



## casebase: An Alternative Framework For Survival Analysis

Sahir Bhatnagar \*  
McGill University

Maxime Turgeon \*  
University of Manitoba

Jesse Islam  
McGill University

James Hanley  
McGill University

Olli Saarela  
University of Toronto

---

### Abstract

The abstract of the article. \* joint co-authors

*Keywords:* keywords, not capitalized, Java.

---

## 1. Introduction

The purpose of the **casebase** package is to provide practitioners with an easy-to-use software tool to predict the risk (or cumulative incidence (CI)) of an event, for a particular patient. The following points should be noted:

1. hazard ratios
2. If, however, the absolute risks are of interest, they have to be recovered using the semi-parametric Breslow estimator
3. Alternative approaches for fitting flexible hazard models for estimating absolute risks, not requiring this two-step approach? Yes! ([Hanley and Miettinen 2009](#))

([Hanley and Miettinen 2009](#)) propose a fully parametric hazard model that can be fit via logistic regression. From the fitted hazard function, cumulative incidence and, thus, risk functions of time, treatment and profile can be easily derived.

## 2. Theoretical details

As discussed in Hanley & Miettinen (2009), the key idea behind case-base sampling is to discretize the study base into an infinite amount of *person moments*. These person moments are indexed by both an individual in the study and a time point, and therefore each person moment has a covariate profile, an exposure status and an outcome status attached to it. We note that there is only a finite number of person moments associated with the event of interest (what Hanley & Miettinen call the *case series*). The case-base sampling refers to the sampling from the base of a representative finite sample called the *base series*.

As shown by Saarela & Arjas (2015) (and further expanded in Saarela (2016)), writing the likelihood arising from this data-generating mechanism using the framework of non-homogeneous Poisson processes, we eventually reach an expression where each person-moment's contribution is of the form

$$\frac{h(t)^{dN(t)}}{\rho(t) + h(t)},$$

where  $N(t)$  is the counting process associated with the event of interest,  $h(t)$  is the corresponding hazard function, and  $\rho(t)$  is the hazard function for the Poisson process associated with case-base sampling. This parametric form suggests that we can readily estimate log-hazards of the form  $\log(h(t)) = g(t; X)$  using logistic regression, where each observation corresponds to a person moment, the function  $g(t; X)$  is linear in a finite number of parameters, and where we treat  $-\log(\rho(t))$  as an offset.

In Hanley & Miettinen (2009), the authors suggest performing case-base sampling *uniformly*, i.e. to sample the base series uniformly from the study base. In terms of Poisson processes, this sampling strategy corresponds essentially to a time-homogeneous Poisson process with intensity equal to  $b/B$ , where  $b$  is the number of sampled observations in the base series, and  $B$  is the total population-time for the study base. More complex examples are also available; see for example Saarela & Arjas (2015), where the probabilities of the sampling mechanism are proportional to the cardiovascular disease event rate given by the Framingham score.

The **casebase** package fits the family of hazard functions of the form

$$h(t; X) = \exp[g(t; X)]$$

where  $t$  denotes time and  $X$ , the individual's covariate profile. Different functions of  $t$  lead to different parametric hazard models. The simplest of these models is the one-parameter exponential distribution which is obtained by taking the hazard function to be constant over the range of  $t$ .

$$h(t; X) = \exp(\beta_0 + \beta_1 X)$$

The instantaneous failure rate is independent of  $t$ , so that the conditional chance of failure in a time interval of specified length is the same regardless of how long the individual has been in the study. This is also known as the *memoryless property* (Kalbfleisch and Prentice, 2002).

The Gompertz hazard model is given by including a linear term for time:

$$h(t; X) = \exp(\beta_0 + \beta_1 t + \beta_2 X)$$

Use of  $\log(t)$  yields the Weibull hazard which allows for a power dependence of the hazard on time (Kalbfleisch and Prentice, 2002):

$$h(t; X) = \exp(\beta_0 + \beta_1 \log(t) + \beta_2 X)$$

For competing-risk analyses with  $J$  possible events, we can show that each person-moment's contribution of the likelihood is of the form

$$\frac{h_j(t)^{dN_j(t)}}{\rho(t) + \sum_{j=1}^J h_j(t)},$$

where  $N_j(t)$  is the counting process associated with the event of type  $j$  and  $h_j(t)$  is the corresponding hazard function. As may be expected, this functional form is similar to the terms appearing in the likelihood function for multinomial regression.<sup>1</sup>

### 3. Existing packages

Survival analysis is an important branch of applied statistics and epidemiology. Accordingly, there is a vast ecosystem of R packages implementing different methods. In this section, we describe how the functionalities of **casebase** compare to these packages.

At the time of writing, a cursory examination of CRAN's task view on survival analysis reveals that there are over 250 packages related to survival analysis (2019). For the purposes of this article, we restricted our description to packages that implement at least one of the following features: parametric modeling, non-proportional hazard models, competing risk analysis, penalized estimation, and cumulative incidence curve estimation. By searching for appropriate keywords in the DESCRIPTION file of these packages, we found 60 relevant packages. These 60 packages were then manually examined to determine which ones are comparable to **casebase**. In particular, we excluded packages that were focused on a different set of problems, such as frailty and multi-state models. The remaining 14 packages appear in Table 1, along with some of the functionalities they offer.

Several packages implement penalized estimation for the Cox model: **glmnet** (2011), **glm-path** (2018), **penalized** (2010), **RiskRegression** (2019). Moreover, some packages also include penalized estimation in the context of Cox models with time-varying coefficients: elastic-net penalization with **CoxRidge** (2015) and **rstpm2** (2019), while **survival** (2015) has an implementation for ridge. **casebase** provides penalized estimation in the context of parametric hazards. To our knowledge, **casebase** and **rsptm2** are the only packages to offer this functionality. **I can dive deeper to determine if this is a fair statement**

Parametric survival models are implemented in a handful of packages: **CFC** (2015), **flexsurv** (2016), **SmoothHazard** (2017), **rsptm2** (2019), **mets** (2014), and **survival**. The types of models they allow vary for each package. For example, **SmoothHazard** is limited to Weibull distributions (2017), whereas both **flexsurv** and **survival** allow users to supply any distribution of their choice. Also, **flexsurv**, **smoothhazard**, **mets** and **rstpm2** also have the ability to model the effect of time using splines, which allows flexible modeling of the hazard function. Moreover,

---

<sup>1</sup>Specifically, it corresponds to the following parametrization:

$$\log \left( \frac{P(Y = j | X)}{P(Y = J | X)} \right) = X^T \beta_j, \quad j = 1, \dots, J - 1$$

**flexsurv** has the ability to estimate both scale and shape parameters for a variety of parametric families (see Table 2). As discussed above, **casebase** can model any parametric family whose log-hazard can be expressed as a linear model of covariates (including time). Therefore, our package allows the user to model the effect of time using splines, and through interaction terms involving covariates and time, it also allows user to fit time-varying coefficient models. However, we do not explicitly model any shape parameter, unlike **flexsurv**.

Of the methods mentioned so far, only **CFC**, **flexsurv**, **mets**, **cmprsk** and **survival** contain implementations for competing risks. The differentiating factor between these packages and **casebase** is that **casebase** handles competing risks through multiple logistic regression. Note that **CFC** has extra features outside the scope of the methods discussed in this paper.

Finally, several packages include functions to estimate the cumulative incidence function. The corresponding methods generally fall into two categories: transformation of the estimated hazard function, and semi-parametric estimation of the baseline hazard. The first category broadly corresponds to parametric survival models, where the full hazard is explicitly modeled. Using this estimate, the survival function and the cumulative incidence function can be obtained using their functional relationships (see Equations 1 and 2 below). Packages including this functionality include **CFC**, **Flexsurv**, **mets**, and **survival**. Our package **casebase** also follows this approach for both single-event and competing-event analyses. The second category outlined above broadly corresponds to Cox models. These models do not model the full hazard function, and therefore the baseline hazard needs to be estimated separately in order to estimate the survival function. This is achieved using semi-parametric estimators (e.g. Breslow’s estimator). Packages that implement this approach include **RiskRegression**, **rstpm2**, and **survival**. As mentioned in the introduction, a key distinguishing factor between these two approaches is that the first category leads to smooth estimates of the cumulative incidence function, whereas the second category produces estimates in the form of step-wise functions. This was one of the main motivations for introducing case-base sampling in survival analysis.

Package	Competing risks	Non-proportional	Penalization	Splines	Parametric	Semi-parametric	Interval/left censoring	Absolute risk
<b>casebase</b> Hanley and Miettinen (2009)	x	x	x	x	x			x
<b>CFC</b> Mahani and Sharabiani (2015)	x	x			x			x
<b>textbfcmprsk</b> Gray (2019)	x					x		x
<b>coxRidge</b> Perperoglou (2015)		x	x			x		
<b>crp</b> Fu (2015)	x		x					
<b>fastcox</b> Yi and Zou (2017)			x			x		
<b>Flexrsurv</b> Clerc-Urmès, Grzebyk, and Hédelin (2017)		x		x	x			x
<b>Flexsurv</b> Jackson (2016)	x	x		x	x			x
<b>glmnet</b> Simon <i>et al.</i> (2011)			x			x		
<b>glmpath</b> Park and Hastie (2018)			x			x		
<b>mets</b> Scheike <i>et al.</i> (2014)	x			x		x		x
<b>penalized</b> Goeman, Meijer, Chaturvedi, and Lueder (2019)			x			x		
<b>RiskRegression</b> Gerds <i>et al.</i> (2019)			x			x		x
<b>Rstpm2</b> Clements <i>et al.</i> (2019)		x	x	x	x	x	x	x
<b>SmoothHazard</b> Touraine <i>et al.</i> (2017)		x		x	x		x	
<b>Survival</b> Therneau (2015)	x	x			x	x	x	x

Table 1: Different features of interest in various survival packages.

## 4. Implementation details

The functions in the casebase package can be divided into two categories: 1) data visualization, in the form of population-time plots; and 2) parametric modeling. We explicitly aimed at being compatible with both **data.frames** and **data.tables**. This is evident in some of the

	Parameters (location in <b>red</b> )	Density R function	dist
Exponential	<b>rate</b>	dexp	"exp"
Weibull (accelerated failure time)	<b>shape</b> , <b>scale</b>	dweibull	"weibull"
Weibull (proportional hazards)	<b>shape</b> , <b>scale</b>	dweibullPH	"weibullPH"
Gamma	<b>shape</b> , <b>rate</b>	dgamma	"gamma"
Log-normal	<b>meanlog</b> , sdlog	dlnorm	"lnorm"
Gompertz	<b>shape</b> , <b>rate</b>	dgomptertz	"gomptertz"
Log-logistic	<b>shape</b> , <b>scale</b>	dllogis	"llogis"
Generalized gamma (Pren- tice 1975)	<b>mu</b> , sigma, Q	dgengamma	"gengamma"
Generalized gamma (Stacy 1962)	<b>shape</b> , <b>scale</b> , k	dgengamma.orig	"gengamma.orig"
Generalized F (stable)	<b>mu</b> , sigma, Q, P	dgenf	"genf"
Generalized F (original)	<b>mu</b> , sigma, s1, s2	dgenf.orig	"genf.orig"

Table 2: Built-in parametric survival distributions in **flexsurv**.

coding choices we made, and it is also reflected in our unit tests.

#### 4.1. Population-time plots

#### 4.2. Parametric modeling

The parametric modeling step was separated into three parts:

1. case-base sampling;
2. estimation of the smooth hazard function;
3. calculation of the risk function.

By separating the sampling and estimation functions, we allowed the possibility of users implementing more complex sampling scheme, as described in Saarela (2016).

The sampling scheme selected for **sampleCaseBase** was described in Hanley and Miettinen (2009): we first sample along the “person” axis, proportional to each individual’s total follow-up time, and then we sample a moment uniformly over their follow-up time. This sampling scheme is equivalent to the following picture: imagine representing the total follow-up time of all individuals in the study along a single dimension, where the follow-up time of the next individual would start exactly when the follow-up time of the previous individual ends. Then the base series could be sampled uniformly from this one-dimensional representation of the overall follow-up time. In any case, the output is a dataset of the same class as the input, where each row corresponds to a person-moment. The covariate profile for each such person-moment is retained, and an offset term is added to the dataset. This output could then be used to fit a smooth hazard function, or for visualization of the base series.

The fitting function `fitSmoothHazard` starts by looking at the class of the dataset: if it was generated from `sampleCaseBase`, it automatically inherited the class `cbData`. If the dataset supplied to `fitSmoothHazard` does not inherit from `cbData`, then the fitting function starts by calling `sampleCaseBase` to generate the base series. In other words, the occasional user can bypass `sampleCaseBase` altogether and only worry about the fitting function `fitSmoothHazard`.

The fitting function retains the familiar formula interface of `glm`. The left-hand side of the formula should be the name of the column corresponding to the event type. The right-hand side can be any combination of the covariates, along with an explicit functional form for the time variable. Note that non-proportional hazard models can be achieved at this stage by adding an interaction term involving time. The offset term does not need to be specified by the user, as it is automatically added to the formula.

To fit the hazard function, we provide several approaches that are available via the `family` parameter. These approaches are:

- `glm`: This is the familiar logistic regression.
- `glmnet`: This option allows for variable selection using Lasso or elastic-net. This functionality is provided through the `glmnet` package (Friedman, Hastie, and Tibshirani 2010).
- `gam`: This option provides support for *Generalized Additive Models* via the `gam` package (Hastie and Tibshirani 1987).
- `gbm`: This option provides support for *Gradient Boosted Trees* via the `gbm` package. This feature is still experimental.

In the case of multiple events, the hazard is fitted via multinomial regression as performed by the **VGAM** package. This package was selected for its ability to fit multinomial regression models with an offset.

Once a model-fit object has been returned by `fitSmoothHazard`, all the familiar summary and diagnostic functions are available: `print`, `summary`, `predict`, `plot`, etc. Our package provides one more functionality: it computes risk functions from the model fit. For the case of a single event, it uses the familiar identity

$$S(t) = \exp \left( - \int_0^t h(u; X) du \right). \quad (1)$$

The integral is computed using either the `stats::integrate` function or Monte-Carlo integration. The risk function (or cumulative incidence function) is then defined as

$$CI(t) = 1 - S(t). \quad (2)$$

For the case of a competing-event analysis, the event-specific risk is computed using the following procedure: first, we compute the overall survival function (i.e. for all event types):

$$S(t) = \exp \left( - \int_0^t H(u; X) du \right), \quad H(t; X) = \sum_{j=1}^J h_j(t; X).$$

From this, we can derive the event-specific subdensities:

$$f_j(t) = h_j(t)S(t).$$

Finally, by integrating these subdensities, we obtain the event-specific cumulative incidence functions:

$$CI_j(t) = \int_0^t f_j(u)du.$$

We created `absoluteRisk` as an `S3` generic, with methods for the different types of outputs of `fitSmoothHazard`. The method dispatch system of R then takes care of matching the correct output to the correct methodology for calculating the cumulative incidence function, without the user's intervention.

In the following sections, we illustrate these functionalities in the context of three case studies.

## 5. Case study 1—European Randomized Study of Prostate Cancer Screening

For the first case study, we will use of the European Randomized Study of Prostate Cancer Screening data. This dataset is available through the `casebase` package:

```
R> data(ERSPC)
R> ERSPC$ScrArm <- factor(ERSPC$ScrArm,
R>                        levels = c(0, 1),
R>                        labels = c("Control group", "Screening group"))
```

The results of this study were published by (Schröder, Hugosson, Roobol, Tammela, Ciatto, Nelen, Kwiatkowski, Lujan, Lilja, Zappa *et al.* 2009), and the dataset itself was obtained using the approach described in (Liu, Rich, and Hanley 2014).

Population time plots can be extremely informative graphical displays of survival data. They should be the first step in an exploratory data analysis. We facilitate this task in the `casebase` package using the `popTime` function. We first create the necessary dataset for producing the population time plots, and we can generate the plot by using the corresponding `plot` method:

```
R> pt_object <- casebase::popTime(ERSPC, event = "DeadOfPrCa")
R> plot(pt_object)
```

We can also create exposure stratified plots by specifying the `exposure` argument in the `popTime` function:

```
R> pt_object_strat <- casebase::popTime(ERSPC,
R>                                     event = "DeadOfPrCa",
R>                                     exposure = "ScrArm")
R> plot(pt_object_strat)
```

We can also plot them side-by-side using the `ncol` argument:

```
R> plot(pt_object_strat, ncol = 2)
```

ADD PARAGRAPH ABOUT WHAT WE CAN CONCLUDE FROM THESE GRAPHS.

Next, we investigate the differences between the control and the screening arms. A common choice for this type of analysis is to use a Cox regression model and estimate the hazard ratio for the screening group (relative to the control group). In R, it can be done as follows:

```
R> cox_model <- survival::coxph(Surv(Follow.Up.Time, DeadOfPrCa) ~ ScrArm,
R>                               data = ERSPC)
R> summary(cox_model)
```

We can also plot the cumulative incidence function (CIF) for each group. However, this involves estimating the baseline hazard, which Cox regression treated as a nuisance parameter.

```
R> new_data <- data.frame(ScrArm = c("Control group", "Screening group"),
R>                          ignore = 99)
R>
R> plot(survfit(cox_model, newdata = new_data),
R>       xlab = "Years since Randomization",
R>       ylab = "Cumulative Incidence (%)",
R>       fun = "event",
R>       xlim = c(0,15), conf.int = FALSE, col = c("red","blue"),
R>       main = sprintf("Estimated Cumulative Incidence (risk) of Death from Prostate
R>                      Cancer Screening group Hazard Ratio: %.2g (%.2g, %.2g)",
R>                      exp(coef(cox_model)),
R>                      exp(confint(cox_model))[1],
R>                      exp(confint(cox_model))[2]),
R>       ylim = c(0, 2), yscale = 100)
R> legend("topleft",
R>       legend = c("Control group", "Screening group"),
R>       col = c("red","blue"),
R>       lty = c(1, 1),
R>       bg = "gray90")
```

Next, we show how case-base sampling can be used to carry out a similar analysis. We fit several models that differ in how we choose to model time.

First, we will fit an exponential model. Recall that this corresponds to excluding time from the linear predictor.

```
R> casebase_exponential <- casebase::fitSmoothHazard(DeadOfPrCa ~ ScrArm,
R>                                                    data = ERSPC,
R>                                                    ratio = 100)
R>
R> summary(casebase_exponential)
R> exp(coef(casebase_exponential)[2])
R> exp(confint(casebase_exponential)[2,])
```



We can then use the `absoluteRisk` function to get an estimate of the cumulative incidence curve for a specific covariate profile. In the plot below, we overlay the estimated CIF from the exponential model on the Cox model CIF:

```
R> smooth_risk_exp <- casebase::absoluteRisk(object = casebase_exponential,
R>                                           time = seq(0,15,0.1),
R>                                           newdata = new_data)
R>
R> plot(survfit(cox_model, newdata = new_data),
R>       xlab = "Years since Randomization",
R>       ylab = "Cumulative Incidence (%)",
R>       fun = "event",
R>       xlim = c(0,15), conf.int = FALSE, col = c("red","blue"),
R>       main = sprintf("Estimated Cumulative Incidence (risk) of Death from Prostate
R>                     Cancer Screening group Hazard Ratio: %.2g (%.2g, %.2g)",
R>                     exp(coef(cox_model)),
R>                     exp(confint(cox_model))[1],
R>                     exp(confint(cox_model))[2])),
R>       ylim = c(0, 2), yscale = 100)
R> lines(smooth_risk_exp[,1], smooth_risk_exp[,2], col = "red", lty = 2)
R> lines(smooth_risk_exp[,1], smooth_risk_exp[,3], col = "blue", lty = 2)
R>
R>
R> legend("topleft",
R>       legend = c("Control group (Cox)", "Control group (Casebase)",
R>                 "Screening group (Cox)", "Screening group (Casebase)"),
R>       col = c("red", "red", "blue", "blue"),
R>       lty = c(1, 2, 1, 2),
R>       bg = "gray90")
```

As we can see, the exponential model gives us an estimate of the hazard ratio that is similar to the one obtained using Cox regression. However, the CIF estimates are quite different. Based on what we observed in the population time plot, where more events are observed later on in time, we do not expect a constant hazard would be a good description for this data. In other words, this poor fit for the exponential hazard is expected. A constant hazard model would overestimate the cumulative incidence earlier on in time, and underestimate it later on; this is what we see on the cumulative incidence plot. This example demonstrates the benefits of population time plots as an exploratory analysis tool.

For the next model, we include time as a linear term. Recall that this corresponds to a Weibull model.

```
R> casebase_time <- fitSmoothHazard(DeadOfPrCa ~ Follow.Up.Time + ScrArm,
R>                                   data = ERSPC,
R>                                   ratio = 100)
R>
R> summary(casebase_time)
```

```
R> exp(coef(casebase_time))
R> exp(confint(casebase_time))
```

Again, the estimate of the hazard ratio is similar to that obtained from Cox regression. We then look at the estimate of the CIF:

```
R> smooth_risk_time <- casebase::absoluteRisk(object = casebase_time,
R>                                           time = seq(0,15,0.1),
R>                                           newdata = new_data)
R>
R> plot(survfit(cox_model, newdata = new_data),
R>      xlab = "Years since Randomization",
R>      ylab = "Cumulative Incidence (%)",
R>      fun = "event",
R>      xlim = c(0,15), conf.int = FALSE, col = c("red","blue"),
R>      main = sprintf("Estimated Cumulative Incidence (risk) of Death from Prostate
R>                    Cancer Screening group Hazard Ratio: %.2g (%.2g, %.2g)",
R>                    exp(coef(cox_model)),
R>                    exp(confint(cox_model))[1],
R>                    exp(confint(cox_model))[2])),
R>      ylim = c(0, 2), yscale = 100)
R> lines(smooth_risk_time[,1], smooth_risk_time[,2], col = "red", lty = 2)
R> lines(smooth_risk_time[,1], smooth_risk_time[,3], col = "blue", lty = 2)
R>
R> legend("topleft",
R>       legend = c("Control group (Cox)", "Control group (Casebase)",
R>                  "Screening group (Cox)", "Screening group (Casebase)"),
R>       col = c("red", "red", "blue", "blue"),
R>       lty = c(1, 2, 1, 2),
R>       bg = "gray90")
```

We see that the Weibull model leads to a better fit of the CIF than the Exponential model, when compared with the semi-parametric approach.

Finally, we can model the hazard as a smooth function of time using the `splines` package:

```
R> casebase_splines <- fitSmoothHazard(DeadOfPrCa ~ bs(Follow.Up.Time) + ScrArm,
R>                                     data = ERSPC,
R>                                     ratio = 100)
R>
R> summary(casebase_splines)
R> exp(coef(casebase_splines))
R> exp(confint(casebase_splines))
```

Once again, we can see that the estimate of the hazard ratio is similar to that from the Cox regression. We then look at the estimate of the CIF:

```

R> smooth_risk_splines <- absoluteRisk(object = casebase_splines,
R>                                     time = seq(0,15,0.1),
R>                                     newdata = new_data)
R>
R> plot(survfit(cox_model, newdata = new_data),
R>       xlab = "Years since Randomization",
R>       ylab = "Cumulative Incidence (%)",
R>       fun = "event",
R>       xlim = c(0,15), conf.int = FALSE, col = c("red","blue"),
R>       main = sprintf("Estimated Cumulative Incidence (risk) of Death from Prostate
R>                     Cancer Screening group Hazard Ratio: %.2g (%.2g, %.2g)",
R>                     exp(coef(cox_model)),
R>                     exp(confint(cox_model))[1],
R>                     exp(confint(cox_model))[2])),
R>       ylim = c(0, 2), yscale = 100)
R> lines(smooth_risk_splines[,1], smooth_risk_splines[,2], col = "red", lty = 2)
R> lines(smooth_risk_splines[,1], smooth_risk_splines[,3], col = "blue", lty = 2)
R>
R> legend("topleft",
R>       legend = c("Control group (Cox)", "Control group (Casebase)",
R>                 "Screening group (Cox)", "Screening group (Casebase)"),
R>       col = c("red", "red", "blue", "blue"),
R>       lty = c(1, 2, 1, 2),
R>       bg = "gray90")

```

Qualitatively, the combination of splines and case-base sampling produces the parametric model that most closely resembles the Cox model. In the following table, we see that the confidence intervals are also similar across all four models. In other words, this reinforces the idea that, under proportional hazards, we do not need to model the full hazard to obtain reliable estimates of the hazard ratio. Of course, different parametric models for the hazards give rise to qualitatively different estimates for the CIF.

As noted above, the usual asymptotic results hold for likelihood ratio tests built using case-base sampling models. Therefore, we can easily test the null hypothesis that the exponential model is just as good as the larger (in terms of number of parameters) splines model.

```

R> anova(casebase_exponential, casebase_splines, test = "LRT")

```

As expected, we see that splines model provides a better fit. Similarly, we can compare the AIC for all three case-base models:

```

R> c("Exp." = AIC(casebase_exponential),
R>    "Weib." = AIC(casebase_time),
R>    "Splines" = AIC(casebase_splines))

```

Once again, we have evidence that the splines model provides the best fit.

With this first case study, we saw how population-time plots can provide a useful summary of the incidence for our event of interests. In particular, we saw how it can provide evidence

against an exponential hazard model. We then explored how different parametric models can provide different estimates of the CIF, even though they all give similar estimates for the hazard ratios.

## 6. Case study 2—Bone-marrow transplant

In the next case study, we show how case-base sampling can also be used in the context of a competing risk analysis. For illustrative purposes, we will use the same data that was used in Scrucca *et al* (2010). The data was downloaded from the first author’s website, and it is also available as part of the **casebase** package.

```
R>
R> data(bmtcrr)
```

The data contains information on 177 patients who received a stem-cell transplant for acute leukemia. The event of interest is relapse, but other competing causes (e.g. transplant-related death) were also recorded. Several covariates were captured at baseline: sex, disease type (acute lymphoblastic or myeloblastic leukemia, abbreviated as ALL and AML, respectively), disease phase at transplant (Relapse, CR1, CR2, CR3), source of stem cells (bone marrow and peripheral blood, coded as BM+PB, or only peripheral blood, coded as PB), and age. A summary of these baseline characteristics appear in Table ???. We note that the statistical summaries were generated differently for different variable types: for continuous variables, we gave the range, followed by the mean and standard deviation; for categorical variables, we gave the counts for each category.

First, we can look at a population-time plot to visualize the incidence density of both relapse and the competing events. In Figure ??, failure times are highlighted on the plot using red points for the event of interest (panel A) and blue dots for competing events (panel B). In both panels, we see evidence of a non-constant hazard function: the density of points is larger at the beginning of follow-up than at the end.

Our main objective is to compute the absolute risk of relapse for a given set of covariates. First, we fit a smooth hazard to the data using a linear term for time:

```
R> model_cb <- fitSmoothHazard(
R>   Status ~ ftime + Sex + D + Phase + Source + Age,
R>   data = bmtcrr,
R>   ratio = 100,
R>   time = "ftime")
```

We will compare our hazard ratio estimates to that obtained from a Cox regression. To do so, we need to treat the competing event as censoring.

```
R> library(survival)
R> # Treat competing event as censoring
R> model_cox <- coxph(Surv(ftime, Status == 1) ~ Sex + D + Phase + Source + Age,
R>   data = bmtcrr)
```

From the fit object, we can extract both the hazard ratios and their corresponding confidence intervals. These quantities appear in Table ???. As we can see, the only significant hazard ratio identified by case-base sampling is the one associated with the phase of the disease at transplant. More precisely, being in relapse at transplant is associated with a hazard ratio of 3.89 when compared to CR1.

Given the estimate of the hazard function obtained using case-base sampling, we can compute the absolute risk curve for a fixed covariate profile. We perform this computation for a 35 year old woman who received a stem-cell transplant from peripheral blood at relapse. We compared the absolute risk curve for such a woman with acute lymphoblastic leukemia with that for a similar woman with acute myeloblastic leukemia. We will estimate the curve from 0 to 60 months.

```
R> # Pick 100 equidistant points between 0 and 60 months
R> time_points <- seq(0, 60, length.out = 50)
R>
R> # Data.frame containing risk profile
R> newdata <- data.frame("Sex" = factor(c("F", "F"),
R>                                levels = levels(bmtcrr[, "Sex"])),
R>                                "D" = c("ALL", "AML"), # Both diseases
R>                                "Phase" = factor(c("Relapse", "Relapse"),
R>                                levels = levels(bmtcrr[, "Phase"])),
R>                                "Age" = c(35, 35),
R>                                "Source" = factor(c("PB", "PB"),
R>                                levels = levels(bmtcrr[, "Source"])))
R>
R> # Estimate absolute risk curve
R> risk_cb <- absoluteRisk(object = model_cb, time = time_points,
R>                          method = "numerical", newdata = newdata)
```

We will compare our estimates to that obtained from a corresponding Fine-Gray model ([Fine and Gray 1999](#)). The Fine-Gray model is a semiparametric model for the cause-specific *subdistribution*, i.e. the function  $f_k(t)$  such that

$$CI_k(t) = 1 - \exp\left(-\int_0^t f_k(u)du\right),$$

where  $CI_k(t)$  is the cause-specific cumulative incidence. The Fine-Gray model allows to directly assess the effect of a covariate on the subdistribution, as opposed to the hazard. For the computation, we will use the **timereg** package ([Scheike and Zhang 2011](#)).

```
R> library(timereg)
R> model_fg <- comp.risk(Event(ftime, Status) ~ const(Sex) + const(D) +
R>                        const(Phase) + const(Source) + const(Age),
R>                        data = bmtcrr, cause = 1, model = "fg")
R>
R> # Estimate absolute risk curve
R> risk_fg <- predict(model_fg, newdata, times = time_points)
```

Figure ?? shows the absolute risk curves for both case-base sampling and the Fine-Gray model. As we can see, the two approaches agree quite well for acute myeloblastic leukemia; however, there seems to be a difference of about 5% between the two curves for acute lymphoblastic leukemia. This difference does not appear to be significant: the curve from case-base sampling is contained within a 95% confidence band around the Fine-Gray absolute risk curve (figure not shown).

## 7. Case study 3—SUPPORT Data

In the first two case studies, we described the basic functionalities of the **casebase** package: creating population-time plots, fitting parametric models for hazard functions, and estimating the corresponding cumulative incidence curves. In the next two case studies, we explore more advanced features.

For the third case study, we will introduce regularization into the estimation of the hazard function. This can be performed using a combination of case-base sampling and regularized logistic regression. To demonstrate this capability, we examine the Study to Understand Prognoses Preferences Outcomes and Risks of Treatment (SUPPORT) dataset (Knaus, Harrell, Lynn, Goldman, Phillips, Connors, Dawson, Fulkerson, Califf, Desbiens *et al.* 1995). The SUPPORT dataset tracks death for individuals who are considered seriously ill within a hospital. Imputation was conducted using the **mice** package in R with default settings (van Buuren and Groothuis-Oudshoorn 2011). Before imputation, there were a total of 1000 individuals and 35 variables. charges, Totcst, totmcst and sfdm2 were removed as they are response variables we are not interested in for the sake of this analysis. After imputation, we have 690 individuals. We lose individuals due to colinearity between variables, preventing imputation. If a pair of covariates are deemed colinear to a certain degree, they are excluded from any imputations. Here, only edu income avtisst\* race ph\* glucose\* adlp\* were imputed. Of these 690 individuals, 66.09% experienced the event of interest. A description of each variable can be found in Table 3 and a breakdown of each categorical variable in Table 4. The original data was obtained from the Department of Biostatistics at Vanderbilt University. The imputed dataset is available as part of the **casebase** package.

As before, the first step is to visualize the incidence density. The population-time plot for the SUPPORT data appears in Figure 1.

Before we proceed with the analysis, we first remove hospital death as a covariate, as this is directly informative of death. All covariates will be used to estimate the hazard function. The hazard function will be used to estimate the cumulative incidence of each individual in the study. Then, the average cumulative incidence curve for all observations will be used for our comparison.

Our main objective is to compute the absolute risk of death for a given set of covariates, using regularization when fitting our model. First, we fit a smooth hazard to the data using a weibull distribution. As casebase makes use of the glmnet package, we interact with the fitsmoothhazard.fit function using a matrix interface, where y contains the time and event variables and x contains all other variables we would like to include in the model. For the SUPPORT dataset, all factor variables were converted into indicator variables. Here, we used lasso penalization by setting alpha to 1.

R>

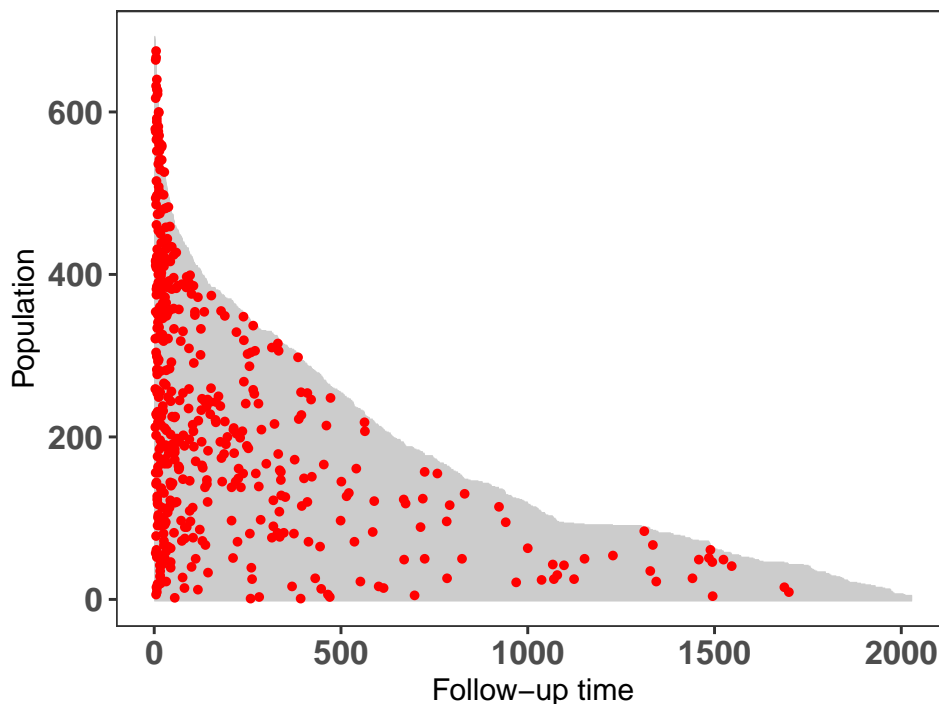


Figure 1: Population-time plot for the SUPPORT data. Red points represent the case series. The gray space demonstrates the base series.

```
R> cb.Model=casebase::fitSmoothHazard.fit(x,y,family="glmnet",time="d.time",event="death",
R+                                     formula_time = ~log(d.time),alpha=1,ratio=ratio)
```

Then, we take our model and use the `absoluteRisk` function to integrate and retrieve our cumulative incidence curves.

We will compare this curve to a regularized version of cox regression, and a kaplan-meier survival curve. The cox regression absolute risk curve with lasso penalization required two extra steps to retrieve. Coxnet from the `glmnet` package has tools to fit a regularized cox model, but requires a survival object skeleton of the selected variables, so that the survival package tools can be used to retrieve the absolute risk. We handle this by fitting a regularized cox model with the `glmnet` package, and a cox model from the `survival` package with the selected variables from the regularized fit. This second model serves as a skeleton that permits the use of survival package functions. The coefficients from the regularized model will replace the coefficients in our skeleton. the resulting model will have the correct coefficients when integrating, but will have an incorrect standard error. For the purposes of this case study, we are only interested in the absolute risk curve itself and not the standard error.

We used the `survival` package to calculate cox regression's cumulative incidence without penalization and the kaplan-meier curve for incidence.

```
R>
R> coxFitup <- survival::coxph(Surv(time=V1,event=V2) ~ ., data = Cox)
R> abcoxFitup<-survival::survfit(coxFitup,newdata=Cox)
R>
```

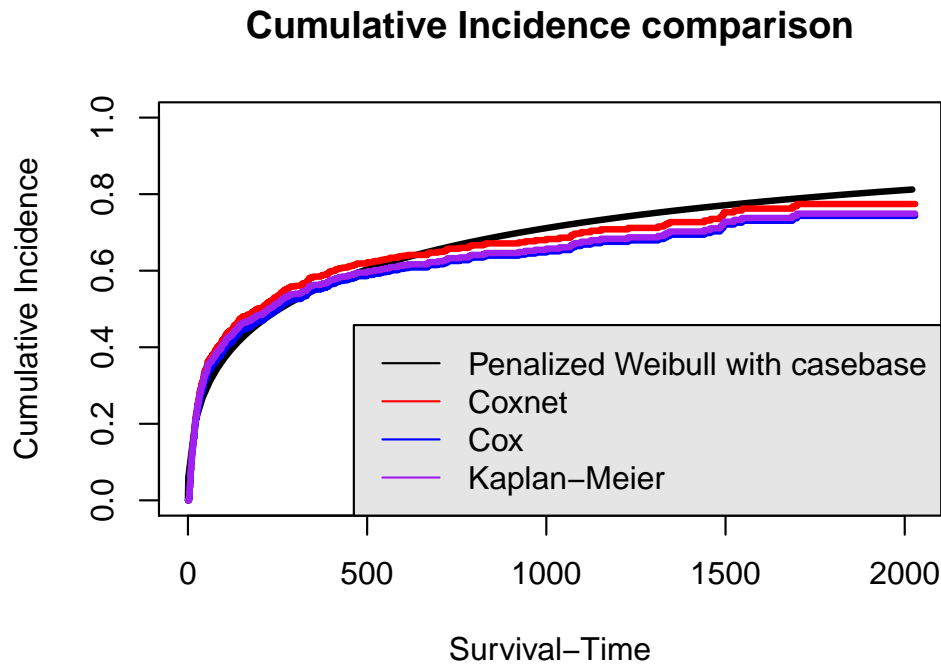


Figure 2: Compares regularized casebase in black, regularized cox in red and Kaplan meier in blue.

```
R> #creating a surv object to be used to fit an unadjusted absolute risk curve.
R> # (Kaplan-Meier risk curve)
R> km <- survival::Surv(time = completeData$d.time, event = completeData$death)
R> abKm<-survival::survfit(km~1,type='kaplan-meier',conf.type='log')
```

Now that our four cumulative incidence estimations have been calculated, we compare them all in Figure 2. Both coxnet and casebase decrease the absolute risk by the end of the study, in comparison to kaplan-meier.

```
R>
R> #get covariates excluding d.time
R> estimatescb=setDT(as.data.frame(coef(cb.Model)[-c(1,2),1]), keep.rownames = TRUE)[]
R> estimatescb$Model="casebase"
R> colnames(estimatescb)<-c("Coefs","Estimate","Model")
R> estimatescox=setDT(as.data.frame(coef(coxglmFit)[-c(1),1]), keep.rownames = TRUE)[]
R> estimatescox$Model="Cox"
R> colnames(estimatescox)<-c("Coefs","Estimate","Model")
R> estimates=rbind(estimatescb, estimatescox)
R> estimates$Estimate[estimates$Estimate==0]=NA# make it so unchosen covariates appear emp
R> library(ggstance)# for vertical dodge
R> ggplot(estimates, aes(x = Estimate, y = Coefs, color = Model)) +
R>   geom_segment(aes(x = 0, y = Coefs, xend = Estimate, yend = Coefs),size=0.5 ) +
R>   geom_point(position=ggstance::position_dodgev(height=0.2))
R>
```



## 8. Case study 4—Stanford Heart Transplant Data

Although the previous case studies provide a broad overview of the capabilities of **casebase**, they all have two properties in common:

- All covariates are fixed, i.e. they do not change over time;
- The estimated hazard functions satisfy the "proportional hazard" assumption.

Case-base sampling has been used in the literature to study vaccination safety, where the exposure period was defined as the week following vaccination (Saarela and Hanley 2015). Hence, the main covariate of interest, i.e. exposure to the vaccine, was changing over time. In this context, case-base sampling offers an efficient alternative to nested case-control designs or self-matching.

The purpose of the next case study is to show how **casebase** can also be used for survival analysis problems that do not share the properties above. To this end, we will use the Stanford Heart Transplant Data (Clark, Stinson, Griep, Schroeder, Shumway, and Harrison 1971, Crowley and Hu (1977)); we will use the version available in the **survival** package.

Recall the setting of this study: patients were admitted to the Stanford program after meeting with their physician and determining that they were unlikely to respond to other forms of treatment. After enrollment, the program searches for a suitable donor for the patient, which can take anywhere between a few days to almost a year. We are interested in the effect of a heart transplant on survival; therefore, the patient is considered exposed only after the transplant has occurred.

As before, we can look at the population-time plot for a graphical summary of the event incidence. As we can see, most events occur early during the follow-up period, and therefore we do not expect the hazard to be constant.

```
R> library(survival)
R> library(casebase)
R>
R> stanford_popTime <- popTime(jasa, time = "fuptime",
R>                             event = "fustat")
R> plot(stanford_popTime)
```

Since the exposure is time-dependent, we need to manually define the exposure variable *after* case-base sampling and *before* fitting the hazard function. For this reason, we will use the **sampleCaseBase** function directly.

```
R> library(tidyverse)
R> library(lubridate)
R>
R> cb_data <- sampleCaseBase(jasa, time = "fuptime",
R>                           event = "fustat", ratio = 10)
```

Next, we will compute the number of days from acceptance into the program to transplant, and we use this variable to determine whether each population-moment is exposed or not.

```

R> # Define exposure variable
R> cb_data <- mutate(cb_data,
R>                   txtime = time_length(accept.dt %--% tx.date,
R>                                       unit = "days"),
R>                   exposure = case_when(
R>                     is.na(txtime) ~ 0L,
R>                     txtime > futime ~ 0L,
R>                     txtime <= futime ~ 1L
R>                   ))

```

Finally, we can fit the hazard using various linear predictors.

```

R> library(splines)
R> # Fit several models
R> fit1 <- fitSmoothHazard(fustat ~ exposure,
R>                        data = cb_data, time = "futile")
R> fit2 <- fitSmoothHazard(fustat ~ exposure + futime,
R>                        data = cb_data, time = "futile")
R> fit3 <- fitSmoothHazard(fustat ~ exposure + bs(futile),
R>                        data = cb_data, time = "futile")
R> fit4 <- fitSmoothHazard(fustat ~ exposure*bs(futile),
R>                        data = cb_data, time = "futile")

```

Note that the fourth model (i.e. `fit4`) includes an interaction term between exposure and follow-up time. In other words, this model no longer exhibit proportional hazards. The evidence of non-proportionality of hazards in the Stanford Heart Transplant data has been widely discussed ([Arjas 1988](#)).

We can then compare the goodness of fit of these four models using the Akaike Information Criterion (AIC).

```

R> # Compute AIC
R> c("Model1" = AIC(fit1),
R>   "Model2" = AIC(fit2),
R>   "Model3" = AIC(fit3),
R>   "Model4" = AIC(fit4))

```

As we can, the best fit is the fourth model. By visualizing the hazard functions for both exposed and unexposed individuals, we can more clearly see how the hazards are no longer proportional.

```

R> # Compute hazards---
R> # First, create a list of time points for both exposure status
R> hazard_data <- expand.grid(exposure = c(0, 1),
R>                            futime = seq(0, 1000,
R>                                       length.out = 100))
R> # Set the offset to zero

```

```

R> hazard_data$offset <- 0
R> # Use predict to get the fitted values, and exponentiate to
R> # transform to the right scale
R> hazard_data$hazard = exp(predict(fit4, newdata = hazard_data,
R>                                type = "link"))
R> # Add labels for plots
R> hazard_data$Status = factor(hazard_data$exposure,
R>                             labels = c("NoTrans", "Trans"))
R>
R> ggplot(hazard_data, aes(futime, hazard, colour = Status)) +
R>   geom_line() +
R>   theme_minimal() +
R>   theme(legend.position = 'top') +
R>   ylab('Hazard') + xlab('Follow-up time')

```

The non-proportionality seems to be more pronounced at the beginning of follow-up than the end. Finally, we can turn these estimate of the hazard function into estimates of the cumulative incidence functions.

```

R> # Compute absolute risk curves
R> newdata <- data.frame(exposure = c(0, 1))
R> absrisk <- absoluteRisk(fit4, newdata = newdata,
R>                        time = seq(0, 1000, length.out = 100))
R>
R> colnames(absrisk) <- c("Time", "NoTrans", "Trans")
R>
R> # Rearrange the data
R> absrisk <- gather(as.data.frame(absrisk),
R>                  "Status", "Risk", -Time)
R>
R> ggplot(absrisk, aes(Time, Risk, colour = Status)) +
R>   geom_line() +
R>   theme_minimal() +
R>   theme(legend.position = 'top') +
R>   expand_limits(y = 0.5) +
R>   xlab('Follow-up time') + ylab('Cum. Incidence')

```

Note that we can easily adapt the code above to the situation where a patient receives a heart transplant at a point in time of interest, for example after 30 days.

```

R> # Compute hazards---
R> # First, create a list of time points for both exposure status
R> one_yr_haz <- data.frame(futime = seq(0, 365,
R>                                     length.out = 100))
R> one_yr_haz <- mutate(one_yr_haz,
R>                      offset = 0,
R>                      exposure = if_else(futime < 30, 0, 1))

```

```
R> one_yr_haz$hazard = exp(predict(fit4, newdata = one_yr_haz,
R>                               type = "link"))
R> ggplot(one_yr_haz, aes(futime, hazard)) +
R>   geom_line() +
R>   theme_minimal() +
R>   theme(legend.position = 'top') +
R>   ylab('Hazard') + xlab('Follow-up time') +
R>   geom_vline(xintercept = 30, linetype = 'dashed')
```

We can then compare the 1-year mortality risk without transplant and with transplant at 30 days.

```
R> absoluteRisk(fit4, newdata = data.frame(exposure = 0),
R>             time = 365)
R>
R> # Use the trapezoidal rule to estimate the cumulative hazard
R> min_hazard <- min(one_yr_haz$hazard)
R> max_hazard <- max(one_yr_haz$hazard)
R> increment <- 365/99
R> cumhaz_est <- 0.5*increment*(sum(2*one_yr_haz$hazard) - min_hazard - max_hazard)
R> 1 - exp(-cumhaz_est)
```

As we can see, the risk estimate at 1-year is about 30% lower if the patient receives a heart transplant at 30 days.

## 9. Discussion

In this article, we presented the R package **casebase** that provides functions to fit smooth parametric hazards and estimate cumulative incidence functions using case-base sampling. We outlined the theoretical underpinnings of the approach, we provided details about our implementation, and we illustrated the merits of the approach and the package through three case studies.

**SmoothHazard**, **survival** and **Rstpm2** can handle interval censoring, whereas **casebase** currently cannot. Interval censoring is a simple feature to implement and is a potential improvement to add to **casebase**.

## 10. Environment Details

This report was generated on 2020-02-06 21:26:01 using the following computational environment and dependencies:

```
R> # which R packages and versions?
R> devtools::session_info()
```

```
#> - Session info -----
#> setting value
```

```

#> version R version 3.6.0 (2019-04-26)
#> os      macOS High Sierra 10.13.6
#> system  x86_64, darwin15.6.0
#> ui      X11
#> language (EN)
#> collate en_CA.UTF-8
#> ctype   en_CA.UTF-8
#> tz      America/Montreal
#> date    2020-02-06
#>
#> - Packages -----
#> package      * version      date      lib
#> assertthat    0.2.1        2019-03-21 [1]
#> backports     1.1.5        2019-10-02 [1]
#> broom         0.5.2        2019-04-07 [1]
#> callr         3.4.1        2020-01-24 [1]
#> casebase      * 0.2.1.9001 2020-01-30 [1]
#> cellranger    1.1.0        2016-07-27 [1]
#> cli           2.0.1        2020-01-08 [1]
#> codetools     0.2-16       2018-12-24 [1]
#> colorspace    1.4-1        2019-03-18 [1]
#> crayon        1.3.4        2017-09-16 [1]
#> data.table    * 1.12.8     2019-12-09 [1]
#> DBI           1.0.0        2018-05-02 [1]
#> dbplyr        1.4.2        2019-06-17 [1]
#> desc          1.2.0        2018-05-01 [1]
#> devtools      2.2.1        2019-09-24 [1]
#> digest        0.6.23       2019-11-23 [1]
#> dplyr         * 0.8.3       2019-07-04 [1]
#> ellipsis      0.3.0        2019-09-20 [1]
#> evaluate      0.14         2019-05-28 [1]
#> fansi         0.4.1        2020-01-08 [1]
#> farver        2.0.3        2020-01-16 [1]
#> forcats      * 0.4.0       2019-02-17 [1]
#> foreach       1.4.7        2019-07-27 [1]
#> fs            1.3.1        2019-05-06 [1]
#> generics      0.0.2        2018-11-29 [1]
#> ggplot2      * 3.2.1       2019-08-10 [1]
#> glmnet        * 3.0-2       2019-12-11 [1]
#> glue          1.3.1        2019-03-12 [1]
#> gtable        0.3.0        2019-03-25 [1]
#> haven         2.2.0        2019-11-08 [1]
#> hms           0.5.2        2019-10-30 [1]
#> htmltools     0.4.0        2019-10-04 [1]
#> httr          1.4.1        2019-08-05 [1]
#> iterators     1.0.12       2019-07-26 [1]
#> jsonlite      1.6          2018-12-07 [1]

```

```

#> kableExtra * 1.1.0      2019-03-16 [1]
#> knitr      1.27        2020-01-16 [1]
#> labeling   0.3         2014-08-23 [1]
#> lattice    0.20-38     2018-11-04 [1]
#> lazyeval   0.2.2       2019-03-15 [1]
#> lifecycle  0.1.0       2019-08-01 [1]
#> lubridate  1.7.4       2018-04-11 [1]
#> magrittr   * 1.5       2014-11-22 [1]
#> Matrix     * 1.2-17    2019-03-22 [1]
#> memoise    1.1.0       2017-04-21 [1]
#> mgcv        1.8-31     2019-11-09 [1]
#> modelr     0.1.5       2019-08-08 [1]
#> munsell    0.5.0       2018-06-12 [1]
#> nlme       3.1-142     2019-11-07 [1]
#> pillar     1.4.3       2019-12-20 [1]
#> pkgbuild   1.0.6       2019-10-09 [1]
#> pkgconfig  2.0.3       2019-09-22 [1]
#> pkgload    1.0.2       2018-10-29 [1]
#> prettyunits 1.1.1     2020-01-24 [1]
#> processx   3.4.1       2019-07-18 [1]
#> ps         1.3.0       2018-12-21 [1]
#> purrr      * 0.3.3     2019-10-18 [1]
#> R6         2.4.1       2019-11-12 [1]
#> Rcpp       1.0.3       2019-11-08 [1]
#> readr      * 1.3.1     2018-12-21 [1]
#> readxl     1.3.1       2019-03-13 [1]
#> remotes    2.1.0       2019-06-24 [1]
#> reprex     0.3.0       2019-05-16 [1]
#> rlang      0.4.4       2020-01-28 [1]
#> rmarkdown  2.1         2020-01-20 [1]
#> rprojroot  1.3-2       2018-01-03 [1]
#> rstudioapi 0.10        2019-03-19 [1]
#> rticles    0.12        2019-11-12 [1]
#> rvest      0.3.5       2019-11-08 [1]
#> scales     1.1.0       2019-11-18 [1]
#> sessioninfo 1.1.1     2018-11-05 [1]
#> shape      1.4.4       2018-02-07 [1]
#> stringi    1.4.5       2020-01-11 [1]
#> stringr    * 1.4.0     2019-02-10 [1]
#> survival   * 3.1-7     2019-11-09 [1]
#> testthat   2.3.1       2019-12-01 [1]
#> tibble     * 2.1.3     2019-06-06 [1]
#> tidyr      * 1.0.0     2019-09-11 [1]
#> tidyselect 0.2.5       2018-10-11 [1]
#> tidyverse  * 1.3.0     2019-11-21 [1]
#> usethis    1.5.1       2019-07-04 [1]
#> vctrs      0.2.2       2020-01-24 [1]

```

[illegible]

[illegible]



```
#> CRAN (R 3.6.0)
#> CRAN (R 3.6.0)
#> CRAN (R 3.6.0)
#>
#> [1] /Library/Frameworks/R.framework/Versions/3.6/Resources/library
```

The current Git commit details are:

```
R> # what commit is this file at?
R> git2r::repository(here::here())

#> Local:      master /Users/jislam/Desktop/Gits/cbpaper
#> Remote:     master @ origin (https://github.com/Jesse-Islam/cbpaper.git)
#> Head:       [c26f96c] 2020-02-07: code clean up to fix conflits
```

## References

- Allignol A, Latouche A (2019). “CRAN Task View: Survival Analysis.” URL <https://cran.r-project.org/web/views/Survival.html>.
- Arjas E (1988). “A graphical method for assessing goodness of fit in Cox’s proportional hazards model.” *Journal of the American Statistical Association*, **83**(401), 204–212.
- Clark DA, Stinson EB, Griep RB, Schroeder JS, Shumway NE, Harrison D (1971). “Cardiac transplantation in man.” *Annals of Internal Medicine*, **75**(1), 15–21.
- Clements M, Liu XR, Lambert P, Jakobsen LH, Gasparini A, Smyth G, Alken P, Wood S, Ulerich R (2019). “Smooth Survival Models, Including Generalized Survival Models [R package rstpm2 version 1.5.1].” URL <https://cran.r-project.org/web/packages/rstpm2/index.html>.
- Clerc-Urmès I, Grzebyk M, Hédelin G (2017). *flexrsurv: An R package for relative survival analysis*. R package version 1.4.1, URL <https://CRAN.R-project.org/package=flexrsurv>.
- Crowley J, Hu M (1977). “Covariance analysis of heart transplant survival data.” *Journal of the American Statistical Association*, **72**(357), 27–36.
- Fine JP, Gray RJ (1999). “A proportional hazards model for the subdistribution of a competing risk.” *Journal of the American statistical association*, **94**(446), 496–509.
- Friedman J, Hastie T, Tibshirani R (2010). “Regularization Paths for Generalized Linear Models via Coordinate Descent.” *Journal of Statistical Software*, **33**(1). ISSN 1548-7660. doi:10.18637/jss.v033.i01.
- Fu Z (2015). “Package crrp.” URL <https://cran.r-project.org/web/packages/crrp/index.html>.

- Gerds TA, Blanche P, Mortersen R, Tollenaar N, Mogensen UB, Ozenne B (2019). “Risk Regression Models and Prediction Scores for Survival Analysis with Competing Risks [R package riskRegression version 2019.11.03].” URL <https://CRAN.R-project.org/package=riskRegression>.
- Goeman J, Meijer R, Chaturvedi N, Lueder M (2019). “Package penalized.” URL <https://cran.r-project.org/web/packages/penalized/index.html>.
- Goeman JJ (2010). “L1 penalized estimation in the Cox proportional hazards model.” *Biometrical Journal*, (52), –14.
- Gray B (2019). “Package cmprsk.” URL <https://cran.r-project.org/web/packages/cmprsk/index.html>.
- Hanley JA, Miettinen OS (2009). “Fitting smooth-in-time prognostic risk functions via logistic regression.” *The International Journal of Biostatistics*, **5**(1).
- Hastie T, Tibshirani R (1987). “Generalized additive models: some applications.” *Journal of the American Statistical Association*, **82**(398), 371–386.
- Jackson C (2016). “flexsurv: A Platform for Parametric Survival Modeling in R.” *Journal of Statistical Software*, **70**(8), 1–33. doi:10.18637/jss.v070.i08.
- Knaus WA, Harrell FE, Lynn J, Goldman L, Phillips RS, Connors AF, Dawson NV, Fulkerson WJ, Califf RM, Desbiens N, *et al.* (1995). “The SUPPORT prognostic model: objective estimates of survival for seriously ill hospitalized adults.” *Annals of internal medicine*, **122**(3), 191–203.
- Liu Z, Rich B, Hanley JA (2014). “Recovering the raw data behind a non-parametric survival curve.” *Systematic reviews*, **3**(1), 151.
- Mahani AS, Sharabiani MT (2015). *Bayesian, and Non-Bayesian, Cause-Specific Competing-Risk Analysis for Parametric and Nonparametric Survival Functions: The R Package CFC*. Ph.D. thesis, Ph. D. Dissertation.
- Park MY, Hastie T (2018). “Package glmpath.” URL <https://CRAN.R-project.org/package=glmpath>.
- Perperoglou A (2015). “Package CoxRidge.” URL <https://CRAN.R-project.org/package=CoxRidge>.
- Saarela O (2016). “A case-base sampling method for estimating recurrent event intensities.” *Lifetime data analysis*, **22**(4), 589–605.
- Saarela O, Arjas E (2015). “Non-parametric Bayesian Hazard Regression for Chronic Disease Risk Assessment.” *Scandinavian Journal of Statistics*, **42**(2), 609–626.
- Saarela O, Hanley JA (2015). “Case-base methods for studying vaccination safety.” *Biometrics*, **71**(1), 42–52.
- Scheike TH, Holst KK, Hjelmberg JB (2014). “Estimating twin concordance for bivariate competing risks twin data.” *Statistics in medicine*, **33**(7), 1193–1204.

- Scheike TH, Zhang MJ (2011). “Analyzing Competing Risk Data Using the R `timereg` Package.” *Journal of Statistical Software*, **38**(2), 1–15. URL <http://www.jstatsoft.org/v38/i02/>.
- Schröder FH, Hugosson J, Roobol MJ, Tammela TL, Ciatto S, Nelen V, Kwiatkowski M, Lujan M, Lilja H, Zappa M, *et al.* (2009). “Screening and prostate-cancer mortality in a randomized European study.” *New England Journal of Medicine*, **360**(13), 1320–1328.
- Scrucca L, Santucci A, Aversa F (2010). “Regression modeling of competing risk using R: an in depth guide for clinicians.” *Bone marrow transplantation*, **45**(9), 1388.
- Simon N, Friedman J, Hastie T, Tibshirani R (2011). “Regularization Paths for Cox’s Proportional Hazards Model via Coordinate Descent.” *Journal of Statistical Software*, **39**(5), 1–13. URL <http://www.jstatsoft.org/v39/i05/>.
- Therneau TM (2015). *A Package for Survival Analysis in S*. Version 2.38, URL <https://CRAN.R-project.org/package=survival>.
- Touraine C, Gerds TA, Joly P (2017). “SmoothHazard: An R Package for Fitting Regression Models to Interval-Censored Observations of Illness-Death Models.” *Journal of Statistical Software*, **79**(7), 1–22. doi:10.18637/jss.v079.i07.
- van Buuren S, Groothuis-Oudshoorn K (2011). “mice: Multivariate Imputation by Chained Equations in R.” *Journal of Statistical Software*, **45**(3), 1–67. URL <https://www.jstatsoft.org/v45/i03/>.
- Yi Y, Zou H (2017). “Package `fastcox`.” URL <https://cran.r-project.org/web/packages/fastcox/index.html>.

## Affiliation:

Sahir Bhatnagar \*

McGill University

1020 Pine Avenue West Montreal, QC, Canada H3A 1A2

E-mail: [sahir.bhatnagar@mail.mcgill.ca](mailto:sahir.bhatnagar@mail.mcgill.ca)

URL: <http://sahirbhatnagar.com/>

Maxime Turgeon \*

University of Manitoba

186 Dysart Road Winnipeg, MB, Canada R3T 2N2

E-mail: [max.turgeon@umanitoba.ca](mailto:max.turgeon@umanitoba.ca)

URL: <https://maxturgeon.ca/>

Jesse Islam

McGill University

1020 Pine Avenue West Montreal, QC, Canada H3A 1A2

E-mail: [jesse.islam@mail.mcgill.ca](mailto:jesse.islam@mail.mcgill.ca)

James Hanley

McGill University

1020 Pine Avenue West Montreal, QC, Canada H3A 1A2

E-mail: [james.hanley@mcgill.ca](mailto:james.hanley@mcgill.ca)

URL: <http://www.medicine.mcgill.ca/epidemiology/hanley/>

Olli Saarela

University of Toronto

Dalla Lana School of Public Health, 155 College Street, 6th floor, Toronto, Ontario M5T 3M7, Canada

E-mail: [olli.saarela@utoronto.ca](mailto:olli.saarela@utoronto.ca)

URL: <http://individual.utoronto.ca/osaarela/>

Name	Labels	Levels	Storage	NAs
age	Age		double	0
death	Death at any time up to NDI date:31DEC94		double	0
sex		2	integer	0
hospdead	Death in Hospital		double	0
slos	Days from Study Entry to Discharge		double	0
d.time	Days of Follow-Up		double	0
dzgroup		8	integer	0
dzclass		4	integer	0
num.co	number of comorbidities		double	0
edu	Years of Education		double	202
income		4	integer	349
scoma	SUPPORT Coma Score based on Glasgow D3		double	0
charges	Hospital Charges		double	25
totcst	Total RCC cost		double	105
totmcst	Total micro-cost		double	372
avtisst	Average TISS, Days 3-25		double	6
race		5	integer	5
meanbp	Mean Arterial Blood Pressure Day 3		double	0
wblc	White Blood Cell Count Day 3		double	24
hrt	Heart Rate Day 3		double	0
resp	Respiration Rate Day 3		double	0
temp	Temperature (Celsius) Day 3		double	0
pafi	PaO2/(.01*FiO2) Day 3		double	253
alb	Serum Albumin Day 3		double	378
bili	Bilirubin Day 3		double	297
crea	Serum creatinine Day 3		double	3
sod	Serum sodium Day 3		double	0
ph	Serum pH (arterial) Day 3		double	250
glucose	Glucose Day 3		double	470
bun	BUN Day 3		double	455
urine	Urine Output Day 3		double	517
adlp	ADL Patient Day 3		double	634
adls	ADL Surrogate Day 3		double	310
sfdm2		5	integer	159
adlsc	Imputed ADL Calibrated to Surrogate		double	0

Table 3: A description of each variable in the SUPPORT dataset.

Variable	Levels
sex	female
	male
dzgroup	ARF/MOSF w/Sepsis
	COPD
	CHF
	Cirrhosis
	Coma
	Colon Cancer
	Lung Cancer
	MOSF w/Malig
dzclass	ARF/MOSF
	COPD/CHF/Cirrhosis
	Coma
	Cancer
income	under \$11k
	11–25k
	25–50k
	>\$50k
race	white
	black
	asian
	other
	hispanic
sfdm2	no(M2 and SIP pres)
	adl $\geq$ 4 ( $\geq$ 5 if sur)
	SIP $\geq$ 30
	Coma or Intub
	<2 mo. follow-up

Table 4: A description of each level within each categorical variable in the dataset.