

PH Models continued

Model Assessment

Is It Any Good?

Always want to know how “well” our model did.

Due to censoring as well as Cox regression making relative predictions, not easy/intuitive to evaluate.

Concordance is a popular method to assess model performance:

- For all possible event and non-event pairs we want to assign the higher predicted value to the subject that had the event.
- Survival analysis spin → assign a higher “risk” (hazard) to the subject that had the event **first**
- How well does model rank who will have the event sooner?

Concordance

What is “risk” in this context?

- Risk: $\hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_k x_{i,k}$
- Piece of the model dealing with the predictors

Example:

- Person 1: event at $t = 3$ and $\hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_k x_{i,k} = 1.5$
- Person 2: event (or censored) at $t = 7$ and $\hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_k x_{i,k} = 0.3$
- Concordant pair since person with higher “risk” score had the event first.

Ties, Incomparable, Indeterminate Pairs

If both people have the same event time or censoring time, then the pair is **incomparable** and we don't count it.

Censoring can still mess up pairs:

- Person 1: **censored** at $t = 3$ and $\hat{\beta}_1 x_{i,1} + \dots + \hat{\beta}_k x_{i,k} = 1.5$
- Person 2: event (or censored) at $t = 7$ and $\hat{\beta}_1 x_{i,1} + \dots + \hat{\beta}_k x_{i,k} = 0.3$
- **Indeterminate** pair since no way to know which person had event first
→ not counted.

If both people have the same predicted “risk” then this pair is tied and counted as 0.5.

Concordance – R

```
concordance(recid.ph)
```

```
## Call:
## concordance.coxph(object = recid.ph)
##
## n= 432
## Concordance= 0.6403 se= 0.02666
## discordant concordant      tied.x      tied.y      tied.xy
##      27242      15291         49        111         0
```

Diagnostics

RESIDUALS

Assumptions

Wait...!?!?! I thought you said there were no distributional assumptions!

Still other assumptions we need to check:

- Linearity (maybe higher powers of x are better?)
- Proportional hazards (no interactions with time)

Will deal with these **NOW**

These assumptions can be checked with the help of residuals!

Survival Analysis Residuals

There are four kinds of residuals for survival models, all with various uses:

- Martingale (check linearity, check PH, detect outliers)
- Schoenfeld (check PH)
- Deviance (check linearity, detect outliers)
- Score (detect influential observations)

R will calculate all of these for you.

Survival Analysis Residuals

There are four kinds of residuals for survival models, all with various uses:

- Martingale (check linearity, check PH, detect outliers)
- Schoenfeld (check PH)
- Deviance (check linearity, detect outliers)
- Score (detect influential observations)

Focus here

R will calculate all of these for you.

Martingale Residuals

Martingale residuals are the difference between the observed number of events and the expected number of events at a specific point in time (indicated by the model) “integrated over the time for which that subject was at risk”.

These are **not** symmetrical around zero!

Schoenfeld Residuals

Schoenfeld residuals are calculated for each variable for each individual.

They are the difference between the actual value of the variable and the expected value for someone who had the event occur at that time.

Diagnostics

LINEARITY

Residual Plots

Martingale residual plots in **R** are useful for checking linearity of predictors by plotting them vs. the predictor.

- Similar to looking for residual patterns in linear regression revealing lack of linearity.
- The visreg package also calculates partial residuals which provides similar information here.

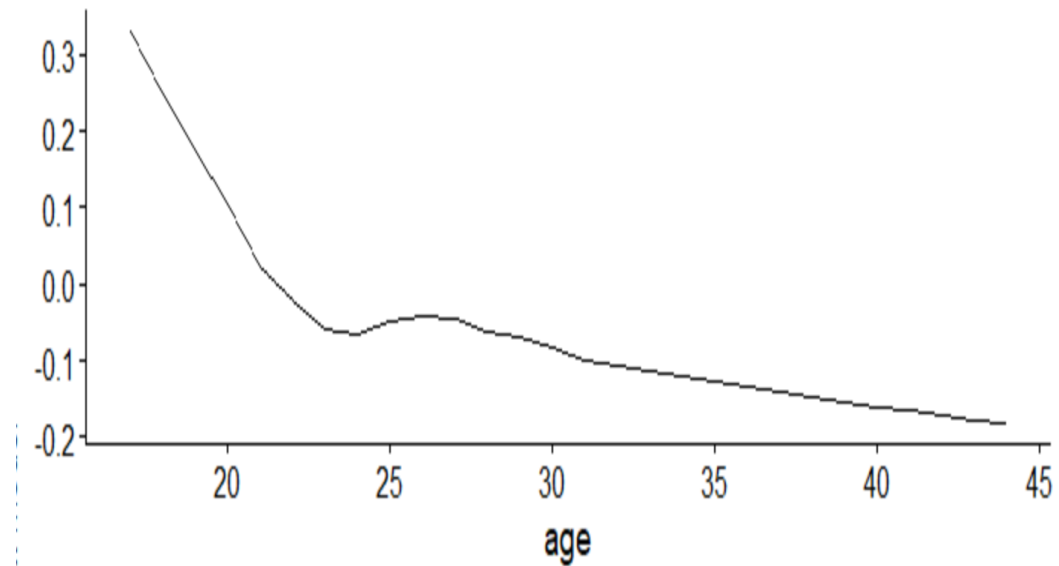
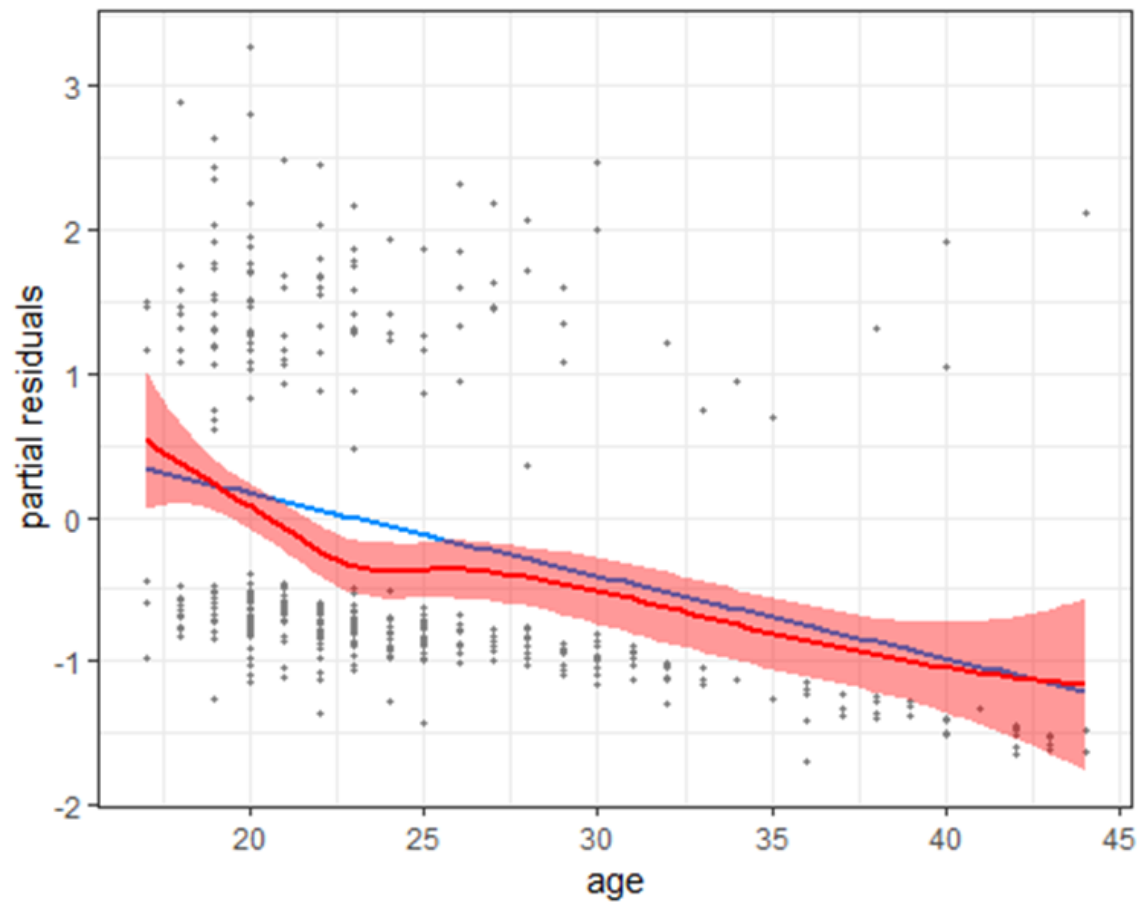
Linearity – R

```
visreg(recid.ph, "age", xlab = "age", ylab = "partial residuals",  
       gg = TRUE, band = FALSE) +  
  geom_smooth(col = "red", fill = "red") + theme_bw()
```

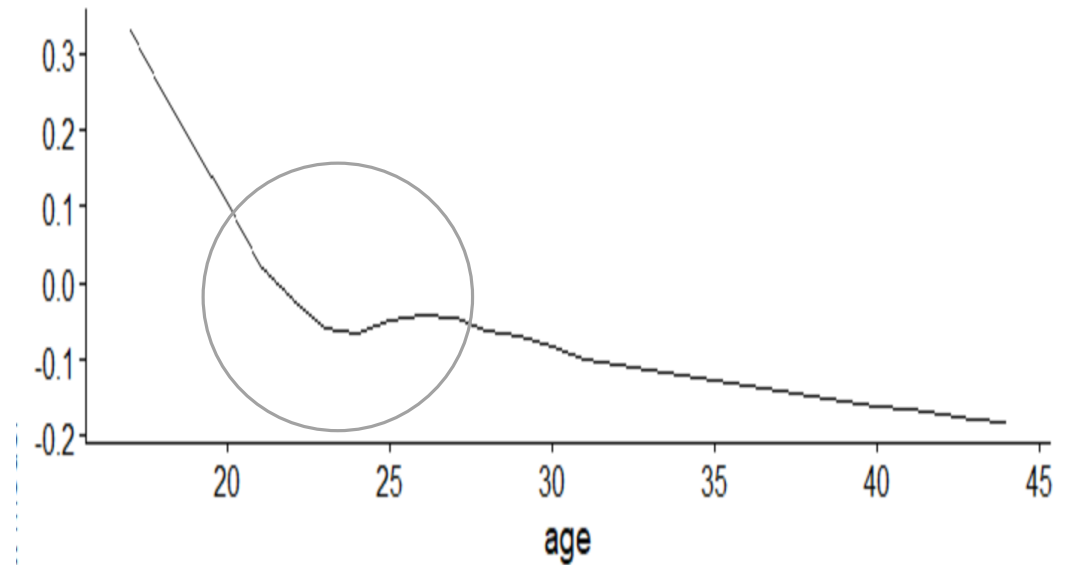
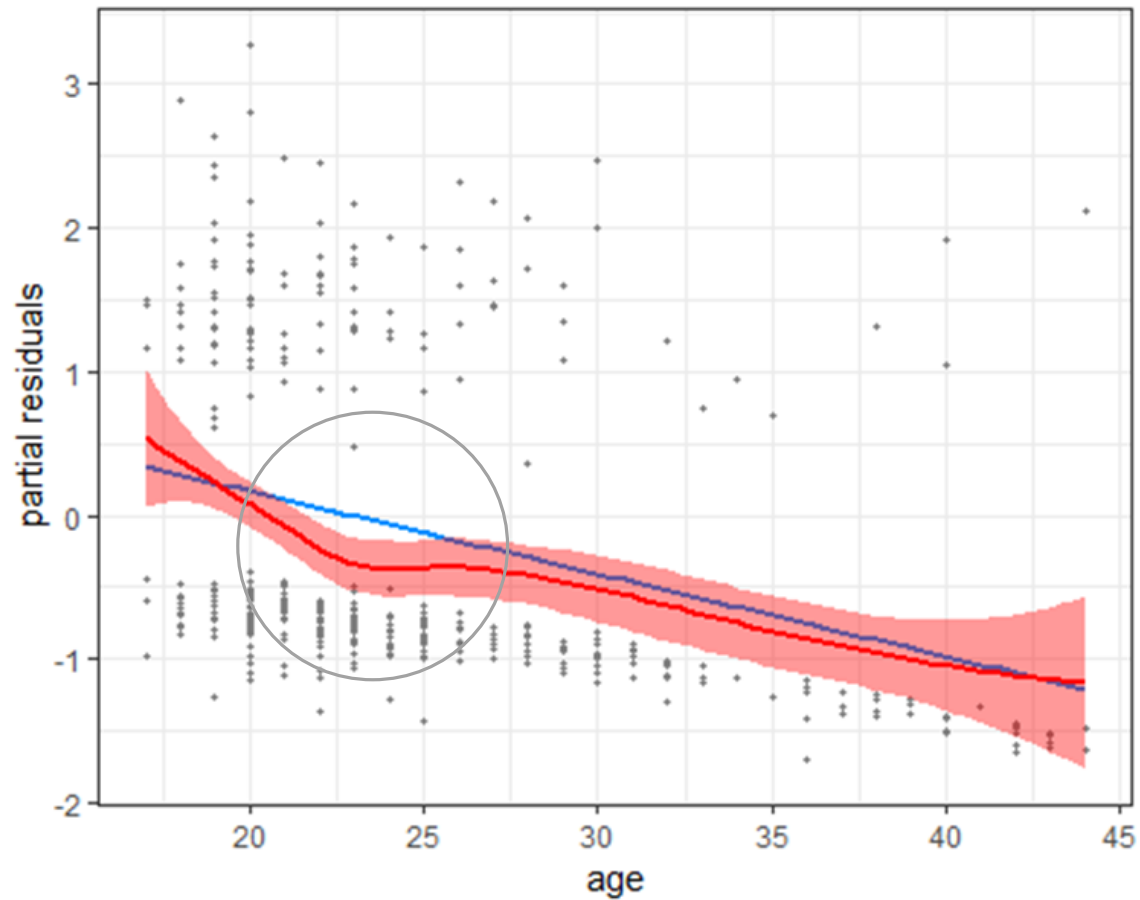
```
visreg(recid.ph, "prio", xlab = "#prior convictions",  
       ylab = "partial residuals",  
       gg = TRUE, band = FALSE) +  
  geom_smooth(col = "red", fill = "red") + theme_bw()
```

```
recid.lin <- coxph(Surv(week, arrest) ~ age + prio, data = recid)  
survminer::ggcoxfunctional(recid.lin, data=recid)
```

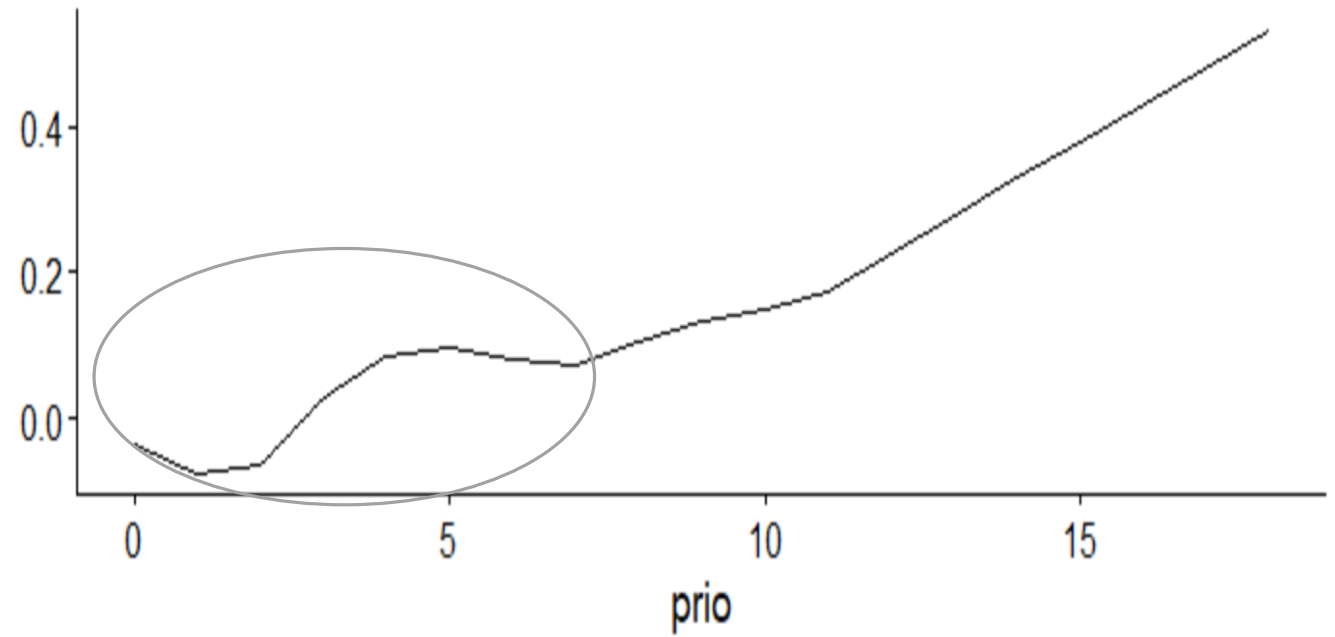
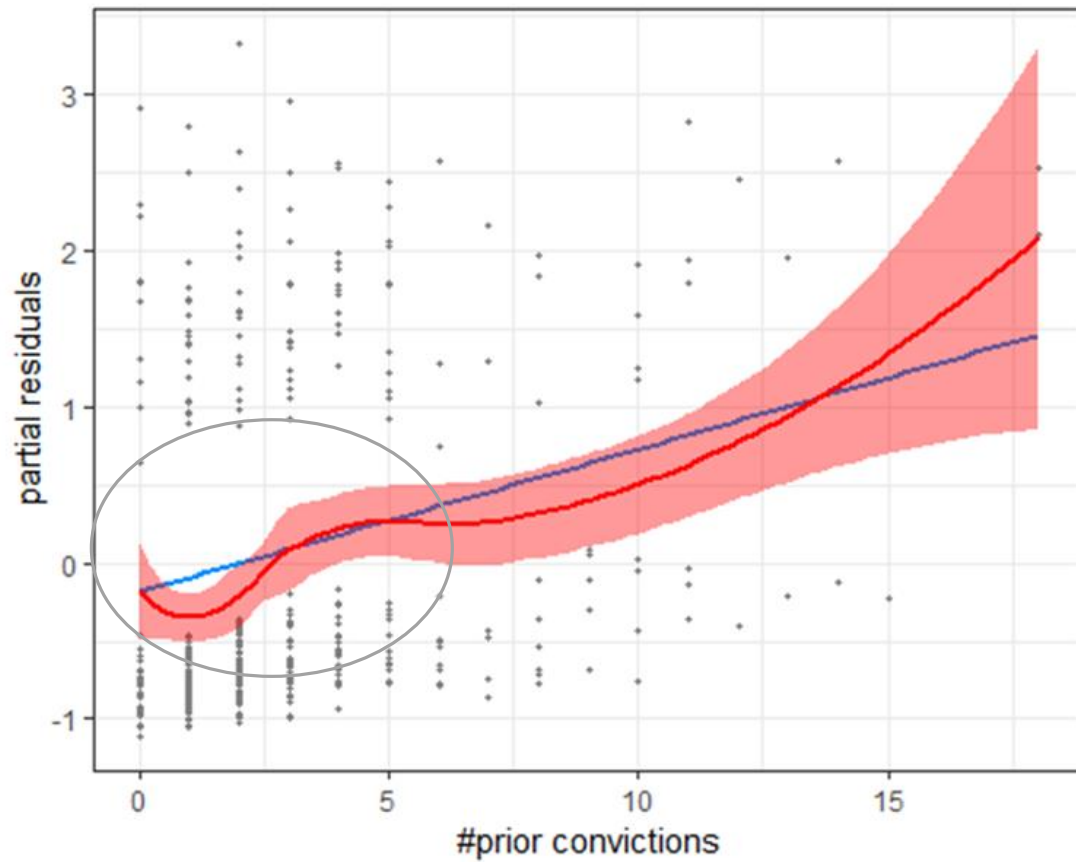
Linearity – R



Linearity – R

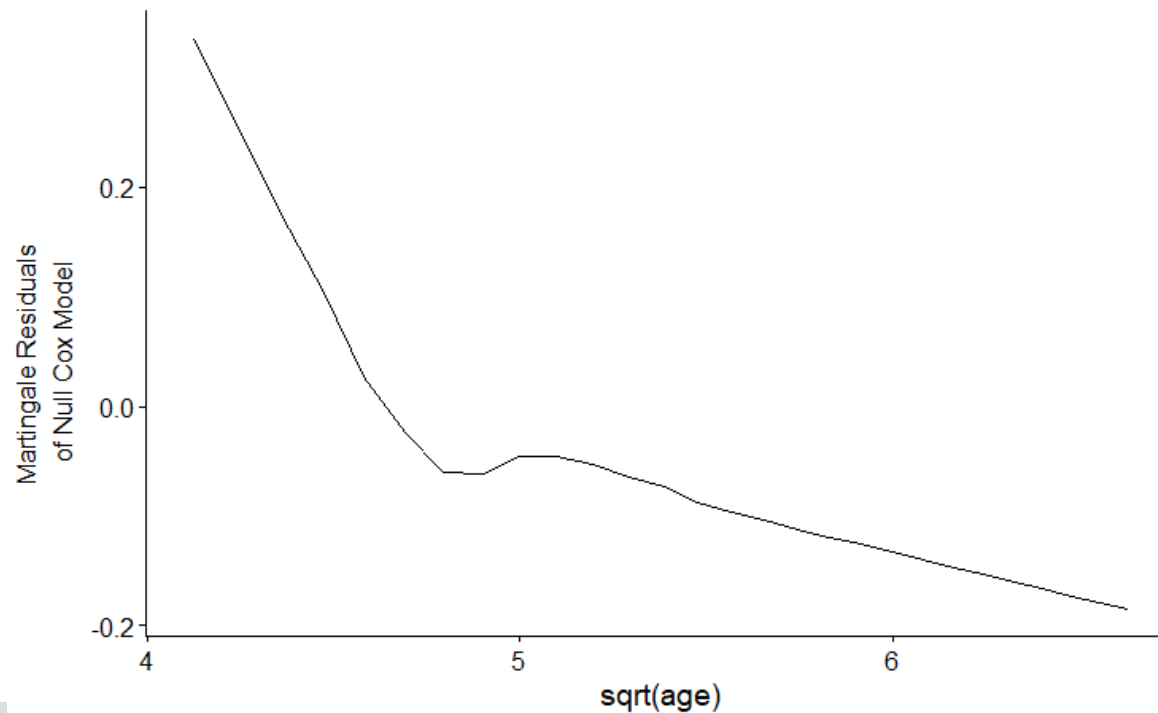


Linearity – R



Can try transformations!

```
recid.lin <- coxph(Surv(week, arrest) ~ sqrt(age) , data = recid)  
survminer::ggcoxfunctional(recid.lin,data=recid)
```



Non-Proportional Hazard Models

TESTS FOR PROPORTIONAL HAZARDS

Schoenfeld residuals for PH

Take a look at the time-dependent coefficients.

If coefficients do NOT depend upon time (i.e. PH holds....constant throughout time), then graphs should be a horizontal line

There is a score test that tests $H_0: \beta = 0$ versus $H_A: \beta \neq 0$ (we want to fail to reject H_0 to assume there is no relationship with time)

```
# Proportional Hazard Test - Schoenfeld Residuals
recid.ph.zph <- cox.zph(recid.ph, transform = "identity")
recid.ph.zph

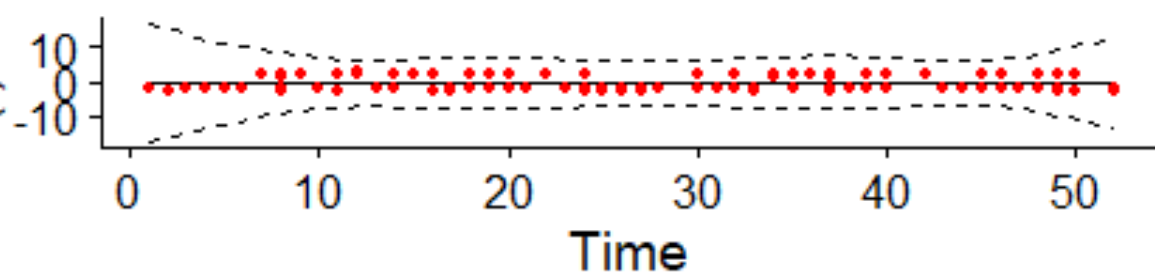
ggcoxzph(recid.ph.zph)
```

	chisq	df	p
fin	0.00878	1	0.9254
age	6.55134	1	0.0105
wexp	3.94780	1	0.0469
mar	1.04080	1	0.3076
paro	0.02280	1	0.8800
prio	0.42929	1	0.5123
GLOBAL	16.85230	6	0.0098

Global Schoenfeld Test p: 0.009842

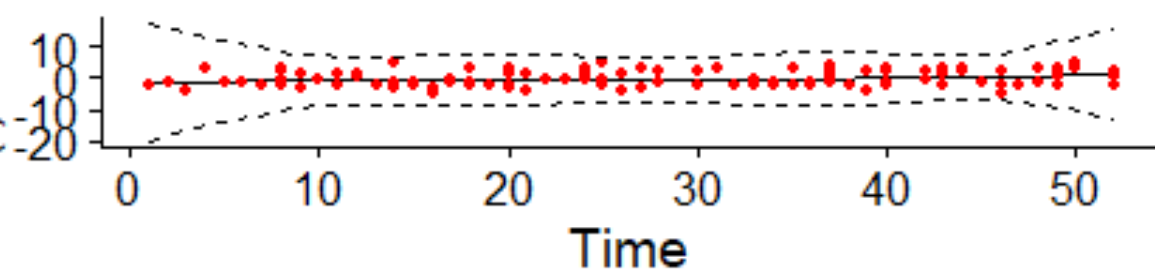
Beta(t) for fin

Schoenfeld Individual Test p: 0.9254



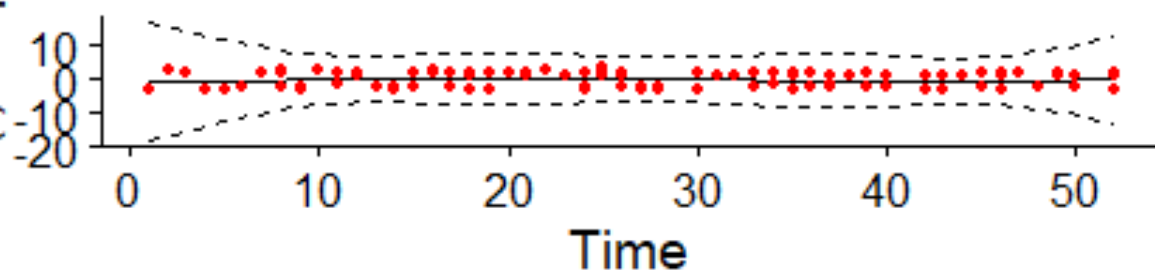
Beta(t) for wexp

Schoenfeld Individual Test p: 0.0469



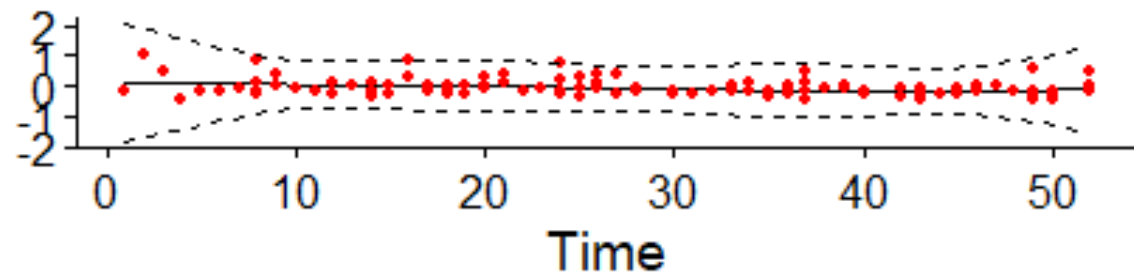
Beta(t) for paro

Schoenfeld Individual Test p: 0.88



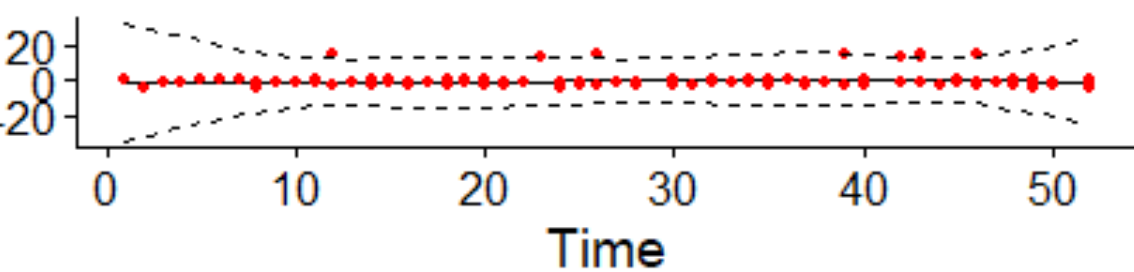
Beta(t) for age

Schoenfeld Individual Test p: 0.0105



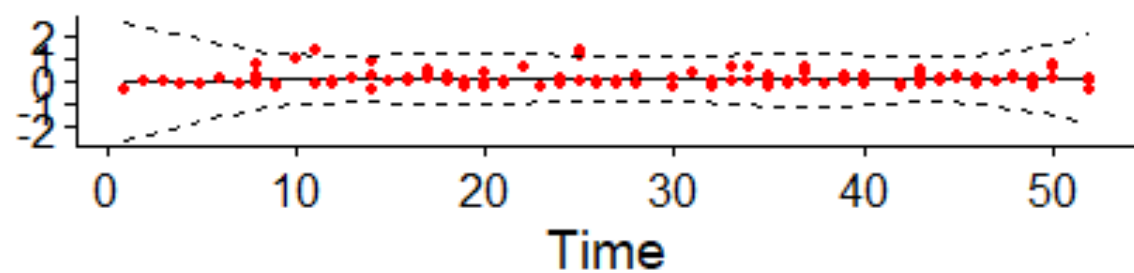
Beta(t) for mar

Schoenfeld Individual Test p: 0.3076



Beta(t) for prio

Schoenfeld Individual Test p: 0.5123



Non-Proportional Hazard Models

TIME-DEPENDENT COEFFICIENTS

Time Dependent Coefficients

Models up until this point have assumed that predictors have a constant effect, β , on the target variable.

In PH models, we assume effects are **constant over time**, so that the hazard ratio is independent of time.

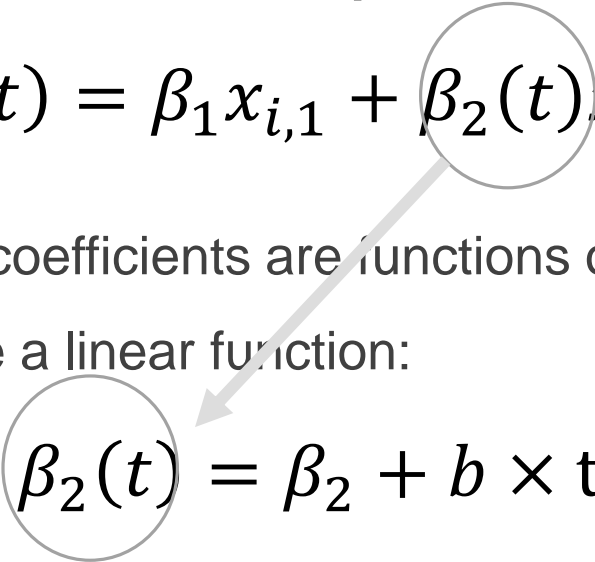
What if this didn't hold true and the effect of the predictor variable could change across time?

- Example: Does age have a constant effect throughout the study?

These effects, $\beta(t)$, are called **time-dependent coefficients**.

Time Dependent Coefficients

These effects, $\beta(t)$, are called **time-dependent coefficients**:

$$\log h(t) = \beta_1 x_{i,1} + \beta_2(t) x_{i,2}$$


These time-dependent coefficients are functions of time.

For example, it could be a linear function:

$$\beta_2(t) = \beta_2 + b \times \text{time}$$

If $b = 0$, then the effect doesn't depend on time (PH assumption satisfied).

If $b \neq 0$, then the effect **does** depend on time (PH assumption **not** satisfied).

Schoenfeld Residuals Again!

Schoenfeld residuals are best used for investigating relationships with time for predictor variables since they are calculated on a per variable basis.

You can plot these residuals against **functions** of time or the more popular technique would be to test the correlation (slope) between these residuals and **functions** of time.

Which functions?

- Common examples: t , $\log(t)$, K-M estimate, etc.

Proportional Hazard Test – R

```
recid.ph.zph <- cox.zph(recid.ph, transform = ...)  
recid.ph.zph
```

Fill with one of: “km”, “identity”, “log”, or “rank”

Proportional Hazard Test – R

		chisq	df	p
“km”	age	5.9161	1	0.015

		chisq	df	p
“identity”	age	6.55134	1	0.0105

		chisq	df	p
“log”	age	8.1533	1	0.0043

Proportional Hazard Test – R

“km”

		chisq	df	p
	age	5.9161	1	0.015

“identity”

		chisq	df	p
	age	6.55134	1	0.0105

“log”

		chisq	df	p
	age	8.1533	1	0.0043

Time Dependent Coefficients

If your software of choice tells you that you need one of these, what do you do?

Need to add these time-dependent coefficients, but luckily R can easily do this for you.

$$\log h(t) = \beta_1 x_{i,1} + \beta_2(t) x_{i,2}$$

Time Dependent Coefficients – R

```
recid.ph.tdc <- coxph(Surv(week, arrest) ~ fin + prio +  
                      wexp + mar + paro + age + tt(age),  
                      data = recid,  
                      tt = function(x, time, ...){x*log(time)})  
  
summary(recid.ph.tdc)
```


Time Dependent Coefficients – R

	coef	exp(coef)	se(coef)	z	Pr(> z)	
fin	-0.36527	0.69401	0.19087	-1.914	0.05566	.
wexp	-0.13317	0.87531	0.21247	-0.627	0.53080	
mar	-0.45279	0.63585	0.38041	-1.190	0.23394	
paro	-0.08490	0.91860	0.19534	-0.435	0.66382	
prio	0.09177	1.09611	0.02880	3.186	0.00144	**
age	0.12174	1.12946	0.06535	1.863	0.06249	.
tt(age)	-0.05931	0.94242	0.02182	-2.718	0.00658	**

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

exp(coef) exp(-coef) lower .95 upper .95

Time Dependent Coefficients – R

	coef	exp(coef)	se(coef)	z	Pr(> z)
fin	-0.36527	0.69401	0.19087	-1.914	0.05566 .
wexp	-0.13317	0.87531	0.21247	-0.627	0.53080
mar	-0.45279	0.63585	0.38041	-1.190	0.23394
paro	-0.08490	0.91860	0.19534	-0.435	0.66382
prio	0.09177	1.09611	0.02880	3.186	0.00144 **
age	0.12174	1.12946	0.06535	1.863	0.06249 .
tt(age)	-0.05931	0.94242	0.02182	-2.718	0.00658 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Interpretation

Let's use our example with age having a time-dependent coefficient:

$$\beta_{\text{age}}(t) = 0.122 - 0.059 \times \log(\text{week})$$

Initially, age has an increasing affect on the hazard (as age increases, the hazard is increasing since the coefficient for age is positive 0.122). This is true up to week 7.

However, as time goes on (week 8 and beyond), this effect becomes negative and being older decreases the hazard of recidivism.

WARNING!

This is **NOT** like creating a standard interaction with time for your predictor variable.

The interaction must be constructed in a way that updates **at each time**.

Trust R to do this for you instead of trying to create this yourself in the data sets.

Non-Proportional Hazard Models

TIME-DEPENDENT COVARIATES

A solid gray horizontal bar spanning the width of the slide at the bottom.

Time Dependent Variables

Similar to time-dependent coefficients, **time-dependent variables** have the actual value of the predictor variable (rather than its effect) change over time.

Time *independent* variable examples:

- Age (at entry)
- Number of prior convictions (at entry)

Time *dependent* variable examples:

- Employment status
- Blood pressure

Time-Dependent Variables

The following equation has one fixed variable and one time-dependent variable:

$$\log h(t) = \beta_1 x_{i,1} + \beta_2 x_{i,2}(t)$$

Prisoner Recidivism Data:

- EMP1 ~ EMP52 variables
 - Measure the full-time employment status during that week.
- Variables measured at same regular interval as response variable week of recapture.

Time-Dependent Variables

The following equation has one fixed variable and one time-dependent variable:

$$\log h(t) = \beta_1 x_{i,1} + \beta_2 x_{i,2}(t)$$

Prisoner Recidivism Data:

- EMP1 ~ EMP52 variables
 - Measure the full-time employment status during that week.
- Variables measured at same regular interval as response variable week of recapture.

Time-Dependent Variables

The following equation has one fixed variable and one time-dependent variable:

$$\log h(t) = \beta_1 x_{i,1} + \beta_2 x_{i,2}(t)$$

Prisoner Recidivism Data:

- EMP1 ~ EMP52 variables
 - Measure the full-time employment status during that week.
- Variables measured at same regular interval as response variable week of recapture.

Coding Time-Dependent Variables

Most important thing to remember with time-dependent variables →
FUTURE DATA CANNOT BE USED TO PREDICT THE PAST

Obvious right?!?!?!?

- So common it has its own name: **Immortal Time Bias**

Just make sure to make sure to structure data appropriately in all the following steps we learn.

Counting Process Structure

For time-dependent variables, it is necessary to split the *time* column of your data set into separate *start* and *stop* columns.

This is known as the **counting process** structure/layout to your data.

This is NEEDED for R to do the analysis.

Counting Process Example

Person 1 has an event at time = 9, but their value of x changes after time = 5.

Observe Person 1 until end of time = 5, after which they are censored:

Person	Start	Stop	x	Event
1	0	5	3	0

Create a “new” person starting after time = 5 who is the *exact same* as Person 1, but with new x value:

Person	Start	Stop	x	Event
1	0	5	3	0
1	5	9	7	1

Counting Process Example

Create a “new” person starting after time = 5 who is the *exact same* as Person 1, but with new x value:

Person	Start	Stop	x	Event
1	0	5	3	0
1	5	9	7	1

We observe this “new” person until either x changes again or their tenure ends (whichever comes first).

Fitting the Model

Most difficult part of modeling time-dependent variables is the formatting of the data correctly.

- Tedious, but usually straight-forward.
- Always print out some of the observations to make sure things look correct!

Everything else in modeling is essentially the same!

Estimates are not effected.

Time-Dependent Variables – R

```
recid_long.ph <- coxph(Surv(start, stop, arrested) ~ fin  
  + age + prio + employed, data = recid_long)  
  
summary(recid_long.ph)
```

Time-Dependent Variables – R

	coef	exp(coef)	se(coef)	z	Pr(> z)	
fin	-0.33051	0.71856	0.19012	-1.738	0.08214	.
age	-0.04977	0.95145	0.02053	-2.424	0.01537	*
prio	0.08364	1.08724	0.02775	3.014	0.00258	**
employed	-1.34815	0.25972	0.24928	-5.408	6.37e-08	***

Time-Dependent Covariates

There are some potential problems with time-dependent variables:

- Variables measured at different regular intervals than response variable.
- Variables measured at irregular time intervals.
- Variables that are undefined for certain intervals of time.

Typically, basic intuition is used for these calculations.