

ConditionalInferenceTrees & Kaplan-Maeier Predict Heart Failure Survival Time

Antonio Pano

11/4/2022

<https://archive.ics.uci.edu/ml/datasets/Heart+failure+clinical+records>
(<https://archive.ics.uci.edu/ml/datasets/Heart+failure+clinical+records>)

- 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)
```

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

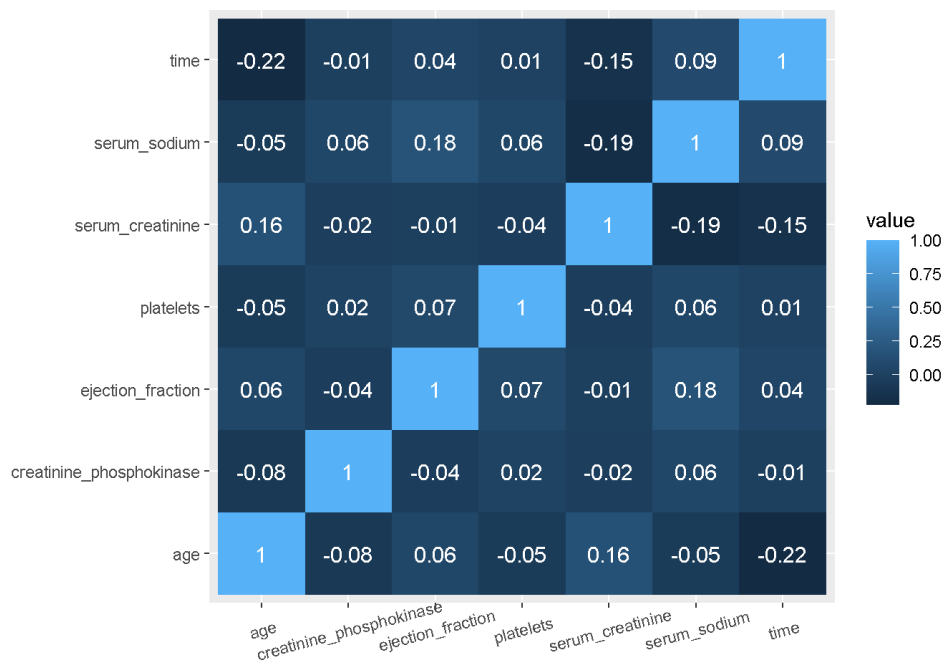
skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
age	0	1	60.83	11.89	40.0	51.0	60.0	70.0	95.0	
creatinine_phosphokinase	0	1	581.84	970.29	23.0	116.5	250.0	582.0	7861.0	
ejection_fraction	0	1	38.08	11.83	14.0	30.0	38.0	45.0	80.0	
platelets	0	1	263358.03	97804.24	25100.0	212500.0	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)

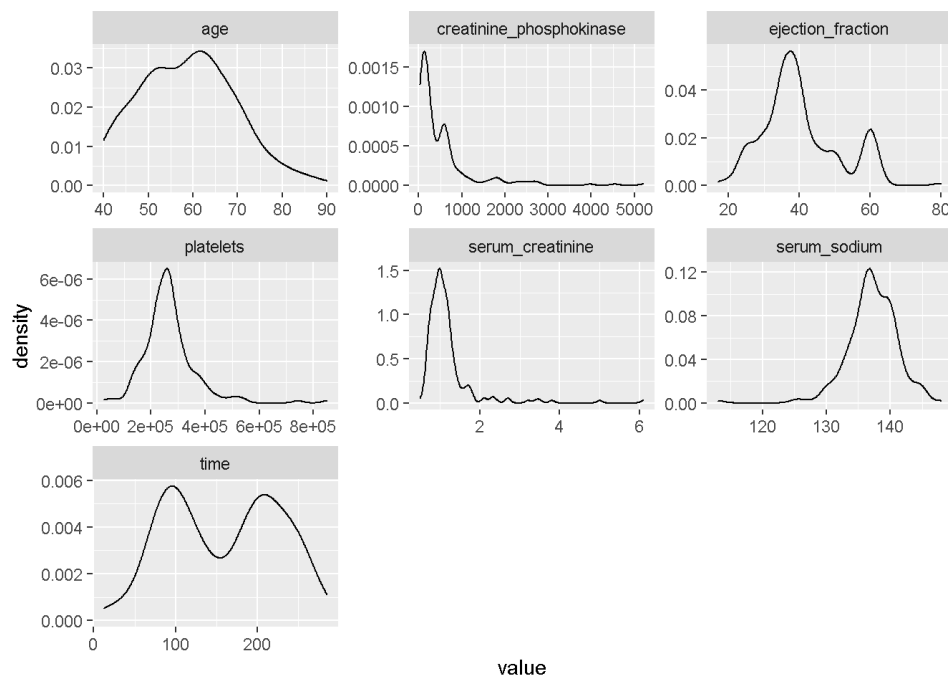
ggplot(data = melted_cormat, aes(x=Var1, y=Var2, fill=value)) +
  geom_tile() +
  geom_text(aes(Var2, Var1, label = value), color = "white", size = 4) +
  theme(axis.title.x=element_blank(),
        axis.title.y=element_blank(),
        axis.text.x = element_text(angle = 15, vjust = 0.8)
  )
```



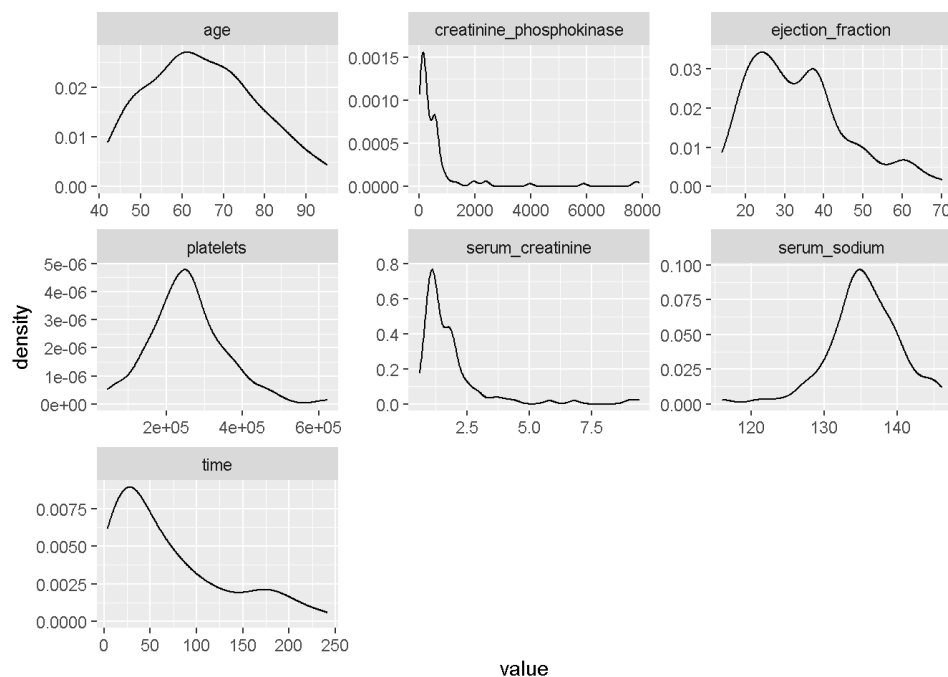
Choosing to grab distributions based on having hypertension– what's traditionally seen as a good indicator of heart failure.

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()
```



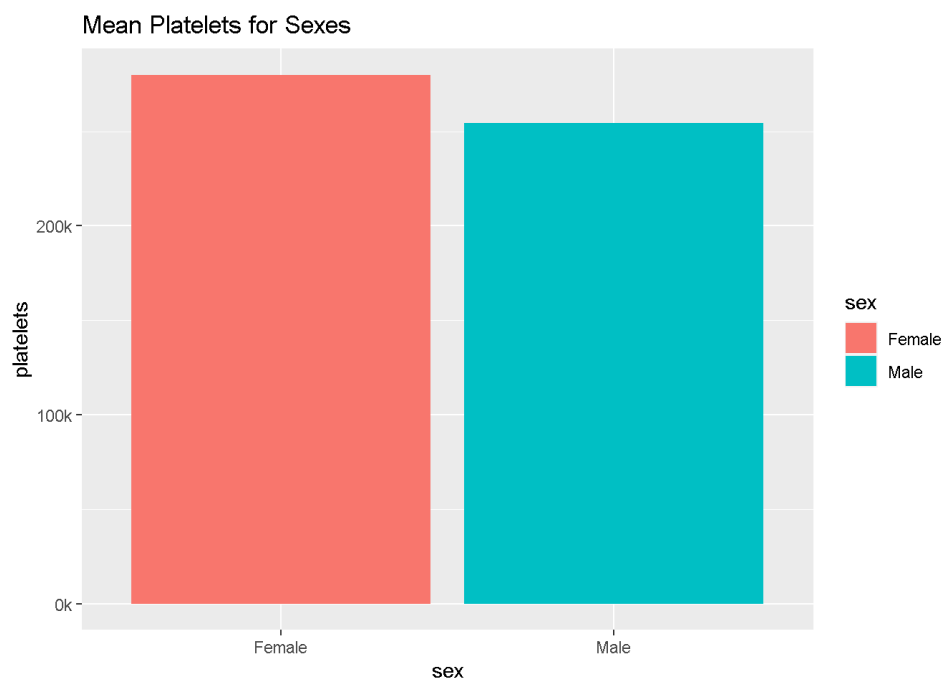
```
HF %>% filter(DEATH_EVENT==1) %>%
  purrr::keep(is.numeric) %>%
  gather() %>%
  ggplot(aes(value)) +
    facet_wrap(~ key, scales = "free") +
    geom_density()
```



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.

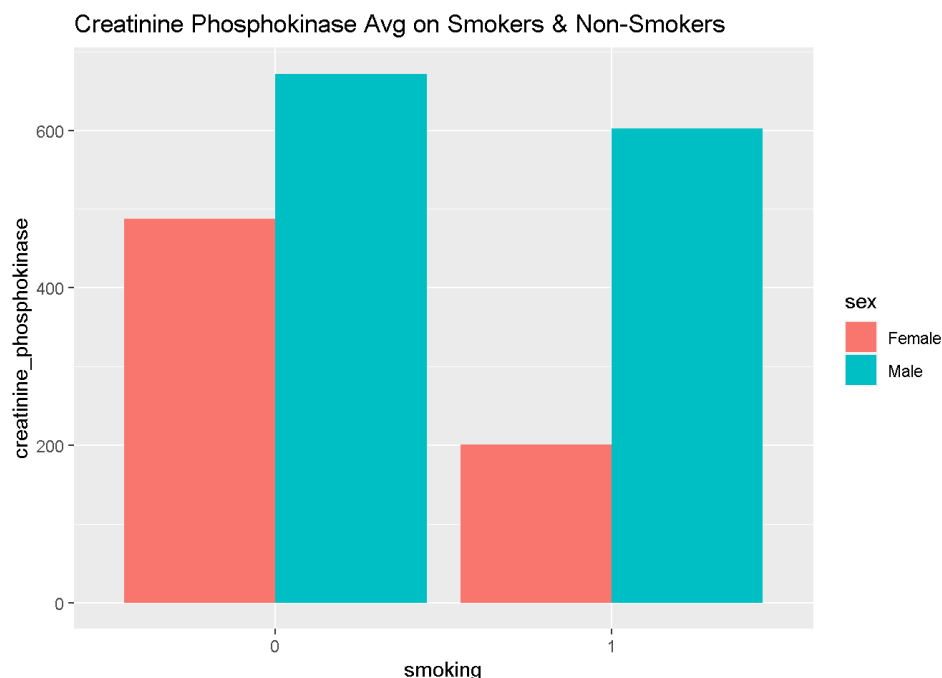
```
ggplot(HF, aes(x=sex, y=platelets, fill=sex)) +
  geom_bar(position = "dodge", stat="summary", fun="mean") +
  scale_y_continuous(labels = scales::label_number(suffix = "k", scale = 1e-3)) +
  ggtitle("Mean Platelets for Sexes")
```



```
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
```

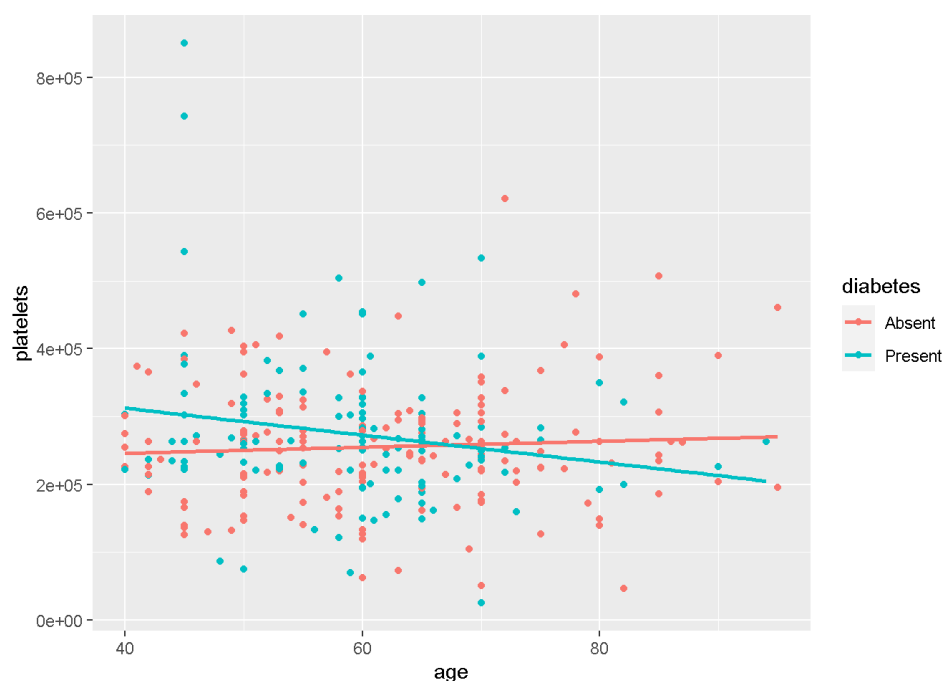
```
ggplot(HF, aes(x=smoking, y=creatinine_phosphokinase, fill=sex)) +
  geom_bar(position = "dodge", stat="summary", fun="mean") +
  ggtitle("Creatinine Phosphokinase Avg on Smokers & Non-Smokers")
```



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

Platelets are incredibly important. Having too few platelets 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)
```



Kaplan - Maeier Curve

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)

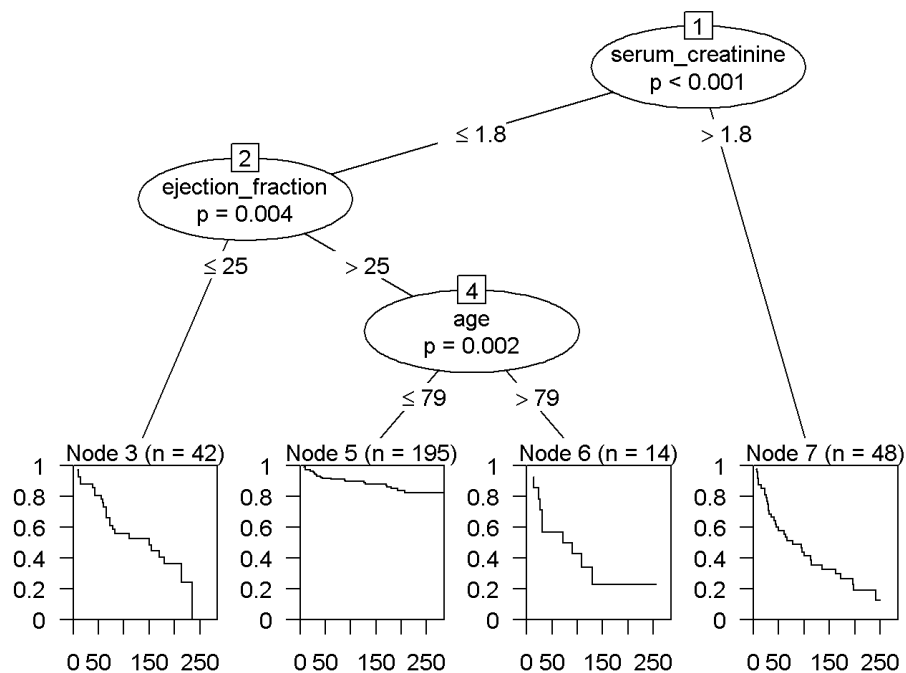
# Won't directly go from factor to numeric. Needed for Survival Analysis.
HF$DEATH_EVENT = as.numeric(as.character(HF$DEATH_EVENT))

# 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.

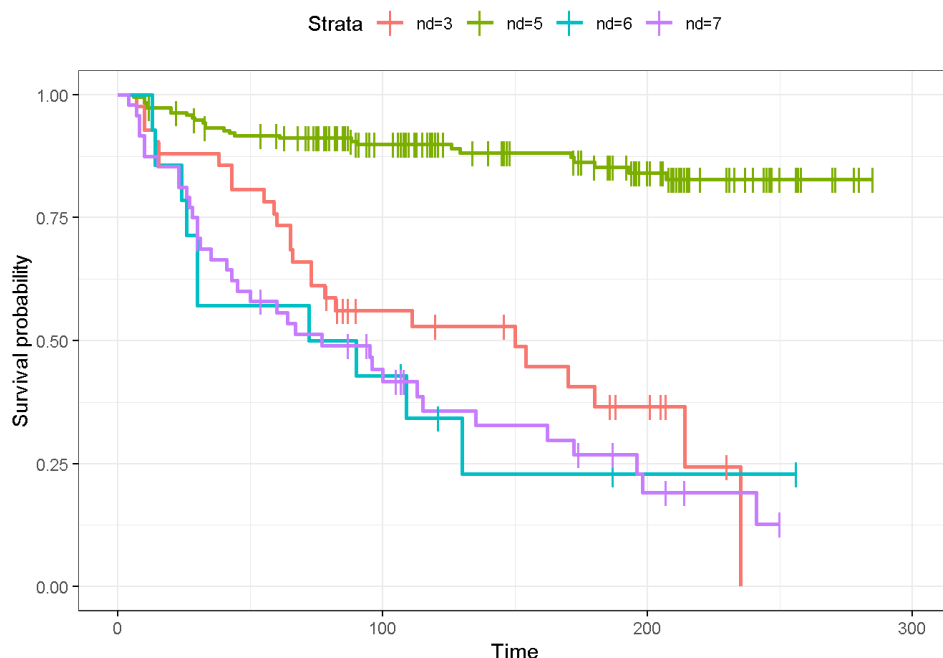
```

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())

```



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 of letting R know that we aren't using any x variables. Just time and whether event occurred which are both y variables.
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 aged 70 or under is 77%.

```
age_subset <- HF %>% filter(age <= 70)
pred_k_age <- survfit(Surv(time, DEATH_EVENT) ~ 1, data = age_subset)

summary(pred_k_age, times = c(150, 200, 250))
```

```
## Call: survfit(formula = Surv(time, DEATH_EVENT) ~ 1, data = age_subset)
##
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##   150   103     53   0.766  0.0288   0.712   0.825
##   200    72      8   0.702  0.0343   0.638   0.773
##   250    17      4   0.639  0.0442   0.558   0.732
```

I'M DONE WITH EVERYTHING THAT CAME BEFORE THIS LINE!

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)
```

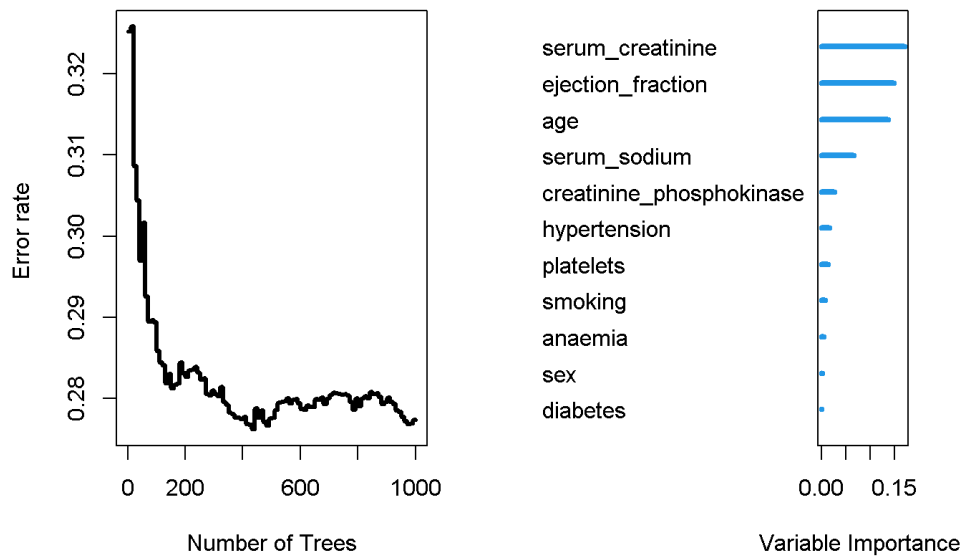
```
fit
```

```
##              Sample size: 299
##              Number of deaths: 96
##              Number of trees: 1000
##              Forest terminal node size: 15
##              Average no. of terminal nodes: 12.966
## No. of variables tried at each split: 2
##              Total no. of variables: 11
##              Resampling used to grow trees: swor
##              Resample size used to grow trees: 189
##              Analysis: RSF
##              Family: surv
##              Splitting rule: logrank *random*
##              Number of random split points: 3
##              (OOB) CRPS: 0.15266934
##              (OOB) Requested performance error: 0.27736247
```

```
summary(fit)
```


##	Length	Class	Mode
## call	7	-none-	call
## family	1	-none-	character
## n	1	-none-	numeric
## ntree	1	-none-	numeric
## nimpute	1	-none-	numeric
## mtry	1	-none-	numeric
## nodesize	1	-none-	numeric
## nodedepth	1	-none-	numeric
## nsplit	1	-none-	numeric
## yvar	2	data.frame	list
## yvar.names	2	-none-	character
## xvar	11	data.frame	list
## xvar.names	11	-none-	character
## event.info	7	-none-	list
## subj	0	-none-	NULL
## subj.names	0	-none-	NULL
## xvar.wt	11	-none-	numeric
## split.wt	11	-none-	numeric
## cause.wt	1	-none-	numeric
## leaf.count	1000	-none-	numeric
## proximity	0	-none-	NULL
## forest	49	rfsrc	list
## forest.wt	0	-none-	NULL
## distance	0	-none-	NULL
## membership	0	-none-	NULL
## tdc.membership	0	-none-	NULL
## splitrule	1	-none-	character
## inbag	0	-none-	NULL
## var.used	0	-none-	NULL
## imputed.indv	0	-none-	NULL
## imputed.data	0	-none-	NULL
## split.depth	0	-none-	NULL
## node.stats	0	-none-	NULL
## ensemble	1	-none-	character
## holdout.array	0	-none-	NULL
## block.size	1	-none-	numeric
## holdout.blk	0	-none-	NULL
## empr.risk	0	-none-	NULL
## oob.empr.risk	0	-none-	NULL
## ctime.internal	1	-none-	numeric
## ctime.external	5	proc_time	numeric
## chf	20631	-none-	numeric
## chf.oob	20631	-none-	numeric
## predicted	299	-none-	numeric
## predicted.oob	299	-none-	numeric
## hazard	0	-none-	NULL
## hazard.oob	0	-none-	NULL
## survival	20631	-none-	numeric
## survival.oob	20631	-none-	numeric
## cif	0	-none-	NULL
## cif.oob	0	-none-	NULL
## err.rate	1000	-none-	numeric
## err.block.rate	100	-none-	numeric
## importance	11	-none-	numeric
## time.interest	69	-none-	numeric
## ndead	1	-none-	numeric

```
plot(fit)
```



```
##
##
```

	Importance	Relative Imp
## serum_creatinine	0.1709	1.0000
## ejection_fraction	0.1492	0.8726
## age	0.1364	0.7981
## serum_sodium	0.0660	0.3861
## creatinine_phosphokinase	0.0262	0.1532
## hypertension	0.0162	0.0950
## platelets	0.0143	0.0839
## smoking	0.0080	0.0468
## anaemia	0.0056	0.0327
## sex	0.0019	0.0111
## diabetes	0.0005	0.0028

```
#calculating variable importance
vimp(fit, importance = "permute")$importance
```

```
##
```

##	age	anaemia	creatinine_phosphokinase
##	0.0374981443	0.0009344516	-0.0013431569
##	diabetes	ejection_fraction	platelets
##	-0.0006854058	0.0420353821	-0.0012737068
##	serum_creatinine	serum_sodium	sex
##	0.0557343406	0.0058497476	0.0006679992
##	smoking	hypertension	
##	0.0025739190	0.0062435021	

```
vimp(fit, importance = "random")$importance
```

```
##
```

##	age	anaemia	creatinine_phosphokinase
##	3.971156e-02	7.544084e-04	-4.813566e-04
##	diabetes	ejection_fraction	platelets
##	-1.064589e-03	4.80002e-02	-2.733320e-03
##	serum_creatinine	serum_sodium	sex
##	5.481908e-02	4.610997e-03	4.388239e-05
##	smoking	hypertension	
##	1.790876e-03	5.035410e-03	

```
vimp(fit, importance = "permute.joint")$importance
```

```
## joint
## 0.1705511
```

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.
coxMod2 <- coxph(Surv(time, DEATH_EVENT) ~ diabetes, data=HF)
summary(coxMod2)
```

```
## 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
```

```
# sex also not stat significant.
coxMod3 <- coxph(Surv(time, DEATH_EVENT) ~ age, data=HF)
summary(coxMod3)
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ age, data = HF)
##
##      n= 299, number of events= 96
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## age 0.042211  1.043115 0.008568 4.927 8.36e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## age      1.043      0.9587      1.026      1.061
##
## Concordance= 0.628 (se = 0.031 )
## Likelihood ratio test= 23.52 on 1 df,  p=1e-06
## Wald test              = 24.27 on 1 df,  p=8e-07
## Score (logrank) test = 24.7 on 1 df,  p=7e-07
```

```
coxMod1 <- coxph(Surv(time, DEATH_EVENT) ~ platelets + age, data=HF)
coxMod1
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ platelets + age, data = HF)
##
##              coef exp(coef) se(coef)      z      p
## platelets -7.506e-07  1.000e+00  1.078e-06 -0.696   0.486
## age      4.253e-02  1.043e+00  8.665e-03  4.909 9.16e-07
##
## Likelihood ratio test=24.01 on 2 df, p=6.11e-06
## n= 299, number of events= 96
```

```
summary(coxMod1)
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ platelets + age, data = HF)
##
##      n= 299, number of events= 96
##
##              coef exp(coef)  se(coef)      z Pr(>|z|)
## platelets -7.506e-07  1.000e+00  1.078e-06 -0.696    0.486
## age       4.253e-02  1.043e+00  8.665e-03  4.909 9.16e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## platelets    1.000      1.0000      1.000    1.000
## age         1.043      0.9584      1.026    1.061
##
## Concordance= 0.628 (se = 0.031 )
## Likelihood ratio test= 24.01 on 2 df,  p=6e-06
## Wald test            = 24.28 on 2 df,  p=5e-06
## Score (logrank) test = 24.87 on 2 df,  p=4e-06
```

```
# Do the Likelihood-Ratio Test
# Try to find combination that may be insightful
anova(coxMod3, coxMod1, test = "LRT")
```

```
## Analysis of Deviance Table
## Cox model: response is Surv(time, DEATH_EVENT)
## Model 1: ~ age
## Model 2: ~ platelets + age
##      loglik  Chisq Df P(>|Chi|)
## 1 -497.45
## 2 -497.20 0.4961 1    0.4812
```

```
# keep only variables that are significant. do manual stepwise, basically, and see what u get.
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+serum_sodium+hypertension, data=HF))
```

```
## Call:
## coxph(formula = Surv(time, DEATH_EVENT) ~ age + anaemia + creatinine_phosphokinase +
##   ejection_fraction + serum_creatinine + serum_sodium + hypertension,
##   data = HF)
##
## n= 299, number of events= 96
##
##               coef exp(coef) se(coef)      z Pr(>|z|)
## age            4.357e-02 1.045e+00 8.831e-03 4.934 8.05e-07 ***
## anaemia1       4.460e-01 1.562e+00 2.150e-01 2.074 0.0380 *
## creatinine_phosphokinase 2.101e-04 1.000e+00 9.825e-05 2.138 0.0325 *
## ejection_fraction -4.747e-02 9.536e-01 1.027e-02 -4.621 3.82e-06 ***
## serum_creatinine 3.139e-01 1.369e+00 6.895e-02 4.552 5.31e-06 ***
## serum_sodium    -4.569e-02 9.553e-01 2.336e-02 -1.956 0.0505 .
## hypertensionPresent 4.965e-01 1.643e+00 2.137e-01 2.324 0.0201 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## age            1.0445      0.9574      1.0266      1.063
## anaemia1       1.5621      0.6402      1.0249      2.381
## creatinine_phosphokinase 1.0002      0.9998      1.0000      1.000
## ejection_fraction 0.9536      1.0486      0.9346      0.973
## serum_creatinine 1.3688      0.7306      1.1957      1.567
## serum_sodium    0.9553      1.0468      0.9126      1.000
## hypertensionPresent 1.6430      0.6086      1.0808      2.498
##
## Concordance= 0.738 (se = 0.027 )
## Likelihood ratio test= 80.58 on 7 df,  p=1e-14
## Wald test            = 88.43 on 7 df,  p=3e-16
## Score (logrank) test = 87.66 on 7 df,  p=4e-16
```

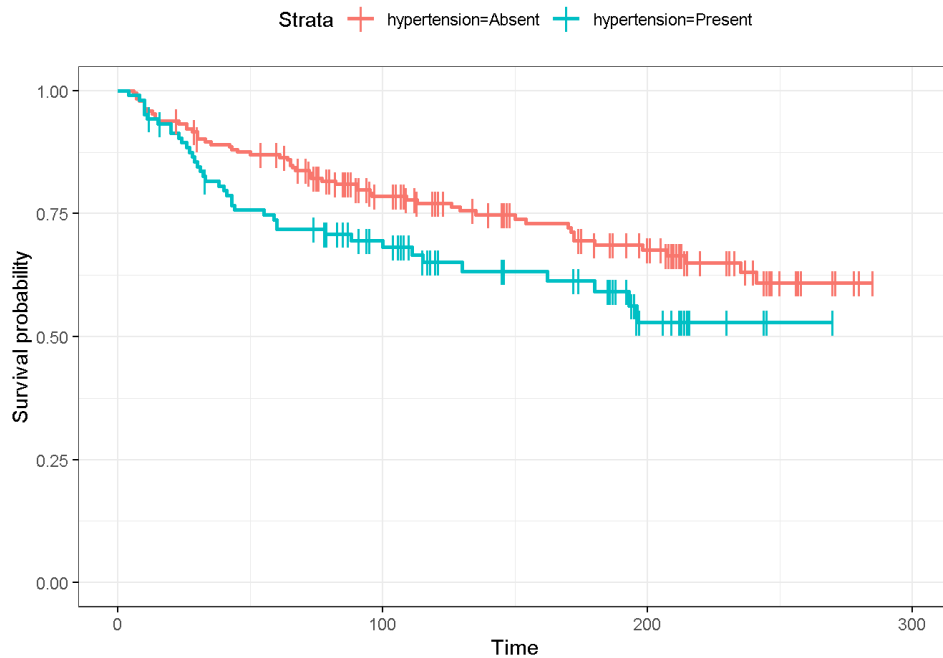
```
summary(coxph(Surv(time, DEATH_EVENT) ~ age+anaemia+creatinine_phosphokinase+ejection_fraction+
  serum_creatinine+hypertension, data=HF))
```

```
## 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
```

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.

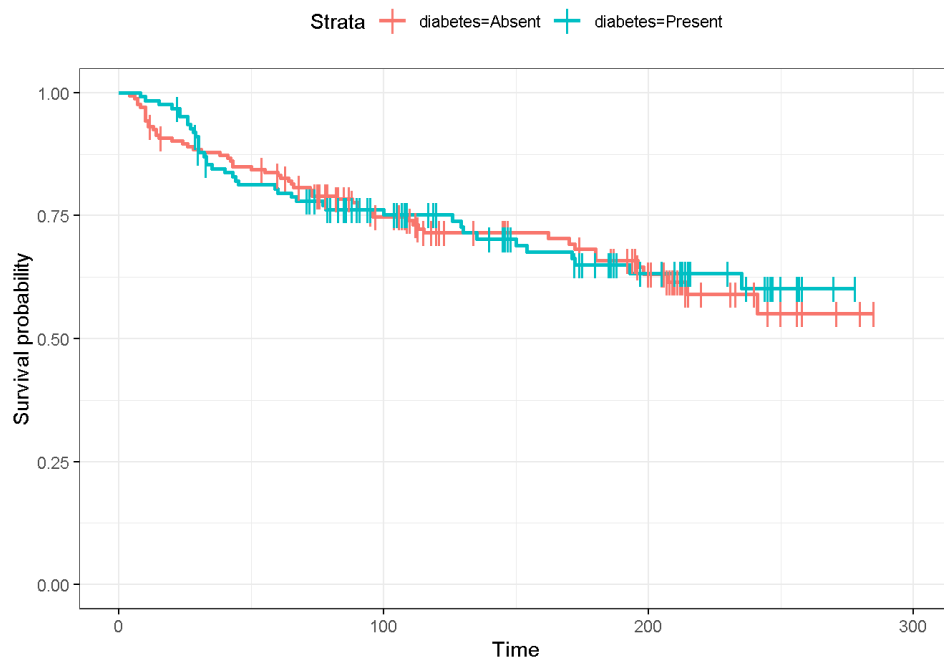
```
#Hypertension
ggsurvplot(survfit(Surv(time,DEATH_EVENT) ~ hypertension, data=HF),
  data = HF,
  censor.shape="|",
  conf.int = FALSE, #surv.median.line = "hv",
  ggtheme = theme_bw())
```



```
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",
  ggtheme = theme_bw())
```



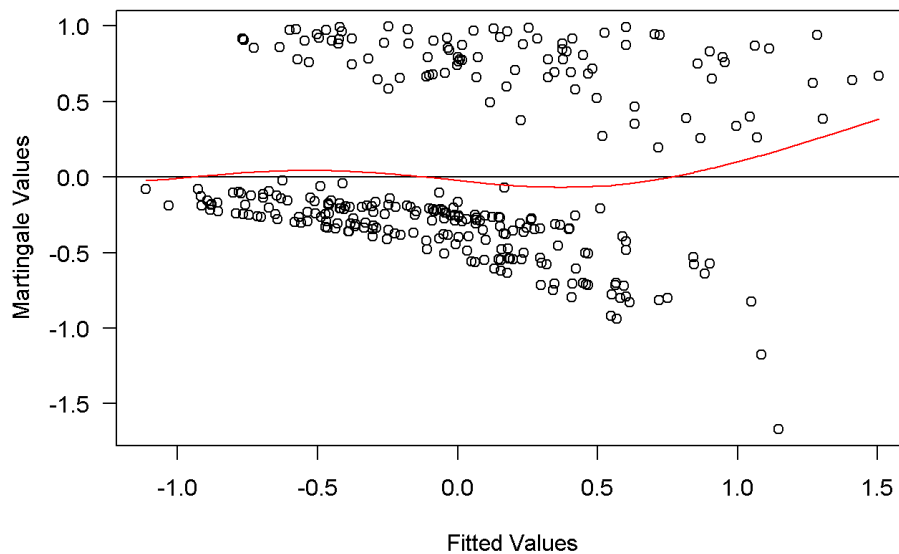
```
survdiff(Surv(time,DEATH_EVENT) ~ diabetes, data=HF)
```

```
## 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
```

Checking Linearity of Model DO THIS AT THE VERY END!!!

```
# Checking Linearity using MARTINGALE residuals.
# Should be as linear as possible.
plot(predict(coxMod1), residuals(coxMod1, type = "martingale"), xlab = "Fitted Values",
      ylab = "Martingale Values", main = "Residual Plot", las = 1) +
  abline(h=0) +
  lines(smooth.spline(predict(coxMod1), residuals(coxMod1, type="martingale")), col="red")
```

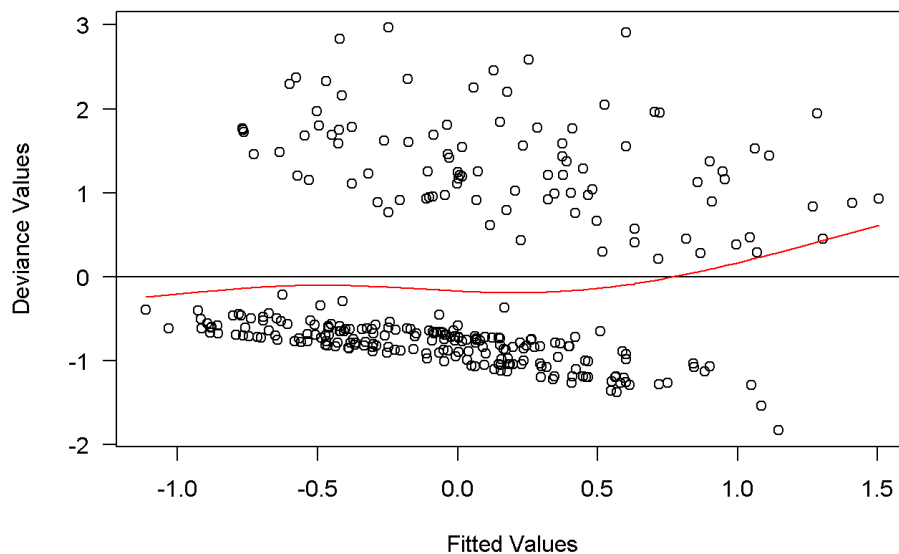

Residual Plot



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

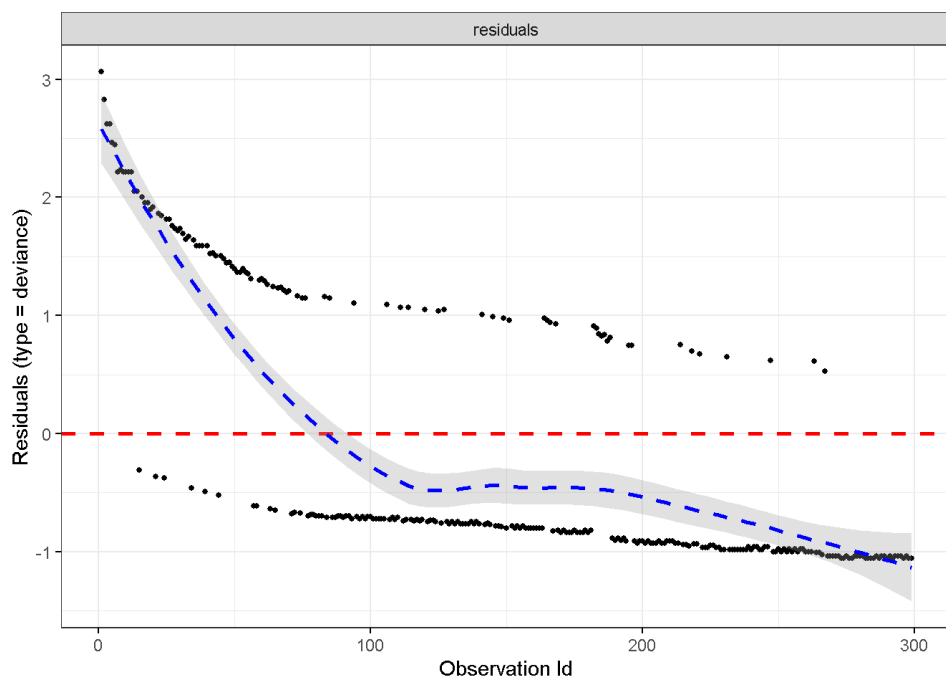
```
# Same thing using deviance residuals
plot(predict(coxMod1), residuals(coxMod1, type = "deviance"), xlab = "Fitted Values",
     ylab = "Deviance Values", main = "Residual Plot", las = 1) +
  abline(h=0) +
  lines(smooth.spline(predict(coxMod1), residuals(coxMod1, type="deviance")), col="red")
```

Residual Plot



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

```
#NO IDEA WHAT THIS DOES BUT IT'S USEFUL, I THINK
ggcoxdiagnostics(coxMod2, type = "deviance",
  linear.predictions = FALSE, ggtheme = theme_bw())
```



NEXT STEPS: CHECK PRESENCE OF LINEARITY FOR COX REGRESSIONS. BUILD A LOGISTIC MODEL FOR PREDICITONS RIGHT AFTER, PERHAPS.