

BIOS522_S8
des1*BIOS 522: Survival Analysis Methods***Lecture 1:****Introduction to time-to-event data****Welcome to BIOS 522!**

- Survival analysis is the branch of statistics that deals with times to events. It has many important applications in clinical research and epidemiology. The goal of this course is to give you a solid understanding of survival analysis and its applications.

2

Your instructor

- Dr. Natalie Dean
 - Assistant Professor in the BIOS Department
 - Research on emerging infectious diseases, vaccine study design
- Office: GCR 336
- Contact: nataliedean@emory.edu or Canvas

3

Your teaching assistant

- Emily Wu
- Contact: emily.wu2@emory.edu or Canvas

4

Course structure

- Before class
 - Read the lecture notes
 - Take a short Canvas quiz
 - One other pre-class responsibility (e.g. homework, discussion, review computing handout)
- In-class
 - Non-exhaustive review of concepts... building content!
 - Examples from the literature
 - Small group activities
- *Active learning!*

5

Review syllabus

Final
is
take-home
project
w/ short answer

6

Course philosophy

- Emphasis on literacy, understanding, and context
- Building a map of concepts >>>> Covering every detail
- Participation >>>> perfection
- *I am glad you are joining me this semester to learn about this important topic. Your feedback throughout the semester is valued. Please do not hesitate to contact me with questions or concerns.*

7

Today's learning objectives

- *Identify examples of time-to-event analyses in practice*
- *Define the time-to-event data format*
- *Define right censoring and censoring notation*
- *Identify the time origin, event, or time scale from an example*
- *Diagnose the limitations of analyzing data as continuous, binary, or as incidence rates*

8

Survival analysis

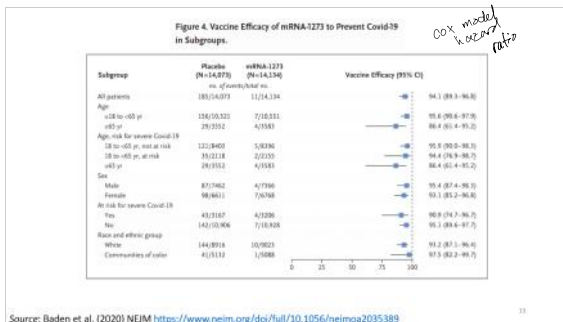
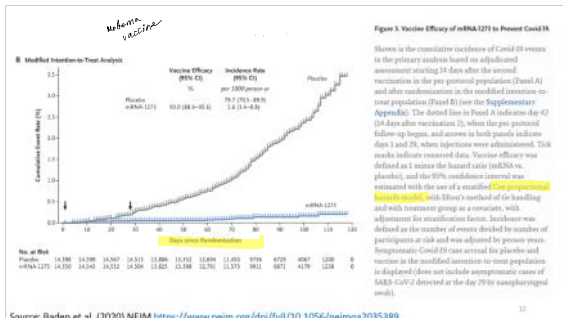
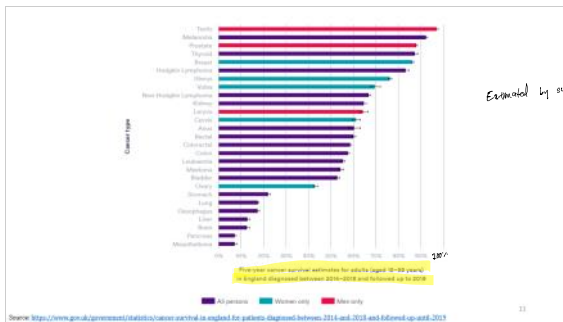
- Survival analysis is the branch of statistics that deals with times to events.
 - Survival time, failure time, occurrence time, event time, time-to-event
- Examples of survival times in clinical research:
 - Time from diagnosis until death
 - Time from infection until diagnosis
 - Time from treatment until suppression of symptoms
 - Length of stay in a hospital

9

Survival analysis examples

- Survival analysis in action

10



Survival analysis overview

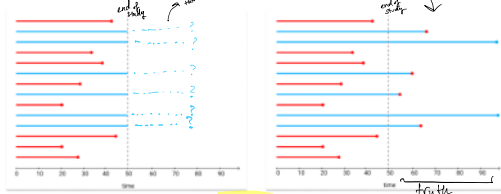
In order to analyze time-to-event data, it is necessary to define:

1. The **time origin**: the **beginning** of the survival time
2. The **failure time**: the **end** of the survival time

Survival time T measures the time elapsed from the origin ("time zero") until the event of interest

Setting	Time origin	Event	Time scale
Human mortality	Birth	Death	Age
Clinical trial of treatment	Randomization	Stroke or cardiovascular death	Time since start of treatment
Pregnancy cohort	12 weeks gestation	Fetal death	Gestational age
Hospital study	Admission	Discharge	Time in hospital
Surgical study	Surgery	Death or complication	Time since surgery
Cancer cohort	Diagnosis	Tumor recurrence	Time since diagnosis
Ebola survival study	Date of symptom onset	Death due to Ebola	Time since symptom onset
Influenza study	Start of flu season (October 1, 2010)	Influenza symptom onset	Calendar time

Not everyone experiences the event during follow-up

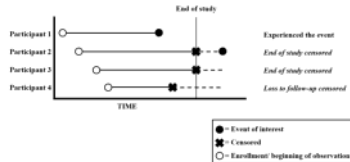


Right censoring

censoring to the right of the line

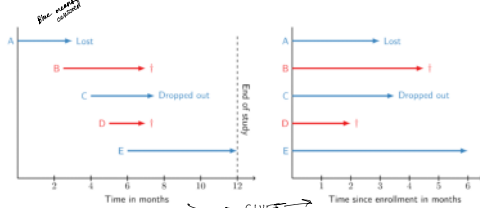
Source: <https://tiny.cc/meywmg>

Some may be censored earlier Enrollment may be rolling over time



17

Align participants at common time origin



Source: <https://tiny.cc/meywmg>

Censoring notation

- We imagine that everyone in our study has two random variables:
 - T is their **failure time**
 - C is their **censoring time** (end of study, date when they will move)

- In practice, we only observe whichever comes first
- We introduce new notation for the observed data
 - $T^o = \min(T, C)$ is the observed time (failure or censoring)
 - $\delta = I(T \leq C)$ is an indicator (0 or 1) for observing a failure time

Underlying data	Observed data
T_i	T_i^o
3	3
6	1
5	5
6	1
8	6
	0

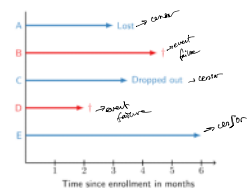
truth

what we see

for observed time of failure or censoring first
if not event or failure or censoring

18

Censoring notation



ID	T_i^o	δ_i
A	3	0
B	4.5	1
C	3.5	0
D	2	1
E	6	0

Using shorthand:

3+, 4.5, 3.5+, 2, 6+

2, 3+, 3.5+, 4.5, 6+

⊕ means censoring time

19

So what do survival data look like?

Patient ID	Date of surgery	Date of death	Date of censoring
1	January 1, 2020		February 1, 2020
2	January 5, 2020	January 13, 2020	
3	January 7, 2020		February 7, 2020
4	January 10, 2020		February 10, 2020
5	January 20, 2020	January 31, 2020	
6	January 22, 2020		February 22, 2020

Patient ID	Time to event	Event indicator (1=death, 0=censoring)
1	31 days	0
2	8 days	1
3	31 days	0
4	31 days	0
5	11 days	1
6	31 days	0

Using shorthand:
31+, 8, 31+, 31+, 11, 31+

Using shorthand, sorted:
8, 11, 31+, 31+, 31+, 31+

22

Example: Worcester Heart Attack Study

id	age	gender	admission	date	status
1	70	1	01/01/1997	12/31/2002	0
174	72	1	01/02/1999	08/11/2002	1
163	64	1	01/03/1999	12/27/1999	1
280	63	0	01/03/1999	12/31/2002	0
20	73	0	01/03/1997	03/01/1997	1
322	61	0	01/03/1999	12/31/2002	0
498	57	1	01/06/2001	12/31/2002	0
36	68	0	01/07/1997	08/11/1998	1
273	47	0	01/07/1999	12/31/2002	0
871	66	1	01/07/1999	12/31/2002	0
168	70	1	01/08/1999	01/10/1999	1
162	66	1	01/12/1999	12/31/2002	0
6	66	0	01/11/1997	08/10/2001	1
177	60	1	01/11/1999	01/10/1999	1
164	66	1	01/12/1999	04/21/2000	1
171	69	0	01/12/1999	12/31/2002	0
1	69	0	01/12/1997	12/31/2002	0
400	77	0	01/15/2001	04/27/2002	1

→ read
admission
for
date
in days

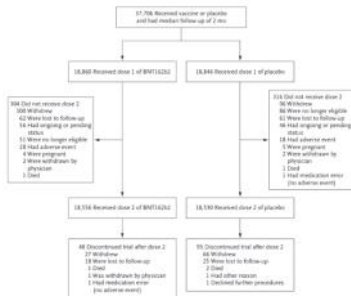
23

Source: Hosmer, Lemeshow, May (2008) John Wiley and Sons Inc.

Other reasons for censoring

- Lost to follow-up
- Move out of the study region
- No longer meet study eligibility criteria
- No longer "at-risk" (e.g., death from unrelated cause)

24



25

Source: Pollock et al. (2020) N Engl J Med 383(25):2543-57

Why survival analysis?

- Why do we need specialized methods for time-to-event data when we have methods for analyzing:
 - Continuous data
 - Binary data
 - Incidence rate data (person-time)

26

Hypothetical study

- We design a study to estimate survival for women diagnosed with stage II breast cancer.
- We can imagine forming a cohort of newly diagnosed women and following them prospectively in time to track survival outcomes.

26

Continuous data

- Imagine that our cohort was comprised of **elderly women** (≥ 75 years), and we continued our study for 20+ years until all women had died. For each woman we can calculate the event time, which is the **time elapsed from diagnosis to death**.
- We can report the **mean or median survival time**.
- What if our cohort was comprised of **middle-aged women** tracked for 5 years only?
- *Can we calculate the mean or median survival time?*

27

Why not continuous data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 3 survive to the end of the 5-year study (survival time is "**censored**")
- Calculate mean time:

$$\frac{1}{n} \sum_{i=1}^n T_i = \frac{3 + 4 + 5 + 5 + 5}{5} = 4.4 \text{ years}$$

MEAN FOLLOW-UP TIME
not mean diagnosis with death

28

Why not continuous data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 3 survive to the end of the 5-year study (survival time is "**censored**")
- Calculate mean SURVIVAL time:

$$\frac{1}{n} \sum_{i=1}^n T_i = \frac{3 + 4 + ? + ? + ?}{5} = ?? \text{ years}$$

29

Why not continuous data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 3 survive to the end of the 5-year study (survival time is "**censored**")
- Median survival time: $(3, 4, ?, ?, ?) = ??$

30

Why not continuous data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 1 woman dies 4.5 years after diagnosis
 - 2 survive to the end of the 5-year study (survival time is "censored")

- Median survival time: $(3, 4, 4.5, ?, ?) = 4.5$ years

these times are > 5

Admission censored =
censoring at end
of study

33

Why not continuous data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 1 woman moves away 2 years after diagnosis (censored)
 - 3 survive to the end of the 5-year study (survival time is "censored")

- Median survival time: $(?, 3, 4, ?, ?) = ??$

32

Why not binary data?

- Another data summary is **5-year mortality** (yes/no, binary)

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 3 survive to the end of the 5-year study (survival time is censored)

- 5-year mortality

$$\frac{1}{n} \sum_{i=1}^n I[T_i \leq 5] = \frac{2}{5}$$

probability
that they
will die
5 years before

35

Why not binary data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis $T_i = 3 \rightarrow T_i^* = 3, \delta_i = 1$
 - 1 woman dies 4 years after diagnosis $T_i = 4 \rightarrow T_i^* = 4, \delta_i = 1$
 - 3 survive to the end of the 5-year study $C_i = 5 \rightarrow T_i^* = 5, \delta_i = 0$

- For the 3 women censored at 5 years, we know $T_i > 5$

34

Why not binary data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - 1 woman moves away at 2 years after diagnosis
 - 2 survive to the end of the 5-year study (survival time is censored)

35

Why not binary data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis $T_i = 3 \rightarrow T_i^* = 3, \delta_i = 1$
 - 1 woman dies 4 years after diagnosis $T_i = 4 \rightarrow T_i^* = 3, \delta_i = 1$
 - **1 woman moves away at 2 years after diagnosis $C_i = 2 \rightarrow T_i^* = 2, \delta_i = 0$**
 - 2 survive to the end of the 5-year study $C_i = 5 \rightarrow T_i^* = 5, \delta_i = 0$

- 5-year mortality

$$\frac{1}{n} \sum_{i=1}^n I[T_i \leq 5] = \frac{?}{5}$$

36

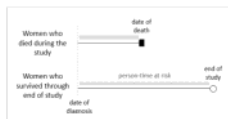
Why not binary data?

- Necessary to pick a single time point for analysis
- Is 10 years more important than 5 years?

37

Why not incidence rate data?

- A **rate-based analysis** can be used when there are differing lengths of follow-up.
- For each woman, we can calculate the length of time during which the event could have occurred and would have been counted in the population, known as the **person-time**.



38

Why not incidence rate data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis
 - 1 woman dies 4 years after diagnosis
 - **3 survive to the end of the 5-year study (survival time is censored)**

Yearly mortality rate

$$\text{incidence rate} = \frac{\# \text{ of women who died in our study}}{\text{total person-years at risk after diagnosis}}$$

39

Why not incidence rate data?

- $n = 5$ women
 - 1 woman dies 3 years after diagnosis $T_i = 3 \rightarrow T_i^* = 3, \delta_i = 1$
 - 1 woman dies 4 years after diagnosis $T_i = 4 \rightarrow T_i^* = 3, \delta_i = 1$
 - 3 survive to the end of the 5-year study $C_i = 5 \rightarrow T_i^* = 5, \delta_i = 0$

Yearly mortality rate

$$\text{incidence rate} = \frac{\sum_{i=1}^n \delta_i}{\sum_{i=1}^n T_i}$$

$$= \frac{3 + 4 + 5 + 5 + 5}{2}$$

$$= 0.09 \text{ deaths/person-year}$$

but what if you wanted it to

vary over time?

40

Why not incidence rate data?

- **Rate-based analyses** are related to simple survival analysis methods.
- **Rate-based analyses assume that the event rate is constant.**
- The majority of survival analysis methods that we will learn about in this class **allow the event rate to vary over time.**
- For example, mortality rates for women newly diagnosed with breast cancer may be **initially high**, as some women may have aggressive or difficult-to-treat forms. Women who survive more than 5 years after diagnosis, though, may have mortality rates closer to the general population.

Why not incidence rate data?

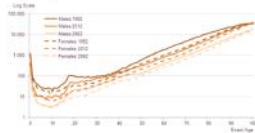
- **Rate-based analyses** are related to simple survival analysis methods.
- **Rate-based analyses assume that the event rate is constant.**
- The majority of survival analysis methods that we will learn about in this class **allow the event rate to vary over time.**
- For example, mortality rates for women newly diagnosed with breast cancer may be **initially high**, as some women may have aggressive or difficult-to-treat forms. Women who survive more than 5 years after diagnosis, though, may have mortality rates closer to the general population.

42

Human mortality

Figure 3: 2013-based Period Mortality Rates (per 100,000), United Kingdom, 1962-2012, 2062

Principal Projection



Source: Office for National Statistics

43

Today's activity

- Word problems
- Work in assigned small groups

44