

Notes

Data Example

Using Longitudinal Data

DEM 7223 - Event History Analysis - Cox Proportional Hazards Model Part 1

Corey S. Sparks, PhD (<https://coreysparks.github.io>)

The University of Texas at San Antonio (<https://hcap.utsa.edu/demography>)

September 28, 2020

Notes

Parametric model specifications

When considering a parametric hazard model, we saw that the choice of the specified distribution function is key * If we expect the hazard (or pdf) to take an exponential form, we use that model, same for the Weibull or log-normal, etc.

- So by saying this, we force our data to correspond to the distribution we specify.
- What if, however, the distribution of the durations do not necessarily follow one of the parametric families? We are then left with the most heinous of statistical quandaries: model mis-specification :(
- So in considering the use of hazard models, we need to also consider the case where we cannot (or adequately) specify an appropriate parametric model, this is the reasoning behind the use of Cox's (1972) semi-parametric modeling approach.

The Cox Proportional Hazards Model

- Cox (1972) suggested a more widely applicable model to be used in situations where a suitable parametric distribution is unavailable
- Also, it allows the analyst the freedom to explore the theoretical connections between the covariates and the hazard rate, free of the parametric assumption.
- We are still modeling the effect of individual characteristics on the hazard of an event outcome, in the same way as with the parametric proportional hazards model. We just leave the baseline hazard rate **unspecified in terms of structural parameters**

Model form

The Cox model has the familiar form:

$$h(t) = h_0(t) \exp(x' \beta)$$

- Which is the same form as the parametric proportional hazards models we saw last week.
- The key difference is in the value of h_0
- In the Cox model the baseline hazard rate, h_0 is the observed empirical hazard rate for individuals in the baseline, or reference group of the sample.
- For any two individuals with different values of a covariate, x , the *hazard ratio* is:

$$\frac{h_i(t)}{h_0(t)} = \exp(\beta(x_i - x_0))$$

so when $x_i = 0, x_0 = 0$, the hazard ratio is just $\exp(\beta)$

This is called a *semi-parametric model* because, while the baseline hazard rate h_0 does not have any structural parameters, the regression effects, β are estimated.

Partial likelihood for the Cox Model

- Unlike the parametric models, where we expected the hazard to be a function of time, and where the time between events actually contributed some information to the estimation, the Cox model handles things differently
- The Cox model works, much like the Kaplan-Meier estimator, in terms of *ordered event times*.
- So we actually only get an estimate of the function at the observed failure times, vs the parametric models where we get an estimate of the risk at all possible failure times (which is one of the main reasons for using them!)

Partial likelihood estimation

- First, we sort all observed failure times, in ascending order:

$$t_1 < t_2 < \dots < t_n$$

for all individuals in the data. Now we assume all events have unique durations. In actuality, ties exist and there are a variety of ways to handle these.

- Each observation has its censoring indicator δ_i , which tells if the individual is observed or censored at each time.
- These observations are then modeled in terms of their relative hazards.
- The partial likelihood is constructed by taking the cumulative product of the hazard, for the *Risk set* at time t . The probability that a case j will fail at time t is :

$$Pr(t_j = T_i | R(t_i)) = \frac{h(t_{ij})}{\sum_{j \in R(t_i)} h(t_{ij})}$$

- Where the denominator in this equation is summing over all individuals at risk at time t_i
- The partial likelihood function in terms of the regression parameters is:

$$L_p = \prod_{\delta=1} \frac{h(t_{ij}) \exp(x' \beta)}{\sum_{j \in R(t_i)} h(t_{ij}) \exp(x' \beta)}$$

* Since both the numerator and denominator contain the overall hazard, it cancels, giving:

$$L_p = \prod_{\delta=1} \frac{\exp(x' \beta)}{\sum_{j \in R(t_i)} \exp(x' \beta)}$$

- Which says nothing about the baseline hazard function or its shape.
- By maximizing this partial likelihood, estimates of the β 's are found
- This is called *Maximum Partial likelihood estimation*, and is not a true likelihood.
- This is because we have not directly included the survival times of the censored cases, instead these are handled by modifying the risk set, $R(t_{ij})$, but not explicitly in the numerator
- Much in the way Kaplan-Meier treats censored cases
- Cox in later papers demonstrated that the same properties (efficiency, asymptotic normality) of the estimates still hold.
- This allows us to use our standard likelihood ratio tests, and Wald parameter tests in interpreting and comparing models.
- Ties in the data may be handled by modification of the partial likelihood function
- Ties are simply events that have the same event time and are very common in demographic work
- We have seen this repeatedly in our child mortality and birth interval analyses.
- The likelihood must be modified to incorporate these "simultaneously" occurring event times
- The ability to modify the likelihood function for the Cox model also exemplifies the flexibility of the model over the parametric forms, which are not specified to handle tied event times.
- The big issue with tied observations is related to determining the risk set at a particular time, but also in determining the future risk sets

Handling ties

- Breslow's Method
 - Assumes the same risk set for all tied events
 - This is weakest if there are a large number of ties at a particular time point

- Efron's Method
 - Uses a different risk set at any time point for tied observations
 - This is done by considering all possible orderings of failure times for the tied observations
- Other methods exist, but Efron's method is the most widely used.

Interpreting the Cox model

- We have already seen how the proportional hazards model is generally interpreted from the Weibull and Exponential cases
- This is done via the $\exp(\beta)$, or the hazards ratio
- If the regression coefficient, β , is positive (the hazard is increasing), $\exp(\beta)$ will be >1 and this indicates that an individual with a value of $x=1$ will have a $1 - \exp(\beta)$ higher hazard rate, compared with an individual with $x=0$
- If the regression coefficient, β , is negative (the hazard is decreasing), $\exp(\beta)$ will be <1 and this indicates that an individual with a value of $x=1$ will have a $1 - \exp(\beta)$ lower hazard rate, compared with an individual with $x=0$

Good variable construction habits

- In order to facilitate the interpretability of the model hazard ratios, in demography, we typically create binary dummy variables for things like age via recoding
- i.e. construct a set of dummy variables for 5 year age intervals between 15 and 50 with the reference group being 30-35.

if age \geq 15 & age $<$ 20 age1=1, else age1=0

if age \geq 20 & age $<$ 25 age2=1, else age2=0

if age \geq 25 & age $<$ 30 age3=1, else age3=0

30-35 is reference group without a covariate constructed

if age \geq 35 & age $<$ 40 age4=1, else age4=0

if age \geq 40 & age $<$ 45 age5=1, else age5=0

if age \geq 45 & age $<$ 50 age6=1, else age6=0

if age \geq 50 age7=1, else age7=0

- You could also use a factor variable with the appropriate level as the reference category.
- This approach is also useful for coding incomes, but is done in terms of the income distribution's quantiles

Good variable construction: continuous case

- Often if our covariate is continuous (like weight, height, maybe income if we're treating it that

way) we construct a z-score for the variable

- The z-score is called the standard score, and centers the covariate around it's mean, so the new mean is 0 and each individual's value represents their departure from the mean
- i.e. if a person's weight z-score was -5, then they are 5 pounds below the average weight
- In R, the `scale()` function does this.

Confidence intervals for hazard ratios

- Because the partial likelihood estimates of the β 's have the same asymptotic properties as mle's of β , we can construct $1 - \alpha\%$ confidence intervals for both β and the $\exp(\beta)$, or hazard ratio.
- These are often reported in output in tabels.
- To find the lower 95% ci for $\exp(\beta)$, we do

$$\text{Lower } 1 - \alpha \% \text{CI} = \exp(\beta - z * s. e. (\beta))$$

$$\text{Upper } 1 - \alpha \% \text{CI} = \exp(\beta + z * s. e. (\beta))$$

- For 95% confidence intervals, z would be 1.96

Risk Scores

- Often we are interested in how “at risk” a certain individual is or at least someone with a particular set of covariates
- *Risk scores* represent the linear combination of all the individual's covariates on their hazard
- If none of our covariates vary with time (which we assume for the present), the “risk score” would be:

$$h_i = h_0 * \exp(\beta_1 * x_1 + \beta_2 * x_2 + \dots + \beta_k * x_k)$$

Since our “baseline hazard” is h_0 , the risk score for an individual with a particular set of covariates is just:

$$h_i = \exp(\beta_1 * x_1 + \beta_2 * x_2 + \dots + \beta_k * x_k)$$

- Which represents their personal risk relative to the baseline.
- The “baseline” is just the hazard when all covariate values are zero!

Visualizing the Cox model

- We can recover the hazard and survival functions from the Cox model
- If we fit the Cox model with no predictors, the estimates of $h(t)$ and $S(t)$ are EXACTLY the same as the Kaplan-Meier estimates
- The baseline hazards and survival functions are just the Kaplan-Meier estimates, for

individuals with the reference level for all predictors

- i.e. all 0's for all x's
- This is because their risk score is:

$$\text{Risk Score}_i = \exp(\beta_1 * x_1 + \beta_2 * x_2 + \dots + \beta_k * x_k)$$

- If all x's are 0, then the risk score is 1: $\exp(0) = 1$
- By turning “off” and “on” different x's, we can build different risk scores for different prototypical individuals
 - Remember, you're only as unique as your covariate vector!
- A hazard for an individual with a risk score, y is:

$$\hat{H}(t_{ij}) = \hat{H}_0(t_j) * y_i$$

- which is just a multiplicative effect on the cumulative hazard function.
- To show this in terms of survival, we do:

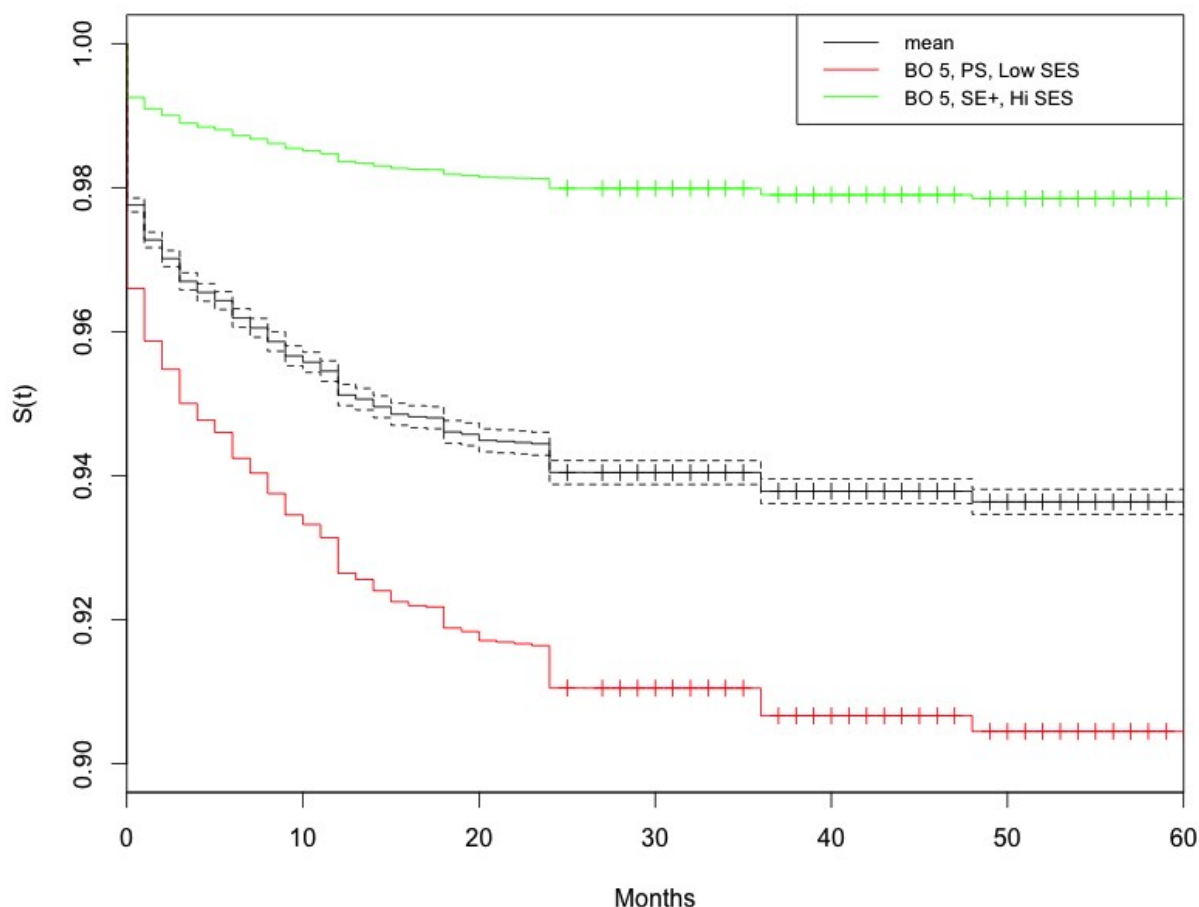
$$\hat{H} = -\log S(t_{ij}) = -\log S_0(t_j) * y_i$$

or

$$\hat{S}(t_{ij}) = \hat{S}(t_j)^{y_i}$$

Which says that the survival function for an individual with risk score y is a power of the baseline survival rate.

- So all we need to do is estimate the K-M functions and multiply them by prototypical risk scores, or prototypical “people”, and we can recover a hazard or survival function estimate for those kinds of people
- This lets us visualize the results very effectively.



[[C:/Users/ozd504/Documents/Github/DEM7223/images/H_ex.png)]

Data Example

This example will illustrate how to fit the Cox Proportional hazards model to continuous duration data (i.e. person-level data) and a discrete-time (longitudinal) data set.

The first example uses longitudinal data from the ECLS-K (<http://nces.ed.gov/ecls/kinderdatainformation.asp>). Specifically, we will examine the transition into poverty between kindergarten and third grade.

In the second example, I use the *time between the first and second birth* for women in the data as the *outcome variable*. The data for this example come from the DHS Model data file Demographic and Health Survey for 2012 (<https://t.co/tM8LfJhomf>) individual recode file. This file contains information for all women sampled in the survey between the ages of 15 and 49.

Using Longitudinal Data

As in the other examples, I illustrate fitting these models to data that are longitudinal, instead of person-duration. In this example, we will examine how to fit the Cox model to a longitudinally collected data set.

First we load our data

```
#Load required libraries
library(foreign)
library(survival)
library(car)
library(survey)
library(eha)
library(dplyr)
options(survey.lonely.psu = "adjust")

eclskk5<-readRDS("C:/Users/ozd504/OneDrive - University of Texas at San Anto
nio/classes/dem7223/dem7223_20//data/eclskk5.rds")
names(eclskk5)<-tolower(names(eclskk5))
#get out only the variables I'm going to use for this example
myvars<-c( "childid", "x_chsex_r", "x_raceth_r",
           "x1kage_r", "x4age", "x5age", "x6age",
           "x7age", "x2povty", "x4povty_i", "x6povty_i",
           "x8povty_i", "x12parled_i", "s2_id", "w6c6p_6psu",
           "w6c6p_6str", "w6c6p_20")
eclskk5<-eclskk5[,myvars]

eclskk5$age1<-ifelse(eclskk5$x1kage_r== -9, NA, eclskk5$x1kage_r/12)
eclskk5$age2<-ifelse(eclskk5$x4age== -9, NA, eclskk5$x4age/12)
#for the later waves, the NCES group the ages into ranges of months, so 1=
<105 months, 2=105 to 108 months. So, I fix the age at the midpoint of the i
nterval they give, and make it into years by dividing by 12
eclskk5$age3<-ifelse(eclskk5$x5age== -9, NA, eclskk5$x5age/12)

eclskk5$pov1<-ifelse(eclskk5$x2povty==1,1,0)
eclskk5$pov2<-ifelse(eclskk5$x4povty_i==1,1,0)
eclskk5$pov3<-ifelse(eclskk5$x6povty_i==1,1,0)

#Recode race with white, non Hispanic as reference using dummy vars
eclskk5$race_rec<-Recode (eclskk5$x_raceth_r, recodes="1 = 'nhwhite';2='nhbl
ack';3:4='hispanic';5='nhasian'; 6:8='other';-9=NA", as.factor = T)
eclskk5$race_rec<-relevel(eclskk5$race_rec, ref = "nhwhite")
eclskk5$male<-Recode(eclskk5$x_chsex_r, recodes="1=1; 2=0; -9=NA")
eclskk5$mlths<-Recode(eclskk5$x12parled_i, recodes = "1:2=1; 3:9=0; else = N
A")
eclskk5$mgths<-Recode(eclskk5$x12parled_i, recodes = "1:3=0; 4:9=1; else =N
A")
```

Now, I need to form the transition variable, this is my event variable, and in this case it will be 1 if a child enters poverty between the first wave of the data and the third grade wave, and 0 otherwise.

NOTE I need to remove any children who are already in poverty age wave 1, because they are not at risk of experiencing **this particular** transition. Again, this is called forming the *risk set*


```
eclskk5<-subset(eclskk5, is.na(pov1)==F&is.na(pov2)==F&
               is.na(pov3)==F&is.na(age1)==F&
               is.na(age2)==F&is.na(age3)==F&pov1!=1)
```

Now we do the entire data set. To analyze data longitudinally, we need to reshape the data from the current “wide” format (repeated measures in columns) to a “long” format (repeated observations in rows). The `reshape()` function allows us to do this easily. It allows us to specify our repeated measures, time varying covariates as well as time-constant covariates.

```
e.long<-reshape(data.frame(eclskk5), idvar="childid",
                    varying=list(c("age1","age2"),
                                c("age2","age3")),
                    v.names=c("age_enter","age_exit"),
                    times=1:2, direction="long" )
e.long<-e.long[order(e.long$childid, e.long$time),]

e.long$povtran<-NA

e.long$povtran[e.long$pov1==0&e.long$pov2==1&e.long$time==1]<-1
e.long$povtran[e.long$pov2==0&e.long$pov3==1&e.long$time==2]<-1

e.long$povtran[e.long$pov1==0&e.long$pov2==0&e.long$time==1]<-0
e.long$povtran[e.long$pov2==0&e.long$pov3==0&e.long$time==2]<-0

#find which kids failed in earlier time periods
#and remove them from the second & third period risk set

failed1<-which(is.na(e.long$povtran)==T)
e.long<-e.long[~failed1,]

e.long$age1r<-round(e.long$age_enter, 0)
e.long$age2r<-round(e.long$age_exit, 0)
head(e.long, n=10)
```

	childid <chr>	x_chsex_r <dbl>	x_raceth_r <dbl>	x1kage_r <dbl>	x4... <dbl>	x5... <dbl>	x6age <dbl>	x7age <dbl>
10000014.1	10000014	1	1	67.82	85.94	91.73	97.51	106.85
10000014.2	10000014	1	1	67.82	85.94	91.73	97.51	106.85
10000020.1	10000020	2	5	68.38	88.57	93.37	100.34	111.12
10000020.2	10000020	2	5	68.38	88.57	93.37	100.34	111.12
10000022.1	10000022	2	8	68.61	87.68	92.98	99.19	110.99
10000022.2	10000022	2	8	68.61	87.68	92.98	99.19	110.99
10000029.1	10000029	2	1	69.40	86.86	92.68	99.32	110.40

	childid <chr>	x_chsex_r <dbl>	x_raceth_r <dbl>	x1kage_r <dbl>	x4... <dbl>	x5... <dbl>	x6age <dbl>	x7age <dbl>
10000029.2	10000029	2	1	69.40	86.86	92.68	99.32	110.40
10000034.1	10000034	1	2	76.24	93.30	99.55	105.96	115.10
10000034.2	10000034	1	2	76.24	93.30	99.55	105.96	115.10

1-10 of 10 rows | 1-10 of 31 columns

Cox regression model

Compared to the parametric models we saw last week (http://rpubs.com/corey_sparks/209276), the Cox model (http://www.jstor.org/stable/pdf/2985181.pdf?seq=1#page_scan_tab_contents) See also the partial likelihood paper (<http://hydra.usc.edu/pm518b/literature/cox-75.pdf>), does not specify a parametric form for the baseline hazard rate. The model still looks the same as the other proportional hazards models:

$$h(t) = h_0 \exp(x'\beta)$$

but h_0 is not a parametric function. Instead, the baseline hazard rate is the empirically observed hazard rate, and the covariates shift it up or down, proportionally.

Using age as the time variable:

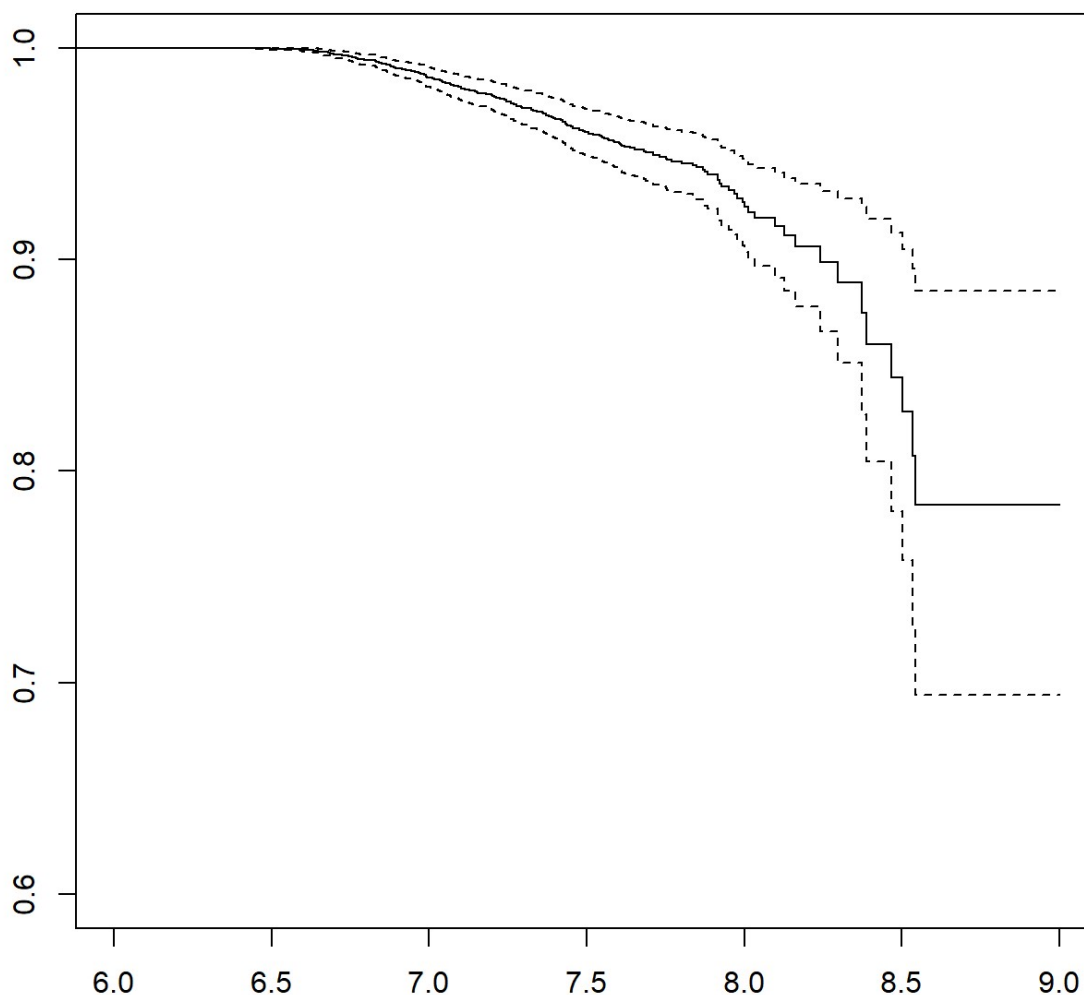
Here I use age of the child as the time variable, this should show how children experience poverty during school.

```
#Cox Model
#interval censored
e.long<-e.long%>%
  filter(complete.cases(age_enter, age_exit, povtran,
                        mlths, mgths, race_rec, w6c6p_6psu))
fit11<-coxreg(Surv(time = age_enter,time2=age_exit, event = povtran)~mlths+m
gths+race_rec,
              data=e.long)
summary(fit11)
```

```
## Call:
## coxreg(formula = Surv(time = age_enter, time2 = age_exit, event = povtra
n) ~
##      mlths + mgths + race_rec, data = e.long)
##
## Covariate      Mean      Coef      Rel.Risk    S.E.      Wald p
## mlths          0.056     0.513     1.669     0.191     0.007
## mgths          0.789    -1.148     0.317     0.161     0.000
## race_rec
##      nhwhite    0.561      0         1 (reference)
##      hispanic   0.221     1.245     3.475     0.179     0.000
##      nhasian    0.083     0.558     1.748     0.309     0.070
##      nhblack    0.065     1.163     3.200     0.255     0.000
##      other      0.070     0.761     2.141     0.291     0.009
##
## Events                221
## Total time at risk    3946.2
## Max. log. likelihood  -1444.4
## LR test statistic      224.58
## Degrees of freedom     6
## Overall p-value        0
```

```
plot(survfit(fitl1), ylim=c(.6,1), xlim=c(6, 9),
      main="Survivorship Function for Cox Regression model - Average Child")
```

Survivorship Function for Cox Regression model - Average Child



The model results (`Rel.Risk`) show that children with mom's who have less than a high school education have 1.67 times higher risk of going into poverty during this period, while children whose mother have more than a high school education are 0.68 % less likely to enter poverty, compared to children whose mothers had a high school education. Likewise, Hispanic, Non-Hispanic black, Native American children all face higher risk of entering poverty, compared to Non-Hispanic whites.

Risk scores

The Cox model generates a "*risk score*" for each individual. This is just $\exp(x'\beta)$, or the exponent of the linear predictor. These are not absolute values that have any real direct interpretation, but they are interpretable in a **relative** sense.

Risk scores >1 indicate that a person has higher risk than the baseline category, while risk scores <1 have lower relative risk, compared to the baseline. If you want to interpret these, it's necessary to have the baseline category be a meaningful "type" of person. In our example above, the baseline group would be non-Hispanic whites, with a mother who had a high school education. i.e. all x 's are 0.

```
e.long$risk<-exp(fit11$linear.predictors)

#highest risk child
e.long[which.max(e.long$risk),c("childid", "age_enter", "mlths", "mgths","race_rec", "risk")]
```

	childid <chr>	age_enter <dbl>	mlths <dbl>	mgths <dbl>	race_rec <fctr>	risk <dbl>
10000744.1	10000744	5.263333	1	0	hispanic	9.002887

1 row

```
#lowest risk child
e.long[which.min(e.long$risk),c("childid", "age_enter", "mlths", "mgths","race_rec", "risk")]
```

	childid <chr>	age_enter <dbl>	mlths <dbl>	mgths <dbl>	race_rec <fctr>	risk <dbl>
10000046.1	10000046	5.9175	0	1	nhwhite	0.4925511

1 row

So, the first of these children has a risk score of 9.07 which means their risk was almost nine times that of the baseline child. Likewise, the lowest risk child had a risk score of .495 that means their score is almost 50% less than the baseline category.

Fitting the Cox model with survey design information

Now we fit the Cox model using full survey design. In the ECLS-K, I use the longitudinal weight for waves 1-5, as well as the associated psu and strata id's for the longitudinal data from these waves from the parents of the child, since no data from the child themselves are used in the outcome.

```
des2<-svydesign(ids = ~w6c6p_6psu,
               strata = ~w6c6p_6str, weights=~w6c6p_20,
               data=e.long, nest=T)

#Fit the model
fit11<-svycoxph(Surv(time =age_enter,
                    time2 = age_exit,
                    event = povtran)~mlths+mgths+race_rec,
               design=des2)

summary(fit11)
```

```
## Stratified 1 - level Cluster Sampling design (with replacement)
## With (123) clusters.
## svydesign(ids = ~w6c6p_6psu, strata = ~w6c6p_6str, weights = ~w6c6p_20,
##      data = e.long, nest = T)
```

```
## Call:
## svycoxph(formula = Surv(time = age_enter, time2 = age_exit, event = povtr
an) ~
##      mlths + mgths + race_rec, design = des2)
##
##      n= 3938, number of events= 221
##
##              coef exp(coef) se(coef) robust se      z Pr(>|z|)
## mlths          0.4620    1.5872  0.1829    0.1931  2.392   0.0168 *
## mgths         -1.1422    0.3191  0.1511    0.2191 -5.213  1.86e-07 ***
## race_rechispanic 1.0563    2.8758  0.1624    0.2091  5.052  4.36e-07 ***
## race_recnhasian  0.4110    1.5084  0.3779    0.3986  1.031   0.3025
## race_recnhblack  1.0686    2.9114  0.2273    0.2325  4.595  4.32e-06 ***
## race_recother    0.4820    1.6193  0.2806    0.2343  2.057   0.0397 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## mlths          1.5872    0.6300    1.0870    2.3176
## mgths          0.3191    3.1336    0.2077    0.4903
## race_rechispanic 2.8758    0.3477    1.9090    4.3324
## race_recnhasian  1.5084    0.6630    0.6905    3.2947
## race_recnhblack  2.9114    0.3435    1.8456    4.5924
## race_recother    1.6193    0.6175    1.0230    2.5633
##
## Concordance= 0.745 (se = 0.023 )
## Likelihood ratio test= NA on 6 df,  p=NA
## Wald test            = 104.6 on 6 df,  p=<2e-16
## Score (logrank) test = NA on 6 df,  p=NA
##
##      (Note: the likelihood ratio and score tests assume independence of
##      observations within a cluster, the Wald and robust score tests do no
##      t).
```

```

plot(survfit(fitl1, conf.int = F), ylab="S(t)",
      xlab="Child Age", xlim=c(6, 9))

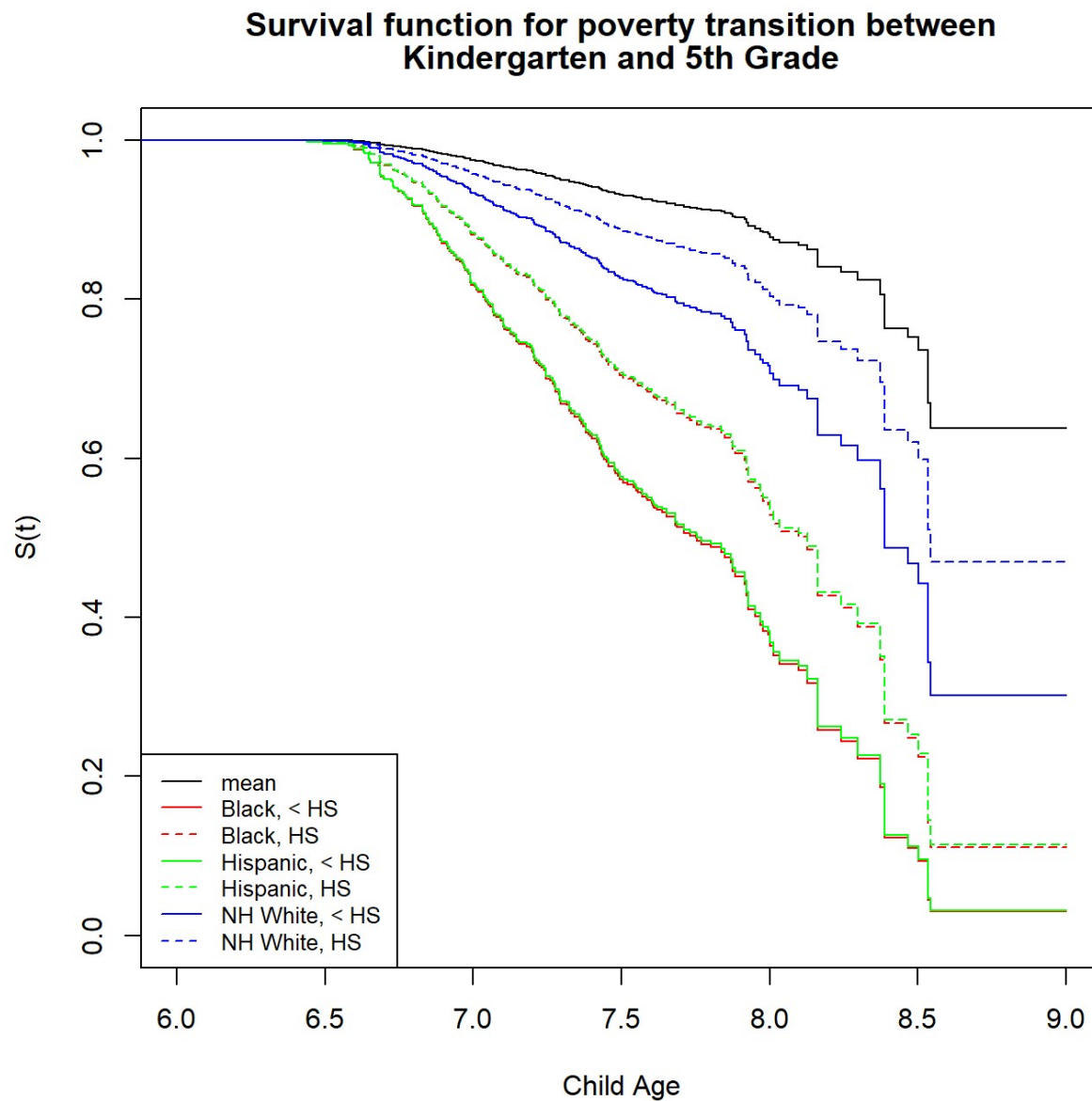
lines(survfit(fitl1, newdata = data.frame(mlths=1, mgths=0, race_rec="nhblac
k"),
          conf.int=F) ,col="red", lty=1)
lines(survfit(fitl1, newdata = data.frame(mlths=0, mgths=0, race_rec="nhblac
k"),
          conf.int=F) ,col="red", lty=2)

lines(survfit(fitl1, newdata = data.frame(mlths=1, mgths=0, race_rec="hispan
ic"),
          conf.int=F) ,col="green", lty=1)
lines(survfit(fitl1, newdata = data.frame(mlths=0, mgths=0, race_rec="hispan
ic"),
          conf.int=F) ,col="green", lty=2)

lines(survfit(fitl1, newdata = data.frame(mlths=1, mgths=0, race_rec="nhwhit
e"),
          conf.int=F) ,col="blue", lty=1)
lines(survfit(fitl1, newdata = data.frame(mlths=0, mgths=0, race_rec="nhwhit
e"),
          conf.int=F) ,col="blue", lty=2)

title(main=c("Survival function for poverty transition between",
             "Kindergarten and 5th Grade"))
legend("bottomleft",
      legend=c("mean", "Black, < HS ", "Black, HS","Hispanic, < HS ",
               "Hispanic, HS","NH White, < HS ", "NH White, HS"),
      col=c(1,"red","red", "green","green", "blue", "blue"),
      lty=c(1,1,2,1,2,1,2), cex=.8)

```



Use time instead of age

Next, I will use the `time` variable we created in `e.long` as the time axis. This model will not focus on the age of the children, but on the probability of experiencing the transition between waves.

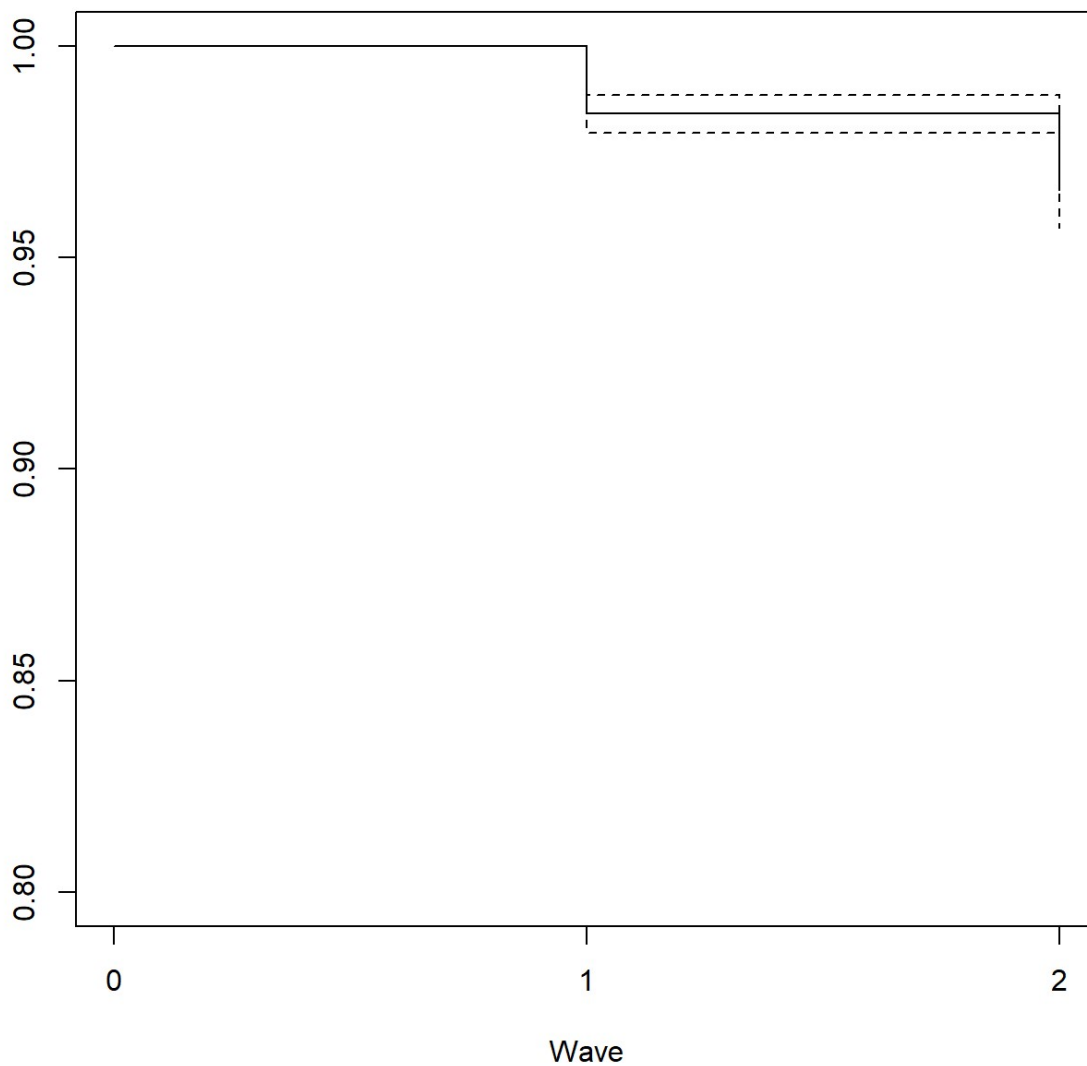
```
#Cox Model
#interval censored
fitl2<-coxreg(Surv(time = time, event = povtran)~mlths+mgths+race_rec, data=
e.long)
summary(fitl2)
```



```
## Call:
## coxreg(formula = Surv(time = time, event = povtran) ~ mlths +
##       mgths + race_rec, data = e.long)
##
## Covariate           Mean           Coef           Rel.Risk    S.E.        Wald p
## mlths                0.051         0.449         1.567        0.189        0.018
## mgths                0.798        -1.226         0.293        0.160        0.000
## race_rec
##       nhwhite        0.570         0           1 (reference)
##       hispanic        0.213         1.159         3.188        0.179        0.000
##       nhasian         0.083         0.424         1.528        0.308        0.169
##       nhblack         0.063         1.134         3.107        0.253        0.000
##       other           0.071         0.783         2.189        0.291        0.007
##
## Events                221
## Total time at risk      5834
## Max. log. likelihood   -1659.6
## LR test statistic       221.90
## Degrees of freedom      6
## Overall p-value        0
```

```
plot(survfit(fit12), ylim=c(.8,1), xlab="Wave",xaxt="n",
      main="Survivorship Function for Cox Regression model - Average Child")
axis(1, at=c(0,1,2))
```

Survivorship Function for Cox Regression model - Average Child



The model results (`Rel.Risk`) show that children with mom's who have less than a high school education have 2.1 times higher risk of going into poverty during this period, while children whose mother have more than a high school education are 67% less likely to enter poverty, compared to children whose mothers had a high school education. Likewise, Hispanic, Non-Hispanic black, Native American and Asian children all face higher risk of entering poverty, compared to Non-Hispanic whites.

Now we fit the Cox model using full survey design. In the ECLS-K, I use the longitudinal weight for waves 1-5, as well as the associated psu and strata id's for the longitudinal data from these waves from the parents of the child, since no data from the child themselves are used in the outcome.

```
#Fit the model
fitl2s<-svycoxph(Surv(time = time, event = povtran)~mlths+mgths+race_rec, de
sign=des2)
summary(fitl2s)
```

```
## Stratified 1 - level Cluster Sampling design (with replacement)
## With (123) clusters.
## svydesign(ids = ~w6c6p_6psu, strata = ~w6c6p_6str, weights = ~w6c6p_20,
##       data = e.long, nest = T)
```

```
## Call:
## svycoxph(formula = Surv(time = time, event = povtran) ~ mlths +
##       mgths + race_rec, design = des2)
##
## n= 3938, number of events= 221
##
##               coef exp(coef) se(coef) robust se      z Pr(>|z|)
## mlths           0.4252   1.5299  0.1841   0.2121  2.004 0.045018 *
## mgths          -1.2290   0.2926  0.1511   0.2195 -5.599 2.16e-08 ***
## race_rechispanic 0.9157   2.4984  0.1658   0.2929  3.127 0.001768 **
## race_recnhasian  0.2686   1.3082  0.3773   0.4022  0.668 0.504237
## race_recnhblack  0.9278   2.5290  0.2251   0.2751  3.373 0.000744 ***
## race_recother    0.4383   1.5501  0.2801   0.2870  1.527 0.126760
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## mlths           1.5299   0.6536   1.0095   2.3186
## mgths           0.2926   3.4178   0.1903   0.4499
## race_rechispanic 2.4984   0.4003   1.4073   4.4355
## race_recnhasian  1.3082   0.7644   0.5947   2.8777
## race_recnhblack  2.5290   0.3954   1.4750   4.3360
## race_recother    1.5501   0.6451   0.8831   2.7208
##
## Concordance= 0.739 (se = 0.022 )
## Likelihood ratio test= NA on 6 df,  p=NA
## Wald test            = 84.51 on 6 df,  p=4e-16
## Score (logrank) test = NA on 6 df,  p=NA
##
## (Note: the likelihood ratio and score tests assume independence of
##       observations within a cluster, the Wald and robust score tests do not).
```

```

plot(survfit(fitl2s, conf.int = F),
     ylab="S(t)", xlab="Wave", xaxt="n",
     ylim=c(.2,1))
axis(1, at=c(0,1,2))

lines(survfit(fitl2s,
              newdata = data.frame(mlths=1, mgths=0, race_rec="nhblack"),
              conf.int=F) ,
      col="red", lty=1)
lines(survfit(fitl2s,
              newdata = data.frame(mlths=0, mgths=0, race_rec="nhblack"),
              conf.int=F) ,
      col="red", lty=2)

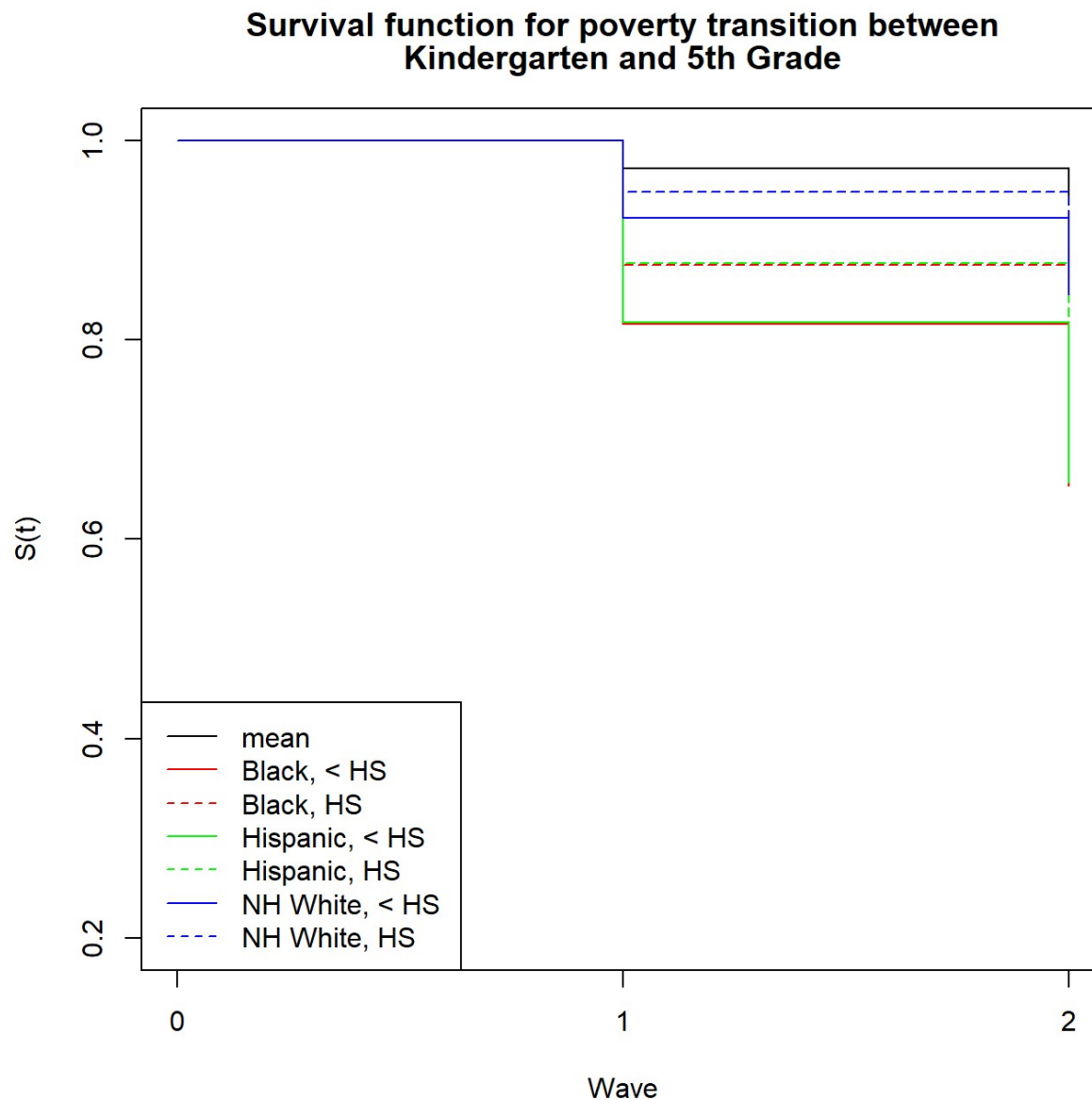
lines(survfit(fitl2s,
              newdata = data.frame(mlths=1, mgths=0, race_rec="hispanic"), c
onf.int=F) ,
      col="green", lty=1)

lines(survfit(fitl2s,
              newdata = data.frame(mlths=0, mgths=0, race_rec="hispanic"), c
onf.int=F),
      col="green", lty=2)

lines(survfit(fitl2s,
              newdata = data.frame(mlths=1, mgths=0, race_rec="nhwhite"),
              conf.int=F) ,
      col="blue", lty=1)
lines(survfit(fitl2s,
              newdata = data.frame(mlths=0, mgths=0, race_rec="nhwhite"),
              conf.int=F) ,
      col="blue", lty=2)

title(main=c("Survival function for poverty transition between", "Kindergarte
n and 5th Grade"))
legend("bottomleft",
      legend=c("mean", "Black, < HS ", "Black, HS", "Hispanic, < HS ",
               "Hispanic, HS", "NH White, < HS ", "NH White, HS"),
      col=c(1,"red","red", "green","green", "blue", "blue"),
      lty=c(1,1,2,1,2,1,2))

```



We see similar results as we did in the age based analysis, but now, we are treating time discretely, and separating it from the child's age entirely.

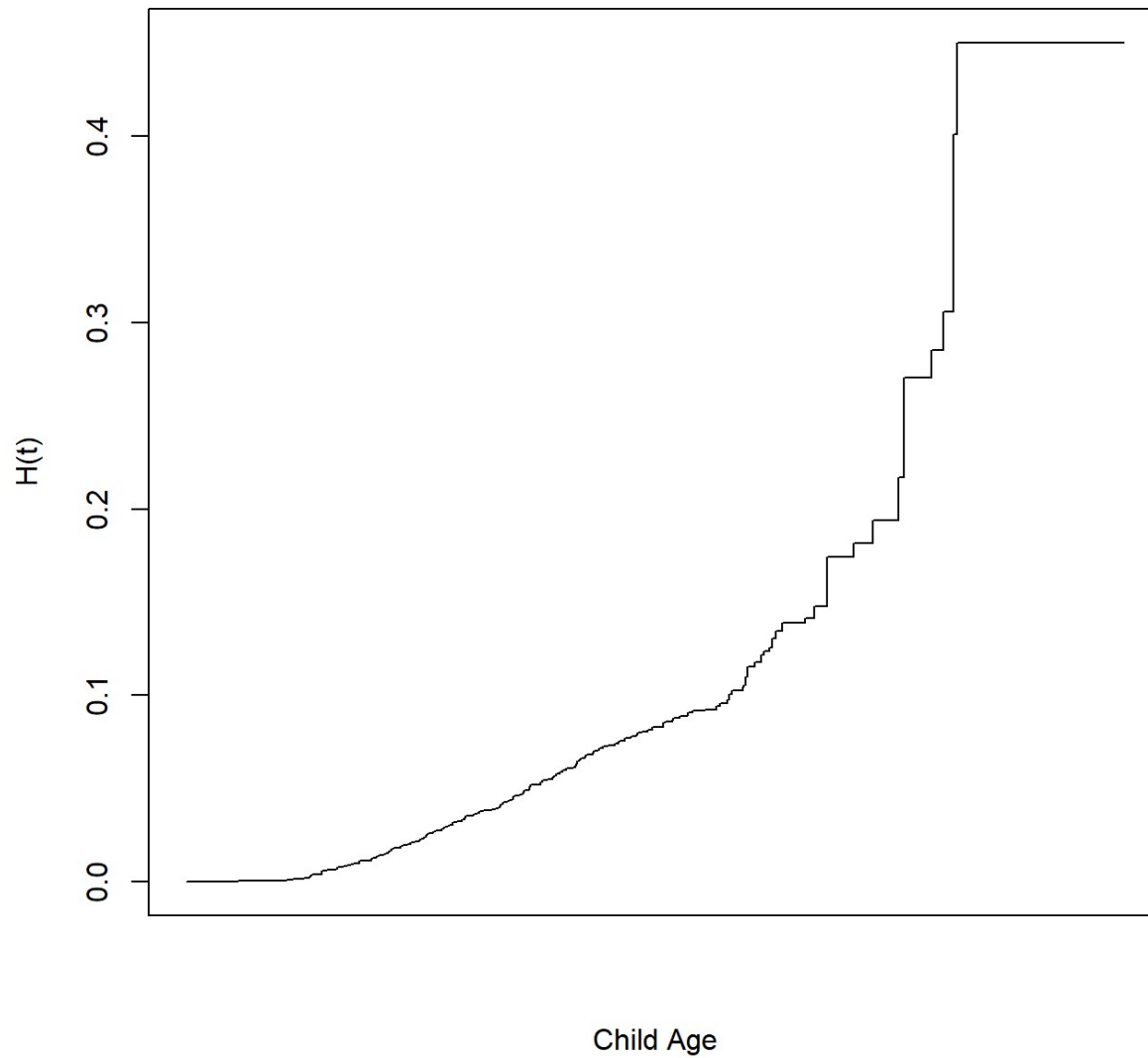
Functions of survival time

Here is an example of calculating the functions of survival time. These aren't very exciting for this example

```
sfl<-survfit(fit11)
H1<--log(sfl$surv)
times<-sfl$time

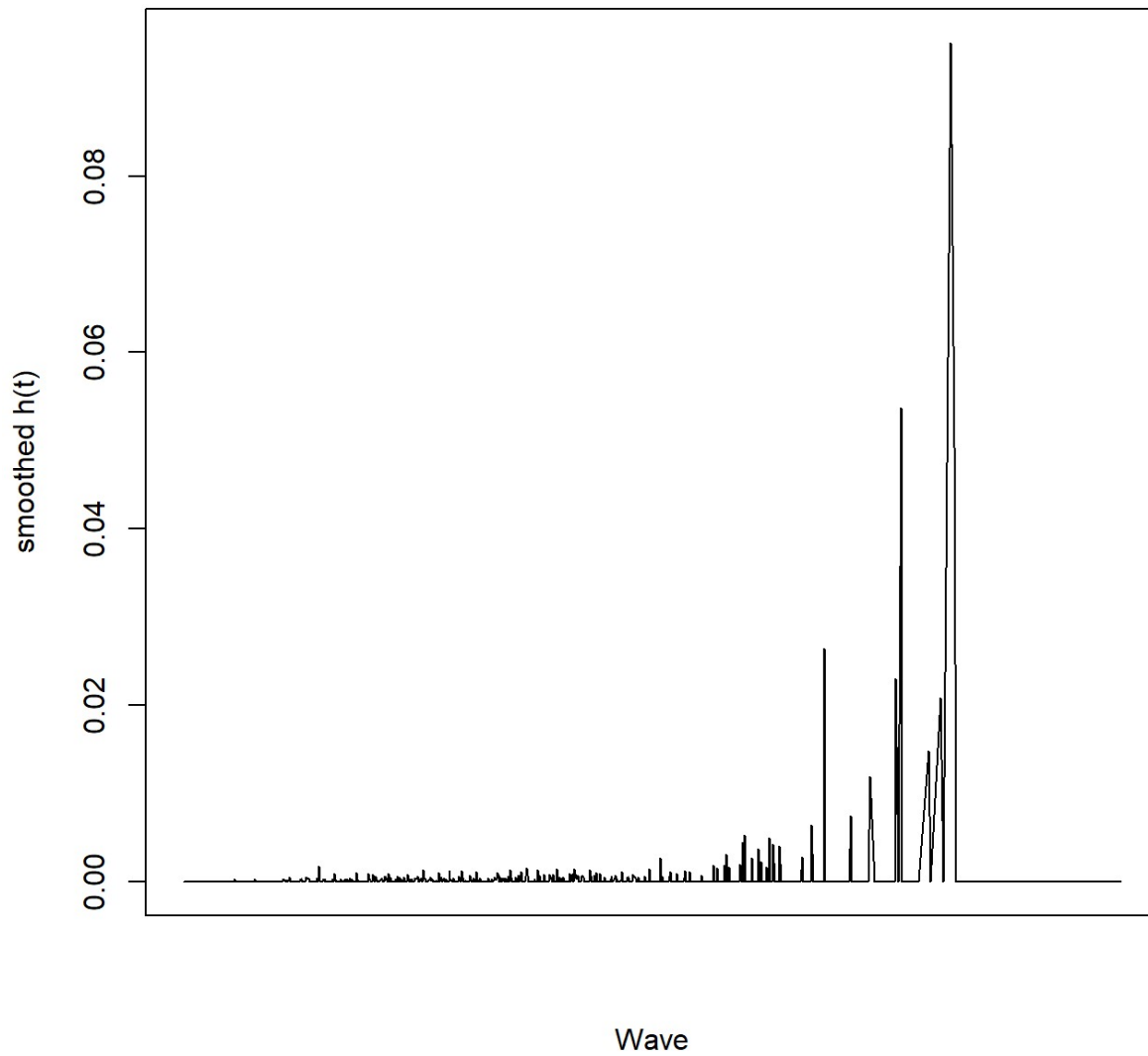
plot(H1~times, type="s", ylab="H(t)",xlab="Child Age",
     xaxt="n", main="Cumulative Hazard plot")
axis(1, at=c(0,1,2))
```

Cumulative Hazard plot



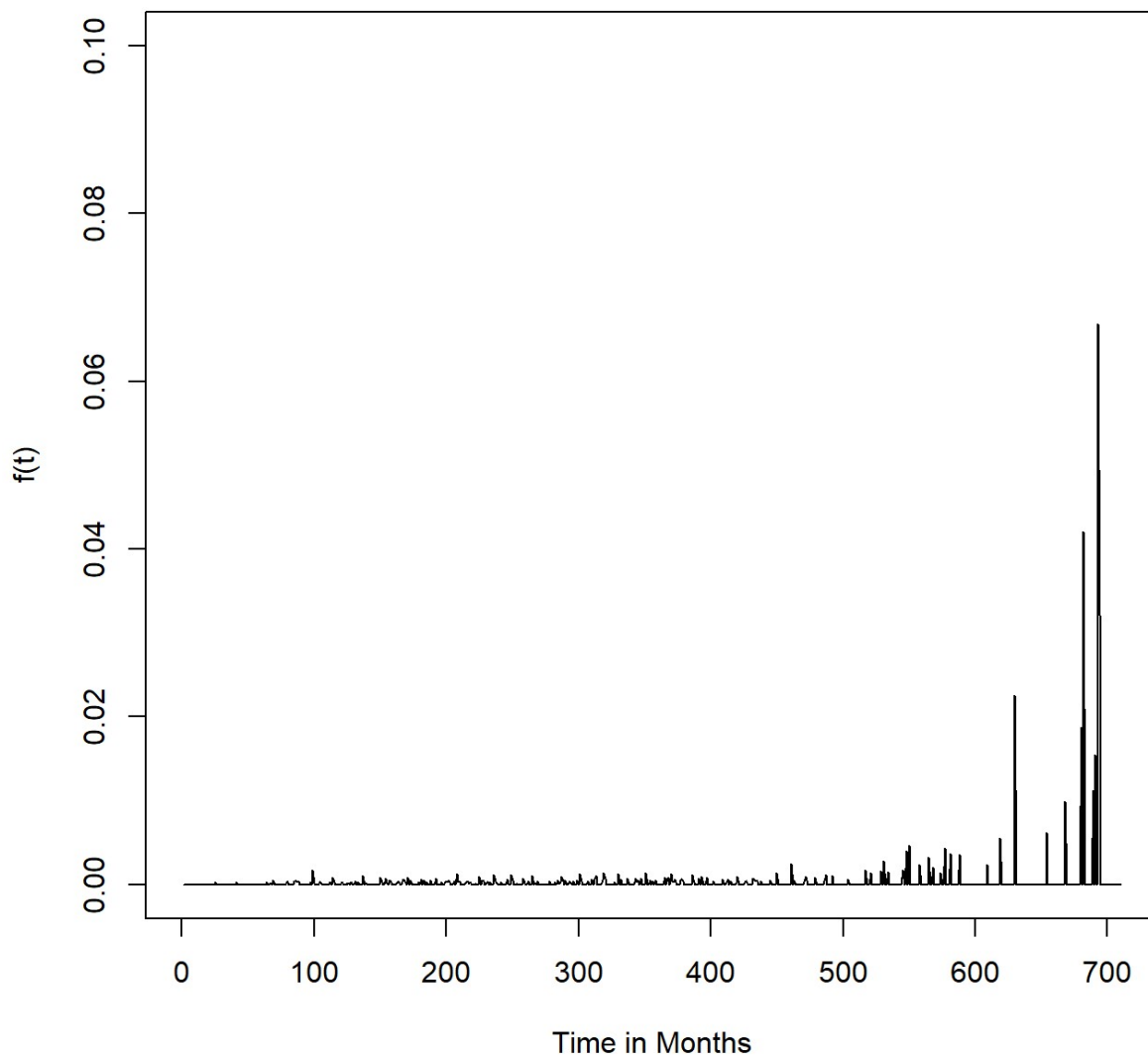
```
hs1<-diff(c(0,H1))  
plot(hs1~times,type="l", ylab="smoothed h(t)", xlab="Wave",  
      xaxt="n", main="Hazard plot")  
axis(1, at=c(0,1,2))
```

Hazard plot



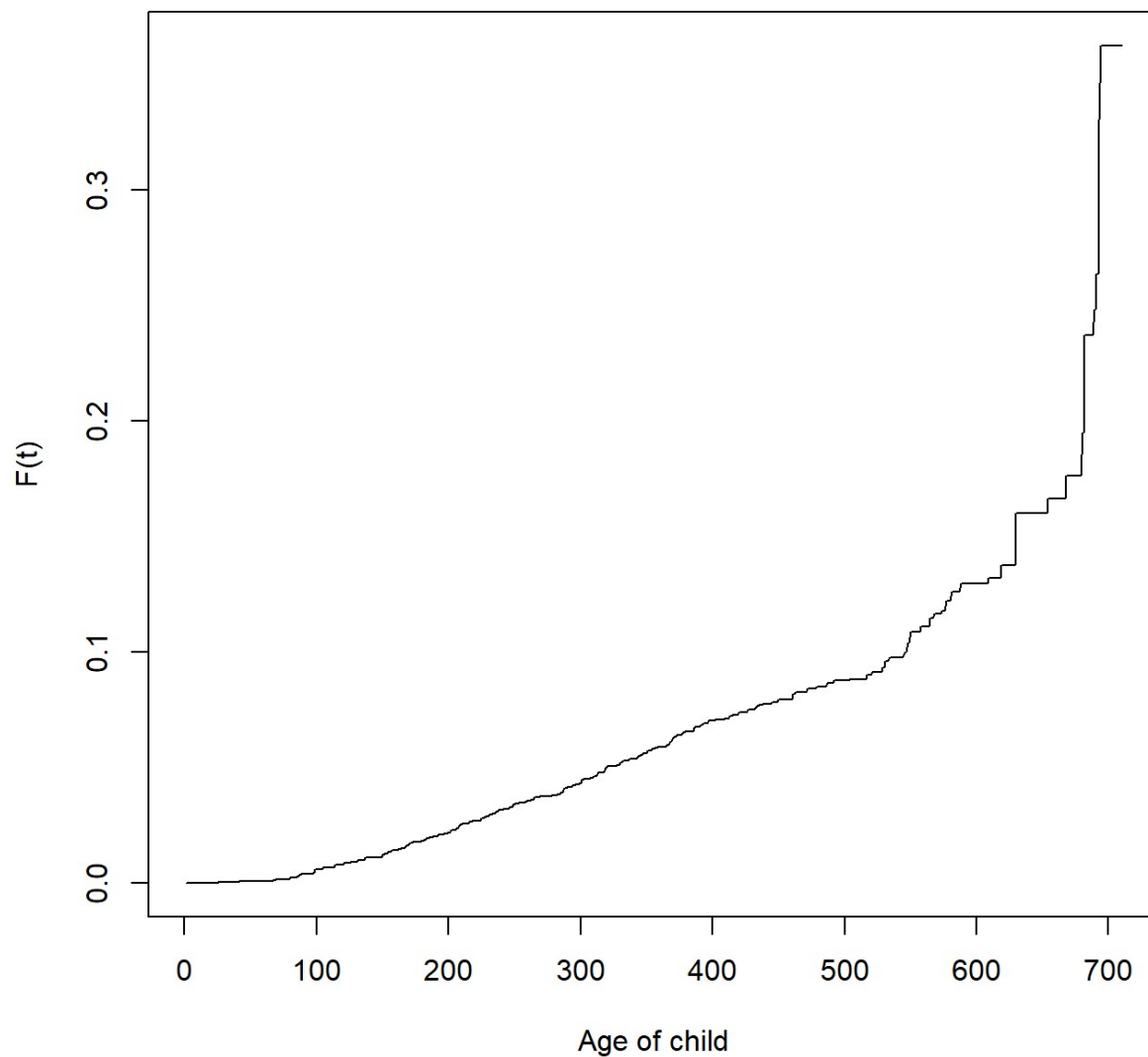
```
ft<- -diff(sfl$surv)
plot(ft,ylim=c(0, .1), type="l",
      ylab="f(t)",xlab="Time in Months",
      main="Probability Density Function")
```

Probability Density Function



```
#here is the cumulative distribution function  
Ft<-cumsum(ft)  
plot(Ft, type="l", ylab="F(t)",xlab="Age of child", main="Cumulative Distrib  
ution Function")
```


Cumulative Distribution Function



DHS data example

```
library(haven)
#load the data
model.dat<-read_dta("https://github.com/coreysparks/data/blob/master/ZZIR62F
L.DTA?raw=true")
model.dat<-zap_labels(model.dat)
```

In the DHS individual recode file, information on every live birth is collected using a retrospective birth history survey mechanism.

Since our outcome is time between first and second birth, we must select as our risk set, only women who have had a first birth.

The `bidx` variable indexes the birth history and if `bidx_01` is not missing, then the woman should be at risk of having a second birth (i.e. she has had a first birth, i.e. `bidx_01==1`).

I also select only non-twin births (`b0 == 0`).

The DHS provides the dates of when each child was born in Century Month Codes.

To get the interval for women who *actually had* a second birth, that is the difference between the CMC for the first birth `b3_01` and the second birth `b3_02` , but for women who had not had a second birth by the time of the interview, the censored time between births is the difference between `b3_01` and `v008` , the date of the interview.

We have 6161 women who are at risk of a second birth.

```
table(is.na(model.dat$bidx_01))
```

```
##
## FALSE TRUE
## 6161 2187
```

```
#now we extract those women
sub<-subset(model.dat, model.dat$bidx_01==1&model.dat$b0_01==0)

#Here I keep only a few of the variables for the dates,
#and some characteristics of the women, and details of the survey

sub2<-data.frame(CASEID=sub$caseid,
                  int.cmc=sub$v008,
                  fbir.cmc=sub$b3_01,
                  sbir.cmc=sub$b3_02,
                  marr.cmc=sub$v509,
                  rural=sub$v025,
                  educ=sub$v106,
                  age=sub$v012,
                  partneredu=sub$v701,
                  partnerage=sub$v730,
                  weight=sub$v005/1000000,
                  psu=sub$v021, strata=sub$v022)

sub2$agefb = (sub2$age - (sub2$int.cmc - sub2$fbir.cmc)/12)
```

Now I need to calculate the birth intervals, both observed and censored, and the event indicator (i.e. did the women *have* the second birth?)

```
sub2$secbi<-ifelse(is.na(sub2$sbir.cmc)==T,
                  ((sub2$int.cmc)-(sub2$fbir.cmc)),
                  (sub2$fbir.cmc-sub2$sbir.cmc))

sub2$b2event<-ifelse(is.na(sub2$sbir.cmc)==T,0,1)
```

Create covariates

Here, we create some predictor variables: Woman's education (secondary +, vs < secondary),

Woman's age², Partner's education (> secondary school)

```
sub2$educ.high<-ifelse(sub2$educ %in% c(2,3), 1, 0)
sub2$age2<-(sub2$age)^2
sub2$partnerhiedu<-ifelse(sub2$partneredu<3,0,
                           ifelse(sub2$partneredu%in%c(8,9),NA,1 ))

options(survey.lonely.psu = "adjust")
des<-svydesign(ids=~psu, strata=~strata,
              data=sub2[sub2$secbi>0,], weight=~weight )
```

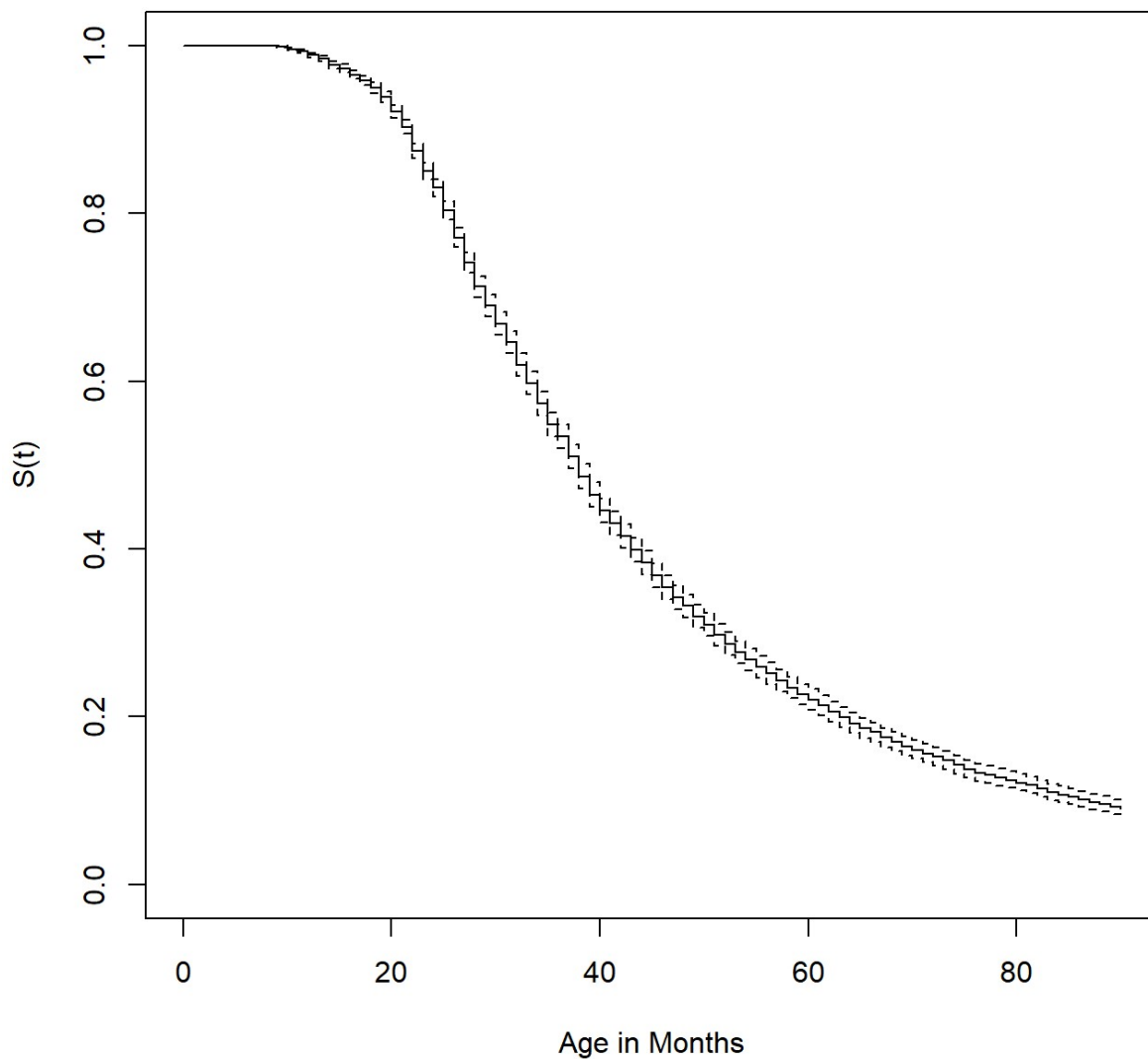
Fit the model

```
#using coxph in survival library
fit.cox2<-coxph(Surv(secbi,b2event)~educ.high+partnerhiedu+age+age2 ,
               data=sub2)
summary(fit.cox2)
```

```
## Call:
## coxph(formula = Surv(secbi, b2event) ~ educ.high + partnerhiedu +
##       age + age2, data = sub2)
##
##      n= 5289, number of events= 4527
##      (737 observations deleted due to missingness)
##
##              coef    exp(coef)    se(coef)      z  Pr(>|z|)
## educ.high    -0.3923429    0.6754725    0.0475997  -8.243   <2e-16 ***
## partnerhiedu -0.1948393    0.8229669    0.0687517  -2.834    0.0046 **
## age          -0.0194407    0.9807471    0.0159741  -1.217    0.2236
## age2          0.0001017    1.0001017    0.0002302    0.442    0.6587
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## educ.high          0.6755      1.4804      0.6153      0.7415
## partnerhiedu        0.8230      1.2151      0.7192      0.9417
## age                0.9807      1.0196      0.9505      1.0119
## age2               1.0001      0.9999      0.9997      1.0006
##
## Concordance= 0.547 (se = 0.005 )
## Likelihood ratio test= 142 on 4 df,  p=<2e-16
## Wald test              = 131.2 on 4 df,  p=<2e-16
## Score (logrank) test = 132.5 on 4 df,  p=<2e-16
```

```
plot(survfit(fit.cox2), xlim=c(0,90),
     ylab="S(t)", xlab="Age in Months")
title(main="Survival Function for Second Birth Interval")
```

Survival Function for Second Birth Interval



```
#use survey design
des<-svydesign(ids=~psu, strata = ~strata ,
              weights=~weight, data=sub2)

cox.s<-svycoxph(Surv(secbi,b2event)~educ.high+partnerhiedu+I(age/5)+age2,
                design=des)
summary(cox.s)
```

```
## Stratified 1 - level Cluster Sampling design (with replacement)
## With (217) clusters.
## svydesign(ids = ~psu, strata = ~strata, weights = ~weight, data = sub2)
```

```
## Call:
## svycoxph(formula = Surv(secbi, b2event) ~ educ.high + partnerhiedu +
##          I(age/5) + age2, design = des)
##
##      n= 5289, number of events= 4527
##      (737 observations deleted due to missingness)
##
##              coef    exp(coef)    se(coef)  robust se         z Pr(>|z|)
## educ.high    -0.4178960  0.6584307  0.0459306  0.0569837 -7.334 2.24e-13
## ***
## partnerhiedu -0.2352044  0.7904093  0.0680935  0.0912254 -2.578  0.00993
## **
## I(age/5)      -0.1970892  0.8211174  0.0811498  0.1041589 -1.892  0.05846
## .
## age2          0.0003923  1.0003924  0.0002343  0.0003023  1.298  0.19430
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## educ.high      0.6584      1.5188    0.5889    0.7362
## partnerhiedu    0.7904      1.2652    0.6610    0.9452
## I(age/5)        0.8211      1.2179    0.6695    1.0071
## age2            1.0004      0.9996    0.9998    1.0010
##
## Concordance= 0.555 (se = 0.006 )
## Likelihood ratio test= NA on 4 df,  p=NA
## Wald test              = 98.89 on 4 df,  p=<2e-16
## Score (logrank) test = NA on 4 df,  p=NA
##
## (Note: the likelihood ratio and score tests assume independence of
## observations within a cluster, the Wald and robust score tests do not).
```

```
dat<-expand.grid(educ.high=c(0,1), partnerhiedu=c(0,1), age=seq(20,40,5), b2
event=1)
dat$age2<-dat$age^2
head(dat)
```

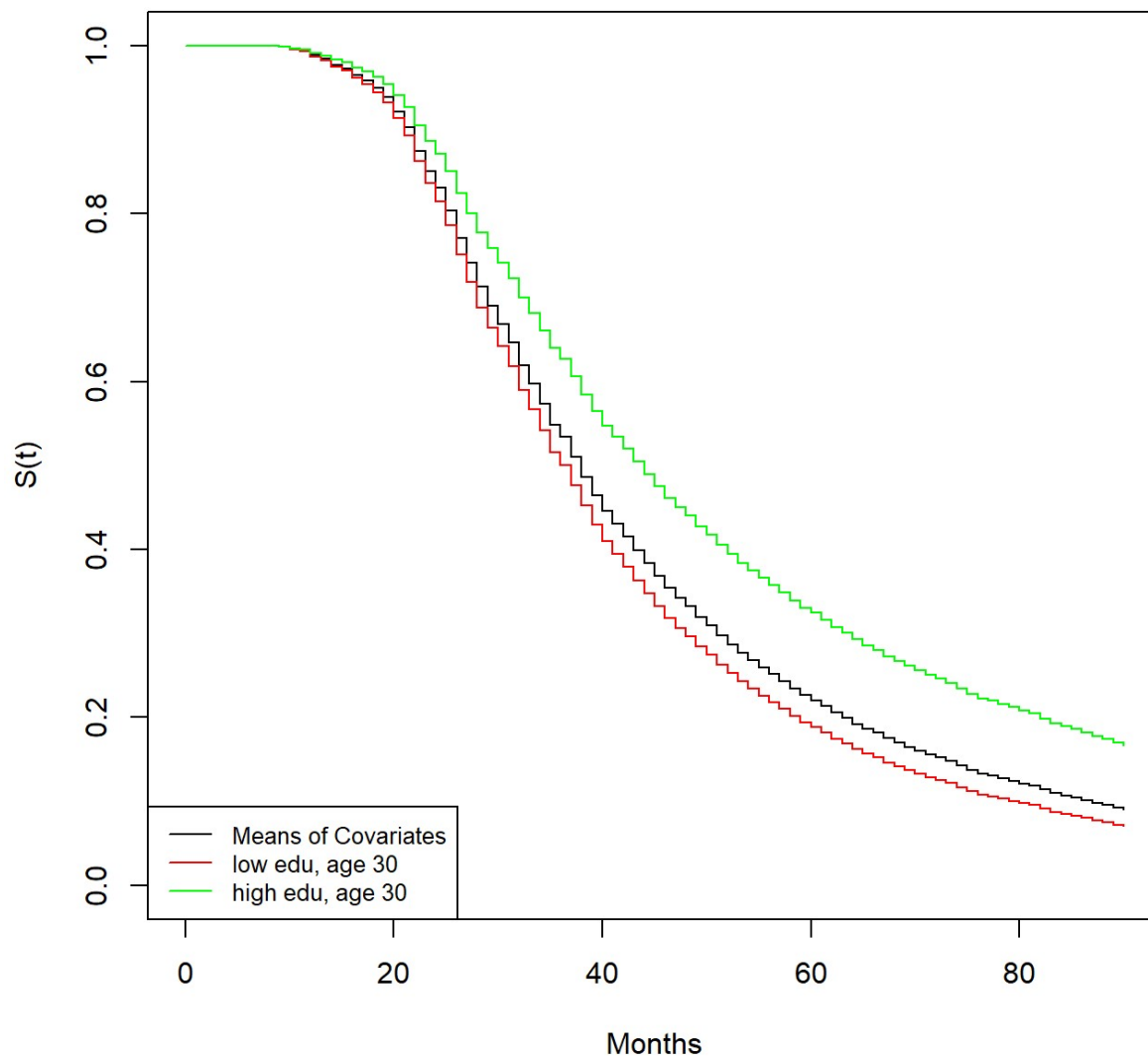
	educ.high <dbl>	partnerhiedu <dbl>	age <dbl>	b2event <dbl>	age2 <dbl>
1	0	0	20	1	400
2	1	0	20	1	400
3	0	1	20	1	400
4	1	1	20	1	400
5	0	0	25	1	625

	educ.high <dbl>	partnerhiedu <dbl>	age <dbl>	b2event <dbl>	age2 <dbl>
6	1	0	25	1	625
6 rows					

```
#plot some survival function estimates for various types of children
plot(survfit(fit.cox2, conf.int = F), xlim=c(0, 90), ylab="S(t)",
      xlab="Months")
title (main = "Survival Plots for for Second Birth Interval")

lines(survfit(fit.cox2,
              newdata=data.frame(educ.high=0, partnerhiedu=0, age=30, age2=9
00),
              conf.int = F), col="red")
lines(survfit(fit.cox2,
              newdata=data.frame(educ.high=1, partnerhiedu=0, age=30, age2=9
00),
              conf.int = F), col="green")
legend("bottomleft", legend=c("Means of Covariates", "low edu, age 30", "hig
h edu, age 30"), lty=1, col=c(1,"red", "green"), cex=.8)
```

Survival Plots for for Second Birth Interval



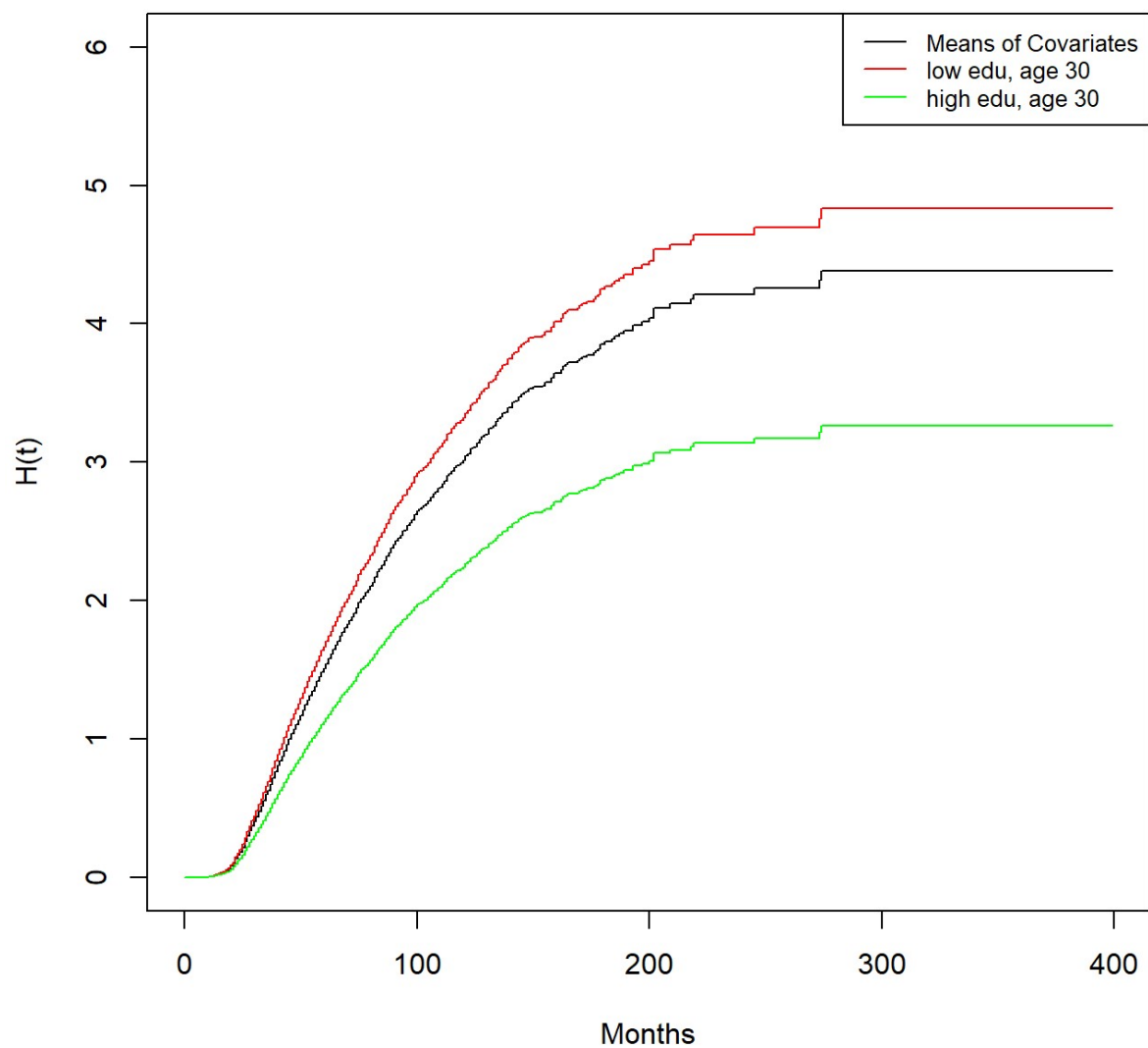
Now we look at some more plots we can examine from the models

```
#Now we look at the cumulative hazard functions
sf1<-survfit(fit.cox2)
sf2<-survfit(fit.cox2,
             newdata=data.frame(educ.high=0, partnerhiedu=0, age=30, age2=90
0))
sf3<-survfit(fit.cox2,
             newdata=data.frame(educ.high=1, partnerhiedu=0, age=30, age2=90
0))

H1<--log(sf1$surv)
H2<--log(sf2$surv)
H3<--log(sf3$surv)

times<-sf1$time

plot(H1~times, type="s", ylab="H(t)",ylim=c(0, 6), xlab="Months")
lines(H2~times,col="red", type="s")
lines(H3~times,col="green", type="s")
legend("topright",
      legend=c("Means of Covariates", "low edu, age 30", "high edu, age 3
0"),
      lty=1, col=c(1,"red", "green"), cex=.8)
```

```
#and the hazard function
hs1<-loess(diff(c(0,H1))~times, degree=1, span=.25)
hs2<-loess(diff(c(0,H2))~times, degree=1, span=.25)
hs3<-loess(diff(c(0,H3))~times, degree=1, span=.25)

plot(predict(hs1)~times,type="l", ylab="smoothed h(t)", xlab="Months",
      ylim=c(0, .04))
title(main="Smoothed hazard plots")
lines(predict(hs2)~times, type="l", col="red")
lines(predict(hs3)~times, type="l", col="green")
legend("topright",
      legend=c("Means of Covariates", "low edu, age 30", "high edu, age 30"),
      lty=1, col=c(1,"red", "green"), cex=.8)
```

