

# Survival Analysis Project: HIV Clinical Trial

*Juste Simanauskaite & Patricia Rivera*

## Contents

<b>Introduction</b>	<b>1</b>
<b>Data-Set Analysis</b>	<b>3</b>
<b>Results</b>	<b>6</b>
Survival Analysis . . . . .	6
Kaplan-Meier Curves . . . . .	6
<b>Patricia's "Something New"</b>	<b>14</b>
1. What is the topic? . . . . .	14
2. How it is relevant? How it relates to survival analysis/analysis at hand? . . . . .	14
3. Resources to learn about the topic. . . . .	14
4. What will be challenging about learning something new? . . . . .	15
Power Analysis code and simulation . . . . .	15
<b>Juste's "Something New": The Schoenfeld Residuals for the Cox PH model</b>	<b>18</b>
3. Resources to learn about the topic. . . . .	18
Explanation of the Theory Behind Schoenfeld Residuals . . . . .	18

## Introduction

HIV (Human Immunodeficiency Virus) is a disease known as an immune system disorder, which causes severe destruction of white blood cells that are responsible for fighting infection. The presence of this disorder is a lead-in for a human to be more prone to infections and cancer diseases. AIDS is the final stage of HIV, which is not always developed in HIV patients. Zidovudine (AZT) is known as antiretroviral medication for prevention of HIV/AIDS, whereas lamivudine (3TC) is an inhibitor medication that works in decreasing HIV and hepatitis B. Previously, it has been founded that three-drug combinations, in particular, with a previous exposure to AZT, have shown the most significant resulted in reducing HIV-1 RNA concentrations. Therefore, this study used indinavir sulfate (a synthetic antiviral agent that inhibits HIV protease activity) in combination with AZT and 3TC as well as variation of placebo treatments to determine the potency of triple drug therapy in the cases of advanced HIV-1 patients. The study hypothesized that a three-drug combination, including a HIV-protease inhibitor and two nucleoside analogues (AZT and 3TC) would alter the progression of the HIV-1 disease. The study was successful in reaching significant data of the clinical superiority of a three-drug approach with inidavor over a treatment containing only a two-drug combination.

The current analysis of the data from a study conducted by Hammer et al. in 1997 considers the response variable to be *time*, which here describes the amount of time in days for the time of death, AIDS diagnosis, or the termination of the study. Another important variable used for the analysis is *sensor*, which indicates the paarticipants of the study that survived till the termination of the study without dying or being diagnosed AIDS. The study explored the influence of the explanatory varibale *tx*. referring to the treatment group that was differentiated into: a control (placebo group) and a treatment group that included IDV (indinavir) # Methods

The study was a randomized, double-blind, and a placebo-controlled trial that compared a three-drug treatment of indinavir (Crixivan), zidovudine (AZT) and lamivudine (3TC) with a two-drug treatment. Patients were selected based on the factor that they had no more than 200 CD4 cells per cubic millimetear at least 3 months prior to AZT therapy. The patients had to be more than 16 years old, with a diagnostic documentation of HIV-1 infection, having no more than 1 week of prior lamivudine treatment, and a Karnofsky score of at least 70.

The approved patients received 200mg of open-label zidovudine three times daily and 150mg of lamuvidine two times daily and were randomly assigned to a placebo or a treatment of 800mg of indinavir every eight hours.

Some modifications were made to the protocol. In October of 1996 prior exposure to AZT was reduced to at least 3 months and permitted patients with no tolerance for this drug to enter the study with stavudine as a substitute.

Patients diagnosed with AIDS-defining events were offered an open-label assignment of the indinavir treatment with nor reveal of their initial treatment assignments. All of these cases had to be reviewed via a blind procedure by the study chair.

Follow ups were made at weeks 4,8, and 16 and every eight weeks afterwards. CD4 cell counts and Plasma HIV-1 RNA concentrations were measured twice at baseline and at weeks 4,8,24, and 40.

The statistical analysis methods used to interpret results were Kaplan-Meier estimates, log-rank tests, and proportional hazards models. The p-values, estimates of treatment differences and 95% confidence intervals were not adjusted for repeated analysis.

The data and results have been reviewed again in 2019 and have been analyzed via statistic methhods such as Cox Proportional Hazards test, Kaplan-Meier estimates, Aalen model, Power analysis, and Schoenfeld residuals.

```
knitr::opts_chunk$set(message=FALSE, warning=FALSE, fig.height=3, fig.width=5,
                        fig.align="center")

library(tidyverse)
library(broom)
library(plyr)
library(survival)
library(survminer)
library(coxed)

aids <- read.csv( "http://pages.pomona.edu/~jsh04747/courses/math150/AIDSdata.csv")
dim(aids)

## [1] 851 16

summary(aids)
```

```
##      id      time      censor      time_d
## Min.   : 1.0   Min.   : 1.0   Min.   :0.00000   Min.   : 1.0
## 1st Qu.:287.5  1st Qu.:179.5  1st Qu.:0.00000   1st Qu.:199.5
## Median :581.0  Median :257.0  Median :0.00000   Median :266.0
## Mean   :579.5  Mean   :231.8  Mean   :0.08108   Mean   :243.4
## 3rd Qu.:873.0  3rd Qu.:300.0  3rd Qu.:0.00000   3rd Qu.:306.0
## Max.   :1156.0 Max.   :362.0  Max.   :1.00000   Max.   :362.0
##      censor_d      tx      txgrp      strat2
## Min.   :0.0000   Min.   :0.0000   Min.   :1.000   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:1.000   1st Qu.:0.0000
## Median :0.0000   Median :1.0000   Median :2.000   Median :1.0000
## Mean   :0.0235   Mean   :0.5041   Mean   :1.504   Mean   :0.6157
## 3rd Qu.:0.0000   3rd Qu.:1.0000   3rd Qu.:2.000   3rd Qu.:1.0000
## Max.   :1.0000   Max.   :1.0000   Max.   :2.000   Max.   :1.0000
##      sex      raceth      ivdrug      hemophil
## Min.   :1.000   Min.   :1.000   Min.   :1.000   Min.   :0.00000
## 1st Qu.:1.000   1st Qu.:1.000   1st Qu.:1.000   1st Qu.:0.00000
## Median :1.000   Median :1.000   Median :1.000   Median :0.00000
## Mean   :1.157   Mean   :1.706   Mean   :1.317   Mean   :0.03408
## 3rd Qu.:1.000   3rd Qu.:2.000   3rd Qu.:1.000   3rd Qu.:0.00000
## Max.   :2.000   Max.   :5.000   Max.   :3.000   Max.   :1.00000
##      karnof      cd4      priorzdv      age
## Min.   : 70.00   Min.   : 0.00   Min.   : 3.00   Min.   :15.00
## 1st Qu.: 90.00   1st Qu.:22.25   1st Qu.:11.00   1st Qu.:33.00
```

```
## Median : 90.00   Median : 75.00   Median : 21.00   Median :38.00
## Mean    : 91.34   Mean      : 86.45   Mean      : 30.63   Mean      :38.81
## 3rd Qu.:100.00   3rd Qu.:135.75   3rd Qu.: 44.00   3rd Qu.:44.00
## Max.    :100.00   Max.      :348.00   Max.      :288.00   Max.      :73.00
```

## Data-Set Analysis

The data set contains a sample size equal to 851 participants and 16 different variables. Out of these participants 782 were considered as uncensored data point, which indicates that these patients survived through the course of the study without diagnosis of AIDS and/or death. 69 were found to be censored meaning that either there was an occurrence of death or AIDS diagnosis, out of which it is known that 20 patients died throughout the course of the study.

```
#Survival Analysis
#mutation of age
aids <- read.csv( "http://pages.pomona.edu/~jsh04747/courses/math150/AIDSdata.csv")
aids <- aids %>%
  mutate(age = ifelse(age <= 20, "under20",
                      ifelse(age <=30, "20-30",
                              ifelse(age <= 40, "30-40",
                                      ifelse(age <=50, "40-50",
                                              ifelse(age <=60, "50-60",
                                                      ifelse(age <=70, "60-70",
                                                                "over70")))))))) %>%

  mutate(age = factor(age,
                      levels = c("under20", "20-30", "30-40", "40-50", "50-60", "60-70", "over70")),
         sex = ifelse(sex == 2, "male", "female"))

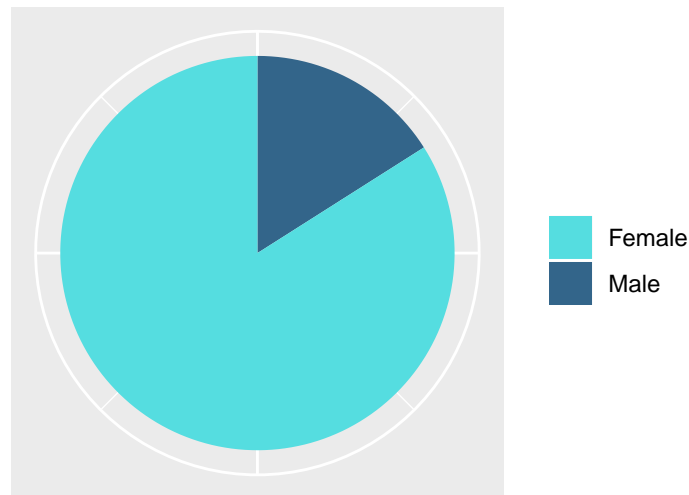
aids <- aids %>%
  mutate(cd4 = ifelse(cd4 <=50, "0-50",
                     ifelse(cd4 <=100, "50-100",
                             ifelse(cd4 <= 150, "100-150",
                                     ifelse(cd4 <=200, "150-200",
                                             ifelse(cd4 <=250, "200-250",
                                                    ifelse(cd4 <=300, "300-350", "350+"))))))))
```

Since there are many values of the explanatory variable *age* in the original data, we've decided to mutate the variable into age categories from under 20 to over 70 in increments of 10 years. Similar modifications have been made to the baseline *CD4* count, just in increments of 50 up until 350+. Furthermore, we changed the labeling and representation of *sex* into "male" and "female" instead of "1" and "2" in the data.

```
library(plotrix)
male<-sum(aids$sex=="male")
female<-sum(aids$sex=="female")
slices <- c(male, female)
lbls <- c("Male", "Female")
pct <- round(slices/sum(slices)*100)
df = data.frame(slices = pct, labels = lbls)
sexplot<- ggplot(df, aes(x = factor(1), y=pct, fill = labels)) +
  geom_bar(stat="identity", width = 1)+
  coord_polar(theta = "y")+
  theme(axis.line = element_blank(),
        axis.text = element_blank(),
        axis.ticks = element_blank(),
        axis.title = element_blank())
print(sexplot + ggtitle("Gender Distribution"))
```

```
scale_fill_manual(values=c("#55DDE0", "#33658A",
                           "#2F4858"))+
labs(x = NULL, y = NULL, fill = NULL))
```

Gender Distribution

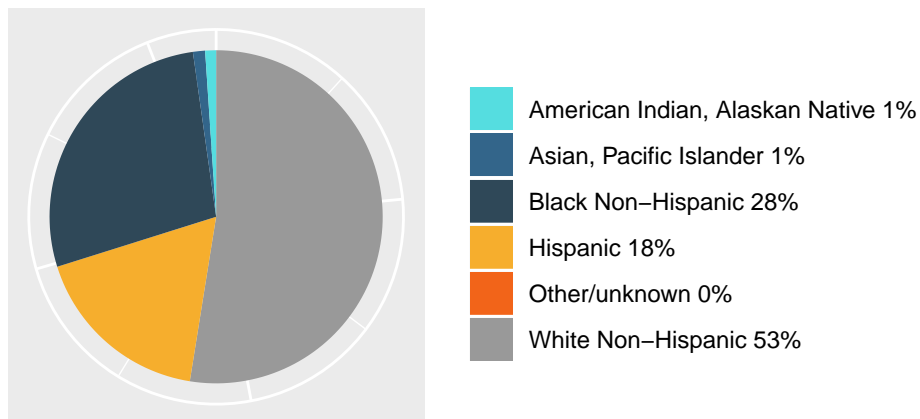


The Pie Chart represents the gender distribution in the sample, with 84% male and 16% female. This shows the potential for the data to not be able to correctly represent the difference of the data variance by gender, if there were to be one. Therefore, gender is something to look into in future data analysis.

```
wnh<-sum(aids$raceth==1)
bnh<-sum(aids$raceth==2)
h<-sum(aids$raceth==3)
api<-sum(aids$raceth==4)
aian<-sum(aids$raceth==5)
oth<-sum(aids$raceth==6)
slices <- c(wnh,bnh,h,api,aian,oth)
lbls <- c("White Non-Hispanic", "Black Non-Hispanic", "Hispanic","Asian, Pacific Islander",
          "American Indian, Alaskan Native", "Other/unknown")
pct <- round(slices/sum(slices)*100)
lbls <- paste(lbls, pct)
lbls <- paste(lbls,"%",sep="")
df = data.frame(slices = slices,labels = lbls)

ethplot<- ggplot(df,aes(x = factor(1),y=slices, fill = labels)) +
  geom_bar(stat="identity", width = 1)+
  coord_polar(theta = "y")+
  theme(axis.line = element_blank(),
        axis.text = element_blank(),
        axis.ticks = element_blank(),
        axis.title = element_blank())
print(ethplot + ggtitle("Race/Ethnicity Distribution among participants")+
  scale_fill_manual(values=c("#55DDE0", "#33658A",
                              "#2F4858", "#F6AE2D", "#F26419",
                              "#999999"))+ labs(x = NULL, y = NULL,
                                                fill = NULL))
```

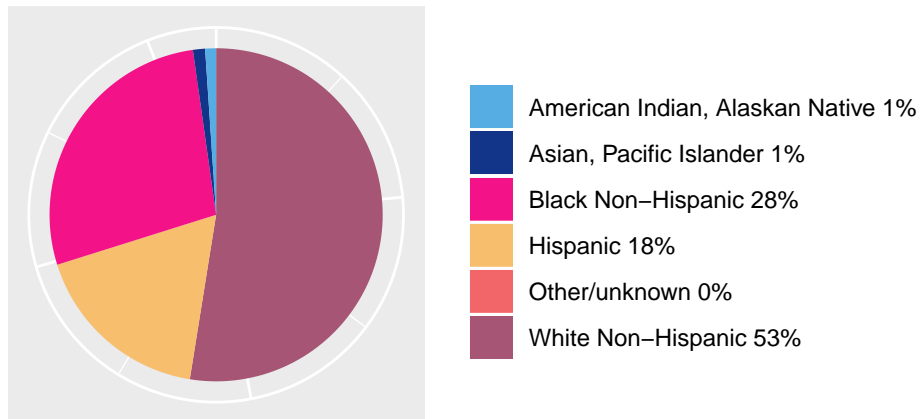
## Race/Ethnicity Distribution among participants



The distribution of race/ethnicity shows that the greatest number of participants consists of white non-Hispanic identifying individuals, with black non-Hispanic following and Hispanic as the 3rd largest represented group.

```
never<-sum(aids$ivdrug==1)
cur<-sum(aids$ivdrug==2)
prev<-sum(aids$ivdrug==3)
slices <- c(never,cur,prev)
lbls <- c("Never", "Currently", "Previously")
pct <- round(slices/sum(slices)*100)
lbls <- paste(lbls, pct)
lbls <- paste(lbls,"%",sep="")
df = data.frame(slices = slices,labels = lbls)
ivplot<- ggplot(df,aes(x = factor(1),y=slices, fill = labels)) +
  geom_bar(stat="identity", width = 1)+
  coord_polar(theta = "y")+
  theme(axis.line = element_blank(),
        axis.text = element_blank(),
        axis.ticks = element_blank(),
        axis.title = element_blank())
print(ethplot + ggtitle("IV Drug Use History")+
  scale_fill_manual(values=c("#56ADE3", "#12358A",
                             "#F41288", "#F7BE6D", "#F26569",
                             "#A65674"))+
  labs(x = NULL, y = NULL, fill = NULL))
```

## IV Drug Use History



From this chart we see that most of the participants (84%) have never used IV drugs, whereas 16% of participants have some type of history of usage and none of the participants reported to be currently using the drugs.

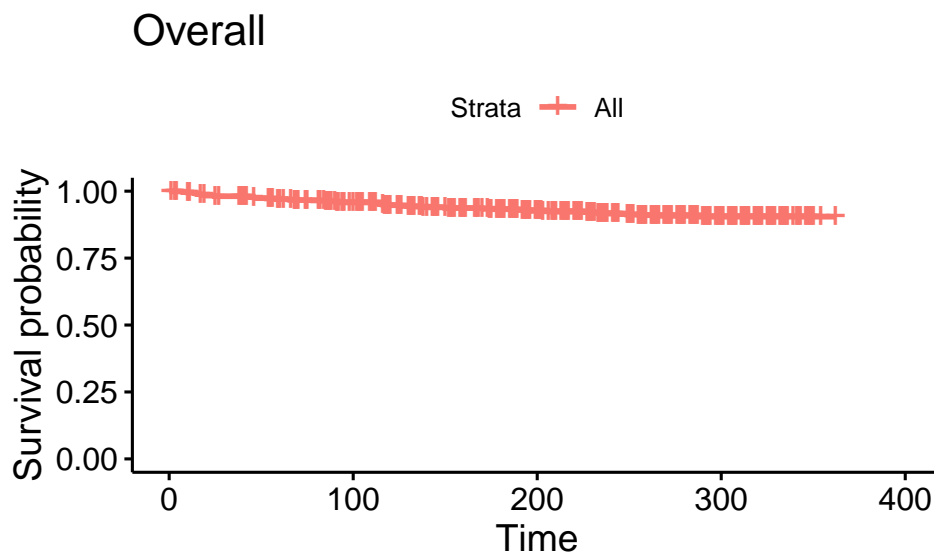
## Results

### Survival Analysis

#### Kaplan-Meier Curves

The following graph is a representation of a Kaplan Meier Curve for all participants in the study, we can see that only a few participants dies or were diagnosed with AIDS during the study as the slope of the curve is not experiencing a high decrease.

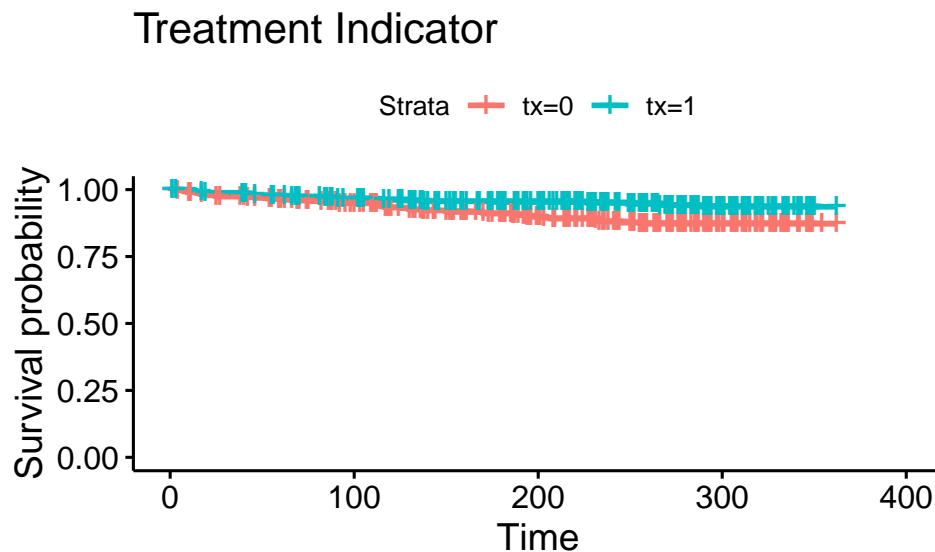
```
fit <- survfit(Surv(time,censor)~1, data = aids)
ggsurvplot(fit,data = aids,conf.int = FALSE) + ggtitle("Overall")
```



The following graph is a representation of the Kaplan-Meier survival probability based on the treatment indicator. In this case  $tx=0$  was the control group and  $tx=1$  the treatment group that was given IDV. Already, we can see a

trend in the graph that the control group shows a lower survival probability with time. According to the log-rank test, we see that the p-value for the test statistic is equal to 0.002 ( $<0.05$ ), thus we can reject the null hypothesis that the two population survival functions are the same, and the alternative is accepted, which says that the survival curves are different. The Wilcoxon test also provides us with a small p-value of 0.002, which again rejects the null and goes in agreement with our primary conclusion.

```
fit1 <- survfit(Surv(time,censor)~tx, data = aids)
ggsurvplot(fit1,data = aids,conf.int = FALSE) + ggtitle("Treatment Indicator")
```



'Log-Rank'

```
## [1] "Log-Rank"
```

```
survdif(Surv(time,censor)~tx, data = aids, rho=0)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ tx, data = aids, rho = 0)
```

```
##
```

```
##      N Observed Expected (O-E)^2/E (O-E)^2/V
```

```
## tx=0 422      46      33.3      4.83      9.35
```

```
## tx=1 429      23      35.7      4.51      9.35
```

```
##
```

```
## Chisq= 9.3 on 1 degrees of freedom, p= 0.002
```

'Wilcoxon'

```
## [1] "Wilcoxon"
```

```
survdif(Surv(time,censor)~tx, data = aids, rho=1)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ tx, data = aids, rho = 1)
```

```
##
```

```
##      N Observed Expected (O-E)^2/E (O-E)^2/V
```

```
## tx=0 422     44.0      31.9      4.57      9.24
```

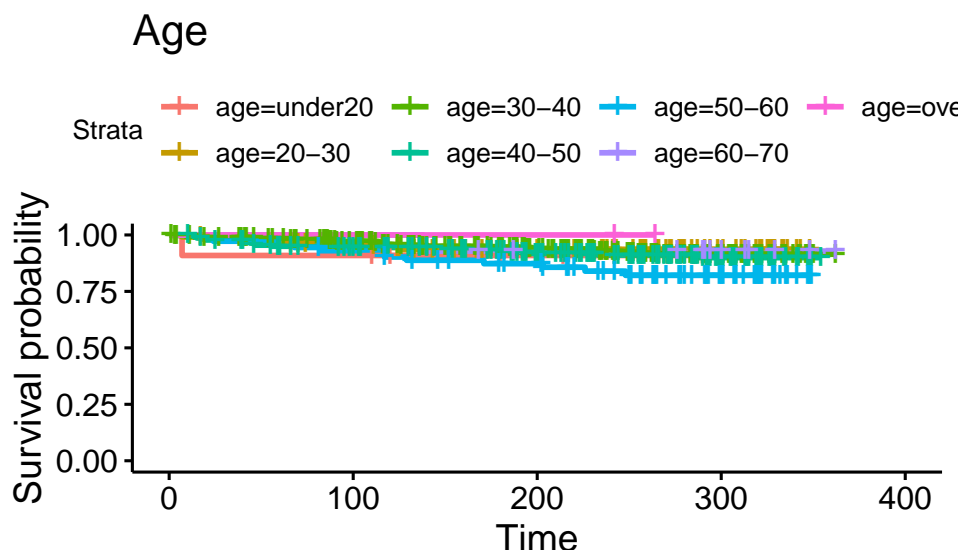
```
## tx=1 429     22.1      34.2      4.28      9.24
```

```
##
```

```
## Chisq= 9.2 on 1 degrees of freedom, p= 0.002
```

```
fit2 <- survfit(Surv(time,censor)~age, data = aids)
```

```
ggsurvplot(fit2,data = aids,conf.int = FALSE) + ggtitle("Age")
```



```
'Log-Rank'
```

```
## [1] "Log-Rank"
```

```
survdif(Surv(time,censor)~age, data = aids, rho=0)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ age, data = aids, rho = 0)
```

```
##
```

	N	Observed	Expected	(O-E) <sup>2</sup> /E	(O-E) <sup>2</sup> /V
age=under20	11	1	0.824	0.0378	0.0383
age=20-30	111	7	9.065	0.4703	0.5418
age=30-40	416	29	33.543	0.6152	1.1980
age=40-50	225	19	17.995	0.0561	0.0759
age=50-60	72	12	6.098	5.7134	6.2712
age=60-70	14	1	1.294	0.0669	0.0682
age=over70	2	0	0.182	0.1817	0.1823

```
##
```

```
## Chisq= 7.1 on 6 degrees of freedom, p= 0.3
```

```
'Wilcoxon'
```

```
## [1] "Wilcoxon"
```

```
survdif(Surv(time,censor)~age, data = aids, rho=1)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ age, data = aids, rho = 1)
```

```
##
```

	N	Observed	Expected	(O-E) <sup>2</sup> /E	(O-E) <sup>2</sup> /V
age=under20	11	0.999	0.790	0.0553	0.0583
age=20-30	111	6.812	8.684	0.4037	0.4850
age=30-40	416	27.640	32.137	0.6292	1.2782
age=40-50	225	18.280	17.238	0.0629	0.0889
age=50-60	72	11.423	5.836	5.3495	6.1293
age=60-70	14	0.941	1.236	0.0706	0.0753
age=over70	2	0.000	0.174	0.1738	0.1821

```
##
```

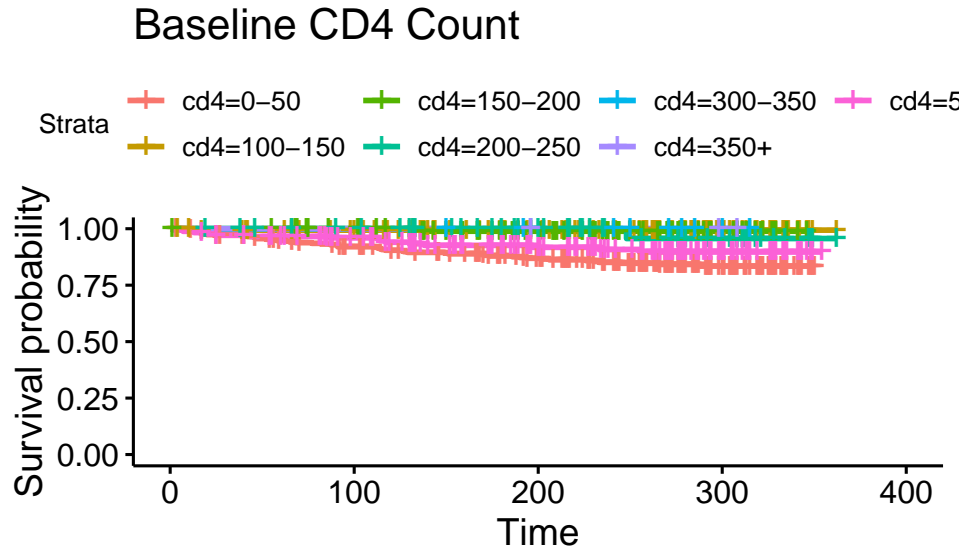
```
## Chisq= 7 on 6 degrees of freedom, p= 0.3
```

The following graph is a representation of the Kaplan-Meier survival probability based on the Baseline CD4. Due to



the close proximity of the curves it's harder to see the significance of the differences. However, according to the log-rank test, we see that the p-value for the test statistic is equal to  $5e-07$  ( $<0.05$ ), thus we can reject the null hypothesis. The Wilcoxon test again provides us with a small p-value of  $5e-07$ , which again rejects the null and goes in agreement with our primary conclusion that the curves are significantly different.

```
fit3 <- survfit(Surv(time,censor)~cd4, data = aids)
ggsurvplot(fit3,data = aids,conf.int = FALSE) + ggtitle("Baseline CD4 Count")
```



'Log-Rank'

```
## [1] "Log-Rank"
```

```
survdif(Surv(time,censor)~cd4, data = aids, rho=0)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ cd4, data = aids, rho = 0)
```

```
##
```

	N	Observed	Expected	$(O-E)^2/E$	$(O-E)^2/V$
cd4=0-50	346	51	28.279	18.256	30.974
cd4=100-150	162	1	13.444	11.519	14.315
cd4=150-200	104	1	8.467	6.585	7.511
cd4=200-250	51	1	4.047	2.294	2.439
cd4=300-350	10	0	0.882	0.882	0.894
cd4=350+	3	0	0.275	0.275	0.276
cd4=50-100	175	15	13.606	0.143	0.178

```
##
```

```
## Chisq= 40 on 6 degrees of freedom, p= 5e-07
```

'Wilcoxon'

```
## [1] "Wilcoxon"
```

```
survdif(Surv(time,censor)~cd4, data = aids, rho=1)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ cd4, data = aids, rho = 1)
```

```
##
```

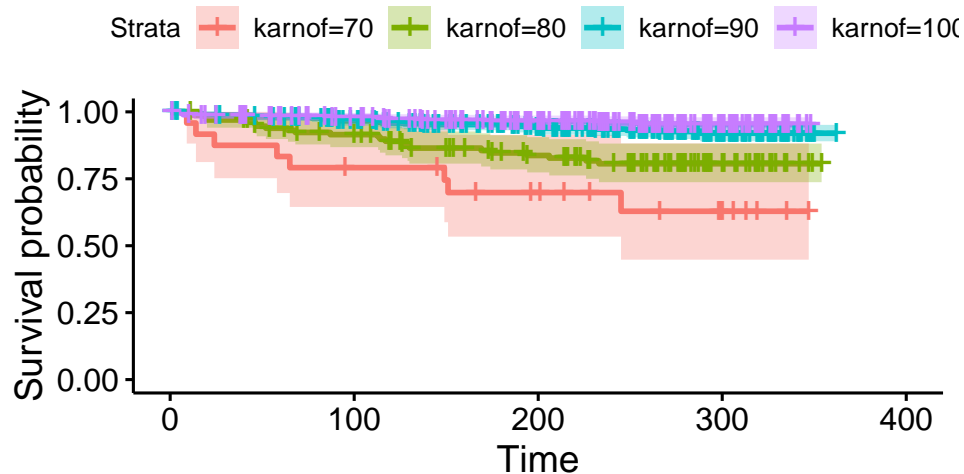
	N	Observed	Expected	$(O-E)^2/E$	$(O-E)^2/V$
cd4=0-50	346	48.834	27.072	17.494	30.954
cd4=100-150	162	0.995	12.872	10.959	14.210
cd4=150-200	104	0.942	8.114	6.339	7.540

```
## cd4=200-250 51 0.914 3.881 2.268 2.513
## cd4=300-350 10 0.000 0.844 0.844 0.893
## cd4=350+ 3 0.000 0.263 0.263 0.276
## cd4=50-100 175 14.409 13.049 0.142 0.184
##
## Chisq= 40 on 6 degrees of freedom, p= 5e-07
```

The final explanatory variable we're investigating as a part of our model is the Karnofsky Performance Scale. The Kaplan-Meier curves for this variable present a higher amplitude of distributions across the survival scale. The p-values of both log-rank and the Wilcoxon Test again present with the same significant p-value of  $5e-10$  ( $<0.05$ ), thus we can reject the null hypothesis of no difference between the curves and consider this a significant variable in the construction of our model.

```
aids_fit_time_k <- survfit(Surv(time, censor) ~ karnof, data=aids)
ggsurvplot(aids_fit_time_k, data=aids, conf.int = TRUE) +
  ggtitle("Karnofsky Performance Score")
```

## Karnofsky Performance Score



'Log-Rank'

```
## [1] "Log-Rank"
```

```
survdif(Surv(time,censor)~karnof, data = aids, rho=0)
```

```
## Call:
```

```
## survdif(formula = Surv(time, censor) ~ karnof, data = aids,
```

```
## rho = 0)
```

```
##
```

```
##      N Observed Expected (O-E)^2/E (O-E)^2/V
```

```
## karnof=70  24      8     1.7    23.43    24.04
```

```
## karnof=80 133     23    10.4    15.26    17.98
```

```
## karnof=90 399     26    32.8     1.41     2.68
```

```
## karnof=100 295     12    24.1     6.08     9.36
```

```
##
```

```
## Chisq= 46.2 on 3 degrees of freedom, p= 5e-10
```

'Wilcoxon'

```
## [1] "Wilcoxon"
```

```
survdif(Surv(time,censor)~karnof, data = aids, rho=1)
```

```
## Call:
```

```
## survdiff(formula = Surv(time, censor) ~ karnof, data = aids,
##      rho = 1)
##
##              N Observed Expected (O-E)^2/E (O-E)^2/V
## karnof=70    24      7.71      1.63      22.69      24.24
## karnof=80   133     22.02      9.97      14.57      17.90
## karnof=90   399     24.82     31.40       1.38       2.75
## karnof=100  295     11.55     23.09       5.77       9.26
##
##  Chisq= 46.3  on 3 degrees of freedom, p= 5e-10
library(survival)
library(survminer)
library(ggplot2)
library(broom)

#### COX PH MODEL USING BACKWARDS SELECTION ####

#full model
cp_full<- coxph(Surv(time_d,censor_d)~.-time -censor, data = aids)
cp_full$loglik

## [1] -129.9983 -107.4059
#cp_full

#reduced model 1
cp_red1<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug, data=aids)
cp_red1$loglik

## [1] -129.9983 -107.4574
#cp_red1

#likelihood ratio test and p-value
s1 <- 2*(cp_full$loglik[2]-cp_red1$loglik[2])
1-pchisq(s1,1)

## [1] 0.7481855
#reduced model 2
cp_red2<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdvd, data=aids)
cp_red2$loglik

## [1] -129.9983 -107.6963
#cp_red2

#likelihood ratio test and p-value
s2 <- 2*(cp_red1$loglik[2]-cp_red2$loglik[2])
1-pchisq(s2,1)

## [1] 0.4894393
#reduced model 3
cp_red3<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdvd -raceth, data=aids)
cp_red3$loglik

## [1] -129.9983 -108.1315
#cp_red3
```

```

#likelihood ratio test and p-value
s3 <- 2*(cp_red2$loglik[2]-cp_red3$loglik[2])
1-pchisq(s3,1)

## [1] 0.3508171

#reduced model 4
cp_red4<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdv -raceth -strat2, data=aids)
cp_red4$loglik

## [1] -129.9983 -109.4399

#cp_red4

#likelihood ratio test and p-value
s4 <- 2*(cp_red3$loglik[2]-cp_red4$loglik[2])
1-pchisq(s4,1)

## [1] 0.1057475

#reduced model 5
cp_red5<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdv -raceth -strat2 -hemophil, d
cp_red5$loglik

## [1] -129.9983 -109.8614

#cp_red5

#likelihood ratio test and p-value
s5 <- 2*(cp_red4$loglik[2]-cp_red5$loglik[2])
1-pchisq(s5,1)

## [1] 0.3585359

#reduced model 6
cp_red6<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdv -raceth -strat2 -hemophil -s
cp_red6$loglik

## [1] -129.9983 -110.5482

#cp_red6

#likelihood ratio test and p-value
s6 <- 2*(cp_red5$loglik[2]-cp_red6$loglik[2])
1-pchisq(s6,1)

## [1] 0.2411968

#reduced model 7
cp_red7<- coxph(Surv(time_d,censor_d)~.-time -censor -txgrp -ivdrug -priorzdv -raceth -strat2 -hemophil -s
cp_red7$loglik

## [1] -129.9983 -111.6966

cp_red7

## Call:
## coxph(formula = Surv(time_d, censor_d) ~ . - time - censor -
##       txgrp - ivdrug - priorzdv - raceth - strat2 - hemophil -
##       sex - id, data = aids)
##
##               coef exp(coef)    se(coef)      z      p
## tx           -8.749e-01  4.169e-01  4.904e-01 -1.784 0.07438

```

```
## karnof      -7.677e-02  9.261e-01  2.561e-02 -2.998 0.00272
## cd4100-150 -1.875e+01  7.201e-09  5.457e+03 -0.003 0.99726
## cd4150-200 -1.861e+01  8.301e-09  6.599e+03 -0.003 0.99775
## cd4200-250 -2.390e-01  7.874e-01  1.068e+00 -0.224 0.82304
## cd4300-350 -1.878e+01  6.980e-09  1.782e+04 -0.001 0.99916
## cd4350+    -1.843e+01  9.912e-09  4.346e+04  0.000 0.99966
## cd450-100  -4.014e-01  6.694e-01  5.968e-01 -0.673 0.50123
## age20-30    1.729e+01  3.212e+07  1.978e+04  0.001 0.99930
## age30-40    1.779e+01  5.299e+07  1.978e+04  0.001 0.99928
## age40-50    1.817e+01  7.814e+07  1.978e+04  0.001 0.99927
## age50-60    1.899e+01  1.759e+08  1.978e+04  0.001 0.99923
## age60-70    1.944e+01  2.762e+08  1.978e+04  0.001 0.99922
## ageover70   1.563e+00  4.772e+00  5.257e+04  0.000 0.99998
##
## Likelihood ratio test=36.6 on 14 df, p=0.0008469
## n= 851, number of events= 20
#likelihood ratio
s7 <- 2*(cp_red6$loglik[2]-cp_red7$loglik[2])
1-pchisq(s7,1)
```

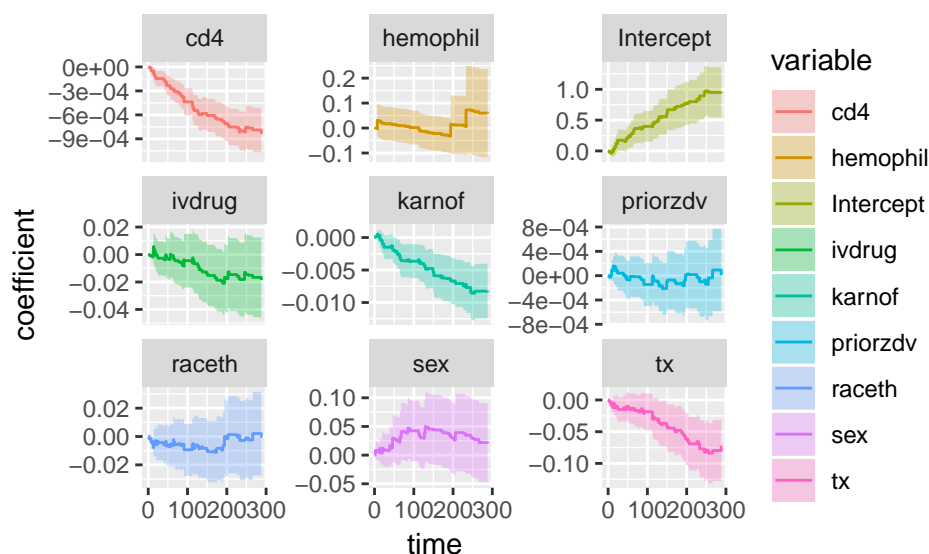
```
## [1] 0.129639
```

To better understand how much the model fit changes with each different explanatory variable, we can use the graphical representation of the Aalen additive regression model. The Aalen model allows for time-varying covariate effects, while the Cox model allows only a common time-dependence through the baseline. In the Aalen model, we have the weighted comparisons of the crude estimate of the hazard rate of each group as compared to a baseline group, which here is defined as the estimate. As we can see, the selected explanatory variables in our model all have an inverse coefficient correlation with the baseline intercept. The slope of an estimated cumulative regression function is positive when covariate increases and this fact correspond to an increasing hazard rate. On the other hand, if the slope is negative while the covariate increases, then this fact points to a decreasing hazard rate.

```
library(ggfortify)
aids <- read.csv( "http://pages.pomona.edu/~jsh04747/courses/math150/AIDSdata.csv")

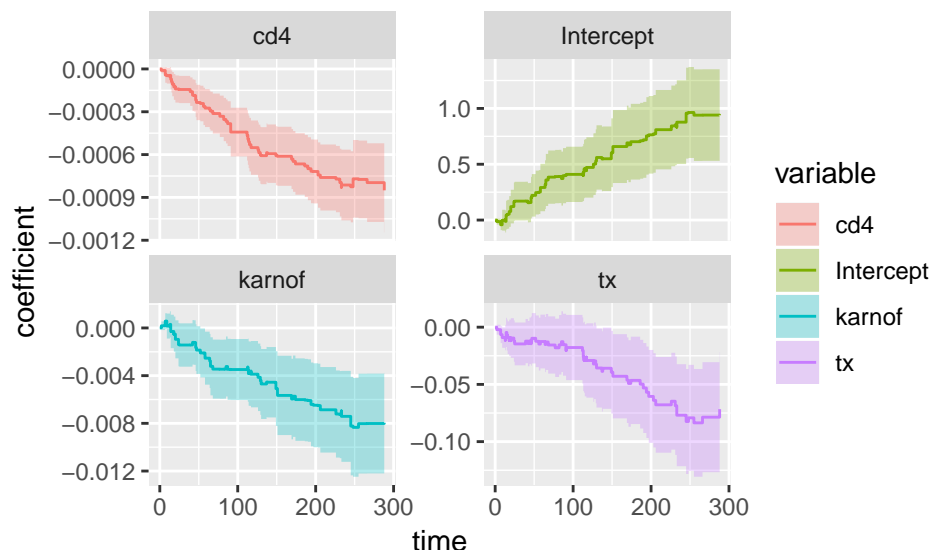
aa_fit <- aareg(Surv(time, censor) ~ cd4 + karnof + priorzdv + hemophil + raceth + sex + tx + ivdrug,
  data = aids)

autoplot(aa_fit, xlab="Coefficient", ylab="Time") + labs(x = "time", y = "coefficient")
```



```
aa_fit2 <-aareg(Surv(time, censor) ~ cd4 + karnof+ tx , data = aids)

autoplot(aa_fit2) + labs(x = "time", y = "coefficient")
```



The Aalen model assumes that the cumulative hazard  $H(t)$  for a subject can be expressed as  $a(t) + X B(t)$ , where  $a(t)$  is a time-dependent intercept term,  $X$  is the vector of covariates for the subject possibly time-dependent, and  $B(t)$  is a time-dependent matrix of coefficients.

The plots show how the effects of the covariates change over time.

## Patricia's “Something New”

I will be doing a power analysis by simulating survival analysis curves

### 1. What is the topic?

The topic is using `sim.survdata` in R to simulate survival data. Using that simulated data, we will make that the alternative and control for the coefficient beta by setting it equal to some value. Then using power analysis, we will see how many times we reject  $H_0$ .

### 2. How it is relevant? How it relates to survival analysis/analysis at hand?

Power analysis relates to survival analysis because if power is large after comparing our data to the simulated survival data, this tells us that there is a high chance that we would reject the null in favor of the alternative (control versus treatment?)

### 3. Resources to learn about the topic.

Below are some of the resources I have begun to use to learn about creating simulations of survival curves and performing power analysis:

a). [https://cran.r-project.org/web/packages/coxed/vignettes/simulating\\_survival\\_data.html](https://cran.r-project.org/web/packages/coxed/vignettes/simulating_survival_data.html) b). [http://www.icssc.org/documents/advbiosgoa/tab%2026.00\\_survss.pdf](http://www.icssc.org/documents/advbiosgoa/tab%2026.00_survss.pdf)

#### 4. What will be challenging about learning something new?

Learning something new will be challenging because in this case, the concept of power analysis is something I just recently learned in Intro to Statistics. So learning to apply this concept in the context of survival analysis curves will be a challenge for me to learn. Learning how to simulate survival curves will also be challenging because I will have to learn how to use and interpret new functions in R.

#### Power Analysis code and simulation

```
simdata <- sim.survdata(N=1000, T=100, num.data.frames=1, beta = c(0.01,0.07,0.3))
head(simdata$data,10)
```

```
##           X1           X2           X3  y failed
## 1  0.1409743 -0.215684124  1.280523331 61   TRUE
## 2  1.2374709  0.249824775 -0.479113239 89   TRUE
## 3 -2.3205264 -0.007826475 -1.904721679 95   TRUE
## 4  0.5333092  0.274866987 -0.670431962 55   TRUE
## 5  0.7080959 -0.517292417 -0.523659530 87  FALSE
## 6 -0.5866109  0.053407869 -1.443545836 87   TRUE
## 7  0.9143624 -0.812643458 -1.336865539 82   TRUE
## 8 -0.7963430 -1.069340109  0.005582356 95   TRUE
## 9  0.8204734 -0.025629577 -1.604662148 88   TRUE
## 10 0.9232624 -1.540258973 -0.050942135 94   TRUE
```

```
simdata$betas
```

```
##      [,1]
## [1,] 0.01
## [2,] 0.07
## [3,] 0.30
```

```
head(simdata$baseline,10)
```

```
##      time  failure.PDF  failure.CDF  survivor      hazard
## 1      1 2.947010e-08 2.947010e-08 1.0000000 2.947010e-08
## 2      2 2.062907e-07 2.357608e-07 0.9999998 2.062907e-07
## 3      3 5.599319e-07 7.956928e-07 0.9999992 5.599321e-07
## 4      4 1.090394e-06 1.886087e-06 0.9999981 1.090395e-06
## 5      5 1.797676e-06 3.683763e-06 0.9999963 1.797680e-06
## 6      6 2.681779e-06 6.365542e-06 0.9999936 2.681789e-06
## 7      7 3.742703e-06 1.010825e-05 0.9999899 3.742727e-06
## 8      8 4.980447e-06 1.508869e-05 0.9999849 4.980498e-06
## 9      9 6.395012e-06 2.148370e-05 0.9999785 6.395109e-06
## 10    10 7.986398e-06 2.947010e-05 0.9999705 7.986569e-06
```

```
#ggsurvplot(survfit(Surv(y,failed) ~ X1 + X2 + X3, data = simdata$data))
```

```
model <- coxph(Surv(y, failed) ~ X1 + X2 + X3, data = simdata$data)
```

```
library(dplyr)
```

```
library(broom)
```

```
model %>% tidy()
```

```
## # A tibble: 3 x 7
```

```
##   term estimate std.error statistic  p.value conf.low conf.high
##   <chr>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1 X1    -0.0367    0.0344    -1.07 2.86e- 1  -0.104    0.0307
## 2 X2     0.0524    0.0339     1.54 1.23e- 1  -0.0141    0.119
## 3 X3     0.287     0.0338     8.51 1.74e-17    0.221     0.354
```

```

set.seed(1234)
n.reps <- 100
simoutput <- c()
for(i in 1:n.reps){
  simdata <- sim.survdata(N=851, T=100, num.data.frames=1, censor= 0.9764,xvars=4, mu=m, sd=s, beta = c(-0
  model <- coxph(Surv(y, failed) ~ X1 + X2 + X3 +X4, data = simdata$data)
  simoutput <- rbind(simoutput, cbind(rep = rep(i, 4), model %>% tidy()))
}

#simoutput

#sum(which(simoutput$p.value < 0.05))
sum(simoutput$p.value < 0.05)

## [1] 65

#simoutput%>%filter(term=="X1")%>%summarize(sum(p.value<0.05))

simoutput%>%dplyr::filter(term=="X1")%>%dplyr::summarize(sum(p.value<0.05))

##    sum(p.value < 0.05)
## 1                      53

simoutput%>%dplyr::filter(term=="X2")%>%dplyr::summarize(sum(p.value<0.05))

##    sum(p.value < 0.05)
## 1                      3

simoutput%>%dplyr::filter(term=="X3")%>%dplyr::summarize(sum(p.value<0.05))

##    sum(p.value < 0.05)
## 1                      3

simoutput%>%dplyr::filter(term=="X4")%>%dplyr::summarize(sum(p.value<0.05))

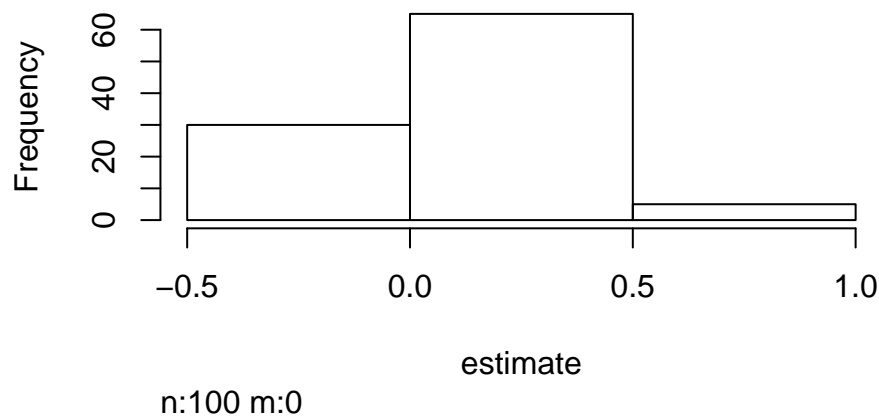
##    sum(p.value < 0.05)
## 1                      6

m <-c(mean(aids$tx), mean(aids$karnof), mean(aids$cd4), mean(aids$age))
s <-c(sd(aids$tx), sd(aids$karnof), sd(aids$cd4), sd(aids$age))

simoutput%>%dplyr::filter(term=="X4")%>%dplyr::select(estimate)%>%hist(breaks=200)

```





```
library(FDRsampsiz)
library(powerSurvEpi)

### POWER ANALYSIS AND SIMULATION ###
set.seed(1234)
n.reps <- 100
simoutput2 <- c()
for(i in 1:n.reps){
  simdata2 <- sim.survdata(N=851, T=362, num.data.frames=1, censor= 0.2,
                           X=aids[,c(6,13)], #,13,14,16) ,
                           beta = c(-0.867409, -0.071620))#, -0.016659, 0.073674))
  model <- coxph(Surv(y, failed) ~ tx + karnof, data = simdata2$data)
  simoutput2 <- rbind(simoutput2, cbind(rep = rep(i, 2), model %>% tidy()))
}

simoutput2%>%dplyr::filter(term=="tx")%>%dplyr::summarize(sum(p.value<0.05))

## sum(p.value < 0.05)
## 1 2

simoutput2%>%dplyr::filter(term=="karnof")%>%dplyr::summarize(sum(p.value<0.05))

## sum(p.value < 0.05)
## 1 4

simoutput2%>%dplyr::filter(term=="cd4")%>%dplyr::summarize(sum(p.value<0.05))

## sum(p.value < 0.05)
## 1 0

simoutput2%>%dplyr::filter(term=="age")%>%dplyr::summarize(sum(p.value<0.05))

## sum(p.value < 0.05)
## 1 0
```

# Juste's "Something New": The Schoenfeld Residuals for the Cox PH model

Cox proportional hazards (PH) model is considered a great way to identify combined effects of several covariates on the relative risk (hazard). This model assumes that the hazards of the different strata formed by the levels of the covariates are proportional. This proportional hazards assumption is particularly important and can be tested via three different classes of tests. The first class is focused on the piece-wise estimation of models for subsets of data defined by stratification of time. The second one considers the interactions between covariates and some function of time. Final, third one is based on examinations of regression residuals. The Schoenfeld Residuals are a part of the third class of proportional hazard assumption testing and I will be exploring it in order to be able to eradicate a method for testing for the PH assumption in the current and future data set analyses. This topic is particularly important in relation to survival analysis since it provides an idea of whether the model is appropriate for the data set at hand and whether some covariates should be considered as variants of time in order to supply the best model for prediction of proportional hazards. Taking a completely new model of analyzing survival data is particularly difficult since the mathematical derivations and notations are also very varied from what we have seen in class. Although, I do remember some of the ideas behind parametric functions, their applications to statistical models are much more challenging than I have expected. Therefore, it will require me a lot of time and extensive research to be able to understand and learn how to apply this model to our data and other instances of survival analysis.

### 3. Resources to learn about the topic.

I have been researching articles and scientific journals that provide insights into the Schoenfeld residuals and their use in the Cox PH model. Sources include:

1. <https://onlinelibrary.wiley.com/doi/full/10.1111/ajps.12176>
2. [https://rstudio-pubs-static.s3.amazonaws.com/39354\\_34153ff19e624116bd2fbdec7d2534aa.html](https://rstudio-pubs-static.s3.amazonaws.com/39354_34153ff19e624116bd2fbdec7d2534aa.html)

### Explanation of the Theory Behind Schoenfeld Residuals

Let  $z_{ij}(t)$  be the  $j^{th}$  covariate of the  $i^{th}$  unit, where  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, p$

This notation indicates that  $z_{ij}$  is allowed to vary as a function of the time scale.

- 1) As we know from lecture, the Cox PH model assumes that  $h(t)$  of the  $i^{th}$  individual satisfies:
  - $h_i(t) = h_0(t)e^{z_i(t)\beta}$  where:
  - $h_0 \rightarrow$  baseline hazard
  - $z_i(t) \rightarrow 1 \times p$  vector of covariates for unit  $i$  each of which can be time fixed or time-varying.
- 2) However, another possibility has been presented by Therneau and Granbsh in 2000, where they proposed an idea that there could be an alternative to the current Cox model, where the coefficient of the estimate could also be varying as a function of time.

The new hazard function would look like this:  $h_i(t) = h_0(t)e^{z_i(t)\beta(t)}$

Therefore, in order to examine thee two models in a case when  $\beta = \beta(t)$  requires a residual analysis that could indicate whether a model should consider a covariate as a variable with time.

---

Due to the fact that that some observations might be censored and in particular, regarding the Cox PH model, the baseline hazard is not estimated, in order to analyse the residuals a particular score process. The risk score for unit  $i$  at time  $t$  is thought to be  $r_i(t) = e^{z_i(t)\beta}$ , where  $Y_i(t)$  is the indicator function and  $Y_i(t) = 1$  indicates a point in which  $i$  is under risk and thus observation and it is equal to 0 in other occasions.

---

The Schoenfeld residuals are given by the equations:

1.  $s_k = Z_{(k)} - \frac{\sum_i Y_i(t_k)r_i(t_k)Z_i(t_k)}{\sum_i Y_i(t_k)r_i(t_k)}$

$$2. s_k = Z_{(k)} - \bar{z}(\hat{\beta}, t_k)$$

In this case, the  $Z(k)$  is the covariate vector of the particular unit that is experiencing the event at time  $k$ ;  $\hat{\beta}$  is the estimate of  $\beta$  and  $\bar{z}(\hat{\beta}, t_k)$  is the weighted mean of covariate values.

Furthermore, the weighted variance can be represented by the derived equation at the  $k^{th}$  time as

$$V(\beta, t_k) = \sum_i Y_i(t_k) r_i(t_k) Z_i(t_k) - \bar{z}(\hat{\beta}, t_k)' Z_i(t_k) - \frac{\bar{z}(\hat{\beta}, t_k)}{\sum_i Y_i(t_k) r_i(t_k)}$$

From this, we can scale the Schoenfeld residuals by  $V(\beta, t_k)$  of  $X$  at  $t_k$  via the equation:

$$s_k^* = V^{-1}(\hat{\beta}, t_k) s_k$$

The scaled Schoenfeld residuals can also be defined as follows:

$$s_k^* = m \sum_{k=1}^d V(\hat{\beta}, t_k) s_k$$

here,  $m$  is the total number of deaths in the data set.

Following the calculations, the residuals are plotted against time in order to test the proportional hazards assumption. If the assumption is correct, the residuals should be fitting around the line centered at zero ( $y=0$ ). The further away this predicted line is from the horizontal of ( $y=0$ ) the more likely one is to call the PH assumption to question and determine whether it is met through the model.

---

To go a little deeper into the analysis of the residual calculation, one can look at the calculations of the test statistic for this residual model.

By producing a least squares slope of regression and assuming a relationship between  $s_{kj}^*$  and  $t_{kj}$  or some function  $g(t_k)$  allows to derive a test statistic for the proportional hazards assumption in regards to the  $j^{th}$  covariate, which is given by:

$$T_j = \frac{[\sum_{k=1}^d (g(t_k) - \hat{g}) s_{kj}^*]^2}{d I_{jj} \sum_{k=1}^d (g(t_k) - \hat{g})^2}$$

Here, the distribution is asymptotically as  $X^2(1)$  stating the null hypothesis that the relationship between the covariate, in this case  $j$  and the event time follows the assumption of PH.

---

Interpretation of Schoenfeld Residuals from plots in R and the p-values presented.

The y-axis of the Schoenfeld residuals graph can be interpreted as the log of the hazard ratio for the explanatory variable— the coefficient in Cox's model if it were allow to vary over time. If the graph is flat, then the PH assumption is adequate. Furthermore, the Schoenfeld residuals are independent of time. A plot that shows a non-random pattern against time is evidence of violation of the PH assumption. The PH assumption is supported when there's a non-significant relationship between residuals and time. ### HIV Data Cox PH model analysis using Schoenfeld Residuals

Schoenfeld Residuals applied to our best Cox PH model for AIDS data where, we have an additive model of explanatory variables: baseline CD4 count, iv drug use history, and karnofsky performance scale score:

```
cph_r10 <- coxph(Surv(time,censor)~.-priorzdv -id -hemophil -raceth -time_d -strat2
               -sex -txgrp -age -tx -censor_d, data = aids)
cph_r10
```

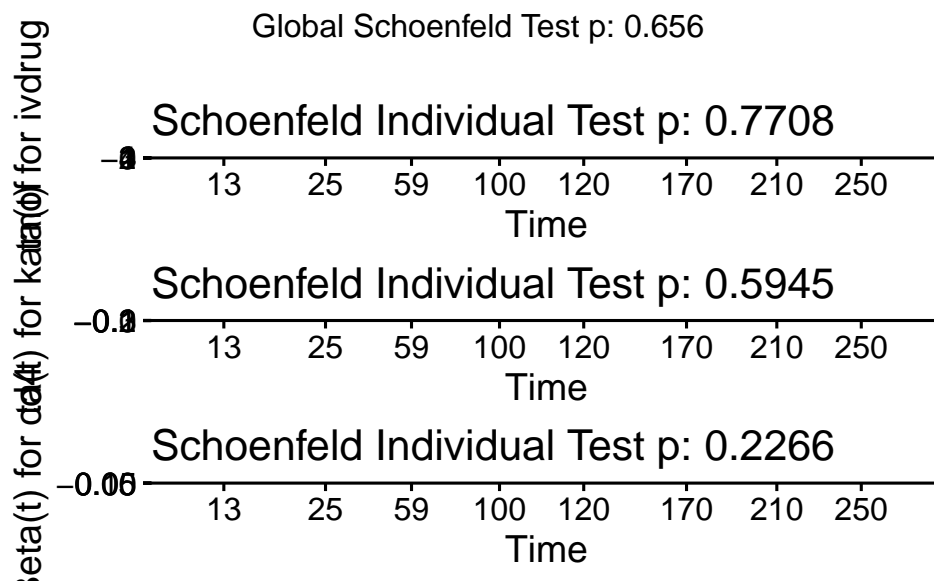
```
## Call:
## coxph(formula = Surv(time, censor) ~ . - priorzdv - id - hemophil -
##       raceth - time_d - strat2 - sex - txgrp - age - tx - censor_d,
##       data = aids)
##
##              coef exp(coef)    se(coef)      z        p
## ivdrug  -0.216832   0.805065   0.180491 -1.201     0.23
## karnof  -0.061043   0.940783   0.014157 -4.312 1.62e-05
```

```
## cd4      -0.015127  0.984987  0.003076 -4.917 8.77e-07
##
## Likelihood ratio test=69.33  on 3 df, p=5.947e-15
## n= 851, number of events= 69
```

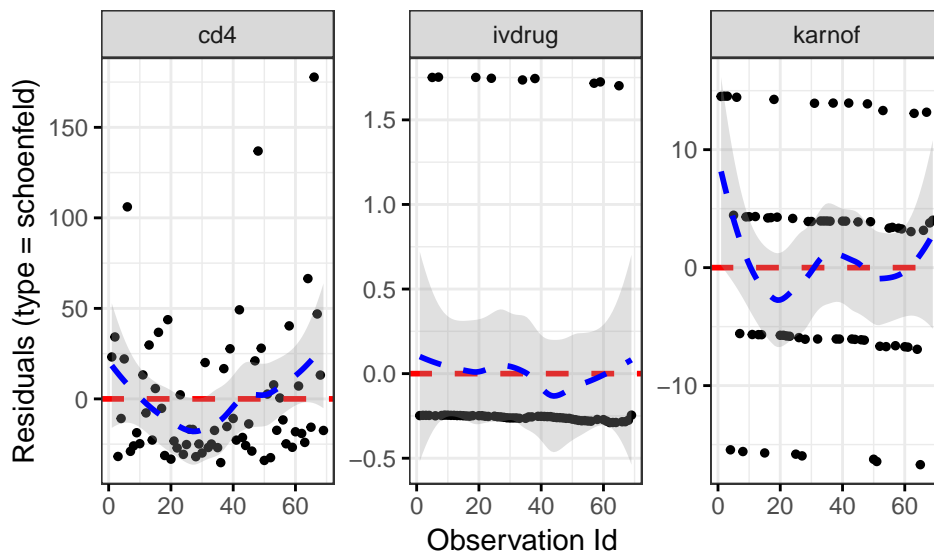
```
zph_r10 <- cox.zph(cph_r10)
zph_r10
```

```
##           rho  chisq      p
## ivdrug -0.0348 0.0849 0.771
## karnof -0.0630 0.2834 0.595
## cd4      0.1524 1.4618 0.227
## GLOBAL      NA 1.6150 0.656
```

```
ggcoxzph(zph_r10)
```



```
ggcoxdiagnostics(cph_r10, type="schoenfeld")
```



Using the best determined Cox PH model for our data, we can look at the Schoenfeld residuals to determine if the PH assumption is met. Via the function “ggcoxzph()”, which produces, for each covariate, graphs of the scaled Schoenfeld residuals against the transformed time. Here, the solid line is a smoothing spline fit to the plot, with the

dashed lines representing a  $\pm 2$ -standard-error. from these graphs, we don't see any patterns or significance of the residual fit regarding the graphs of the covariates with time. Therefore, the assumption of proportional hazards seems to be supported for the covariates: baseline CD4 count, iv drug use history, and karnofsky performance scale score.

Using the `ggcoxdiagnostics()` function we can provide another graphic representation of the residual distribution in regards to the covariates with time. Here, we also see that there's no particular pattern of the residuals around the line of fit, therefore again, we can state that the PH assumption has been met.

References: 1. [http://www.ukm.my/jsm/pdf\\_files/SM-PDF-46-3-2017/15%20Aditif%20Aalen.pdf](http://www.ukm.my/jsm/pdf_files/SM-PDF-46-3-2017/15%20Aditif%20Aalen.pdf)