i) \times 1$
Survival outcome:	
7:	Generate n random samples from $\mathcal{U}(30, 70)$ for the individual's baseline age, age_i : \mathbf{age} . $n \times 1$
8:	Generate n random samples from Bern(0.5) for the individual's sex, sex_i : \mathbf{sex} . $n \times 1$
9:	Generate n random samples from $\mathcal{U}(0, 1)$, u_i : \mathbf{u} . $n \times 1$
10:	Define $H_i(t) = \int_0^t h_i(s) ds$, where $h_i(t) = \underbrace{\phi \sigma_0 t^{\sigma_0 - 1}}_{h_0(t)} \exp\{\gamma_{\text{sex}} \text{sex}_{\text{male}, i} + \gamma_{\text{age}} \text{age}_i + \eta_i(t) \alpha\}$
11:	Solve numerically $\exp(-H_i(t_i^*)) = u_i$ for t_i^* (Bender et al. 2005) to obtain the individual true event times: \mathbf{t}^* . $n \times 1$
12:	The observed event times $t_i = \min(t_i^*, t_{\text{max}})$, where t_{max} is the deterministic maximum follow-up time: \mathbf{t} . $n \times 1$
13:	Define the type I censoring indicator as $\delta_i = \begin{cases} 1 & t_i \leq t_{\text{max}} \\ 0 & t_i > t_{\text{max}} \end{cases}$.
Longitudinal outcome (2/2):	
14:	Remove all y_{ij} for $t_{ij} > t_i$.

Table S4: Characteristics of the simulated datasets.

Number of individuals (n)	500	1,000	2,500	5,000
Number of replicas	200	200	200	200
Total number of observations				
Median (IQR)	2961.50 (2897.50, 3032.25)	5927.50 (5849.25, 6026.00)	14811.50 (14664.75, 14951.50)	29658.50 (29462.75, 29870.50)
Number of observations/ind.				
2.5 th PCT, Median (IQR)	1.00 (1.00, 1.00)	1.00 (1.00, 1.00)	1.00 (1.00, 1.00)	1.00 (1.00, 1.00)
Median, Median (IQR)	5.00 (5.00, 5.00)	5.00 (5.00, 5.00)	5.00 (5.00, 5.00)	5.00 (5.00, 5.00)
97.5 th PCT, Median (IQR)	15.00 (15.00, 15.00)	15.00 (15.00, 15.00)	15.00 (15.00, 15.00)	15.00 (15.00, 15.00)
				
Aggregated follow-up time				
Median (IQR)	1308.38 (1274.51, 1332.12)	2612.78 (2581.94, 2655.62)	6524.46 (6460.77, 6585.95)	13058.13 (12965.20, 13176.31)
Total follow-up time/ind.				
2.5 th PCT, Median (IQR)	0.00 (0.00, 0.00)	0.00 (0.00, 0.00)	0.00 (0.00, 0.00)	0.00 (0.00, 0.00)
Median, Median (IQR)	2.38 (2.32, 2.46)	2.39 (2.35, 2.43)	2.38 (2.36, 2.41)	2.39 (2.37, 2.41)
97.5 th PCT, Median (IQR)	6.77 (6.71, 6.81)	6.77 (6.73, 6.81)	6.78 (6.76, 6.80)	6.78 (6.76, 6.79)
				
Terminal-event time				
2.5 th PCT, Median (IQR)	0.56 (0.50, 0.60)	0.54 (0.50, 0.58)	0.54 (0.51, 0.56)	0.54 (0.52, 0.55)
Median, Median (IQR)	2.65 (2.59, 2.71)	2.65 (2.61, 2.69)	2.65 (2.63, 2.68)	2.65 (2.64, 2.68)
97.5 th PCT, Median (IQR)	6.18 (6.04, 6.30)	6.19 (6.10, 6.29)	6.20 (6.15, 6.26)	6.20 (6.16, 6.23)
				
Censoring time				
Median (IQR)	7.00 (7.00, 7.00)	7.00 (7.00, 7.00)	7.00 (7.00, 7.00)	7.00 (7.00, 7.00)
Terminal event				
%, Median (IQR)	94.20 (93.40, 94.80)	94.00 (93.60, 94.60)	94.12 (93.84, 94.49)	94.14 (93.90, 94.34)
Sex (Female)				
%, Median (IQR)	50.40 (49.60, 51.20)	50.10 (49.50, 51.00)	50.34 (49.96, 50.84)	50.36 (50.04, 50.62)
Age at baseline				
Median (IQR)	50.14 (49.67, 50.79)	50.24 (49.74, 50.72)	50.13 (49.92, 50.47)	50.23 (50.06, 50.45)

IQR: interquartile range; PCT: percentile.

Table S5: Computing time (hours) against the number of individuals and the number of data splits, using 63 cores (200 replicas). The best performance is underlined.

	Number of individuals (n)							
	500		1,000		2,500		5,000	
	Median	IQR	Median	IQR	Median	IQR	Median	IQR
All data	0.655	(0.648, 0.664)	1.137	(1.123, 1.159)	3.489	(3.478, 3.500)	9.035	(9.016, 9.051)
2 splits	<u>0.607</u>	<u>(0.587, 0.626)</u>	<u>0.809</u>	<u>(0.796, 0.827)</u>	1.811	(1.796, 1.825)	3.960	(3.947, 3.975)
5 splits	0.996	(0.962, 1.039)	1.067	(1.031, 1.103)	<u>1.365</u>	<u>(1.334, 1.396)</u>	<u>1.996</u>	<u>(1.957, 2.038)</u>
10 splits	1.699	(1.641, 1.748)	2.171	(2.115, 2.239)	2.191	<u>(2.126, 2.250)</u>	2.335	<u>(2.294, 2.365)</u>

IQR: interquartile range.

39

Table S6: Computing time (hours) against the number of individuals and the number of data splits, using 7 cores (200 replicas). The best performance is underlined.

	Number of individuals (n)							
	500		1,000		2,500		5,000	
	Median	IQR	Median	IQR	Median	IQR	Median	IQR
All data	0.583	(0.579, 0.589)	1.057	(1.052, 1.116)	3.426	(3.342, 3.441)	8.861	(8.823, 8.987)
2 splits	<u>0.438</u>	<u>(0.432, 0.442)</u>	<u>0.661</u>	<u>(0.656, 0.665)</u>	<u>1.647</u>	<u>(1.642, 1.652)</u>	3.709	(3.677, 3.802)
5 splits	0.859	(0.850, 0.870)	1.081	(1.071, 1.089)	1.840	(1.832, 1.853)	3.384	(3.366, 3.428)
10 splits	1.265	(1.257, 1.275)	1.473	(1.460, 1.483)	2.093	(2.076, 2.114)	<u>3.187</u>	<u>(3.170, 3.208)</u>

IQR: interquartile range.

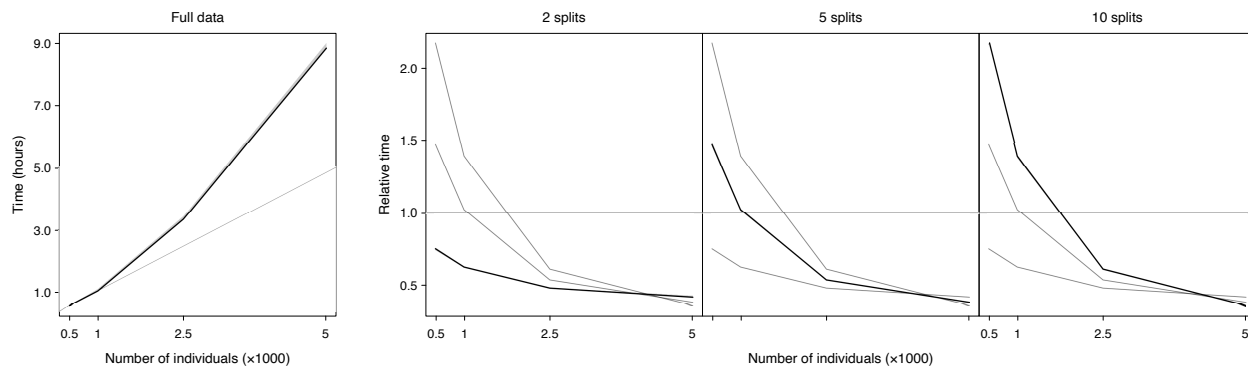


Figure S7: Left: Median computing time, with associated IQR, against the number of individuals, when using all data together, using 7 cores. The gray diagonal line shows a linear evolution. Right: Median and IQR computing time, relative to the time required to fit all data together, against the number of individuals, in the scenario with an unlimited number of cores, using 7 cores. The gray lines show the median time from the remaining panels.

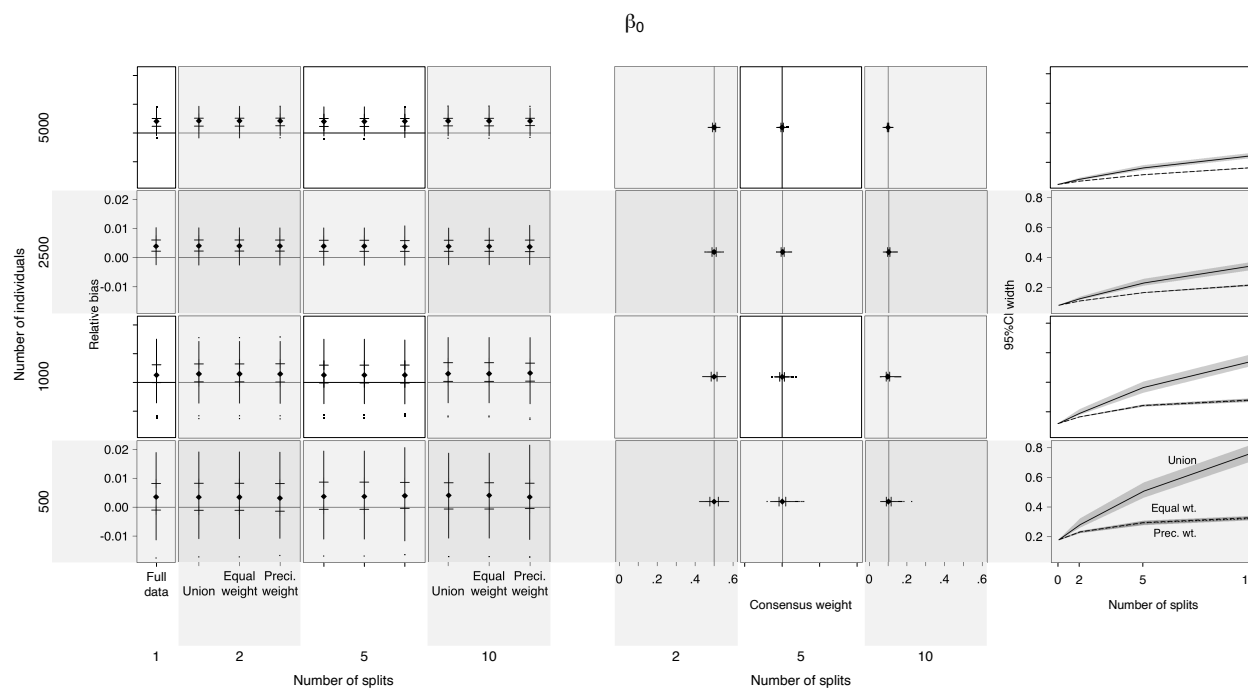


Figure S8: Left: Box plot of the relative bias for the β_0 estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the β_0 estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weights. Right: Median width of the 95% credible interval, with the associated IQR, of the β_0 estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

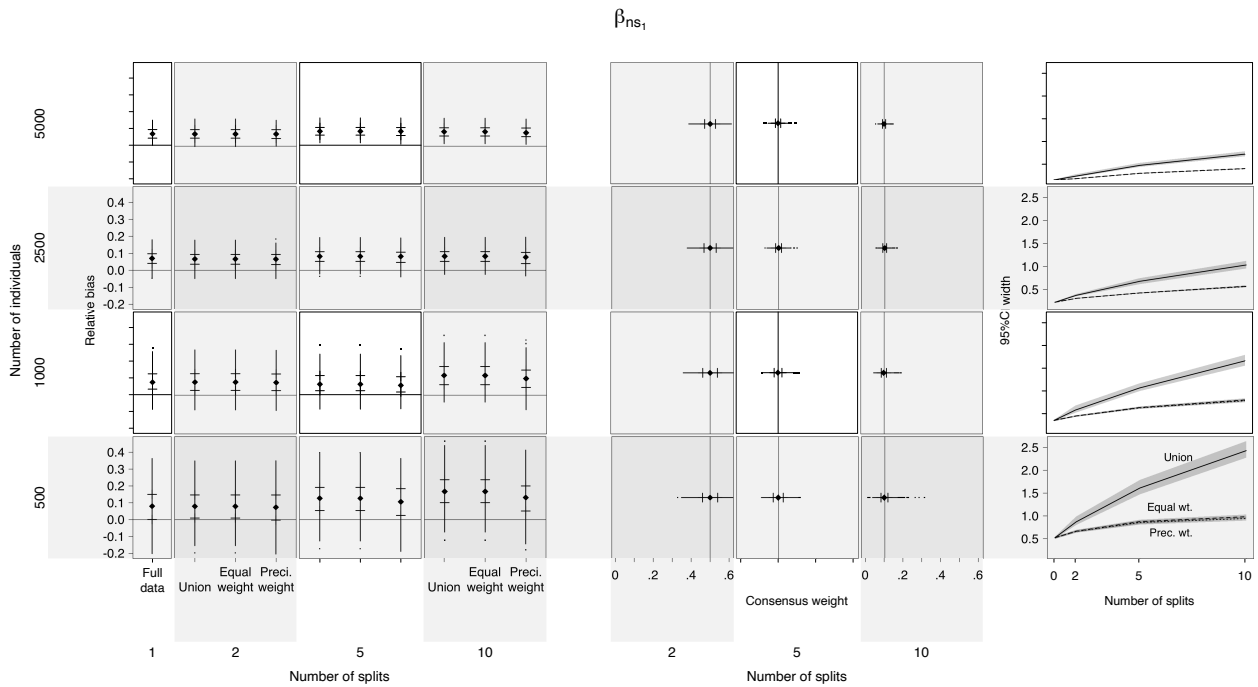


Figure S9: Left: Box plot of the relative bias for the β_{ns_1} estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the β_{ns_1} estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weights. Right: Median width of the 95% credible interval, with the associated IQR, of the β_{ns_1} estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

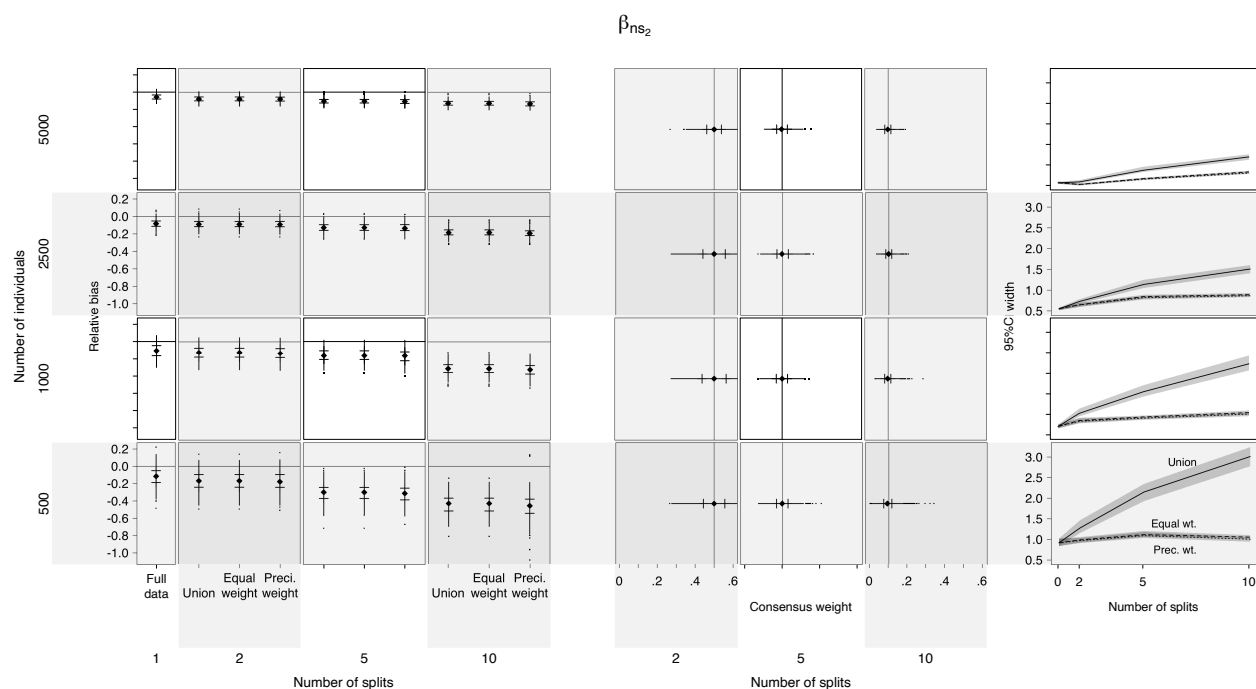


Figure S10: Left: Box plot of the relative bias for β_{ns_2} estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the β_{ns_2} estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weight corresponding for that number of data. Right: Median width of the 95% credible interval, with the associated IQR, of the β_{ns_2} estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

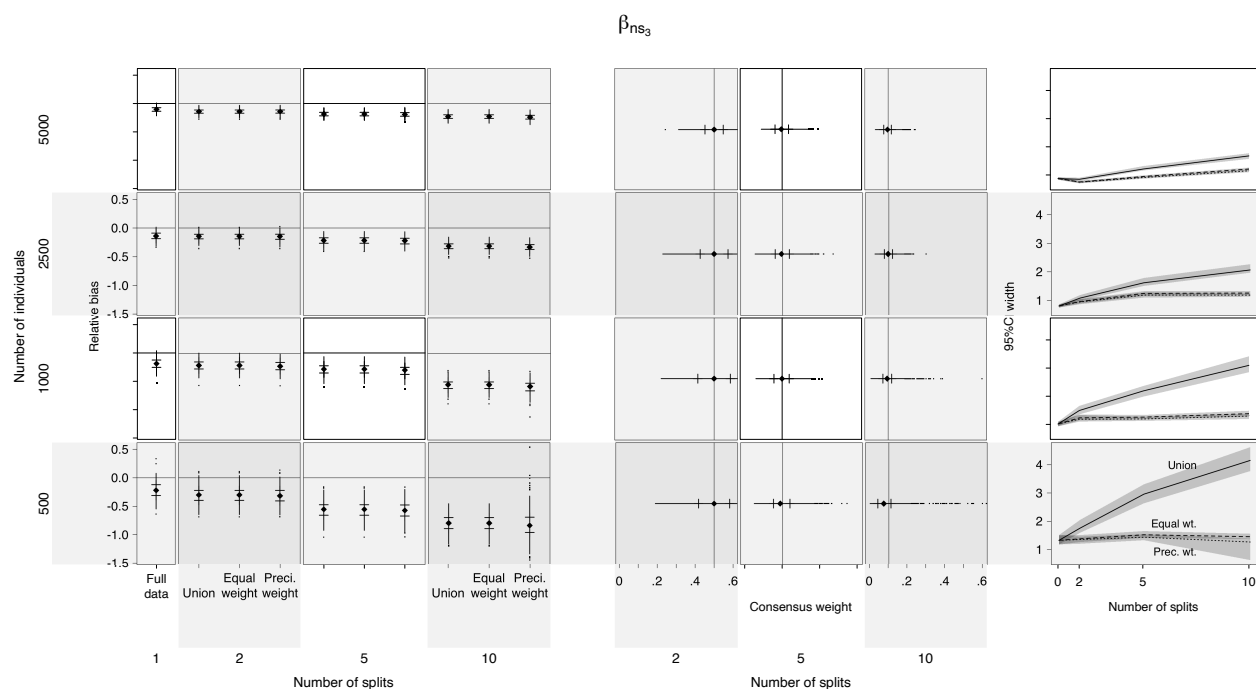


Figure S11: Left: Box plot of the relative bias for β_{ns_3} estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the β_{ns_3} estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weight corresponding for that number of data. Right: Median width of the 95% credible interval, with the associated IQR, of the β_{ns_2} estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

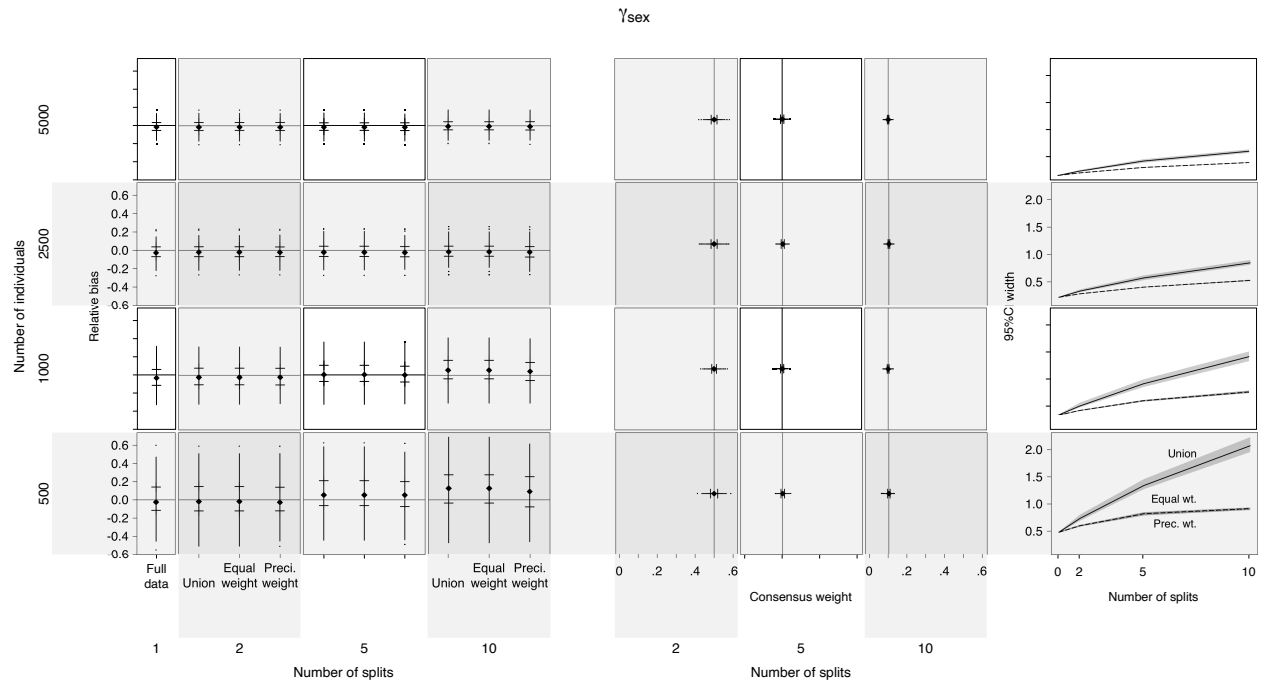


Figure S12: Left: Box plot of the relative bias for the γ_{sex} estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the γ_{sex} estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weights. Right: Median width of the 95% credible interval, with the associated IQR, of the γ_{sex} estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

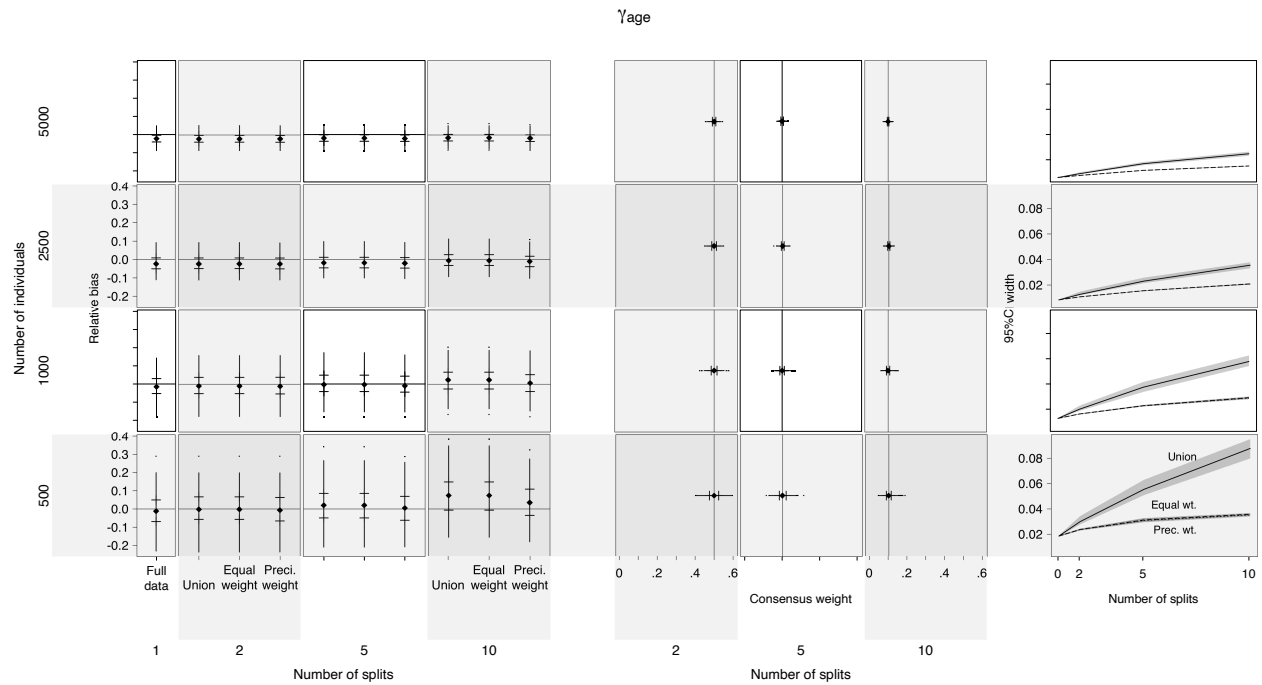


Figure S13: Left: Box plot for the relative bias for γ_{age} estimate for different sample sizes and numbers of data splits. Center: Consensus (standardized) precision weights for the γ_{age} estimate for different sample sizes and numbers of data splits. The vertical lines represent the equal weights. Right: Median width of the 95% credible interval, with the associated IQR, of the γ_{age} estimate—1st, 2nd, and 3rd quartiles—against the number of splits for different sample sizes.

C Application study

Table S7: Estimated posterior means and 95% credible interval for the joint model coefficients obtained from the different consensus methods.

Param.	Consensus methods					
	Union		Equal weight		Precision weight	
	Mean	95% CI	Mean	95% CI	Mean	95% CI
LME						
β_0	92.097	(90.012, 93.915)	92.097	(91.293, 92.913)	92.107	(91.309, 92.919)
β_{ns1}	-288.629	(-298.948, -274.531)	-288.629	(-290.853, -286.377)	-288.644	(-290.858, -286.396)
β_{ns2}	-459.129	(-480.053, -428.993)	-459.129	(-463.062, -455.213)	-459.181	(-463.108, -455.269)
β_{pa}	-1.342	(-1.479, -1.224)	-1.342	(-1.393, -1.290)	-1.342	(-1.393, -1.290)
β_{ins}	-0.154	(-0.361, 0.012)	-0.154	(-0.214, -0.092)	-0.158	(-0.219, -0.097)
β_{dpx}	-0.611	(-2.453, 1.540)	-0.611	(-1.116, -0.090)	-0.632	(-1.136, -0.112)
$\beta_{[88,93]}$	-0.376	(-2.513, 1.910)	-0.376	(-1.389, 0.642)	-0.373	(-1.380, 0.640)
$\beta_{[93,98]}$	-0.045	(-2.056, 2.081)	-0.045	(-1.071, 1.005)	-0.054	(-1.072, 0.991)
$\beta_{[98,11]}$	3.809	(2.021, 5.663)	3.809	(2.993, 4.619)	3.826	(3.017, 4.629)
β_{sex}	1.529	(0.507, 2.570)	1.529	(0.868, 2.177)	1.529	(0.868, 2.176)
β_{htz}	2.427	(1.073, 3.859)	2.427	(1.720, 3.145)	2.435	(1.731, 3.153)
β_{oth}	3.218	(0.807, 6.330)	3.218	(2.258, 4.183)	3.248	(2.287, 4.213)
σ_y	8.844	(8.805, 8.881)	8.844	(8.819, 8.863)	8.844	(8.819, 8.863)
PH						
γ_{sex}	-0.425	(-0.628, -0.252)	-0.425	(-0.519, -0.332)	-0.420	(-0.514, -0.328)
α	-0.120	(-0.132, -0.110)	-0.120	(-0.124, -0.114)	-0.119	(-0.123, -0.113)

CI: credible interval; LME: linear mixed-effects model; PH: proportional hazards model.

Table S8: Follow-up, demographic, social, and clinical characteristics of the CF individuals analyzed.

		Subsamples					
		Full data	1	2	3	4	5
Number of individuals		35,153	7,030	7,031	7,031	7,031	7,030
Birth cohort							
	[1924, 1988]	14,473 (41.17%)	2,905 (41.32%)	2,896 (41.19%)	2,925 (41.60%)	2,885 (41.03%)	2,862 (40.71%)
	[1988, 1993]	4,762 (13.55%)	934 (13.29%)	943 (13.41%)	961 (13.67%)	967 (13.75%)	957 (13.61%)
	[1993, 1998]	4,630 (13.17%)	915 (13.02%)	948 (13.48%)	907 (12.90%)	929 (13.21%)	931 (13.24%)
	[1998, 2011]	11,288 (32.11%)	2,276 (32.38%)	2,244 (31.91%)	2,238 (31.83%)	2,250 (32.00%)	2,280 (32.43%)
Baseline age (years)	Median (IQR)	8.92 (6.23, 18.56)	8.98 (6.23, 18.77)	9.06 (6.23, 18.49)	8.93 (6.23, 19.13)	8.86 (6.23, 18.26)	8.77 (6.22, 18.20)
Terminal event							
	Censoring	26,902 (76.53%)	5,383 (76.57%)	5,370 (76.38%)	5,326 (75.71%)	5,373 (76.42%)	5,450 (77.53%)
	Death or transplantation	8,251 (23.47%)	1,647 (23.43%)	1,661 (23.62%)	1,705 (24.25%)	1,658 (23.58%)	1,580 (22.48%)
Age at end of follow-up (years)							
	Censoring, median (IQR)	21.33 (14.12, 30.94)	21.32 (14.04, 31.49)	21.33 (14.17, 30.63)	21.33 (14.06, 30.86)	21.43 (14.20, 30.88)	21.30 (14.12, 30.96)
	Death/transplantation, median (IQR)	27.12 (21.36, 35.99)	26.69 (21.20, 36.16)	27.61 (21.52, 36.44)	27.47 (21.55, 36.20)	26.82 (21.16, 35.38)	27.10 (21.33, 35.44)
Genotype (F508del)							
	Homozygous	15,656 (44.54%)	3,116 (44.32%)	3,130 (44.52%)	3,165 (45.02%)	3,114 (44.29%)	3,131 (44.54%)
	Heterozygous	14,002 (39.83%)	2,820 (40.11%)	2,774 (39.45%)	2,790 (39.68%)	2,815 (40.04%)	2,803 (39.87%)
	Neither	5,495 (15.63%)	1,094 (15.53%)	1,127 (16.03%)	1,076 (15.30%)	1,102 (15.67%)	1,096 (15.59%)
Sex	Female	16,992 (48.34%)	3,392 (48.25%)	3,375 (48.00%)	3,378 (48.04%)	3,376 (48.02%)	3,471 (49.37%)
Number of ppFEV₁ measurements		1,523,406	304,119	305,358	302,592	304,541	306,796
Number of ppFEV₁ measurements/ind.	Median (IQR)	36.00 (15.00, 64.00)	36.00 (15.00, 63.75)	36.00 (15.00, 64.00)	36.00 (15.00, 62.00)	36.00 (15.00, 64.00)	36.00 (16.00, 63.00)
Total follow-up duration (years)		372,365.79	74,084.40	74,478.90	74,411.97	74,393.33	74,997.19
Follow-up duration/ind. (years)	Median (IQR)	10.28 (4.59, 16.78)	10.26 (4.47, 16.70)	10.33 (4.50, 16.78)	10.16 (4.64, 16.79)	10.19 (4.62, 16.72)	10.43 (4.71, 16.90)
Baseline ppFEV₁	Median (IQR)	86.00 (65.50, 100.70)	86.10 (66.40, 100.80)	86.40 (65.50, 100.80)	85.40 (64.90, 100.25)	86.30 (65.20, 100.80)	85.70 (65.60, 100.90)
Medicaid insurance possession							
	At baseline	16,350 (46.51%)	3,318 (47.20%)	3,233 (45.98%)	3,240 (46.08%)	3,259 (46.35%)	3,300 (46.94%)
	Throughout follow-up	8,677 (24.68%)	1,778 (25.29%)	1,698 (24.15%)	1,708 (24.29%)	1,770 (25.17%)	1,723 (24.51%)
	Sometime during follow-up	28,632 (81.45%)	5,707 (81.18%)	5,714 (81.27%)	5,692 (80.96%)	5,727 (81.45%)	5,792 (82.39%)
<i>Pseudomonas aeruginosa</i>							
	At baseline	3,621 (10.30%)	769 (10.94%)	695 (9.89%)	735 (10.45%)	704 (10.01%)	718 (10.21%)
	Throughout follow-up	235 (0.67%)	66 (0.94%)	48 (0.68%)	43 (0.61%)	36 (0.51%)	42 (0.60%)
	Sometime during follow-up	25,970 (73.88%)	5,208 (74.08%)	5,162 (73.42%)	5,181 (73.69%)	5,171 (73.55%)	5,248 (74.65%)

IQR: interquartile range.