

Part III: Survival Models

Load Packages

Again, we must load the packages that will be used in the first part of this workshop.

```
library(pastecs, quietly = TRUE)
library(lm.beta, quietly = TRUE)
library(lmtest, quietly = TRUE)
library(foreign, quietly = TRUE)
library(lattice, quietly = TRUE)
library(lme4, quietly = TRUE)
library(nlme, quietly = TRUE)
library(survival, quietly = TRUE)
library(dplyr, quietly = TRUE)
library(ggfortify, quietly = TRUE)
library(survminer, quietly = TRUE)
library(rms, quietly = TRUE)
library(MASS, quietly = TRUE)
```

Introduction

Survival models concerns the analysis of the time for an event to occur. The response variable is the time for the event to occur. The event is generally called “death.”

Definitions

Survival models involve two functions:

- $S(t)$: the survival function. $S(t)$ is probability that death has not occurred until after time, t .
- $\lambda(t)$: the hazard function.

$$\lambda(t) = \frac{\text{probability of dying in at time, } t}{\text{probability of survival until time, } t} \approx \frac{\text{number of people who died at time, } t}{\text{number of people who lived until time, } t}$$

$\lambda(t)$ measures the likelihood of death in a very small time interval, t and $t + dt$. It is a measure of *risk*.

- $\Lambda(t)$: the cumulative hazard. It is total hazard from 0 to time, t .

Note that these two functions are related:

- $\lambda(t) \leftrightarrow S(t)$

$$\lambda(t) = -\frac{d}{dt} \log S(t)$$

- $\lambda(t) \leftrightarrow \Lambda(t)$

$$\Lambda(t) = \int_0^t \lambda(t) dt$$

- $S(t) \leftrightarrow \Lambda(t)$ and $S(t) \leftrightarrow \lambda(t)$

$$S(t) = \exp(-\Lambda(t)) = \exp\left(-\int_0^t \lambda(t) dt\right).$$

Censoring

Like most models, survival models are susceptible to imperfect data. Let's say a subject is recorded for a study up until a time, t^* . After time t^* , the subject may decide not to continue with study or it is not possible to locate the subject. Many things could have caused a lack of follow up. This subject is called *censored*. While it may be reasonable to discard this data point, the censored data actually contains information that we know the event has not occurred prior to t^* . This gives more information to our model about time prior to t^* than if we were to discard the censored data.

Data

Description

We will be working the `colon` data set. This data comes from one of the first successful trials of a drug for colon cancer. The recurrence and death times are recorded for all patients in the study.

The `colon` dataset has the following columns:

- **id:** id
- **study:** 1 for all patients
- **rx:** Treatment - Observation, Lev(amisole), Lev(amisole)+5-FU. Levamisole is a low-toxicity compound previously used to treat worm infestations in animals; 5-FU is a moderately toxic (as these things go) chemotherapy agent.
- **sex:** 0 = female, 1 = male
- **age:** age of the patient
- **obstruct:** 0 = if tumour did not obstruct colon, 1 = if tumour obstructed colon
- **perfor:** perforation of colon
- **adhere:** adherence to nearby organs
- **nodes:** number of lymph nodes with detectable cancer
- **time:** days until event or censoring
- **status:** censoring status
- **differ:** differentiation of tumour (1=well, 2=moderate, 3=poor)
- **extent:** Extent of local spread (1=submucosa, 2=muscle, 3=serosa, 4=contiguous structures)
- **surg:** time from surgery to registration (0=short, 1=long)
- **node4:** more than 4 positive lymph nodes
- **etype:** event type: 1=recurrence, 2=death

```
attach(colon)
head(colon)
```

```
##   id study      rx sex age obstruct perfor adhere nodes status differ
## 1  1     1 Lev+5FU  1  43         0      0      0      5      1      2
## 2  1     1 Lev+5FU  1  43         0      0      0      5      1      2
## 3  2     1 Lev+5FU  1  63         0      0      0      1      0      2
## 4  2     1 Lev+5FU  1  63         0      0      0      1      0      2
## 5  3     1    Obs  0  71         0      0      1      7      1      2
## 6  3     1    Obs  0  71         0      0      1      7      1      2
##   extent surg node4 time etype
## 1      3     0     1 1521     2
## 2      3     0     1  968     1
## 3      3     0     0 3087     2
## 4      3     0     0 3087     1
## 5      2     0     1  963     2
## 6      2     0     1  542     1
```

Subsetting data and converting data

We will be studying the recurrence event of colon cancer.

```
colon_subset_recurrence = colon[colon$etype==1,]
```

Some survival models can only handle variables encoded in 0 and 1. We need to convert continuous variables, such as age and nodes, to 0 and 1.

```
colon_subset_recurrence$age.ds = sapply(colon_subset_recurrence$age,  
                                         function(x) ifelse(x > 60, 1, 0))
```

If the binary variables are stored as `numeric` variables, the survival models will treat the explanatory variables as continuous variables rather than as discrete variables.

```
sapply(colon,class)
```

```
##      id      study      rx      sex      age obstruct  perfor  
## "numeric" "numeric" "factor" "numeric" "numeric" "numeric" "numeric"  
##   adhere    nodes    status   differ    extent    surg    node4  
## "numeric" "numeric" "numeric" "numeric" "numeric" "numeric" "numeric"  
##    time    etype  
## "numeric" "numeric"
```

Many discrete variables are stored as `numeric` variables. We have to convert these columns to `factor`.

The `factor` takes as arguments:

- the discrete data in the first argument
- `level` is current coding the discrete data. This is an optional argument.
- `label` is the encoding that you would like to change to discrete data. This is an optional argument. Use this argument if you would to change the labeling of the discrete data.

```
colon_subset_recurrence$age.ds <- factor(colon_subset_recurrence$age.ds,  
                                         levels= c("0","1"),  
                                         labels=c("<60",">60"))
```

```
colon_subset_recurrence$node4 <- factor(colon_subset_recurrence$node4,  
                                         levels= c("0","1"),  
                                         labels=c("<4",">4"))
```

```
colon_subset_recurrence$sex <- factor(colon_subset_recurrence$sex,  
                                       levels= c("0","1"), labels=c("F","M"))
```

```
colon_subset_recurrence$obstruct <- factor(colon_subset_recurrence$obstruct,  
                                           levels= c("0","1"),  
                                           labels=c("no obstruct","obstruct"))
```

```
colon_subset_recurrence$adhere <- factor(colon_subset_recurrence$adhere,  
                                         levels= c("0","1"),  
                                         labels=c("no adhere","adhere"))
```

```
colon_subset_recurrence$perfor <- factor(colon_subset_recurrence$perfor,  
                                         levels= c("0","1"),  
                                         labels=c("no perfor","perfor"))
```

```
colon_subset_recurrence$differ <- factor(colon_subset_recurrence$differ,  
                                         levels= c("1","2","3"),  
                                         labels=c("well","mod","poor"))
```

```
colon_subset_recurrence$extent <- factor(colon_subset_recurrence$extent,  
                                         levels= c("1","2","3","4"),
```

```

labels=c("submucosa", "muscle", "serosa", "contiguous"))
colon_subset_recurrence$surg <- factor(colon_subset_recurrence$surg,
levels= c("0","1"),
labels=c("short","long"))

```

Now, let's take a look at the data.

```
head(colon_subset_recurrence)
```

```

##      id study      rx sex age   obstruct   perfor   adhere nodes status
## 2      1      1 Lev+5FU  M  43 no obstruct no perfor no adhere     5      1
## 4      2      1 Lev+5FU  M  63 no obstruct no perfor no adhere     1      0
## 6      3      1      Obs  F  71 no obstruct no perfor   adhere     7      1
## 8      4      1 Lev+5FU  F  66   obstruct no perfor no adhere     6      1
## 10     5      1      Obs  M  69 no obstruct no perfor no adhere    22      1
## 12     6      1 Lev+5FU  F  57 no obstruct no perfor no adhere     9      1
##      differ extent  surg node4 time etype age.ds
## 2      mod serosa short    >4  968      1    <60
## 4      mod serosa short    <4 3087      1    >60
## 6      mod muscle short    >4  542      1    >60
## 8      mod serosa long    >4  245      1    >60
## 10     mod serosa long    >4  523      1    >60
## 12     mod serosa short    >4  904      1    <60

```

Surv Object

The `Surv` function takes as input the time and censoring status (0 or 1) of a data point. It returns a object that packages together time and censoring status.

```

surv <-with(colon_subset_recurrence, Surv(time,status))
head(surv)

```

```
## [1] 968 3087+ 542 245 523 904
```

The + at the end of the time indicates that the data point was censored.

Kalpan-Meier Estimator

First, let t_i be the i th recorded time in the data. That is, t_1 is the 1st recorded time, t_2 is the 2nd recorded time, \dots , t_{20} is the 20th recorded, etc.

Kalpan-Meier assumes that the survival function can be estimated as

$$\hat{S}(t) = \prod_{\text{for } i: t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

where d_i is the number of persons that “died” after time t_i and n_i is the number of uncensored persons that have lived up to t_i .

Kalpan-Meier Estimator for the entire data

To fit $\hat{S}(t) = \prod_{\text{for } i: t_i \leq t} 1 - \frac{d_i}{n_i}$ to the entire data, we use the command below.

```
km_fit <- survfit(surv~1, data=colon_subset_recurrence)
```

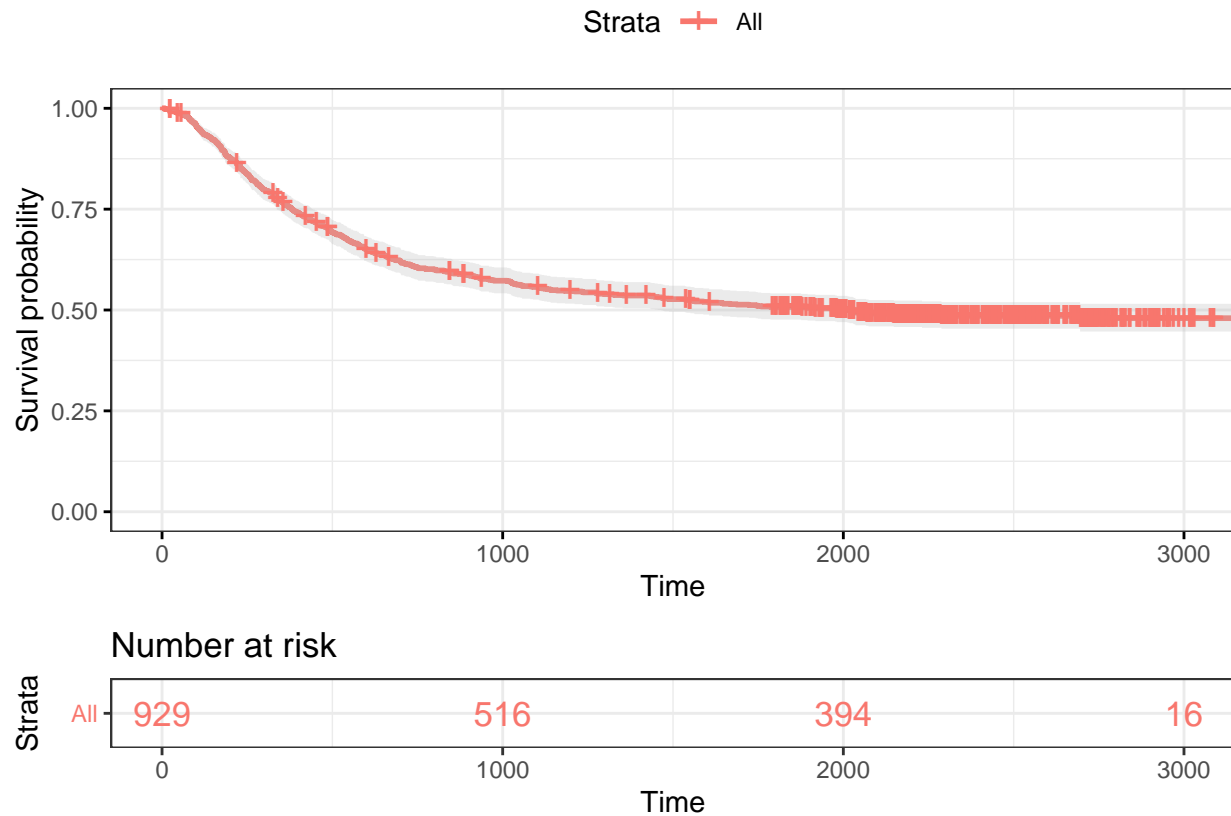
We can return a summary of the $\hat{S}(t)$ at certain time points. `summary(km_fit)` will return a summary `km_fit` for all time points in the data. Since the data set is large, I do not run the command, `summary(km_fit)`. However, you are free to do this on your own.

```
summary(km_fit, times=c(1, 10, 20, 30, 40, 50))
```

```
## Call: survfit(formula = surv ~ 1, data = colon_subset_recurrence)
##
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##    1     929      0   1.000 0.00000    1.000    1.000
##   10     927      2   0.998 0.00152    0.995    1.000
##   20     926      2   0.996 0.00215    0.991    1.000
##   30     922      1   0.995 0.00240    0.990    0.999
##   40     919      4   0.990 0.00322    0.984    0.997
##   50     914      3   0.987 0.00371    0.980    0.994
```

There is a convenience function `ggsurvplot` that generates a plot for a `survfit` object. `conf.int = TRUE` shows the confidence interval around the estimate. `risk.table = TRUE` shows a tabulation of risk below $\hat{S}(t)$.

```
ggsurvplot(km_fit, data = colon_subset_recurrence,
            conf.int = TRUE, risk.table = TRUE,
            ggtheme = theme_bw(),
            risk.table.col = "strata")
```



Kalpan-Meier Estimator for the data divided into obstruct and no obstruct

colon_subset_recurrence can be divided two data sets by the `obstruct` column. Those patients whose colons are obstructed by the tumour and those whose colons aren't. We can fit to each data partition to a Kalpan-Meier Estimator

$$\hat{S}_{\text{obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{obstruct}_i = \text{obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

$$\hat{S}_{\text{no obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{obstruct}_i = \text{no obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

to the entire data. To do this, we use the command below.

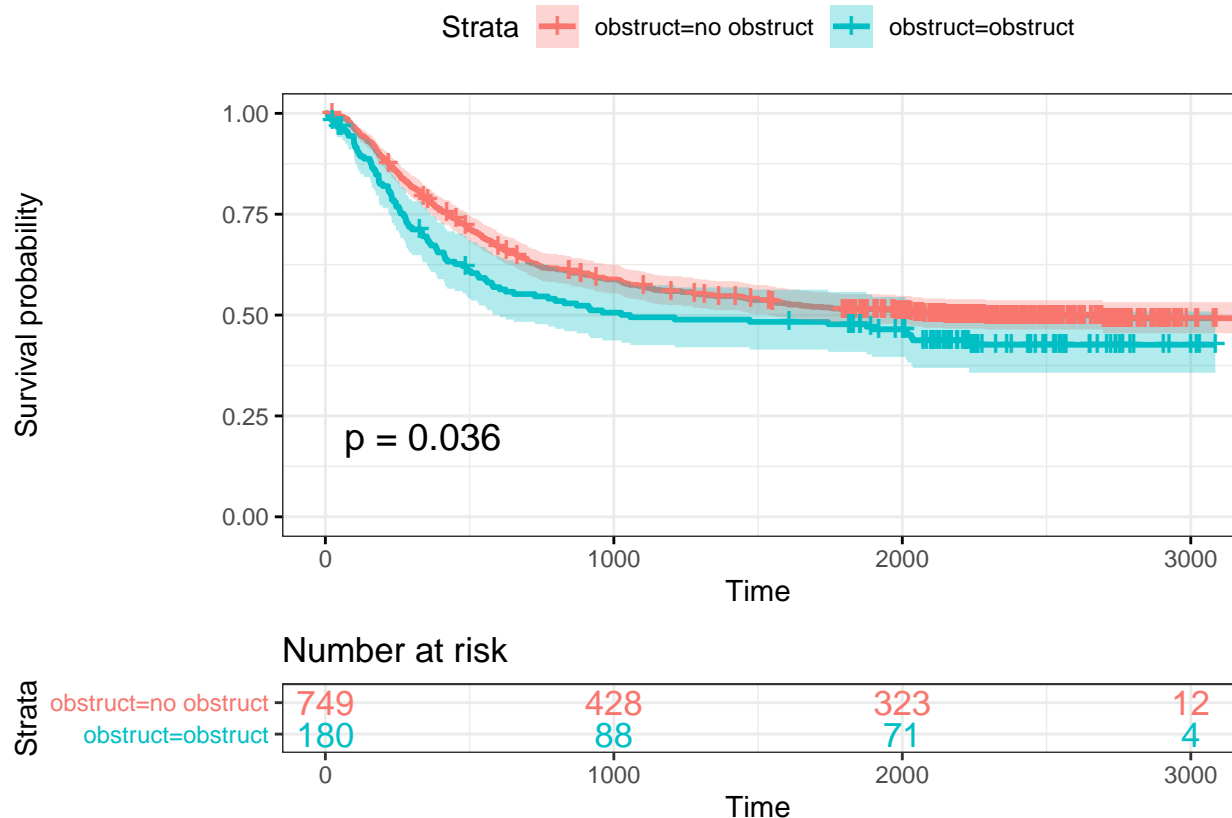
```
km_fit <- survfit(surv~obstruct, data=colon_subset_recurrence)
```

```
summary(km_fit, times=c(1, 10, 20, 30, 40, 50))
```

```
## Call: survfit(formula = surv ~ obstruct, data = colon_subset_recurrence)
##
##               obstruct=no obstruct
##  time  n.risk  n.event  survival  std.err  lower 95% CI  upper 95% CI
##    1      749       0    1.000  0.00000    1.000    1.000
##   10      748       1    0.999  0.00133    0.996    1.000
##   20      748       0    0.999  0.00133    0.996    1.000
##   30      746       1    0.997  0.00189    0.994    1.000
##   40      745       1    0.996  0.00231    0.991    1.000
##   50      742       3    0.992  0.00326    0.986    0.998
```

```
##
##          obstruct=obstruct
## time n.risk n.event survival std.err lower 95% CI upper 95% CI
##    1    180      0   1.000 0.00000      1.000      1.000
##   10    179      1   0.994 0.00554      0.984      1.000
##   20    178      2   0.983 0.00954      0.965      1.000
##   30    176      0   0.983 0.00954      0.965      1.000
##   40    174      3   0.967 0.01342      0.941      0.993
##   50    172      0   0.967 0.01342      0.941      0.993
```

```
ggsurvplot(km_fit, data = colon_subset_recurrence,
  pval = TRUE, conf.int = TRUE,
  risk.table = TRUE, ggtheme = theme_bw(),
  risk.table.col = "strata")
```



The p-value in the plot comes from the log-rank hypothesis test which allows us to compare a set of Kaplan-Meier estimators. The null hypothesis is that there is no significant difference between the Kaplan-Meier estimators. Since $p < 0.05$, we reject the null hypothesis.

We also do the log-rank hypothesis test using the `survdif` function.

```
p_value <- survdiff(surv~obstruct, data=colon_subset_recurrence)
print(p_value)
```

```
## Call:
## survdiff(formula = surv ~ obstruct, data = colon_subset_recurrence)
##
##          N Observed Expected (O-E)^2/E (O-E)^2/V
## obstruct=no obstruct 749      369      386.2      0.768      4.4
## obstruct=obstruct   180       99       81.8      3.628      4.4
```

```
##
## Chisq= 4.4 on 1 degrees of freedom, p= 0.04
survdifff(surv~adhere + obstruct,data=colon_subset_recurrence)

## Call:
## survdifff(formula = surv ~ adhere + obstruct, data = colon_subset_recurrence)
##
##
##              N Observed Expected (O-E)^2/E
## adhere=no adhere, obstruct=no obstruct 642      306   335.6      2.61
## adhere=no adhere, obstruct=obstruct   152       80    69.5      1.59
## adhere=adhere, obstruct=no obstruct   107       63    50.6      3.02
## adhere=adhere, obstruct=obstruct      28       19    12.3      3.68
##              (O-E)^2/V
## adhere=no adhere, obstruct=no obstruct      9.24
## adhere=no adhere, obstruct=obstruct      1.86
## adhere=adhere, obstruct=no obstruct      3.39
## adhere=adhere, obstruct=obstruct      3.79
##
## Chisq= 10.9 on 3 degrees of freedom, p= 0.01
```

Kalpan-Meier Estimator for the data divided into adhere and no adhere

colon_subset_recurrence can be divided two data sets by the adhere column. Those patients whose colons are obstructed by the tumour and those whose colons aren't. We can fit to each data partition to a Kalpan-Meier Estimator

$$\hat{S}_{\text{adhere}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{adhere}}} \left(1 - \frac{d_i}{n_i}\right)$$

$$\hat{S}_{\text{no adhere}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{no adhere}}} \left(1 - \frac{d_i}{n_i}\right)$$

to the entire data. To do this, we use the command below.

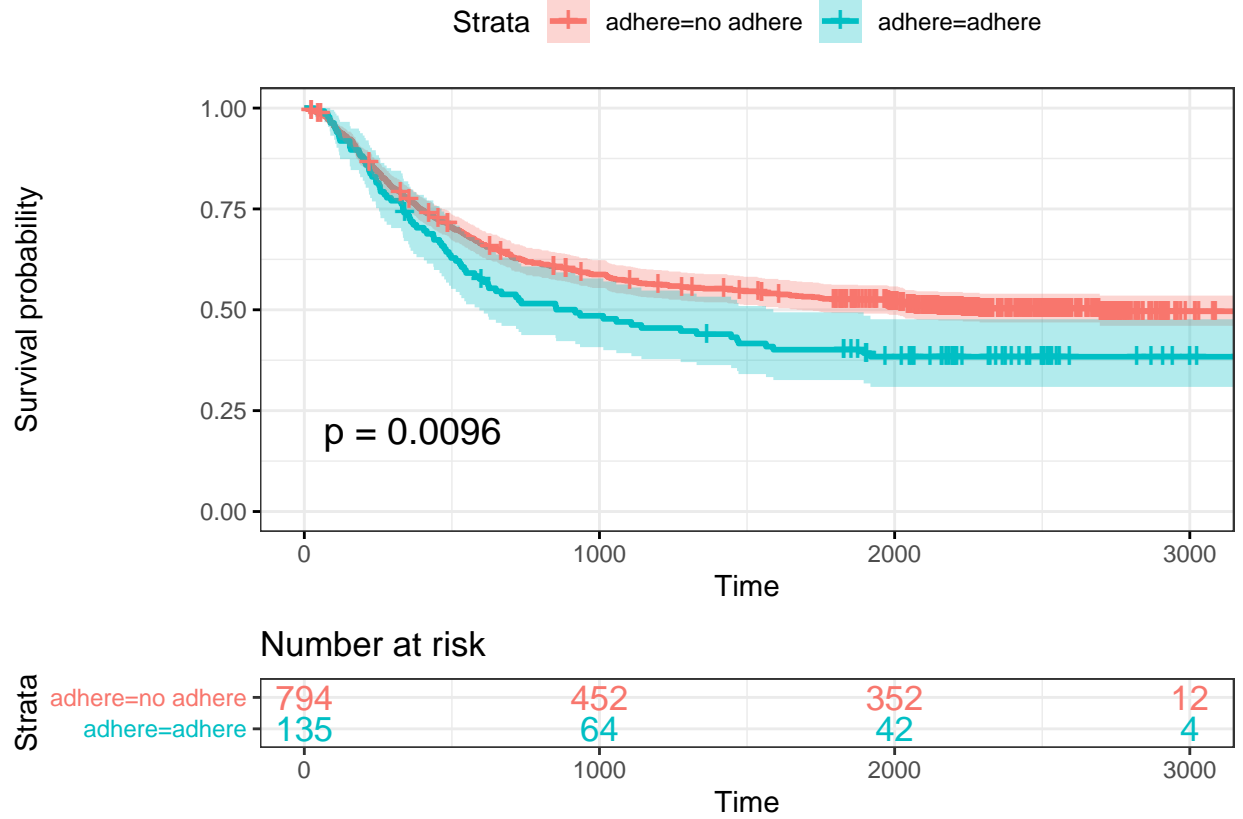
```
km_fit <- survfit(surv~adhere, data=colon_subset_recurrence)
```

```
summary(km_fit, times=c(1,10,20,30,40,50))
```

```
## Call: survfit(formula = surv ~ adhere, data = colon_subset_recurrence)
##
##              adhere=no adhere
## time n.risk n.event survival std.err lower 95% CI upper 95% CI
##    1    794      0   1.000 0.00000      1.000      1.000
##   10    792      2   0.997 0.00178      0.994      1.000
##   20    791      2   0.995 0.00251      0.990      1.000
##   30    787      1   0.994 0.00281      0.988      0.999
##   40    785      3   0.990 0.00355      0.983      0.997
##   50    780      3   0.986 0.00416      0.978      0.994
##
##              adhere=adhere
## time n.risk n.event survival std.err lower 95% CI upper 95% CI
##    1    135      0   1.000 0.00000      1.000      1
##   10    135      0   1.000 0.00000      1.000      1
##   20    135      0   1.000 0.00000      1.000      1
```


##	30	135	0	1.000	0.00000	1.000	1
##	40	134	1	0.993	0.00738	0.978	1
##	50	134	0	0.993	0.00738	0.978	1

```
ggsurvplot(km_fit, data = colon_subset_recurrence,
  pval = TRUE, conf.int = TRUE,
  risk.table = TRUE, ggtheme = theme_bw(),
  risk.table.col = "strata")
```



Kalpan-Meier Estimator for the data divided into (adhere, obstruct), (adhere, no obstruct), (no adhere, obstruct) and (no adhere, no obstruct)

colon_subset_recurrence can be divided in any amount by the explanatory variables Let's consider breaking up the data based on a patient's obstruction and adherence status. We can fit to each data partition to a Kalpan-Meier Estimator

$$\hat{S}_{\text{adhere,obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{adhere} \\ \text{obstruct}_i = \text{obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

$$\hat{S}_{\text{no adhere,obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{no adhere} \\ \text{obstruct}_i = \text{obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

$$\hat{S}_{\text{adhere,no obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{adhere} \\ \text{obstruct}_i = \text{no obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

$$\hat{S}_{\text{no adhere, no obstruct}}(t) = \prod_{\substack{\text{for } i: t_i \leq t \\ \text{adher}_i = \text{no adhere} \\ \text{obstruct}_i = \text{no obstruct}}} \left(1 - \frac{d_i}{n_i}\right)$$

to the entire data. To do this, we use the command below.

```
km_fit <- survfit(surv~adhere + obstruct, data=colon_subset_recurrence)
```

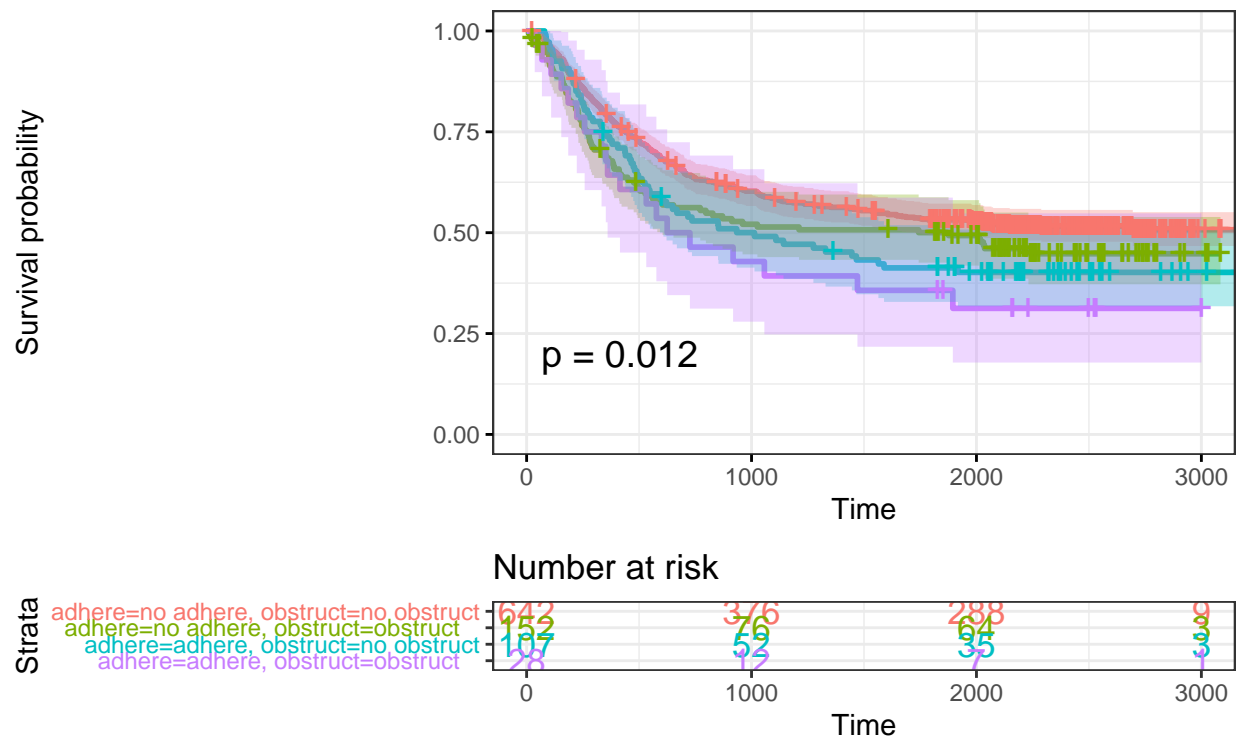
```
summary(km_fit, times=c(1,10,20,30,40,50))
```

```
## Call: survfit(formula = surv ~ adhere + obstruct, data = colon_subset_recurrence)
```

```
##
##               adhere=no adhere, obstruct=no obstruct
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##     1     642      0    1.000 0.00000      1.000      1.000
##    10     641      1    0.998 0.00156      0.995      1.000
##    20     641      0    0.998 0.00156      0.995      1.000
##    30     639      1    0.997 0.00220      0.993      1.000
##    40     638      1    0.995 0.00269      0.990      1.000
##    50     635      3    0.991 0.00380      0.983      0.998
##
##               adhere=no adhere, obstruct=obstruct
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##     1     152      0    1.000 0.00000      1.000      1.000
##    10     151      1    0.993 0.00656      0.981      1.000
##    20     150      2    0.980 0.01128      0.958      1.000
##    30     148      0    0.980 0.01128      0.958      1.000
##    40     147      2    0.967 0.01451      0.939      0.996
##    50     145      0    0.967 0.01451      0.939      0.996
##
##               adhere=adhere, obstruct=no obstruct
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##     1     107      0      1      0      1      1
##    10     107      0      1      0      1      1
##    20     107      0      1      0      1      1
##    30     107      0      1      0      1      1
##    40     107      0      1      0      1      1
##    50     107      0      1      0      1      1
##
##               adhere=adhere, obstruct=obstruct
##   time n.risk n.event survival std.err lower 95% CI upper 95% CI
##     1      28      0    1.000 0.0000      1.000      1
##    10      28      0    1.000 0.0000      1.000      1
##    20      28      0    1.000 0.0000      1.000      1
##    30      28      0    1.000 0.0000      1.000      1
##    40      27      1    0.964 0.0351      0.898      1
##    50      27      0    0.964 0.0351      0.898      1
```

```
ggsurvplot(km_fit, data = colon_subset_recurrence,
            pval = TRUE, conf.int = TRUE,
            risk.table = TRUE, ggtheme = theme_bw(),
            risk.table.col = "strata")
```

adhere=no adhere, obstruct=no obstruct + adhere=no adhere, obstruct=obstruct + adhere=adhere, obstruct=no obstruct + adhere=adhere, obstruct=obstruct



```
survdif(surv~adhere + obstruct,data=colon_subset_recurrence)
```

```
## Call:
## survdif(formula = surv ~ adhere + obstruct, data = colon_subset_recurrence)
##
##              N Observed Expected (O-E)^2/E
## adhere=no adhere, obstruct=no obstruct 642      306    335.6      2.61
## adhere=no adhere, obstruct=obstruct   152       80     69.5      1.59
## adhere=adhere, obstruct=no obstruct   107       63     50.6      3.02
## adhere=adhere, obstruct=obstruct       28       19     12.3      3.68
##              (O-E)^2/V
## adhere=no adhere, obstruct=no obstruct      9.24
## adhere=no adhere, obstruct=obstruct       1.86
## adhere=adhere, obstruct=no obstruct       3.39
## adhere=adhere, obstruct=obstruct       3.79
##
## Chisq= 10.9 on 3 degrees of freedom, p= 0.01
```

Cox Proportional Hazard

In the limit of large data, the Kaplan-Meier estimator converges to true survival function. However, the Kaplan-Meier has two disadvantages:

- it cannot effectively accomodate continuous data
- it is non-parametric – this means that given a data point, we cannot predict their life trajectory from data. This will be seen more clearly later in this section.

Rather than estimating survival function at each time interval, the *Cox Proportional Hazard* assumes that hazard function is an exponentiated linear function of explanatory variables. That is,

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_1 X_{1i} + \dots + \beta_n X_{ni}).$$

$\lambda_0(t)$ is called the baseline function. $\lambda(t) = \lambda_0(t)$ when $X_{1i} = X_{2i} = \dots = X_{ni} = 0$.

The Cox Proportional Hazard models the effects of the covariates on the baseline function. It assumes that the ratio of hazards are independent of time. The baseline function is generally unknown. However, the effects of the covariates can still be determined regardless of the baseline function. The β_i 's is calculated using *partial maximum likelihood*. Avoiding the estimation of $\lambda_0(t)$ prevents accumulation of errors in a unknown function.

Note that the Cox Proportional Hazard does not solve all the problems of the Kaplan-Meier estimator. Cox Proportional Hazard has one (or 1/2) disadvantage:

- it is semi-parametric. Given a data point, we can estimate the effect of a covariate on the baseline function. However, we cannot predict the life trajectory of data point unless we know $\lambda_0(t)$.

Cox Proportional Hazard for $X_1 = \text{surg}$

Given only one covariate, our Cox Proportional Hazard function takes the form

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_1 X_{1i}).$$

where

$$X_{1i} = \begin{cases} 1 & \text{if surgery time of ith data point is long} \\ 0 & \text{otherwise} \end{cases}.$$

Learning Cox Proportional Hazard model

We fit the Cox Proportional Hazard model accordingly.

```
cox <- coxph(surv ~ surg,
             data=colon_subset_recurrence)
```

```
summary(cox)
```

```
## Call:
## coxph(formula = surv ~ surg, data = colon_subset_recurrence)
##
##      n= 929, number of events= 468
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## surglong 0.2549    1.2903   0.1008 2.529   0.0114 *
## ---
```

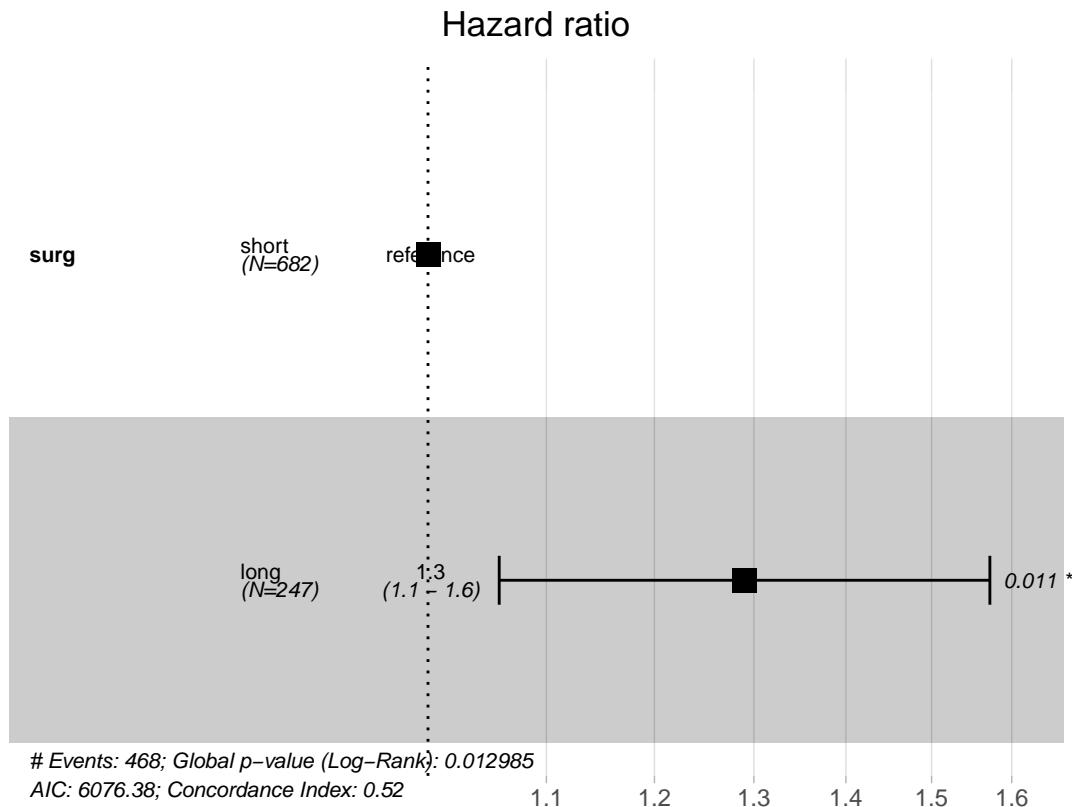
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      exp(coef) exp(-coef) lower .95 upper .95
## surglong      1.29      0.775    1.059    1.572
##
## Concordance= 0.523  (se = 0.01 )
## Likelihood ratio test= 6.17  on 1 df,   p=0.01
## Wald test            = 6.39  on 1 df,   p=0.01
## Score (logrank) test = 6.43  on 1 df,   p=0.01
```

```
coef(cox)
```

```
## surglong
## 0.2548703
```

```
ggforest(cox, data = colon_subset_recurrence)
```

```
## Warning: Removed 1 rows containing missing values (geom_errorbar).
```



Testing Proportionality Assumption

The Cox proportionality hazard model assumes that ratio of the hazards are constant over time. If ratio of the hazards are constant over time, then covariates and their effects must also be constant over time. If this assumption is violated, then one might get strange results (such as the crossing of Kaplan-Meier curves).

To test for proportionality hazard assumption, we use the `cox.zph` function. `cox.zph` takes a `coxph` model as input and returns a p-value to determine whether the proportionality hazard assumption was violated for each covariate. `cox.zph` tests the null hypothesis that there are no time dependent relationships in the covariates and their effects.

```
test.ph <- cox.zph(cox)
test.ph
```

```
##           rho chisq      p
## surglong 0.0503  1.18 0.277
```

Since the p value is greater than 0.05, we fail to reject the null hypothesis

Model Selection

```
anova(cox)
```

```
## Analysis of Deviance Table
## Cox model: response is surv
## Terms added sequentially (first to last)
##
##      loglik  Chisq Df Pr(>|Chi|)
## NULL -3040.3
## surg -3037.2 6.1712  1    0.01299 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Cox Proportional Hazard for $X_1 = \text{surg}$, $X_2 = \text{adher}$

Given only two covariate, our Cox Proportional Hazard function takes the form

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_1 X_{1i} + \beta_2 X_{2i}).$$

where

$$X_{1i} = \begin{cases} 1 & \text{if surgery time of } i \text{ th data point is long} \\ 0 & \text{otherwise} \end{cases},$$

$$X_{2i} = \begin{cases} 1 & \text{if the } i \text{ th data point has adherence to other organs} \\ 0 & \text{otherwise} \end{cases}.$$

Learning Cox Proportional Hazard model

We fit the Cox Proportional Hazard model accordingly.

```
cox <- coxph(surv ~ surg + adhere,
             data=colon_subset_recurrence)
```

```
summary(cox)
```

```
## Call:
## coxph(formula = surv ~ surg + adhere, data = colon_subset_recurrence)
##
##      n= 929, number of events= 468
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## surglong      0.2481   1.2816   0.1008 2.460  0.0139 *
## adhereadhere 0.3053   1.3570   0.1217 2.508  0.0121 *
## ---
```

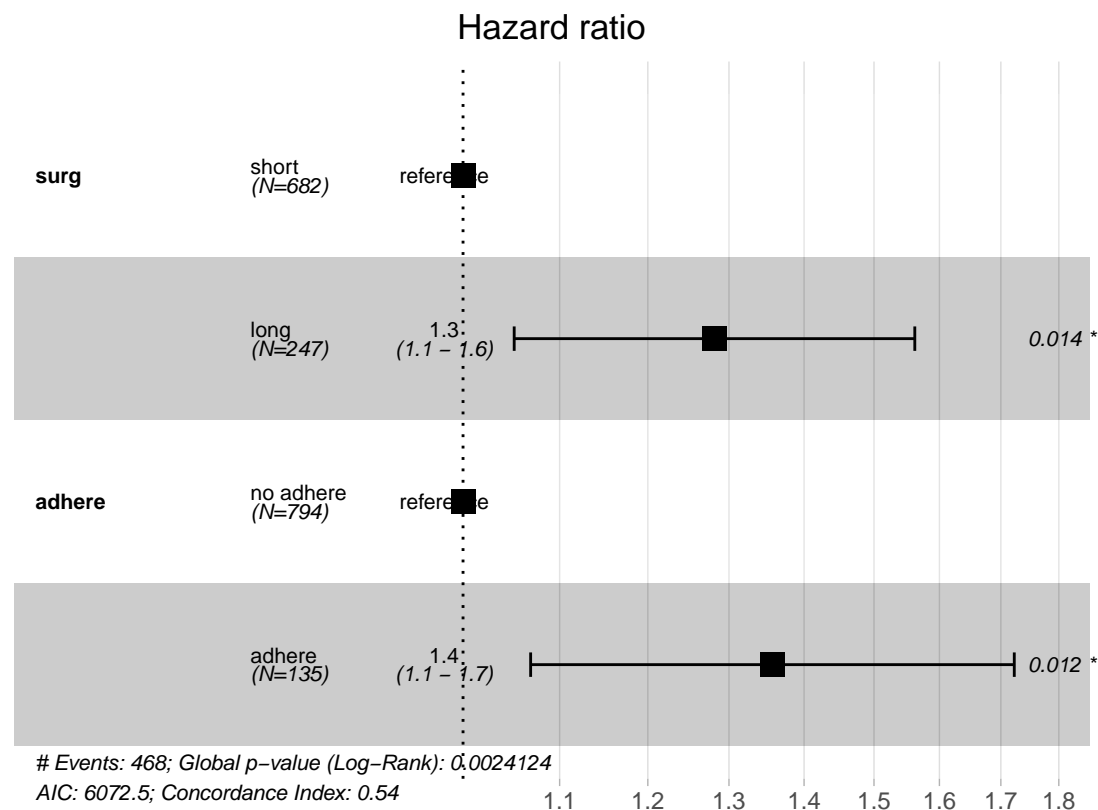
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## surglong      1.282     0.7803     1.052     1.562
## adhereadhere  1.357     0.7369     1.069     1.723
##
## Concordance= 0.538 (se = 0.012 )
## Likelihood ratio test= 12.05 on 2 df,  p=0.002
## Wald test            = 12.72 on 2 df,  p=0.002
## Score (logrank) test = 12.81 on 2 df,  p=0.002
```

```
coef(cox)
```

```
##      surglong adhereadhere
##      0.2480832    0.3052677
```

```
ggforest(cox, data = colon_subset_recurrence)
```

```
## Warning: Removed 2 rows containing missing values (geom_errorbar).
```



Testing Proportionality Assumption

```
test.ph <- cox.zph(cox)
test.ph
```

```
##           rho chisq    p
## surglong   0.0486 1.110 0.292
## adhereadhere 0.0403 0.764 0.382
## GLOBAL      NA  1.931 0.381
```

Since the p value is greater than 0.05, we fail to reject the null hypothesis

Model Selection

```
anova(cox)

## Analysis of Deviance Table
## Cox model: response is surv
## Terms added sequentially (first to last)
##
##          loglik  Chisq Df Pr(>|Chi|)
## NULL      -3040.3
## surg      -3037.2 6.1712  1    0.01299 *
## adhere    -3034.2 5.8831  1    0.01529 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Cox Proportional Hazard for $X_1 = \text{surg}$, $X_2 = \text{adher}$, $X_3 = \text{nodes}$

Given only three covariate, our Cox Proportional Hazard function takes the form

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i}).$$

where

$$X_{1i} = \begin{cases} 1 & \text{if surgery time of } i \text{ th data point is long} \\ 0 & \text{otherwise} \end{cases},$$

$$X_{2i} = \begin{cases} 1 & \text{if the } i \text{ th data point has adherence to other organs} \\ 0 & \text{otherwise} \end{cases}$$

and X_{3i} is number of nodes of the i th data point.

Learning Cox Proportional Hazard model

We fit the Cox Proportional Hazard model accordingly.

```
cox <- coxph(surv ~ surg + adhere + nodes,
             data=colon_subset_recurrence)

summary(cox)

## Call:
## coxph(formula = surv ~ surg + adhere + nodes, data = colon_subset_recurrence)
##
##      n= 911, number of events= 456
##      (18 observations deleted due to missingness)
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## surglong      0.27318   1.31414  0.10273  2.659  0.00783 **
## adhereadhere  0.30821   1.36099  0.12317  2.502  0.01234 *
## nodes         0.08562   1.08939  0.00888  9.642 < 2e-16 ***
## ---
```



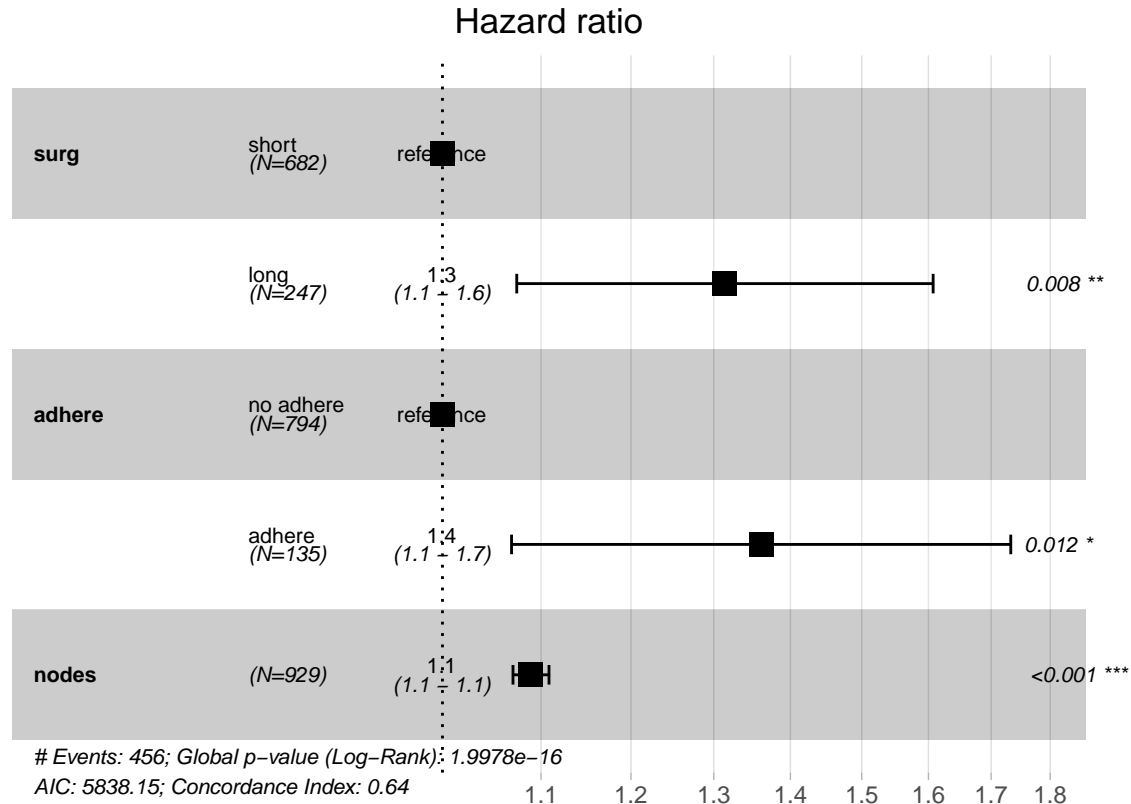
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## surglong      1.314     0.7610     1.074     1.607
## adhereadhere  1.361     0.7348     1.069     1.733
## nodes         1.089     0.9179     1.071     1.109
##
## Concordance= 0.636 (se = 0.013 )
## Likelihood ratio test= 76.21 on 3 df,  p=<2e-16
## Wald test            = 102.7 on 3 df,  p=<2e-16
## Score (logrank) test = 105.5 on 3 df,  p=<2e-16
```

```
coef(cox)
```

```
##      surglong adhereadhere      nodes
## 0.27318255 0.30821401 0.08561674
```

```
ggforest(cox, data = colon_subset_recurrence)
```

```
## Warning: Removed 2 rows containing missing values (geom_errorbar).
```



Testing Proportionality Assumption

```
test.ph <- cox.zph(cox)
test.ph
```

```
##      rho chisq    p
## surglong 0.0590 1.593 0.207
## adhereadhere 0.0307 0.430 0.512
```

```
## nodes          -0.0460 0.601 0.438
## GLOBAL          NA 2.815 0.421
```

Since the p value is greater than 0.05, we fail to reject the null hypothesis

Model Selection

```
anova(cox)

## Analysis of Deviance Table
## Cox model: response is surv
## Terms added sequentially (first to last)
##
##          loglik   Chisq Df Pr(>|Chi|)
## NULL      -2954.2
## surg      -2951.6  5.1725  1    0.02295 *
## adhere    -2948.6  6.0414  1    0.01397 *
## nodes     -2916.1 64.9922  1   7.519e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Estimating Survival Curve

It is possible to estimate the survival curve for the Cox Proportional Model as long as we have some estimate for $\lambda_0(t)$. One way to estimate $\lambda_0(t)$ from data is to use formula:

$$\lambda_0(t_i) \approx \frac{d_i}{\sum_{s \in R_i} \exp(\beta_1 X_{1s} + \dots + \beta_n X_{ns})}$$

where d_i is the number of deaths in at time t_i , R_i is set of persons alive after t_i and X_{ij} is the i th explanatory variable of the j th person.

Now let's create some data point. This data point will have the **surg** set to **short**, **adhere** set to **no adhere**, **nodes** set to 5 and **extent** set to **serosa**.

```
subject_one <- data.frame(surg = factor('short'),
                          adhere = factor('adhere'), nodes = 5)
```

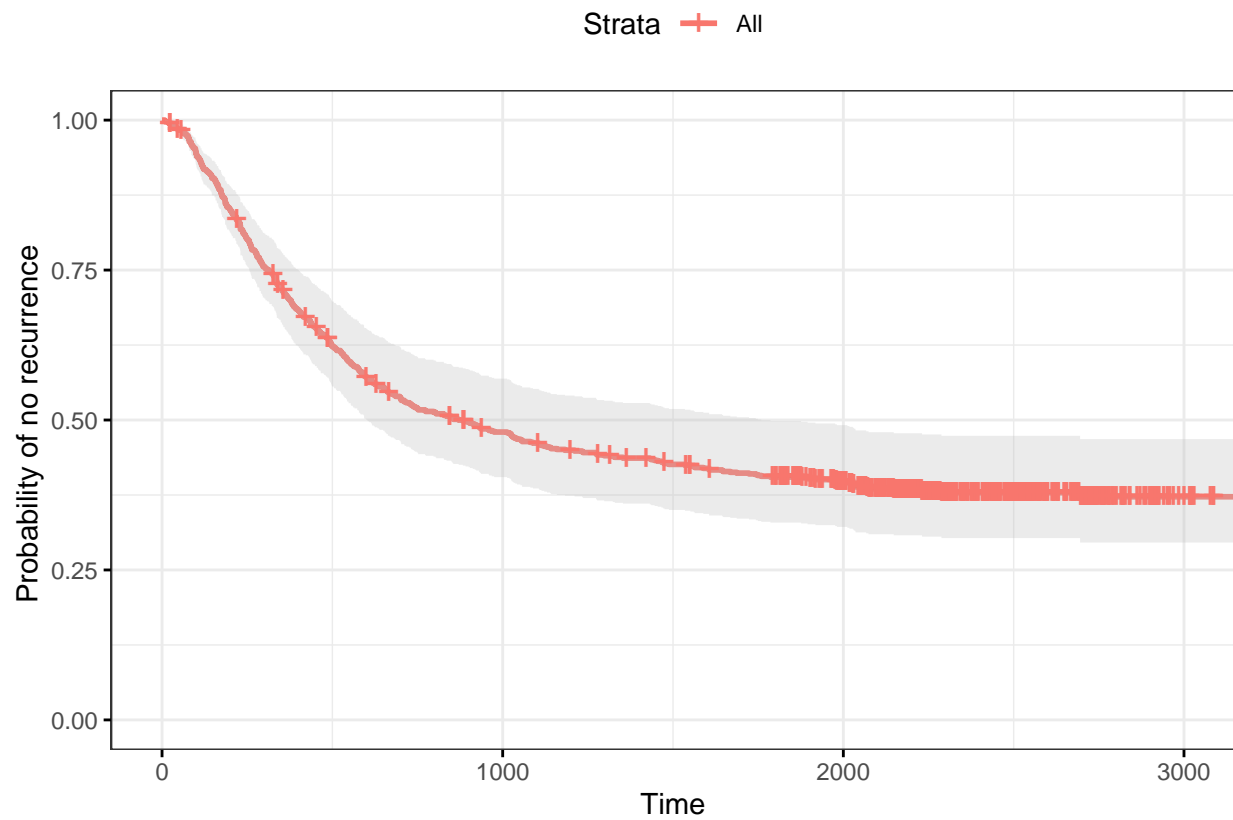
Using the **survfit** function, we can generate an object which will be used for plotting. **survfit** takes as argument:

- first argument: cox proportional hazard model fit with **coxph**
- second argument: the data point in question. It must have the same explanatory variables as the model in the first argument
- data: the data set used to fit the **coxph** object.

```
prediction_one <- survfit(cox, subject_one,
                         data = colon_subset_recurrence)
```

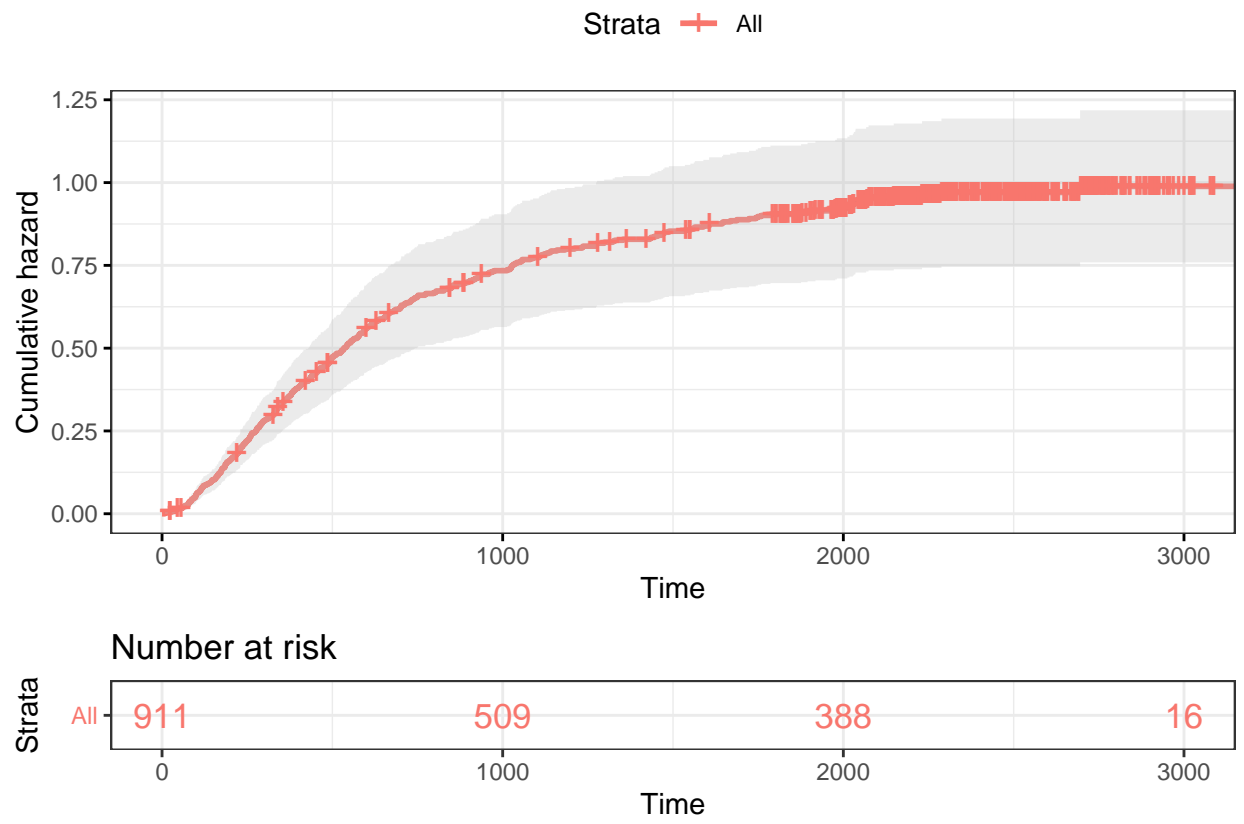
We then use the **ggsurvplot** function to plot the estimate of the survival curve from **survfit** fit object.

```
ggsurvplot(prediction_one, ylab = "Probability of no recurrence ",
            conf.int = TRUE,
            ggtheme = theme_bw())
```



We can also use the `ggsurvplot` function to plot the estimate of the cumulative hazard curve from `survfit` fit object

```
ggsurvplot(prediction_one, fun="cumhaz",  
            conf.int = TRUE, risk.table = TRUE,  
            ggtheme = theme_bw(),  
            risk.table.col = "strata")
```



Accelerated failure time models

Accelerated failure time model assume that the log time for an event to occur is a function of the covariates of the data. That is,

$$\log T_i = \beta_1 X_{1i} + \cdots + \beta_n X_{ni} + \varepsilon_i$$

where ε is a random error term that follows a distribution.

This is called an **accelerated** failure model since covariates can scale the base time distribution, T_0 , by their effects.

$$T_i = T_0 \exp(\beta_1 X_{1i} + \cdots + \beta_n X_{ni})$$

where $T_0 = \exp(\varepsilon_i)$.

There is difference between proportional hazard models (PH) and accelerated failure time models (AFT). The effect of the covariates in PH models act multiplicately on the base hazard. However, in AFT models, these effects act multiplicately on the base time. Despite this difference, it is possible that AFT models are also PH models.

We use the function, **survreg**, to fit accelerated failure time models. The argument, **dist**, specifies the distribution which implies the form of $\lambda_0(t)$. We will be considering:

- exponential models, **dist**="exponential"
- weibull models, **dist**="weibull"
- lognormal models, **dist**="lognormal"

These are fully parameteric model and are thus a suitable alternative Kaplan-Meier estimators and Cox Proportional Hazard models. However, AFT assume the function form T_0 and, thus, the baseline hazard and the baseline survival functions. Incorrect assumptions introduce errors in our modeling.

Exponential models

Exponential accelerated failure time models are also proportional hazard models. Exponential accelerated failure time models assume that T_0 follows a exponential distribution with parameter λ .

From our definitions of terms and with some probability theory (not covered), the hazard and survival function of an exponential AFT models are

$$\lambda_i(t) = \lambda \exp(\beta_1 X_{1i} + \cdots + \beta_n X_{ni})$$

$$\text{and } S_i(t) = \exp(\beta_1 X_{1i} + \cdots + \beta_n X_{ni}) S_0(t), \quad S_0(t) = \exp(-\lambda t).$$

As proportional hazard model, exponential accelerated failure time models assumes that the baseline hazard is constant, $\lambda_0(t) = \lambda$.

Learning Exponential models

survreg learns the parameter value, λ , and the regression coefficients. As an example, we will be consider the model: **surv ~ 1 + surg + adhere + nodes**.

```

survregExp <- survreg(surv ~ 1 + surg + adhere + nodes,
                      dist="exponential",data=colon_subset_recurrence)
summary(survregExp)

##
## Call:
## survreg(formula = surv ~ 1 + surg + adhere + nodes, data = colon_subset_recurrence,
##         dist = "exponential")
##               Value Std. Error      z      p
## (Intercept)   8.45944    0.07338 115.29 <2e-16
## surglong     -0.32521    0.10289  -3.16 0.0016
## adhereadhere -0.34689    0.12321  -2.82 0.0049
## nodes        -0.10446    0.00878 -11.90 <2e-16
##
## Scale fixed at 1
##
## Exponential distribution
## Loglik(model)= -4024.3  Loglik(intercept only)= -4078.5
##  Chisq= 108.33 on 3 degrees of freedom, p= 2.5e-23
## Number of Newton-Raphson Iterations: 5
## n=911 (18 observations deleted due to missingness)

```

Therefore, $\lambda = \exp(8.45944)$.

Estimating Survival Curve

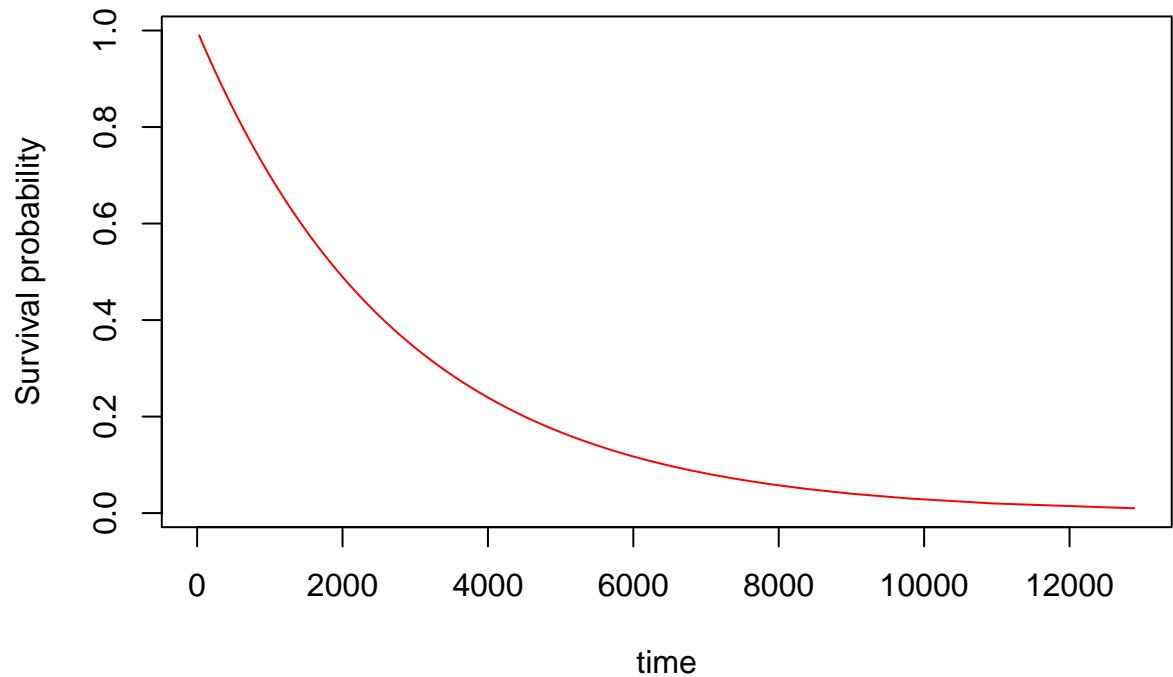
```

subject_two = list(surg = factor('short'), adhere = factor('no adhere'), nodes = 5)

plot(predict(survregExp, newdata=subject_two,
             type="quantile",p=seq(.01,.99,by=.01)),
     seq(.99,.01,by=-.01), col="red",type='l',xlab='time',
     ylab='Survival probability',main='Exponential AFT Model')

```

Exponential AFT Model



##

Weibull models

Weibull accelerated failure time models are also proportional hazard models. Weibull accelerated failure time models assume that T_0 follows a Weibull distribution with parameters, λ and γ .

From our definitions of terms and with some probability theory (not covered), the hazard and survival function of a Weibull AFT models are

$$\lambda_i(t) = \lambda \gamma t^{\gamma-1} \exp(\beta_1 X_{1i} + \dots + \beta_n X_{ni}).$$

$$\text{and } S_i(t) = \exp(\beta_1 X_{1i} + \dots + \beta_n X_{ni}) S_0(t), \quad S_0(t) = \exp(-(\lambda t)^\gamma).$$

As proportional hazard model, Weibull accelerated failure time models assumes that the baseline hazard is constant, $\lambda \gamma t^{\gamma-1}$. One can see that exponential accelerated failure time models are a special case of Weibull accelerated failure time models with $\gamma = 1$.

Learning Weibull models

`survreg` learns the parameter value, λ and γ , and the regression coefficients.

As an example, we will consider the model: `surv ~ 1 + surg + adhere + nodes` for all the accelerated time models.

```
survregWeibull = survreg(surv ~ 1 + surg + adhere + nodes,
                        dist="weibull", data=colon_subset_recurrence)
summary(survregWeibull)
```

##

Call:

`survreg(formula = surv ~ 1 + surg + adhere + nodes, data = colon_subset_recurrence,`

```
##      dist = "weibull")
##              Value Std. Error      z      p
## (Intercept)   8.7993    0.1155  76.21 <2e-16
## surglong      -0.4156    0.1454  -2.86 0.0042
## adhereadhere -0.4629    0.1743  -2.66 0.0079
## nodes         -0.1315    0.0131 -10.04 <2e-16
## Log(scale)    0.3432    0.0414   8.30 <2e-16
##
## Scale= 1.41
##
## Weibull distribution
## Loglik(model)= -3984.6   Loglik(intercept only)= -4028.2
##  Chisq= 87.17 on 3 degrees of freedom, p= 8.9e-19
## Number of Newton-Raphson Iterations: 5
## n=911 (18 observations deleted due to missingness)
```

Therefore,

$$\gamma = \exp(0.3432)$$

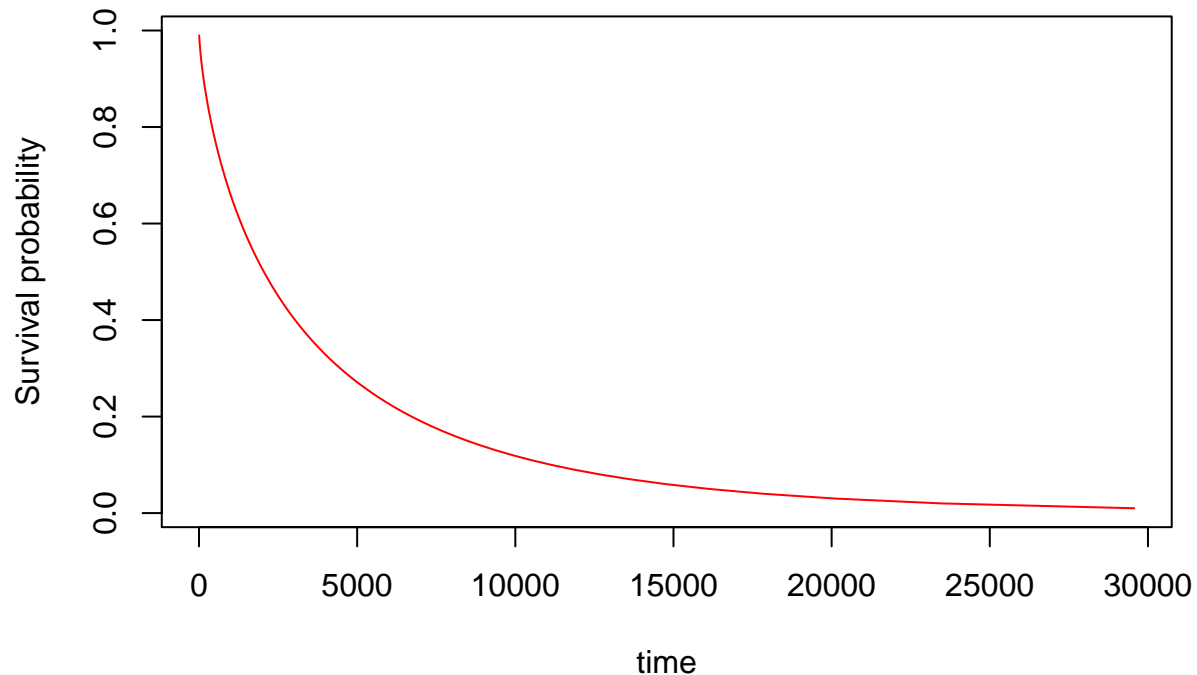
$$\lambda = \exp(8.7993 \times \gamma)$$

Estimating Survival Curve

```
subject_two = list(surg = factor('short'),
                   adhere = factor('no adhere'),
                   nodes = 5)

plot(predict(survregWeibull, newdata=subject_two,
            type="quantile", p=seq(.01,.99,by=.01)),
     seq(.99,.01,by=-.01), col="red", type='l', xlab='time',
     ylab='Survival probability', main='Weibull AFT Model')
```


Weibull AFT Model



Log-normal models

Log-normal accelerated failure time models assume that T_0 follows a log normal distribution with a scale parameter. Log-normal accelerated failure time models are not proportional hazard models.

The hazard and survival function of a log-normal AFT models are a bit complicated so they will not be shown here.

Learning Log-normal models

`survreg` learns the parameter value, λ and γ , and the regression coefficients.

As an example, we will consider the model: `surv ~ 1 + surg + adhere + nodes` for all the accelerated time models.

```
survregLogNormal = survreg(surv ~ 1 + surg + adhere + nodes,  
                           dist="lognormal", data=colon_subset_recurrence)  
summary(survregLogNormal)
```

```
##  
## Call:  
## survreg(formula = surv ~ 1 + surg + adhere + nodes, data = colon_subset_recurrence,  
##       dist = "lognormal")  
##           Value Std. Error      z      p  
## (Intercept)  8.3066      0.1223 67.93 <2e-16  
## surglong    -0.3521      0.1525 -2.31  0.021  
## adhereadhere -0.4685      0.1876 -2.50  0.013  
## nodes       -0.1613      0.0182 -8.86 <2e-16  
## Log(scale)   0.6127      0.0371 16.49 <2e-16  
##
```

```
## Scale= 1.85
##
## Log Normal distribution
## Loglik(model)= -3940.5   Loglik(intercept only)= -3983.6
##  Chisq= 86.09 on 3 degrees of freedom, p= 1.5e-18
## Number of Newton-Raphson Iterations: 3
## n=911 (18 observations deleted due to missingness)
```

Estimating Survival Curve

```
subject_two = list(surg = factor('short'),
                   adhere = factor('no adhere'),
                   nodes = 5)

plot(predict(survregLogNormal, newdata=subject_two,
            type="quantile",p=seq(.01,.99,by=.01)),
     seq(.99,.01,by=-.01), col="red",type='l',xlab='time',
     ylab='Survival probability',main='Log Normal AFT Model')
```

Log Normal AFT Model

