# Survival analysis with `R` (Part II)

*Juan R Gonzalez*

## Contents

## 1 Introduction

**Objectives**

- Understand the concept of survival analysis
- Learn how to perform survival analysis using Cox proportional hazard models with `R`
- Peform data analyses where the scientific question is to determine factors associated with time until event considering different covariates

## 2 Cox proportional hazard model

### 2.1 Single model

The survival experience of the cohort of patients depends on several variables, whose values have been recorded for each patient at the time origin. The proportional hazards model like other regression models allows explore the relationship between this survival experience (or hospital admission-free time of a patient) and explanatory variables. The focus is modeling the hazard function (hospital admission hazard or risk of death) at any time after the time origin of the study . The Cox regression model specifies the hazard for individual $i$ as:

$$\lambda_i(t) = \exp(\beta_1 x_{1i} + \beta_2 x_{2i} + \cdots \beta_p x_{pi})\lambda_0(t)$$

where $\lambda_0(t)$ is called *baseline hazard function* and $\beta_1, \beta_2, \cdots, \beta_p$ are coefficients to be determined.

Let us illustrate how to perform such an analysis by using Chronic Obstructive Pulmonary Disease (COPD) data that can be loaded from our repository.

```
copd <- read.table("../data/copd.txt", head=TRUE)
head(copd)
##   id lobenr enum phys.act fev smoke      age status.readmission
## 1  1     19    1        2 2.0     1 57.16359                  0
## 2  2     22    1        1 2.0     0 56.16427                  0
## 3  3     25    1        1 4.1     2 54.16838                  0
## 4  4     28    1        1 2.4     2 53.16359                  0
## 5  5     29    1        0 1.7     0 53.17180                  0
## 6  6     32    1        1 2.4     2 51.16496                  0
##   time.readmission status.death time.death
## 1            16141            0      16141
## 2            16116            1      16116
## 3            16141            0      16141
## 4            15073            1      15073
## 5            16141            0      16141
## 6            13479            1      13479
```

Our COPD database includes detailed information of a cohort of 2226 adults recruited from the general population, with repeated examinations every 5-10 years. During the follow-up period, hospital admissions may occur more than once for a given subject (recurrent even) and another kind of event like death (terminating event) can be observed as well. Despite the multivariate nature of these data, we start using univariate failure time methods with failure taken to be first occurrence of the recurring event (admission to hospital), where this event regarded as censored if the terminating event (death) occurs before any occurrence of the recurring event. In the next sessions this single approach will be compared with the multivariate one.

Four variables were considered for this study and collected at the moment of the first examination: `physical activity` (`phys.act`) categorized in three levels: low, moderate and high, a lung function test, `forced expiratory volume` (`fev`) measured with an electronic spirometer, the `smoking situation` (`smoke`), categorized in never, ex and current, and finally the age (`age`) of the patient. The pair of variables (`time.readmission`, `status.readmission`) depicts the censoring information about the first hospital admissions where death is regarded as censor, and the pair (`time.death`, `status.death`) describes the censoring information about the terminating event (death).

Let us then start by illustrating how survival analysis can be perfomed when our event of interest is hospital admission. We can fit the Cox proportional hazards model in `R` with the `coxph()` function (available in `survival` library). Let us investigate whether physical activity is associated with the time until hospital readmission. Next model can be considered as single analysis and it is similar to perform Kaplan-Meier estimates and log-rank test.

```
library(survival)
cox.pa <- coxph(Surv(time.readmission, status.readmission) ~
                as.factor(phys.act), data=copd)
cox.pa
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##     as.factor(phys.act), data = copd)
##
##                        coef exp(coef) se(coef)     z       p
## as.factor(phys.act)1 -0.444     0.642    0.110 -4.03 5.6e-05
## as.factor(phys.act)2 -0.671     0.511    0.126 -5.33 9.7e-08
##
## Likelihood ratio test=28.1  on 2 df, p=7.77e-07
## n= 2223, number of events= 495
##    (3 observations deleted due to missingness)
```

We can observe that those patients doing moderate and high physical activity (codes 1 and 2 respectively) have less risk of being readmitted than those doing low level of exercise. In all cases the p-value associated to each category are statistically significant at 5% level.
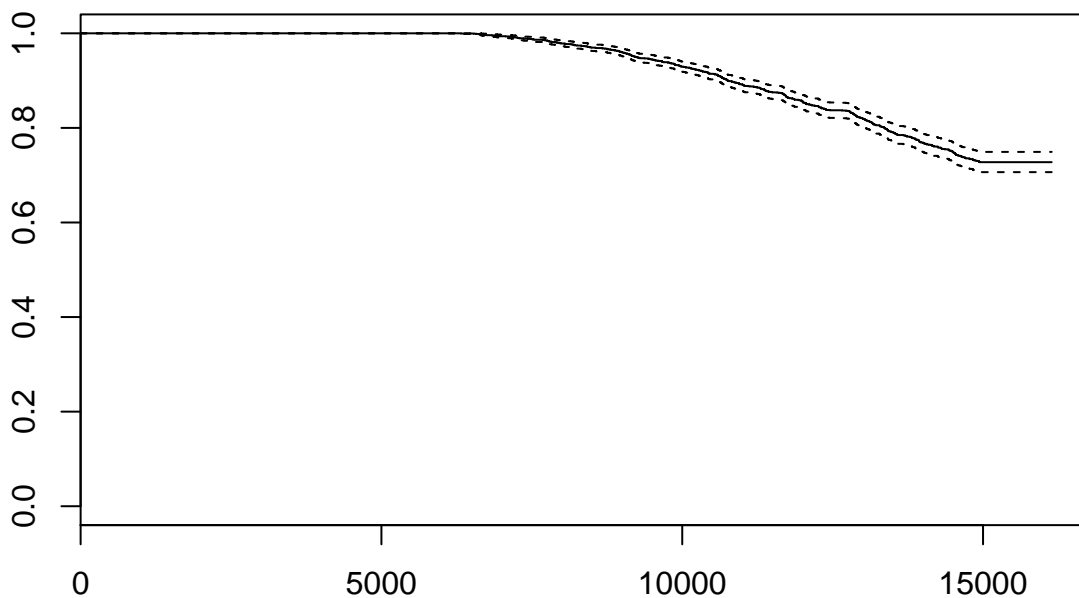
The exponentiated coefficients are interpretable as multiplicative effects on the hazard. `z` records the Wald statistic (ratio of each regression coefficient to its standard error) which is asymptotically standard normal under the hypothesis that the corresponding $\beta$ is zero. The likelihood ratio test, at the bottom of the output, is a overall test for the null hypothesis that all of the $\beta$'s are zero. Other equivalent overall tests (Wald test and Score test) could be obtained using `summary()` function.

```
summary(cox.pa)
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##     as.factor(phys.act), data = copd)
##
##   n= 2223, number of events= 495
##    (3 observations deleted due to missingness)
##
##                         coef exp(coef) se(coef)      z Pr(>|z|)
## as.factor(phys.act)1 -0.4438    0.6416   0.1102 -4.027 5.64e-05 ***
## as.factor(phys.act)2 -0.6709    0.5112   0.1258 -5.331 9.75e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##                      exp(coef) exp(-coef) lower .95 upper .95
## as.factor(phys.act)1    0.6416      1.559    0.5169    0.7962
## as.factor(phys.act)2    0.5112      1.956    0.3995    0.6543
##
## Concordance= 0.563  (se = 0.012 )
## Rsquare= 0.013   (max possible= 0.963 )
## Likelihood ratio test= 28.14  on 2 df,   p=7.768e-07
## Wald test            = 29.77  on 2 df,   p=3.44e-07
## Score (logrank) test = 30.53  on 2 df,   p=2.349e-07
```

The fit shows that our variable has an impact on the hospital admission-free time of the patients. In particular, the fitted model estimates that that people with moderate and high level of physical activity reduces the risk of hospital admission by 36% ($= 1 - 0.64$) and 49% ($= 1 - 0.51$) respectively compared to people with low level of physical activity.

Having fit a Cox model to the data, it is often of interest to examine the estimated distribution of survival times. The function `survfit()` estimates $S(t)$ and its confidence intervals.

```
plot(survfit(cox.pa))
```

## 2.2   Multivariate model

In some ocassion one might be interested in adjusting the results by other covariates. This ends up with fitting multivariate models that can be estimated as following:

```
cox.pa.adj <- coxph(Surv(time.readmission, status.readmission) ~
                    as.factor(phys.act) + age + fev +
                    as.factor(smoke), data=copd)
cox.pa.adj
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##     as.factor(phys.act) + age + fev + as.factor(smoke), data = copd)
##
##
##                        coef exp(coef) se(coef)      z       p
## as.factor(phys.act)1 -0.33589   0.71470  0.11117  -3.02  0.0025
## as.factor(phys.act)2 -0.29579   0.74394  0.12861  -2.30  0.0215
## age                   0.02435   1.02465  0.00518   4.70 2.6e-06
## fev                  -0.99877   0.36833  0.07646 -13.06 < 2e-16
## as.factor(smoke)1     1.50586   4.50802  0.28295   5.32 1.0e-07
## as.factor(smoke)2     2.02998   7.61390  0.25582   7.94 2.1e-15
##
## Likelihood ratio test=389  on 6 df, p=0
## n= 2206, number of events= 492
##    (20 observations deleted due to missingness)
```

Here, one may see that the effect of doing physical activity is not so strong as in the univariate case. The adjusted model

shows that the hazard ratio of doing high level of activity is 0.74 that means that the probability of being readmitted to the hospital is reduced in only 26% with respect to those people who do low level of physical activity. This risk is much lower than the one estimated using the single model probably due to the fact that the single HR was somehow confused by age, fev or skoming status. This makes perfect sense since smoking is associated with copd and probably those people who smoke normally do less exercise.
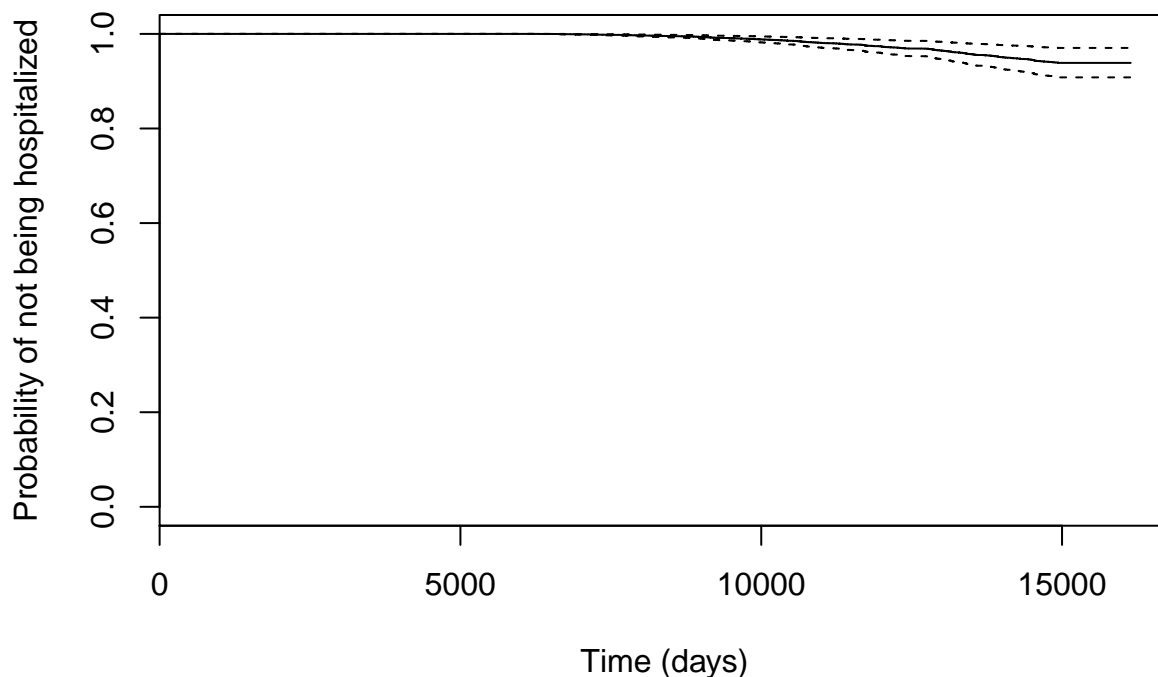
In that case one may also be interested in examining the estimated distribution of survival times. The function `survfit()` estimates $S(t)$ and confidence intervals, by default at the mean values of the covariates. This is obtained with the follow command lines:

```r
copd$smoke2 <- 1*(copd$smoke==1)
copd$smoke3 <- 1*(copd$smoke==2)
copd$phys.act2 <- 1*(copd$phys.act==1)
copd$phys.act3 <- 1*(copd$phys.act==2)

cox2.pa.adj <- coxph(Surv(time.readmission, status.readmission) ~
                         age + fev + smoke2 + smoke3 + phys.act2 +
                         phys.act3, data=copd)

newdata <- data.frame(age = mean(copd$age),
                         fev = mean(copd$fev, na.rm=TRUE),
                         smoke2 = 0, smoke3=0,
                         phys.act2=0, phys.act3=0)

plot(survfit(cox2.pa.adj, newdata= newdata),
          xlab="Time (days)",
          ylab="Probability of not being hospitalized")
```

This figure represents the estimated survival function for reference patients: non-smokers, with low level of physical activity, aged 54 years old and 2.4 of level of forced expiratory volume

# 3   Model selection

Model selection can be carried out by using an automatic medhos by using the function `stepAIC()` from `MASS` library. We start by considering the model having all of our variables of interest

```
library(MASS)
cox.pa.adj
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##     as.factor(phys.act) + age + fev + as.factor(smoke), data = copd)
##
##
##                         coef exp(coef) se(coef)      z        p
## as.factor(phys.act)1 -0.33589   0.71470  0.11117  -3.02   0.0025
## as.factor(phys.act)2 -0.29579   0.74394  0.12861  -2.30   0.0215
## age                   0.02435   1.02465  0.00518   4.70 2.6e-06
## fev                  -0.99877   0.36833  0.07646 -13.06 < 2e-16
## as.factor(smoke)1     1.50586   4.50802  0.28295   5.32 1.0e-07
## as.factor(smoke)2     2.02998   7.61390  0.25582   7.94 2.1e-15
##
## Likelihood ratio test=389  on 6 df, p=0
## n= 2206, number of events= 492
##    (20 observations deleted due to missingness)
```

Then, the procedure is performed by executing

```
mod <- stepAIC(cox.pa.adj)
## Start:  AIC=6870.59
## Surv(time.readmission, status.readmission) ~ as.factor(phys.act) +
##     age + fev + as.factor(smoke)
## Error in dropterm.default(fit, scope$drop, scale = scale, trace = max(0, : number of rows in use has cha
```

Here you see an error message that is due to the fact that missing data are present. One solution can be to impute data using another package like `mice`. Here, let's illustrate how to perform stepwise using complete cases.

```
copd.complete <- copd[complete.cases(copd),]
cox.complete <- coxph(Surv(time.readmission, status.readmission) ~
                      as.factor(phys.act) + age + fev +
                      as.factor(smoke), data=copd.complete)
mod <- stepAIC(cox.complete)
## Start:  AIC=6870.59
## Surv(time.readmission, status.readmission) ~ as.factor(phys.act) +
##     age + fev + as.factor(smoke)
##
##                        Df    AIC
## <none>                     6870.6
## - as.factor(phys.act)  2 6875.7
## - age                  1 6891.2
## - as.factor(smoke)     2 6992.0
## - fev                  1 7055.5
```

We observe that all variables are statistically significant after performing stepwise procedure. This method basically compares nested models by using likelihood ratio test (LRT). Here you can see an example of LRT to compare a model

having *age*, *phys.act* and *smoke* with the comple one ( *age*, *phys.act* , *smoke* and *fev* )

```
cox1 <- coxph(Surv(time.readmission, status.readmission) ~
                    as.factor(phys.act) + age + as.factor(smoke),
                    data=copd.complete)
cox2 <- coxph(Surv(time.readmission, status.readmission) ~
                    as.factor(phys.act) + age + as.factor(smoke) +
                    fev, data=copd.complete)
anova(cox1, cox2)
## Analysis of Deviance Table
##  Cox model: response is  Surv(time.readmission, status.readmission)
##  Model 1: ~ as.factor(phys.act) + age + as.factor(smoke)
##  Model 2: ~ as.factor(phys.act) + age + as.factor(smoke) + fev
##    loglik  Chisq Df P(>|Chi|)
## 1 -3522.7
## 2 -3429.3 186.86  1 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Here we see that the p-value obtained with LRT is pretty similart to the one obtained using Wald test (e.g. using z-score)

```
cox2
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##     as.factor(phys.act) + age + as.factor(smoke) + fev, data = copd.complete)
##
##
##                          coef exp(coef) se(coef)       z        p
## as.factor(phys.act)1 -0.33589   0.71470  0.11117  -3.02  0.0025
## as.factor(phys.act)2 -0.29579   0.74394  0.12861  -2.30  0.0215
## age                   0.02435   1.02465  0.00518   4.70 2.6e-06
## as.factor(smoke)1     1.50586   4.50802  0.28295   5.32 1.0e-07
## as.factor(smoke)2     2.02998   7.61390  0.25582   7.94 2.1e-15
## fev                  -0.99877   0.36833  0.07646 -13.06 < 2e-16
##
## Likelihood ratio test=389  on 6 df, p=0
## n= 2206, number of events= 492
```

# 4   Validation and diagnostic

The aim of this section is to determine if the fitted model adequately describes the data. The model-checking procedure involves four kinds of diagnostics: for assessment of model fit, for checking the functional form of covariates, for identification of influential observations and for violation of the assumption of proportional hazards. All these model-checking procedures are based on the follow *residuals*: the martingale, deviance, score, Schoenfeld, dfbeta and scaled Schoenfeld residuals. This section gives a brief definition of these residuals and how use them for the model-checking procedure.

## 4.1   Assessment of model fit

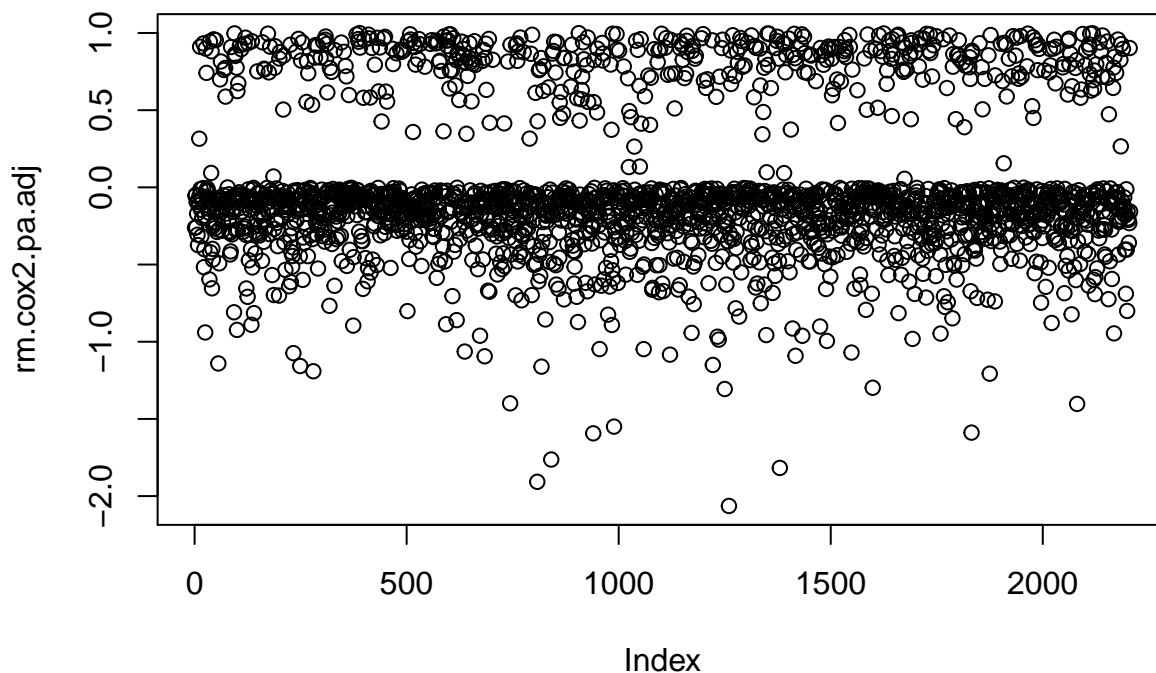The *martingale* residual for the $i$th individual is given by the expression,

$$r_{M_i} = \delta_i - \hat{\Lambda}(t_i)$$

where $\delta_i$ takes the value 0 if the observation is censored and the value 1 if it is a failure. *Martingale* residuals may be interpreted as the difference between the observed and the expected number of failures in the time interval $(0, t_i)$. So,

a plot of these residuals will highlight those individuals with a bigger difference, and, consequently, the residuals will highlight individuals whose survival time is not been well fitted by the model (*outliers*). Plots of these residuals against explanatory variables can be interesting for indicate whether there are values of explanatory variables where the model does not fit well.

The martingale residuals are the default output of `residuals()` on a `coxph()` fit,

```
rm.cox2.pa.adj<- resid(cox2.pa.adj)
plot(rm.cox2.pa.adj)
```
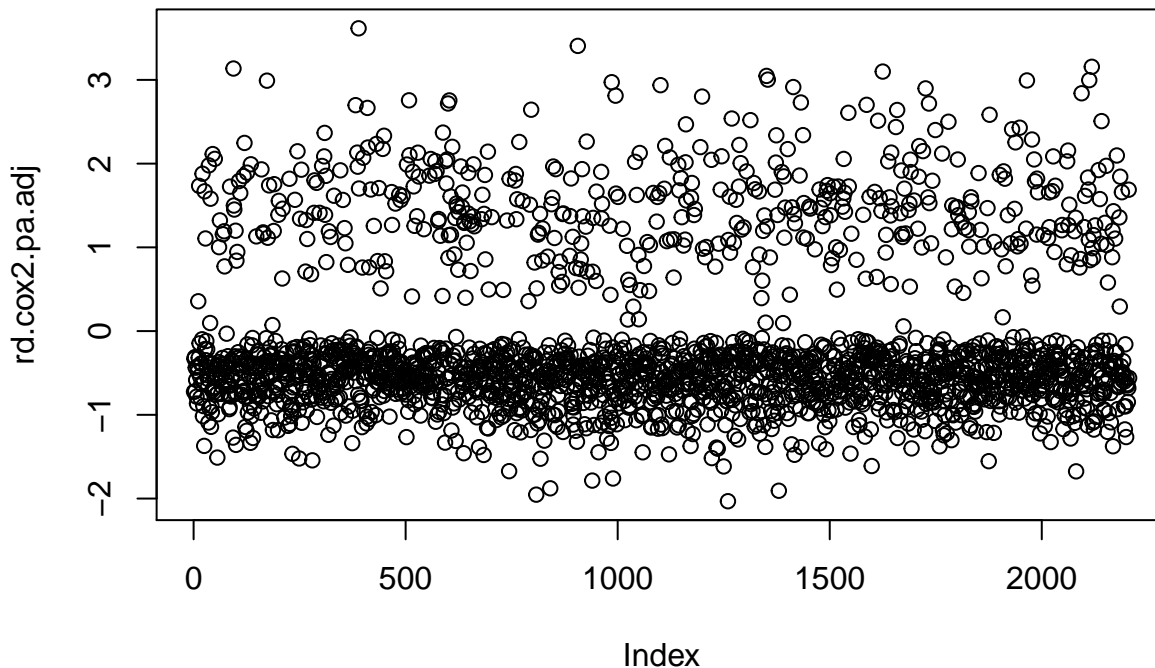


The *Deviance* residual for the $i$th individual is defined by:

$$r_{D_i} = \text{sign}(r_{M_i})[-2\{r_{M_i} + \delta_i \log{(\delta_i - r_{M_i})}\}]^{\frac{1}{2}},$$

where $\text{sgn}(\cdot)$ is the sign function, which takes the value $+1$ if its argument is positive and the value -1 if it is negative.

This kind of residuals are a transformation of *Martingale* residuals and generate values that are symmetric around zero when the fitted model is appropriate. They are also useful to detect *outliers*. Plotting *deviance* residuals against *Risk Score*, we may also detect those individuals with risk of failure below the mean value (*Risk Score* very negative), and those above it (high *Risk Score*). Using the function `resid()` and the option `type="deviance` we can obtain these residuals from the fitted model. For martingale and deviance residuals, the returned object is a vector with one element for each subject,

```
rd.cox2.pa.adj<- resid(cox2.pa.adj, type="deviance")
plot(rd.cox2.pa.adj)
```
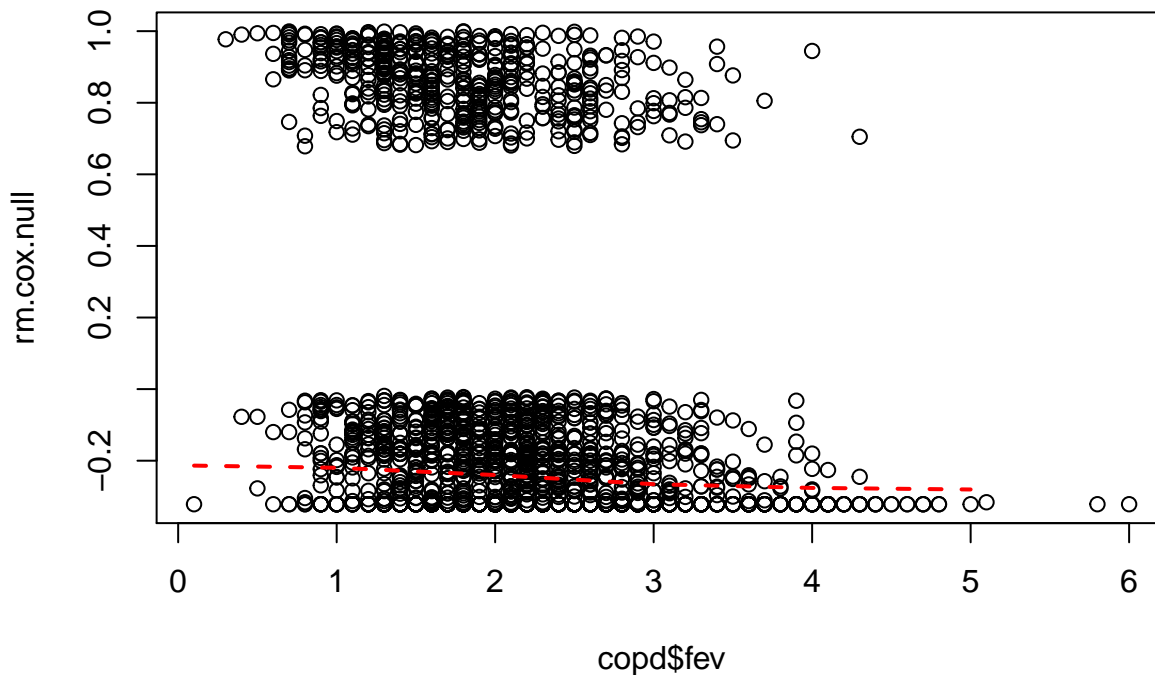
In both plots we observe as there is no outliers and that all the residuals are constant.

## 4.2   Funcional form of covariates

An improvement in the fit of a model may well be obtained by using some transformation of the values of a variable instead of the original values. The simplest approach is one examined by Therneau and Grambsch, who suggested plotting the martingale residuals obtained from fitting the null model (model that contains no variables) against the variable, and superimposing a smoothed scatterplot. This plot should display the functional form required for the variable. In particular, a straight line plot indicates that a linear term is needed.

The following code creates null residual plots for the forced expiratory volume variable:

```
cox.null <- coxph(Surv(time.readmission, status.readmission) ~ 1,
                  data = copd)
rm.cox.null <- resid(cox.null)
plot(copd$fev, rm.cox.null)
smooth <- lowess(copd$fev, rm.cox.null, delta=1)
lines(smooth, lty=2, lwd=2, col="red")
```

We observe that there is a linear relationship between residuals and our variable of interest, and hence, no further transformation are needed.

## 4.3   Testing Proportional Hazards

We know that hazards are said to be proportional if ratios of hazards are independent of time. If there are one or more explanatory variables in the model whose coefficients vary with time, or if there are explanatory variables that are time-dependent, the proportional hazards assumption will be violated. So, it is required a method to detect this possibility: if there is some form of time dependency in particular variables.

The tests and graphical diagnostics for proportional hazards are based on the *scaled Schoenfeld residuals*, $r^*_{Pji}\}$, and are useful in evaluating the assumption of proportional hazards after fitting a Cox regression model.

Therneau and Grambsch show that the expected value of the $i$th *scaled Schoenfeld residual* is given by $E\left(r^*_{Pji}\right) \approx \hat{\beta}_j\left(t_i\right) - \hat{\beta}_j$, and so a plot of the values of $r^*_{Pji} + \hat{\beta}_j$ against the event times should give information about the form of the time-dependent coefficient of $X_j$, $\beta_j\left(t\right)$.

The interpretation of these graphs is greatly facilitated by smoothing shown on each graph by a solid line. An horizontal line in each graph indicates no suggestion of non-proportional hazards and that the coefficients of these variables are constant (see next figure).

This graphical diagnostic is supplemented by a test for each variable, along with a global test for the model as a whole. These tests for the proportional-hazards assumption are obtained from `cox.zph()`,
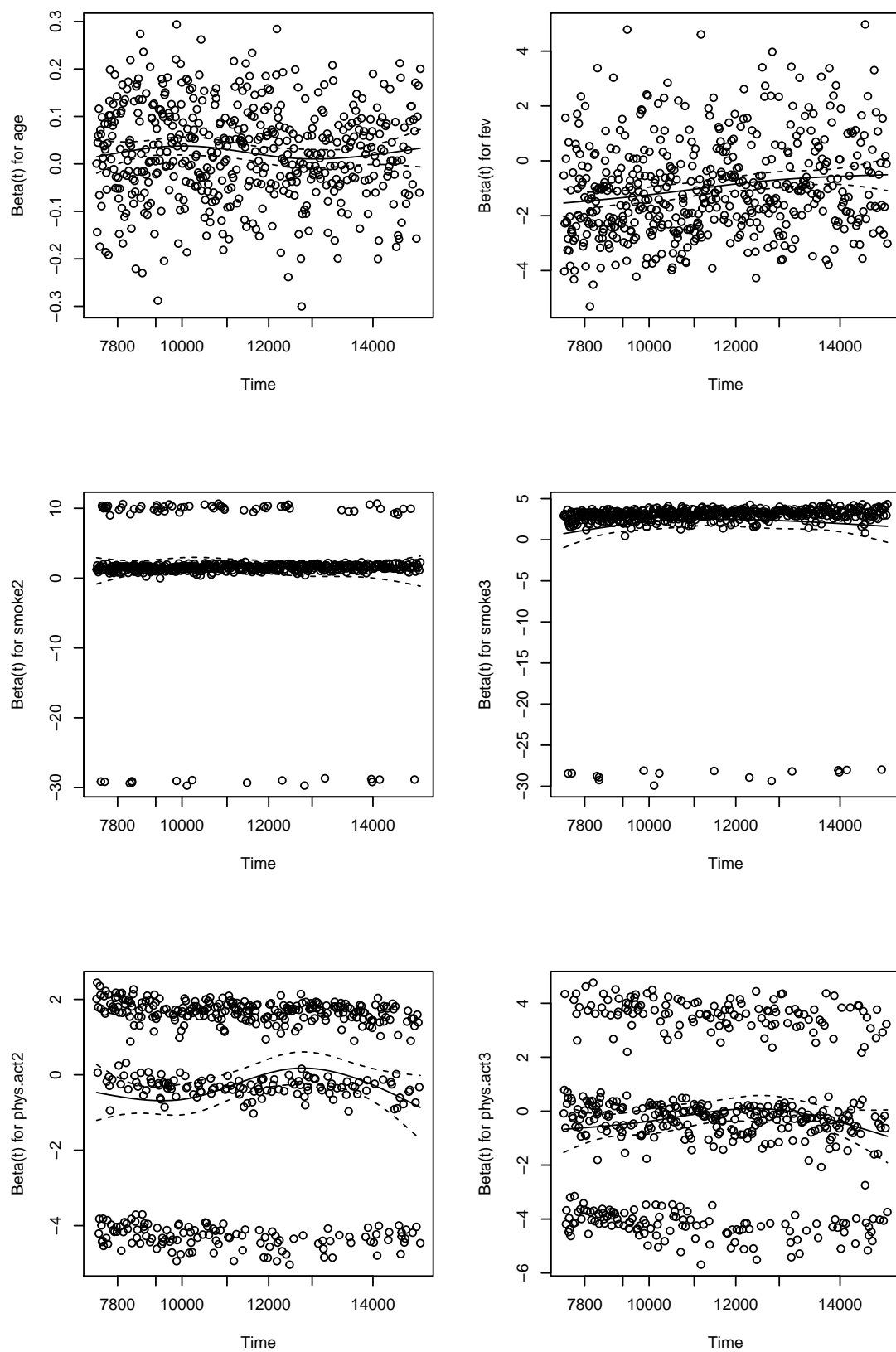
```
cox.pa.adj.zph <- cox.zph(cox2.pa.adj)
cox.pa.adj.zph
##              rho  chisq        p
## age      -0.0417  0.654 4.19e-01
```

```
## fev          0.1910 18.642 1.58e-05
## smoke2      -0.0234  0.269 6.04e-01
## smoke3       0.0368  0.661 4.16e-01
## phys.act2 0.0642  2.022 1.55e-01
## phys.act3 0.0264  0.341 5.59e-01
## GLOBAL            NA 32.061 1.59e-05
```

Here `rho` is the Pearson product-moment correlation between the *scaled Schoenfeld residuals* and time for each variable. The column `chisq` gives the tests statistics for each variable and the last row GLOBAL gives the global test for a $\chi^2$ of 6 degree of freedom. There is strong evidence of non proportionality hazard for *fev*, the GLOBAL test gives also strong evidence of non-proportionality of **at least** one covariate.

Plotting the object returned by `cox.zph` produces graphs of the scaled Schoenfeld residuals against time,

```
par(mfrow=c(3,2))
plot(cox.pa.adj.zph)
```

We observe as *fev* is not constant across time (we see and increase in the beta value when time increases - in other words, the risk is increasing with time) and the Cox model assumes that the risk is constant over time.

The solution to this problem is to include a time-dependent variable into Cox model, or to stratify the initial model according to the variable that violates the assumption of proportional hazard. In that case the variable *fev* is continuos. Therefore, we can create categories by dissecting the variable into intervals. We decided to create four intervals or categories defined from the quartiles of the variable,

```
copd$fev4 <- cut(copd$fev, 4)
table(copd$fev4)
##
## (0.0941,1.58]    (1.58,3.05]    (3.05,4.53]    (4.53,6.01]
##          502           1444            253             11
```

A stratified Cox regression model is fit by including the `strata` term on the right hand side of the model formula. The strata divide the individuals into disjoint groups, each of which has a distinct baseline hazard function but common values for the coefficient vector $\beta$,

```
cox.pa.strata <- coxph(Surv(time.readmission, status.readmission) ~
                  age + as.factor(smoke) + as.factor(phys.act) +
                  strata(fev4), data=copd)
cox.pa.strata
## Call:
## coxph(formula = Surv(time.readmission, status.readmission) ~
##      age + as.factor(smoke) + as.factor(phys.act) + strata(fev4),
##      data = copd)
##
##
##                         coef exp(coef) se(coef)      z       p
## age                  0.02746   1.02784  0.00518   5.31 1.1e-07
## as.factor(smoke)1    1.44096   4.22474  0.28300   5.09 3.5e-07
## as.factor(smoke)2    1.96300   7.12067  0.25566   7.68 1.6e-14
## as.factor(phys.act)1 -0.35796   0.69910  0.11111  -3.22  0.0013
## as.factor(phys.act)2 -0.37952   0.68419  0.12891  -2.94  0.0032
##
## Likelihood ratio test=147  on 5 df, p=0
## n= 2206, number of events= 492
##      (20 observations deleted due to missingness)
```

And now the test of proportionality is doing well:

```
cox.pa.strata.zph <- cox.zph(cox.pa.strata)
cox.pa.strata.zph
##                         rho chisq      p
## age                  -0.0424 0.678 0.4104
## as.factor(smoke)1    -0.0252 0.311 0.5769
## as.factor(smoke)2     0.0333 0.538 0.4634
## as.factor(phys.act)1  0.0635 1.984 0.1589
## as.factor(phys.act)2  0.0241 0.287 0.5923
## GLOBAL                   NA 9.262 0.0991
```

There is no evidence of non-proportional hazards for the remaining covariates. An advantage of this approach is that it gives most general adjustment for a confounding variable. A disadvantage is that no estimate of the effects of the stratifying covariate is produced.

Further information about survival data analysis with `R` can be found in this tutorial Tutorial Survival Analysis.

# 5    Exercise (to deliver)

---

Data for exercises are in the repository https://github.com/isglobal-brge/TeachingMaterials/tree/master/Longitudinal_data_analysis/data

File *pulmon.sav* contains data about a survival study about lung cancer (NOTE: data can be loaded into `R` by using `read.spss` function available at `foreign` library - use argument *to.data.frame=TRUE* ). Colums contain this information:

- TIEMPO Supervivencia (meses)
- ESTADO: 0 VIVO, 1 MORT
- EDAD4 Age at diagnosis in years (quartiles)
- SEXO: HOMBRES, MUJERES
- ESTCLIN Estadio clinico: EST 0/I, EST II, EST IIIA, EST IIIB, EST IV
- IK Indice de estado general (100 estado perfecto, 0 muerte)
- CIRUGIA: 1 No operado, 2 Cirugia no radical, 3 Cirugia Radical
- QUIMIO: 1 No Quimio, 2 Platino
- RADIOTER: 1 No RT, 2 <60 Gy, 3 >60 Gy

**Exercise 1:**

- Estimate a separate Cox model of each variable (univariate analysis)
- Select those variables that are statistically significant in the previous step and fit a multivariate model including all of them
- Is this a good model to predict survival in lung cancer?
- If not, estimate the best model by using an automatic method (e.g. stepwise)
- Does this model hold Cox model assumption?

---

# 6    References

- The [survival] package (https://cran.r-project.org/web/packages/survival/)

# 7    Session information

```
## R version 3.3.2 (2016-10-31)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 14393)
##
## locale:
## [1] LC_COLLATE=Spanish_Spain.1252  LC_CTYPE=Spanish_Spain.1252
## [3] LC_MONETARY=Spanish_Spain.1252 LC_NUMERIC=C
## [5] LC_TIME=Spanish_Spain.1252
##
## attached base packages:
## [1] stats     graphics  grDevices utils     datasets  methods   base
##
## other attached packages:
## [1] MASS_7.3-45    survival_2.40-1 knitr_1.15.1    BiocStyle_2.2.1
##
## loaded via a namespace (and not attached):
```

```
##  [1] Rcpp_0.12.9     lattice_0.20-34 digest_0.6.11   rprojroot_1.2
##  [5] grid_3.3.2      backports_1.0.5 magrittr_1.5    evaluate_0.10
##  [9] stringi_1.1.2   Matrix_1.2-7.1  rmarkdown_1.3   splines_3.3.2
## [13] tools_3.3.2     stringr_1.2.0   yaml_2.1.14     htmltools_0.3.5
```