

simon-5502-08-slides

Topics to be covered

- What you will learn
 - A simple example of survival data
 - Overall Kaplan-Meier curve
 - The log rank test
 - The hazard function
 - The Cox regression model
 - Assumptions and data management

Survival analysis

- Time to event models
 - Death
 - Relapse
 - Rehospitalization
 - Failure of medical device
 - Pregnancy
- Not every patient experiences the event
 - These are censored observations

Speaker notes

Survival analysis models are more properly called time to event models. You follow a group of patients from a certain time point and note the amount of time until they die.

Or the amount of time until they relapse. Or the amount of time until they need to be rehospitalized. Or the amount time until a medical device that you implanted in them fails.

I should note that while almost all of the events in a time to event model are bad, there are a few exceptions. In a study of couples with fertility problems, you might use a time to event model to study the time to pregnancy, a very happy event for any couple with fertility problems.

Mortality is the context under which time to event models were derived, so the term survival analysis has been used even when the event is different.

A key feature of survival analysis is that not every patient experiences the event. You should be glad that not everyone that you recruit for a clinical trial dies, but this adds a layer of complexity to the analysis.

The reasons for not experiencing the event can vary. Everyone dies, but maybe the event is death from cancer and if you patient gets hit by a bus, that patient does not experience the event. If the event is death from any cause, you still have to end the study in some time frame, and not everyone will die in that time frame. Not every patient gets rehospitalized, not every patient relapses, not every medical device fails.

A patient may drop out of a study, and you no longer are able to tell from that point onward whether that patient would have experienced the event sometime during the rest of your study.

When your patient does not die during the study, this is not a missing value. You have partial information. You know that the patient was alive for a certain amount of time. When you end the study within a certain time frame, you know that your patient dies at a time beyond the end date of your study. If your patient drops out after six months, you know that the patient survived for more than six months.

First fruit fly experiment, 1

`data_dictionary: fly1.txt`

`description: |`

`This dataset provides a simple example of what survival and censoring. It provides an intuitive explanation of estimation of survival probabilities.`

`vars:`

`day:`

`label: Time until death`

`unit: days`

Speaker notes

The following data represents survival time for a group of fruit flies and is a subset of a larger data set found at the Data and Story Library (DASL). The data set has been slightly modified to simplify some of these explanations.

There are 25 flies in the sample, with the first fly dying on day 37 and the last fly dying on day 96.

First fruit fly experiment, 2

37, 40, 43, 44, 45, 47, 49, 54, 56, 58, 59, 60, 61, 62, 68, 70, 71, 72, 73,
75, 77, 79, 89, 94, 96

Speaker notes

If you wanted to estimate the survival probability for this data, you would draw a curve that decreases by 4% ($1/25$) every time a fly dies.

First fruit fly experiment, 3

	day	p
1	37	96%
2	40	92%
3	43	88%
4	44	84%
5	45	80%
6	47	76%
7	49	72%
8	54	68%
9	56	64%

	day	p
10	58	60%
11	59	56%
12	60	52%
13	61	48%
14	62	44%
15	68	40%
16	70	36%
17	71	32%
18	72	28%

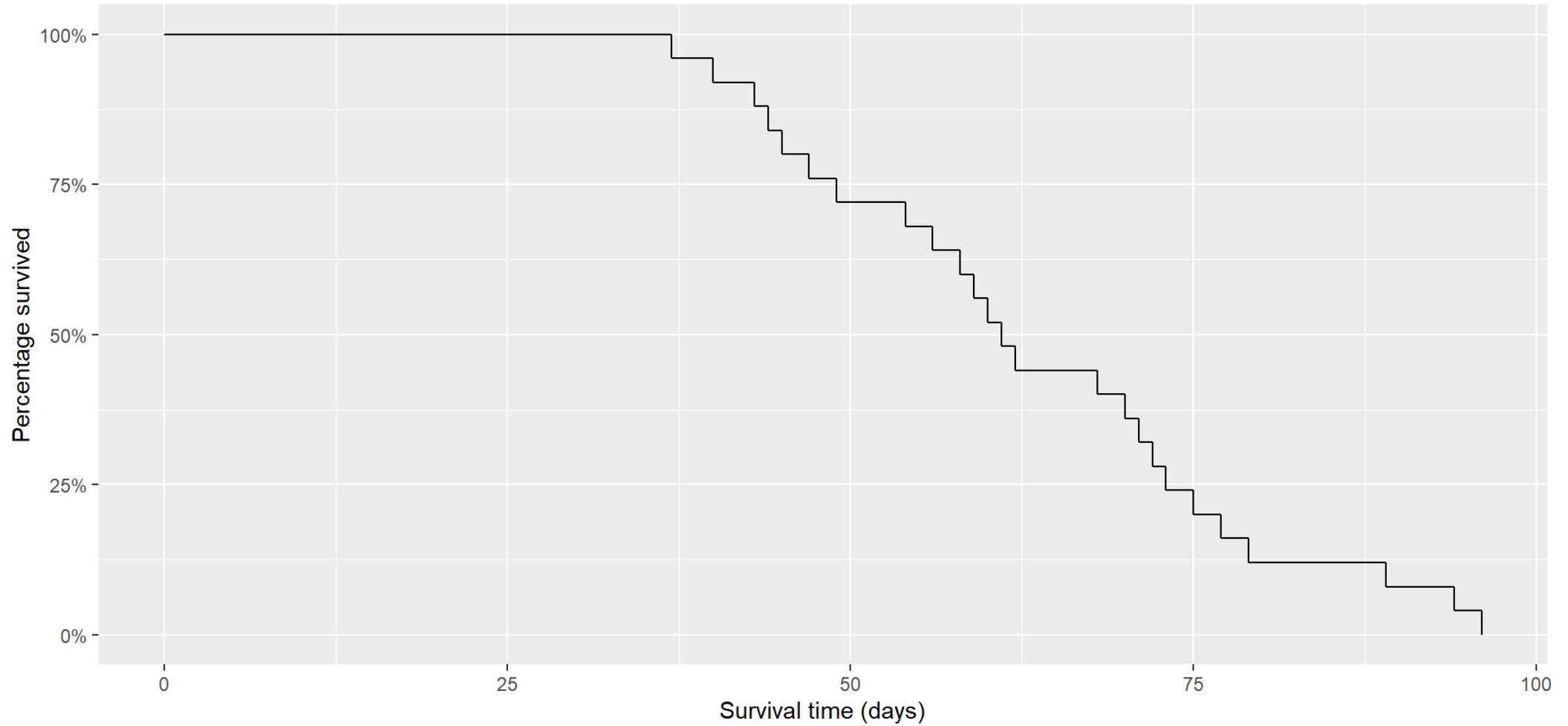
	day	p
19	73	24%
20	75	20%
21	77	16%
22	79	12%
23	89	8%
24	94	4%
25	96	0%

Speaker notes

The probability of survival drops by 4% ($1/25$) at each day of death.

First fruit fly experiment, 4

Graph drawn by Steve Simon on 2025-03-10



Speaker notes

Here's a graph of these probabilities over time.

By tradition and for some rather technical reasons, you should use a stair step pattern rather than a diagonal line to connect adjacent survival probabilities.

Second fruit fly experiment, 1

37, 40, 43, 44, 45, 47, 49, 54, 56, 58, 59, 60, 61, 62, 68, ??, ??, ??, ??,
??, ??, ??, ??, ??, ??

Speaker notes

Now let's alter the experiment. Suppose that totally by accident, a technician leaves the screen cover open on day 70 and all the flies escape. This includes the fly who was going to die on the afternoon of the 70th day anyway. Oh the sadness of it all; the poor fly has the briefest of tastes of freedom then ends up shriveled up on a window sill.

You're probably worried that the whole experiment has been ruined. But don't be so pessimistic. You still have complete information on survival of the fruit flies up to their 70th day of life.

Second fruit fly experiment, 2

	day	event
1	37	1
2	40	1
3	43	1
4	44	1
5	45	1
6	47	1
7	49	1
8	54	1
9	56	1

	day	event
10	58	1
11	59	1
12	60	1
13	61	1
14	62	1
15	68	1
16	70	0
17	70	0
18	70	0

	day	event
19	70	0
20	70	0
21	70	0
22	70	0
23	70	0
24	70	0
25	70	0

Speaker notes

Here's how you would code the data for importing into SPSS or any other software.

Second fruit fly experiment, 3

	day	event	p
1	37	1	96%
2	40	1	92%
3	43	1	88%
4	44	1	84%
5	45	1	80%
6	47	1	76%
7	49	1	72%
8	54	1	68%
9	56	1	64%

	day	event	p
10	58	1	60%
11	59	1	56%
12	60	1	52%
13	61	1	48%
14	62	1	44%
15	68	1	40%
16	70	0	
17	70	0	
18	70	0	

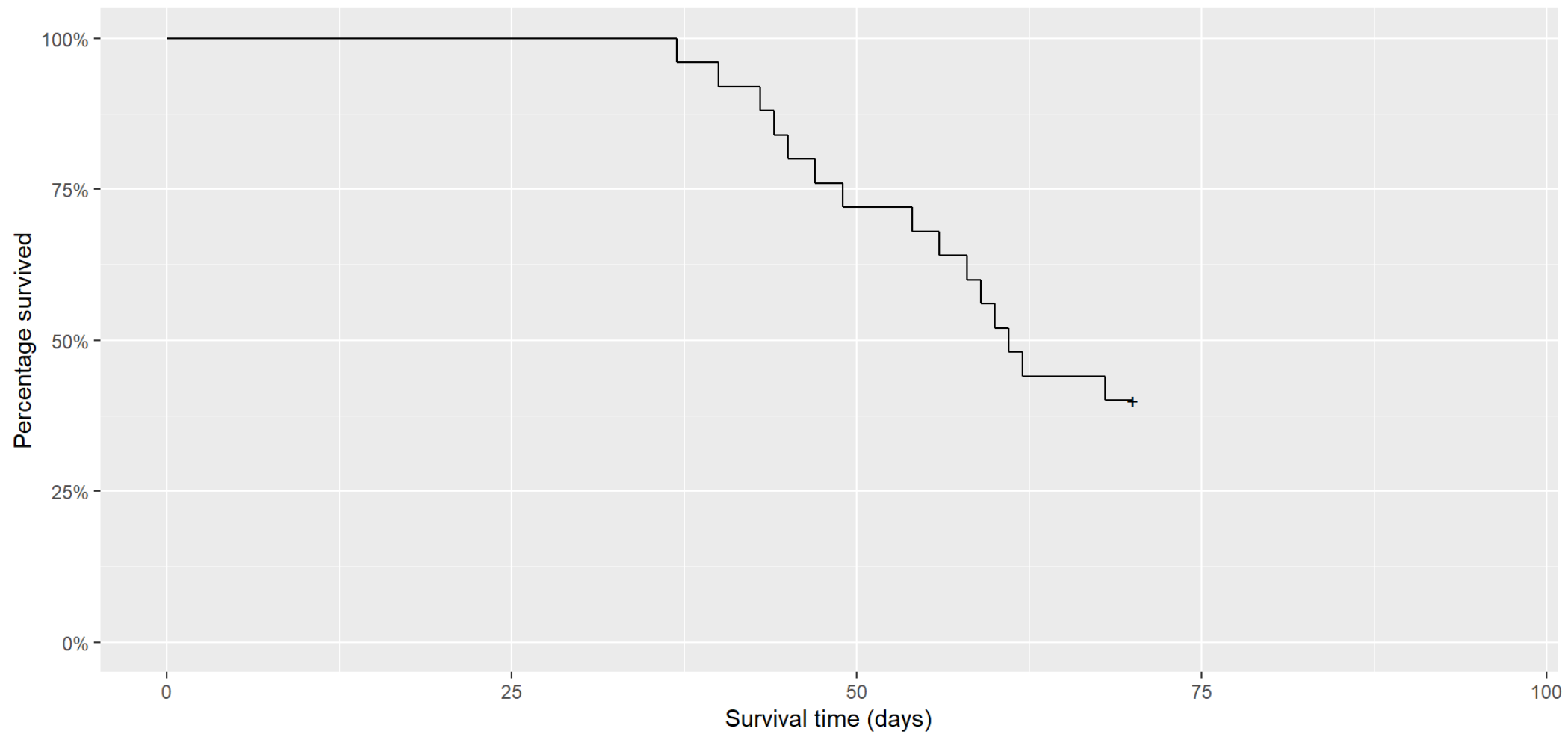
	day	event	p
19	70	0	
20	70	0	
21	70	0	
22	70	0	
23	70	0	
24	70	0	
25	70	0	

Speaker notes

We clearly have enough data to make several important statements about survival probability. For example, the median survival time is 61 days because roughly half of the flies had died before this day.

Second fruit fly experiment, 4

Graph drawn by Steve Simon on 2025-03-10



Speaker notes

Here is a graph of the survival probabilities of the second experiment. This graph is identical to the graph in the first experiment up to day 70 after which you can no longer estimate survival probabilities.

By the way, you might be tempted to ignore the ten flies who escaped. But that would seriously bias your results. All of these flies were tough and hardy flies who lived well beyond the median day of death. If you pretended that they didn't exist, you would seriously underestimate the survival probabilities. The median survival time, for example, of the 15 flies who did not escape, for example, is only 54 days which is much smaller than the actual median.

Third fruit fly experiment, 1

37, 40, 43, 44, 45, 47, 49, 54, 56, 58, 59, 60, 61, 62, 68, ??, 71, ??, ??,
75, ??, ??, 89, ??, 96

Speaker notes

Let's look at a third experiment, where the screen cover is left open and all but four of the remaining flies escape. It turns out that those four remaining flies who didn't bug out will allow us to still get reasonable estimates of survival probabilities beyond 70 days.

Third fruit fly experiment, 2

	day	event
1	37	1
2	40	1
3	43	1
4	44	1
5	45	1
6	47	1
7	49	1
8	54	1
9	56	1

	day	event
10	58	1
11	59	1
12	60	1
13	61	1
14	62	1
15	68	1
16	70	0
17	71	1
18	70	0

	day	event
19	70	0
20	75	1
21	70	0
22	70	0
23	89	1
24	70	0
25	96	1

Speaker notes

Here is how you would code the data for importing into SPSS.

Third fruit fly experiment, 3

	day	event	p
1	37	1	96%
2	40	1	92%
3	43	1	88%
4	44	1	84%
5	45	1	80%
6	47	1	76%
7	49	1	72%
8	54	1	68%
9	56	1	64%

	day	event	p
10	58	1	60%
11	59	1	56%
12	60	1	52%
13	61	1	48%
14	62	1	44%
15	68	1	40%
16	70	0	
17	71	1	30%
18	70	0	

	day	event	p
19	70	0	
20	75	1	20%
21	70	0	
22	70	0	
23	89	1	10%
24	70	0	
25	96	1	0%

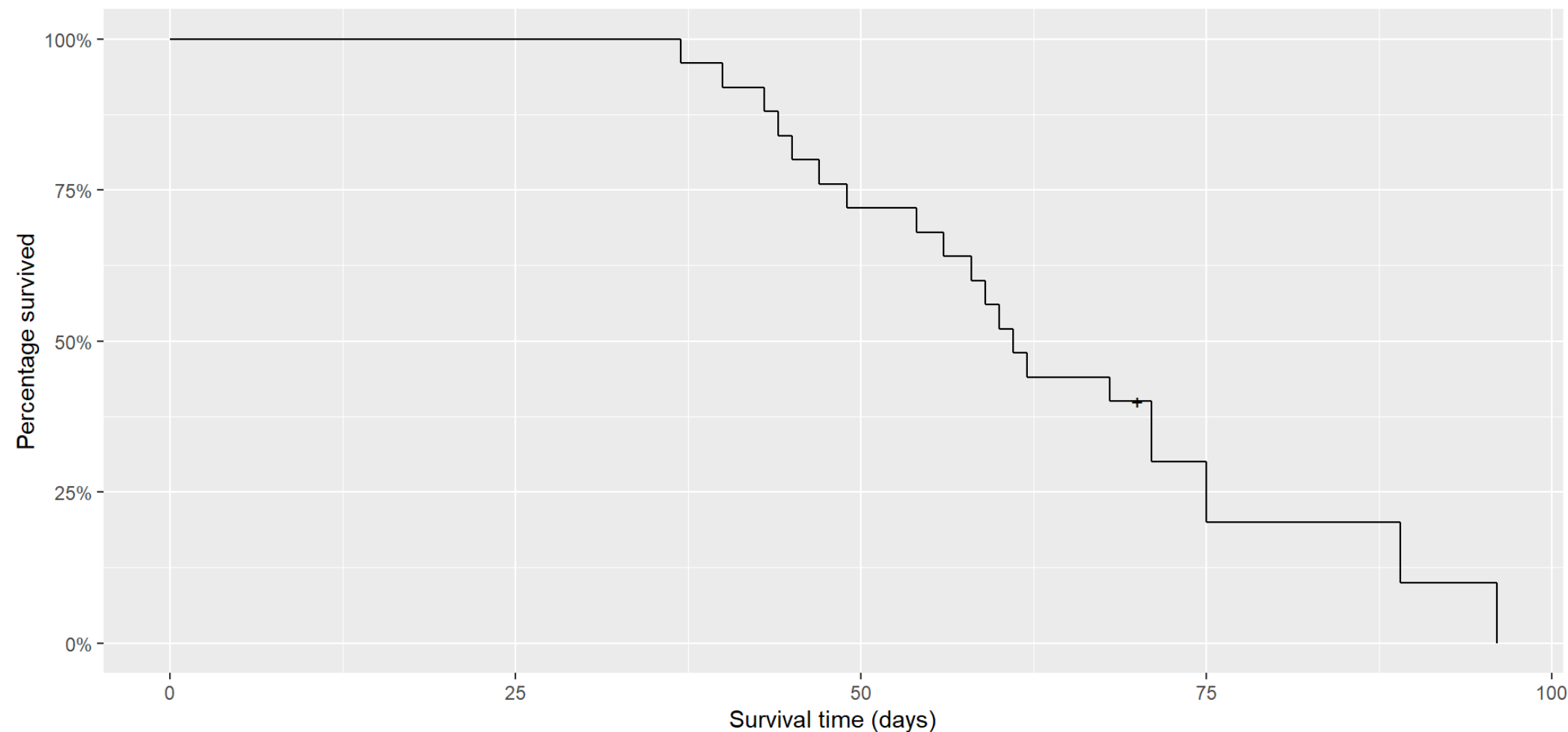
Speaker notes

What you do with the six escaped flies is to allocate their survival probabilities equally among the four flies who didn't bug out. This places a great responsibility among each of those four remaining flies since each one is now responsible for 10% of the remaining survival probability, their original 4% plus 6% more which represents a fourth of the 24% survival probability that was lost with the six escaping flies.

Another way of looking at this is that the six flies who escaped influence the denominator of the survival probabilities up to day 70 and then totally drop out of the calculations for any further survival probabilities. Because the denominator has been reduced, the jumps at each remaining death are much larger.

Third fruit fly experiment, 4

Graph drawn by Steve Simon on 2025-03-10



Speaker notes

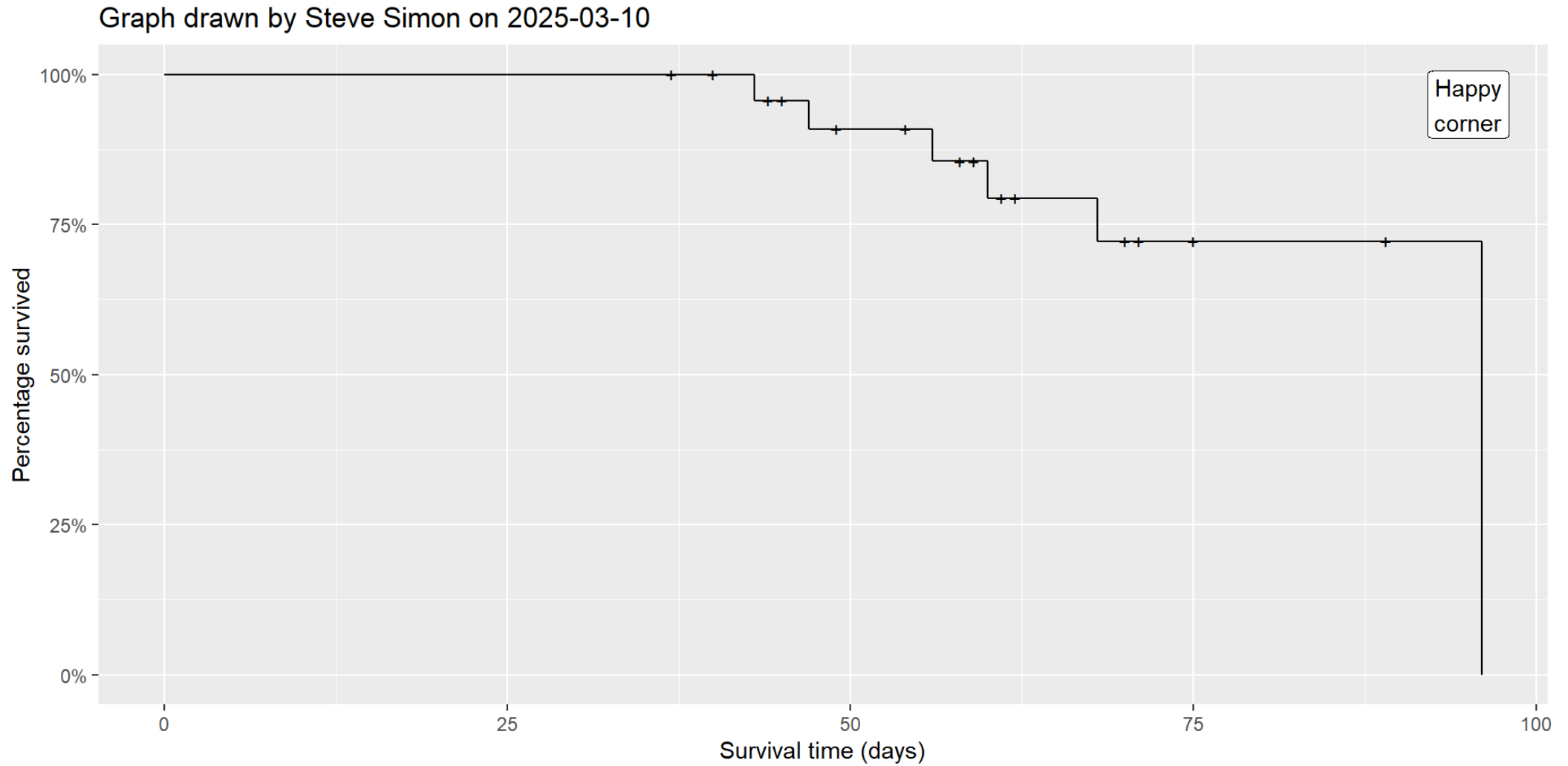
Here is a graph of the survival probability estimates from the third experiment.

These survival probabilities differ only slightly from the survival probabilities in the original experiment. This works out because the mechanism that caused us to lose information on six of the fruit flies was independent of their ultimate survival.

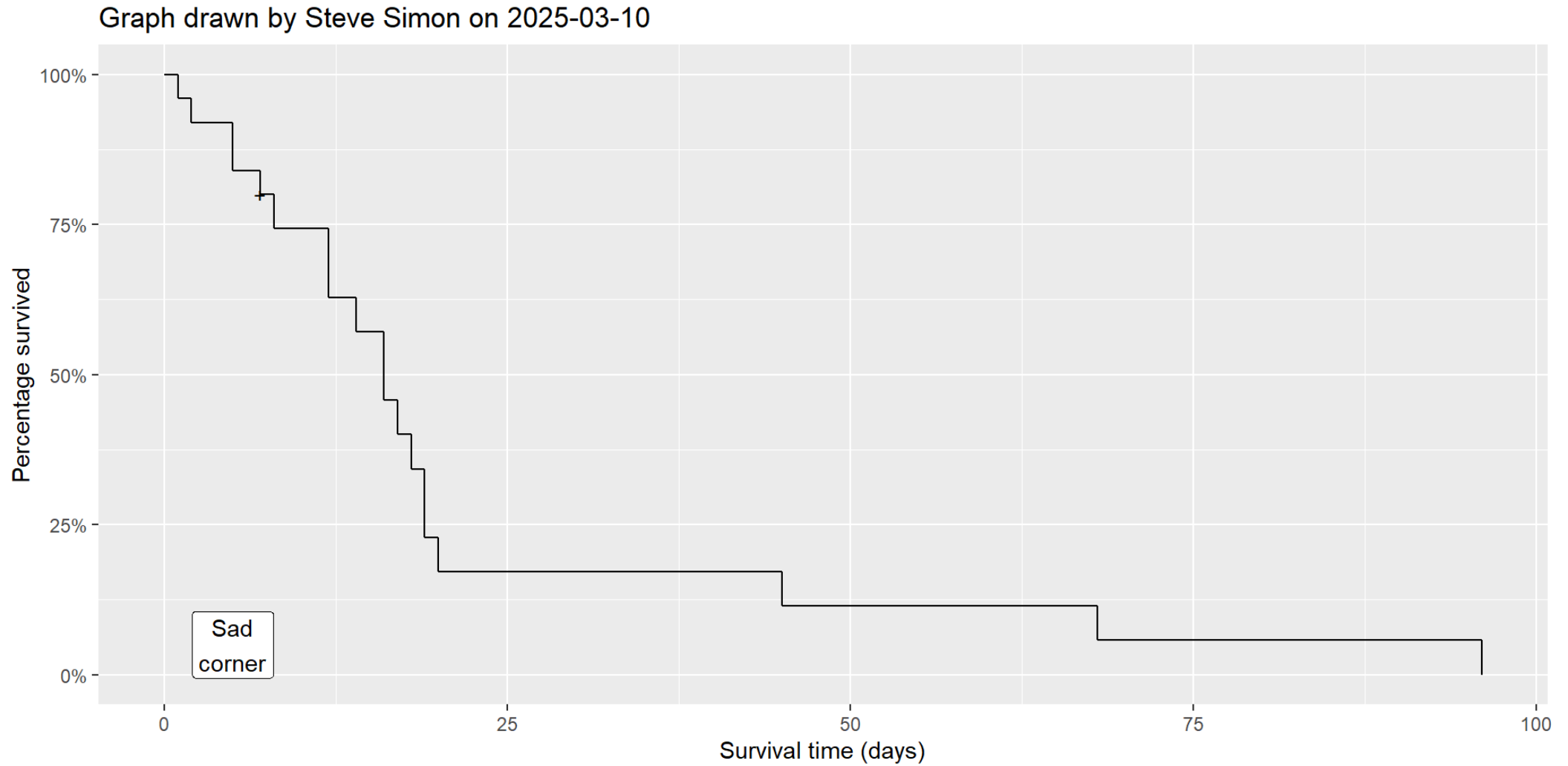
If the censoring mechanism were somehow related to survival prognosis, then you would have the possibility of serious bias in your estimates. Suppose for example, that only the toughest of flies (those with the most days left in their short lives) would have been able to escape. The flies destined to die on days 70, 71, 72, and 73, were already on their deathbeds and unable to fly at all, much less make a difficult escape. Then these censored values would not be randomly interspersed among the remaining survival times, but would constitute some of the larger values. But since these larger values would remain unobserved, you would underestimate survival probabilities beyond the 70th day.

This is known as informative censoring, and it happens more often than you might expect. Suppose someone drops out of a cancer mortality study because they are abandoning the drugs being studied in favor of laetrile treatments down in Mexico. Usually, this is a sign that the current drugs are not working well, so a censored observation here might represent a patient with a poorer prognosis. Excluding these patients would lead to an overestimate of survival probabilities.

Interpreting Kaplan-Meier plots, 1



Interpreting Kaplan-Meier plots, 2



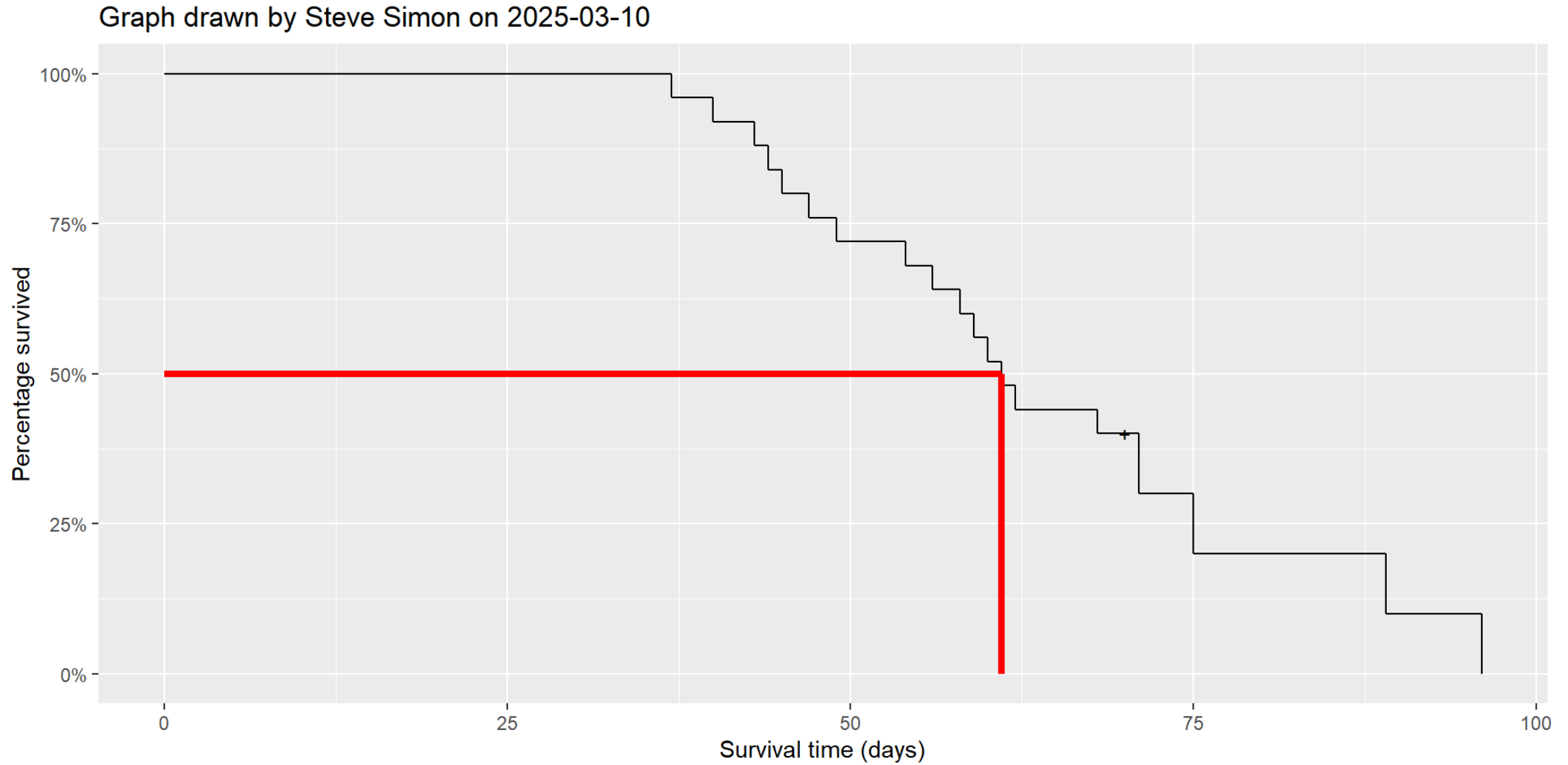
Speaker notes

When you see a survival curve in a research paper, there are three ways to interpret it.

First, presuming that the event in question is a sad event (such as death, relapse), then the upper right hand corner is the happy corner. Most of your patients go for a very long time with only a small proportion suffering the negative event.

In contrast, the lower left corner is the sad corner. Most of your patients experience the bad event, and they experience it very quickly.

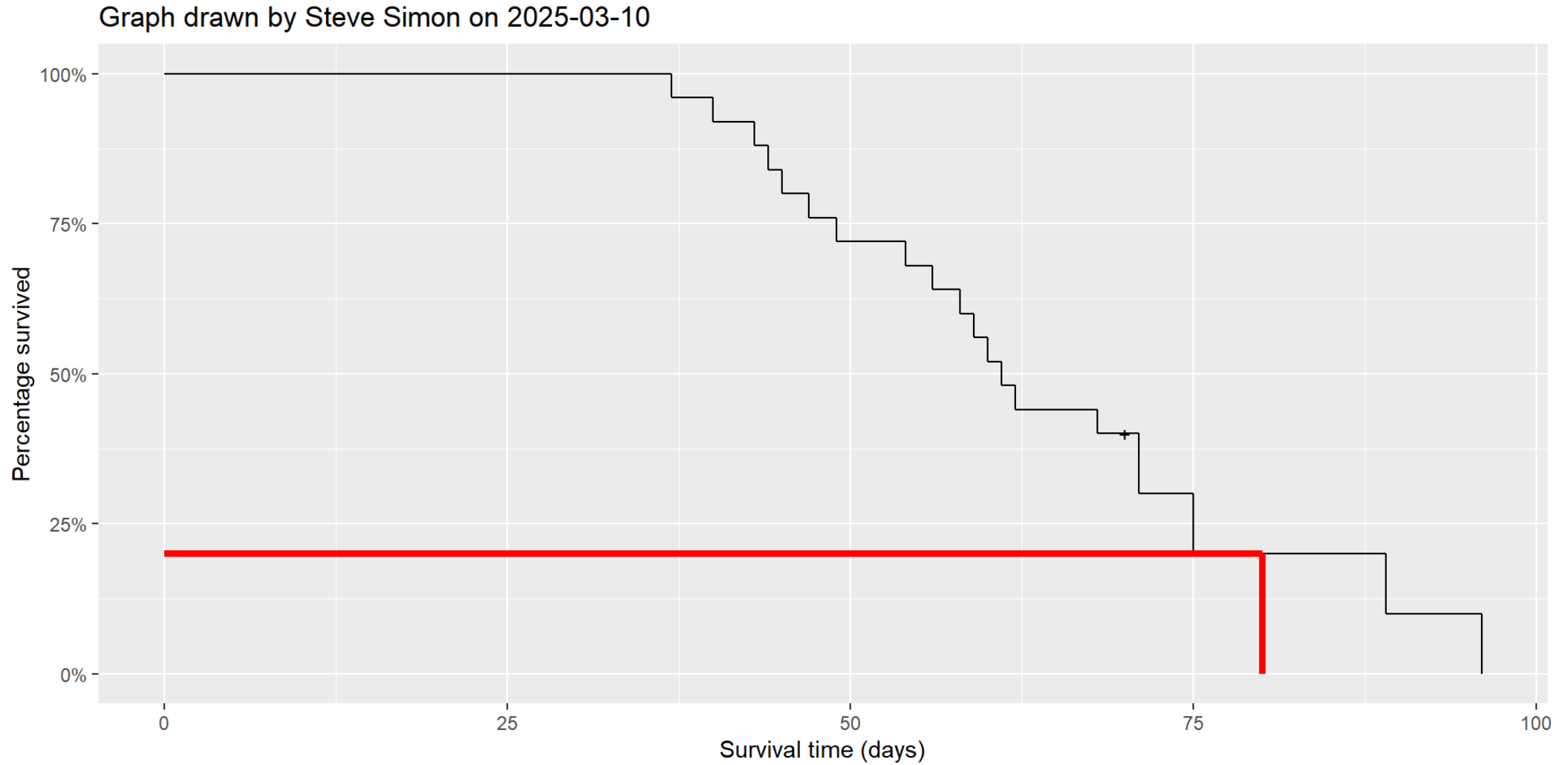
Interpreting Kaplan-Meier plots, 3



Speaker notes

Next, you can get an estimate of the median (or other percentiles) by projecting horizontally until you intersect with the survival curve and then head down to get your estimate. In the survival curve we have just looked at, you would estimate the median survival as slightly more than 60 days.

Interpreting Kaplan-Meier plots, 4



Speaker notes

You can also estimate probabilities for survival at any given time by projecting up from the time and then moving to the left to estimate the probability. In the example below, you can see that the 80 day survival probability is a little bit less than 25%.

Break #1

- What you have learned
 - A simple example of survival data
- What's coming next
 - Overall Kaplan-Meier curve

Worcester Heart Attack Study, 1

`data_dictionary: whas500.dat`

`description: The data represents survival times for a 500 patient subset of data from`

`the Worcester Heart Attack Study. You can find more information about this data`

`set in Chapter 1 of Hosmer, Lemeshow, and May.`

Worcester Heart Attack Study, 2

```
id:  
  label: a sequential code from 1 to 100  
age:  
  scale: ratio  
  label: Age at Admission  
  unit: years  
gender:  
  scale: binary  
  value:  
    Male: 0  
    Female: 1
```


Worcester Heart Attack Study, 3

hr:

scale: ratio

label: Initial Heart Rate

unit: Beats per minute

sysbp:

scale: ratio

label: Initial Systolic Blood Pressure

unit: mmHg

diasbp:

scale: ratio

label: Initial Diastolic Blood Pressure

unit: mmHg

Worcester Heart Attack Study, 4

```
bmi:  
  scale: ratio  
  label: Body Mass Index  
  unit: kg/m^2  
cvd:  
  scale: binary  
  label: History of Cardiovascular Disease  
value:  
  'FALSE': 0  
  'TRUE': 1
```


Worcester Heart Attack Study, 5

```
afb:
  scale: binary
  label: Atrial Fibrillation
  value:
    'FALSE': 0
    'TRUE': 1
sho:
  scale: binary
  label: Cardiogenic Shock
  value:
    'FALSE': 0
    'TRUE': 1
```


Worcester Heart Attack Study, 6

```
chf:
  scale: binary
  label: Congestive Heart Complications
  value:
    'FALSE': 0
    'TRUE': 1
av3:
  scale: binary
  label: Complete Heart Block
  value:
    'FALSE': 0
    'TRUE': 1
```


Worcester Heart Attack Study, 7

```
miord:
  scale: binary
  label: MI Order
  value:
    First: 0
    Recurrent: 1
mitype:
  scale: binary
  label: MI Type
  value:
    non Q-wave: 0
    Q-wave: 1
```


Worcester Heart Attack Study, 8

```
year:  
  scale: ordinal  
  label: Cohort Year  
  value:  
    yr1997: 1  
    yr1999: 2  
    yr2001: 3  
admitdate:  
  label: Admission Date  
  format: mm/dd/yyyy  
disdate:  
  label: Hospital Discharge Date  
  format: mm/dd/yyyy
```


Worcester Heart Attack Study, 9

fdate:

label: Date of last Follow Up

format: mm/dd/yyyy

los:

scale: ratio

label: Length of Hospital Stay

unit: Days

dstat:

scale: binary

label: Discharge Status from Hospital

value:

Alive: 0

Dead: 1

Worcester Heart Attack Study, 9

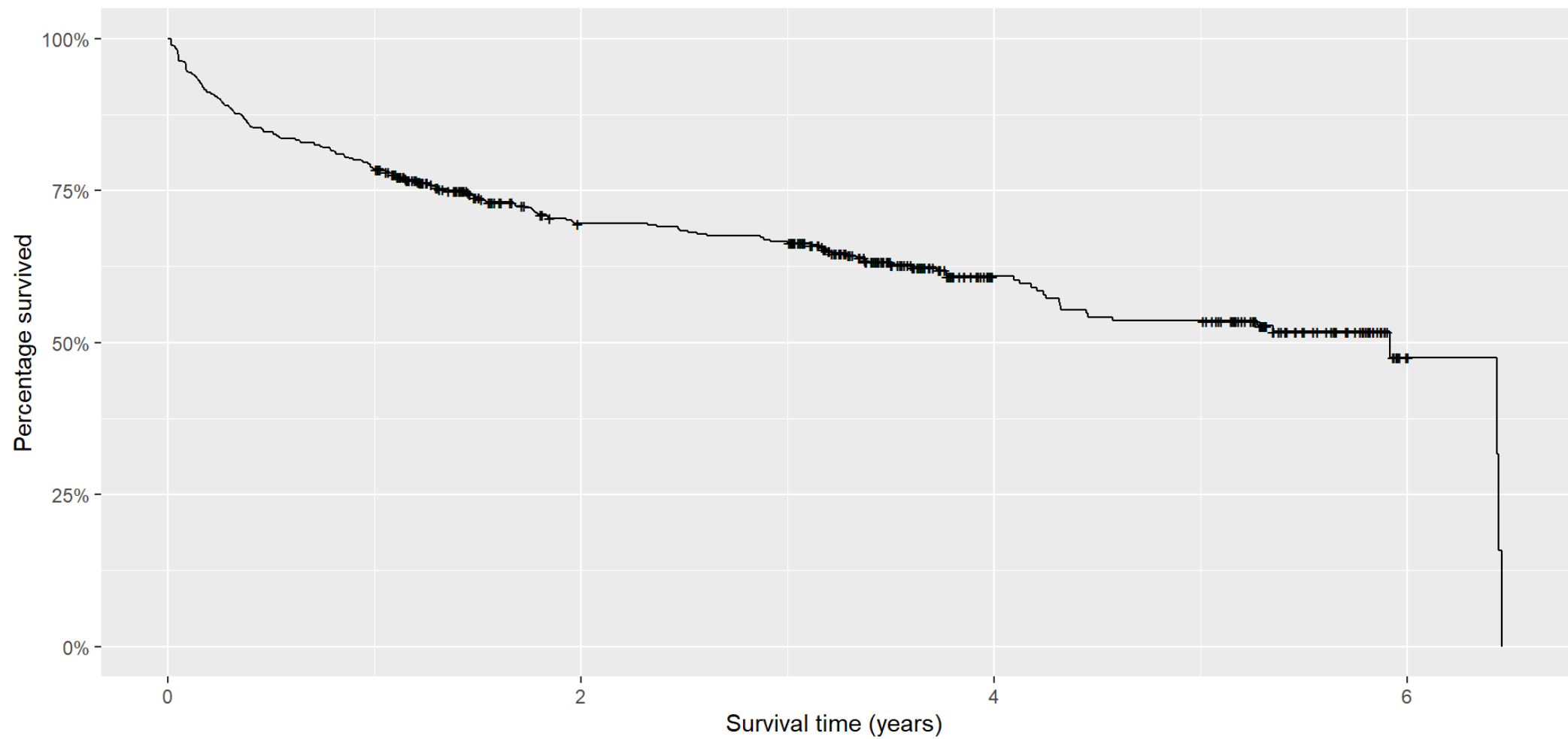
```
lenfol:
  scale: ratio
  label: Follow Up Time
  unit: days
fstat:
  scale: binary
  label: Vital Satus
value:
  Alive: 0
  Dead: 1
```


Event count

```
# A tibble: 2 × 2
  fstat      n
  <dbl> <int>
1     0    285
2     1    215
```


Overall Kaplan-Meier curve

Graph drawn by Steve Simon on 2025-03-10



Live demo, Overall Kaplan-Meier curve

Break #2

- What you have learned
 - Overall Kaplan-Meier curve
- What's coming next
 - The log rank test

Cox regression for gender

Call:

```
coxph(formula = Surv(lenfol, fstat) ~ gender, data = whas_4)
```

	coef	exp(coef)	se(coef)	z	p
gender	0.3417	1.4074	0.1526	2.24	0.0251

Likelihood ratio test=4.95 on 1 df, p=0.02606

n= 461, number of events= 176

Covariate imbalance

```
# A tibble: 2 × 2
  gender age_mean
  <dbl>   <dbl>
1     0     65.9
2     1     74.0
```

Cox regression adjusted for age

Call:

```
coxph(formula = Surv(lenfol, fstat) ~ gender, data = whas_4)
```

	coef	exp(coef)	se(coef)	z	p
gender	0.3417	1.4074	0.1526	2.24	0.0251

Likelihood ratio test=4.95 on 1 df, p=0.02606

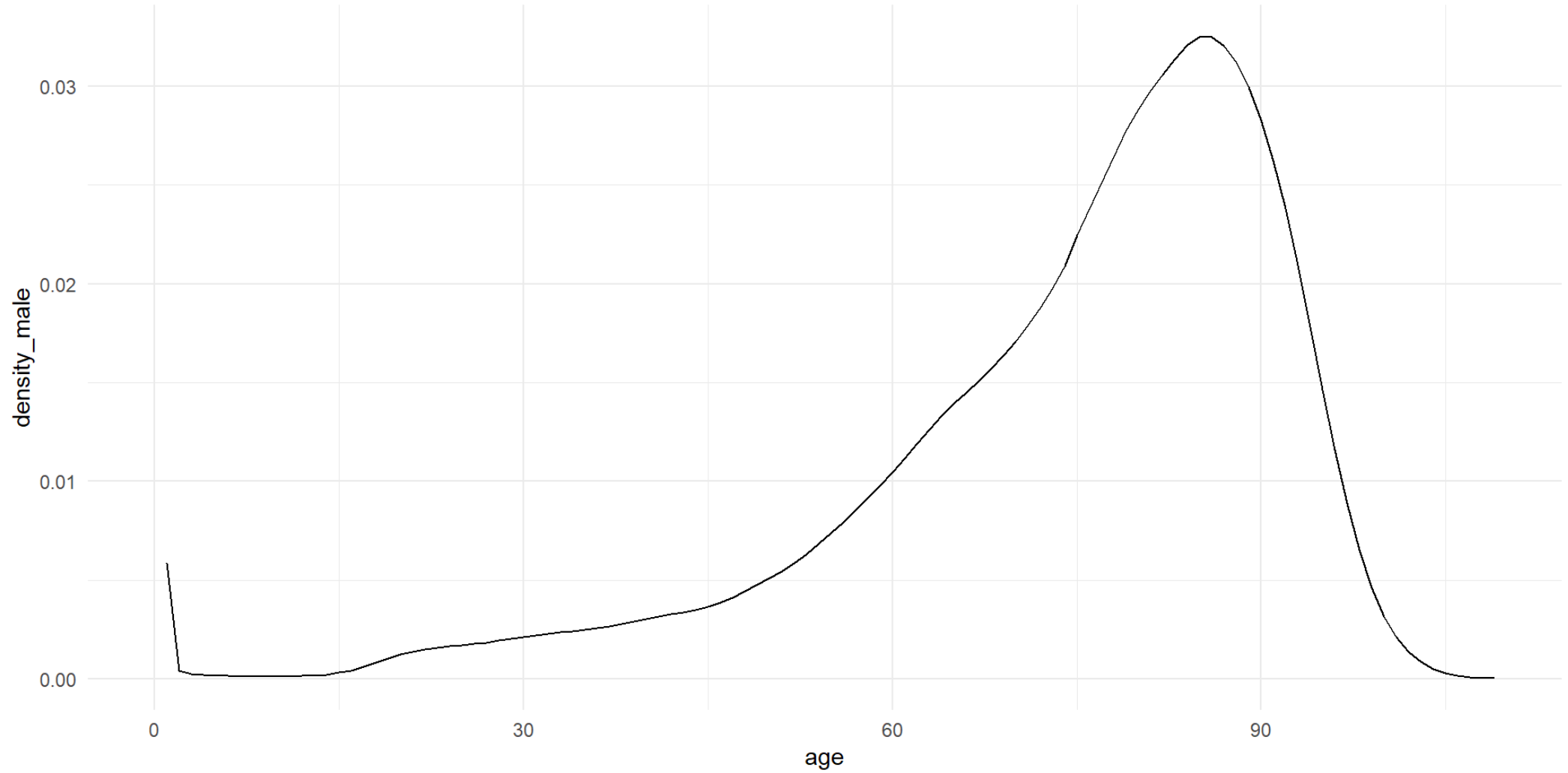
n= 461, number of events= 176

Live demo, The log rank test

Break #3

- What you have learned
 - The log rank test
- What's coming next
 - The hazard function

Life insurance example

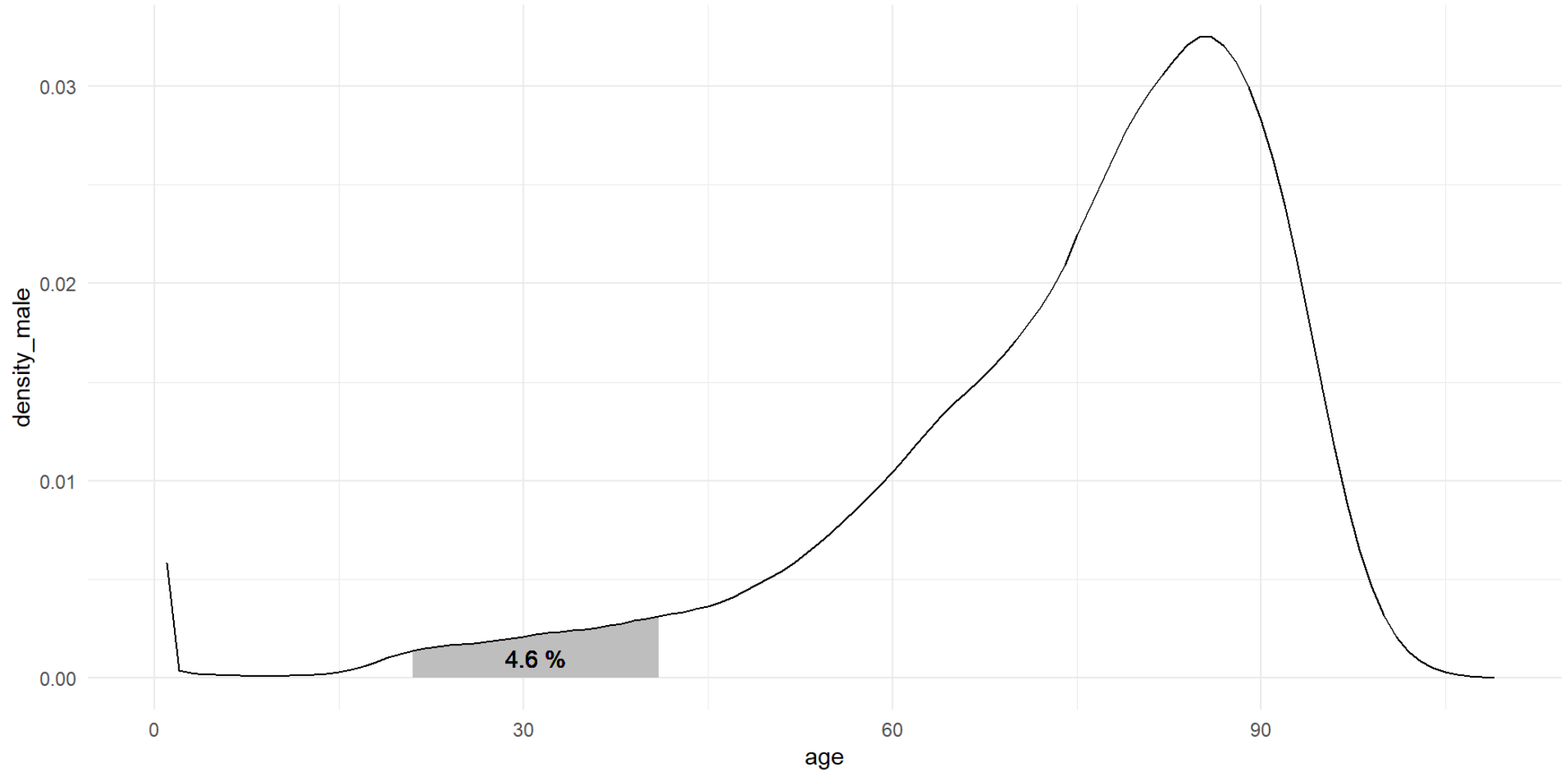


Speaker notes

I found some data on mortality from the Social Security website and plotted an approximation to the probability density function. There is an unusual early peak in this function because the first year of your life is one of the most dangerous ones you will have to face.

Imagine yourself working in life insurance sales. You want to price your policies so that you only ask for low payments on the policy when the risk of death is low. So let's calculate some probabilities.

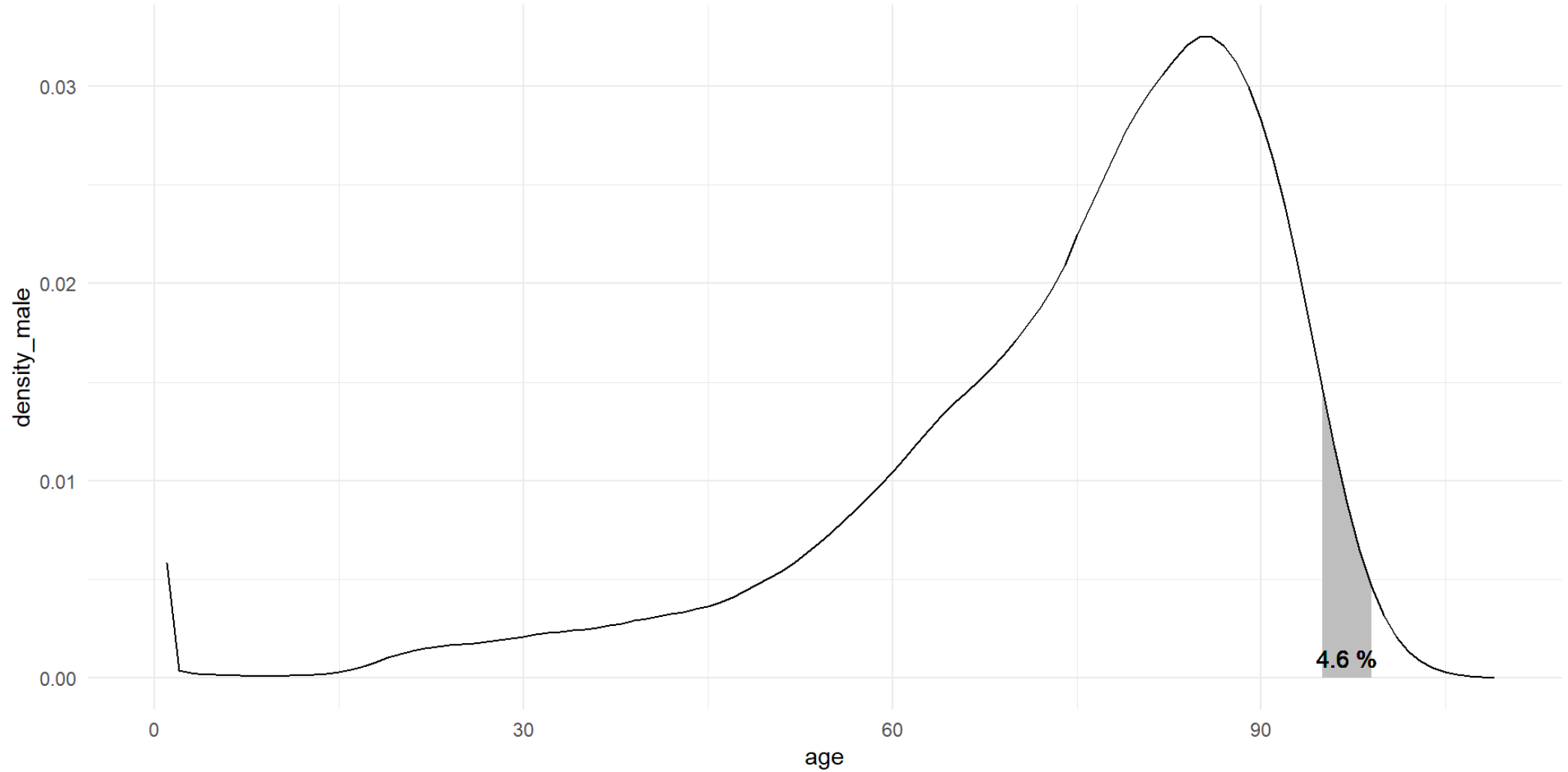
Probabilities for ages 21 through 41



Speaker notes

The probability of a potential customer dying between the ages of 21 and 41 is 0.04638.

Probabilities for ages 95 through 99



Speaker notes

The probability of a potential customer dying between the ages of 95 and 99 is about the same, 0.04626. So should you charge the same amount for an insurance policy for someone 21 years old and someone 95 years old?

Why are these probabilities not comparable?

- Unequal time intervals
 - Fix by computing a rate
- Non-uniform probabilities over the interval
 - Fix by looking at narrow interval
- No adjustment for survivorship
 - Fix by dividing by survival probability

Speaker notes

Obviously not. There are three things you need to fix first.

The most obvious flaw is the unequal time intervals, 20 years for the first probability and 4 years for the second probability.

You can fix this by computing a rate. You get the rate by dividing the probability by the width of the time interval.

The second flaw is that the probability changes over the interval, increasing in the first case and decreasing in the second case.

You can fix this by shrinking the width of the time interval.

The third flaw is a bit more subtle. The probability of dying between the ages of 95 and 99 are probabilities computed from the perspective of a newborn child. That probability is small not because the chances of dying are small at that age, but because so many have died before their 95th birthday.

If you are in insurance sales, you do not sell policies to newborn infants. You sell to people who have survived to a certain age. No one rises from their grave on their 95th birthday and asks for an insurance policy. First, because zombies aren't real, and second the zombie who died prior to year 95 would not be able to collect on an insurance policy that paid off for a death between 95 and 99.

You can fix this by dividing by the survivor probability.

Hazard function, definition

-

-

$$h(t) = \frac{f(t)}{S(t)}$$

- where f is the density function, and

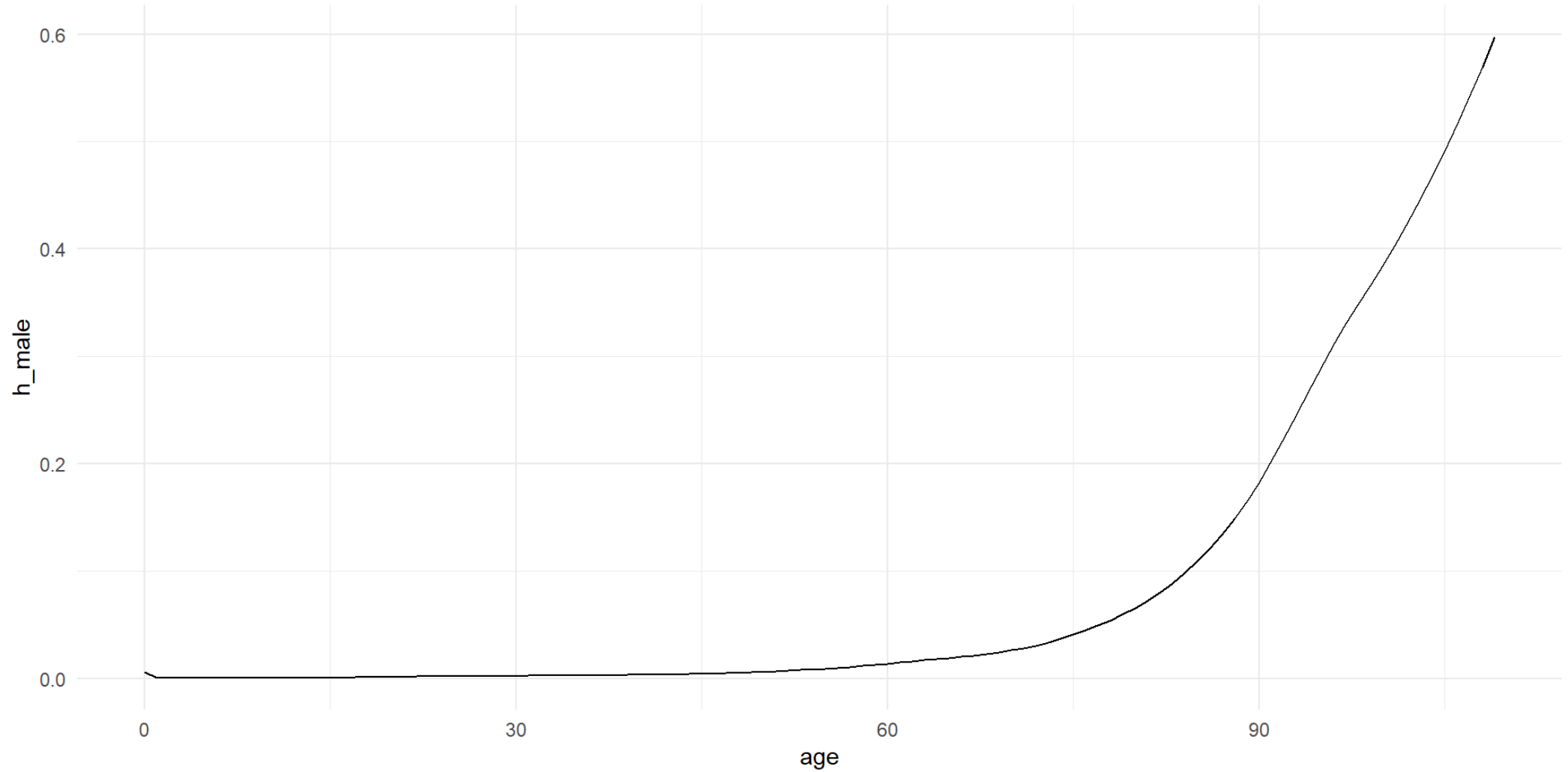
- S is the survival function ()

$$S(t) = 1 - F(t)$$

Speaker notes

The hazard function addresses all three of the concerns mentioned above. It computes a rate by dividing by Δt . It shrinks the interval but using a limit. And it adjusts for survivorship by dividing by the survivor probability.

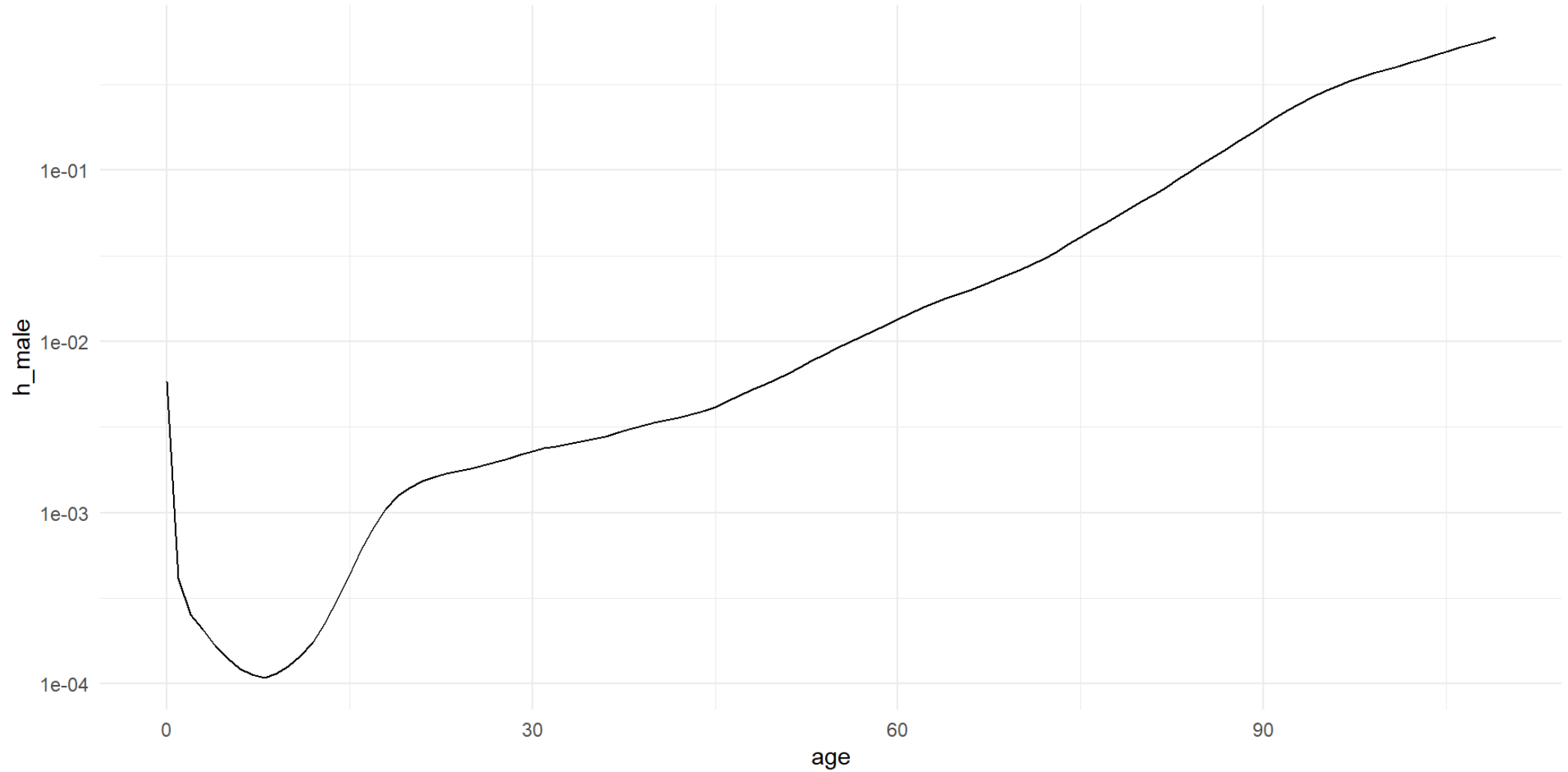
Hazard function, example



Speaker notes

This is what the hazard function for mortality data looks like.

Hazard function on a log scale



Speaker notes

The pattern becomes a bit clearer when you look at the hazard function on a log scale. The risk of death is high early in your life, but drops. There is a safe period during your pre-teen and early teen years, but then the risk rises because of an increase in deaths associated with things like driving, alcohol, and other drugs. Some of that fades as you mature but other risks increase because of the unavoidable aging of your body.

Break #4

- What you have learned
 - The hazard function
- What's coming next
 - The Cox regression model

Mean ages for men and women

Report

age			
gender	Mean	N	Std. Deviation
Male	66.60	300	14.943
Female	74.72	200	12.301
Total	69.85	500	14.491

Unadjusted and adjusted Cox regression models for gender

Variables in the Equation

	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
							Lower	Upper
gender	.381	.138	7.679	1	.006	1.464	1.118	1.917

Variables in the Equation

	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
							Lower	Upper
gender	-.066	.141	.217	1	.641	.937	.711	1.234
age	.067	.006	116.401	1	<.001	1.069	1.056	1.082

Live demo, The Cox regression model

Break #5

- What you have learned
 - The Cox regression model
- What's coming next
 - Assumptions and data management

Assumptions of the log rank test

- Independence
 - From one patient to another
 - Of censoring mechanism

Speaker notes

There is only one assumption in the log rank test, independence. There are, however, two dimensions of independence.

First, you have to assume that the chances of an event occurring for one patient is independent of the chances of that event for another patient. There is a famous example where this was not true. Carrie Fisher, an actor and author, died on December 27, 2016. Her mother, Debbie Reynolds, herself a famous actor, died on December 28, 2016. The two events were probably related, the stress of dealing with Carrie Fisher's death contributed to the death of Debbie Reynolds.

Often you can only assess the assumption of independence qualitatively, using your knowledge of the setting. Certain situations might make you concerned about the independence assumption, an infectious disease for example. If your data has clusters (family units, litter mates), then independence within a cluster might be questionable.

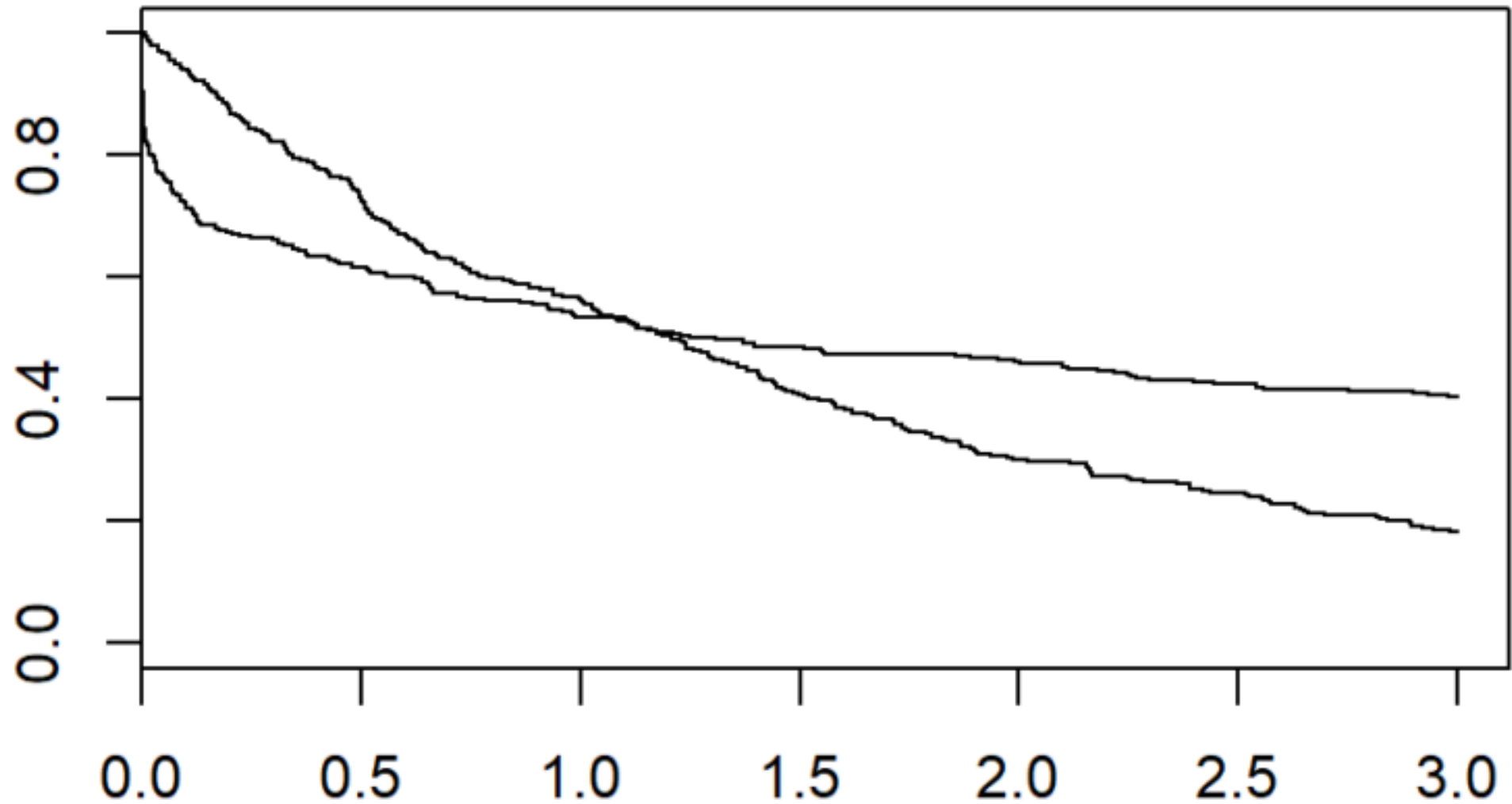
You also have to assume that when an event is censored, it is for a reason unrelated to what that event time might have been.

A patient may drop out of a study, for example, and you no longer know if or when the event occurred, except that it had to occur sometime after the date they dropped out. If a cancer patient drops out of a study because they went down to Mexico for an laetrile treatments, that might cause you to question the independence assumption. No one makes the trip to Mexico for an unapproved and unproven treatment if their current therapy was going well. The act of dropping out is informative as it indicates, at least in some cases, that death, relapse, or any other bad event is likely to happen sooner rather than later.

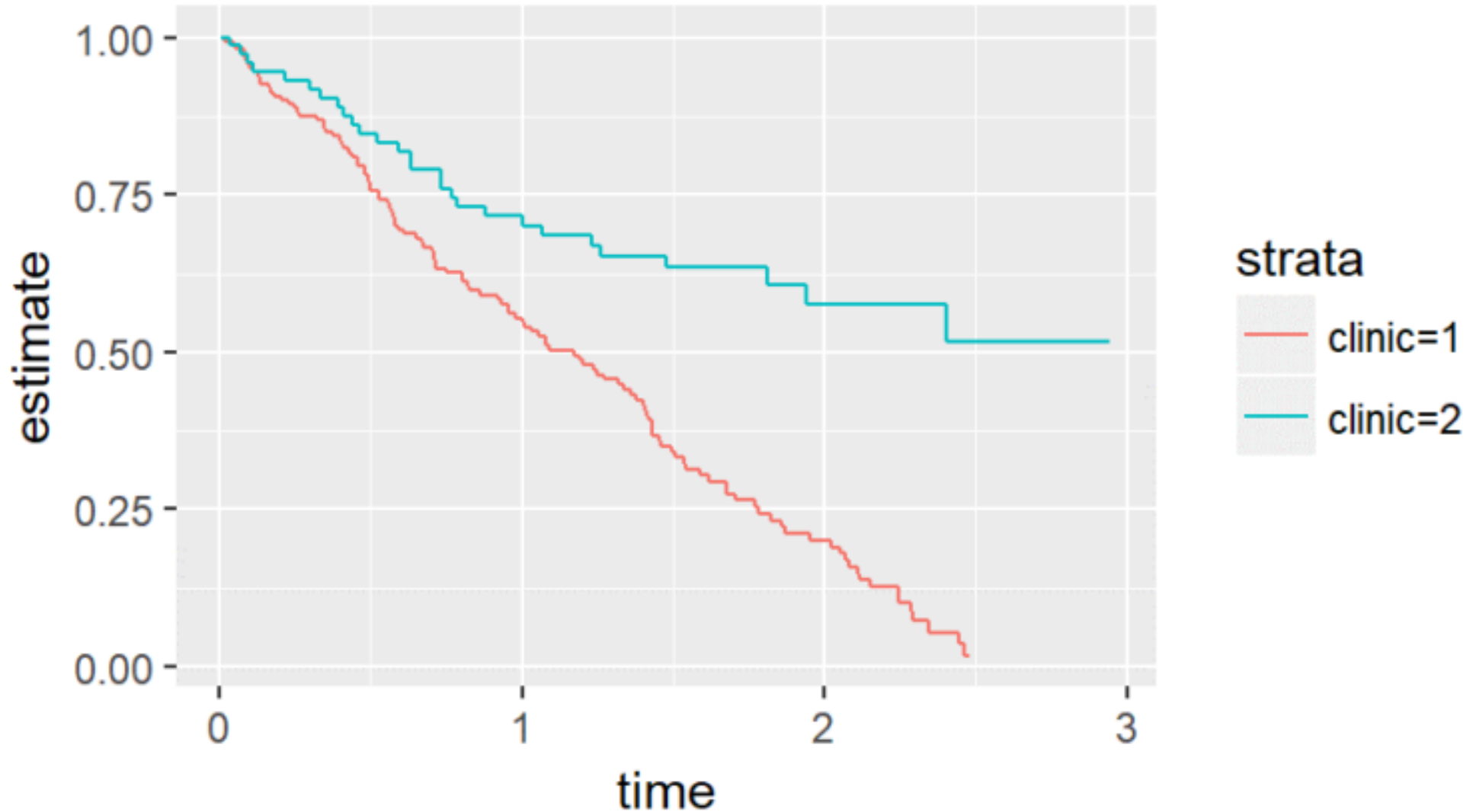
Assumptions of the Cox regression model

- Independence
- Proportional hazards assumption
- Possible violations of proportional hazards
 - Survival curves that cross
 - One curve flattening out over time
 - Curves diverge only at later times

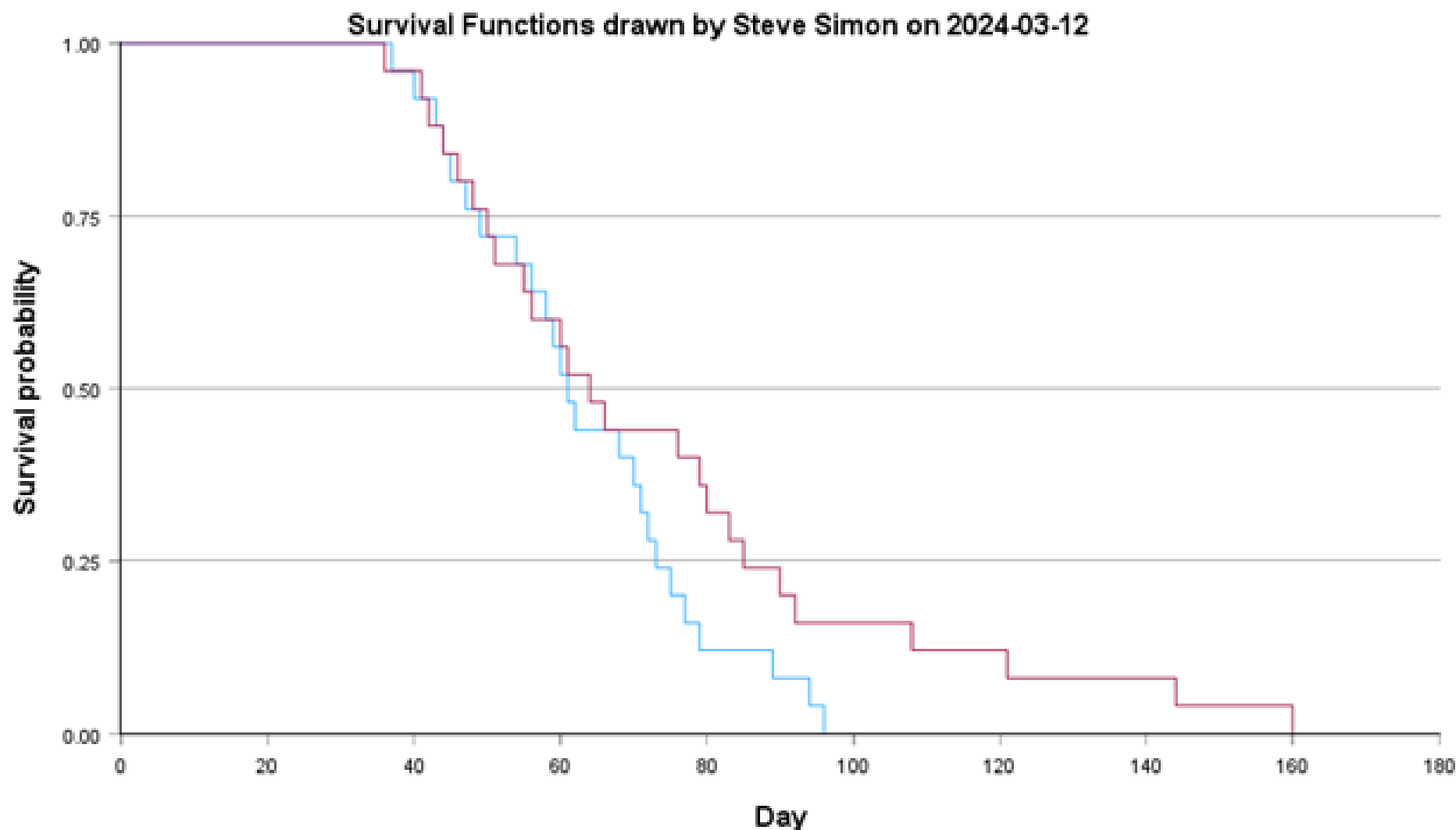
Survival curves that cross



One curve flattening out over time



Curves diverge only at later times



Sample size issues

- Rule of 50
- Rule of 15

Use ISO format for dates


PUBLIC SERVICE ANNOUNCEMENT:

OUR DIFFERENT WAYS OF WRITING DATES AS NUMBERS CAN LEAD TO ONLINE CONFUSION. THAT'S WHY IN 1988 ISO SET A GLOBAL STANDARD NUMERIC DATE FORMAT.

THIS IS *THE* CORRECT WAY TO WRITE NUMERIC DATES:

2013-02-27

THE FOLLOWING FORMATS ARE THEREFORE DISCOURAGED:

02/27/2013 02/27/13 27/02/2013 27/02/13
20130227 2013.02.27 27.02.13 27-02-13
27.2.13 2013.II.27. $27\frac{1}{2}$ -13 2013.158904109
MMXIII-II-XXVII MMXIII $\frac{\text{LVII}}{\text{CCCLXV}}$ 1330300800
 $((3+3) \times (111+1) - 1) \times 3 / 3 - 1 / 3^3$ 2013
10/11011/1101 02/27/20/13 $\begin{matrix} 2 & 3 & 1 & 4 \\ 0 & 1 & 2 & 3 & 7 \\ & 5 & 6 & 7 & 8 \end{matrix}$ 

Understand the internal storage system for dates

The screenshot shows the IBM SPSS Statistics help page for 'Dates and Times in IBM SPSS Statistics'. The page has a dark blue header with navigation links: 'Home > SPSS Statistics 24.0.0 > ... > Date and Time Wizard >'. On the right, there are 'Previous' and 'Next' buttons. The main title 'Dates and Times in IBM SPSS Statistics' is prominently displayed. Below the title, there are links for 'Table of contents' and 'Change version'. A search bar is located on the right side of the header. The main content area is white and contains a paragraph explaining that date/time variables are numeric with specific display formats. It distinguishes between date variables (e.g., mm/dd/yyyy) and date/time variables (e.g., dd-mmm-yyyy hh:mm:ss). A section titled 'Date and date/time variables' explains that these are stored as seconds from October 14, 1582. A bulleted list notes that both two- and four-digit year specifications are recognized. On the right side of the page, there is a 'More topics' section with links to related topics like 'Create a Date/Time Variable from a String' and 'Create a Date/Time Variable from a Set of Variables'. At the bottom right, there are icons for 'Print this topic', 'PDF download page [Beta]', and 'PDF download section', along with a 'Feedback' button.

Home > SPSS Statistics 24.0.0 > ... > Date and Time Wizard > Previous Next

Dates and Times in IBM SPSS Statistics

Search in all products

Table of contents Change version

Search in this pro...

Variables that represent dates and times in IBM® SPSS® Statistics have a variable type of numeric, with display formats that correspond to the specific date/time formats. These variables are generally referred to as date/time variables. Date/time variables that actually represent dates are distinguished from those that represent a time duration that is independent of any date, such as 20 hours, 10 minutes, and 15 seconds. The latter are referred to as duration variables and the former as date or date/time variables. For a complete list of display formats, see [Date and Time Formats](#).

> **Date and date/time variables.** Date variables have a format representing a date, such as mm/dd/yyyy. Date/time variables have a format representing a date and time, such as dd-mmm-yyyy hh:mm:ss. Internally, date and date/time variables are stored as the number of seconds from October 14, 1582. Date and date/time variables are sometimes referred to as date-format variables.

- Both two- and four-digit year specifications are recognized. By default, two-digit years represent a range beginning 69 years prior to the current date and ending 30 years after the current date. This range is determined by your Options settings and is configurable (from the Edit menu, choose **Options** and click the **Data** tab).

More topics

- [Create a Date/Time Variable from a String](#)
- [Create a Date/Time Variable from a Set of Variables](#)
- [Add or Subtract Values from Date/Time Variables](#)
- [Extract Part of a Date/Time Variable](#)

Print this topic

PDF download page [Beta]

PDF download section

Feedback

Date management

- The three dates you need
 - the date of origin,
 - the date of the event (if it occurred),
 - the date of last contact with the patient.

The date of origin

- Rehospitalization
 - use date of first discharge.
- Failure of a mechanical device
 - use date of implant.
- Divorce
 - use date of marriage.
- Loan default
 - use date of loan contract.
- Infectious disease
 - use date of first exposure.

Speaker notes

You have various choices for the date of origin. It depends a lot on the context of the research.

The date of the event

- Define your event precisely
 - All-cause mortality
 - Mortality related to the health condition
 - Composite endpoints (e.g., death or relapse)
 - Requires comparing the earlier of two dates.
- If the event did NOT happen, leave this field blank/missing.

The date of last contact

- If event did not occur
 - Must be specified
 - Typically last medical exam or last telephone contact
- If event did occur
 - Make same as event date, or
 - Leave blank

Speaker notes

Every patient will have a date of last contact. It will be the last time that you have been able to contact the patient and confirm that the event has not yet occurred.

If the event has occurred, you have several reasonable choices: put the event date in this field also, or leave the field blank/missing.

A date after the event date may represent a data error!

Survival calculations, 1 of 2

- $\text{Time} = \max(\text{Date of event}, \text{Date of last contact}) - \text{Date of origin}$
- Censoring variable = 0 if Date of event is missing, 1 if not

Survival calculations, 2 of 2

patient	origin	event	last_contact	time	censored		
1	2024-01-02		2024-01-30	28	0		
2	2024-01-04	2024-01-12		8	1		
3	2024-01-06		2024-01-30	24	0		
4	2024-01-08		2024-01-30	22	0		
5	2024-01-10		2024-01-30	20	0		
6	2024-01-12		2024-01-30	18	0		
7	2024-01-14	2024-01-15		1	1		
8	2024-01-16	2024-01-28		12	1		
9	2024-01-18	2024-01-25		7	1		

Summary

- What you have learned
 - A simple example of survival data
 - Overall Kaplan-Meier curve
 - The log rank test
 - The hazard function
 - The Cox regression model
 - Assumptions and data management

