# flexsurv: flexible parametric survival modelling in R. Supplementary examples

**Christopher H. Jackson**

MRC Biostatistics Unit, Cambridge, UK

chris.jackson@mrc-bsu.cam.ac.uk

---

**Abstract**

This vignette of examples supplements the main **flexsurv** user guide.

*Keywords*: survival.

---

# 1. Examples of custom distributions

## 1.1. Proportional hazards generalized gamma model

Crowther and Lambert (2013) discuss using the **stgenreg** Stata package to construct a proportional hazards parameterisation of the three-parameter generalised gamma distribution. A similar trick can be used in **flexsurv**. A four-parameter custom distribution is created by defining its hazard (and cumulative hazard) functions. These are obtained by multiplying the built-in functions `hgengamma` and `Hgengamma` by an extra dummy parameter, which is used as the location parameter of the new distribution. The intercept of this parameter is fixed at 1 when calling `flexsurvreg`, so that the new model is no more complex than the generalized gamma AFT model `fs3`, but covariate effects on the dummy parameter are now interpreted as hazard ratios.

```
R> library(flexsurv)
R> hgengammaPH <- function(x, dummy, mu=0, sigma=1, Q){
+      dummy * hgengamma(x=x, mu=mu, sigma=sigma, Q=Q)
+ }
R> HgengammaPH <- function(x, dummy, mu=0, sigma=1, Q){
+      dummy * Hgengamma(x=x, mu=mu, sigma=sigma, Q=Q)
+ }
R> custom.gengammaPH <- list(name="gengammaPH",
+                       pars=c("dummy","mu","sigma","Q"), location="dummy",
+                       transforms=c(log, identity, log, identity),
+                       inv.transforms=c(exp, identity, exp, identity),
+                       inits=function(t){
+                           lt <- log(t[t>0])
+                           c(1, mean(lt), sd(lt), 0)
```

```
+                              })
R> fs7 <- flexsurvreg(Surv(recyrs, censrec) ~ group, data=bc,
+                     dist=custom.gengammaPH, fixedpars=1)
```

## 2. Examples of custom model summaries

### 2.1. Plotting a hazard ratio against time

The following code plots the hazard ratio (Medium versus Good prognostic group) against time for both the proportional hazards model `fs7` and the better-fitting accelerated failure time model `fs2`. It illustrates the use of the following functions.

`summary.flexsurvreg` for generating the estimated hazard at a series of times, for particular covariate categories.

`normboot.flexsurvreg` for generating a bootstrap-style sample from the sampling distribution of the parameter estimates, for particular covariate categories.

`do.call` for constructing a function call by supplying a list containing the function's arguments. This is used throughout the source of **flexsurv**.

```
R> fs2 <- flexsurvreg(Surv(recyrs, censrec) ~ group + sigma(group),
+                     data=bc, dist="gengamma")
R> B <- 5000
R> t <- seq(0.1, 8, by=0.1)
R> hrAFT.est <-
+     summary(fs2, t=t, type="hazard",
+             newdata=data.frame(group="Medium"),ci=FALSE)[[1]][,"est"] /
+     summary(fs2, t=t, type="hazard",
+             newdata=data.frame(group="Good"),ci=FALSE)[[1]][,"est"]
R> pars <- normboot.flexsurvreg(fs2, B=B, newdata=data.frame(group=c("Good","Medium")))
R> hrAFT <- matrix(nrow=B, ncol=length(t))
R> for (i in seq_along(t)){
+     haz.medium.rep <- do.call(hgengamma, c(list(t[i]), as.data.frame(pars[[2]])))
+     haz.good.rep <- do.call(hgengamma, c(list(t[i]), as.data.frame(pars[[1]])))
+     hrAFT[,i] <- haz.medium.rep / haz.good.rep
+ }
R> hrAFT <- apply(hrAFT, 2, quantile, c(0.025, 0.975))
R> hrPH.est <-
+     summary(fs7, t=t, type="hazard",
+             newdata=data.frame(group="Medium"),ci=FALSE)[[1]][,"est"] /
+     summary(fs7, t=t, type="hazard",
+             newdata=data.frame(group="Good"),ci=FALSE)[[1]][,"est"]
R> pars <- normboot.flexsurvreg(fs7, B=B, newdata=data.frame(group=c("Good","Medium")))
R> hrPH <- matrix(nrow=B, ncol=length(t))
```
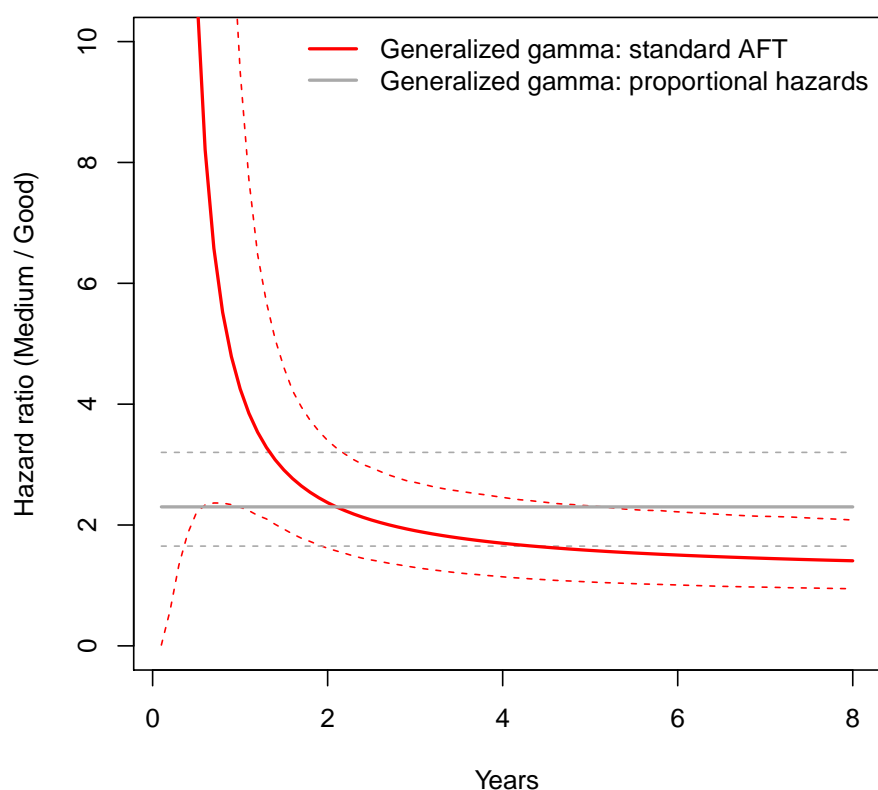
```
R> for (i in seq_along(t)){
+     haz.medium.rep <- do.call(hgengammaPH, c(list(t[i]), as.data.frame(pars[[2]])))
+     haz.good.rep <- do.call(hgengammaPH, c(list(t[i]), as.data.frame(pars[[1]])))
+     hrPH[,i] <- haz.medium.rep / haz.good.rep
+ }
R> hrPH <- apply(hrPH, 2, quantile, c(0.025, 0.975))
R> plot(t, hrAFT[1,], type="l", ylim=c(0, 10), col="red", xlab="Years",
+      ylab="Hazard ratio (Medium / Good)", lwd=1, lty=2)
R> lines(t, hrAFT[2,], col="red", lwd=1, lty=2)
R> lines(t, hrPH[1,], col="darkgray", lwd=1, lty=2)
R> lines(t, hrPH[2,], col="darkgray", lwd=1, lty=2)
R> lines(t, hrAFT.est, col="red", lwd=2)
R> lines(t, hrPH.est, col="darkgray", lwd=2)
R> legend("topright", lwd=c(2,2), col=c("red","darkgray"), bty="n",
+        c("Generalized gamma: standard AFT", "Generalized gamma: proportional hazards"))
```



## 2.2. Restricted mean survival

The expected survival up to time $t$, from a model with cumulative distribution $F(t|\alpha)$, is

$$E(T|T < t) = \int_0^t 1 - F(u|\alpha)du$$

An estimate and confidence interval for this, for a specified covariate value, can be computed using a custom summary function as follows. (Note that `summary.flexsurvreg` can be abbreviated to `summary`). As in the `median.weibull` example in the user guide vignette, the summary function is independent of time, so any value can be specified for `t` in the call to `summary`. The time horizon up to which to compute the mean is specified by the default value of the `horizon` argument to the custom function. The mean survival is computed here up to 100 years. Setting `horizon=Inf` is theoretically also possible for an unrestricted mean, but the integral does not converge in this example.

```
R> mean.gengamma <- function(mu, sigma, Q, horizon=100, ...){
+     surv <- function(t, ...) {  1 - pgengamma(q=t, mu=mu, sigma=sigma, Q=Q, ...) }
+     integrate(surv, 0, horizon, ...)$value
+ }
R> summary(fs2, newdata=list(group="Good"), t=1, fn=mean.gengamma)


group=Good
  time      est      lcl      ucl
1    1 21.50102 13.53971 31.73238


R> summary(fs2, newdata=list(group="Medium"), t=1, fn=mean.gengamma)


group=Medium
  time      est      lcl      ucl
1    1 12.00025 8.082643 17.74195


R> summary(fs2, newdata=list(group="Poor"), t=1, fn=mean.gengamma)


group=Poor
  time      est     lcl      ucl
1    1 5.539196 3.82227 9.299211
```

Note that the (unrestricted) median is more stable, and less than the restricted mean due to the skewness of this distribution.

```
R> median.gengamma <- function(mu, sigma, Q) {
+     qgengamma(0.5, mu=mu, sigma=sigma, Q=Q)
+ }
R> summary(fs2, newdata=list(group="Good"), t=1, fn=median.gengamma)


group=Good
  time      est      lcl      ucl
1    1 9.736555 7.352547 13.91304
```

# 3. Spline models

## 3.1. Prognostic model for the German breast cancer data

The regression model III in Sauerbrei and Royston (1999) used to create the prognostic group from the breast cancer data (supplied as `bc` in **flexsurv** and `GBSG2` in **TH.data**) can be reproduced as follows. Firstly, the required fractional polynomial transformations of the covariates are constructed. `progc` implements the Cox model used by Sauerbrei and Royston (1999), and `prog3` is a flexible fully-parametric alternative, implemented as a spline with three internal knots. The number of knots was chosen to minimise AIC. The covariate effects are very similar.

After fitting the model, the prognostic index can then be derived from categorising observations in three groups according to the tertiles of the linear predictor in each model. The indices produced by the Cox model (`progc`) and the spline-based model (`progf`) agree exactly.

```
R> if (require("TH.data")){
+
+ GBSG2 <- transform(GBSG2,
+                    X1a=(age/50)^-2,
+                    X1b=(age/50)^-0.5,
+                    X4=tgrade %in% c("II","III"),
+                    X5=exp(-0.12*pnodes),
+                    X6=(progrec+1)^0.5
+                    )
+ (progc <- coxph(Surv(time, cens) ~ horTh + X1a + X1b + X4 +
+                 X5 + X6, data=GBSG2))
+ (prog3 <- flexsurvspline(Surv(time, cens) ~ horTh + X1a + X1b + X4 +
+                          X5 + X6, k=3, data=GBSG2))
+ predc <- predict(progc, type="lp")
+ progc <- cut(predc, quantile(predc, 0:3/3))
+ predf <- model.matrix(prog3) %*% prog3$res[-(1:5),"est"]
+ progf <- cut(predf, quantile(predf, 0:3/3))
+ table(progc, progf)
+
+ }
```

|                  | progf            |                  |                  |
|------------------|------------------|------------------|------------------|
| progc            | (-9.91,-7.76]    | (-7.76,-7.11]    | (-7.11,-3.36]    |
| (-2.52,-0.359]   | 228              | 0                | 0                |
| (-0.359,0.285]   | 0                | 228              | 0                |
| (0.285,4]        | 0                | 0                | 229              |

# 4. Right truncation: retrospective ascertainment

Suppose we want to estimate the distribution of the time from onset of a disease to death, but have only observed cases known to have died by the current date. In this case, times from

onset to death for individuals in the data are *right-truncated* by the current date minus the onset date. Predicted survival times for new cases can then be described by an un-truncated version of the fitted distribution. Denote the time from onset to death as the *delay* time.

This is illustrated in the following simulated example. Suppose individual onset times are uniformly distributed between 0 and 30 days. Their delay times are generated from a Gamma distribution.

```
R> set.seed(1)
R> nsim <- 10000
R> onsetday <- runif(nsim, 0, 30)
R> deathday <- onsetday + rgamma(nsim, shape=1.5, rate=1/10)
```

The data are examined at 40 days. Therefore we do not observe people who have died after this time. For each individual in the observed data, their delay time is right-truncated by 40 days minus their onset day, since their delay times cannot be greater than this if they are included in the data. The right-truncation point is specified by the variable `rtrunc` in the data.

```
R> datt <- data.frame(delay = deathday - onsetday,
+                     event = rep(1, nsim),
+                     rtrunc = 40 - onsetday)
R> datt <- datt[datt$delay < datt$rtrunc, ]
```

The truncated Gamma model is fitted with `flexsurvreg` by specifying individual-specific truncation points in the argument `rtrunc`. The fitted model reproduces the gamma parameters that were used to simulate the data. After fitting the model, we can use the fitted model to predict the mean time to death - this is approximately the shape / rate of the untruncated gamma distribution.

```
R> fitt <- flexsurvreg(Surv(delay, event) ~ 1, data=datt, rtrunc = rtrunc, dist="gamma")
R> fitt

Call:
flexsurvreg(formula = Surv(delay, event) ~ 1, data = datt, rtrunc = rtrunc,
    dist = "gamma")

Estimates:
        est      L95%     U95%     se
shape   1.50886  1.45760  1.56194  0.02661
rate    0.10067  0.09439  0.10737  0.00331

N = 7728,  Events: 7728,  Censored: 0
Total time at risk: 82028.85
Log-likelihood = -24252.16, df = 2
AIC = 48508.32

R> summary(fitt, t=1, fn = mean_gamma)
```

```
  time      est      lcl      ucl
1    1 14.98784 14.40418 15.62192
```

# References

Crowther MJ, Lambert PC (2013). "**stgenreg**: A Stata Package for General Parametric Survival Analysis." *Journal of Statistical Software*, **53**, 1–17.

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.