

Cox Proportional-Hazards Regression for Survival Data

Appendix to *An R and S-PLUS Companion to Applied Regression*

John Fox

15 June 2008 (small corrections)

1 Introduction

Survival analysis examines and models the time it takes for events to occur. The prototypical such event is death, from which the name ‘survival analysis’ and much of its terminology derives, but the ambit of application of survival analysis is much broader. Essentially the same methods are employed in a variety of disciplines under various rubrics — for example, ‘event-history analysis’ in sociology. In this appendix, therefore, terms such as *survival* are to be understood generically.

Survival analysis focuses on the distribution of survival times. Although there are well known methods for estimating unconditional survival distributions, most interesting survival modeling examines the relationship between survival and one or more predictors, usually termed *covariates* in the survival-analysis literature. The subject of this appendix is the Cox proportional-hazards regression model (introduced in a seminal paper by Cox, 1972), a broadly applicable and the most widely used method of survival analysis. Although I will not discuss them here, the `survival` library in R and S-PLUS also contains all of the other commonly employed tools of survival analysis.¹

As is the case for the other appendices to *An R and S-PLUS Companion to Applied Regression*, I assume that you have read the main text and are therefore familiar with S. In addition, I assume familiarity with Cox regression. I nevertheless begin with a review of basic concepts, primarily to establish terminology and notation. The second section of the appendix takes up the Cox proportional-hazards model with time-independent covariates. Time-dependent covariates are introduced in the third section. A fourth and final section deals with diagnostics.

There are many texts on survival analysis: Cox and Oakes (1984) is a classic (if now slightly dated) source. Allison (1995) presents a highly readable introduction to the subject based on the SAS statistical package, but nevertheless of general interest. The major example in this appendix is adapted from Allison. A book by Therneau and Grambsch (2000) is also worthy of mention here because Therneau is the author of the `survival` library for S. Extensive documentation for the `survival` library may be found in Therneau (1999).

2 Basic Concepts and Notation

Let T represent survival time. We regard T as a random variable with cumulative distribution function $P(t) = \Pr(T \leq t)$ and probability density function $p(t) = dP(t)/dt$. The more optimistic *survival function* $S(t)$ is the complement of the distribution function, $S(t) = \Pr(T > t) = 1 - P(t)$. A fourth representation of the distribution of survival times is the *hazard function*, which assesses the instantaneous risk of demise

¹The `survival` library is a standard part of S-PLUS 2000 and 6.0 for Windows, and need not be attached via the `library` function. In R, the `survival` library is among the recommended packages, and is included with the standard Windows distribution; it must be attached prior to use, however.

at time t , conditional on survival to that time:

$$\begin{aligned} h(t) &= \lim_{\Delta t \rightarrow 0} \frac{\Pr[(t \leq T < t + \Delta t) | T \geq t]}{\Delta t} \\ &= \frac{f(t)}{S(t)} \end{aligned}$$

Modeling of survival data usually employs the hazard function or the log hazard. For example, assuming a constant hazard, $h(t) = \nu$, implies an exponential distribution of survival times, with density function $p(t) = \nu e^{-\nu t}$. Other common hazard models include

$$\log h(t) = \nu + \rho t$$

which leads to the *Gompertz distribution* of survival times, and

$$\log h(t) = \nu + \rho \log(t)$$

which leads to the *Weibull distribution* of survival times. (See, for example, Cox and Oakes, 1984: Sec. 2.3, for these and other possibilities.) In both the Gompertz and Weibull distributions, the hazard can either increase or decrease with time; moreover, in both instances, setting $\rho = 0$ yields the exponential model.

A nearly universal feature of survival data is *censoring*, the most common form of which is *right-censoring*: Here, the period of observation expires, or an individual is removed from the study, before the event occurs — for example, some individuals may still be alive at the end of a clinical trial, or may drop out of the study for various reasons other than death prior to its termination. An observation is *left-censored* if its initial time at risk is unknown. Indeed, the same observation may be both right and left-censored, a circumstance termed *interval-censoring*. Censoring complicates the likelihood function, and hence the estimation, of survival models.

Moreover, conditional on the value of any covariates in a survival model and on an individual's survival to a particular time, censoring must be independent of the future value of the hazard for the individual. If this condition is not met, then estimates of the survival distribution can be seriously biased. For example, if individuals tend to drop out of a clinical trial shortly before they die, and therefore their deaths go unobserved, survival time will be over-estimated. Censoring that meets this requirement is *noninformative*. A common instance of noninformative censoring occurs when a study terminates at a predetermined date.

3 The Cox Proportional-Hazards Model

As mentioned, survival analysis typically examines the relationship of the survival distribution to covariates. Most commonly, this examination entails the specification of a linear-like model for the log hazard. For example, a parametric model based on the exponential distribution may be written as

$$\log h_i(t) = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$$

or, equivalently,

$$h_i(t) = \exp(\alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik})$$

that is, as a linear model for the log-hazard or as a multiplicative model for the hazard. Here, i is a subscript for observation, and the x 's are the covariates. The constant α in this model represents a kind of log-baseline hazard, since $\log h_i(t) = \alpha$ [or $h_i(t) = e^\alpha$] when all of the x 's are zero. There are similar parametric regression models based on the other survival distributions described in the preceding section.²

The *Cox model*, in contrast, leaves the baseline hazard function $\alpha(t) = \log h_0(t)$ unspecified:

$$\log h_i(t) = \alpha(t) + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$$

²The `survreg` function in the `survival` library fits the exponential model and other parametric *accelerated failure time* models. Because the Cox model is now used much more frequently than parametric survival regression models, I will not describe `survreg` in this appendix. Enter `help(survreg)` and see Therneau (1999) for details.

or, again equivalently,

$$h_i(t) = h_0(t) \exp(\beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik})$$

This model is *semi-parametric* because while the baseline hazard can take any form, the covariates enter the model linearly. Consider, now, two observations i and i' that differ in their x -values, with the corresponding linear predictors

$$\eta_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$$

and

$$\eta_{i'} = \beta_1 x_{i'1} + \beta_2 x_{i'2} + \cdots + \beta_k x_{i'k}$$

The hazard ratio for these two observations,

$$\begin{aligned} \frac{h_i(t)}{h_{i'}(t)} &= \frac{h_0(t)e^{\eta_i}}{h_0(t)e^{\eta_{i'}}} \\ &= \frac{e^{\eta_i}}{e^{\eta_{i'}}} \end{aligned}$$

is independent of time t . Consequently, the Cox model is a *proportional-hazards model*.

Remarkably, even though the baseline hazard is unspecified, the Cox model can still be estimated by the *method of partial likelihood*, developed by Cox (1972) in the same paper in which he introduced the Cox model. Although the resulting estimates are not as efficient as maximum-likelihood estimates for a correctly specified parametric hazard regression model, not having to make arbitrary, and possibly incorrect, assumptions about the form of the baseline hazard is a compensating virtue of Cox's specification. Having fit the model, it is possible to extract an estimate of the baseline hazard (see below).

3.1 The coxph Function

The Cox proportional-hazards regression model is fit in S with the `coxph` function (located in the `survival` library in R):

```
> library(survival) # R only
> args(coxph)
function (formula = formula(data), data = sys.frame(sys.parent()),
        weights, subset, na.action, init, control, method = c("efron",
        "breslow", "exact"), singular.ok = T, robust = F, model = F,
        x = F, y = T, ...)
NULL
```

Most of the arguments to `coxph`, including `formula`, `data`, `weights`, `subset`, `na.action`, `singular.ok`, `model`, `x` and `y`, are familiar from `lm` (see Chapter 4 of the text, especially Section 4.7), although the `formula` argument requires special consideration: The right-hand side of the model formula for `coxph` is the same as for a linear model.³ The left-hand side is a survival object, created by the `Surv` function. In the simple case of right-censored data, the call to `Surv` takes the form `Surv(time, event)`, where `time` is either the event time or the censoring time, and `event` is a dummy variable coded 1 if the event is observed or 0 if the observation is censored. See the on-line help for `Surv` for other possibilities.

Among the remaining arguments to `coxph`:

- `init` (initial values) and `control` are technical arguments: See the on-line help for `coxph` for details.
- `method` indicates how to handle observations that have tied (i.e., identical) survival times. The default "efron" method is generally preferred to the once-popular "breslow" method; the "exact" method is much more computationally intensive.

³There are, however, special functions `cluster` and `strata` that may be included on the right side of the model formula: The `cluster` function is used to specify non-independent observations (such as several individuals in the same family); the `strata` function may be used to divide the data into sub-groups with potentially different baseline hazard functions, as explained in Section 5.1.

- If `robust` is `TRUE`, `coxph` calculates robust coefficient-variance estimates. The default is `FALSE`, unless the model includes non-independent observations, specified by the `cluster` function in the model formula. I do not describe Cox regression for clustered data in this appendix.

3.2 An Illustration: Recidivism

The file `Rossi.txt` contains data from an experimental study of recidivism of 432 male prisoners, who were observed for a year after being released from prison (Rossi, Berk, and Lenihan, 1980).⁴ The following variables are included in the data; the variable names are those used by Allison (1995), from whom this example and variable descriptions are adapted:

- `week`: week of first arrest after release, or censoring time.
- `arrest`: the event indicator, equal to 1 for those arrested during the period of the study and 0 for those who were not arrested.
- `fin`: a dummy variable, equal to 1 if the individual received financial aid after release from prison, and 0 if he did not; financial aid was a randomly assigned factor manipulated by the researchers.
- `age`: in years at the time of release.
- `race`: a dummy variable coded 1 for blacks and 0 for others.
- `wexp`: a dummy variable coded 1 if the individual had full-time work experience prior to incarceration and 0 if he did not.
- `mar`: a dummy variable coded 1 if the individual was married at the time of release and 0 if he was not.
- `paro`: a dummy variable coded 1 if the individual was released on parole and 0 if he was not.
- `prio`: number of prior convictions.
- `educ`: education, a categorical variable, with codes 2 (grade 6 or less), 3 (grades 6 through 9), 4 (grades 10 and 11), 5 (grade 12), or 6 (some post-secondary).
- `emp1 – emp52`: dummy variables coded 1 if the individual was employed in the corresponding week of the study and 0 otherwise.

After changing to the directory containing the data, I read the data file into a data frame, and print the first few observations (omitting the variables `emp1 – emp52`, which are in columns 11–62 of the data frame):

```
> Rossi <- read.table('Rossi.txt', header=T)
> Rossi[1:5, 1:10]
  week arrest fin age race wexp mar paro prio educ
1    20     1  0  27   1    0  0    1    3    3
2    17     1  0  18   1    0  0    1    8    4
3    25     1  0  19   0    1  0    1   13    3
4    52     0  1  23   1    1  1    1    1    5
5    52     0  0  19   0    1  0    1    3    3
```

Thus, for example, the first individual was arrested in week 20 of the study, while the fourth individual was never rearrested, and hence has a censoring time of 52.

Following Allison, a Cox regression of time to rearrest on the time-constant covariates is specified as follows:

⁴The data file `Rossi.txt` is available at <http://socserv.mcmaster.ca/jfox/Books/Companion/Rossi.txt>.

```
> mod.allison <- coxph(Surv(week, arrest) ~ fin + age + race + wexp + mar + paro + prio,
+   data=Rossi)
> mod.allison
Call:
coxph(formula = Surv(week, arrest) ~ fin + age + race + wexp +
      mar + paro + prio, data = Rossi)
```

	coef	exp(coef)	se(coef)	z	p
fin	-0.3794	0.684	0.1914	-1.983	0.0470
age	-0.0574	0.944	0.0220	-2.611	0.0090
race	0.3139	1.369	0.3080	1.019	0.3100
wexp	-0.1498	0.861	0.2122	-0.706	0.4800
mar	-0.4337	0.648	0.3819	-1.136	0.2600
paro	-0.0849	0.919	0.1958	-0.434	0.6600
prio	0.0915	1.096	0.0286	3.195	0.0014

Likelihood ratio test=33.3 on 7 df, p=2.36e-05 n= 432

The summary method for Cox models produces a more complete report:

```
> summary(mod.allison)
Call:
coxph(formula = Surv(week, arrest) ~ fin + age + race + wexp +
      mar + paro + prio, data = Rossi)
```

n= 432

	coef	exp(coef)	se(coef)	z	p
fin	-0.3794	0.684	0.1914	-1.983	0.0470
age	-0.0574	0.944	0.0220	-2.611	0.0090
race	0.3139	1.369	0.3080	1.019	0.3100
wexp	-0.1498	0.861	0.2122	-0.706	0.4800
mar	-0.4337	0.648	0.3819	-1.136	0.2600
paro	-0.0849	0.919	0.1958	-0.434	0.6600
prio	0.0915	1.096	0.0286	3.195	0.0014

	exp(coef)	exp(-coef)	lower .95	upper .95
fin	0.684	1.461	0.470	0.996
age	0.944	1.059	0.904	0.986
race	1.369	0.731	0.748	2.503
wexp	0.861	1.162	0.568	1.305
mar	0.648	1.543	0.307	1.370
paro	0.919	1.089	0.626	1.348
prio	1.096	0.913	1.036	1.159

```
Rsquare= 0.074 (max possible= 0.956 )
Likelihood ratio test= 33.3 on 7 df, p=2.36e-05
Wald test = 32.1 on 7 df, p=3.86e-05
Score (logrank) test = 33.5 on 7 df, p=2.11e-05
```

- The column marked **z** in the output records the ratio of each regression coefficient to its standard error, a Wald statistic which is asymptotically standard normal under the hypothesis that the corresponding β is zero. The covariates **age** and **prio** (prior convictions) have highly statistically significant coefficients, while the coefficient for **fin** (financial aid — the focus of the study) is marginally significant.

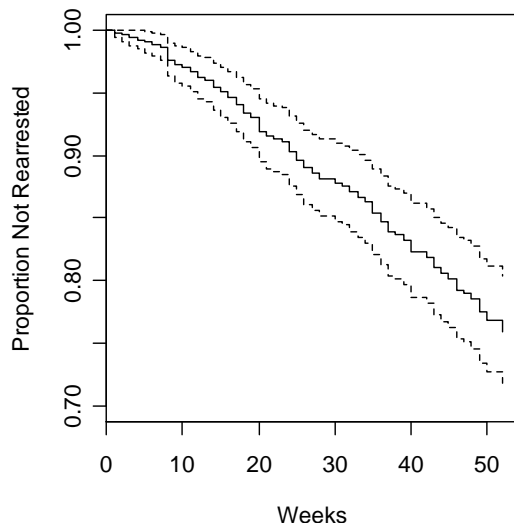


Figure 1: Estimated survival function $\hat{S}(t)$ for the Cox regression of time to rearrest on several predictors. The broken lines show a point-wise 95-percent confidence envelope around the survival function.

- The exponentiated coefficients in the second column of the first panel (and in the first column of the second panel) of the output are interpretable as multiplicative effects on the hazard. Thus, for example, holding the other covariates constant, an additional year of age reduces the weekly hazard of rearrest by a factor of $e^{b_2} = 0.944$ on average — that is, by 5.6 percent. Similarly, each prior conviction increases the hazard by a factor of 1.096, or 9.6 percent.
- The likelihood-ratio, Wald, and score chi-square statistics at the bottom of the output are asymptotically equivalent tests of the omnibus null hypothesis that all of the β 's are zero. In this instance, the test statistics are in close agreement, and the hypothesis is soundly rejected.

Having fit a Cox model to the data, it is often of interest to examine the estimated distribution of survival times. The `survfit` function estimates $S(t)$, by default at the mean values of the covariates. The `plot` method for objects returned by `survfit` graphs the estimated survival function, along with a point-wise 95-percent confidence band. For example, for the model just fit to the recidivism data:

```
> plot(survfit(mod.allison), ylim=c(.7, 1), xlab='Weeks',
+      ylab='Proportion Not Rearrested')
>
```

This command produces Figure 1. [The limits for the vertical axis, set by `ylim=c(.7, 1)`, were selected after examining an initial plot.]

Even more cogently, we may wish to display how estimated survival depends upon the value of a covariate. Because the principal purpose of the recidivism study was to assess the impact of financial aid on rearrest, let us focus on this covariate. I construct a new data frame with two rows, one for each value of `fin`; the other covariates are fixed to their average values. (In the case of a dummy covariate, such as `race`, the average value is the proportion coded 1 in the data set — in the case of `race`, the proportion of blacks). This data frame is passed to `survfit` via the `newdata` argument:

```
> attach(Rossi)
> Rossi.fin <- data.frame(fin=c(0,1), age=rep(mean(age),2), race=rep(mean(race),2),
```

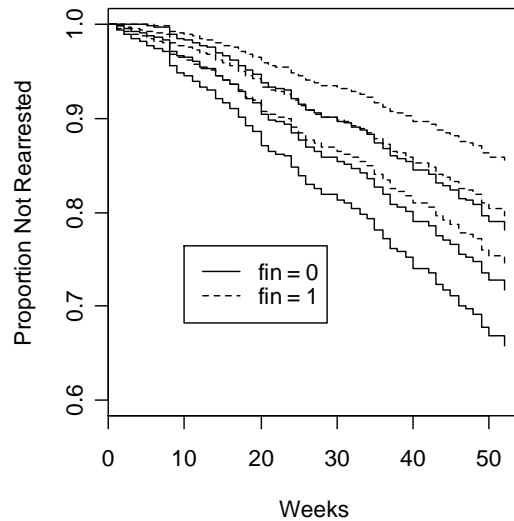


Figure 2: Estimated survival functions for those receiving (`fin = 1`) and not receiving (`fin = 0`) financial aid. Other covariates are fixed at their average values. Each estimate is accompanied by a point-wise 95-percent confidence envelope.

```

+      wexp=rep(mean(wexp),2), mar=rep(mean(mar),2), paro=rep(mean(paro),2),
+      prio=rep(mean(prio),2))
> detach()
> plot(survfit(mod.allison, newdata=Rossi.fin), conf.int=T,
+      lty=c(1,2), ylim=c(.6, 1))
> legend(locator(1), legend=c('fin = 0', 'fin = 1'), lty=c(1,2))
>

```

I specified two additional arguments to `plot`: `lty=c(1,2)` indicates that the survival function for the first group (i.e., for `fin = 0`) will be plotted with a solid line, while that for the second group (`fin = 1`) will be plotted with a broken line; `conf.int=T` requests that confidence envelopes be drawn around each estimated survival function (which is not the default when more than one survival function is plotted). Notice, as well, the use of the `legend` function (along with `locator`) to place a legend on the plot: Click the left mouse button to position the legend.⁵ The resulting graph, which appears in Figure 2, shows the higher estimated ‘survival’ of those receiving financial aid, but the two confidence envelopes overlap substantially, even after 52 weeks.

4 Time-Dependent Covariates

The `coxph` function handles time-dependent covariates by requiring that each time period for an individual appear as a separate observation — that is, as a separate row (or *record*) in the data set. Consider, for example, the `Rossi` data frame, and imagine that we want to treat weekly employment as a time-dependent predictor of time to rearrest. As is often the case, however, the data for each individual appears as a single row, with the weekly employment indicators as 52 columns in the data frame, with names `emp1` through `emp52`; for example, for the first person in the study:

⁵The `plot` method for `survfit` objects can also draw a legend on the plot, but separate use of the `legend` function provides greater flexibility. Legends, line types, and other aspects of constructing graphs in S are described in Chapter 7 of the text.

```

> Rossi[1,]
  week arrest fin age race wexp mar paro prio educ emp1 emp2
1   20     1  0  27   1   0   0   1   3   3   0   0
  emp3 emp4 emp5 emp6 emp7 emp8 emp9 emp10 emp11 emp12 emp13
1    0   0   0   0   0   0   0   0   0   0   0   0
  emp14 emp15 emp16 emp17 emp18 emp19 emp20 emp21 emp22 emp23
1    0   0   0   0   0   0   0   0   NA   NA   NA
  emp24 emp25 emp26 emp27 emp28 emp29 emp30 emp31 emp32 emp33
1   NA   NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
  emp34 emp35 emp36 emp37 emp38 emp39 emp40 emp41 emp42 emp43
1   NA   NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
  emp44 emp45 emp46 emp47 emp48 emp49 emp50 emp51 emp52
1   NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
>

```

Notice that the employment indicators are missing after week 20, when individual 1 was rearrested.

To put the data in the requisite form, we need to write one row for each non-missing period of observation. Here is a simple sequence of commands that accomplishes this purpose:

- First, noting that the employment indicators are in columns 11 through 62, I calculate the number of non-missing records that will be required:

```

> sum(!is.na(Rossi[,11:62])) # record count
[1] 19809

```

- Next, I create a matrix, initially filled with 0's, to hold the data. This matrix has 19,809 rows, one for each record, and 14 columns, to contain the first 10 variables in the original data frame; the `start` and `stop` times of each weekly record; a time-dependent indicator variable (`arrest.time`) set to 1 if a rearrest occurs during the current week, or 0 otherwise; and another indicator (`employed`) set to 1 if the individual was employed during the current week, and 0 if he was not. This last variable is the time-dependent covariate. If there were more than one time-dependent covariate, then there would be a column in the new data set for each.

```

> Rossi.2 <- matrix(0, 19809, 14) # to hold new data set
> colnames(Rossi.2) <- c('start', 'stop', 'arrest.time', names(Rossi)[1:10], 'employed')
>

```

- Finally, I loop over the observations in the original data set, and over the weeks of the year within each observation, to construct the new data set.⁶

```

> row <- 0 # set record counter to 0
> for (i in 1:nrow(Rossi)) { # loop over individuals
+   for (j in 11:62) {      # loop over 52 weeks
+     if (is.na(Rossi[i, j])) next # skip missing data
+     else {
+       row <- row + 1      # increment row counter
+       start <- j - 11     # start time (previous week)
+       stop <- start + 1   # stop time (current week)
+       arrest.time <- if (stop == Rossi[i, 1] && Rossi[i, 2] == 1) 1 else 0
+       # construct record:
+       Rossi.2[row,] <- c(start, stop, arrest.time, unlist(Rossi[i, c(1:10, j)]))
+     }
+   }

```

⁶Programming constructs such as `for` loops are described in Chapter 8 of the text.


```

+     }
> Rossi.2 <- as.data.frame(Rossi.2)
> remove(i, j, row, start, stop, arrest.time) # clean up
>

```

This procedure is very inefficient computationally, taking more than four minutes under R 1.4.1 on my Windows 2000 computer, which has an 800 MHz processor and plenty of memory. But the programming was very simple, requiring perhaps five minutes to write and debug: A time expenditure of about 10 minutes is insignificant in preparing data for analysis.⁷

If, however, we often want to perform these operations, it makes sense to encapsulate them in a function, and to spend some programming time to make the computation more efficient. I have written such a function, named `fold`; the function is included with the script file for this appendix, and takes the following arguments:

- **data**: A data frame or numeric matrix (with column names) to be ‘folded.’ For reasons of efficiency, if there are factors in **data** these will be converted to numeric variables in the output data frame.
- **time**: The quoted name of the event/censoring-time variable in **data**.
- **event**: The quoted name of the event/censoring indicator variable in **data**.
- **cov**: A vector giving the column numbers of the time-dependent covariate in **data**, or a list of vectors if there is more than one time-dependent covariate.
- **cov.names**: A character string or character vector giving the name(s) to be assigned to the time-dependent covariate(s) in the output data set.
- **suffix**: The suffix to be attached to the name of the time-to-event variable in the output data set; defaults to `’.time’`.
- **cov.times**: The observation times for the covariate values, including the start time. This argument can take several forms:
 - The default is the vector of integers from 0 to the number of covariate values (i.e., containing one more entry — the initial time of observation — than the length of each vector in **cov**).
 - An arbitrary numerical vector with one more entry than the length of each vector in **cov**.
 - The columns in the input data set that give the (potentially different) covariate observation times for each individual. There should be one more column than the length of each vector in **cov**.
- **common.times**: A logical value indicating whether the times of observation are the same for all individuals; defaults to `TRUE`.
- **lag**: Number of observation periods to lag each value of the time-dependent covariate(s); defaults to 0. The use of **lag** is described later in this section.

To create the same data set as above using `fold`, I enter:

```

> Rossi.2 <- fold(Rossi, time='week',
+   event='arrest', cov=11:62, cov.names='employed')
>

```

This command required less than 16 seconds on my computer — still not impressively efficient, but (not counting programming effort) much better than the brute-force approach that I took previously. Here are the first 50 of the nearly 20,000 records in the data frame `Rossi.2`:

⁷See the discussion of ‘quick and dirty’ programming in Chapter 8 of the text.

```

> Rossi.2[1:50,]
      start stop arrest.time week arrest fin age race wexp mar paro prio educ employed
1.1      0   1         0    20     1  0  27   1  0  0   1   3   3     0
1.2      1   2         0    20     1  0  27   1  0  0   1   3   3     0
. . . .
1.19     18  19         0    20     1  0  27   1  0  0   1   3   3     0
1.20     19  20         1    20     1  0  27   1  0  0   1   3   3     0
2.1      0   1         0    17     1  0  18   1  0  0   1   8   4     0
2.2      1   2         0    17     1  0  18   1  0  0   1   8   4     0
. . . .
2.16     15  16         0    17     1  0  18   1  0  0   1   8   4     0
2.17     16  17         1    17     1  0  18   1  0  0   1   8   4     0
3.1      0   1         0    25     1  0  19   0  1  0   1  13   3     0
3.2      1   2         0    25     1  0  19   0  1  0   1  13   3     0
. . . .
3.13     12  13         0    25     1  0  19   0  1  0   1  13   3     0

```

Once the data set is constructed, it is simple to use `coxph` to fit a model with time-dependent covariates. The right-hand-side of the model is essentially the same as before, but both the start and end times of each interval are specified in the call to `Surv`, in the form `Surv(start, stop, event)`. Here, `event` is the time-dependent version of the event indicator variable, equal to 1 only in the time-period during which the event occurs. For the example:

```

> mod.allison.2 <- coxph(Surv(start, stop, arrest.time) ~
+   fin + age + race + wexp + mar + paro + prio + employed,
+   data=Rossi.2)
> summary(mod.allison.2)
Call:
coxph(formula = Surv(start, stop, arrest.time) ~ fin + age +
      race + wexp + mar + paro + prio + employed, data = Rossi.2)

```

n= 19809

	coef	exp(coef)	se(coef)	z	p
fin	-0.3567	0.700	0.1911	-1.866	6.2e-02
age	-0.0463	0.955	0.0217	-2.132	3.3e-02
race	0.3387	1.403	0.3096	1.094	2.7e-01
wexp	-0.0256	0.975	0.2114	-0.121	9.0e-01
mar	-0.2937	0.745	0.3830	-0.767	4.4e-01
paro	-0.0642	0.938	0.1947	-0.330	7.4e-01
prio	0.0851	1.089	0.0290	2.940	3.3e-03
employed	-1.3282	0.265	0.2507	-5.298	1.2e-07

	exp(coef)	exp(-coef)	lower .95	upper .95
fin	0.700	1.429	0.481	1.018
age	0.955	1.047	0.915	0.996
race	1.403	0.713	0.765	2.574
wexp	0.975	1.026	0.644	1.475
mar	0.745	1.341	0.352	1.579
paro	0.938	1.066	0.640	1.374
prio	1.089	0.918	1.029	1.152
employed	0.265	3.774	0.162	0.433

Rsquare= 0.003 (max possible= 0.066)

```

Likelihood ratio test= 68.7  on 8 df,  p=9.11e-12
Wald test             = 56.1  on 8 df,  p=2.63e-09
Score (logrank) test = 64.5  on 8 df,  p=6.1e-11

```

Setting Up Time-Dependent Data in S-PLUS

The operations described above work in S-PLUS as well as in R. An exception is the `remove` command: In S-PLUS `remove` takes a character vector of names as its argument; moreover, the loop indices `i` and `j` are not global variables in S-PLUS, and hence do not have to be removed.

There are, however, differences in efficiency: In S-PLUS 2000 for Windows, the brute-force approach to constructing the time-dependent data set required nearly nine minutes to execute, while the `fold` function required 27 seconds; in S-PLUS 6.0 for Windows, the first operation required a little more than six minutes, while the second consumed more than two minutes. All of these timings were on the same 800 MHz Windows PC used for R.

4.1 Lagged Covariates

The time-dependent employment covariate therefore has an apparently large effect: The hazard of rearrest was smaller by a factor of $e^{-1.3282} = 0.265$ (i.e., a decline of 73.5 percent) during a week in which the former inmate was employed. As Allison (1995) points out, however, the direction of causality here is ambiguous, since a person cannot work when he is in jail. One way of addressing this problem is to use instead a lagged value of employment, from the previous week for example. The `fold` function can easily provide lagged time-dependent covariates:

```

> Rossi.3 <- fold(Rossi, 'week', 'arrest', 11:62, 'employed', lag=1)
> mod.allison.3 <- coxph(Surv(start, stop, arrest.time) ~
+   fin + age + race + wexp + mar + paro + prio + employed,
+   data=Rossi.3)
> summary(mod.allison.3)

```

Call:

```

coxph(formula = Surv(start, stop, arrest.time) ~ fin + age +
      race + wexp + mar + paro + prio + employed, data = Rossi.3)

```

n= 19377

	coef	exp(coef)	se(coef)	z	p
fin	-0.3513	0.704	0.1918	-1.831	0.06700
age	-0.0498	0.951	0.0219	-2.274	0.02300
race	0.3215	1.379	0.3091	1.040	0.30000
wexp	-0.0477	0.953	0.2132	-0.223	0.82000
mar	-0.3448	0.708	0.3832	-0.900	0.37000
paro	-0.0471	0.954	0.1963	-0.240	0.81000
prio	0.0920	1.096	0.0288	3.195	0.00140
employed	-0.7869	0.455	0.2181	-3.608	0.00031

	exp(coef)	exp(-coef)	lower .95	upper .95
fin	0.704	1.421	0.483	1.025
age	0.951	1.051	0.911	0.993
race	1.379	0.725	0.752	2.528
wexp	0.953	1.049	0.628	1.448

mar	0.708	1.412	0.334	1.501
paro	0.954	1.048	0.649	1.402
prio	1.096	0.912	1.036	1.160
employed	0.455	2.197	0.297	0.698

```

Rsquare= 0.002 (max possible= 0.067 )
Likelihood ratio test= 47.2 on 8 df, p=1.43e-07
Wald test = 43.4 on 8 df, p=7.47e-07
Score (logrank) test = 46.4 on 8 df, p=1.99e-07

```

The coefficient for the now-lagged employment indicator is still highly statistically significant, but the estimated effect of employment is much smaller: $e^{-0.7869} = 0.455$ (or a decrease of 54.5 percent).

5 Model Diagnostics

As is the case for a linear or generalized linear model (see Chapter 6 of the text), it is desirable to determine whether a fitted Cox regression model adequately describes the data. I will briefly consider three kinds of diagnostics: for violation of the assumption of proportional hazards; for influential data; and for nonlinearity in the relationship between the log hazard and the covariates. All of these diagnostics use the **residuals** method for **coxph** objects, which calculates several kinds of residuals (along with some quantities that are not normally thought of as residuals). Details are in Therneau (1999).

5.1 Checking Proportional Hazards

Tests and graphical diagnostics for proportional hazards may be based on the *scaled Schoenfeld residuals*; these can be obtained directly as **residuals(model, "scaledsch")**, where **model** is a **coxph** model object. The matrix returned by **residuals** has one column for each covariate in the model. More conveniently, the **cox.zph** function calculates tests of the proportional-hazards assumption for each covariate, by correlating the corresponding set of scaled Schoenfeld residuals with a suitable transformation of time [the default is based on the *Kaplan-Meier estimate* of the survival function, $K(t)$].

I will illustrate these tests with a scaled-down version of the Cox regression model fit to the recidivism data in Section 3.2, eliminating the covariates whose coefficients were not statistically significant:⁸

```

> mod.allison.4 <- coxph(Surv(week, arrest) ~ fin + age + prio,
+   data=Rossi)
> mod.allison.4
Call:
coxph(formula = Surv(week, arrest) ~ fin + age + prio, data = Rossi)

```

	coef	exp(coef)	se(coef)	z	p
fin	-0.3469	0.707	0.1902	-1.82	0.06800
age	-0.0671	0.935	0.0209	-3.22	0.00130
prio	0.0969	1.102	0.0273	3.56	0.00038

```

Likelihood ratio test=29.1 on 3 df, p=2.19e-06 n= 432

```

Note that the coefficient for financial aid, which is the focus of the study, now has a two-sided p -value greater than .05; a one-sided test is appropriate here, however, since we expect the coefficient to be negative, so there is still marginal evidence for the effect of this covariate on the time of rearrest.

As mentioned, tests for the proportional-hazards assumption are obtained from **cox.zph**, which computes a test for each covariate, along with a global test for the model as a whole:

⁸It is possible that a covariate that is not statistically significant when its effect is, in essence, *averaged* over time nevertheless has a statistically significant *interaction* with time, which manifests itself as nonproportional hazards. I leave it to the reader to check for this possibility using the model fit originally to the recidivism data.

```
> cox.zph(mod.allison.4)
      rho  chisq    p
fin   -0.00657 0.00507 0.9432
age   -0.20976 6.54118 0.0105
prio  -0.08003 0.77263 0.3794
GLOBAL      NA 7.12999 0.0679
```

There is, therefore, strong evidence of non-proportional hazards for age, while the global test (on 3 degrees of freedom) is not quite statistically significant. These tests are sensitive to linear trends in the hazard.

Plotting the object returned by `cox.zph` produces graphs of the scaled Schoenfeld residuals against transformed time (see Figure 3):

```
> par(mfrow=c(2,2))
> plot(cox.zph(mod.allison.4))
Warning messages:
1: Collapsing to unique x values in: approx(xx, xtime,
      seq(min(xx), max(xx), length = 17)[2 * (1:8)])
2: Collapsing to unique x values in: approx(xtime, xx, temp)
```

Interpretation of these graphs is greatly facilitated by smoothing, for which purpose `cox.zph` uses a smoothing spline, shown on each graph by a solid line; the broken lines represent ± 2 -standard-error envelopes around the fit. Systematic departures from a horizontal line are indicative of non-proportional hazards. The assumption of proportional hazards appears to be supported for the covariates `fin` (which is, recall, a dummy variable, accounting for the two bands in the graph) and `prio`, but there appears to be a trend in the plot for `age`, with the `age` effect declining with time; this effect was detected in the test reported above.

One way of accommodating non-proportional hazards is to build interactions between covariates and time into the Cox regression model; such interactions are themselves time-dependent covariates. For example, based on the diagnostics just examined, it seems reasonable to consider a linear interaction of time and `age`; using the previously constructed `Rossi.2` data frame:

```
> mod.allison.5 <- coxph(Surv(start, stop, arrest.time) ~
+   fin + age + age:stop + prio,
+   data=Rossi.2)
> mod.allison.5
Call:
coxph(formula = Surv(start, stop, arrest.time) ~ fin + age +
      age:stop + prio, data = Rossi.2)
```

	coef	exp(coef)	se(coef)	z	p
fin	-0.34855	0.706	0.19023	-1.832	0.06700
age	0.03219	1.033	0.03943	0.817	0.41000
prio	0.09820	1.103	0.02726	3.603	0.00031
age:stop	-0.00383	0.996	0.00147	-2.608	0.00910

```
Likelihood ratio test=36 on 4 df, p=2.85e-07 n= 19809
```

As expected, the coefficient for the interaction is negative and highly statistically significant: The effect of `age` declines with time.⁹ Notice that the model does not require a ‘main-effect’ term for `stop` (i.e., time); such a term would be redundant, since the time effect is the baseline hazard.

An alternative to incorporating an interaction in the model is to divide the data into *strata* based on the value of one or more covariates. Each stratum is permitted to have a different baseline hazard function, while the coefficients of the remaining covariates are assumed to be constant across strata. An advantage of this

⁹That is, initially, `age` has a positive partial effect on the hazard (given by the `age` coefficient, 0.032), but this effect gets progressively smaller with time (at the rate -0.0038 per week), becoming negative after about 10 weeks.

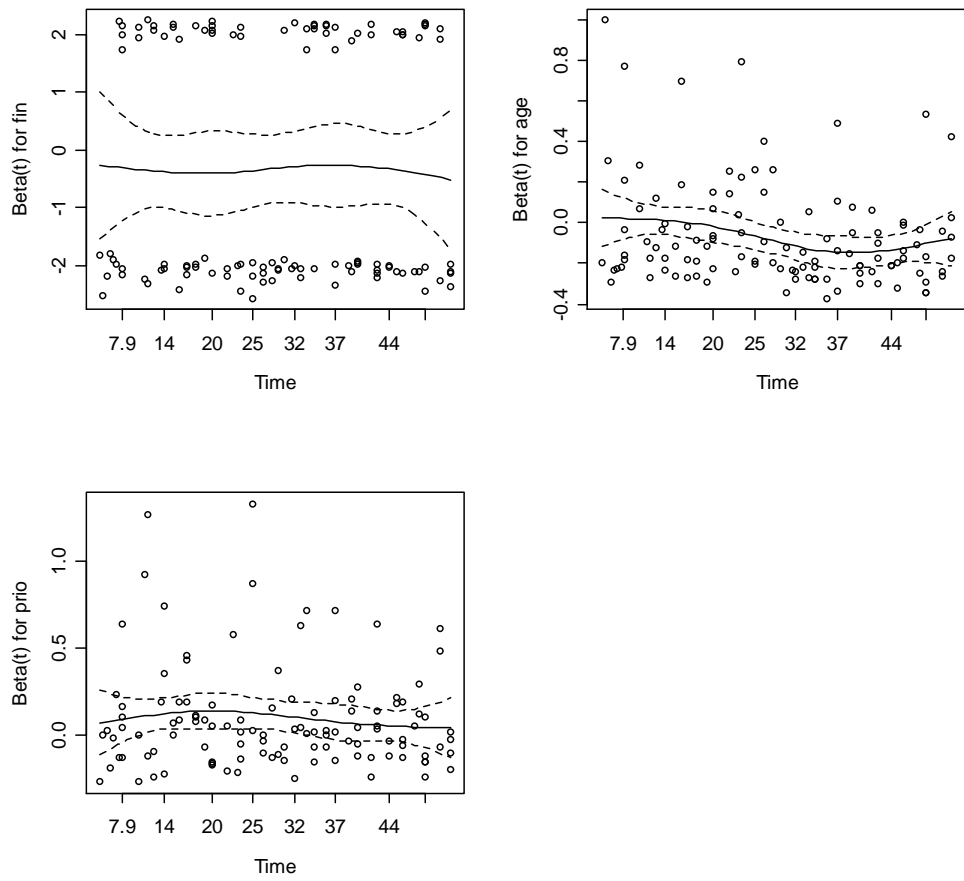


Figure 3: Plots of scaled Schoenfeld residuals against transformed time for each covariate in a model fit to the recidivism data. The solid line is a smoothing-spline fit to the plot, with the broken lines representing a ± 2 -standard-error band around the fit.

approach is that we do not have to assume a particular form of interaction between the stratifying covariates and time. A disadvantage is the resulting inability to examine the effects of the stratifying covariates. Stratification is most natural when a covariate takes on only a few distinct values, and when the effect of the stratifying variable is not of direct interest. In our example, `age` takes on many different values, but we can create categories by arbitrarily dissecting the variable into class intervals. After examining the distribution of `age`, I decided to define four intervals: 19 or younger; 20 to 25; 26 to 30; and 31 or older. I use the `recode` function in the `car` library to categorize `age`:¹⁰

```
> library(car)
. . .
> Rossi$age.cat <- recode(Rossi$age, " lo:19=1; 20:25=2; 26:30=3; 31:hi=4 ")
> table(Rossi$age.cat)

 1  2  3  4
66 236 66 64
```

Most of the individuals in the data set are in the second age category, 20 to 25, but since this is a reasonably narrow range of ages, I did not feel the need to sub-divide the category.

A stratified Cox regression model is fit by including a call to the `strata` function on the right-hand side of the model formula. The arguments to this function are one or more stratifying variables; if there is more than one such variable, then the strata are formed from their cross-classification. In the current illustration, there is just one stratifying variable:

```
> mod.allison.6 <- coxph(Surv(week, arrest) ~
+   fin + prio + strata(age.cat),
+   data=Rossi)
> mod.allison.6
Call:
coxph(formula = Surv(week, arrest) ~ fin + prio + strata(age.cat),
      data = Rossi)
```

	coef	exp(coef)	se(coef)	z	p
fin	-0.341	0.711	0.190	-1.79	0.0730
prio	0.094	1.099	0.027	3.48	0.0005

Likelihood ratio test=13.4 on 2 df, p=0.00122 n= 432

```
> cox.zph(mod.allison.6)
      rho chisq p
fin    -0.0183 0.0392 0.843
prio   -0.0771 0.6858 0.408
GLOBAL      NA 0.7297 0.694
```

There is no evidence of non-proportional hazards for the remaining covariates.

5.2 Influential Observations

Specifying the argument `type=dfbeta` to `residuals` produces a matrix of estimated changes in the regression coefficients upon deleting each observation in turn; likewise, `type=dfbetas` produces the estimated changes in the coefficients divided by their standard errors (cf., Sections 6.1 and 6.6 of the text for similar diagnostics for linear and generalized linear models).

For example, for the model regressing time to rearrest on financial aid, age, and number of prior offenses:

¹⁰An alternative is to use the standard S function `cut: cut(Rossi$age, c(0, 19, 25, 30, Inf))`. See Chapter 2 of the text.

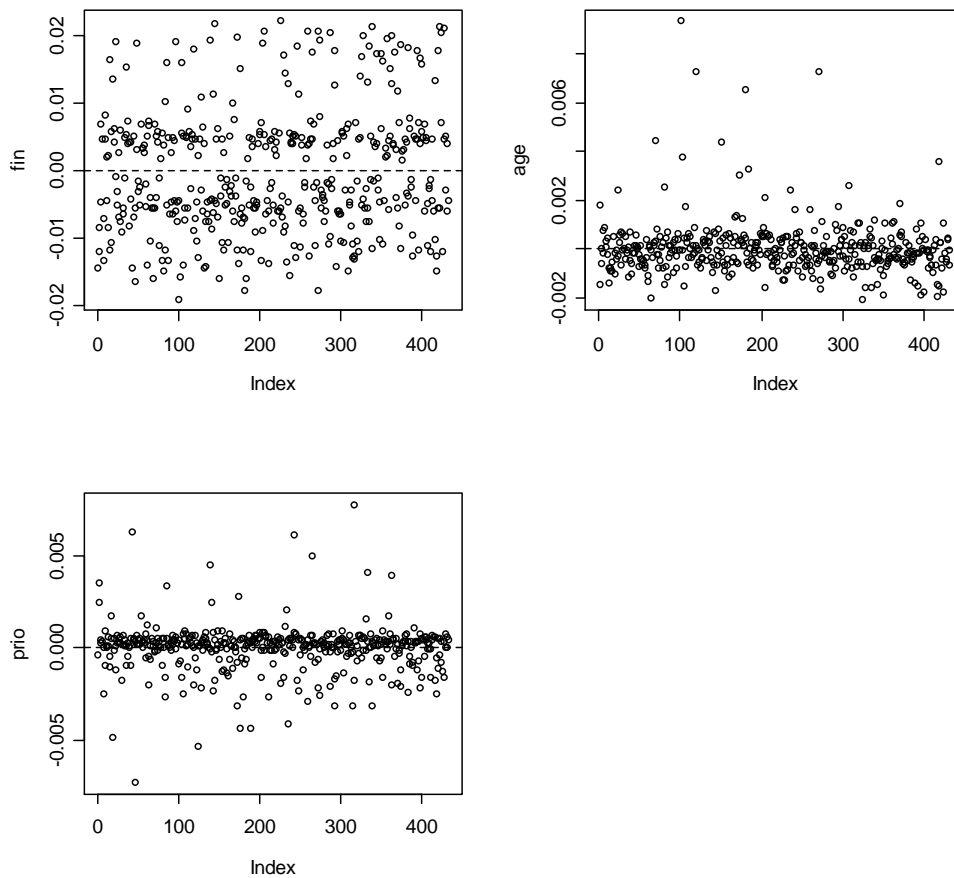


Figure 4: Index plots of dfbeta for the Cox regression of time to rearrest on `fin`, `age`, and `prio`.

```

> dfbeta <- residuals(mod.allison.4, type='dfbeta')
> par(mfrow=c(2,2))
> for (j in 1:3) {
+   plot(dfbeta[,j], ylab=names(coef(mod.allison.4))[j])
+   abline(h=0, lty=2)
+ }
>

```

The index plots produced by these commands appear in Figure 4. Comparing the magnitudes of the largest dfbeta values to the regression coefficients suggests that none of the observations is terribly influential individually (even though some of the dfbeta values for `age` are large compared with the others¹¹).

5.3 Nonlinearity

Nonlinearity — that is, an incorrectly specified functional form in the parametric part of the model — is a potential problem in Cox regression as it is in linear and generalized linear models (see Sections 6.4 and 6.6 of the text). The *martingale residuals* may be plotted against covariates to detect nonlinearity, and may

¹¹As an exercise, the reader may wish to identify these observations and, in particular, examine their ages.

also be used to form component-plus-residual (or partial-residual) plots, again in the manner of linear and generalized linear models.

For the regression of time to rearrest on financial aid, age, and number of prior arrests, let us examine plots of martingale residuals and partial residuals against the last two of these covariates; nonlinearity is not an issue for financial aid, because this covariate is dichotomous:

```
> par(mfrow=c(2,2))
> res <- residuals(mod.allison.4, type='martingale')
> X <- as.matrix(Rossi[,c("age", "prio")]) # matrix of covariates
> par(mfrow=c(2,2))
> for (j in 1:2) { # residual plots
+   plot(X[,j], res, xlab=c("age", "prio")[j], ylab="residuals")
+   abline(h=0, lty=2)
+   lines(lowess(X[,j], res, iter=0))
+ }

> b <- coef(mod.allison.4)[c(2,3)] # regression coefficients
> for (j in 1:2) { # partial-residual plots
+   plot(X[,j], b[j]*X[,j] + res, xlab=c("age", "prio")[j],
+       ylab="component+residual")
+   abline(lm(b[j]*X[,j] + res ~ X[,j]), lty=2)
+   lines(lowess(X[,j], b[j]*X[,j] + res, iter=0))
+ }
>
```

The resulting residual and component-plus-residual plots appear in Figure 5. As in the plots of Schoenfeld residuals, smoothing these plots is also important to their interpretation; The smooths in Figure 5 are produced by local linear regression (using the `lowess` function). Nonlinearity, it appears, is slight here.

References

- Allison, P. D. 1995. *Survival Analysis Using the SAS System: A Practical Guide*. Cary NC: SAS Institute.
- Cox, D. R. 1972. "Regression Models and Life Tables (with Discussion)." *Journal of the Royal Statistical Society, Series B* 34:187–220.
- Cox, D. R. & D. Oakes. 1984. *Analysis of Survival Data*. London: Chapman and Hall.
- Rossi, P. H., R. A. Berk & K. J. Lenihan. 1980. *Money, Work and Crime: Some Experimental Results*. New York: Academic Press.
- Therneau, T. M. 1999. A Package for Survival Analysis in S. Technical Report <<http://www.mayo.edu/hsr/people/therneau/survival.ps>> Mayo Foundation.
- Therneau, T. M. & P. M. Grambsch. 2000. *Modeling Survival Data: Extending the Cox Model*. New York: Springer.

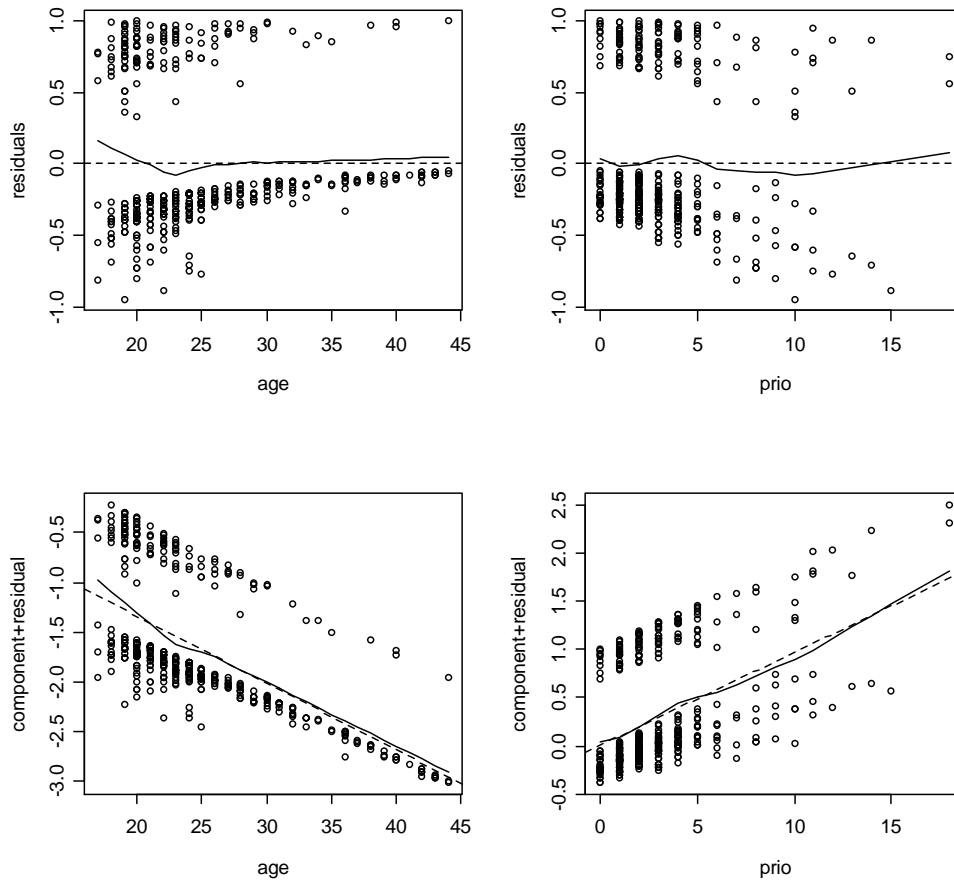


Figure 5: Martingale-residual plots (top) and component-plus-residual plots (bottom) for the covariates **age** and **prio**. The broken lines on the residual plots are at the vertical value 0, and on the component-plus-residual plots are fit by linear least-squares; the solid lines are fit by local linear regression (lowess).