

# flexsurv: A Platform for Parametric Survival Modelling in R

Christopher H. Jackson  
MRC Biostatistics Unit, Cambridge, UK

---

## Abstract

**flexsurv** is an R package for fully-parametric modelling of survival data. Any parametric time-to-event distribution may be fitted if the user supplies a probability density or hazard function, and ideally also their cumulative versions. Standard survival distributions are built in, including the three and four-parameter generalized gamma and F distributions. Any parameter of any distribution can be modelled as a linear or log-linear function of covariates. The package also includes the spline model of Royston and Parmar (2002), in which both baseline survival and covariate effects can be arbitrarily flexible parametric functions of time. The main model-fitting function, **flexsurvreg**, uses the familiar syntax of **survreg** from the standard **survival** package (Therneau 2014). Censoring or left-truncation are specified in **Surv** objects. The models are fitted by maximising the full log-likelihood, and estimates and confidence intervals for any function of the model parameters can be printed or plotted. **flexsurv** also provides functions for fitting and predicting from fully-parametric multi-state models, and connects with the **mstate** package (de Wreede *et al.* 2011). This article explains the methods and design principles of the package, giving several worked examples of its use. *[Note: A version of this vignette is published as Jackson (2016) in Journal of Statistical Software. All content there is included here. There have been no substantial changes in the survival modelling parts since then. Version 2.0 of flexsurv added new features for multi-state modelling, and since that version, multi-state modelling with flexsurv has been described in a separate vignette.]*

*Keywords:* survival, multi-state models, multistate models.

---

## 1. Motivation and design

The Cox model for survival data is ubiquitous in medical research, since the effects of predictors can be estimated without needing to supply a baseline survival distribution that might be inaccurate. However, fully-parametric models have many advantages, and even the originator of the Cox model has expressed a preference for parametric modelling (see Reid 1994). Fully-specified models can be more convenient for representing complex data structures and processes (Aalen *et al.* 2008), e.g. hazards that vary predictably, interval censoring, frailties, multiple responses, datasets or time scales, and can help with out-of-sample prediction. For example, the mean survival  $E(T) = \int_0^\infty S(t)dt$ , used in health economic evaluations (Latimer 2013), needs the survivor function  $S(t)$  to be fully-specified for all times  $t$ , and parametric models that combine data from multiple time periods can facilitate this (Benaglia *et al.* 2014). **flexsurv** for R (R Core Team 2014) allows parametric distributions of arbitrary complexity to be fitted to survival data, gaining the convenience of parametric modelling, while avoiding

the risk of model misspecification. Built-in choices include spline-based models with any number of knots (Royston and Parmar 2002) and 3–4 parameter generalized gamma and F distribution families. Any user-defined model may be employed by supplying at minimum an R function to compute the probability density or hazard, and ideally also its cumulative form. Any parameters may be modelled in terms of covariates, and any function of the parameters may be printed or plotted in model summaries.

**flexsurv** is intended as a general platform for survival modelling in R. The **survreg** function in the R package **survival** (Therneau 2014) only supports two-parameter (location/scale) distributions, though users can supply their own distributions if they can be parameterised in this form. Some other contributed R packages can fit survival models, e.g., **eha** (Broström 2014) and **VGAM** (Yee and Wild 1996), though these are either limited to specific distribution families, or not specifically designed for survival analysis. Others, e.g. **ActuDistns** (Nadarajah and Bakar 2013), contain only the definitions of distribution functions. **flexsurv** enables such functions to be used in survival models.

It is similar in spirit to the Stata packages **stpm2** (Lambert and Royston 2009) for spline-based survival modelling, and **stgenreg** (Crowther and Lambert 2013) for fitting survival models with user-defined hazard functions using numerical integration. Though in **flexsurv**, slow numerical integration can be avoided if the analytic cumulative distribution or hazard can be supplied, and optimisation can also be speeded by supplying analytic derivatives. **flexsurv** also has features for multi-state modelling and interval censoring, and general output reporting. It employs functional programming to work with user-defined or existing R functions.

§2 explains the general model that **flexsurv** is based on. §3 gives examples of its use for fitting built-in survival distributions with a fixed number of parameters, and §4 explains how users can define new distributions. §5 concentrates on classes of models where the number of parameters can be chosen arbitrarily, such as splines. §6 mentions the use of **flexsurv** for fitting and predicting from fully-parametric multi-state models, which is described more fully in a separate vignette. Finally §7 suggests some potential future extensions.

## 2. General parametric survival model

The general model that **flexsurv** fits has probability density for death at time  $t$ :

$$f(t|\mu(\mathbf{z}), \boldsymbol{\alpha}(\mathbf{z})), \quad t \geq 0 \quad (1)$$

The cumulative distribution function  $F(t)$ , survivor function  $S(t) = 1 - F(t)$ , cumulative hazard  $H(t) = -\log S(t)$  and hazard  $h(t) = f(t)/S(t)$  are also defined (suppressing the conditioning for clarity).  $\mu = \alpha_0$  is the parameter of primary interest, which usually governs the mean or *location* of the distribution. Other parameters  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_R)$  are called “ancillary” and determine the shape, variance or higher moments.

**Covariates** All parameters may depend on a vector of covariates  $\mathbf{z}$  through link-transformed linear models  $g_0(\mu(\mathbf{z})) = \gamma_0 + \boldsymbol{\beta}_0^\top \mathbf{z}$  and  $g_r(\alpha_r(\mathbf{z})) = \gamma_r + \boldsymbol{\beta}_r^\top \mathbf{z}$ .  $g()$  will typically be  $\log()$  if the parameter is defined to be positive, or the identity function if the parameter is unrestricted. Suppose that the location parameter, but not the ancillary parameters, depends on covariates. If the hazard function factorises as  $h(t|\boldsymbol{\alpha}, \mu(\mathbf{z})) = \mu(\mathbf{z})h_0(t|\boldsymbol{\alpha})$ , then this is a *proportional*

*hazards* (PH) model, so that the hazard ratio between two groups (defined by two different values of  $\mathbf{z}$ ) is constant over time  $t$ .

Alternatively, if  $S(t|\mu(\mathbf{z}), \boldsymbol{\alpha}) = S_0(\mu(\mathbf{z})t|\boldsymbol{\alpha})$  then it is an *accelerated failure time* (AFT) model, so that the effect of covariates is to speed or slow the passage of time. For example, doubling the value of a covariate with coefficient  $\beta = \log(2)$  would give half the expected survival time.

**Data and likelihood** Let  $t_i : i = 1, \dots, n$  be a sample of times from individuals  $i$ . Let  $c_i = 1$  if  $t_i$  is an observed death time, or  $c_i = 0$  if this is censored. Most commonly,  $t_i$  may be right-censored, thus the true death time is known only to be greater than  $t_i$ . More generally, the survival time may be interval-censored on  $(t_i^{\min}, t_i^{\max})$ .

Also let  $s_i$  be corresponding left-truncation (or delayed-entry) times, meaning that the  $i$ th survival time is only observed conditionally on the individual having survived up to  $s_i$ , thus  $s_i = 0$  if there is no left-truncation. Time-dependent covariates (§3.1) and some multi-state models (§6) can be represented through left-truncation.

With at most right-censoring, the likelihood for the parameters  $\boldsymbol{\theta} = \{\gamma, \beta\}$  in Equation 1, given the corresponding data vectors, is

$$l(\boldsymbol{\theta}|\mathbf{t}, \mathbf{c}, \mathbf{s}) = \left\{ \prod_{i: c_i=1} f_i(t_i) \prod_{i: c_i=0} S_i(t_i) \right\} / \prod_i S_i(s_i) \quad (2)$$

where  $f_i(t_i)$  is shorthand for  $f(t_i|\mu(\mathbf{z}_i), \boldsymbol{\alpha}(\mathbf{z}_i))$ ,  $S_i(t_i)$  is  $S(t_i|\mu(\mathbf{z}_i), \boldsymbol{\alpha}(\mathbf{z}_i))$ , and  $\mu, \boldsymbol{\alpha}$  are related to  $\gamma, \beta$  and  $\mathbf{z}_i$  via the link functions defined above. The log-likelihood also has a concise form in terms of hazards and cumulative hazards, as

$$\log l(\boldsymbol{\theta}|\mathbf{t}, \mathbf{c}, \mathbf{s}) = \sum_{i: c_i=1} \{\log(h_i(t_i)) - H_i(t_i)\} - \sum_{i: c_i=0} H_i(t_i) + \sum_i H_i(s_i)$$

With interval-censoring, the likelihood is

$$l(\boldsymbol{\theta}|\mathbf{t}^{\min}, \mathbf{t}^{\max}, \mathbf{c}, \mathbf{s}) = \left\{ \prod_{i: c_i=1} f_i(t_i) \prod_{i: c_i=0} (S_i(t_i^{\min}) - S_i(t_i^{\max})) \right\} / \prod_i S_i(s_i) \quad (3)$$

These likelihoods assume that the times of censoring are fixed or otherwise distributed independently of the parameters  $\boldsymbol{\theta}$  that govern the survival times (see, e.g. Aalen *et al.* (2008)). The individual survival times are also independent, so that **flexsurv** does not currently support shared frailty, clustered or random effects models (see §7).

The parameters are estimated by maximising the full log-likelihood with respect to  $\boldsymbol{\theta}$ , as detailed further in §3.6.

### 3. Fitting standard parametric survival models

An example dataset used throughout this paper is from 686 patients with primary node positive breast cancer, available in the package as **bc**. This was originally provided with **stpm** (Royston 2001), and analysed in much more detail by Sauerbrei and Royston (1999) and Royston and Parmar (2002)<sup>1</sup>. The first two records are shown by:

<sup>1</sup>A version of this dataset, including more covariates but excluding the prognostic group, is also provided as **GBSG2** in the package **TH.data** (Hothorn 2015).

```
R> library("flexsurv")
```

```
R> head(bc, 2)
```

```
  censrec rectime group  recyrs
1      0    1342  Good 3.676712
2      0    1578  Good 4.323288
```

The main model-fitting function is called `flexsurvreg`. Its first argument is an R *formula* object. The left hand side of the formula gives the response as a survival object, using the `Surv` function from the `survival` package.

```
R> fs1 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc,
+                    dist = "weibull")
```

Here, this indicates that the response variable is `recyrs`. This represents the time (in years) of death or cancer recurrence when `censrec` is 1, or (right-)censoring when `censrec` is 0. The covariate `group` is a factor representing a prognostic score, with three levels "Good" (the baseline), "Medium" and "Poor". All of these variables are in the data frame `bc`. If the argument `dist` is a string, this denotes a built-in survival distribution. In this case we fit a Weibull survival model.

Printing the fitted model object gives estimates and confidence intervals for the model parameters and other useful information. Note that these are the *same parameters* as represented by the R distribution function `dweibull`: the *shape*  $\alpha$  and the *scale*  $\mu$  of the survivor function  $S(t) = \exp(-(t/\mu)^\alpha)$ , and `group` has a linear effect on  $\log(\mu)$ .

```
R> fs1
```

Call:

```
flexsurvreg(formula = Surv(recyrs, censrec) ~ group, data = bc,
            dist = "weibull")
```

Estimates:

	data	mean	est	L95%	U95%	se
shape		NA	1.3797	1.2548	1.5170	0.0668
scale		NA	11.4229	9.1818	14.2110	1.2728
groupMedium	0.3338		-0.6136	-0.8623	-0.3649	0.1269
groupPoor	0.3324		-1.2122	-1.4583	-0.9661	0.1256
	exp(est)		L95%	U95%		
shape		NA	NA	NA		
scale		NA	NA	NA		
groupMedium	0.5414		0.4222	0.6943		
groupPoor	0.2975		0.2326	0.3806		

N = 686, Events: 299, Censored: 387

Total time at risk: 2113.425

Log-likelihood = -811.9419, df = 4

AIC = 1631.884

For the Weibull (and exponential, log-normal and log-logistic) distribution, `flexsurvreg` simply acts as a wrapper for `survreg`: the maximum likelihood estimates are obtained by `survreg`, checked by `flexsurvreg` for optimisation convergence, and converted to `flexsurvreg`'s preferred parameterisation. Therefore the same model can be fitted more directly as

```
R> survreg(Surv(recyrs, censrec) ~ group, data = bc, dist = "weibull")
```

Call:

```
survreg(formula = Surv(recyrs, censrec) ~ group, data = bc, dist = "weibull")
```

Coefficients:

```
(Intercept) groupMedium    groupPoor
      2.4356168    -0.6135892    -1.2122137
```

Scale= 0.7248206

```
Loglik(model)= -811.9    Loglik(intercept only)= -873.2
      Chisq= 122.53 on 2 degrees of freedom, p= <2e-16
n= 686
```

The maximised log-likelihoods are the same, however the parameterisation is different: the first coefficient (`Intercept`) reported by `survreg` is  $\log(\mu)$ , and `survreg`'s `"scale"` is `dweibull`'s (thus `flexsurvreg`)'s  $1 / \text{shape}$ . The covariate effects  $\beta$ , however, have the same “accelerated failure time” interpretation, as linear effects on  $\log(\mu)$ . The multiplicative effects  $\exp(\beta)$  are printed in the output as `exp(est)`.

The same model can be fitted in `eha`, also by maximum likelihood, as

```
R> library(eha)
R> aftreg(Surv(recyrs, censrec) ~ group, data = bc, dist = "weibull")
```

The results are presented in the same parameterisation as `flexsurvreg`, except that the shape and scale parameters are log-transformed, and (unless the argument `param="lifeExp"` is supplied) the covariate effects have the opposite sign.

### 3.1. Additional modelling features

#### *Truncation and time-dependent covariates*

If we also had left-truncation times in a variable called `start`, the response would be `Surv(start, recyrs, censrec)`. Or if all responses were interval-censored between lower and upper bounds `tmin` and `tmax`, then we would write `Surv(tmin, tmax, type = "interval2")`.

Time-dependent covariates can sometimes be represented in “counting process” form — as a series of left-truncated survival times, which may also be right-censored. For each individual there would be multiple records, each corresponding to an interval where the covariate is assumed to be constant. The response would be of the form `Surv(start, stop, censrec)`, where `start` and `stop` are the limits of each interval, and `censrec` indicates whether a

death was observed at `stop`. Care is required however. Whether this is a valid approach in **flexsurv** depends on whether the probability of survival up to the left-truncation time can be represented by a term of the form  $S_i(s_i) = \exp(-H_i(s_i))$  in the likelihood (equation 2). The cumulative hazard  $H$  over the interval from time 0 to time  $s_i$  depends on how the covariates change on this time interval. If the covariates are constant between time 0 and time  $s_i$ , or if covariates are modelled with proportional hazards, then this cumulative hazard is  $H_i(s_i)$ , hence the likelihood used by **flexsurv** is valid. It is not valid in general however for other forms of dependence on covariates, e.g. accelerated failure time models.

In versions of **flexsurv** since April 2020, models with individual-specific right-truncation times are also supported. These are used for situations with “retrospective ascertainment”, where cases are only included in the data if they have died by a specific time. These models are specified through an argument `rtrunc` to `flexsurvreg` that names the variable with the truncation times. See the Supplementary Examples vignette for a worked example.

### *Relative survival*

In relative survival models (Nelson *et al.* 2007), the survivor function is expressed as  $S(t) = S^*(t)R(t)$ , where  $S^*(t)$  is the “expected” or “baseline” survival, and  $R(t)$  is the *relative* survival. Equivalently, the hazard is defined as  $h(t) = h^*(t) + \lambda(t)$ , where  $h^*(t)$  is the baseline hazard function, and  $\lambda(t)$  is the excess mortality rate associated with the disease of interest. The baseline represents a reference population, and is typically obtained from national routinely-collected mortality statistics, adjusted (e.g. by age/sex) to represent the population under study. The parametric model is defined and estimated for  $R(t)$ .

These models are implemented in **flexsurv** by supplying the variable in the data that represents the expected mortality rate  $h^*(t)$  in the `bhazard` argument to `flexsurvreg`. This is only used for the individuals in the data who die, and `bhazard` describes the expected hazard at the death time. The values of `bhazard` for censored individuals are ignored.

Note that the parameters returned in the model fitted by `flexsurvreg` refer to the relative survival  $R(t)$ , rather than the absolute survival. The likelihood returned by `flexsurvreg` here is a *partial* likelihood defined (as in Nelson *et al.* 2007, equations 4–5) by omitting the term  $\sum_i \log(S^*(t_i))$  (summed over all individuals  $i$  in the data, including both censored and uncensored times  $t_i$ ) from the full likelihood. This term is equivalent to minus the sum of the cumulative hazards. It can be omitted from the likelihood for the purpose of estimating the parameters of the relative survival model, since it does not depend on these parameters. Hence if a full likelihood is required, (e.g. for model comparison) then this term should be added to the partial likelihood.

Similarly, the predicted survival or hazard (e.g. as returned by `summary.flexsurvreg`, see Section 3.4) from a relative survival model refers to  $R(t)$  or  $h(t)$ . Hence if the overall survival or hazard is required, the predictions of relative survival should be converted to the “absolute” scale by combining with the baseline, though no specific tools for doing this are provided by **flexsurv**.

### *Weighting and subsetting*

Case weights and data subsets can also be specified, as in standard R modelling functions, using `weights` or `subset` arguments.

### 3.2. Built-in models

**flexsurvreg**'s currently built-in distributions are listed in Table 1. In each case, the probability density  $f()$  and parameters of the fitted model are taken from an existing R function of the same name but beginning with the letter **d**. For the Weibull, exponential (**dexp**), gamma (**dgamma**) and log-normal (**dlnorm**), the density functions are provided with standard installations of R. These density functions, and the corresponding cumulative distribution functions (with first letter **p** instead of **d**) are used internally in **flexsurvreg** to compute the likelihood.

**flexsurv** provides some additional survival distributions, including a Gompertz distribution with unrestricted shape parameter, Weibull with proportional hazards parameterisation, log-logistic, and the three- and four-parameter families described below. For all built-in distributions, **flexsurv** also defines functions beginning with **h** giving the hazard, and **H** for the cumulative hazard.

A package vignette “Distributions reference” lists the survivor function and parameterisation of covariate effects used by each built-in distribution.

**Generalized gamma** This three-parameter distribution includes the Weibull, gamma and log-normal as special cases. The original parameterisation from Stacy (1962) is available as `dist = "gengamma.orig"`, however the newer parameterisation (Prentice 1974) is preferred: `dist = "gengamma"`. This has parameters  $(\mu, \sigma, q)$ , and survivor function

$$\begin{aligned} 1 - I(\gamma, u) & \quad (q > 0) \\ 1 - \Phi(z) & \quad (q = 0) \end{aligned}$$

where  $I(\gamma, u) = \int_0^u x^{\gamma-1} \exp(-x) / \Gamma(\gamma)$  is the incomplete gamma function (the cumulative gamma distribution with shape  $\gamma$  and scale 1),  $\Phi$  is the standard normal cumulative distribution,  $u = \gamma \exp(|q|z)$ ,  $z = (\log(t) - \mu) / \sigma$ , and  $\gamma = q^{-2}$ . The Prentice (1974) parameterisation extends the original one to include a further class of models with negative  $q$ , and survivor function  $I(\gamma, u)$ , where  $z$  is replaced by  $-z$ . This stabilises estimation when the distribution is close to log-normal, since  $q = 0$  is no longer near the boundary of the parameter space. In R notation,<sup>2</sup> the parameter values corresponding to the three special cases are

```
dgengamma(x, mu, sigma, Q=0)      == dlnorm(x, mu, sigma)
dgengamma(x, mu, sigma, Q=1)      == dweibull(x, shape = 1 / sigma,
                                                scale = exp(mu))
dgengamma(x, mu, sigma, Q=sigma) == dgamma(x, shape = 1 / sigma^2,
                                              rate = exp(-mu) / sigma^2)
```

**Generalized F** This four-parameter distribution includes the generalized gamma, and also the log-logistic, as special cases. The variety of hazard shapes that can be represented is discussed by Cox (2008). It is provided here in alternative “original” (`dist = "genf.orig"`) and “stable” parameterisations (`dist = "genf"`) as presented by Prentice (1975). See `help(GenF)` and `help(GenF.orig)` in the package documentation for the exact definitions.

<sup>2</sup>The parameter called  $q$  here and in previous literature is called  $Q$  in **dgengamma** and related functions, since the first argument of a cumulative distribution function is conventionally named  $q$ , for quantile, in R.

	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>	dcompertz	"gompertz"
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 1: Built-in parametric survival distributions in **flexsurv**.

### 3.3. Covariates on ancillary parameters

The generalized gamma model is fitted to the breast cancer survival data. **fs2** is an AFT model, where only the parameter  $\mu$  depends on the prognostic covariate **group**. In a second model **fs3**, the first ancillary parameter **sigma** ( $\alpha_1$ ) also depends on this covariate, giving a model with a time-dependent effect that is neither PH nor AFT. The second ancillary parameter **Q** is still common between prognostic groups.

```
R> fs2 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc,
+                    dist = "gengamma")
R> fs3 <- flexsurvreg(Surv(recyrs, censrec) ~ group + sigma(group), data = bc,
+                    dist = "gengamma")
```

Ancillary covariates can alternatively be supplied using the **anc** argument to **flexsurvreg**. This syntax is required if any parameter names clash with the names of functions used in model formulae (e.g., **factor()** or **I()**).

```
R> fs3 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc,
+                    anc = list(sigma = ~ group), dist = "gengamma")
```

Table 3 compares all the models fitted to the breast cancer data, showing absolute fit to the data as measured by the maximised  $-2 \times \log$  likelihood  $-2LL$ , number of parameters  $p$ , and Akaike's information criterion  $-2LL + 2p$  (AIC). The model **fs2** has the lowest AIC, indicating the best estimated predictive ability.

### 3.4. Plotting outputs

The **plot()** method for **flexsurvreg** objects is used as a quick check of model fit. By default, this draws a Kaplan-Meier estimate of the survivor function  $S(t)$ , one for each combination of



categorical covariates, or just a single “population average” curve if there are no categorical covariates (Figure 1). The corresponding estimates from the fitted model are overlaid. Fitted values from further models can be added with the `lines()` method.

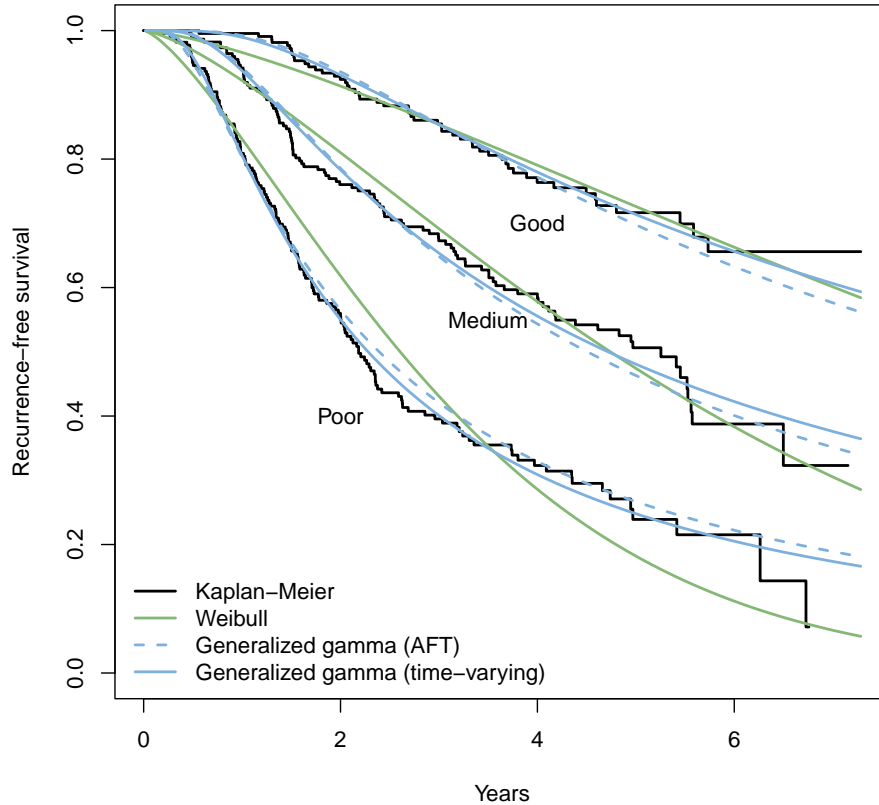


Figure 1: Survival by prognostic group from the breast cancer data: fitted from alternative parametric models and Kaplan-Meier estimates.

The argument `type = "hazard"` can be set to plot hazards from parametric models against kernel density estimates obtained from `muhaz` (Hess 2010; Mueller and Wang 1994). Figure 2 shows more clearly that the Weibull model is inadequate for the breast cancer data: the hazard must be increasing or decreasing — while the generalized gamma can represent the increase and subsequent decline in hazard seen in the data. Similarly, `type = "cumhaz"` plots cumulative hazards.

The numbers plotted are available from the `summary.flexsurvreg()` method. Confidence intervals are produced by simulating a large sample from the asymptotic normal distribution of the maximum likelihood estimates of  $\{\beta_r : r = 0, \dots, R\}$  (Mandel 2013), via the function `normboot.flexsurvreg`. This very general method allows confidence intervals to be obtained for arbitrary functions of the parameters, as described in the next section.

In this example, there is only a single categorical covariate, and the `plot` and `summary` methods return one observed and fitted trajectory for each level of that covariate. For more complicated

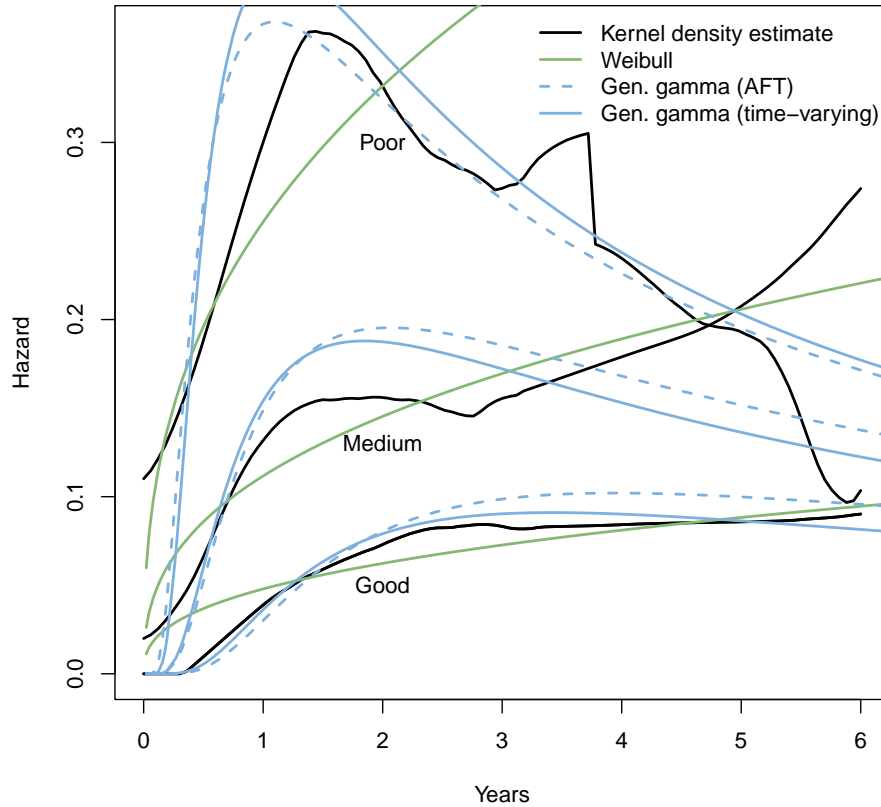


Figure 2: Hazards by prognostic group from the breast cancer data: fitted from alternative parametric models and kernel density estimates.

models, users should specify what covariate values they want summaries for, rather than relying on the default <sup>3</sup>. This is done by supplying the `newdata` argument, a data frame or list containing covariate values, just as in standard R functions like `predict.lm`. Time-dependent covariates are not understood by these functions.

This `plot()` method is only for casual exploratory use. For publication-standard figures, it is preferable to set up the axes beforehand (`plot(..., type = "n")`), and use the `lines()` methods for `flexsurvreg` objects, or construct plots by hand using the data available from `summary.flexsurvreg()`.

### 3.5. Custom model summaries

Any function of the parameters of a fitted model can be summarised or plotted by supplying

<sup>3</sup>If there are only factor covariates, all combinations are plotted. If there are any continuous covariates, these methods by default return a “population average” curve, with the linear model design matrix set to its average values, including the 0/1 contrasts defining factors, which doesn’t represent any specific covariate combination.

the argument `fn` to `summary.flexsurvreg` or `plot.flexsurvreg`. This should be an R function, with optional first two arguments `t` representing time, and `start` representing a left-truncation point (if the result is conditional on survival up to that time). Any remaining arguments must be the parameters of the survival distribution. For example, median survival under the Weibull model `fs1` can be summarised as follows

```
R> median.weibull <- function(shape, scale) {
+   qweibull(0.5, shape = shape, scale = scale)
+ }
R> summary(fs1, fn = median.weibull, t = 1, B = 10000)
```

```
group=Good
  time      est      lcl      ucl
1    1 8.75794 7.091625 10.74015
```

```
group=Medium
  time      est      lcl      ucl
1    1 4.741585 4.119985 5.438689
```

```
group=Poor
  time      est      lcl      ucl
1    1 2.605819 2.308052 2.943888
```

Although the median of the Weibull has an analytic form as  $\mu \log(2)^{1/\alpha}$ , the form of the code given here generalises to other distributions. The argument `t` (or `start`) can be omitted from `median.weibull`, because the median is a time-constant function of the parameters, unlike the survival or hazard.

10000 random samples are drawn to produce a slightly more precise confidence interval than the default — users should adjust this until the desired level of precision is obtained. A useful future extension of the package would be to employ user-supplied (or built-in) derivatives of summary functions if possible, so that the delta method can be used to obtain approximate confidence intervals without simulation.

### 3.6. Computation

The likelihood is maximised in `flexsurvreg` using the optimisation methods available through the standard R `optim` function. By default, this is the "BFGS" method (Nash 1990) using the analytic derivatives of the likelihood with respect to the model parameters, if these are available, to improve the speed of convergence to the maximum. These derivatives are built-in for the exponential, Weibull, Gompertz, log-logistic, and hazard- and odds-based spline models (see §5.1). For custom distributions (see §4), the user can optionally supply functions with names beginning "DLd" and "DLS" respectively (e.g., `DLdweibull`, `DLSweibull`) to calculate the derivatives of the log density and log survivor functions with respect to the transformed baseline parameters  $\gamma$  (then the derivatives with respect to  $\beta$  are obtained automatically). Arguments to `optim` can be passed to `flexsurvreg` — in particular, `control` options, such as convergence tolerance, iteration limit or function or parameter scaling, may need to be adjusted to achieve convergence.

## 4. Custom survival distributions

**flexsurv** is not limited to its built-in distributions. Any survival model of the form (1–3) can be fitted if the user can provide either the density function  $f()$  or the hazard  $h()$ . Many contributed R packages provide probability density and cumulative distribution functions for positive distributions. Though survival models may be more naturally characterised by their hazard function, representing the changing risk of death through time. For example, for survival following major surgery we may want a “U-shaped” hazard curve, representing a high risk soon after the operation, which then decreases, but increases naturally as survivors grow older.

To supply a custom distribution, the **dist** argument to **flexsurvreg** is defined to be an R list object, rather than a character string. The list has the following elements.

**name** Name of the distribution. In the first example below, we use a log-logistic distribution, and the name is “**llogis**”<sup>4</sup>. Then there is assumed to be at least either

- a function to compute the probability density, which would be called **dllogis** here, or
- a function to compute the hazard, called **hllogis**.

There should also be a function called **pllogis** for the cumulative distribution (if **d** is given), or **H** for the cumulative hazard (to complement **h**), if analytic forms for these are available. If not, then **flexsurv** can compute them internally by numerical integration, as in **stgenreg** (Crowther and Lambert 2013). The default options of the built-in R routine **integrate** for adaptive quadrature are used, though these may be changed using the **integ.opts** argument to **flexsurvreg**. Models specified this way will take an order of magnitude more time to fit, and the fitting procedure may be unstable. An example is given in §5.2.

These functions must be *vectorised*, and the density function must also accept an argument **log**, which when **TRUE**, returns the log density. See the examples below.

In some cases, R’s scoping rules may not find the functions in the working environment. They may then be supplied through the **dfns** argument to **flexsurvreg**.

**pars** Character vector naming the parameters of the distribution  $\mu, \alpha_1, \dots, \alpha_R$ . These must match the arguments of the R distribution function or functions, in the same order.

**location** Character: quoted name of the location parameter  $\mu$ . The location parameter will not necessarily be the first one, e.g., in **dweibull** the **scale** comes after the **shape**.

**transforms** A list of functions  $g()$  which transform the parameters from their natural ranges to the real line, for example, **c(log, identity)** if the first is positive and the second unrestricted.<sup>5</sup>

**inv.transforms** List of corresponding inverse functions.

<sup>4</sup>though since version 0.5.1, this distribution is built into **flexsurv** as **dist="llogis"**

<sup>5</sup>This is a *list*, not an *atomic vector* of functions, so if the distribution only has one parameter, we should write **transforms = c(log)** or **transforms = list(log)**, not **transforms = log**.

**inits** A function which provides plausible initial values of the parameters for maximum likelihood estimation. This is optional, but if not provided, then each call to **flexsurvreg** must have an **inits** argument containing a vector of initial values, which is inconvenient. Implausible initial values may produce a likelihood of zero, and a fatal error message (initial value in ‘vmmn’ is not finite) from the optimiser.

Each distribution will ideally have a heuristic for initialising parameters from summaries of the data. For example, since the median of the Weibull is  $\mu \log(2)^{1/\alpha}$ , a sensible estimate of  $\mu$  might be the median log survival time divided by  $\log(2)$ , with  $\alpha = 1$ , assuming that in practice the true value of  $\alpha$  is not far from 1. Then we would define the function, of one argument **t** giving the survival or censoring times, returning the initial values for the Weibull **shape** and **scale** respectively <sup>6</sup>.

```
inits = function(t) c(1, median(t[t > 0]) / log(2))
```

More complicated initial value functions may use other data such as the covariate values and censoring indicators: for an example, see the function **flexsurv.splineinits** in the package source that computes initial values for spline models (§5.1).

**Example: Using functions from a contributed package** The following custom model uses the log-logistic distribution functions (**dllogis** and **pllogis**) available in the package **eha**. The survivor function is  $S(t|\mu, \alpha) = 1/(1 + (t/\mu)^\alpha)$ , so that the log odds  $\log((1 - S(t))/S(t))$  of having died are a linear function of log time.

```
R> custom.llogis <- list(name = "llogis", pars = c("shape", "scale"),
+                       location = "scale",
+                       transforms = c(log, log),
+                       inv.transforms = c(exp, exp),
+                       inits = function(t){ c(1, median(t)) })
R> fs4 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc,
+                   dist = custom.llogis)
```

This fits the breast cancer data better than the Weibull, since it can represent a peaked hazard, but less well than the generalized gamma (Table 3).

**Example: Wrapping functions from a contributed package** Sometimes there may be probability density and similar functions in a contributed package, but in a different format. For example, **eha** also provides a three-parameter Gompertz-Makeham distribution with hazard  $h(t|\mu, \alpha_1, \alpha_2) = \alpha_2 + \alpha_1 \exp(t/\mu)$ . The shape parameters  $\alpha_1, \alpha_2$  are provided to **dmakeham** as a vector argument of length two. However, **flexsurvreg** expects distribution functions to have one argument for each parameter. Therefore we write our own functions that wrap around the third-party functions.

```
R> dmakeham3 <- function(x, shape1, shape2, scale, ...) {
+   dmakeham(x, shape = c(shape1, shape2), scale = scale, ...)
+ }
```

---

<sup>6</sup>though Weibull models in **flexsurvreg** are “initialised” by fitting the model with **survreg**, unless there is left-truncation.

```
R> pmakeham3 <- function(q, shape1, shape2, scale, ...) {
+   pmakeham(q, shape = c(shape1, shape2), scale = scale, ...)
+ }
```

`flexsurvreg` also requires these functions to be *vectorized*, as the standard distribution functions in R are. That is, we can supply a vector of alternative values for one or more arguments, and expect a vector of the same length to be returned. The R base function `Vectorize` can be used to do this here.

```
R> dmakeham3 <- Vectorize(dmakeham3)
R> pmakeham3 <- Vectorize(pmakeham3)
```

and this allows us to write, for example,

```
R> pmakeham3(c(0, 1, 1, Inf), 1, c(1, 1, 2, 1), 1)
```

We could then use `dist = list(name = "makeham3", pars = c("shape1", "shape2", "scale"), ...)` in a `flexsurvreg` model, though in the breast cancer example, the second shape parameter is poorly identifiable.

**Example: Changing the parameterisation of a distribution** We may want to fit a Weibull model like `fs1`, but with the proportional hazards (PH) parameterisation  $S(t) = \exp(-\mu t^\alpha)$ , so that the covariate effects reported in the printed `flexsurvreg` object can be interpreted as hazard ratios or log hazard ratios without any further transformation. Here instead of the density and cumulative distribution functions, we provide the hazard and cumulative hazard. (Note that since version 0.7, the "weibullPH" distribution is built in to `flexsurvreg` — but this example has been kept here for illustrative purposes.)<sup>7</sup>

```
R> hweibullPH <- function(x, shape, scale = 1, log = FALSE){
+   hweibull(x, shape = shape, scale = scale ^ {-1 / shape}, log = log)
+ }
R> HweibullPH <- function(x, shape, scale = 1, log = FALSE){
+   Hweibull(x, shape = shape, scale = scale ^ {-1 / shape}, log = log)
+ }
R> custom.weibullPH <- list(name = "weibullPH",
+   pars = c("shape", "scale"), location = "scale",
+   transforms = c(log, log),
+   inv.transforms = c(exp, exp),
+   inits = function(t){
+     c(1, median(t[t > 0]) / log(2))
+   })
R> fs6 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc,
+   dist = custom.weibullPH)
R> fs6$res["scale", "est"] ^ {-1 / fs6$res["shape", "est"]}
```

<sup>7</sup>The `eha` package may need to be detached first so that `flexsurv`'s built-in `hweibull` is used, which returns `NaN` if the parameter values are zero, rather than failing as `eha`'s currently does.

```
[1] 11.42287
```

```
R> - fs6$res["groupMedium", "est"] / fs6$res["shape", "est"]
```

```
[1] -0.61359
```

```
R> - fs6$res["groupPoor", "est"] / fs6$res["shape", "est"]
```

```
[1] -1.212215
```

The fitted model is the same as `fs1`, therefore the maximised likelihood is the same. The parameter estimates of `fs6` can be transformed to those of `fs1` as shown. The shape  $\alpha$  is common to both models, the scale  $\mu'$  in the AFT model is related to the PH scale  $\mu$  as  $\mu' = \mu^{-1/\alpha}$ . The effects  $\beta'$  on life expectancy in the AFT model are related to the log hazard ratios  $\beta$  as  $\beta' = -\beta/\alpha$ .

A slightly more complicated example is given in the package vignette `flexsurv-examples` of constructing a proportional hazards generalized gamma model. Note that `phreg` in `eha` also fits the Weibull and other proportional hazards models, though again the parameterisation is slightly different.

## 5. Arbitrary-dimension models

`flexsurv` also supports models where the number of parameters is arbitrary. In the models discussed previously, the number of parameters in the model family is fixed (e.g., three for the generalized gamma). In this section, the model complexity can be chosen by the user, given the model family. We may want to represent more irregular hazard curves by more flexible functions, or use bigger models if a bigger sample size makes it feasible to estimate more parameters.

### 5.1. Royston and Parmar spline model

In the spline-based survival model of [Royston and Parmar \(2002\)](#), a transformation  $g(S(t, z))$  of the survival function is modelled as a natural cubic spline function of log time:  $g(S(t, z)) = s(x, \gamma)$  where  $x = \log(t)$ . This model can be fitted in `flexsurv` using the function `flexsurvspline`, and is also available in the `Stata` package `stpm2` ([Lambert and Royston 2009](#)) (historically `stpm`, [Royston \(2001, 2004\)](#)).

Typically we use  $g(S(t, \mathbf{z})) = \log(-\log(S(t, \mathbf{z}))) = \log(H(t, \mathbf{z}))$ , the log cumulative hazard, giving a proportional hazards model.

**Spline parameterisation** The complexity of the model, thus the dimension of  $\gamma$ , is governed by the number of *knots* in the spline function  $s(\cdot)$ . Natural cubic splines are piecewise cubic polynomials defined to be continuous, with continuous first and second derivatives at the knots, and also constrained to be linear beyond boundary knots  $k_{min}, k_{max}$ . As well as the boundary knots there may be up to  $m \geq 0$  *internal* knots  $k_1, \dots, k_m$ . Various spline parameterisations exist — the one used here is from [Royston and Parmar \(2002\)](#).

$$s(x, \gamma) = \gamma_0 + \gamma_1 x + \gamma_2 v_1(x) + \dots + \gamma_{m+1} v_m(x) \quad (4)$$

Model	$g(S(t, \mathbf{z}))$	In <code>flexsurvspline</code>	With $m = 0$
Proportional hazards	$\log(-\log(S(t, \mathbf{z})))$ (log cumulative hazard)	<code>scale = "hazard"</code>	Weibull <code>shape</code> $\gamma_1$ , <code>scale</code> $\exp(-\gamma_0/\gamma_1)$
Proportional odds	$\log(S(t, \mathbf{z})^{-1} - 1)$ (log cumulative odds)	<code>scale = "odds"</code>	Log-logistic <code>shape</code> $\gamma_1$ , <code>scale</code> $\exp(-\gamma_0/\gamma_1)$
Normal / probit	$\Phi^{-1}(S(t, \mathbf{z}))$ (inverse normal CDF, <code>qnorm</code> )	<code>scale = "normal"</code>	Log-normal <code>meanlog</code> $-\gamma_0/\gamma_1$ , <code>sdlog</code> $1/\gamma_1$

Table 2: Alternative modelling scales for `flexsurvspline`, and equivalent distributions for  $m = 0$  (with parameter definitions as in the R `d` functions referred to elsewhere in the paper).

where  $v_j(x)$  is the  $j$ th *basis* function

$$v_j(x) = (x - k_j)_+^3 - \lambda_j(x - k_{\min})_+^3 - (1 - \lambda_j)(x - k_{\max})_+^3, \quad \lambda_j = \frac{k_{\max} - k_j}{k_{\max} - k_{\min}}$$

and  $(x - a)_+ = \max(0, x - a)$ . If  $m = 0$  then there are only two parameters  $\gamma_0, \gamma_1$ , and this is a Weibull model if  $g()$  is the log cumulative hazard. Table 2 explains two further choices of  $g()$ , and the parameter values and distributions they simplify to for  $m = 0$ . The probability density and cumulative distribution functions for all these models are available as `dsurvspline` and `psurvspline`. A model with an absolute time scale ( $x = t$ ) is also available through `timescale="identity"`.

**Covariates on spline parameters** Covariates can be placed on any parameter  $\gamma$  through a linear model (with identity link function). Most straightforwardly, we can let the intercept  $\gamma_0$  vary with covariates  $\mathbf{z}$ , giving a proportional hazards or odds model (depending on  $g()$ ).

$$g(S(t, \mathbf{z})) = s(\log(t), \gamma) + \beta^\top \mathbf{z}$$

The spline coefficients  $\gamma_j : j = 1, 2, \dots$ , the “ancillary” parameters, may also be modelled as linear functions of covariates  $\mathbf{z}$ , as

$$\gamma_j(\mathbf{z}) = \gamma_{j0} + \gamma_{j1}z_1 + \gamma_{j2}z_2 + \dots$$

giving a model where the effects of covariates are arbitrarily flexible functions of time: a non-proportional hazards or odds model.

**Spline models in flexsurv** The argument `k` to `flexsurvspline` defines the number of internal knots  $m$ . Knot locations are chosen by default from quantiles of the log uncensored death times, or users can supply their own locations in the `knots` argument. Initial values for numerical likelihood maximisation are chosen using the method described by Royston and Parmar (2002) of Cox regression combined with transforming an empirical survival estimate. For example, the best-fitting model for the breast cancer dataset identified in Royston and Parmar (2002), a proportional odds model with one internal spline knot, is



```
R> sp1 <- flexsurvspline(Surv(recyrs, censrec) ~ group, data = bc, k = 1,
+                        scale = "odds")
```

A further model where the first ancillary parameter also depends on the prognostic group, giving a time-varying odds ratio, is fitted as

```
R> sp2 <- flexsurvspline(Surv(recyrs, censrec) ~ group + gamma1(group),
+                        data = bc, k = 1, scale = "odds")
```

These models give qualitatively similar results to the generalized gamma in this dataset (Figure 3), and have similar predictive ability as measured by AIC (Table 3). Though in general, an advantage of spline models is that extra flexibility is available where necessary.

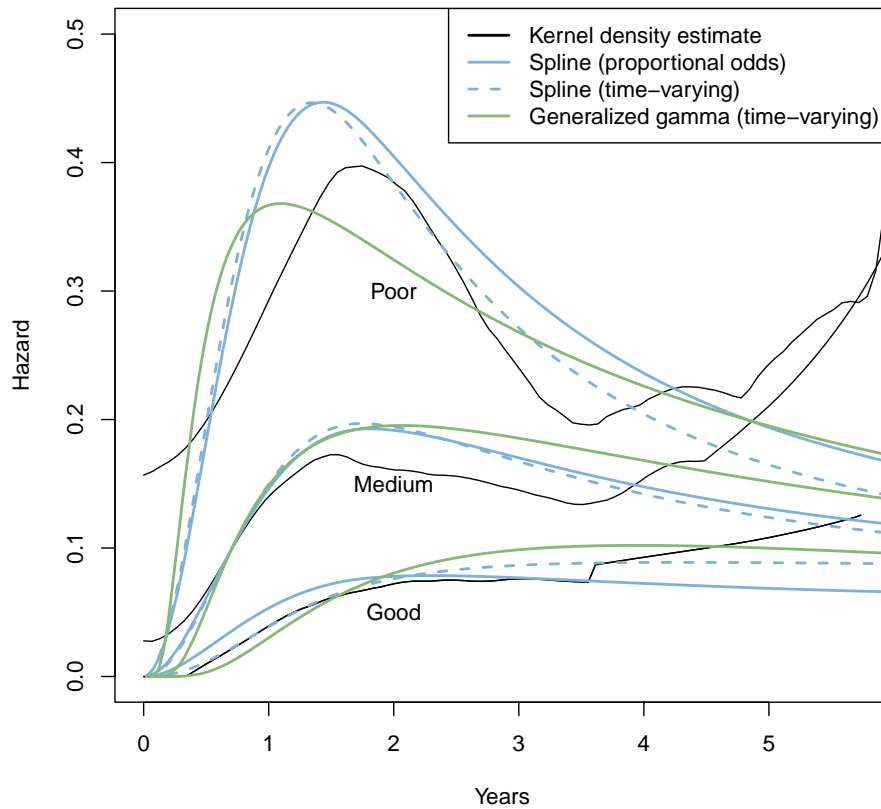


Figure 3: Comparison of spline and generalized gamma fitted hazards for the breast cancer survival data by prognostic group.

In this example, proportional odds models (`scale = "odds"`) are better-fitting than proportional hazards models (`scale = "hazard"`) (Table 3). Note also that under a proportional hazards spline model with one internal knot (`sp3`), the log hazard ratios, and their standard errors, are substantively the same as under a standard Cox model (`cox3`). This illustrates that this class of flexible fully-parametric models may be a reasonable alternative to the

(semi-parametric) Cox model. See [Royston and Parmar \(2002\)](#) for more discussion of these issues.

```
R> sp3 <- flexsurvspline(Surv(recyrs, censrec) ~ group, data = bc, k = 1,
+                        scale = "hazard")
R> sp3$res[c("groupMedium", "groupPoor"), c("est", "se")]
```

```
              est      se
groupMedium 0.8345026 0.1712764
groupPoor   1.6120990 0.1641750
```

```
R> cox3 <- coxph(Surv(recyrs, censrec) ~ group, data = bc)
R> coef(summary(cox3))[ , c("coef", "se(coef)")]
```

```
              coef  se(coef)
groupMedium 0.8401002 0.1713926
groupPoor   1.6180720 0.1645443
```

An equivalent of a “stratified” Cox model may be obtained by allowing *all* the spline parameters to vary with the categorical covariate that defines the strata. In this case, this covariate might be `group`. With  $k=m$  internal knots, the formula should then include `group`, representing  $\gamma_0$ , and  $m+1$  further terms representing the parameters  $\gamma_1, \dots, \gamma_{m+1}$ , named as follows.

```
R> sp4 <- flexsurvspline(Surv(recyrs, censrec) ~ group + gamma1(group) +
+                        gamma2(group), data = bc, k = 1, scale = "hazard")
```

Other covariates might be added to this formula — if placed on the intercept, these will be modelled through proportional hazards, as in `sp1`. If placed on higher-order parameters, these will represent time-varying hazard ratios. For example, if there were a covariate `treat` representing treatment, then

```
R> flexsurvspline(Surv(recyrs, censrec) ~ group + gamma1(group) +
+                gamma2(group) + treat + gamma1(treat),
+                data = bc, k = 1, scale = "hazard")
```

would represent a model stratified by `group`, where the hazard ratio for treatment is time-varying, but the model is not fully stratified by treatment.

```
R> res <- t(sapply(list(fs1, fs2, fs3, sp1, sp2, sp3, sp4),
+                  function(x) rbind(-2 * round(x$loglik, 1), x$npars,
+                                              round(x$AIC, 1))))
R> rownames(res) <- c("Weibull (fs1)", "Generalized gamma (fs2)",
+                  "Generalized gamma (fs3)",
+                  "Spline (sp1)", "Spline (sp2)", "Spline (sp3)",
+                  "Spline (sp4)")
R> colnames(res) <- c("-2 log likelihood", "Parameters", "AIC")
```

```
R> res
```

	-2 log likelihood	Parameters	AIC
Weibull (fs1)	1623.8	4	1631.9
Generalized gamma (fs2)	1575.2	5	1585.1
Generalized gamma (fs3)	1572.4	7	1586.4
Spline (sp1)	1578.0	5	1588.0
Spline (sp2)	1574.8	7	1588.8
Spline (sp3)	1585.8	5	1595.7
Spline (sp4)	1571.4	9	1589.3

Table 3: Comparison of parametric survival models fitted to the breast cancer data.

## 5.2. Implementing new general-dimension models

The spline model above is an example of the general parametric form (Equation 1), but the number of parameters,  $R + 1$  in Equation 1,  $m + 2$  in Equation 4, is arbitrary. **flexsurv** has the tools to deal with any model of this form. **flexsurvspline** works internally by building a custom distribution and then calling **flexsurvreg**. Similar models may in principle be built by users using the same method. This relies on a functional programming trick.

**Creating distribution functions dynamically** The R distribution functions supplied to custom models are expected to have a fixed number of arguments, including one for each scalar parameter. However, the distribution functions for the spline model (e.g., **dsurvspline**) have an argument **gamma** representing the *vector* of parameters  $\gamma$ , whose length is determined by choosing the number of knots. Just as the *scalar parameters* of conventional distribution functions can be supplied as *vector arguments* (as explained in §4), similarly, the vector parameters of spline-like distribution functions can be supplied as *matrix arguments*, representing alternative parameter values.

To convert a spline-like distribution function into the correct form, **flexsurv** provides the utility **unroll.function**. This converts a function with one (or more) vector parameters (matrix arguments) to a function with an arbitrary number of scalar parameters (vector arguments). For example, the 5-year survival probability for the baseline group under the model **sp1** is

```
R> gamma <- sp1$res[c("gamma0", "gamma1", "gamma2"), "est"]
R> 1 - psurvspline(5, gamma = gamma, knots = sp1$knots)

[1] 0.6897025
```

An alternative function to compute this can be built by **unroll.function**. We tell it that the vector parameter **gamma** should be provided instead as three scalar parameters named **gamma0**, **gamma1**, **gamma2**. The resulting function **pfn** is in the correct form for a custom **flexsurvreg** distribution.

```
R> pfn <- unroll.function(psurvspline, gamma = 0:2)
R> 1 - pfn(5, gamma0 = gamma[1], gamma1 = gamma[2], gamma2 = gamma[3],
+       knots = sp1$knots)
```

[1] 0.6897025

Users wishing to fit a new spline-like model with a known number of parameters could just as easily write distribution functions specific to that number of parameters, and use the methods in §4. However the `unroll.function` method is intended to simplify the process of extending the **flexsurv** package to implement new model families, through wrappers similar to `flexsurvspline`.

**Example: splines on alternative scales** An alternative to the Royston-Parmar spline model is to model the log *hazard* as a spline function of (log) time instead of the log cumulative hazard. Crowther and Lambert (2013) demonstrate this model using the Stata `stgenreg` package. An advantage explained by Royston and Lambert (2011) is that when there are multiple time-dependent effects, time-dependent hazard ratios can be interpreted independently of the values of other covariates.

This can also be implemented in `flexsurvreg` using `unroll.function`. A disadvantage of this model is that the cumulative hazard (hence the survivor function) has no analytic form, therefore to compute the likelihood, the hazard function needs to be integrated numerically. This is done automatically in `flexsurvreg` (just as in `stgenreg`) if the cumulative hazard is not supplied.

Firstly, a function must be written to compute the hazard as a function of time `x`, the vector of parameters `gamma` (which can be supplied as a matrix argument so the function can give a vector of results), and a vector of knot locations. This uses **flexsurv**'s function `basis` to compute the natural cubic spline basis (Equation 4), and replicates `x` and `gamma` to the length of the longest one.

```
R> hsurvspline.lh <- function(x, gamma, knots){
+   if(!is.matrix(gamma)) gamma <- matrix(gamma, nrow = 1)
+   lg <- nrow(gamma)
+   nret <- max(length(x), lg)
+   gamma <- apply(gamma, 2, function(x)rep(x, length = nret))
+   x <- rep(x, length = nret)
+   loghaz <- rowSums(basis(knots, log(x)) * gamma)
+   exp(loghaz)
+ }
```

The equivalent function is then created for a three-knot example of this model (one internal and two boundary knots) that has arguments `gamma0`, `gamma1` and `gamma2` corresponding to the three columns of `gamma`,

```
R> hsurvspline.lh3 <- unroll.function(hsurvspline.lh, gamma = 0:2)
```

To complete the model, the custom distribution list is formed, the internal knot is placed at the median uncensored log survival time, and the boundary knots are placed at the minimum and maximum. These are passed to `hsurvspline.lh` through the `aux` argument of `flexsurvreg`.

```
R> custom.hsurvspline.lh3 <- list(
+   name = "survspline.lh3",
```

```

+   pars = c("gamma0", "gamma1", "gamma2"),
+   location = c("gamma0"),
+   transforms = rep(c(identity), 3), inv.transforms = rep(c(identity), 3)
+ )
R> dtime <- log(bc$recyrs)[bc$censrec == 1]
R> ak <- list(knots = quantile(dtime, c(0, 0.5, 1)))

```

Initial values must be provided in the call to `flexsurvreg`, since the custom distribution list did not include an `inits` component. For this example, “default” initial values of zero suffice, but the permitted values of  $\gamma_2$  are fairly tightly constrained (from -0.5 to 0.5 here) using the “L-BFGS-B” bounded optimiser from R’s `optim` (Nash 1990). Without the constraint, extreme values of  $\gamma_2$ , visited by the optimiser, cause the numerical integration of the hazard function to fail.

```

R> sp5 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data = bc, aux = ak,
+   inits = c(0, 0, 0, 0, 0),
+   dist = custom.hsurvspline.lh3,
+   method = "L-BFGS-B", lower = c(-Inf, -Inf, -0.5),
+   upper = c(Inf, Inf, 0.5),
+   control = list(trace = 1, REPORT = 1))

```

This takes around ten minutes to converge, so is not presented here, though the fit is poorer than the equivalent spline model for the cumulative hazard. The 95% confidence interval for  $\gamma_2$  of (0.16, 0.37) is firmly within the constraint. Crowther and Lambert (2014) present a combined analytic / numerical integration method for this model that may make fitting it more stable.

**Other arbitrary-dimension models** Another potential application is to fractional polynomials (Royston and Altman 1994). These are of the form  $\sum_{m=1}^M \alpha_m x^{p_m} \log(x)^n$  where the power  $p_m$  is in the standard set  $\{2, -1, -0.5, 0, 0.5, 1, 2, 3\}$  (except that  $\log(x)$  is used instead of  $x^0$ ), and  $n$  is a non-negative integer. They are similar to splines in that they can give arbitrarily close approximations to a nonlinear function, such as a hazard curve, and are particularly useful for expressing the effects of continuous predictors in regression models. See e.g., Sauerbrei *et al.* (2007), and several other publications by the same authors, for applications and discussion of their advantages over splines. The R package `gamlss` (Rigby and Stasinopoulos 2005) has a function to construct a fractional polynomial basis that might be employed in `flexsurv` models.

Polyhazard models (Louzada-Neto 1999) are another potential use of this technique. These express an overall hazard as a sum of latent cause-specific hazards, each one typically from the same class of distribution, e.g., a *poly-Weibull* model if they are all Weibull. For example, a U-shaped hazard curve following surgery may be the sum of early hazards from surgical mortality and later deaths from natural causes. However, such models may not always be identifiable without external information to fix or constrain the parameters of particular hazards (Demiris *et al.* 2011).

## 6. Multi-state models

A *multi-state model* represents how an individual moves between multiple states in continuous time. Survival analysis is a special case with two states, “alive” and “dead”. *Competing risks* are a further special case, where there are multiple causes of death, that is, one starting state and multiple possible destination states.

Multi-state modelling with **flexsurv** was previously described in this section of the current vignette. Version 2.0 of **flexsurv** added several new features for multi-state modelling, including multi-state modelling using mixtures, and transition-specific distribution families in cause-specific hazards models. These models are now fully described in a separate **flexsurv** vignette, “Flexible parametric multi-state modelling with the flexsurv package”.

## 7. Potential extensions

Multi-state modelling is still an area of ongoing work, and while version 2.0 extended **flexsurv** in this area, more tools and documentation in this area would still be useful. The **msm** package arguably has a more accessible interface for fitting and summarising multi-state models, but it was designed mainly for panel data rather than event time data, and therefore the event time distributions it fits are relatively inflexible.

Models where multiple survival times are assumed to be correlated within groups, sometimes called (shared) frailty models (Hougaard 1995), would also be a useful development. See, e.g., Crowther *et al.* (2014) for a recent application based on parametric models. These might be implemented by exploiting tractability for specific distributions, such as gamma frailties, or by adjusting standard errors to account for clustering, as implemented in **survreg**. More complex random effects models would require numerical integration, for example, Crowther *et al.* (2014) provide Stata software based on Gauss-Hermite quadrature. Alternatively, a probabilistic modelling language such as Stan (Stan Development Team 2014) or BUGS (Lunn *et al.* 2012) would be naturally suited to complex extensions such as random effects on multiple parameters or multiple hierarchical levels.

**flexsurv** is intended as a platform for parametric survival modelling. Extensions of the software to deal with different models may be written by users themselves, through the facilities described in §4 and §5.2. These might then be included in the package as built-in distributions, or at least demonstrated in the package’s other vignette **flexsurv-examples**. Each new class of models would ideally come with

- guidance on what situations the model is useful for, e.g., what shape of hazards it can represent
- some intuitive interpretation of the model parameters, their plausible values in typical situations, and potential identifiability problems. This would also help with choosing initial values for numerical maximum likelihood estimation, ideally through an **inits** function in the custom distribution list (§4).

**flexsurv** is available from <http://CRAN.R-project.org/package=flexsurv>. Development versions are available on <https://github.com/chjackson/flexsurv-dev>, and contributions are welcome.

## Acknowledgements

Thanks to Milan Bouchet-Valat for help with implementing covariates on ancillary parameters, Andrea Manca for motivating the development of the package, the reviewers of the paper, and all users who have reported bugs and given suggestions.

## References

- Aalen O, Borgan O, Gjessing H (2008). *Survival and Event History Analysis: A Process Point of View*. Springer-Verlag.
- Benaglia T, Jackson CH, Sharples LD (2014). “Survival Extrapolation in the Presence of Cause Specific Hazards.” *Statistics in Medicine*. In press.
- Broström G (2014). **eha**: *Event History Analysis*. R package version 2.4-1, URL <http://CRAN.R-project.org/package=eha>.
- Cox C (2008). “The Generalized F Distribution: An Umbrella for Parametric Survival Analysis.” *Statistics in Medicine*, **27**(21), 4301–4312.
- Crowther MJ, Lambert PC (2013). “**stgenreg**: A Stata Package for General Parametric Survival Analysis.” *Journal of Statistical Software*, **53**, 1–17.
- Crowther MJ, Lambert PC (2014). “A General Framework for Parametric Survival Analysis.” *Statistics in Medicine*. Early view, DOI: 10.1002/sim.6300.
- Crowther MJ, Look MP, Riley RD (2014). “Multilevel Mixed Effects Parametric Survival Models Using Adaptive Gauss–Hermite Quadrature With Application to Recurrent Events and Individual Participant Data Meta-Analysis.” *Statistics In Medicine*, **33**(22), 3844–3858.
- de Wreede LC, Fiocco M, Putter H (2011). “**mstate**: An R Package for the Analysis of Competing Risks and Multi-State Models.” *Journal of Statistical Software*, **38**, 1–30.
- Demiris N, Lunn D, Sharples L (2011). “Survival Extrapolation Using the Poly-Weibull Model.” *Statistical Methods in Medical Research*. Early view, DOI: 10.1177/0962280211419645.
- Hess K (2010). **muhaz**: *Hazard Function Estimation in Survival Analysis*. R package version 1.2.5, S original by K. Hess and R port by R. Gentleman, URL <http://CRAN.R-project.org/package=muhaz>.
- Hothorn T (2015). **TH.data**: *TH’s Data Archive*. R package version 1.0-6, URL <http://CRAN.R-project.org/package=TH.data>.
- Hougaard P (1995). “Frailty Models for Survival Data.” *Lifetime Data Analysis*, **1**(3), 255–273.
- 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.



- Lambert PC, Royston P (2009). “Further Development of Flexible Parametric Models for Survival Analysis.” *Stata Journal*, **9**(2), 265.
- Latimer NR (2013). “Survival Analysis for Economic Evaluations Alongside Clinical Trials — Extrapolation with Patient-Level Data, Inconsistencies, Limitations, and a Practical Guide.” *Medical Decision Making*, **33**(6), 743–754.
- Louzada-Neto F (1999). “Polyhazard Models for Lifetime Data.” *Biometrics*, **55**, 1281–1285.
- Lunn D, Jackson C, Best N, Thomas A, Spiegelhalter D (2012). *The BUGS Book: A Practical Introduction to Bayesian Analysis*. CRC Press.
- Mandel M (2013). “Simulation-Based Confidence Intervals for Functions with Complicated Derivatives.” *The American Statistician*, **67**(2), 76–81.
- Mueller HG, Wang JL (1994). “Hazard Rates Estimation Under Random Censoring with Varying Kernels and Bandwidths.” *Biometrics*, **50**, 61–76.
- Nadarajah S, Bakar SAA (2013). “A New R Package for Actuarial Survival Models.” *Computational Statistics*, **28**(5), 2139–2160.
- Nash JC (1990). *Compact Numerical Methods for Computers: Linear Algebra and Function Minimisation*. CRC Press.
- Nelson CP, Lambert PC, Squire IB, Jones DR (2007). “Flexible Parametric Models for Relative Survival, With Application in Coronary Heart Disease.” *Statistics in Medicine*, **26**(30), 5486–5498.
- Prentice RL (1974). “A Log Gamma Model and its Maximum Likelihood Estimation.” *Biometrika*, **61**(3), 539–544.
- Prentice RL (1975). “Discrimination Among Some Parametric Models.” *Biometrika*, **62**(3), 607–614.
- Stan Development Team (2014). *Stan Modeling Language Users Guide and Reference Manual, Version 2.4*. URL <http://mc-stan.org/>.
- R Core Team (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- Reid N (1994). “A Conversation with Sir David Cox.” *Statistical Science*, **9**(3), 439–455.
- Rigby RA, Stasinopoulos DM (2005). “Generalized Additive Models for Location, Scale and Shape (with discussion).” *Journal of the Royal Statistical Society C*, **54**(3), 507–554.
- Royston P (2001). “Flexible Parametric Alternatives to the Cox Model, and More.” *Stata Journal*, **1**(1), 1–28.
- Royston P (2004). “Flexible Parametric Alternatives to the Cox Model: Update.” *The Stata Journal*, **4**(1), 98–101.
- Royston P, Altman DG (1994). “Regression Using Fractional Polynomials of Continuous Covariates: Parsimonious Parametric Modelling.” *Journal of the Royal Statistical Society C*, **43**(3), 429–467.



- Royston P, Lambert PC (2011). “Flexible Parametric Survival Analysis Using Stata: Beyond the Cox Model.” *Stata Press books*.
- Royston P, Parmar M (2002). “Flexible Parametric Proportional-Hazards and Proportional-Odds Models for Censored Survival Data, with Application to Prognostic Modelling and Estimation of Treatment Effects.” *Statistics in Medicine*, **21**(1), 2175–2197.
- Sauerbrei W, Royston P (1999). “Building Multivariable Prognostic and Diagnostic Models: Transformation of the Predictors by Using Fractional Polynomials.” *Journal of the Royal Statistical Society A*, **162**(1), 71–94.
- Sauerbrei W, Royston P, Binder H (2007). “Selection of Important Variables and Determination of Functional Form for Continuous Predictors in Multivariable Model Building.” *Statistics in Medicine*, **26**(30), 5512–5528.
- Stacy EW (1962). “A Generalization of the Gamma Distribution.” *The Annals of Mathematical Statistics*, **33**(3), 1187–92.
- Therneau T (2014). “A Package for Survival Analysis in S.” R package version 2.37-7. <http://CRAN.R-project.org/package=survival>.
- Yee TW, Wild CJ (1996). “Vector Generalized Additive Models.” *Journal of the Royal Statistical Society B*, **58**(3), 481–493.