

Anthony Orso

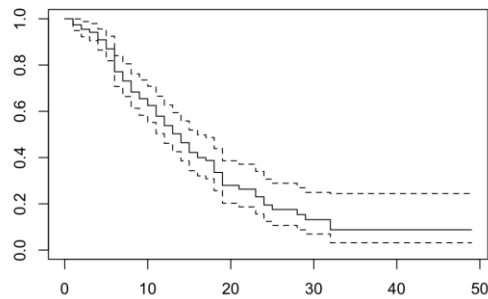
STAT 411

December 18, 2023

Question 1 Report

I chose to analyze the burn dataset, which has 11 covariates and three different event types. I began by printing out the summary statistics. All the variables are categorical, except for the percentage of body covered in burns. Because of this, there are plenty of opportunities to test survival rates in certain categorical variables. I focused on the event and delta indicator for measuring the time to excision. Excision is used in cases of deep burns to reduce mortality. Thus, the presence of excision indicates a serious medical case. I'm curious to see if certain indicator variables signal greater medical intervention, measured as an excision in this case. I am curious to view the survival curve to see how fast events are happening, and I will see if a Nelson-Aalen and Kaplan-Meier estimator are similar. Then I will look at hazard rates with log-rank tests—including trend tests and the use of strata—to see a group-based differences in hazard and survival to make conclusions about what precludes medical intervention and, thus, more serious cases. Then I will use modeling to predict survival, and I will use model fit metrics to evaluate and compare various models.

To begin, a simple survival curve was generated using a Kaplan-Meier estimator.



Overall, the survival curve is consistent in shape. The rate of change seems to be fairly similar across each timepoint, so I believe an AFT model will likely perform well with it since exponential functions have a constant hazard rate. I also tested a random time-point of $t=30$, and $S(x)$ at that time is estimated to be 0.132. Around that time, the survival curve flatlines. In essence, around $t=30$, all events have happened. This tracks with our understanding of burn victims. Burns are very serious injuries that can cause permanent disfiguration, nerve damage, and death. Excisions involve removing flesh that prevents wounds from healing, and while it can be vital for preventing mortality, it comes with significant risk of blood loss and infection. As such, it is a decision that must be made carefully while closely monitoring patient vitals. Thus, it tracks that excisions happen at a consistent pace until plateauing when no more events happen. In these instances, I assume it would be less severe injuries.

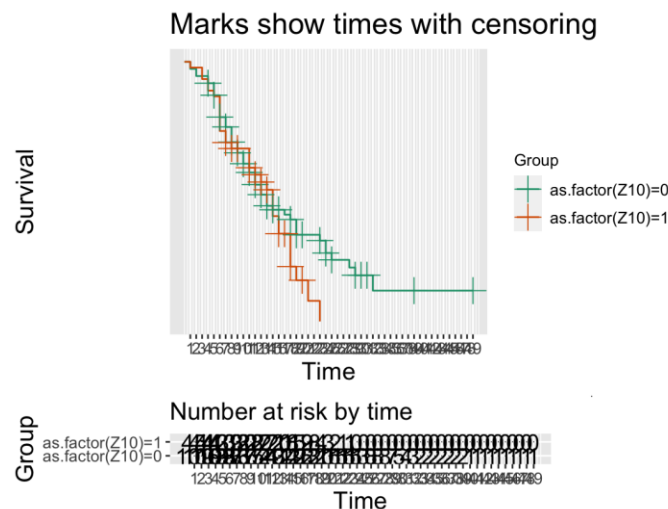
I then manually calculated a Nelson-Aalen estimator and plotted it against the Kaplan-Meier estimator.

Figure 1 is a Kaplan-Meier survival plot. The Y-axis represents Survival Probability, ranging from 0.0 to 1.0. The X-axis represents Time, ranging from 0 to 30. Two curves are plotted: a blue line for Kaplan-Meier and a red line for Nelson-Aalen. Both curves show a decreasing trend in survival probability over time. The Nelson-Aalen curve is slightly higher than the Kaplan-Meier curve in the early stages, but they converge as time increases.

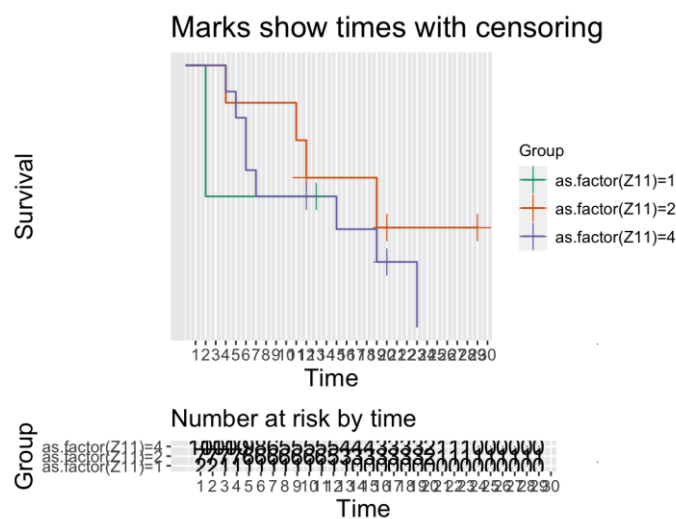
I now want to move into hypothesis testing to see if there are meaningful differences in survival rates based on groupings. I was particularly interested in seeing if different burn types (chemical, scald, flame, and electric) produce different survival, and the log-rank test showed that survival was indeed different with a p-value of 0.04. Since there are four levels, I wanted to see the possibility of conducting a trend test. The autoplot function produced the following plot.

[illegible]

I also saw the respiratory indicator variable and wondered if a burn near the respiratory tract would have meaningfully different survival rates from non-respiratory burns. My instinct is that it would increase survival times, as I can imagine it's not feasible to just cut out a bunch of flesh around the lungs since it is a major organ. However, the log-rank test produced a p-value of 0.2, so we conclude survival is the same regardless of whether the burn impacted the patient's lungs. The autoplot confirms this.



I then wanted to see if using strata would change our understanding of survival. Since previous research already showed meaningful differences among the burn types, I wanted to see if layering various strata would also produce different survival rates. Gender failed to achieve statistical significance, although the p-value was directional at 0.08. When I plotted the survival curves by burn type for both male and female, the survival curves did indeed look quite different. However, one of the genders had no burns of type 3, so I'm guessing this introduced uncertainty into the data. When I tested survival based on burn type and stratified on the race variable (0 for nonwhite, 1 for white), alpha was exactly equal to 0.05. The interpretation is dependent on whether you describe alpha as strictly less than 0.05 or less than or equal to 0.05. In my case, I do not elect for the strictly less than interpretation. The non-white survival curve is below, followed by the white survival curve.



Survival

Group

- as.factor(Z11)=1
- as.factor(Z11)=2
- as
- as

Time

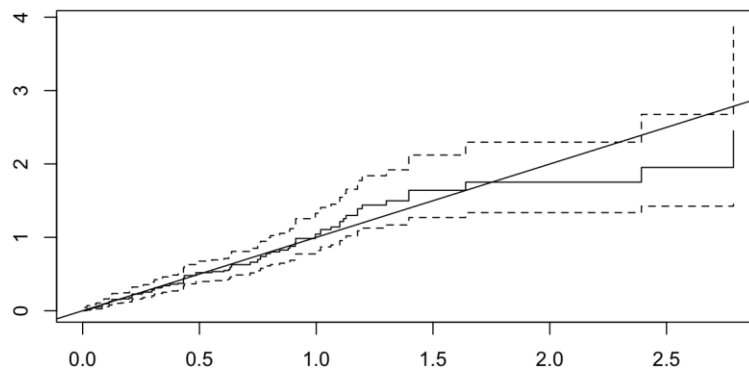
Number at risk by time

Group

Time

After testing various hypotheses with the log-rank test, I now move onto the modeling phase. I will be exploring both Cox Proportional-Hazards and Accelerated Failure Time models. I first started by building a Cox PH with the treatment type, race, gender, respiratory indicator, and burn type variables. The only variables that achieved statistical significance were treatment type and gender. Based off the coefficients for these two variables, we can see what impacts this has on the hazards ratio. Body cleansing increases the risk of excision by 80% compared to routine cleaning. Also, female sex increases the risk of excision by 92% compared to male sex. I suspect women are more likely to have an excision because doctors often consider damage to reproductive organs when working with female patients. A more direct option is women are more likely to be victims of severe burns, possibly due to domestic violence and sexual assault. After doing my research, I did see that women are indeed more likely to be burn victims.

All model fit tests indicate statistical significance of the model compared to a null intercept model. I then ran the model through the `step()` function, and it suggested reducing the model to just sex, cleaning type, and burn type. Although the burn type was not statistically significant, it resulted in the optimal AIC, which showed the usefulness of the parameter that I was able to discern in the hypothesis testing of the report. However, the difference in AIC was minute when removing the burn type parameter, so my final Cox PH model includes just the cleaning type and sex parameter. I also did a residual plot to check model assumptions.



The model overall fits well, although it is less unstable at the tail end. Although this is typical with survival curves, I'm hoping fitting a parametric model will solve this issue.

I wanted to see if an AFT model would work well since the hazard rate in the overall survival curve appeared to be relatively stable, which indicated the data potentially followed a distribution in the exponential family. I tested both a Weibull and exponential distribution, but the Weibull distribution lowered the AIC the greatest. According to the Weibull AFT, survival time was decreased by 32% for body cleansing compared to routine cleaning, and female sex decreased survival time by 39% compared to male sex. These track with the Cox PH model assumptions, but allowing the model to fit based on a

parametric distribution measurably lowered the AIC. For this reason, I elect to use the parametric model to predict time to excision in this data set.

```
#Final Exam
```

```
library(survival)
```

```
library(KMsurv)
```

```
library(tidyverse)
```



```
library(survMisc)
```

```
##Q1
```

```
data(burn)
```

```
summary(burn)
```

```
burnSurv = survfit(Surv(T1,D1)~1,data=burn)
```

```
info = summary(burnSurv)
```

```
plot(burnSurv)
```

```
summary(survfit(Surv(T1,D1)~1,conf.type='plain',data=burn),times=30)
```

```
?burn
```

```
surv_prob <- info$surv
```

```
time_points <- info$time
```

```
n_at_risk <- info$n.risk
```

```
n_events <- info$n.event
```

```
n_censored <- info$n.censor
```

```
summary_df <- data.frame(
```

```
  Time = time_points,
```

```
  Survival_Probability = surv_prob,
```

```
  Number_at_Risk = n_at_risk,
```

```
  Number_of_Events = n_events,
```

```
  Number_Censored = n_censored
```

```
)
```

```
Htilde=cumsum(summary_df$Number_of_Events/summary_df$Number_at_Risk)
```

```
survs = data.frame(time=info$time, survKM = info$surv, survNA = exp(-1*Htilde))
```

```
plot(0, 0, type = "n", xlim = c(0, 30), ylim = c(0, 1),
```

```
  xlab = "Time", ylab = "Survival Probability",
```

```
  main = "Kaplan-Meier and Nelson-Aalen Survival Estimates")
```

```
# Plot Kaplan-Meier estimate
```

```
lines(survs$survKM, col = "blue")
```

```
# Plot Nelson-Aalen estimate
```

```
lines(survvs$urvNA, col = "red")
```

```
# Add legend
```

```
legend("topright", legend = c("Kaplan-Meier", "Nelson-Aalen"), col = c("blue", "red"), lty = 1)
```

```
fit0 = survdiff(Surv(T1,D1)~as.factor(Z11),data=burn)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z11),data=burn),timeTicks='custom', times=1:49,lty=1:4)
```

```
trend= c(3,2,4,1)
```

```
z <- t(trend)%*%(fit0$obs - fit0$exp) / sqrt(t(trend) %*% fit0$var %*% trend)
```

```
z^2
```

```
1-pchisq(z^2,df=1)
```

```
survdiff(Surv(T1,D1)~as.factor(Z10),data=burn)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z10),data=burn),timeTicks='custom', times=1:49,lty=1:2)
```

```
##
```

```
survdiff(Surv(T1,D1)~as.factor(Z11)+strata(as.factor(Z2)),data=burn)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z11),data=burn[burn$Z2==0,]),timeTicks='custom',
```

```
times=1:49,lty=1:4)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z11),data=burn[burn$Z2==1,]),timeTicks='custom',  
times=1:49,lty=1:4)
```

```
survdifff(Surv(T1,D1)~as.factor(Z11)+strata(as.factor(Z3)),data=burn)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z11),data=burn[burn$Z3==0,]),timeTicks='custom',  
times=1:49,lty=1:4)
```

```
autoplot(survfit(Surv(T1,D1) ~ as.factor(Z11),data=burn[burn$Z3==1,]),timeTicks='custom',  
times=1:49,lty=1:4)
```

```
##
```

```
coxFit <-
```

```
coxph(Surv(T1,D1)~as.factor(Z1)+as.factor(Z2)+as.factor(Z3)+as.factor(Z10)+as.factor(Z11),data=burn)
```

```
summary(coxFit)
```

```
step(coxFit)
```

```
coxReduce <- coxph(Surv(T1,D1)~as.factor(Z1)+as.factor(Z2)+as.factor(Z11),data=burn)
```

```
summary(coxReduce)
```

```
anova(coxFit,coxReduce)
```

```
AIC(coxReduce)
```

```
AIC(coxSmaller)
```

```
?burn
```

```
coxSmaller = coxph(Surv(T1,D1)~as.factor(Z1)+as.factor(Z2),data=burn)
```

```
summary(coxSmaller)
```

```
anova(coxFit,coxSmaller)
```

```
AIC(coxSmaller)
```

```
#Z1 class 1 increases the hazards rate by 86% compared to baseline
```

```
#Z2 class 1 increases hazards rate by 104%
```

```
##
```

```
resid = -coxSmaller$residuals + burn$D1
```

```
plot(survfit(Surv(resid,burn$D1)~1,conf.type='none'),fun='cumhaz')
```

```
abline(0,1)
```

```
##
```

```
weib.fit<-survreg(Surv(T1,D1)~as.factor(Z1)+as.factor(Z2),data=burn,dist='weibull')
```

```
summary(weib.fit)
```

```
AIC(weib.fit)
```

```
1-exp(weib.fit$coefficients)
```

```
exp.fit <- survreg(Surv(T1,D1)~as.factor(Z1)+as.factor(Z2),data=burn,dist='exponential')
```

```
summary(exp.fit)
```

```
AIC(exp.fit)
```

```
##
```