

432 Class 20 Slides

github.com/THOMASELOVE/2020-432

2020-04-07

Preliminaries

```
library(here); library(janitor); library(magrittr)
library(rms)
library(survival); library(survminer)
library(tidyverse)

theme_set(theme_bw())

survex <- read_csv(here("data/survex.csv")) %>%
  type.convert()
```

Working with Time-to-Event Data

- The Survival Function, $S(t)$
 - Kaplan-Meier Estimation of the Survival Function
 - Creating Survival Objects in R
 - Drawing Survival Curves
- Testing the difference between Survival Curves
- The Hazard Function and its Estimation
- Getting Started with Cox Proportional Hazards Regression

Working with Time to Event Data

In many medical studies, the main outcome variable is the time to the occurrence of a particular event.

- In a randomized controlled trial of cancer, for instance, surgery, radiation, and chemotherapy might be compared with respect to time from randomization and the start of therapy until death.
 - In this case, the event of interest is the death of a patient, but in other situations it might be remission from a disease, relief from symptoms or the recurrence of a particular condition.
 - Such observations are generally referred to by the generic term survival data even when the endpoint or event being considered is not death but something else.

What Do We Study in a Time-to-Event Study?

Survival analysis is concerned with prospective studies. We start with a cohort of patients and follow them forwards in time to determine some clinical outcome.

- Follow-up continues until either some event of interest occurs, the study ends, or further observation becomes impossible.

The outcomes in a survival analysis consist of the patient's **fate** and **length of follow-up** at the end of the study.

- For some patients, the outcome of interest may not occur during follow-up.
- For such patients, whose follow-up time is *censored*, we know only that this event did not occur while the patient was being followed. We do not know whether or not it will occur at some later time.

Problems with Time to Event Data

The primary problems are *non-normality* and *censoring*...

- 1 Survival data are not symmetrically distributed. They will often appear positively skewed, with a few people surviving a very long time compared with the majority; so assuming a normal distribution will not be reasonable.
- 2 At the completion of the study, some patients may not have reached the endpoint of interest (death, relapse, etc.). Consequently, the exact survival times are not known.
 - All that is known is that the survival times are greater than the amount of time the individual has been in the study.
 - The survival times of these individuals are said to be **censored** (precisely, they are right-censored).

Next, we'll define some special functions to build models that address these concerns.

The Survival Function, $S(t)$

The **survival function**, $S(t)$ (sometimes called the survivor function) is the probability that the survival time, T , is greater than or equal to a particular time, t .

- $S(t)$ = proportion of people surviving to time t or beyond

If there's no censoring, the survival function is easy to estimate

When there is no censoring, this function is easily estimated as ...

$$\hat{S}(t) = \frac{\# \text{ of subjects with survival times } \geq t}{n}$$

but this won't work if there is censoring.

Understanding the Kaplan-Meier Estimator

The survival function $S(t)$ is the probability of surviving until at least time t . It is essentially estimated by the number of patients alive at time t divided by the total number of study subjects remaining at that time.

The Kaplan-Meier estimator first orders the (unique) survival times from smallest to largest, then estimates the survival function at each unique survival time.

- The survival function at the second death time, $t_{(2)}$ is equal to the estimated probability of not dying at time $t_{(2)}$ conditional on the individual being still at risk at time $t_{(2)}$.

The Kaplan-Meier Estimator

- 1 Order the survival times from smallest to largest, where $t_{(j)}$ is the j th largest unique survival time, so we have...

$$t_{(1)} \leq t_{(2)} \leq t_{(3)} \leq \dots t_{(n)}$$

- 2 The Kaplan-Meier estimate of the survival function is

$$\hat{S}(t) = \prod_{j: t_{(j)} \leq t} \left(1 - \frac{d_j}{r_j}\right)$$

where r_j is the number of people at risk just before $t_{(j)}$, including those censored at time $t_{(j)}$, and d_j is the number of people who experience the event at time $t_{(j)}$.

Creating a Survival Object in R

The `Surv` function, part of the `survival` package in R, will create a **survival object** from two arguments:

- ❶ `time` = follow-up time
- ❷ `event` = a status indicator, where
 - `event = 1` or `TRUE` means the event was observed (for instance, the patient died)
 - `event = 0` or `FALSE` means the follow-up time was censored

The survex data frame

The `survex.csv` file on the course website is essentially the same as a file simulated by Frank Harrell and his team¹ to introduce some of the key results from the `cph` function, which is part of the `rms` package in R.

The `survex` data includes 1,000 subjects. . .

- `id` = patient ID (1-1000)
- `age` = patient's age at study entry, years
- `sex` = patient's sex (Male or Female)
- `study.yrs` = patient's years of observed time in study until death or censoring
- `death` = 1 if patient died, 0 if censored.

¹see the `rms` package documentation

A first example: Looking at just 100 observations

```
set.seed(4322020)
ex100 <- sample_n(survex, 100, replace = F)
ex100 %>% select(id, study.yrs, death) %>% summary()
```

id	study.yrs	death
Min. : 23.0	Min. : 0.175	Min. :0.00
1st Qu.:258.2	1st Qu.: 2.122	1st Qu.:0.00
Median :468.0	Median : 4.864	Median :0.00
Mean :479.1	Mean : 6.007	Mean :0.17
3rd Qu.:710.0	3rd Qu.: 9.759	3rd Qu.:0.00
Max. :938.0	Max. :14.817	Max. :1.00

For a moment, let's focus on developing a survival object in this setting.

Relationship between death and study.yrs?

- study.yrs here is follow-up time, in years
- death = 1 if subject had the event (death), 0 if not.

```
ex100 %$% mosaic::favstats(study.yrs ~ death)
```

	death	min	Q1	median	Q3	max	mean	sd
1	0	0.175	2.4775	5.268	10.233	14.817	6.373952	4.464091
2	1	0.641	1.8460	2.641	4.815	13.746	4.213882	3.780889
	n missing							
1	83	0						
2	17	0						

Building a Survival Object

```
surv_100 = ex100 %>% Surv(time = study.yrs, event = death)  
  
head(surv_100, 3)
```

```
[1] 3.047  9.454+ 4.023+
```

- Subject 1 survived 3.047 years and then died.
- Subject 2 survived 9.454 years before being censored.
- Subject 3 survived 4.023 years before being censored.

Remember that 17 of these 100 subjects died, the rest were censored at the latest time where they were seen for follow-up.

On dealing with time-to-event data

You have these three subjects.

- 1 Alice died in the hospital after staying for 20 days.
- 2 Betty died at home on the 20th day after study enrollment, after staying in the hospital for the first ten days.
- 3 Carol left the hospital after 20 days, but was then lost to follow up.

Suppose you plan a time-to-event analysis.

- How should you code “time” and “event” to produce a “time-to-event” object you can model if . . .
 - **death** is your primary outcome
 - **length of hospital stay** is your primary outcome?

Building a Kaplan-Meier Estimate

Remember that `surv_100` is the survival object we created.

```
km_100 <- survfit(surv_100 ~ 1)

print(km_100, print.rmean = TRUE)
```

Call: `survfit(formula = surv_100 ~ 1)`

n	events	*rmean	*se(rmean)	median
100.000	17.000	12.155	0.567	NA
0.95LCL	0.95UCL			
13.746	NA			
* restricted mean with upper limit = 14.8				

- 17 events (deaths) occurred in 100 subjects.
- Restricted mean survival time is 12.16 years (upper limit 14.8?)
- Median survival time is NA (why?) but has a lower bound for 95% CI.

Summary of the Kaplan-Meier Estimate

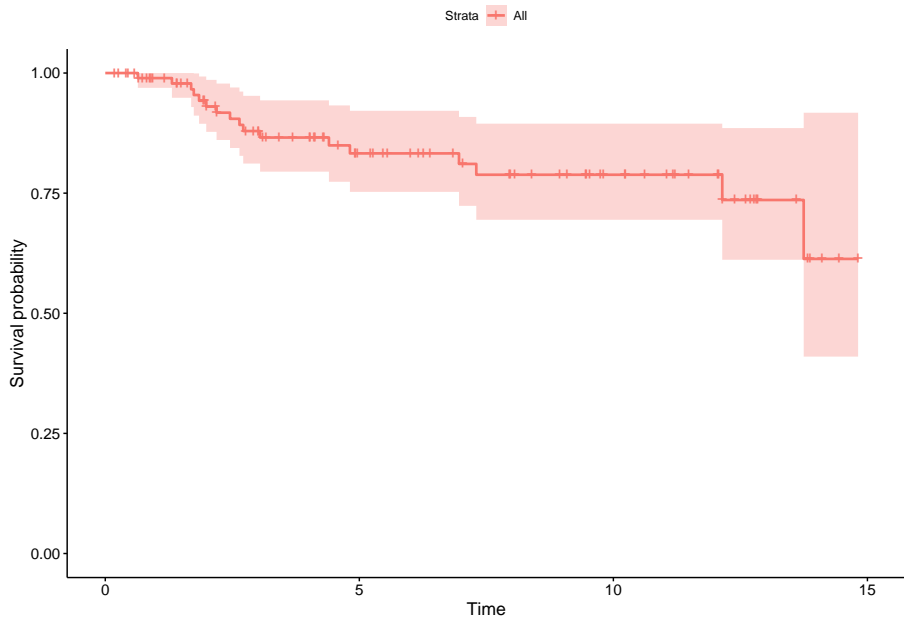
- Up to 0.641 years, no one died, but five people were censored (so 95 were at risk at that time). (Estimated survival probability = 0.989)
- By the time of the next death at 1.312 years, only 87 people were still at risk. (Estimated $\Pr(\text{survival})$ now 0.978)

```
summary(km_100)
```

```
Call: survfit(formula = surv_100 ~ 1)
```

time	n.risk	n.event	survival	std.err	lower 95% CI
0.641	95	1	0.989	0.0105	0.969
1.312	87	1	0.978	0.0153	0.949
1.690	82	1	0.966	0.0192	0.929
1.742	81	1	0.954	0.0224	0.911
1.846	80	1	0.942	0.0251	0.894
1.987	77	1	0.930	0.0276	0.878
2.190	74	1	0.918	0.0299	0.861
2.455	72	1	0.905	0.0321	0.844

Kaplan-Meier Plot, via survminer



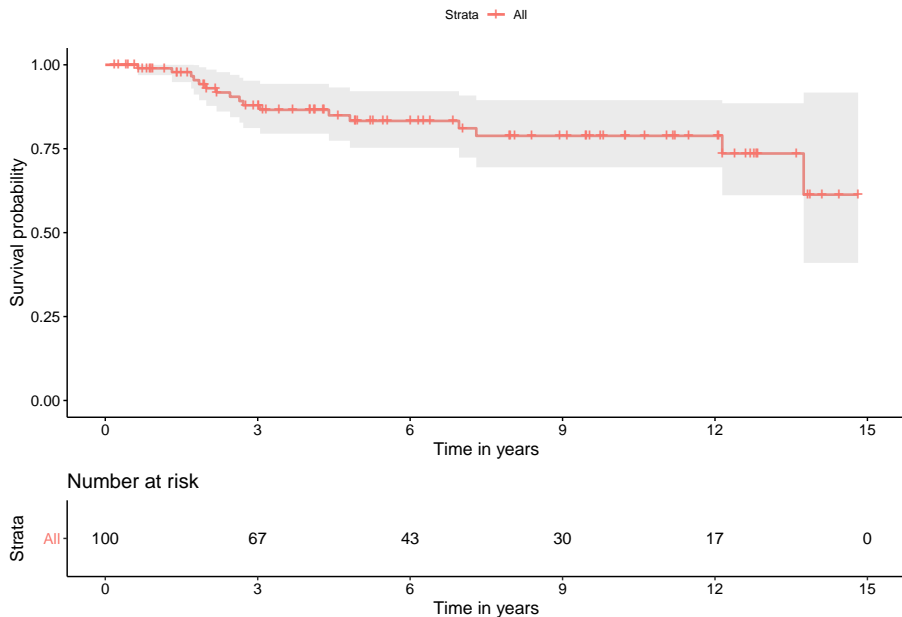
Kaplan-Meier Plot, via survminer (code)

```
ggsurvplot(km_100, data = ex100)
```

- The solid line indicates survival probability at each time point (in years.)
- The crosses indicate time points where censoring has occurred.
- The steps down indicate events (deaths.)
- The shading indicates (by default, 95%) pointwise confidence intervals.

For simultaneous confidence bands, visit the OpenIntro Statistics *Survival Analysis in R* materials, written by David Diez, as posted on our web site.

Adding a Number at Risk Table



Adding a Number at Risk Table (code)

```
ggsurvplot(km_100, data = ex100,  
  conf.int = TRUE,           # Add confidence interval  
  risk.table = TRUE,         # Add risk table  
  xlab = "Time in years",    # Adjust X axis label  
  break.time.by = 3         # X ticks every 3 years  
)
```

Comparing Survival, by Sex

Suppose we want to compare the survival functions for subjects classified by their sex.

- So, for instance, in our sample, 8 of 32 females and 9 of 68 males had the event (died).

```
ex100 %>% tabyl(death, sex) %>% adorn_totals()
```

death	Female	Male
0	24	59
1	8	9
Total	32	68

Summarizing the Survival Function Estimate, by Sex

```
km_100_sex <- survfit(surv_100 ~ ex100$sex)
```

```
print(km_100_sex, print.rmean = TRUE)
```

Call: survfit(formula = surv_100 ~ ex100\$sex)

	n	events	*rmean	*se(rmean)	median	0.95LCL
ex100\$sex=Female	32	8	9.45	1.163	NA	7.3
ex100\$sex=Male	68	9	12.05	0.479	NA	13.7

0.95UCL

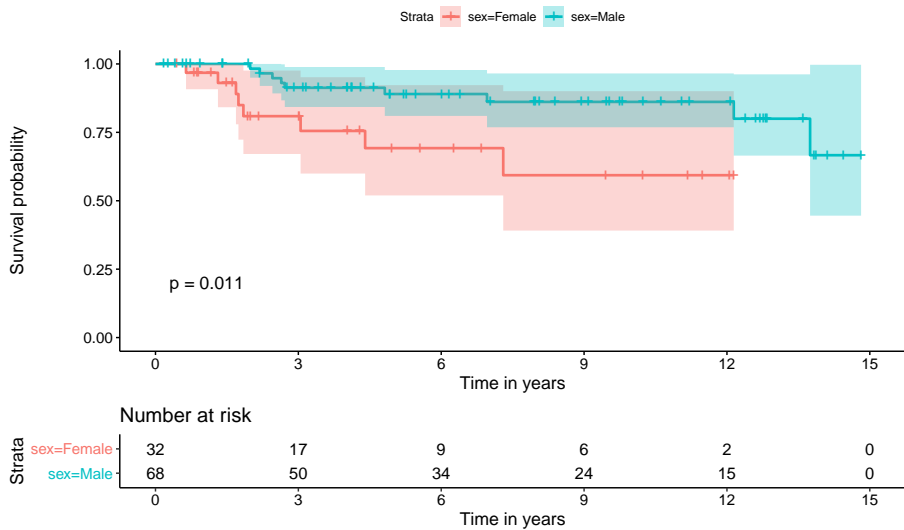
ex100\$sex=Female NA

ex100\$sex=Male NA

* restricted mean with upper limit = 13.5

- Among females, 8 of 32 subjects died, and the estimated restricted mean survival is 9.45 years.
- Among males, 9 of 68 subjects died, and the restricted mean survival estimate is 12.05 years.

Kaplan-Meier Survival Function Estimates, by Sex



Kaplan-Meier Survival Function Estimates, by Sex (code)

```
ggsurvplot(km_100_sex, data = ex100,  
            conf.int = TRUE,  
            xlab = "Time in years",  
            break.time.by = 3,  
            risk.table = TRUE,  
            risk.table.height = 0.25,  
            pval = TRUE)
```

- Note that I turned off the warning for this chunk of code. Otherwise you get the warning:

Vectorized input to `element_text()` is not officially supported. Results may be unexpected or may change in future versions of `ggplot2`.

Testing the difference between 2 survival curves

To obtain a significance test comparing these two survival curves, we turn to a log rank test, which tests the null hypothesis $H_0 : S_1(t) = S_2(t)$ for all t where the two exposures have survival functions $S_1(t)$ and $S_2(t)$.

```
survdif(surv_100 ~ ex100$sex)
```

Call:

```
survdif(formula = surv_100 ~ ex100$sex)
```

	N	Observed	Expected	(O-E) ² /E	(O-E) ² /V
ex100\$sex=Female	32	8	3.75	4.81	6.39
ex100\$sex=Male	68	9	13.25	1.36	6.39

Chisq= 6.4 on 1 degrees of freedom, p= 0.01

At usual α levels, there's a statistically significant difference ($p = 0.011$) between the survival curves stratified by sex.

Alternative log rank tests

An alternative is the *Peto and Peto modification of the Gehan-Wilcoxon test*, which results from adding $\rho=1$ to the `survdif` function ($\rho=0$, the default, yields the log rank test.)

```
survdif(surv_100 ~ ex100$sex, rho = 1)
```

Call:

```
survdif(formula = surv_100 ~ ex100$sex, rho = 1)
```

	N	Observed	Expected	$(O-E)^2/E$	$(O-E)^2/V$
ex100\$sex=Female	32	7.44	3.45	4.62	6.7
ex100\$sex=Male	68	7.79	11.79	1.35	6.7

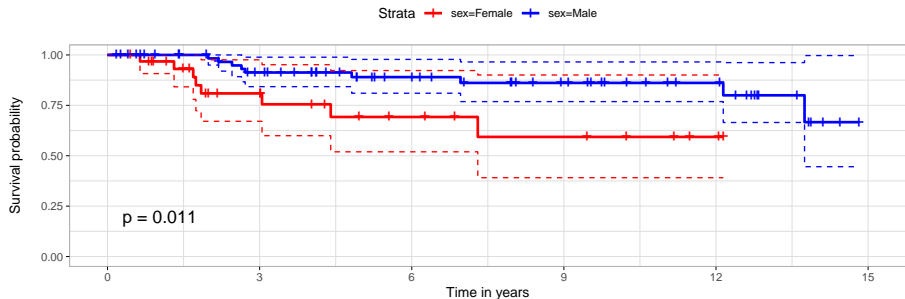
Chisq= 6.7 on 1 degrees of freedom, p= 0.01

Alternative log rank tests

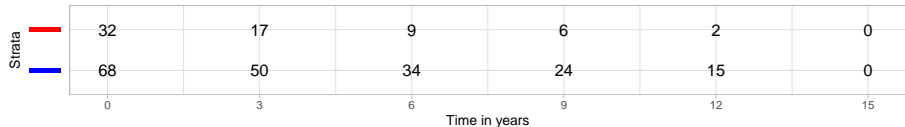
- As compared to the log rank test, this Peto-Peto modification (and others using $\rho > 0$) give greater weight to the left hand (earlier) side of the survival curves.
- To obtain chi-square tests that give greater weight to the right hand (later) side of the survival curves than the log rank test, use $\rho < 0$.

The log rank test generalizes to permit survival comparisons across more than two groups, with the test statistic having an asymptotic chi-squared distribution with one degree of freedom less than the number of patient groups being compared.

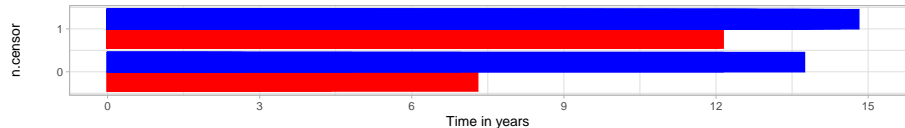
A Highly Customized K-M Plot



Number at risk



Number of censoring



Customizing the K-M Plot Further

See <https://rpkgs.datanovia.com/survminer/> or
<https://github.com/kassambara/survminer/> for many more options.

Comparing Survival Functions, by sex, 1000 observations

```
surv_obj2 <- Surv(time = survex$study.yrs,  
                  event = survex$death)
```

```
km_sex2 <- survfit(surv_obj2 ~ survex$sex)
```

```
survdif(surv_obj2 ~ survex$sex)
```

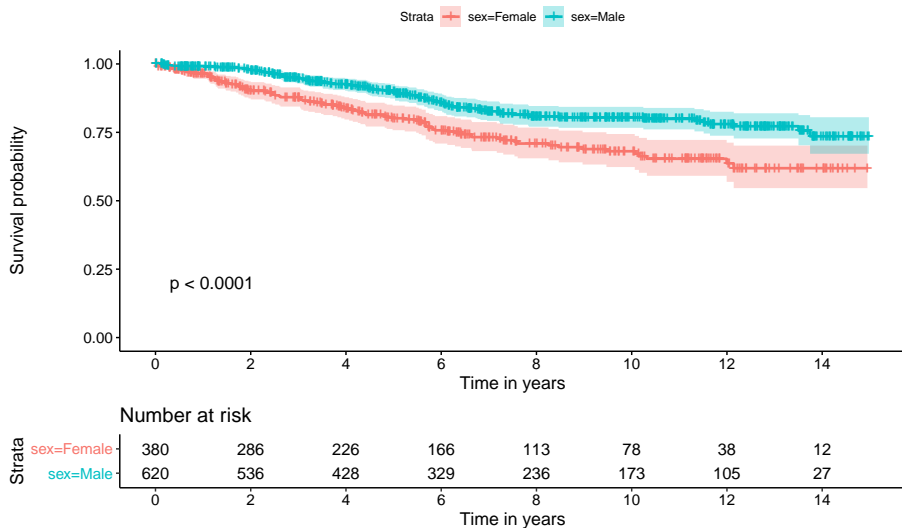
Call:

```
survdif(formula = surv_obj2 ~ survex$sex)
```

	N	Observed	Expected	$(O-E)^2/E$	$(O-E)^2/V$
survex\$sex=Female	380	90	62.7	11.85	18.1
survex\$sex=Male	620	93	120.3	6.18	18.1

Chisq= 18.1 on 1 degrees of freedom, p= 2e-05

Kaplan-Meier Plot of Survival, by Sex (n = 1000)

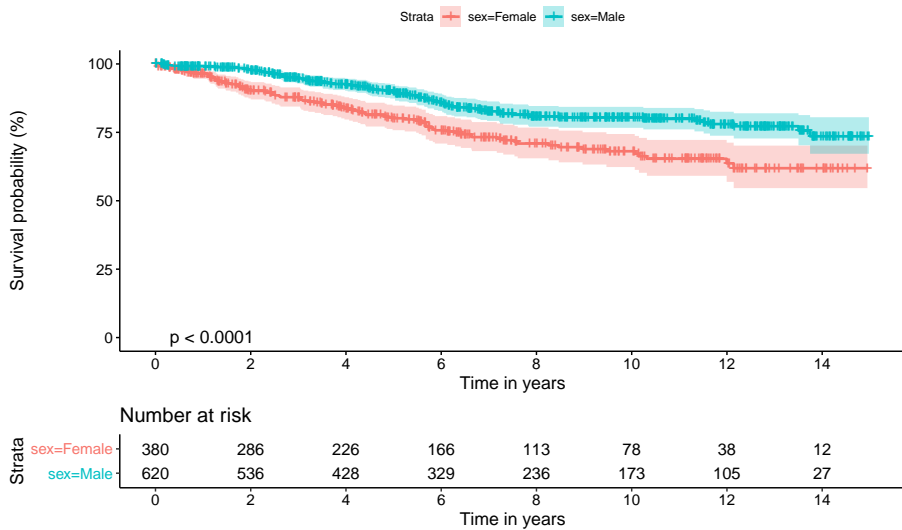


Kaplan-Meier Plot of Survival Percentage, Instead?

Just add `fun = "pct"` to the plot.

```
ggsurvplot(km_sex2, data = survex, fun = "pct",  
            conf.int = TRUE,  
            pval = TRUE,  
            xlab = "Time in years",  
            break.time.by = 2,  
            risk.table = TRUE,  
            risk.table.height = 0.25)
```

Kaplan-Meier Plot of Survival Percentage



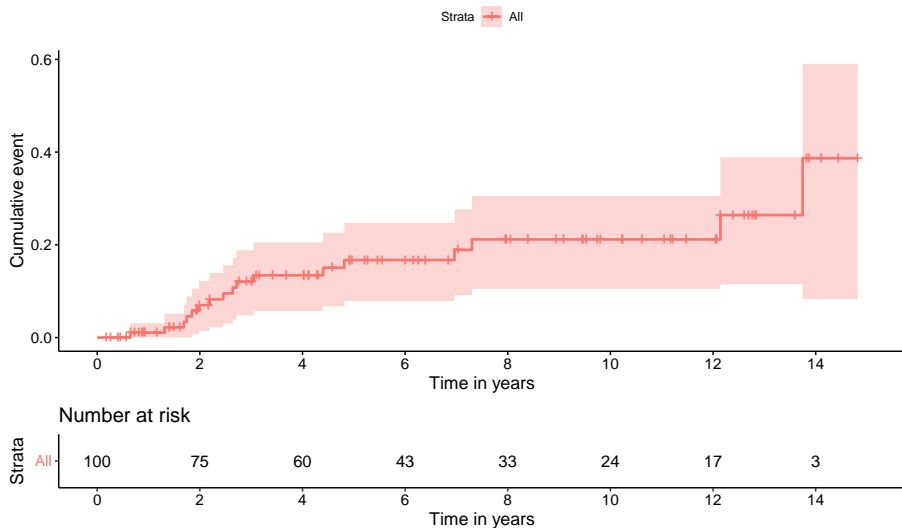
Code to plot Cumulative Event Rate

Let's look at our original km_100 model for 100 observations.

- Add fun = "event" to our ggsurvplot.

```
ggsurvplot(km_100, data = survex, fun = "event",  
           xlab = "Time in years",  
           break.time.by = 2,  
           risk.table = TRUE,  
           risk.table.height = 0.25)
```

Can we plot the cumulative event rate instead?

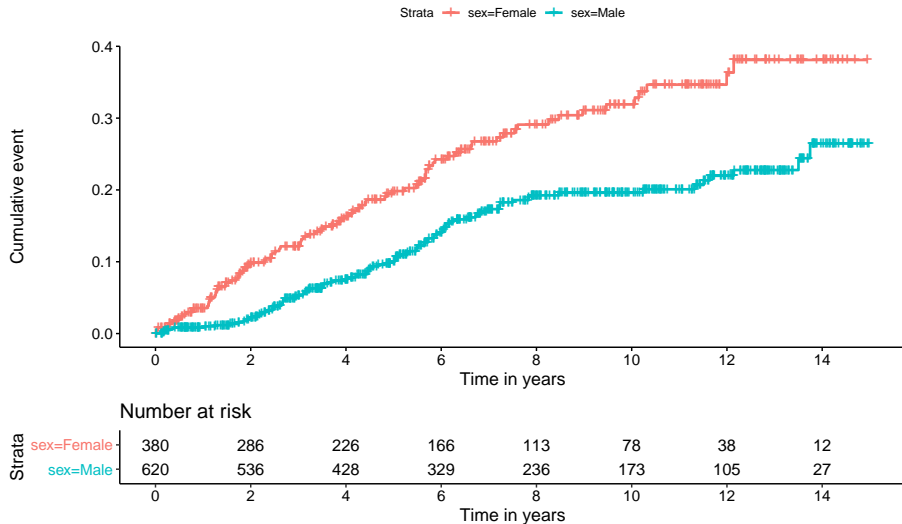


Cumulative Event Rate for km_sex2 model

Let's look at our model for 1000 observations, that includes sex:

```
ggsurvplot(km_sex2, data = survex, fun = "event",  
            xlab = "Time in years",  
            break.time.by = 2,  
            risk.table = TRUE,  
            risk.table.height = 0.25)
```

Cumulative Event Rate for km_sex2 model (Results)



The Hazard Function

To build regression models for time-to-event data, we will need to introduce the **hazard function**.

If $S(t)$ is the survival function, and time t is taken to be continuous, then $S(t) = e^{-H(t)}$ defines the hazard function $H(t)$.

- Note that $H(t) = -\ln(S(t))$.
- The function $H(t)$ is an important analytic tool.
 - It is used to describe the concept of the risk of “failure” in an interval after time t , conditioned on the subject having survived to time t .
 - It is often called the *cumulative hazard function*, to emphasize the fact that its value is the “sum” of the hazard up to time t .

Understanding the Hazard Function

Consider a subject in the survex study who has a survival time of 4 years.

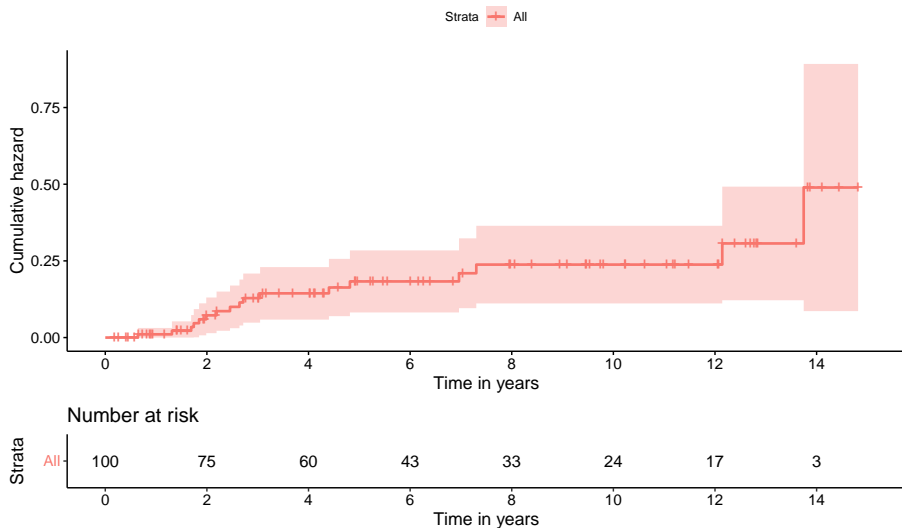
- For this subject to die at 4 years, they had to survive for the first 3 years.
- The subject's hazard at 4 years is the failure rate “per year” conditional on the subject being alive through the first 3 years.

Plotting the Cumulative Hazard Function

For our initial `km_100` fit, we'd use something like this...

```
ggsurvplot(km_100, data = survex, fun = "cumhaz",  
            xlab = "Time in years",  
            break.time.by = 2,  
            risk.table = TRUE,  
            risk.table.height = 0.25)
```

Cumulative Hazard Function for km_100 (Result)



Estimating the Cumulative Hazard Function

There are several different methods to estimate $H(t)$ and I'll discuss two now:

- 1 The inverse Kaplan-Meier estimator
- 2 The Nelson-Aalen estimator

The Inverse Kaplan-Meier Estimator of $H(t)$

I'll create something called H.est1, the inverse K-M estimate...

```
surv_100 <- Surv(ex100$study.yrs, ex100$death)
km_100 <- survfit(surv_100 ~ 1)
Haz1.almost <- -log(km_100$surv)
H_est1 <- c(Haz1.almost, tail(Haz1.almost, 1))
```

Next, we create a tibble of the times and hazard estimates

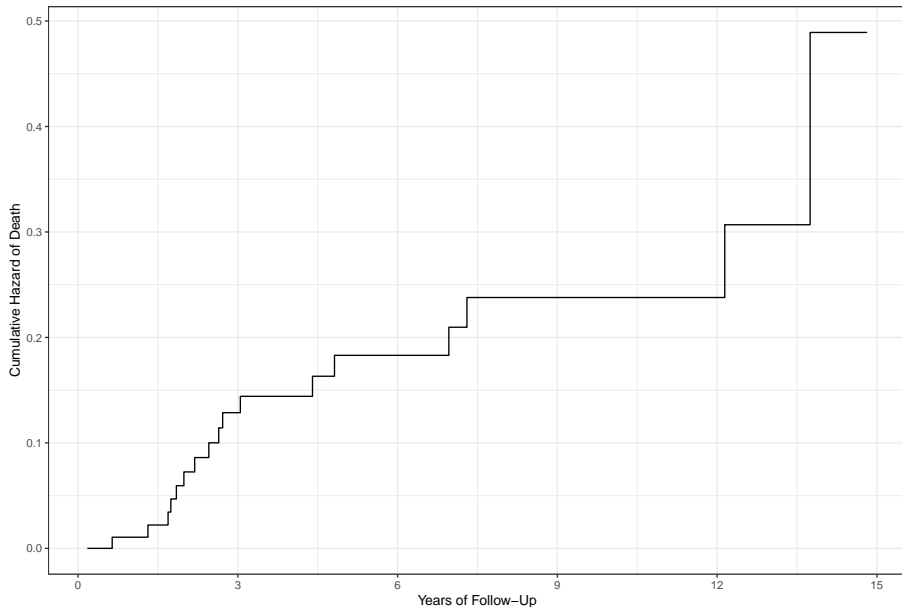
```
haz_frame <- tibble(
  time = c(km_100$time, tail(km_100$time, 1)),
  inverse_KM = H_est1
)
```

Cumulative Hazard Function from Inverse Kaplan-Meier (code)

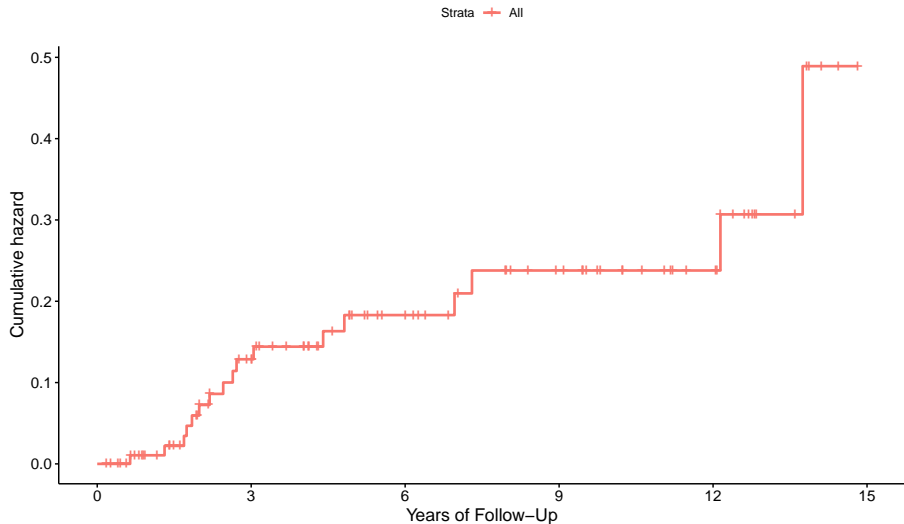
```
ggplot(haz_frame, aes(x = time, y = inverse_KM)) +  
  geom_step() +  
  scale_x_continuous(breaks = c(0, 3, 6, 9, 12, 15)) +  
  labs(x = "Years of Follow-Up",  
       y = "Cumulative Hazard of Death",  
       title = "Cumulative Hazard Function via Inverse K-M")
```

Cumulative Hazard Function (Inverse K-M)

Cumulative Hazard Function via Inverse K-M



Cumulative Hazard Plot via ggsurvplot



Nelson-Aalen Estimator of $H(t)$

We'll create the Nelson-Aalen estimate, `H_est2`.

```
h.sort.of <- km_100$n.event / km_100$n.risk
Haz2.almost <- cumsum(h.sort.of)
H_est2 <- c(Haz2.almost, tail(Haz2.almost, 1))
```

Add Nelson-Aalen Estimate to our Data Frame

```
haz_frame$Nelson_Aalen <- H_est2
```

```
haz_frame %>% slice(., c(5:6, 91:92))
```

```
# A tibble: 4 x 3
```

	time	inverse_KM	Nelson_Aalen
	<dbl>	<dbl>	<dbl>
1	0.569	0	0
2	0.641	0.0106	0.0105
3	13.6	0.307	0.302
4	13.7	0.489	0.469

Convert Wide Data to Long

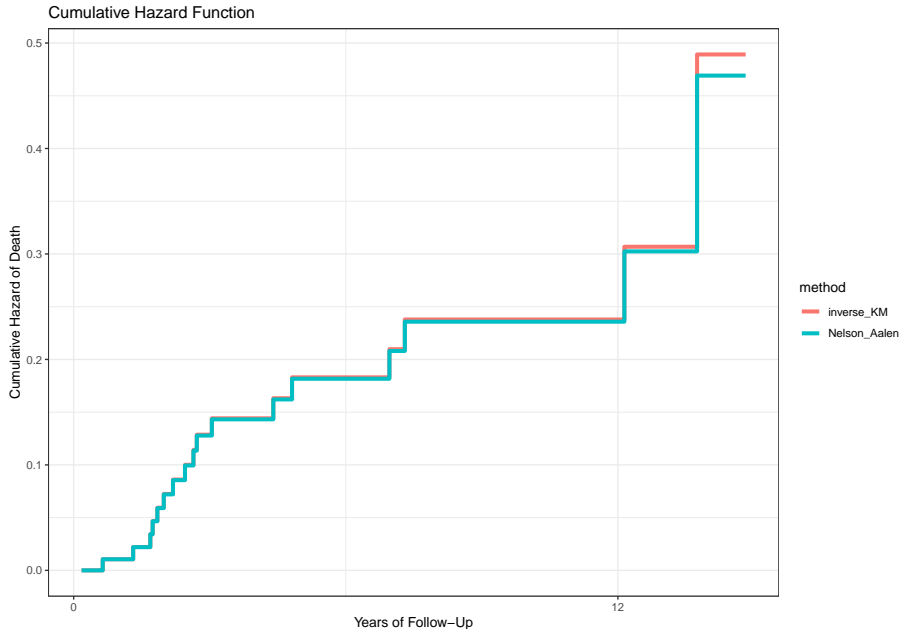
In order to easily plot the two hazard function estimates in the same graph, we'll want to convert these data from wide format to long format, with the `pivot_longer` function.

```
haz_frame_long <- pivot_longer(  
  haz_frame, cols = inverse_KM:Nelson_Aalen,  
  names_to = "method", values_to = "hazardest"  
)
```

```
tail(haz_frame_long)
```

```
# A tibble: 6 x 3  
   time method      hazardest  
<dbl> <chr>         <dbl>  
1  14.4 inverse_KM      0.489  
2  14.4 Nelson_Aalen    0.469  
3  14.8 inverse_KM      0.489  
4  14.8 Nelson_Aalen    0.469
```

Plot Hazard Estimates and Compare



Plot Hazard Estimates and Compare (code)

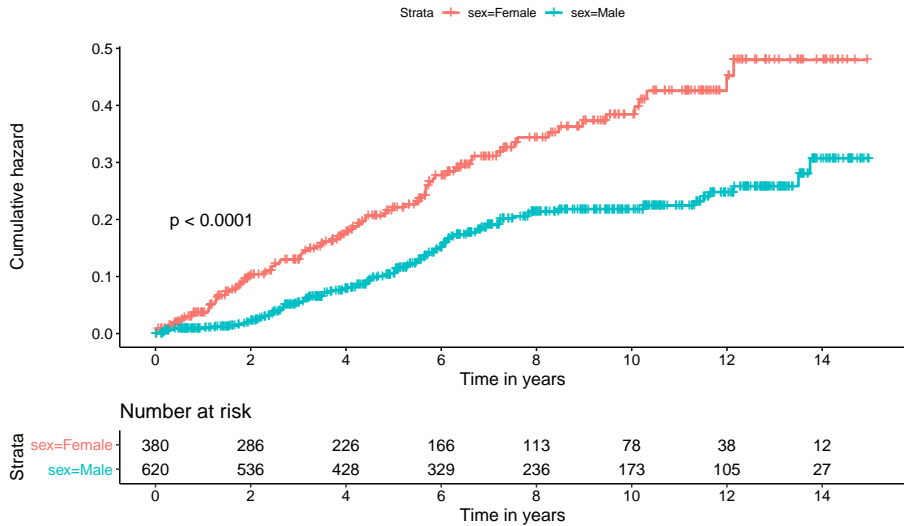
```
ggplot(haz_frame_long, aes(x = time, y = hazardest,  
                           col = method)) +  
  geom_step(size = 1.5) +  
  scale_x_continuous(breaks = c(0, 12, 24, 36, 48)) +  
  labs(x = "Years of Follow-Up",  
       y = "Cumulative Hazard of Death",  
       title = "Cumulative Hazard Function") +  
  theme_bw()
```

Plotting the Cumulative Hazard Function by Sex

For our `km_sex2` fit, we'd use something like this...

```
ggsurvplot(km_sex2, data = survex, fun = "cumhaz",  
            xlab = "Time in years",  
            pval = TRUE,  
            break.time.by = 2,  
            risk.table = TRUE,  
            risk.table.height = 0.25)
```

Cumulative Hazard Function for km_sex2 (Result)



What does this look like in a Cox model?

```
mod_sex <- survex %$%  
  coxph(Surv(study.yrs, death) ~ sex)
```

The Cox proportional hazards model fits survival data with a constant (not varying over time) covariate (here, sex) to a hazard function of the form:

$$h(t|sex) = h_0(t)\exp(\beta_1 sex)$$

where we estimate the unknown value of β_1 and where $h_0(t)$ is the baseline hazard which depends on t but not on sex.

Coefficients of our Cox model

```
mod_sex
```

Call:

```
coxph(formula = Surv(study.yrs, death) ~ sex)
```

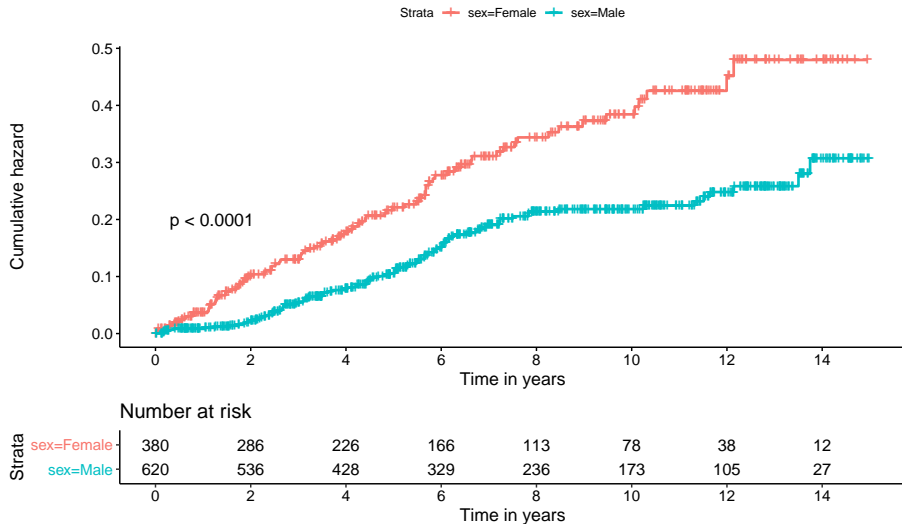
	coef	exp(coef)	se(coef)	z	p
sexMale	-0.6195	0.5382	0.1481	-4.184	2.86e-05

Likelihood ratio test=17.18 on 1 df, p=3.399e-05
n= 1000, number of events= 183

Our hazard ratio estimate is 0.5382 for Males (vs. Females)

- Hazard Ratio < 1 indicates a decrease in hazard for Males as compared to Females
- Does this match our plot?

The ggsurvplot of Cumulative Hazard (km_sex2)



What if we also include Age?

```
mod_age_sex <- coxph(Surv(study.yrs, death) ~ sex + age,  
                     data = survex)
```

Interpreting the Age + Sex model

```
mod_age_sex
```

Call:

```
coxph(formula = Surv(study.yrs, death) ~ sex + age, data = sur
```

	coef	exp(coef)	se(coef)	z	p
sexMale	-0.597528	0.550170	0.148207	-4.032	5.54e-05
age	0.041920	1.042811	0.005571	7.525	5.26e-14

Likelihood ratio test=69.93 on 2 df, p=6.522e-16

n= 1000, number of events= 183

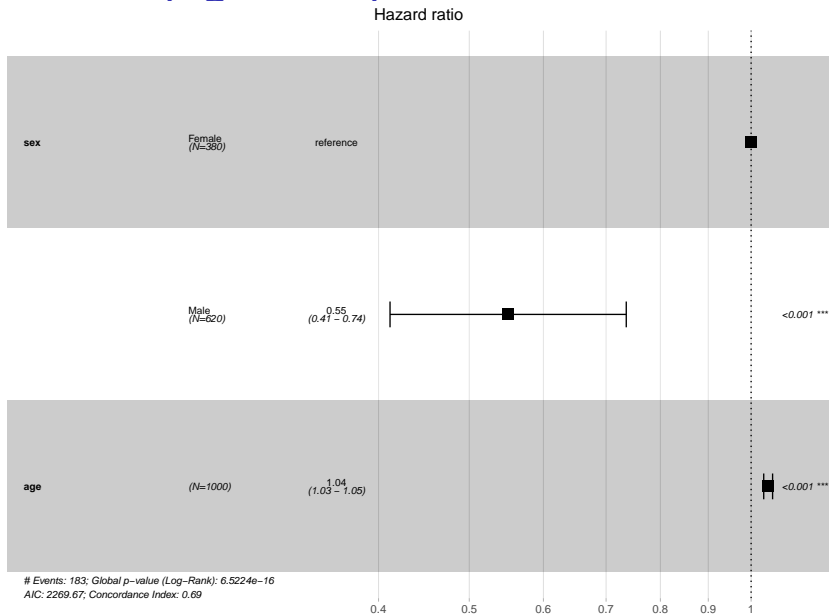
- If Harry is a year older than Steve and both are male, then Harry's hazard of death is 1.04 times that of Steve.
- If Harry (male) and Sally (female) are the same age, then Harry's hazard of death is 0.55 times that of Sally.

Summarizing the Cox Model with ggforest

Here is the code. Result on the next slide...

```
ggforest(mod_age_sex, data = survex)
```

Cox Model (Age + Sex) Coefficients



Next Time

- Diagnostics for Cox Proportional Hazards Regression
- Using `cph` from the `rms` package to fit a Cox model
- An Example with Real Data