# Conditional Inference Trees & Cox Regression to Predict Heart Failure Survival Time

Antonio Pano

11/10/2022

# https://archive.ics.uci.edu/ml/datasets/Heart+failure+clinical+record (https://archive.ics.uci.edu/ml/datasets/Heart+failure+clinical+recor

- All 299 patients had left ventricular systolic dysfunction

Initial Variables:

- age: age of the patient (years)
- anaemia: decrease of red blood cells or hemoglobin since last measure (boolean)
- high blood pressure: if the patient has hypertension (boolean)
- creatinine phosphokinase (CPK): level of the CPK enzyme in the blood (mcg/L)
- diabetes: if the patient has diabetes (boolean)
- ejection fraction: percentage of blood leaving the heart at each contraction (percentage)
- platelets: platelets in the blood (kiloplatelets/mL)
- sex: woman or man (binary)
- serum creatinine: level of serum creatinine in the blood (mg/dL)
- serum sodium: level of serum sodium in the blood (mEq/L)
- smoking: if the patient smokes or not (boolean)
- time: follow-up period (days)
- [target] death event: if the patient deceased during the follow-up period (boolean)

```
library(skimr)
library(ggplot2)
library(dplyr)
library(tidyr)
library(survival)
library(survminer)
library(partykit)
library(coin)
library(survminer)
library(flexsurv)
library(randomForestSRC)
library(broom)
library(gtsummary)
```

Loading in the data

Creating Left Ventricular Ejection Fraction Groups set by Cardiology Experts (https://www.ncbi.nlm.nih.gov/books/NBK459131/). Rounding for averages instead of only using data for men and women.

```
HF <- read.csv("heart_failure_clinical_records_dataset.csv")

HF$anaemia = as.factor(HF$anaemia)
HF$diabetes = factor(HF$diabetes,levels=c(0,1),labels=c("Absent","Present"))
HF$hypertension = factor(HF$high_blood_pressure,levels=c(0,1),labels=c("Absent","Present"))

HF$sex = factor(HF$sex,levels=c(0,1),labels=c("Female","Male"))
HF$smoking = as.factor(HF$smoking)
HF$DEATH_EVENT = as.factor(HF$DEATH_EVENT)

HF <- HF %>%
  mutate(EF_Condition = cut(HF$ejection_fraction, breaks = c(0, 30, 40, 52, Inf),
         labels = c("Severe", "Moderate", "Mild", "Normal"), include.lowest = TRUE))


HF <- select(HF, -high_blood_pressure)

skim(HF)
```

Data summary

| Name | HF |
|---|---|
| Number of rows | 299 |
| Number of columns | 14 |
| _____ | |
| Column type frequency: | |
| factor | 7 |
| numeric | 7 |
| _____ | |
| Group variables | None |

**Variable type: factor**

| skim_variable | n_missing | complete_rate | ordered | n_unique | top_counts |
|---|---|---|---|---|---|
| anaemia | 0 | 1 | FALSE | 2 | 0: 170, 1: 129 |
| diabetes | 0 | 1 | FALSE | 2 | Abs: 174, Pre: 125 |
| sex | 0 | 1 | FALSE | 2 | Mal: 194, Fem: 105 |
| smoking | 0 | 1 | FALSE | 2 | 0: 203, 1: 96 |
| DEATH_EVENT | 0 | 1 | FALSE | 2 | 0: 203, 1: 96 |
| hypertension | 0 | 1 | FALSE | 2 | Abs: 194, Pre: 105 |
| EF_Condition | 0 | 1 | FALSE | 4 | Mod: 126, Sev: 93, Mil: 41, Nor: 39 |

**Variable type: numeric**

| skim_variable | n_missing | complete_rate | mean | sd | p0 | p25 | p50 | p75 | p100 | hist |
|---|---|---|---|---|---|---|---|---|---|---|
| age | 0 | 1 | 60.83 | 11.89 | 40.0 | 51.0 | 60.0 | 70.0 | 95.0 | ▄▆▇▄▁ |
| creatinine_phosphokinase | 0 | 1 | 581.84 | 970.29 | 23.0 | 116.5 | 250.0 | 582.0 | 7861.0 | █▁▁▁▁ |
| ejection_fraction | 0 | 1 | 38.08 | 11.83 | 14.0 | 30.0 | 38.0 | 45.0 | 80.0 | ▃▇▃▁▁ |
| platelets | 0 | 1 | 263358.03 | 97804.24 | 25100.0 | 212500.0 | 262000.0 | 303500.0 | 850000.0 | ▂▇▂▁▁ |
| serum_creatinine | 0 | 1 | 1.39 | 1.03 | 0.5 | 0.9 | 1.1 | 1.4 | 9.4 | █▁▁▁▁ |
| serum_sodium | 0 | 1 | 136.63 | 4.41 | 113.0 | 134.0 | 137.0 | 140.0 | 148.0 | ▁▁▃▇▁ |
| time | 0 | 1 | 130.26 | 77.61 | 4.0 | 73.0 | 115.0 | 203.0 | 285.0 | ▇▆▇▅▆ |

Correlation

Time and Serum_Creatinine have a correlation to Serum_Sodium of 0.15 & 0.19, respectively.

```
cormat <- HF %>% select(where(is.numeric)) %>% cor() %>% round(2)
melted_cormat <- reshape2::melt(cormat)

ggplot(data = melted_cormat, aes(x=Var1, y=Var2, fill=value)) +
  geom_tile() +
  geom_text(aes(Var2, Var1, label = value), color = "white", size = 4) +
   theme(axis.title.x=element_blank(),
       axis.title.y=element_blank(),
       axis.text.x = element_text(angle = 15, vjust = 0.8)
       )
```
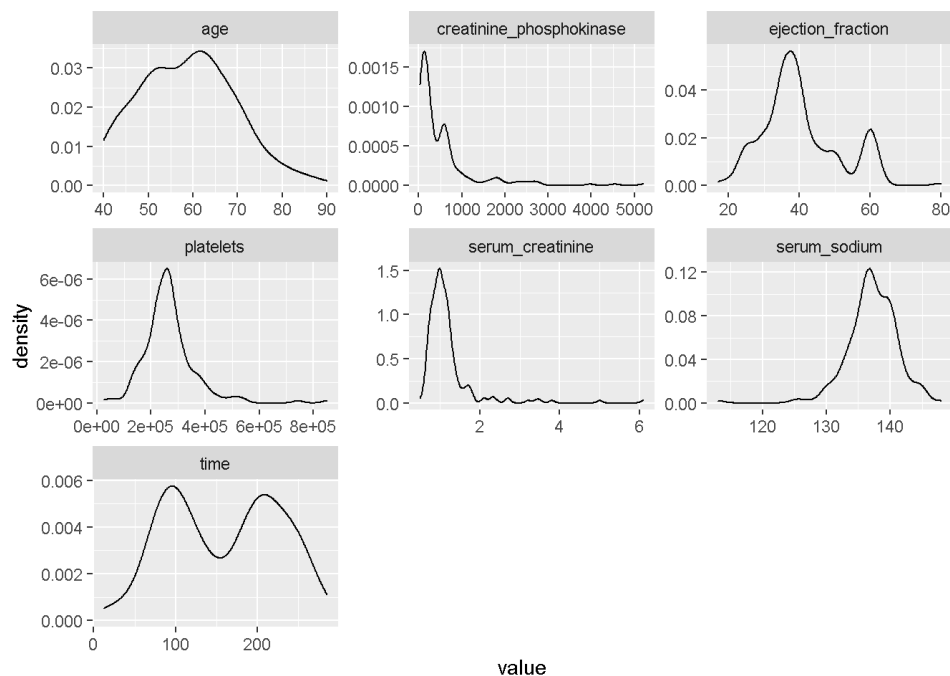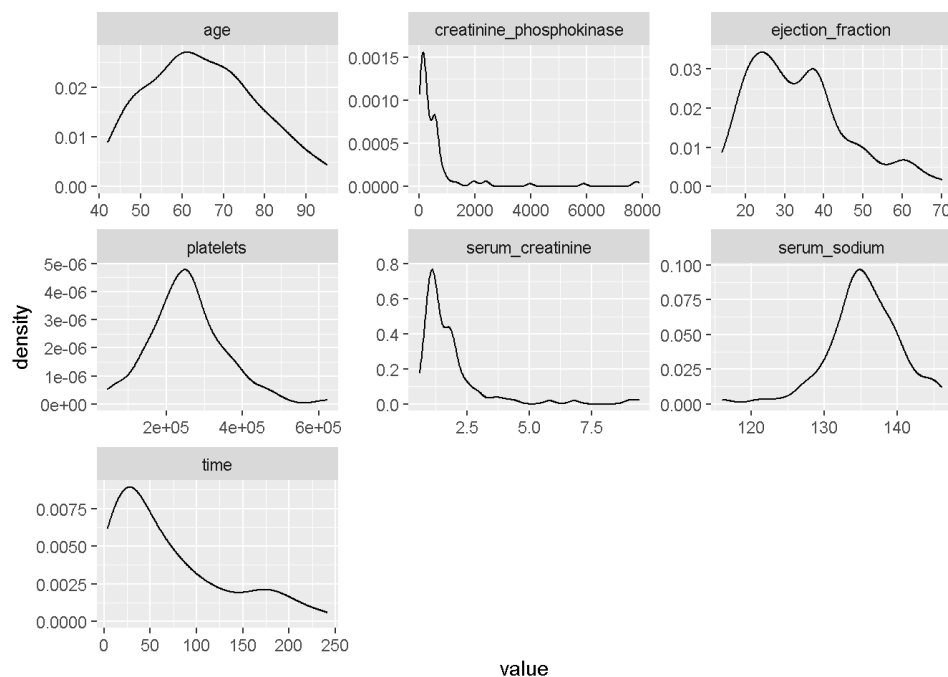
Choosing to grab distributions based on having hypertension– what's traditionally seen as a good indicator of heart failure.

Doing so to look at, specifically, Ejection Fraction right after to see if there is correlation.

```
HF %>% filter(DEATH_EVENT==0) %>%
  purrr::keep(is.numeric) %>%
  gather() %>%
  ggplot(aes(value)) +
    facet_wrap(~ key, scales = "free") +
    geom_density()
```
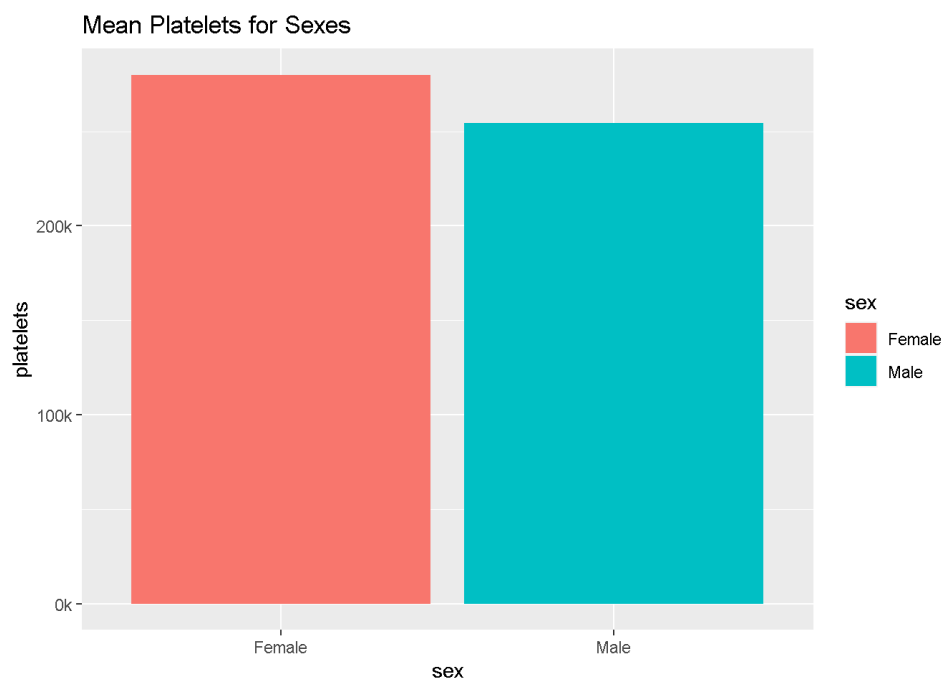


```
HF %>% filter(DEATH_EVENT==1) %>%
  purrr::keep(is.numeric) %>%
  gather() %>%
  ggplot(aes(value)) +
    facet_wrap(~ key, scales = "free") +
    geom_density()
```

Comparing `creatinine_phosphokinase` to Men & Women– those who smoke and those who do not.

- Noticing that the average `creatinine_phosphokinase` is higher for non-smokers.
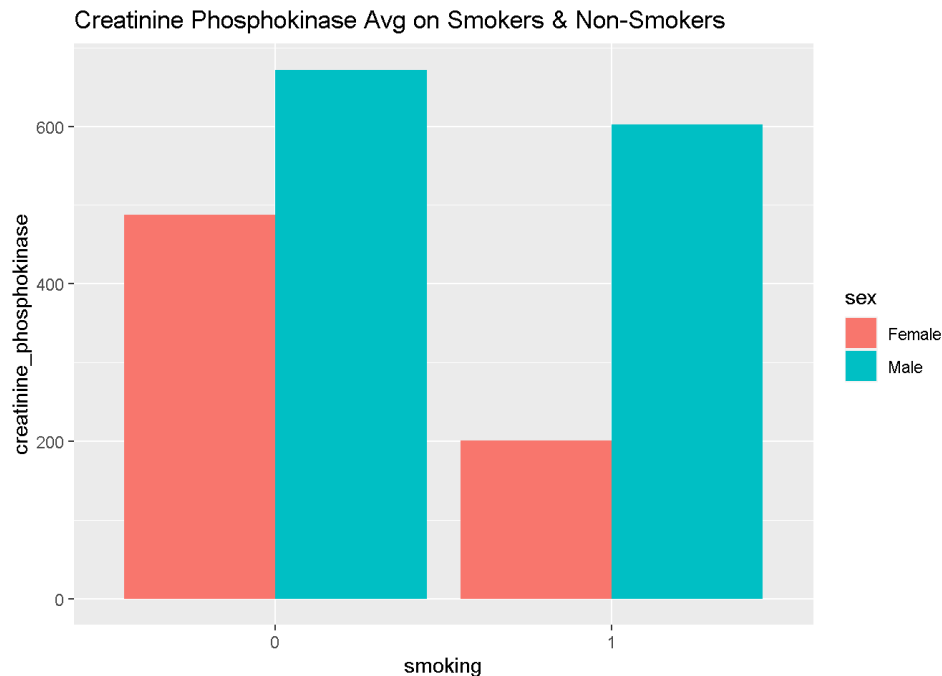
```
ggplot(HF, aes(x=sex, y=platelets, fill=sex)) +
  geom_bar(position = "dodge", stat="summary", fun="mean") +
  scale_y_continuous(labels = scales::label_number(suffix = "k", scale = 1e-3)) +
  ggtitle("Mean Platelets for Sexes")
```



```
HF %>% group_by(sex, DEATH_EVENT) %>%
  summarize(count = n(), .groups="drop")
```

```
## # A tibble: 4 × 3
##   sex    DEATH_EVENT count
##   <fct>  <fct>       <int>
## 1 Female 0              71
## 2 Female 1              34
## 3 Male   0             132
## 4 Male   1              62
```
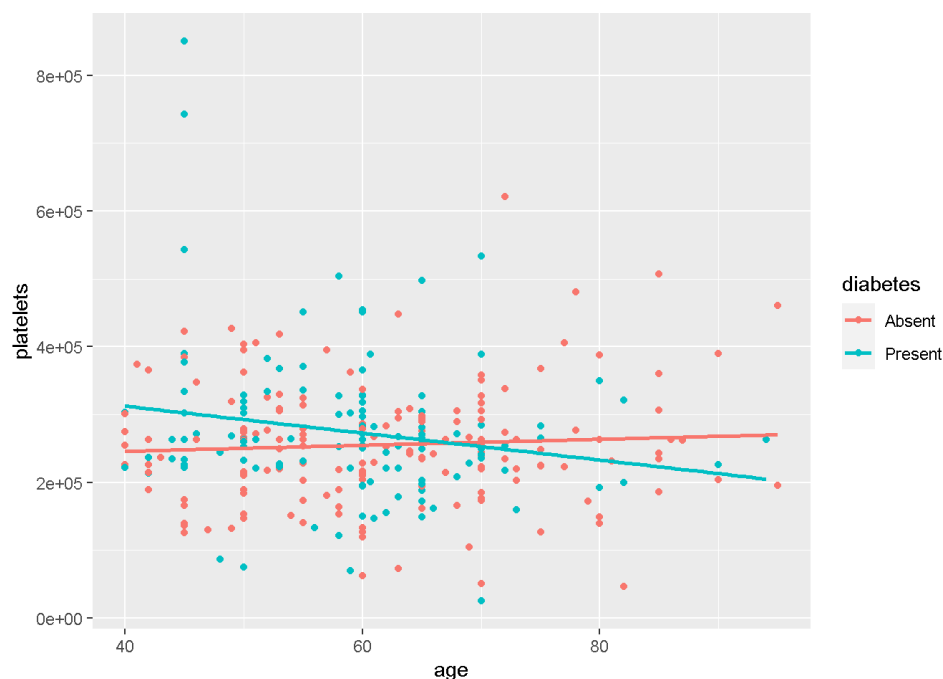
```
ggplot(HF, aes(x=smoking, y=creatinine_phosphokinase, fill=sex)) +
  geom_bar(position = "dodge", stat="summary", fun="mean") +
  ggtitle("Creatinine Phosphokinase Avg on Smokers & Non-Smokers")
```



- Finding out that for those diabetic, plateletes reduce as age increases.
- For those who aren't diabetic, plateletes generally stay the same and potentially, increase by a marginal amount.

Plateletes are incredibly important. Having too few plateletes can lead to internal bleeding in intestines or stroke.

```
ggplot(HF, aes(x=age, y=platelets,color=diabetes)) + geom_point() +
  geom_smooth(method='lm', se = FALSE)
```



# Conditional Inference Trees - Kaplan Maeier Curves

We can see we have remaining cases in which the person did was not declared deceased due to the ending of the curve not dropping down to 0%.

Insights from this graph include: * `Serum Creatinine` is highly significant with the showcased split at 1.8 for survival prediction.
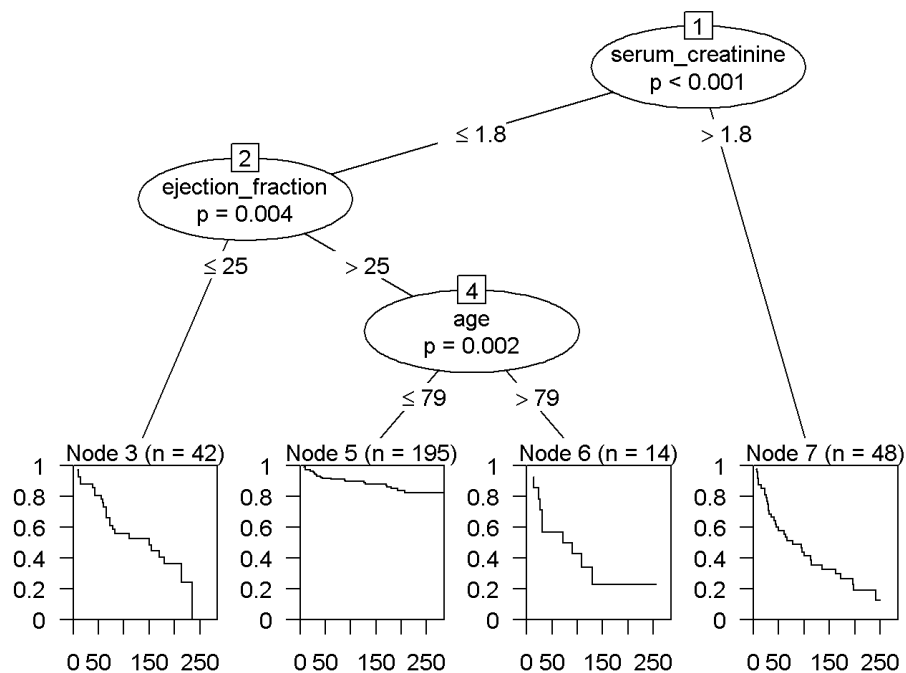
```
set.seed(0)

# Won't directly go from factor to numeric. Needed for Survival Analysis.
HF$DEATH_EVENT = as.numeric(as.character(HF$DEATH_EVENT))

# Dropping categorical Ejection Fraction.
HF <- HF %>% select(-EF_Condition)


# Creating a Conditional Inference Tree for descriptive analytics
CondInfTree <- ctree(Surv(time, DEATH_EVENT) ~ .,
                     data = HF)


plot(CondInfTree)
```



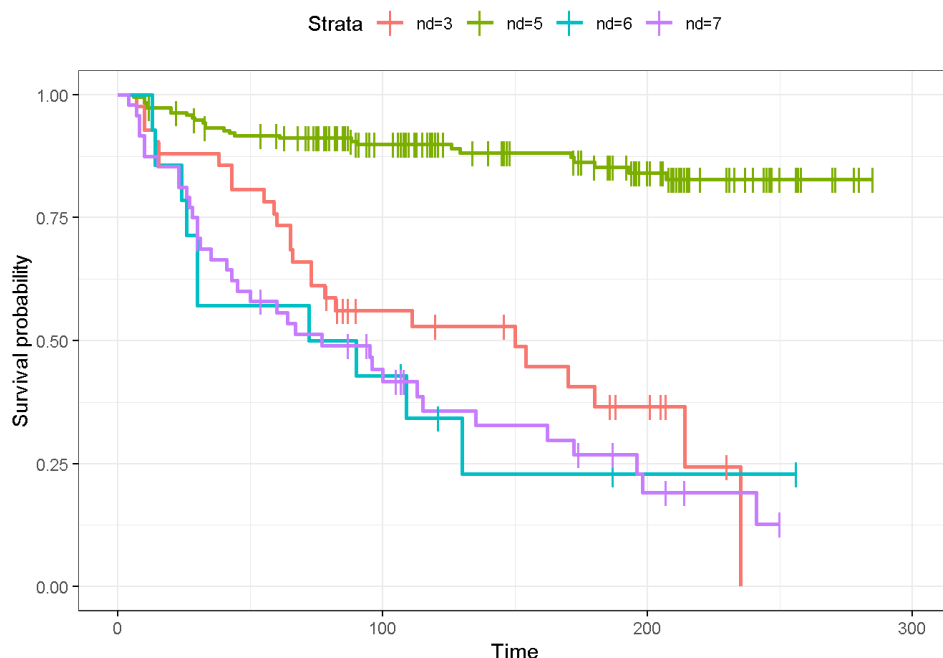Plotting all node distributions/curves in one plot.

```
nd <- factor(predict(CondInfTree, type = "node"))

all_nd <- survfit(Surv(time, DEATH_EVENT) ~ nd, data = HF)

ggsurvplot(all_nd, data = HF,
           censor.shape="|",
           conf.int = FALSE, #surv.median.line = "hv",
           ggtheme = theme_bw())
```

EXTRACTING SURVIVAL CURVE FOR ONLY ONE OBSERVATION/PERSON FROM THE CTREE! PULL OUT AT LEAST ONE INSIGHT. THE 'X' DETERMINES WHICH OBSERVATION YOU'LL LOOK AT. PERHAPS LOOK AT AN OUTLIER TO TALK ABOUT A SPECIAL CASE.

```
#nd1 <- predict(CondInfTree, type = "prob")[[X]]
#summary(nd1, times=c(20, 45, 60, 80, 100, 10*(11:15)))
```

Constructing an exponential curve for previous graph's second node. * 24% probability of survival after t=130 days for patients older than 79, that have less than or equal to 1.8 in serum creatine, and an ejection fraction over 25.

```
K <- HF %>%
  filter(serum_creatinine <= 1.8, ejection_fraction > 25, age > 79)


# This one is best.
# The ~ 1 is our way ofletting R know that we aren't using any x variables. Just time and whether event occured which are both
 y variabes.
pred_k_surv <- survfit(Surv(time, DEATH_EVENT) ~ 1, data = K)

summary(pred_k_surv, times=c(20, 45, 60, 80, 100, 10*(11:15)))
```

```
## Call: survfit(formula = Surv(time, DEATH_EVENT) ~ 1, data = K)
##
##  time n.risk n.event survival std.err lower 95% CI upper 95% CI
##    20     12       2    0.857  0.0935       0.6921        1.000
##    45      8       4    0.571  0.1323       0.3630        0.899
##    60      8       0    0.571  0.1323       0.3630        0.899
##    80      7       1    0.500  0.1336       0.2961        0.844
##   100      6       1    0.429  0.1323       0.2341        0.785
##   110      4       1    0.343  0.1307       0.1624        0.724
##   120      4       0    0.343  0.1307       0.1624        0.724
##   130      3       1    0.229  0.1277       0.0765        0.683
##   140      2       0    0.229  0.1277       0.0765        0.683
##   150      2       0    0.229  0.1277       0.0765        0.683
```

- No pruning was done since most trees found revolve around the same 3 variables.
- Probability of survival after 150 days for those younger than 70 is 77%.
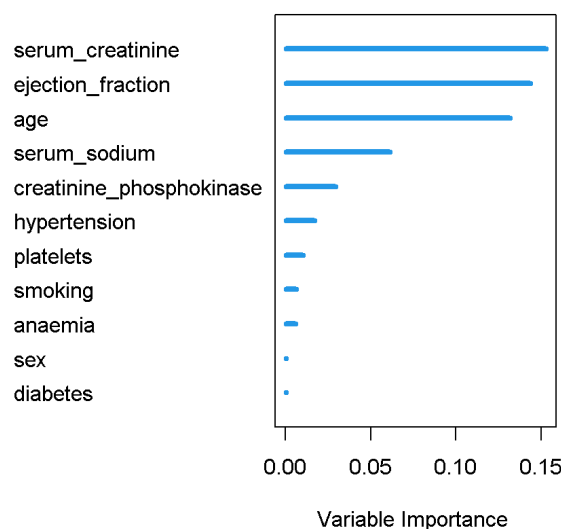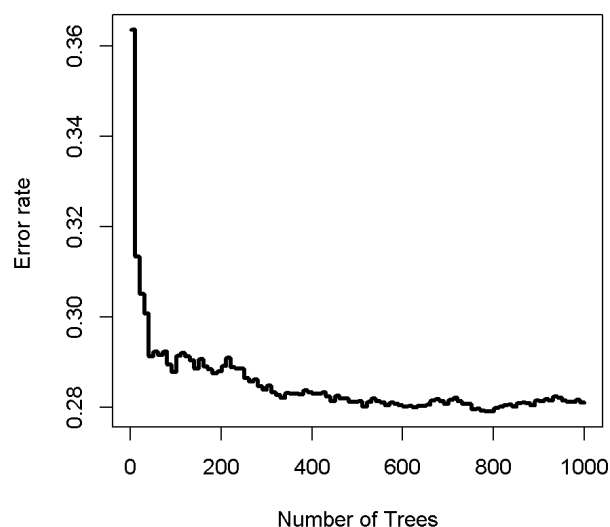- Probability of survival after 200 days for those younger than 70 is 70%.

```
survfit(Surv(time, DEATH_EVENT) ~ 1, data = HF %>% filter(age <= 70)) %>%
  tbl_survfit(
    times = c(150,200),
    label_header = "**{time} Day Survival (95% CI) For Those Younger Than 70**"
    )
```

| Characteristic | 150 Day Survival (95% CI) For Those Younger Than 70 | 200 Day Survival (95% CI) For Those Younger Than 70 |
|---|---|---|
| Overall | 77% (71%, 82%) | 70% (64%, 77%) |

# Random Forest Survival

```
# mtry means how many nodes at each split
fit <- rfsrc(Surv(time, DEATH_EVENT) ~ .,
           data = HF,
           ntree = 1000,
           importance = TRUE,
           nsplit = 3,
           mtry = 2)

#fit
plot(fit)
```



```
##
##                            Importance   Relative Imp
## serum_creatinine              0.1527        1.0000
## ejection_fraction             0.1435        0.9399
## age                           0.1319        0.8635
## serum_sodium                  0.0614        0.4020
## creatinine_phosphokinase      0.0295        0.1929
## hypertension                  0.0170        0.1112
## platelets                     0.0103        0.0674
## smoking                       0.0062        0.0404
## anaemia                       0.0059        0.0385
## sex                           0.0002        0.0014
## diabetes                      0.0001        0.0006
```

# Cox Proportional Hazards Model (Cox Regression)

KM will make the curve based on event & time but that's all. We need to include the rest of the variables.

- At a given instance in time, someone who has hypertension is 0.42 times as likely to die as someone without hypertension adjusting for age.
- At any given instance in time, someone who does *not* have hypertension is 0.65 times as likely to die as someone who does, adjusting for age.
- Concordance: Goodness of fit for survival analysis.

```
# diabetes isn't stat significant.
coxMod1 <- coxph(Surv(time, DEATH_EVENT) ~ diabetes, data=HF)
summary(coxMod1)
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ diabetes, data = HF)
##
##   n= 299, number of events= 96
##
##                     coef exp(coef) se(coef)      z Pr(>|z|)
## diabetesPresent -0.04184   0.95902  0.20728 -0.202     0.84
##
##                 exp(coef) exp(-coef) lower .95 upper .95
## diabetesPresent     0.959      1.043    0.6388      1.44
##
## Concordance= 0.502  (se = 0.027 )
## Likelihood ratio test= 0.04  on 1 df,    p=0.8
## Wald test            = 0.04  on 1 df,    p=0.8
## Score (logrank) test = 0.04  on 1 df,    p=0.8
```
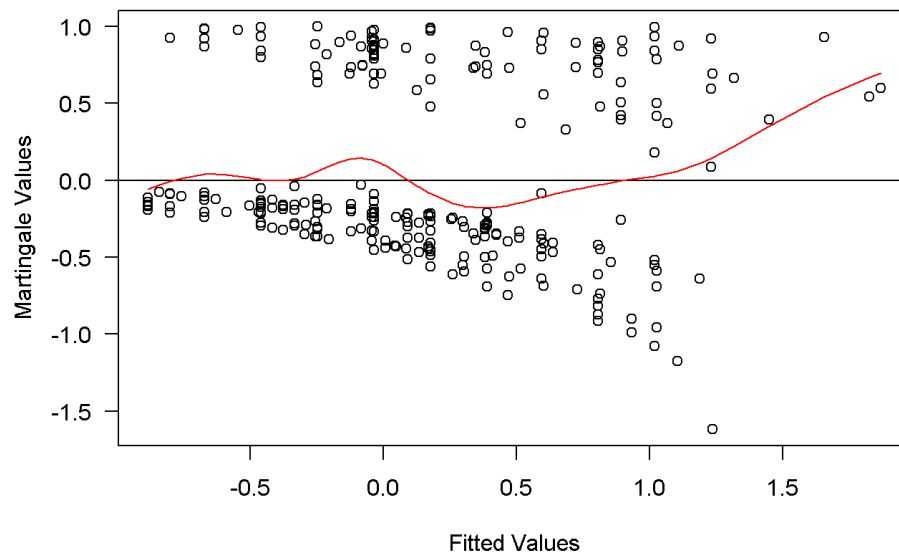
```
# hypertension useful bc tree didn't output it. i paired it w/ age bc why not? historia mejor.
coxMod2 <- coxph(Surv(time, DEATH_EVENT) ~ hypertension + age, data=HF)
summary(coxMod2)
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ hypertension + age,
##     data = HF)
##
##   n= 299, number of events= 96
##
##                         coef exp(coef) se(coef)      z Pr(>|z|)
## hypertensionPresent 0.417717  1.518491 0.209708 1.992   0.0464 *
## age                 0.042424  1.043336 0.008693 4.880 1.06e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##                     exp(coef) exp(-coef) lower .95 upper .95
## hypertensionPresent     1.518     0.6585     1.007     2.290
## age                     1.043     0.9585     1.026     1.061
##
## Concordance= 0.638  (se = 0.031 )
## Likelihood ratio test= 27.36  on 2 df,    p=1e-06
## Wald test            = 27.52  on 2 df,    p=1e-06
## Score (logrank) test = 28.25  on 2 df,    p=7e-07
```

```
# so long as most part of red doesn't stray, it's linear. This one strays a lot at end bc of less values overall so they hold m
ore weight.
plot(predict(coxMod2), residuals(coxMod2, type = "martingale"), xlab = "Fitted Values",
     ylab = "Martingale Values", main = "Residual Plot", las = 1) +
  abline(h=0) +
  lines(smooth.spline(predict(coxMod2), residuals(coxMod2, type="martingale")), col="red")
```

## Residual Plot



```
## integer(0)
```

```
# Do the Likelihood-Ratio Test
# Try to find combination that may be insightful and make them "oh? interesting".
anova(coxMod1, coxMod2, test = "LRT")
```

```
## Analysis of Deviance Table
##  Cox model: response is  Surv(time, DEATH_EVENT)
##  Model 1: ~ diabetes
##  Model 2: ~ hypertension + age
##    loglik  Chisq Df P(>|Chi|)
## 1 -509.18
## 2 -495.52 27.322  1 1.722e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# keep only variables that are significant. do manual stepwise, basically, and see what u get.
# boom, found it.
summary(coxph(Surv(time, DEATH_EVENT) ~ ., data=HF))
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ ., data = HF)
##
##   n= 299, number of events= 96
##
##                              coef  exp(coef)   se(coef)       z Pr(>|z|)
## age                      4.641e-02  1.048e+00  9.324e-03   4.977 6.45e-07 ***
## anaemia1                 4.601e-01  1.584e+00  2.168e-01   2.122   0.0338 *
## creatinine_phosphokinase 2.207e-04  1.000e+00  9.919e-05   2.225   0.0260 *
## diabetesPresent          1.399e-01  1.150e+00  2.231e-01   0.627   0.5307
## ejection_fraction       -4.894e-02  9.522e-01  1.048e-02  -4.672 2.98e-06 ***
## platelets               -4.635e-07  1.000e+00  1.126e-06  -0.412   0.6806
## serum_creatinine         3.210e-01  1.379e+00  7.017e-02   4.575 4.76e-06 ***
## serum_sodium            -4.419e-02  9.568e-01  2.327e-02  -1.899   0.0575 .
## sexMale                 -2.375e-01  7.886e-01  2.516e-01  -0.944   0.3452
## smoking1                 1.289e-01  1.138e+00  2.512e-01   0.513   0.6078
## hypertensionPresent      4.757e-01  1.609e+00  2.162e-01   2.201   0.0278 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##                          exp(coef) exp(-coef) lower .95 upper .95
## age                         1.0475     0.9547    1.0285     1.067
## anaemia1                    1.5843     0.6312    1.0358     2.423
## creatinine_phosphokinase    1.0002     0.9998    1.0000     1.000
## diabetesPresent             1.1501     0.8695    0.7427     1.781
## ejection_fraction           0.9522     1.0502    0.9329     0.972
## platelets                   1.0000     1.0000    1.0000     1.000
## serum_creatinine            1.3786     0.7254    1.2014     1.582
## serum_sodium                0.9568     1.0452    0.9141     1.001
## sexMale                     0.7886     1.2681    0.4816     1.291
## smoking1                    1.1376     0.8790    0.6953     1.861
## hypertensionPresent         1.6092     0.6214    1.0534     2.458
##
## Concordance= 0.741  (se = 0.027 )
## Likelihood ratio test= 81.95  on 11 df,    p=6e-13
## Wald test            = 87.27  on 11 df,    p=6e-14
## Score (logrank) test = 88.39  on 11 df,    p=3e-14
```

```
summary(coxph(Surv(time, DEATH_EVENT) ~ age+anaemia+creatinine_phosphokinase+ejection_fraction+
              serum_creatinine+hypertension, data=HF))
```
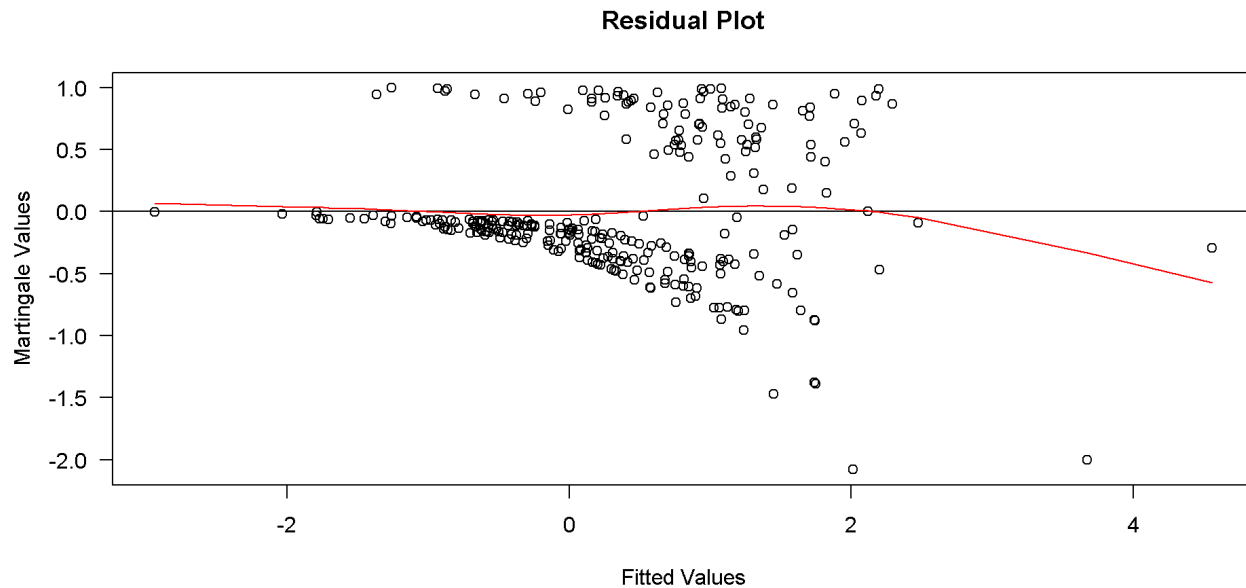
```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ age + anaemia + creatinine_phosphokinase +
##     ejection_fraction + serum_creatinine + hypertension, data = HF)
##
##   n= 299, number of events= 96
##
##                              coef  exp(coef)   se(coef)       z Pr(>|z|)
## age                      4.361e-02  1.045e+00  8.853e-03   4.926 8.41e-07 ***
## anaemia1                 3.933e-01  1.482e+00  2.129e-01   1.847   0.0648 .
## creatinine_phosphokinase 1.965e-04  1.000e+00  9.856e-05   1.993   0.0462 *
## ejection_fraction       -5.179e-02  9.495e-01  1.005e-02  -5.152 2.57e-07 ***
## serum_creatinine         3.483e-01  1.417e+00  6.550e-02   5.318 1.05e-07 ***
## hypertensionPresent      4.668e-01  1.595e+00  2.129e-01   2.192   0.0284 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##                          exp(coef) exp(-coef) lower .95 upper .95
## age                         1.0446     0.9573    1.0266    1.0629
## anaemia1                    1.4818     0.6749    0.9762    2.2493
## creatinine_phosphokinase    1.0002     0.9998    1.0000    1.0004
## ejection_fraction           0.9495     1.0531    0.9310    0.9684
## serum_creatinine            1.4167     0.7059    1.2460    1.6108
## hypertensionPresent         1.5948     0.6270    1.0506    2.4209
##
## Concordance= 0.738  (se = 0.028 )
## Likelihood ratio test= 77.02  on 6 df,    p=1e-14
## Wald test            = 85.82  on 6 df,    p=2e-16
## Score (logrank) test = 83.51  on 6 df,    p=7e-16
```

Checking Linearity of Model * Linearity of the final cox regression is sufficient.

SAY SOMETHING ABOUT THE HAZARD RATIO PLOT

```
sigMod = coxph(Surv(time, DEATH_EVENT) ~ age+anaemia+creatinine_phosphokinase+ejection_fraction+
                serum_creatinine+hypertension, data=HF)

plot(predict(sigMod), residuals(sigMod, type = "martingale"), xlab = "Fitted Values",
     ylab = "Martingale Values", main = "Residual Plot", las = 1) +
  abline(h=0) +
  lines(smooth.spline(predict(sigMod), residuals(sigMod, type="martingale")), col="red")
```
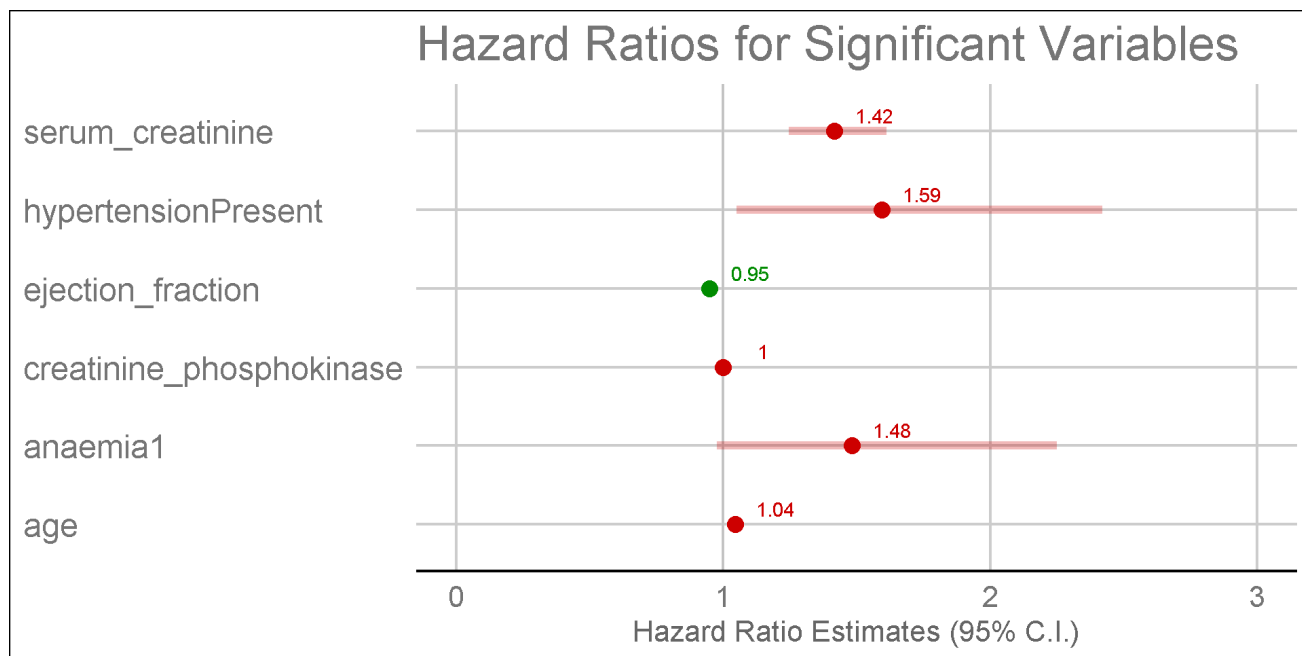
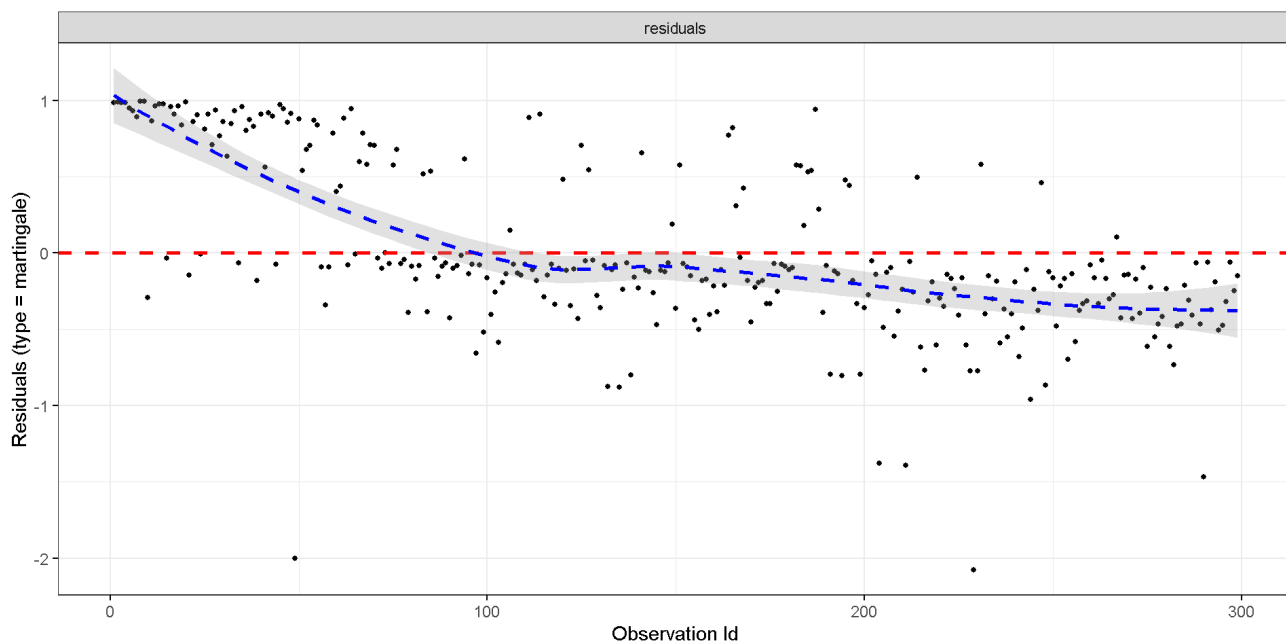## Residual Plot



```
## integer(0)
```

```
#ggforest(sigMod, data = HF)

library(ggthemes)
finMod <- sigMod %>% tidy()

finMod %>%  mutate(upper = estimate + 1.96 * std.error,
         lower = estimate - 1.96 * std.error) %>%
  mutate(across(all_of(c("estimate", "lower", "upper")), exp)) %>%
  ggplot(aes(estimate, term, color = estimate > 1)) +
  geom_vline(xintercept = 1, color = "gray75") +
  geom_linerange(aes(xmin = lower, xmax = upper), size = 2.25, alpha = 0.28) +
  geom_point(size = 4) +
  theme_gdocs(base_size = 16) +
  scale_color_manual(values = c("green4", "red3"), guide = "none") +
  xlim(c(0, 3)) +
  labs(title = "Hazard Ratios for Significant Variables", y = NULL,
       x = "Hazard Ratio Estimates (95% C.I.)") +
  theme(axis.text.y = element_text(hjust = 0, size = 18)) +
    geom_text(label = exp(finMod$estimate) %>% round(2),
              nudge_y = .2, nudge_x = .15)
```
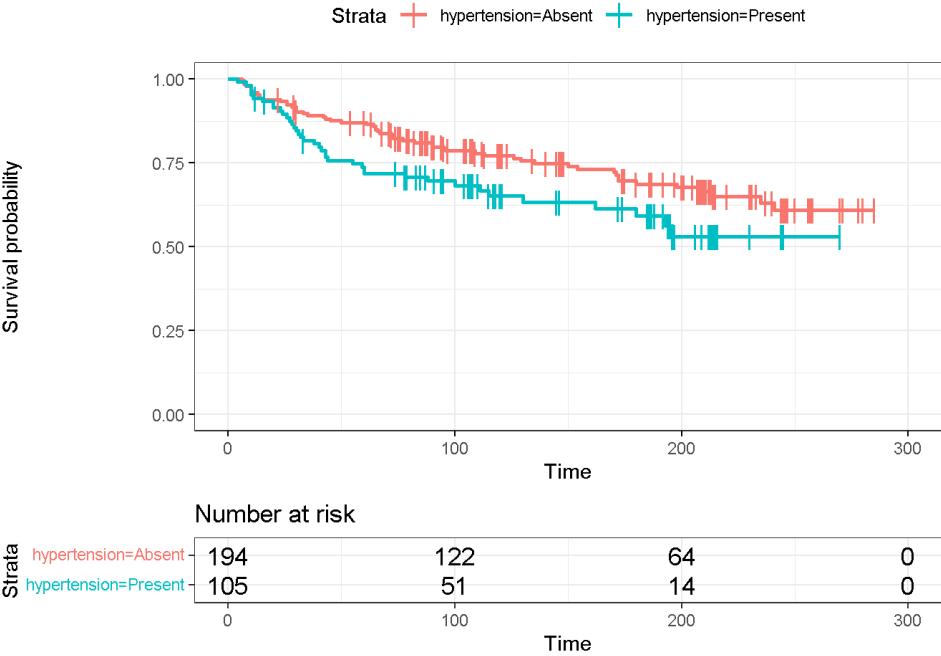
## Hazard Ratios for Significant Variables



```
ggcoxdiagnostics(sigMod, type = "martingale",
                 linear.predictions = FALSE,
                 ggtheme = theme_bw(),
                 )
```



Performing the Log-Rank Test on the `hypertension` & `diabetes`.

- Finding out that the distribution of present hypertension is statistically significant when compared against the distribution of the absence of it.

- The presence of diabetes, however, does not impact survival rate.
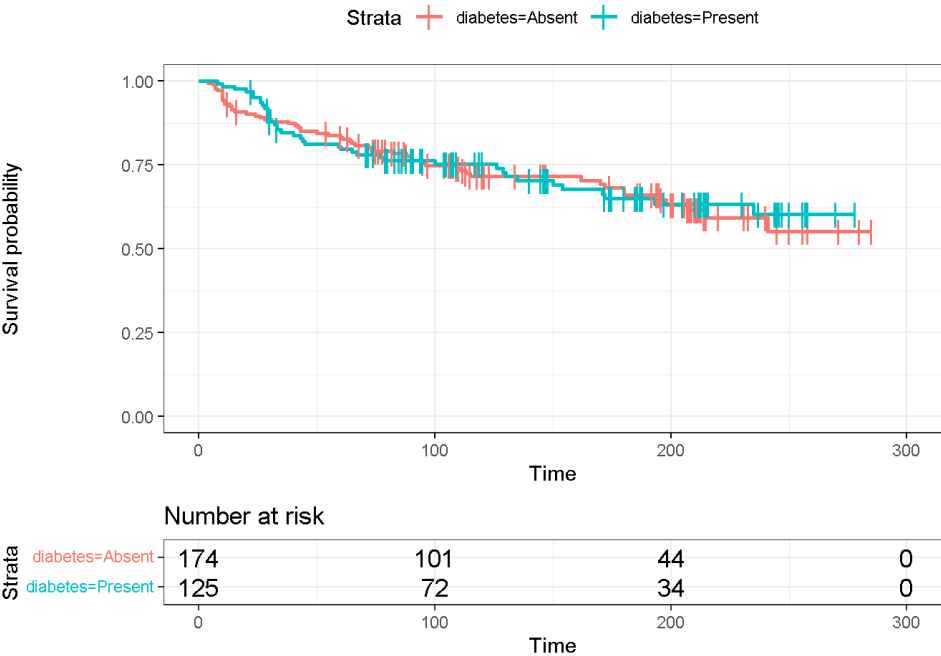
```
#Hypertension
ggsurvplot(survfit(Surv(time,DEATH_EVENT) ~ hypertension, data=HF),
           data = HF,
           censor.shape="|",
           conf.int = FALSE, #surv.median.line = "hv",
           risk.table = TRUE,
           ggtheme = theme_bw())
```

Strata  + hypertension=Absent  + hypertension=Present



### Number at risk

| Strata | | | | |
|---|---|---|---|---|
| hypertension=Absent | 194 | 122 | 64 | 0 |
| hypertension=Present | 105 | 51 | 14 | 0 |

```
survdiff(Surv(time,DEATH_EVENT) ~ hypertension, data=HF)
```

```
## Call:
## survdiff(formula = Surv(time, DEATH_EVENT) ~ hypertension, data = HF)
##
##                          N Observed Expected (O-E)^2/E (O-E)^2/V
## hypertension=Absent    194       57     66.4      1.34      4.41
## hypertension=Present   105       39     29.6      3.00      4.41
##
##  Chisq= 4.4  on 1 degrees of freedom, p= 0.04
```

```
#Diabetes
ggsurvplot(survfit(Surv(time,DEATH_EVENT) ~ diabetes, data=HF),
          data = HF,
          censor.shape="|",
          conf.int = FALSE, #surv.median.line = "hv",
          risk.table = TRUE,
          ggtheme = theme_bw())
```

Strata  + diabetes=Absent  + diabetes=Present



### Number at risk

| Strata | | | | |
|---|---|---|---|---|
| diabetes=Absent | 174 | 101 | 44 | 0 |
| diabetes=Present | 125 | 72 | 34 | 0 |

```
survdiff(Surv(time,DEATH_EVENT) ~ diabetes, data=HF)
```

```
## Call:
## survdiff(formula = Surv(time, DEATH_EVENT) ~ diabetes, data = HF)
##
##                     N Observed Expected (O-E)^2/E (O-E)^2/V
## diabetes=Absent  174       56       55    0.0172    0.0405
## diabetes=Present 125       40       41    0.0231    0.0405
##
##  Chisq= 0  on 1 degrees of freedom, p= 0.8
```

# QUESTIONS TO ASK:

- Probability that person survives longer than time = x? CHECK
- Typical survival rate for person in a certain region? N/A
- Out of 10 random observations, how many do we expect to be deceased after 60 days?

# NEXT STEPS: BUILD A LOGISTIC MODEL FOR PREDICITONS RIGHT AFTER, PERHAPS.

# Findings:

- `Diabetes` isn't a statistically significant predictor of survival time.
- At a given instance in time, someone who has `hypertension` is 0.42 times as likely to die as someone without hypertension adjusting for age.
- At any given instance in time, someone who does *not* have hypertension is 0.65 times as likely to die as someone who does, adjusting for age.
- Probability of survival after 150 days for those younger than 70 is 77%.
- Probability of survival after 200 days for those younger than 70 is 70%.
- 24% probability of survival after t=130 days for patients older than 79, that have less than or equal to 1.8 in serum creatine, and an ejection fraction over 25.
- For those diabetic, plateletes reduce as age increases.
- On average, `creatinine_phosphokinase` is higher for non-smokers.
- Men, on average, have higher `creatinine_phosphokinase`.
- Women, on average, have a higher `platelets` count.