

# History and Ecology of R

Martyn Plummer

University of Warwick, UK

SPE 2023, Tartu



Pre-history  
ooooooo

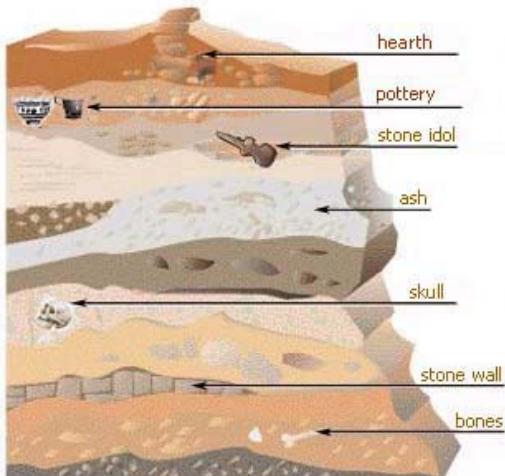
History  
ooooooo

Present  
ooooooo

Future?  
oooo

## Abstract

more recent  
↑  
more ancient



Pre-history  
●ooooooo

History  
ooooooo

Present  
ooooooo

Future?  
oooo

## Pre-history



## The S language

Developed at AT&T Bell laboratories by Rick Becker, John Chambers, Doug Dunn, Paul Tukey, Graham Wilkinson.

Version 1	1976–1980	Honeywell GCOS, Fortran-based
Version 2	1980–1988	Unix; Macros, Interface Language
	1981–1986	QPE (Quantitative Programming Environment)
	1984–	General outside licensing; books
Version 3	1988–1998	C-based; S functions and objects
	1991–	Statistical models; informal classes and methods
Version 4	1998	Formal class-method model; connections; large objects
	1991–	Interfaces to Java, Corba?

Source: Stages in the Evolution of S <http://ect.bell-labs.com/s1/S/history.html>

## The “Blue Book” and the “White Book”



Key features of S version 3 outlined in two books:

- Becker, Chambers and Wilks, *The New S Language: A Programming Environment for Statistical Analysis and Graphics* (1988)
  - Functions and objects
- Chambers and Hastie (Eds), *Statistical Models in S* (1992)
  - Data frames, formulae

These books were later used as a prototype for R.

## Programming with Data

*“We wanted users to be able to begin in an interactive environment, where they did not consciously think of themselves as programming. Then as their needs became clearer and their sophistication increased, they should be able to slide gradually into programming.” – John Chambers, Stages in the Evolution of S*

This philosophy was later articulated explicitly in *Programming With Data* (Chambers, 1998) as a kind of mission statement for S

*To turn ideas into software, quickly and faithfully*

Pre-history  
oooo●ooo

History  
oooooooo

Present  
oooooooo

Future?  
oooo

## The “Green Book”



Key features of S version 4 were outlined in Chambers, *Programming with Data* (1998).

- S as a programming language
- Introduced formal classes and methods, which were later introduced into R by John Chambers himself.

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

Pre-history  
ooooooo●

History  
oooooooo

Present  
oooooooo

Future?  
oooo

## S-PLUS

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

Pre-history  
ooooooo●

History  
oooooooo

Present  
oooooooo

Future?  
oooo

## The Rise and Fall of S-PLUS

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

Pre-history  
ooooooo

History  
●oooooo

Present  
ooooooo

Future?  
oooo

## History

*How R started, and how it turned into an S clone*

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

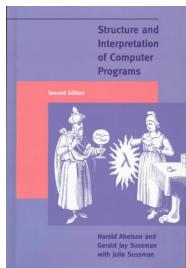
Pre-history  
ooooooo

History  
○●ooooo

Present  
ooooooo

Future?  
oooo

## The Dawn of R



- Ross Ihaka and Robert Gentleman at the University of Auckland
- An experimental statistical environment
- Scheme interpreter with S-like syntax
  - Replaced scalar type with vector-based types of S
  - Added lazy evaluation of function arguments
- Announced to *s-news* mailing list in August 1993.

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

Pre-history  
ooooooo

History  
○○●oooo

Present  
ooooooo

Future?  
oooo

## A free software project

- June 1995. Martin Maechler (ETH, Zurich) persuades Ross and Robert to release R under GNU Public License (GPL)
- March 1996. Mailing list *r-testers* mailing list
  - Later split into three *r-announce*, *r-help*, and *r-devel*.
- Mid 1997. Creation of *core team* with access to central repository (CVS)
  - Doug Bates, Peter Dalgaard, Robert Gentleman, Kurt Hornik, Ross Ihaka, Friedrich Leisch, Thomas Lumley, Martin Maechler, Paul Murrell, Heiner Schwarte, Luke Tierney
- 1997. Adopted by the GNU Project as “GNU S”.

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿

Pre-history  
ooooooo

History  
ooooooo

Present  
ooooooo

Future?  
oooo

## The draw of S

*"Early on, the decision was made to use S-like syntax. Once that decision was made, the move toward being more and more like S has been irresistible"*  
– Ross Ihaka, R: Past and Future History (Interface '98)

R 1.0.0, a complete and stable implementation of S version 3, was released in 2000.



Pre-history  
ooooooo

History  
ooooooo

Present  
ooooooo

Future?  
oooo

## A Souvenir



Pre-history  
ooooooo

History  
ooooooo

Present  
ooooooo

Future?  
oooo

## Packages

- Comprehensive R Archive Network (CRAN) started in 1997
  - Quality assurance tools built into R
  - Increasingly demanding with each new R release
- Recommended packages distributed with R
  - Third-party packages included with R distribution
  - Provide more complete functionality for the R environment
  - Starting with release 1.3.0 (completely integrated in 1.6.0)



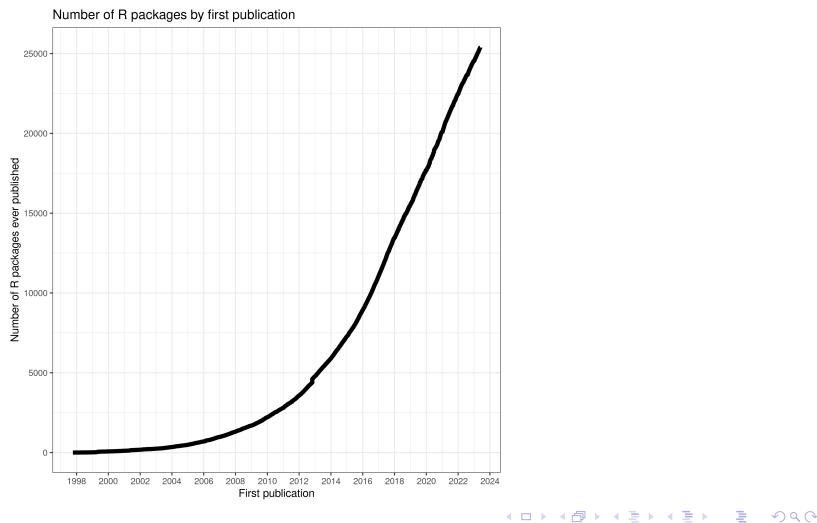
Pre-history  
ooooooo

History  
oooooo●

Present  
oooooooo

Future?  
oooo

## Growth of CRAN



Pre-history  
ooooooo

History  
oooooo●

Present  
●oooooooo

Future?  
oooo

## Community

- useR! Annual conference
  - Toulouse (2019), Online (2020, 2021), Nashville (2022)
- R Journal (<http://journal.r-project.org>)
  - Journal of record, peer-reviewed articles, indexed
  - Journal of Statistical Software (JSS) has many articles dedicated to R packages (<http://jstatsoft.org>)
- Migration to social media
  - Stack Exchange/Overflow, Github, Twitter, Mastodon (#rstats)
  - Follow @R\_Foundation on Twitter, or @R\_Foundation@fosstodon.org on Mastodon

Pre-history  
ooooooo

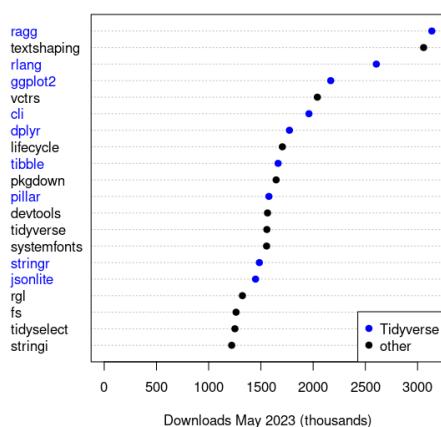
History  
oooooo●

Present  
●oooooooo

Future?  
oooo

## Much important R infrastructure is now in package space

Top 20 packages by downloads



Pre-history  
ooooooo

History  
oooooooo

Present  
ooo●ooo

Future?  
oooo

## Much important R infrastructure is now in package space

Top 100 packages by downloads



◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾

Pre-history  
ooooooo

History  
oooooooo

Present  
ooo●ooo

Future?  
oooo

## The tidyverse

- Many of the popular packages on CRAN come from the company Posit (formerly R Studio).
- These packages are known as the “tidyverse” ([www.tidyverse.org](http://www.tidyverse.org)).
- All packages in the tidyverse have a common design philosophy and work together. Common features are:
  - Non-standard evaluation rules for function calls.
  - Use of the pipe operator `|>` (or `%>%`) to pass data transparently from one function call to another.
- The CRAN meta-package `tidyverse` installs all of these packages.

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾

Pre-history  
ooooooo

History  
oooooooo

Present  
ooo●ooo

Future?  
oooo

## The R Foundation for Statistical Computing

A non-profit organization working in the public interest, founded in 2002 in order to:

- Provide support for the R project and other innovations in statistical computing.
- Provide a reference point for individuals, institutions or commercial enterprises that want to support or interact with the R development community.
- Hold and administer the copyright of R software and documentation (This never happened)

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾

Pre-history  
ooooooo

History  
oooooooo

Present  
oooooo●

Future?  
oooo

## The R Consortium

In 2015, a group of organizations created a consortium to support the R ecosystem.  
Current members (May 2023)

**R Foundation** A statutory member of The R Consortium

**Platinum members** Biogen, Genentech, Microsoft, Posit

**Gold members** ASA, Esri, Google, GSK, Janssen, Lander Analytics, Mango Solutions, Merck

**Silver members** Novo Nordisk, Oracle, Parexel, Pfizer, Procogia, Swiss Re



Pre-history  
ooooooo

History  
oooooooo

Present  
oooooo

Future?  
●ooo

## The Future



Pre-history  
ooooooo

History  
oooooooo

Present  
oooooo

Future?  
●ooo

## Trends

We cannot make predictions, but some long-term trends are very visible:

- Average age of R Core Team?
- Younger R developers more closely associated with industry than academia
- Strong competition from Python



Pre-history  
ooooooo

History  
oooooooo

Present  
ooooooo

Future?  
ooo●○

## What does all of this mean for the course?

- R incorporates over 40 years of ideas in statistical computing from multiple contributors.
- There is usually more than one way to do something in R.
- Some of the peculiarities of the R language are there for historical reasons.
- The course does not cover some of the recent additions to the R ecosystem.



Pre-history  
ooooooo

History  
oooooooo

Present  
ooooooo

Future?  
ooo●○

## Resources

- Chambers J, Stages in the Evolution of S
- Becker, R, A Brief History of S
- Chambers R, Evolution of the S language
- Ihaka, R and Gentleman R, R: A language for Data Analysis and Graphics, *J Comp Graph Stat*, 5, 299–314, 1996.
- Ihaka, R, R: Past and Future History, Interface 98.
- Ihaka, R, Temple Lang, D, Back to the Future: Lisp as a Base for a Statistical Computing System
- Fox, J, Aspects of the Social Organization and Trajectory of the R Project, R Journal, Vol 1/2, 5–13, 2009.



## R: language and basic data management

Krista Fischer

Statistical Practice in Epidemiology, Tartu, 2019  
(initial slides by P. Dalgaard)

### Language

- ▶ R is a programming language – also on the command line
- ▶ (This means that there are *syntax rules*)

On the command line (or a line in a script) one could:

- ▶ Print an object by typing its name
- ▶ Evaluate an expression
- ▶ Call a function, giving the arguments in parentheses – possibly empty
- ▶ Notice `objects` vs. `objects()`

### R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object names
- ▶ Explicit constants
- ▶ Arithmetic operators
- ▶ Function calls
- ▶ Assignment of results to names

## R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object **names**
- ▶ Explicit constants
- ▶ Arithmetic operators
- ▶ Function calls
- ▶ Assignment of results to names

## R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object names
- ▶ Explicit constants
- ▶ Arithmetic operators
- ▶ Function calls
- ▶ Assignment of results to names

## R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object names
- ▶ Explicit constants
- ▶ Arithmetic **operators**
- ▶ Function calls
- ▶ Assignment of results to names

## R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object names
  - ▶ Explicit constants
  - ▶ Arithmetic operators
  - ▶ **Function calls**
  - ▶ Assignment of results to names

## R expressions

```
x <- rnorm(10, mean=20, sd=5)
m <- mean(x)
sum((x - m)^2)
```

- ▶ Object names
  - ▶ Explicit constants
  - ▶ Arithmetic operators
  - ▶ Function calls
  - ▶ **Assignment** of results to names

## Objects

- ▶ The simplest object type is *vector*
  - ▶ Modes: numeric, character, factor, ...
  - ▶ Operations are vectorized: you can add entire vectors with  
 $a + b$
  - ▶ Recycling of objects: If the lengths don't match, the shorter vector is reused

## Example (numeric vectors)

```

> a <- c(2, 8, 3, 1, 0, 7)
> b <- c(3, 4, 1, 4, 5, 2)
> a+b
[1] 5 12 4 5 5 9
> mean(a)
[1] 3.5
> m <- mean(a)
> m
[1] 3.5
> a - m # notice recycling
[1] -1.5 4.5 -0.5 -2.5 -3.5 3.5

> z <- c(1, 2, 3)
> a - z #recycling!
[1] 1 6 0 0 -2 4

```

## Factors

- ▶ **Factors** are used to describe groupings – these are just integer codes plus a set of names, as labels for the *levels*
  - ▶ In model specifications, a factor variable is treated as a classification rather than as a quantitative variable

### Example:

```
> x<-c(1,3,3,2,1,3,1)
> fx<-factor(x,labels=c("bad","average","good"))

> fx
[1] bad      good     good     average  bad      good

> levels(fx)
[1] "bad"    "average" "good"
```

# Lists

- ▶ Lists are vectors where the elements can have different types – thus collections of any elements, gathered into one object
  - ▶ Functions often return lists
  - ▶ 

```
lst <- list(A=rnorm(5), B="hello")
```
  - ▶ Special indexing:
  - ▶ `lst$A`
  - ▶ `lst[[1]]` first element (NB: double brackets)
  - ▶ **Data frames** are special type of lists

## Matrices

- ▶ A **matrix** is a rectangular collection of data. All columns of a matrix should be of the same type.

```

> A<-matrix(c(1, 4, 2, 6, 7, 8), nrow=3, ncol=2,
           byrow=T)
> A
      [,1]  [,2]
[1,]     1     4
[2,]     2     6
[3,]     7     8

```

- ▶ One can also construct a matrix from its columns using `cbind`, whereas joining two matrices with equal no of columns (with the same column names) can be done using `rbind`.

## Data frames

- ▶ Usually a dataset in R is stored in a form of a **data frame**.
  - ▶ While reading in data from text files (using `read.table()`, `read.csv()`), a data frame is created.
  - ▶ A data frame is similar to a matrix, but can have columns (variables) of different types.
  - ▶ A variable can be extracted using `dataframe$variable` (as data frames are lists)

```
(

> D<- data.frame(a=c(8,3,5),b=c("X","Z","Y"))
> D
  a b
1 8 X
2 3 Z
3 5 Y
> D$a
[1] 8 3 5
```

## Matrices or data frames?

- ▶ A (numeric or character) matrix can be converted to a data frame and vice versa (with `as.data.frame(A)` and `as.matrix(B)`).
  - ▶ Most R functions for statistical analysis work with data frames, but in some cases it is useful to have a matrix (incl the occasions where you want to use some matrix algebra).
  - ▶ If you need more dimensions than two, there is also `array`.

## How to access variables in the data frame?

Different ways to tell R to use variable X from data frame D:

- ▶ As mentioned, you can use the `dataframe$variable` notation  
`summary(D$X)`
  - ▶ Use the `with` function  
`with(D, summary(X))`
  - ▶ Use the `data` argument (does not work for all functions)  
`lm(Y~X, data=D)`
  - ▶ Attach the dataframe – **DISCOURAGED!**  
(seems a convenient solution, but can actually make things more complicated, as it creates a temporary copy of the dataset)  
`attach(D)`  
`summary(X)`  
`detach()`

Basics	Objects in R	Data frames and data manipulation	Functions	Graphics	The workspace	Additional topics
oo	oooooooo	ooooooo	oooo	oo	o	oooo

## Data manipulation

To create a new variable `bmi` in the existing data frame `students`, use either of the two:

```
students$bmi <-  
    with(students, weight/(height/100)^2)  
students <-  
    transform(students, bmi=weight/(height/100)^2)
```

(notice: you need an assignment, to save the transformed object)

Basics	Objects in R	Data frames and data manipulation	Functions	Graphics	The workspace	Additional topics
oo	oooooooo	oooo●ooo	oooo	oo	o	oooo

## Indexing – extracting elements from objects

Square brackets [ ] are used for indexing!

## Examples:

- ▶ Elements of vectors: `a[5]` (5th element); `a[5:7]` (5th to 7th elements); `a[-6]` (all elements except the 6th)
  - ▶ Logical index: `a[a < 3]`, `a[b > 2]`, `a[is.na(b)]` (elements of `a` corresponding to missing values of `b`)
  - ▶ In a data frame or matrix – two dimensions, two indexes:  
`students[5, 7]`, `students[1:10, c(2, 5)]`,  
`students[1, ]`, `students[, 3]` (entire row/column)

## Examples of indexing

```
> x<- c(2,7,3,1,5,9,0)
> x[c(1,5,7)]
[1] 2 5 0
> x[x<3]
[1] 2 1 0

> NMRimp[1:2,1:4]    #quick look at a large data
  sample.id XXL.VLDL.P XXL.VLDL.L XXL.VLDL.PL
1      V18566   1.46e-04    0.0313    0.00331
2      V36115   9.00e-05    0.0195    0.00178

> fgsa[is.na(fgsa$height),"age"]
[1] 18 69 52 41 52 44 73 28 66 20 73 63 26
# ages of those with missing height

# equivalent: fgsa$age[is.na(fgsa$height)]
```

## Naming

- ▶ Elements of vectors, rows and columns of matrices and data frames can have names

```
> x <- c(boys=1.2, girls=1.1)
> x
  boys  girls
  1.2   1.1
> x["boys"]
boys
  1.2
> D[, "a"]  # works for matrices and data frames
[1] 8 3 5
```

- ▶ You can extract and set names with `names(x)`; for matrices and data frames also `colnames(x)` and `rownames(x)`;

## Classes, generic functions

- ▶ R objects have *classes*
- ▶ Functions can behave differently depending on the class of an object
- ▶ E.g. `summary(x)` or `print(x)` does different things if `x` is numeric, a factor, or a linear model fit

```
> summary(x)  # a numeric vector
  Min. 1st Qu. Median Mean 3rd Qu. Max.
  1       1       2      2       3       3
> summary(fx) # a factor
  bad average good
  3         1     3
```

## Function calls

Round brackets ( ) are used for function calls!

Lots of things you do with R involve calling functions (you have seen that already!).

For instance

```
mean(x, na.rm=TRUE)
```

The important parts of this are

- ▶ The name of the function
- ▶ Arguments: input to the function
- ▶ Sometimes, we have named arguments

## Function calls

Round brackets ( ) are used for function calls!

Lots of things you do with R involve calling functions (you have seen that already!).

For instance

```
mean(x, na.rm=TRUE)
```

The important parts of this are

- ▶ The **name** of the function
- ▶ Arguments: input to the function
- ▶ Sometimes, we have named arguments

## Function calls

Round brackets ( ) are used for function calls!

Lots of things you do with R involve calling functions (you have seen that already!).

For instance

```
mean(x, na.rm=TRUE)
```

The important parts of this are

- ▶ The name of the function
- ▶ **Arguments**: input to the function
- ▶ Sometimes, we have named arguments

## Function calls

Round brackets ( ) are used for function calls!

Lots of things you do with R involve calling functions (you have seen that already!).

For instance

```
mean(x, na.rm=TRUE)
```

The important parts of this are

- ▶ The name of the function
- ▶ Arguments: input to the function
- ▶ Sometimes, we have **named arguments**

## Function arguments

Examples:

```
rnorm(10, mean=m, sd=s)
hist(x, main="My histogram")
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ Names of R objects
- ▶ Explicit constants
- ▶ Return values from another function call or expression
- ▶ Some arguments have their *default values*.
- ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
- ▶ Quite often – first argument is not named, but the others are named

## Function arguments

Examples:

```
rnorm(10, mean=m, sd=s)
hist(x, main="My histogram")
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ **Names** of R objects
- ▶ Explicit constants
- ▶ Return values from another function call or expression
- ▶ Some arguments have their *default values*.
- ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
- ▶ Quite often – first argument is not named, but the others are named

## Function arguments

## Examples:

```
rnorm(10, mean=m, sd=s)  
hist(x, main="My histogram")  
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ Names of R objects
  - ▶ Explicit **constants**
  - ▶ Return values from another function call or expression
  - ▶ Some arguments have their *default values*.
  - ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
  - ▶ Quite often – first argument is not named, but the others are named

## Function arguments

## Examples:

```
rnorm(10, mean=m, sd=s)  
hist(x, main="My histogram")  
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ Names of R objects
  - ▶ Explicit constants
  - ▶ Return values from another function call or expression
  - ▶ Some arguments have their *default values*.
  - ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
  - ▶ Quite often – first argument is not named, but the others are named

## Function arguments

### Examples:

```
rnorm(10, mean=m, sd=s)  
hist(x, main="My histogram")  
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ Names of R objects
  - ▶ Explicit constants
  - ▶ Return values from another function call or expression
  - ▶ Some arguments have their *default values*.
  - ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
  - ▶ Quite often – first argument is not named, but the others are named

## Function arguments

## Examples:

```
rnorm(10, mean=m, sd=s)  
hist(x, main="My histogram")  
mean(log(x + 1))
```

Items which may appear as arguments:

- ▶ Names of R objects
  - ▶ Explicit constants
  - ▶ Return values from another function call or expression
  - ▶ Some arguments have their *default values*.
  - ▶ Use `help(function)` or `args(function)` to see the arguments (and their order and default values) that can be given to any function.
  - ▶ Quite often – first argument is not named, but the others are named

## Example

From R-help (`help(t.test)`):

```
t.test(x, y = NULL,  
       alternative = c("two.sided", "less", "greater"),  
       mu = 0, paired = FALSE, var.equal = FALSE,  
       conf.level = 0.95, ...)
```

- ▶ The first argument (`x`) does not have a default – you have to provide some data!
  - ▶ The other arguments can be modified, if you need to.

## Example (cont.)

The following lines of code are equivalent:

```
t.test(a, b, alternative="less", paired=TRUE)  
t.test(a, b, paired=TRUE, alt="less")
```

```
t.test(a, b, p=T, a="l")    #not a good style!
```

Order does not matter for named arguments!

Partial keyword matching is possible ("alternative" or "alt" or "a")  
(partial matching is possible)

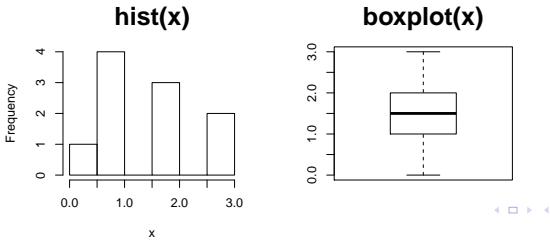
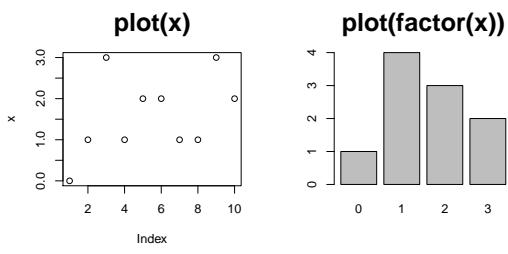
For a readable code, the use of explicit argument names is highly recommended!

## Basic graphics

The `plot()` function is a generic function, producing different plots for different types of arguments. For instance, `plot(x)` produces:

- ▶ a plot of observation index against the observations, when  $x$  is a numeric variable
  - ▶ a bar plot of category frequencies, when  $x$  is a factor variable
  - ▶ a time series plot (interconnected observations) when  $x$  is a time series
  - ▶ a set of diagnostic plots, when  $x$  is a fitted regression model
  - ▶ Similarly, the `plot(x, y)` produces a scatter plot, when  $x$  is a numeric variable and a bar plot of category frequencies, when  $x$  is a factor variable

## Some simple plots:



# The workspace

- ▶ The *global environment* contains R objects created on the command line.
  - ▶ There is an additional *search path* of loaded packages and attached data frames.
  - ▶ When you request an object by name, R looks first in the global environment, and if it doesn't find it there, it continues along the search path.
  - ▶ The search path is maintained by `library()`, `attach()`, and `detach()`
  - ▶ Notice that objects in the global environment may mask objects in packages and attached data frames

## More on factors: the `cut` Function

- ▶ The `cut` function converts a numerical variable into groups (a factor variable) according to a set of break points
  - ▶ The intervals are left-open, right-closed by default (`right=FALSE` changes that)
  - ▶ ... and that the lowest endpoint is *not* included by default (set `include.lowest=TRUE` if it bothers you)

## Example

```

> age <- c(35,20,21,50,46,23,30)
> agegr<-cut(age, c(20,30,40,50))
> table(agegr)
agegr      # the 20-year old is not included!
(20,30] (30,40] (40,50]
      3      1      2
> agegr<-cut(age, c(20,30,40,50),right=FALSE)
> table(agegr)
agegr      # the 50-year old is not included!
[20,30) [30,40) [40,50)
      3      2      1
> agegr<-cut(age, c(20,30,40,50),
+                                     include.lowest=TRUE)
> table(agegr)
agegr
[20,30] (30,40] (40,50]
4       1       2

```

# Working with Dates

- ▶ Dates are usually read as character or factor variables
  - ▶ Use the `as.Date` function to convert them to objects of class "Date"
  - ▶ If data are not in the default format (YYYY-MM-DD) you need to supply a format specification

## Working with Dates

- ▶ Dates are usually read as character or factor variables
  - ▶ Use the `as.Date` function to convert them to objects of class "Date"
  - ▶ If data are not in the default format (YYYY-MM-DD) you need to supply a format specification
    - > `as.Date("11/3/1959", format = "%d/%m-%Y")`  
[1] "1959-03-11"

A set of small, light-gray navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and table of contents.

## Working with Dates

- ▶ Dates are usually read as character or factor variables
  - ▶ Use the `as.Date` function to convert them to objects of class "Date"
  - ▶ If data are not in the default format (YYYY-MM-DD) you need to supply a format specification
    - > `as.Date("11/3-1959", format = "%d/%m-%Y")`

```
[1] "1959-03-11"
```

## Working with Dates

- ▶ Dates are usually read as character or factor variables
  - ▶ Use the `as.Date` function to convert them to objects of class "Date"
  - ▶ If data are not in the default format (YYYY-MM-DD) you need to supply a format specification

```
> as.Date("11/3-1959", format="%d/%m-%Y")
[1] "1959-03-11"
```

- ▶ You can calculate differences between `Date` objects. The result is an object of class "difftime". To get the number of days between two dates, use

```
> as.numeric(as.Date("2017-6-1") -  
           as.Date("1959-3-11"), "days")  
[1] 17607
```

## Creating your own functions

A very simple example:

```
logit <- function(p) log(p/(1-p))
```

The function `logit` requires one argument  $p$  and produces the logit of  $p$ . Try `logit(0.5)`, or `logit(0.25)`, ...

More complex (but still simple):

```
simpsum <- function(x, dec=5) {  
  m <- mean(x, na.rm=TRUE)  
  s <- sd(x, na.rm=TRUE)  
  round(c(mean=m, sd=s), dec) }
```

The function `simpsum` requires one argument  $x$ , but the second argument `dec` (no of decimal points in the output) has a default value 5. Try `simpsum(a)`, or `simpsum(a, dec=2)`.

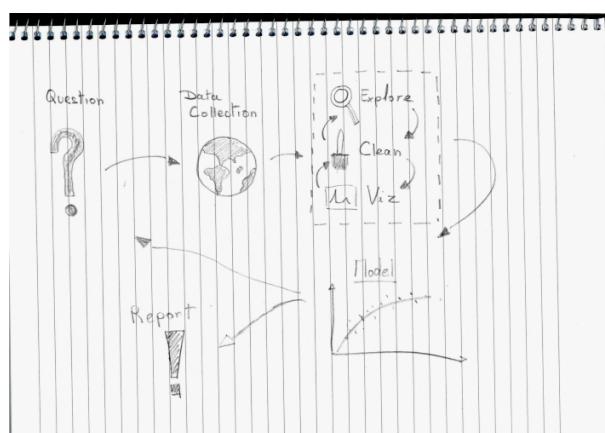
# Data manipulation with dplyr

Damien Georges

International Agency for Research on Cancer

June 2023 - Tartu

## Epidemiological study workflow



## Data manipulation tools



- ▶ R core function
- ▶ dplyr
- ▶ data.table
- ▶ ...

=> The best tool is the one you feel the most comfortable with

## Tidyverse (from [www.tidyverse.org](http://www.tidyverse.org))

### R packages for data science

The tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.



### pipe functions %>%

```
chill(fold(add(melt(add(chocolate, butter)),  
               beat(add(eggs.white, cream))))
```

### pipe functions %>%

```
chill(fold(add(melt(add(chocolate, butter)),  
               beat(add(eggs.white, cream))))
```

```
chocolate %>%  
  add(butter) %>%  
  melt() %>%  
  add(  
    eggs.white %>%  
      add(cream) %>%  
      beat()  
  ) %>%  
  fold() %>%  
  chill()
```

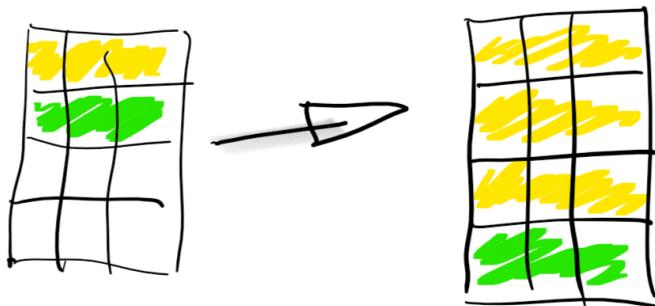
## code as you speak

Data manipulation with dplyr is done using a limited number of **verbs** corresponding to an action to be applied to a table.

- ▶ slice
- ▶ filter
- ▶ arrange
- ▶ select
- ▶ mutate
- ▶ group\_by
- ▶ summarize
- ▶ join
- ▶ ...

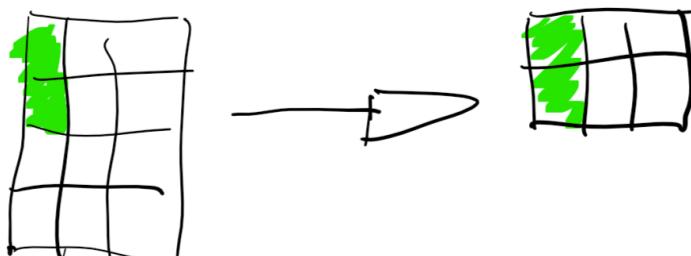
### select rows

```
dat %>% slice(c(1, 1, 1, 2))
```



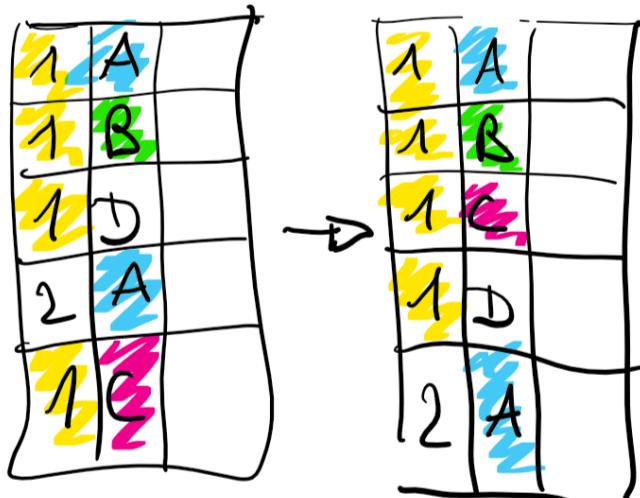
### select columns

```
dat %>% filter(C1 == 'green')
```



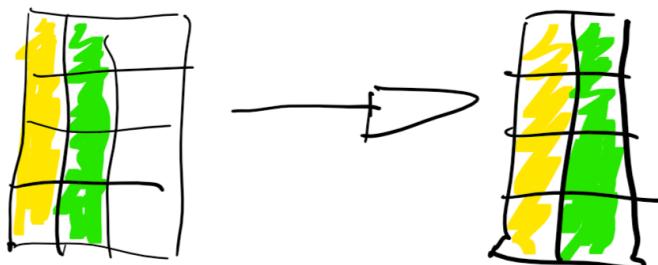
### arrange rows

```
dat %>% arrange(C1, C2)
```



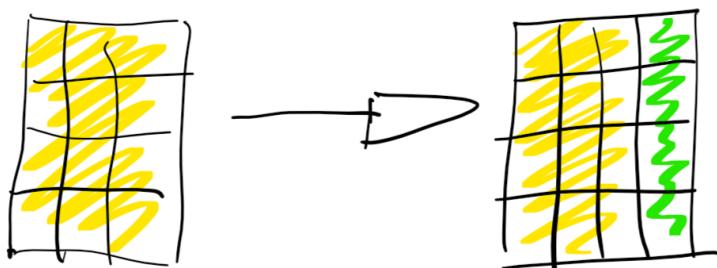
### columns selection

```
dat %>% select(C1, C2)
```



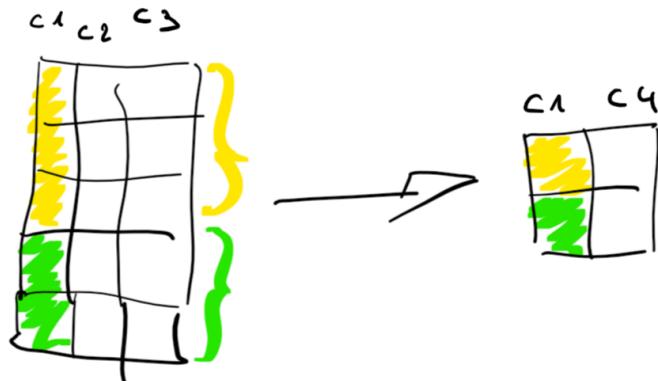
### create/modify columns

```
dat %>% mutate(C4 = C1 + C2 + C3)
```



## group and summarize data

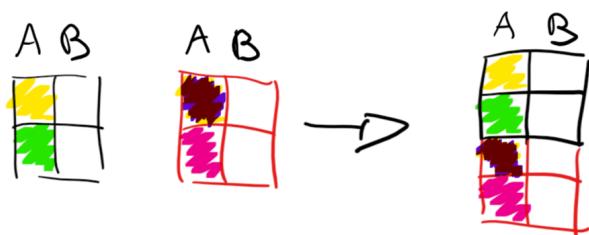
```
dat %>% group_by(C1) %>% summarize(C4 = mean(C2 + C3))
```



**note:** summarise() is an alias for summarize()

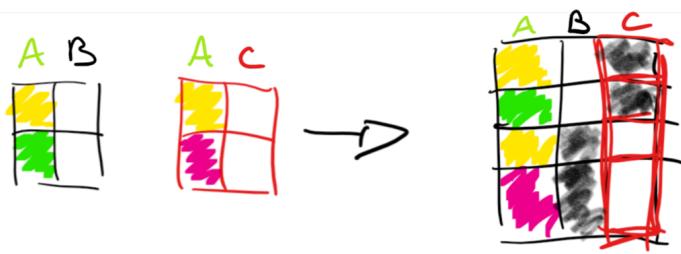
## bind and merge tables

```
dat1 %>% bind_rows(dat2)
```



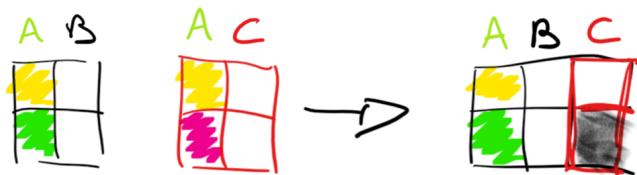
## bind and merge tables

```
dat1 %>% bind_rows(dat2)
```



bind and merge tables

```
dat1 %>% left_join(dat2)
```



**note:** right\_join will keep all rows of dat2

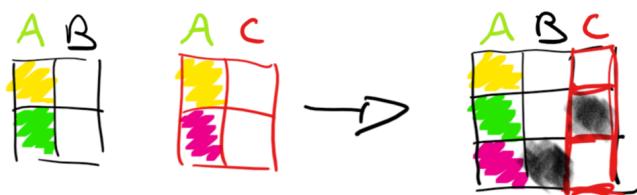
bind and merge tables

```
dat1 %>% inner_join(dat2)
```



bind and merge tables

```
dat1 %>% full_join(dat2)
```



## Poisson and Logistic Regression

Janne Pitkäniemi

Finnish Cancer Registry  
Tampere university

Statistical Practice in Epidemiology (2023, Tartu)

A set of small, light-blue navigation icons typically found in presentation software like Beamer. The icons include symbols for back, forward, search, and table of contents.

1 / 19

## Points to be covered

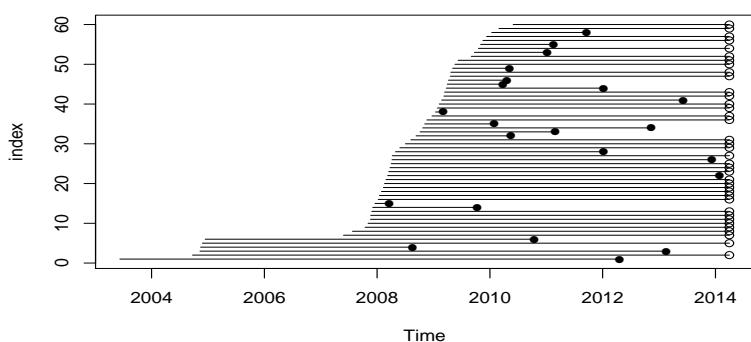
- ▶ Incidence rates, rate ratios and rate differences from *follow-up studies* can be computed by fitting *Poisson regression models*.
  - ▶ Odds ratios can be computed from binary data by fitting *Logistic regression models*.
  - ▶ Both models are special instances of *Generalized linear models*.
  - ▶ There are various ways to do these tasks in R.

2 / 19

2 / 19

## The Estonian Biobank cohort: survival among the elderly

Follow-up of 60 random individuals aged 75-103 at recruitment, until death (●) or censoring (○) in April 2014 (linkage with the Estonian Causes of Death Registry). (time-scale: calendar time).

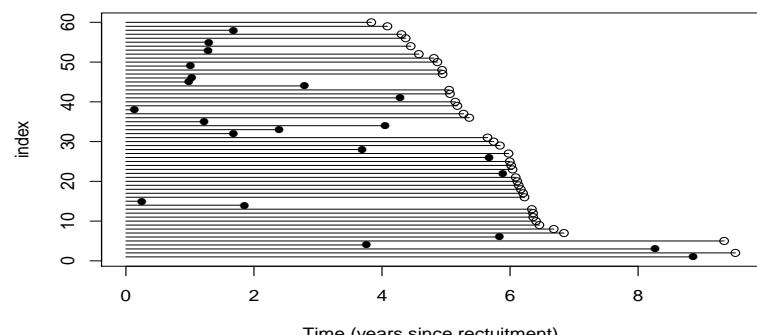


A set of small, light-blue navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and table of contents.

3 / 19

## The Estonian Biobank cohort: survival among the elderly

Follow-up time for 60 random individuals aged 75-103 at recruitment (time-scale: time in study).



## Events, dates and risk time

- ▶ Mortality as the outcome:
    - d: indicator for **status** at exit:
      - 1: death observed
      - 0: censored alive
  - ▶ Dates:

doe = date of Entry to follow-up,  
dox = date of eXit, end of follow-up.

- Follow-up time (years) computed as:

$$y = (\text{dox} - \text{doe}) / 365.25$$



Crude overall rate computed by hand and model

Total no. cases, person-years & rate (/1000 y):

```

> D <- sum( d ); Y <- sum(y) ; R <- D/(Y/1000)
> round( c(D=D, Y=Y, R=R), 2)
      D      Y      R
  884.00 11678.24   75.70

```

## R-implementation of the rate estimation with Poisson regression: A model with `poisreg`-family (Epi package)

```
> glm(cbind(D, Y) ~ 1, family=poisreg)
```

Coefficients :  
(Intercept) -2.581

*Why do we get the same results?*



## Constant hazard — Poisson model

Let  $Y \sim \text{exp}(\lambda)$ , then  $f(y; \lambda) = \lambda e^{-\lambda y} I(y > 0)$

Constant rate model:  $\lambda(y) = \frac{f(y; \lambda)}{S(y; \lambda)} = \lambda$  and observed data  $\{(y_i, \delta_i); i = 1, \dots, n\}$ .

The likelihood  $L(\lambda) = \prod_{i=1}^n \lambda^{\delta_i} e^{-\lambda y_i}$  and

$$\log(L) = \sum_{i=1}^n [\delta_i \log(\lambda) - \lambda y_i]$$

Solving the *score equations*:

$$\frac{\partial \log L(\lambda)}{\partial \lambda} = \sum \left[ \frac{\delta_i}{\lambda} - y_i \right] = \frac{D}{\lambda} - Y = 0 \text{ and } D - \lambda Y = 0$$

→ **maximum likelihood estimator (MLE)** of  $\lambda$ :

$$\hat{\lambda} = \frac{D}{Y} = \frac{\text{number of cases}}{\text{total person-time}} = \text{empirical rate!}$$

## offset term — Poisson model

- ▶ Previous model without offset: Intercept 6.784 =  $\log(884)$
- ▶ We should use an offset if we suspect that the underlying **population sizes (person-years) differ** for each of the observed counts – For example varying person-years by sex, age, treatment group, ...
- ▶ We need a term in the model that "scales" the likelihood, but does not depend on model parameters (include a **term with reg. coef. fixed to 1**) – offset term is  $\log(y)$
- ▶ This is all taken care of by family = poisson – recommend to use

$$\begin{aligned}\log\left(\frac{\mu}{y}\right) &= \beta_0 + \beta_1 x_1 \\ \log(\mu) &= 1 \times \log(y) + \beta_0 + \beta_1 x_1\end{aligned}$$

## Comparing rates: The Thorotrast Study

- ▶ Cohort of seriously ill patients in Denmark on whom angiography of brain was performed.
- ▶ Exposure: contrast medium used in angiography,
  1. thor = thorotrast (with  $^{232}\text{Th}$ ), used 1935-50
  2. ctrl = other medium (?), used 1946-63
- ▶ Outcome of interest: death

`doe` = date of Entry to follow-up,  
`dox` = date of eXit, end of follow-up.

- ▶ `data(thoro)` in the Epi package.

## Tabulating rates: thorotrast vs. control

Tabulating cases, person-years & rates by group

```
> stat.table( contrast,
+               list ( N = count(),
+                     D = sum(d),
+                     Y = sum(y),
+                     rate = ratio(d,y,1000) ) )
-----
contrast      N      D      Y    rate
-----
ctrl        1236  797.00 30517.56  26.12
thor         807   748.00 19243.85  38.87
-----
```

10 / 19

## Rate ratio estimation with Poisson regression

- ▶ Include `contrast` as the explanatory variable (factor).
- ▶ Insert person years in units that you want rates in

```
> m2 <- glm( cbind(d,y/1000) ~ contrast,family = poisreg(link="log") )
> round( summary(m2)$coef, 4)[, 1:2]
```

	Estimate	Std. Error
(Intercept)	3.2626	0.0354
contrast thor	0.3977	0.0509

- ▶ Rate ratio and CI?

Call function `ci.exp()` in `Epi`

```
> round( ci.exp( m2 ), 3 )
      exp(Est.) 2.5% 97.5%
(Intercept) 26.116 24.364 27.994
contrast thor 1.488 1.347 1.644
```

11 / 19

## Rates in groups with Poisson regression

- ▶ Include `contrast` as the explanatory variable (factor).
- ▶ Remove the intercept (-1)
- ▶ Insert person-years in units that you want rates in

```
> m3 <- glm( cbind(d,y/1000) ~ factor(contrast)-1,family = poisreg)
> round( summary(m3)$coef, 4)[, 1:2]
```

	Estimate	Std. Error
contrast ctrl	3.2626	0.0354
contrast thor	3.6602	0.0366

```
> round( ci.exp( m3 ), 3 )

      exp(Est.) 2.5% 97.5%
contrast ctrl 26.116 24.364 27.994
contrast thor 38.870 36.181 41.757
```

12 / 19

## Rate difference estimation with Poisson regression

- The approach with d/y enables additive rate models too:

```
> contrast<-c(0,1)
> m5 <-glm(cbind(d,y/1000) ~contrast,
           family=poisreg(link="identity") )
> round( ci.exp(m5,Exp=F), 3 )
```

```
Estimate 2.5% 97.5%
(Intercept) 26.116 24.303 27.929
contrast thor 12.753 9.430 16.077
```

## Binary data: Treatment success Y/N

85 diabetes-patients with foot-wounds:

- Dalteparin (Dal)
- Placebo (Pl)

Treatment/Placebo given to diabetes patients, the design is prospective and outcome is measured better(Y)/worse(N). Is the probability of outcome more than 15% – yes, then use the risk difference or risk ratio (RR)

	Treatment group	
	Dalteparin	Placebo
Better	29	20
Worse	14	22
Total	43	42

$$\hat{p}_{\text{Dal}} = \frac{29}{43} = 67\% \quad \hat{p}_{\text{Pl}} = \frac{20}{42} = 47\%$$

## Binary data: Crosstabulation analysis of 2x2 table

```
> library(Epi)
> dlt <- rbind( c(29,14), c(20,22) )
> colnames( dlt ) <- c("Better","Worse")
> rownames( dlt ) <- c("Dal","Pl")
> kable(twoby2( dlt ),"latex")
```

```
2 by 2 table analysis:
  Better Worse P(Better) 95% conf. interval
Dal    29     14     0.6744   0.5226   0.7967
Pl     20     22     0.4762   0.3316   0.6249
                                         95% conf. interval
                                         Relative Risk: 1.4163  0.9694  2.0692
                                         Sample Odds Ratio: 2.2786  0.9456  5.4907
                                         Conditional MLE Odds Ratio: 2.2560  0.8675  6.0405
                                         Probability difference: 0.1982 -0.0110  0.3850

Exact P-value: 0.0808
Asymptotic P-value: 0.0665
```

## Binary regression – estimation of odds ratio

For grouped binary data, the response is a two-column matrix with columns (successes,failures).

```
> library(Epi)
> library(xtable)
> dlt <- data.frame(rbind( c(29,14),c(20,22) ))
> colnames( dlt ) <- c("Better","Worse")
> dlt$trt <- c(1,0)
> b2<-glm(cbind(Better,Worse)~trt,
+           family=binomial(link="logit"),
+           data=dlt)
> xtable(round( ci.exp( b2 ), digits=6 ))
```

	exp(Est.)	2.5%	97.5%
(Intercept)	0.91	0.50	1.67
trt	2.28	0.95	5.49

- ▶ The default parameters in logistic regression are **odds** (the intercept:  $20/22 = 0.9090$ ) and the **odds-ratio**  $((29/14)/(20/22) = 2.28)$ .
- ▶ This is **NOT** what you want, because odds ratio is biased estimate of the risk ratio.(recall if  $p>10\% \frac{p}{1-p} \not\approx p$ )

## Binary regression - Estimation of risk ratio (Relative risk)

```
> library(Epi)
> library(xtable)
> dlt <- data.frame(rbind( c(29,14),c(20,22) ))
> colnames( dlt ) <- c("Better","Worse")
> dlt$trt <- c(1,0)
> b2<-glm(cbind(Better,Worse)~trt,
+           family=binomial(link="log"),
+           data=dlt)
> xtable(round( ci.exp( b2 ), digits=6 ))
```

	exp(Est.)	2.5%	97.5%
(Intercept)	0.48	0.35	0.65
trt	1.42	0.97	2.07

Diabetics with Dalterapin treatment are 1.4 times likely to get better than those treated with placebo

## Binary regression - Estimation of risk difference

```
> library(Epi)
> library(xtable)
> dlt <- data.frame(rbind( c(29,14),c(20,22) ))
> colnames( dlt ) <- c("Better","Worse")
> dlt$trt <- c(1,0)
> b2<-glm(cbind(Better,Worse)~trt,
+           family=binomial(link="identity"),
+           data=dlt)
> xtable(round( ci.exp( b2,Exp=F ), digits=6 ))
```

	Estimate	2.5%	97.5%
(Intercept)	0.48	0.33	0.63
trt	0.20	-0.01	0.40

Twenty percent more of the Diabetics with Dalterapin treatment are getting better compared to Diabetics treated with placebo

## Conclusion: What did we learn?

- ▶ Rates, their ratio and difference can be analysed by Poisson regression
- ▶ In Poisson models the response can be either:
  - ▶ case indicator d with offset =  $\log(y)$ , or
  - ▶ case and person-years cbind(d,y) with poisreg-family (Epi-package)
- ▶ Both may be fitted on either grouped data, or individual records.
- ▶ Binary outcome can be modeled with binary regression.

# Linear and generalized linear models

Saturday 24 August, 2019

Esa Läärä

Statistical Practice in Epidemiology with R

23 to 28 August, 2019

University of Tartu, Estonia

## Outline

- ▶ Simple linear regression.
- ▶ Fitting a model and extracting results.
- ▶ Predictions and diagnostics.
- ▶ Categorical factors and contrast matrices.
- ▶ Main effects and interactions.
- ▶ Generalized linear models.
- ▶ Modelling curved effects.

## Variables in generalized linear models

- ▶ The **outcome** or **response** variable must be numeric.
- ▶ Main types of response variables are
  - Metric or continuous (a measurement with units)
  - Binary (two values coded 0/1)
  - Failure (does the subject fail at end of follow-up)
  - Count (aggregated failure data, number of cases)
- ▶ **Explanatory** variables or **regressors** can be
  - Numeric or quantitative variables
  - Categorical factors, represented by class indicators or contrast matrices.

## The births data in Epi

id: Identity number for mother and baby.  
bweight: Birth weight of baby.  
lowbw: Indicator for birth weight less than 2500 g.  
gestwks: Gestation period in weeks.  
preterm: Indicator for gestation period less than 37 weeks.  
matage: Maternal age.  
hyp: Indicator for maternal hypertension (0 = no, 1 = yes).  
sex: Sex of baby (1 = male, 2 = female).

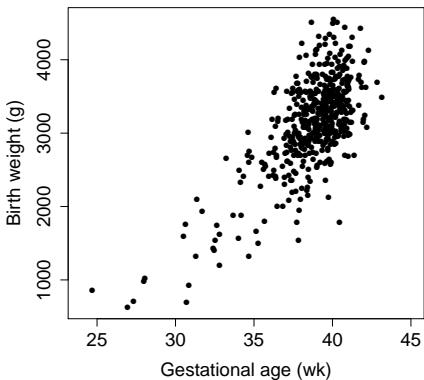
Declaring and transforming some variables as factors:

```
> library(Epi) ; data(births)
> births <- transform(births,
+   hyp = factor(hyp, labels=c("N", "H")),
+   sex = factor(sex, labels=c("M", "F")),
+   gest4 = cut(gestwks,breaks=c(20, 35, 37, 39, 45), right=FALSE) )
> births <- subset(births, !is.na(gestwks))
```

Linear and generalized linear models

3 / 22

## Birth weight and gestational age



```
> with(births, plot(bweight ~ gestwks, xlim = c(24,45), pch = 16, cex.axis=1.5, cex.lab = 1.5,
+   xlab= "Gestational age (wk)", ylab= "Birth weight (g)" ) )
```

Linear and generalized linear models

4 / 22

## Metric response, numeric explanatory variable

Roughly linear relationship btw bweight and gestwks

→ Simple **linear regression model** fitted.

```
> m <- lm(bweight ~ gestwks, data=births)
```

- ▶ `lm()` is the function that fits linear regression models, assuming **Gaussian** distribution for **error** terms.
- ▶ `bweight ~ gestwks` is the **model formula**
- ▶ `m` is a **model object** belonging to **class "lm"**.

```
> coef(m) – Printing the estimated regression coefficients
```

(Intercept)	gestwks
-4489.1	197.0

Interpretation of **intercept** and **slope**?

Linear and generalized linear models

5 / 22

## Model object and extractor functions

Model object = **list** of different elements, each being separately accessible. – See `str(m)` for the full list.

Functions that extract results from the fitted model object

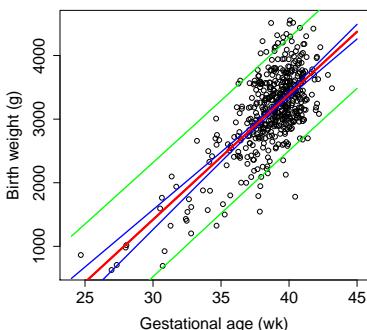
- ▶ `summary(m)` – lots of output
  - ▶ `coef(m)` – beta-hats only (see above)
  - ▶ `ci.lin(m) [,c(1,5,6)]` –  $\hat{\beta}_j$ s plus confidence limits
    - Estimate 2.5% 97.5%
    - (Intercept) -4489.1 -5157.3 -3821.0
    - gestwks 197.0 179.7 214.2
- This function is in Epi package
- ▶ `anova(m)` – Analysis of Variance Table

## Other extractor functions, for example

- ▶ `fitted(m)`, `resid(m)`, `vcov(m)`, ...
- ▶ `predict(m, newdata = ..., interval=...)`
  - Predicted responses for desired combinations of new values of the regressors – `newdata`
  - Argument `interval` specifies whether **confidence** intervals for the *mean* response or **prediction** intervals for *individual* responses are returned.
- ▶ `plot(m)` – produces various diagnostic plots based on residuals (raw or standardized)

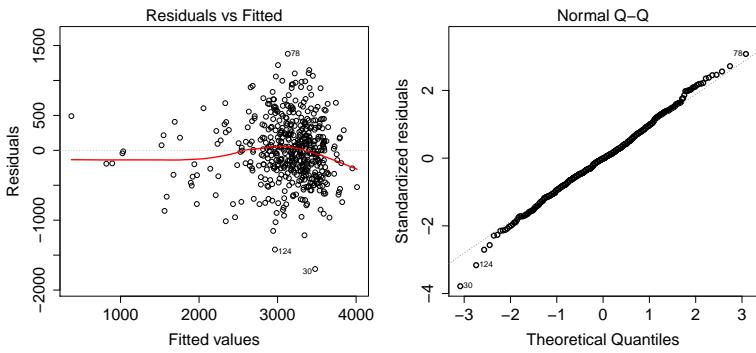
Many of these are special **methods** for certain **generic functions**, aimed at acting on objects of class “`lm`”.

## Fitted values, confidence & prediction intervals



```
> nd <- data.frame( gestwks = seq(24, 45, by = 0.25) )
> pr.c1 <- predict( m, newdata=nd, interval="conf" )
> pr.p1 <- predict( m, newdata=nd, interval="pred" )
> with(births, plot(bweight ~ gestwks, xlim = c(24,45), cex.axis=1.5, cex.lab = 1.5, xlab = 'Gestation
> matlines( nd$gestwks, pr.c1, lty=1, lwd=c(3,2,2), col=c('red','blue','blue')))
> matlines( nd$gestwks, pr.p1, lty=1, lwd=c(3,2,2), col=c('red','green','green')))
```

## A couple of diagnostic plots



```
> par(mfrow=c(1,2))
> plot(m, 1:2, cex.lab = 1.5, cex.axis=1.5, cex.caption=1.5, lwd=2)
```

- ▶ Some deviation from linearity?
- ▶ Reasonable agreement with Gaussian error assumption?

## Factor as an explanatory variable

- ▶ How bweight depends on maternal hypertension?

```
> mh <- lm( bweight ~ hyp, data=births)
```

	Estimate	2.5%	97.5%
(Intercept)	3198.9	3140.2	3257.6
hypH	-430.7	-585.4	-275.9

- ▶ Removal of intercept → mean bweights by hyp:

```
> mh2 <- lm( bweight ~ -1 + hyp, data = births)
> coef(mh2)
```

	hypN	hypH
	3198.9	2768.2

- ▶ Interpretation:  $-430.7 = 2768.2 - 3198.9$  = difference between level 2 vs. reference level 1 of hyp

## Additive model with both gestwks and hyp

- ▶ Joint effect of hyp and gestwks under additivity is modelled e.g. by updating a simpler model:

```
> mhg <- update(mh, . ~ . + gestwks)
Estimate      2.5%    97.5%
(Intercept) -4285.0 -4969.7 -3600.3
hypH         -143.7  -259.0   -28.4
gestwks       192.2   174.7   209.8
```

- ▶ The effect of hyp: H vs. N is attenuated (from  $-430.7$  to  $-143.7$ ).
- ▶ This suggests that much of the effect of hypertension on birth weight is mediated through a shorter gestation period among hypertensive mothers.

## Model with interaction of hyp and gestwks

- ▶ 

```
mhgi <- lm(bweight ~ hyp + gestwks +
               hyp:gestwks, data = births)
```
- ▶ Or with shorter formula: `bweight ~ hyp * gestwks`

	Estimate	2.5%	97.5%
(Intercept)	-3960.8	-4758.0	-3163.6
hypH	-1332.7	-2841.0	175.7
gestwks	183.9	163.5	204.4
hypH:gestwks	31.4	-8.3	71.1
- ▶ Estimated slope: 183.9 g/wk in reference group N and  $183.9 + 31.4 = 215.3$  g/wk in hypertensive mothers.
- ⇒ For each additional week the difference in mean bweight between H and N group increases by 31.4 g.
- ▶ *Interpretation of Intercept and “main effect” hypH?*

## Model with interaction (cont'd)

More interpretable parametrization obtained if `gestwks` is **centered** at some reference value, using e.g. the **insulate** operator `I()` for explicit transformation of an original term.

- ▶ 

```
mi2 <- lm(bweight ~ hyp*I(gestwks-40), ...)
```

	Estimate	2.5%	97.5%
(Intercept)	3395.6	3347.5	3443.7
hypH	-77.3	-219.8	65.3
I(gestwks - 40)	183.9	163.5	204.4
hypH:I(gestwks - 40)	31.4	-8.3	71.1
- ▶ Main effect of `hyp` =  $-77.3$  is the difference between H and N at `gestwks = 40`.
- ▶ Intercept = 3395.6 is the estimated mean `bweight` at the reference value 40 of `gestwks` in group N.

## Factors and contrasts in R

- ▶ A categorical explanatory variable or **factor** with  $L$  **levels** will be represented by  $L - 1$  linearly independent columns in the **model matrix** of a linear model.
- ▶ These columns can be defined in various ways implying alternative **parametrizations** for the effect of the factor.
- ▶ Parametrization is defined by given type of **contrasts**.
- ▶ Default: **treatment** contrasts, in which 1st class is the **reference**, and regression coefficient  $\beta_k$  for class  $k$  is interpreted as  $\beta_k = \mu_k - \mu_1$
- ▶ Own parametrization may be tailored by function `C()`, with the pertinent **contrast matrix** as argument.
- ▶ Or, use `ci.lin(mod, ctr.mat = CM)` after fitting.

## Two factors: additive effects

- ▶ Factor  $X$  has 3 levels,  $Z$  has 2 levels – Model:

$$\mu = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \gamma_1 Z_1 + \gamma_2 Z_2$$

- ▶  $X_1$  (reference),  $X_2, X_3$  are the indicators for  $X$ ,

- ▶  $Z_1$  (reference),  $Z_2$  are the indicators for  $Z$ .

- ▶ Omitting  $X_1$  and  $Z_1$  the model for mean is:

$$\mu = \alpha + \beta_2 X_2 + \beta_3 X_3 + \gamma_2 Z_2$$

with predicted means  $\mu_{jk}$  ( $j = 1, 2, 3; k = 1, 2$ ):

	$Z = 1$	$Z = 2$
	$\mu_{11} = \alpha$	$\mu_{11} = \alpha + \gamma_2$
$X$	$\mu_{21} = \alpha + \beta_2$	$\mu_{22} = \alpha + \beta_2 + \gamma_2$
	$\mu_{31} = \alpha + \beta_3$	$\mu_{32} = \alpha + \beta_3 + \gamma_2$

## Two factors with interaction

- ▶ Effect of  $Z$  differs at different levels of  $X$ :

	$Z = 1$	$Z = 2$
	$\mu_{11} = \alpha$	$\mu_{12} = \alpha + \gamma_2$
$X$	$\mu_{21} = \alpha + \beta_2$	$\mu_{22} = \alpha + \beta_2 + \gamma_2 + \delta_{22}$
	$\mu_{31} = \alpha + \beta_3$	$\mu_{32} = \alpha + \beta_3 + \gamma_2 + \delta_{32}$

- ▶ How much the effect of  $Z$  (level 2 vs. 1) changes when the level of  $X$  is changed from 1 to 3:

$$\begin{aligned}\delta_{32} &= (\mu_{32} - \mu_{31}) - (\mu_{12} - \mu_{11}) \\ &= (\mu_{32} - \mu_{12}) - (\mu_{31} - \mu_{11}),\end{aligned}$$

= how much the effect of  $X$  (level 3 vs. 1) changes when the level of  $Z$  is changed from 1 to 2.

- ▶ See the exercise: interaction of hyp and gest4.

## Contrasts in R

- ▶ All contrasts can be implemented by supplying a suitable **contrast function** giving the **contrast matrix** e.g:

```
> contr.cum(3)           > contr.sum(3)
 1 0 0                  1   1   0
 2 1 0                  2   0   1
 3 1 1                  3  -1  -1
```

- ▶ In model formula factor name faktori can be replaced by expression like `C(faktori, contr.cum)`.

- ▶ Function `ci.lin()` has an option for calculating CI's for linear functions of the parameters of a fitted model `mall` when supplied by a relevant contrast matrix

```
> ci.lin(mall, ctr.mat = CM)[ , c(1,5,6)]
```

→ No need to specify contrasts in model formula!

## From linear to generalized linear models

- ▶ An alternative way of fitting our 1st Gaussian model:

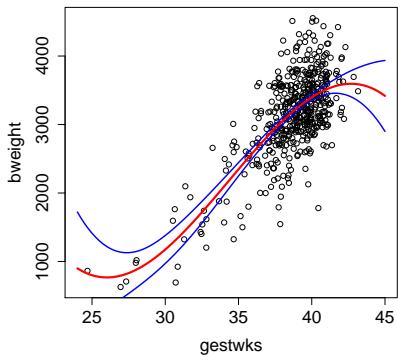
```
> m <- glm(bweight ~ gestwks, family=gaussian, data=births)
```
- ▶ Function `glm()` fits **generalized linear models** (GLM).
- ▶ Requires specification of the
  - **family** – i.e. the assumed “error” distribution for  $Y_i$ s,
  - **link** function – a transformation of the expected  $Y_i$ .
- ▶ Covers common models for other types of response variables and distributions, too, e.g. **logistic** regression for binary responses and **Poisson** regression for counts.
- ▶ Fitting: method of **maximum likelihood**.
- ▶ Many extractor functions for a `glm` object similar to those for an `lm` object.

## More about numeric regressors

What if dependence of  $Y$  on  $X$  is non-linear?

- ▶ **Categorize** the values of  $X$  into a factor.
  - Continuous effects violently discretized by often arbitrary cutpoints. – Inefficient.
- ▶ Fit a low-degree (e.g. 2 to 4) **polynomial** of  $X$ .
  - Tail behaviour may be problematic.
- ▶ Use **fractional polynomials**.
  - Invariance problems. Only useful if  $X = 0$  is well-defined.
- ▶ Use a **spline** model: smooth function  $s(X; \beta)$ . – See Martyn’s lecture
  - More flexible models that act locally.
  - Effect of  $X$  reported by graphing  $\hat{s}(X; \beta)$  & its CI

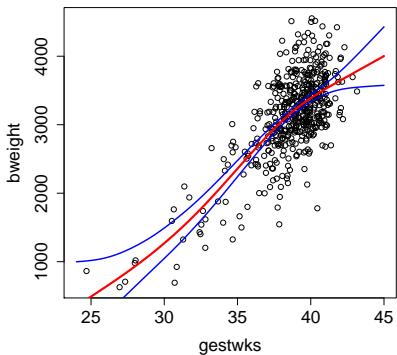
## Mean bweight as 3rd order polynomial of gestwks



```
> mp3 <- update( m, . ~ . - gestwks + poly(gestwks, 3) )
```

- ▶ The model is linear in parameters with 4 terms & 4 df.
- ▶ Otherwise good, but the tails do not behave well.

## Penalized spline model with cross-validation



```
> library(mgcv)
> mpen <- gam( bweight ~ s(gestwks), data = births)
```

- ▶ Looks quite nice.
- ▶ Model df ≈ 4.2; close to 4, as in the 3rd degree polynomial model.

## What was covered

- ▶ A wide range of models from simple linear regression to splines.
- ▶ R functions fitting linear and generalized models:  
`lm()` and `glm()`.
- ▶ Parametrization of categorical explanatory factors; contrast matrices.
- ▶ Extracting results and predictions:  
`ci.lin()`, `fitted()`, `predict()`, ... .
- ▶ Model diagnostics:  
`resid()`, `plot.lm()`, ... .

# Everything you ever wanted to know about splines but were too afraid to ask

Martyn Plummer

University of Warwick

03 June 2023



Categorization and its discontents  
ooooo

Join the dots  
oooooooooooo

Smoothing splines  
oooo

Splines in R  
ooooo

## Overview

Categorization and its discontents

Join the dots

Smoothing splines

Splines in R



Categorization and its discontents  
ooooo

Join the dots  
oooooooooooo

Smoothing splines  
oooo

Splines in R  
ooooo

## Introduction

- Splines are a flexible class of models that can be helpful for representing dose-response relationships in epidemiology
- In this course we will be using spline models extensively.
- However, spline models are widely misunderstood.
- The purpose of this lecture is to give a conceptual background on where spline models come from.



## Outline

## Categorization and its discontents



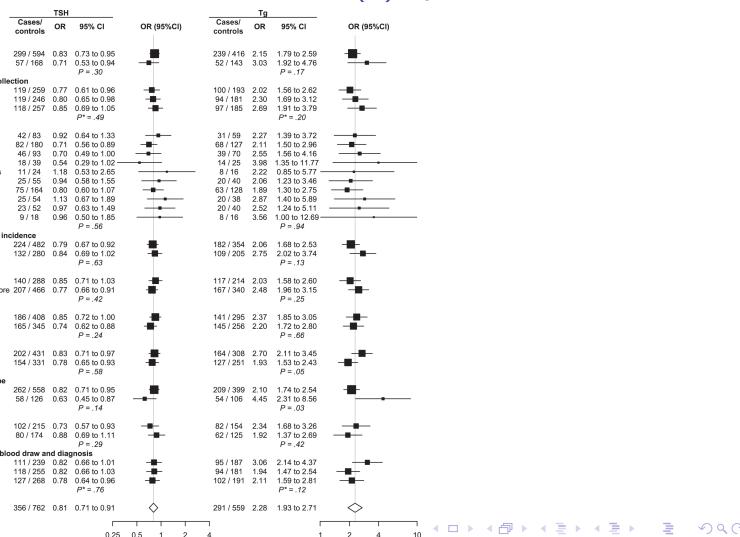
## Categorization and its discontents

Join the dots  
oooooooooooo

## Smoothing splines

Splines in R  
ooooo

Rinaldi et al, JNCI. 2014 Jun;106(6):dju097



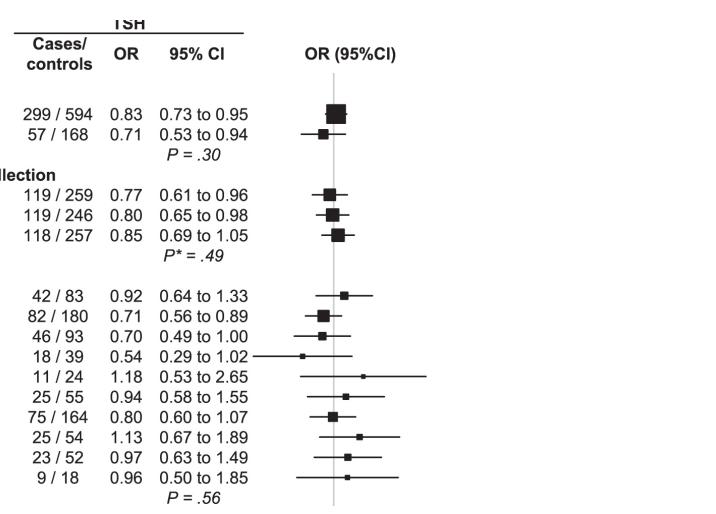
Categorization and its discontents

Join the dots

Smoothing splines

Splines in R

Rinaldi et al. JNCI 2014 Jun;106(6):diu097



## Statisticians against categorization

- Greenland S (1995) Avoiding power loss associated with categorization and ordinal scores in dose-response and trend analysis, *Epidemiology*, **6**, 450–454.
- Senn S (2005) Dichotomania: an obsessive compulsive disorder that is badly affecting the quality of analysis of pharmaceutical trials.
- Bennette C, and Vickers A, (2012), Against quantiles: categorization of continuous variables in epidemiologic research, and its discontents. *BMC Medical Research Methodology* 12:21

## Epidemiologists against categorization

Rose, G. (1992) The Strategy of Preventive Medicine

- Many diseases are not discrete. Instead there is an underlying continuum of increasing severity (e.g. hypertension).
- In medicine, we tend to conflate a clinical action (treat vs. do not treat) with the presence/absence of disease.
- Disease prevention efforts are best targeted at shifting the distribution of risk for the whole population instead of trying to identify and target a “high risk” group.

## Outline

Categorization and its discontents

Join the dots

Smoothing splines

Splines in R

Categorization and its discontents

Join the dots

Smoothing splines

Splines in R

# Join the dots

The figure displays a scatter plot with numerous points, each labeled with a number. The points are interconnected by lines, forming a complex, winding path. The numbers are as follows:

- Top row: 81, 82, 83, 84, 85, 86, 87, 88, 89, 80, 79.
- Second row: 78, 77, 76, 75, 74, 73, 72, 71, 70, 69, 68, 67, 66, 65, 64, 63, 62, 61, 60, 59, 58, 57, 56, 55, 54, 53, 52, 51, 50, 49, 48, 47, 46, 45, 44, 43, 42, 41, 40, 39, 38, 37, 36, 35, 34, 33, 32, 31, 30, 29, 28, 27, 26, 25, 24, 23, 22, 21, 20, 19, 18, 17, 16, 15, 14, 13, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1.
- Bottom row: 109, 108, 107, 106, 105, 104, 103, 102, 101, 100, 99, 98, 97, 96, 95, 94, 93, 92, 91, 90, 89, 88, 87, 86, 85, 84, 83, 82, 81, 80, 79, 78, 77, 76, 75, 74, 73, 72, 71, 70, 69, 68, 67, 66, 65, 64, 63, 62, 61, 60, 59, 58, 57, 56, 55, 54, 53, 52, 51, 50, 49, 48, 47, 46, 45, 44, 43, 42, 41, 40, 39, 38, 37, 36, 35, 34, 33, 32, 31, 30, 29, 28, 27, 26, 25, 24, 23, 22, 21, 20, 19, 18, 17, 16, 15, 14, 13, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1.

The figure illustrates two different ways to connect a series of data points. On the left, labeled "Join the dots", the points are connected by straight line segments, forming a polygon. On the right, labeled "Smoothing splines", the points are connected by a single, continuous smooth curve. The points are numbered from 1 to 100.

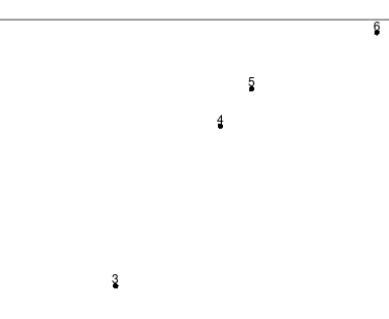
Categorization and its discontents  
ooooo

Join the dots  
oo●oooooooooooo

Smoothing splines  
oooo

Splines in R  
oooooo

## Linear interpolation

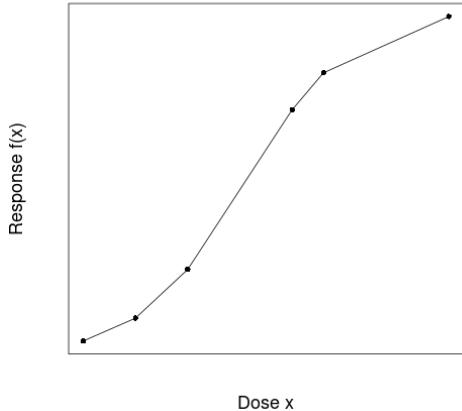


The figure is a scatter plot with 'Dose x' on the horizontal axis and 'Response f(x)' on the vertical axis. Six data points are plotted and labeled with integers 1 through 6. Point 1 is at a low dose and low response. Point 2 is at a slightly higher dose and response. Point 3 is at a medium dose and response. Point 4 is at a medium-high dose and response. Point 5 is at a high dose and response. Point 6 is at a very high dose and response. The plot area is enclosed in a light gray border.

- Suppose a dose response curve is known exactly at certain points
- We can fill in the gaps (interpolate) by drawing a straight (linear) line between adjacent points
- This creates a mathematical function  $f()$  which gives a response value  $f(x)$  for every dose value  $x$ .

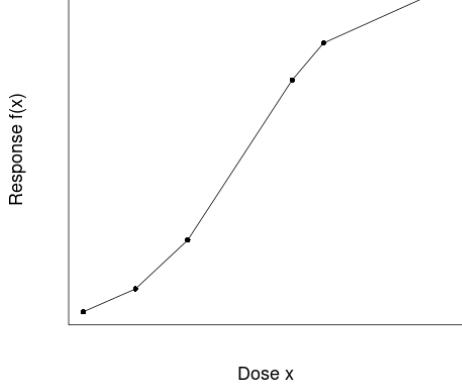
- Suppose a dose response curve is known exactly at certain points
  - We can fill in the gaps (interpolate) by drawing a straight (linear) line between adjacent points
  - This creates a mathematical function  $f()$  which gives a response value  $f(x)$  for every dose value  $x$ .

## Linear interpolation



- Suppose a dose response curve is known exactly at certain points
- We can fill in the gaps (interpolate) by drawing a straight (linear) line between adjacent points
- This creates a mathematical function  $f()$  which gives a response value  $f(x)$  for every dose value  $x$ .

## Linear interpolation



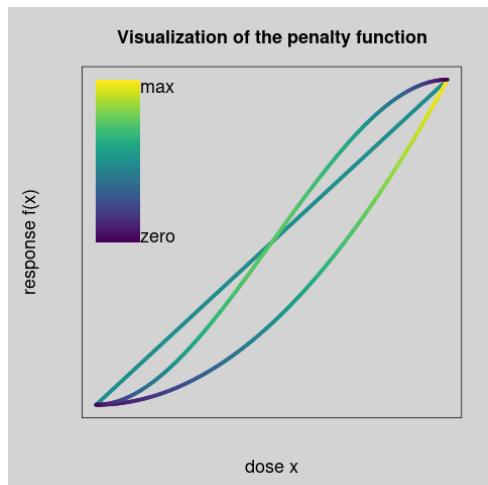
- Suppose a dose response curve is known exactly at certain points
- We can fill in the gaps (interpolate) by drawing a straight (linear) line between adjacent points
- This creates a mathematical function  $f()$  which gives a response value  $f(x)$  for every dose value  $x$ .

## Why linear interpolation?

Out of all possible curves that go through the observed points, linear interpolation is the one that minimizes the penalty function

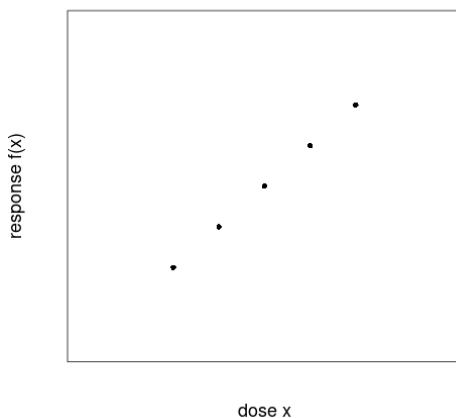
$$\int \left( \frac{\partial f}{\partial x} \right)^2 dx$$

## What does the penalty mean?



- The contribution to the penalty at each point depends on the steepness of the curve (represented by a colour gradient)
- Any deviation from a straight line between the two fixed points will incur a higher penalty overall.

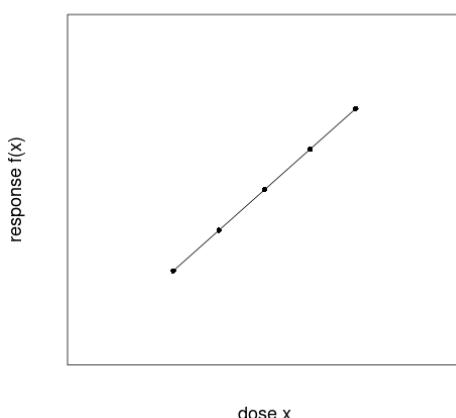
## Extrapolation



- Linear interpolation fits a linear dose-response curve exactly.
- But it breaks down when we extrapolate.

## Extrapolation

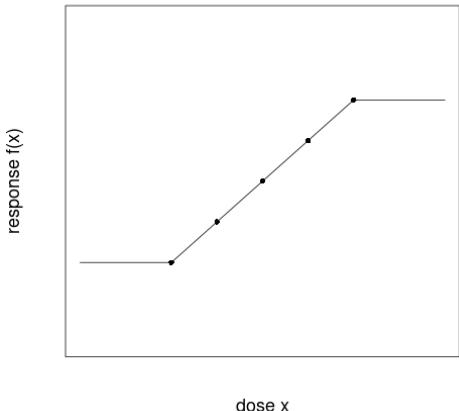
Linear interpolation



- Linear interpolation fits a linear dose-response curve exactly.
- But it breaks down when we extrapolate.

## Extrapolation

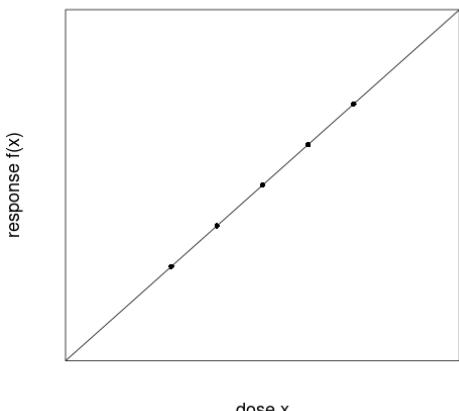
Extrapolation - not what we want



- Linear interpolation fits a linear dose-response curve exactly.
- But it breaks down when we extrapolate.

## Extrapolation

We want this



- Linear interpolation fits a linear dose-response curve exactly.
- But it breaks down when we extrapolate.

## Why does linear interpolation break down?

- The penalty function

$$\int \left( \frac{\partial f}{\partial x} \right)^2 dx$$

penalizes the steepness of the curve

- Minimizing the penalty function gives us the “flattest” curve that goes through the points.

- In between two observations the flattest curve is a straight line.
- Outside the range of the observations the flattest curve is completely flat.

## A roughness penalty

- If we want a fitted curve that extrapolates a linear trend then we want to minimize the curvature.

$$\int \left( \frac{\partial^2 f}{\partial x^2} \right)^2 dx$$

- Like the first penalty function but uses the second derivative of  $f$  (i.e. the curvature).
- This is a roughness penalty.

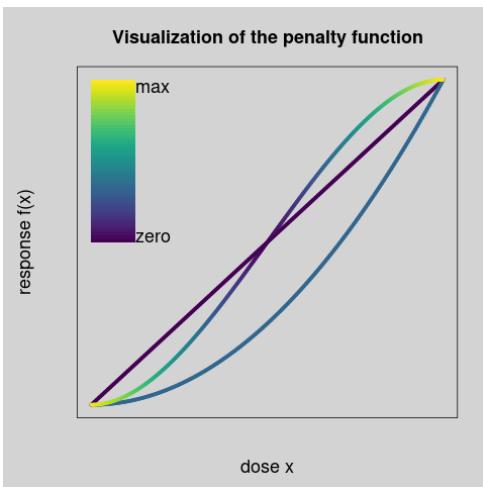
## A roughness penalty

- If we want a fitted curve that extrapolates a linear trend then we want to minimize the **curvature**.

$$\int \left( \frac{\partial^2 f}{\partial x^2} \right)^2 dx$$

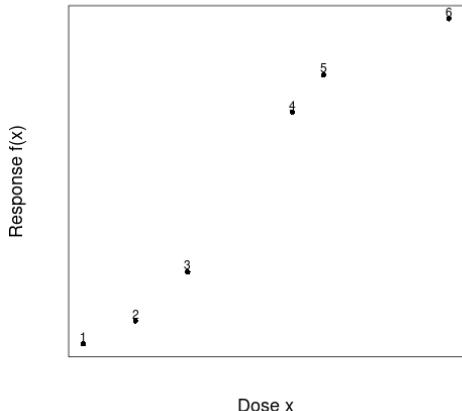
- Like the first penalty function but uses the **second derivative** of  $f$  (i.e. the curvature).
- This is a roughness penalty.

## What does the roughness penalty mean?



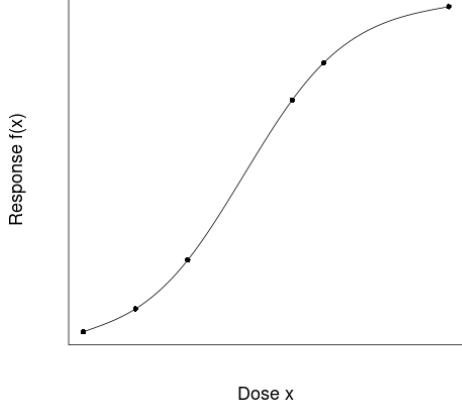
- The contribution to the penalty at each point depends on the curvature (represented by a colour gradient)
- A straight line has no curvature, hence zero penalty.
- Sharp changes in the slope are heavily penalized.

## An interpolating cubic spline



- The smoothest curve that goes through the observed points is a cubic spline.

## An interpolating cubic spline



- The smoothest curve that goes through the observed points is a cubic spline.

## What is a cubic spline?

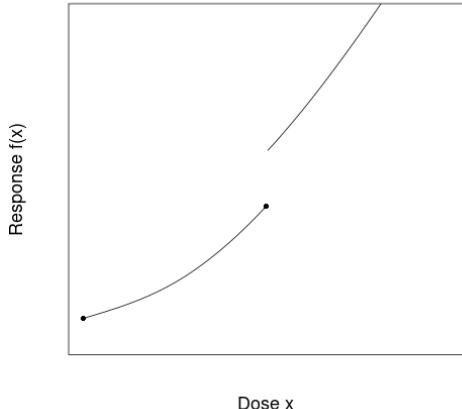
Splines are piecewise cubic curves

- Every observed point is a knot.
- The knots divide the curve into sections
- Each section is a cubic function

$$f(x) = a + bx + cx^2 + dx^3$$

- The parameters  $a, b, c, d$  are different for different sections

## Boundary conditions

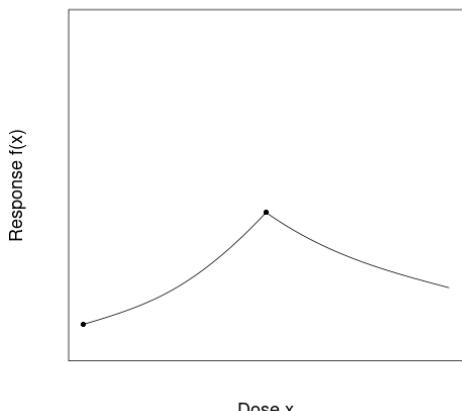


Sections need to join up smoothly.

- Both sides must go through the knot.
- The slope cannot change at a knot
- The curvature cannot change at a knot



## Boundary conditions

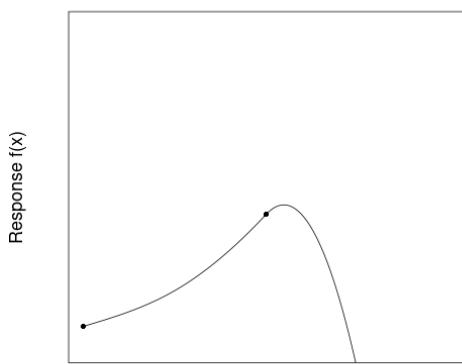


Sections need to join up smoothly.

- Both sides must go through the knot.
- The slope cannot change at a knot
- The curvature cannot change at a knot



## Boundary conditions



Sections need to join up smoothly.

- Both sides must go through the knot.
- The slope cannot change at a knot
- The curvature cannot change at a knot



## Outline

Categorization and its discontents

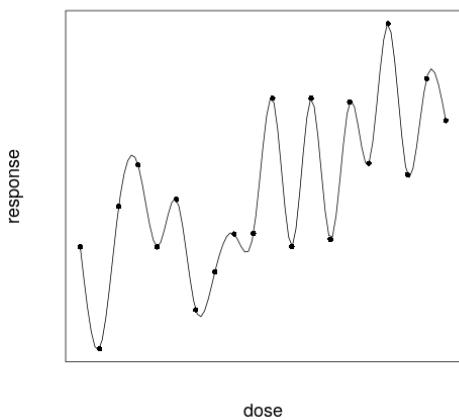
Join the dots

Smoothing splines

Splines in R

## Dose response with error

Perfect fit

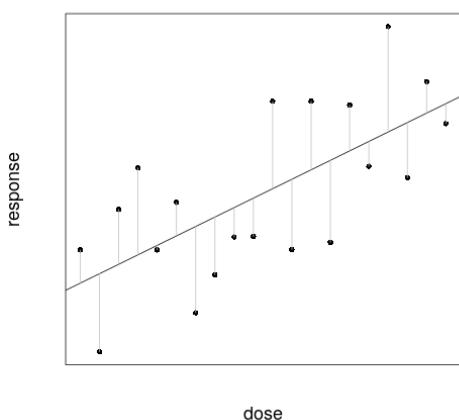


In practice we never know the dose response curve exactly at any point but always measure with error. A spline model is then a compromise between

- Model fit
- Smoothness of the spline

## Dose response with error

Perfectly smooth



In practice we never know the dose response curve exactly at any point but always measure with error. A spline model is then a compromise between

- Model fit
- Smoothness of the spline

## Fitting a smoothing spline

Minimize

$$\sum_i [y_i - f(x_i)]^2 + \lambda \int \left( \frac{\partial^2 f}{\partial x^2} \right)^2 dx$$

Or, more generally

Deviance +  $\lambda \times$  Roughness penalty

Size of tuning parameter  $\lambda$  determines compromise between model fit (small  $\lambda$ ) and smoothness (large  $\lambda$ ).

## Smoothing and degrees of freedom

Software will choose the smoothing parameter  $\lambda$  for you automatically using cross-validation.

The smoothing parameter is adapted to the data.

Smoothness of the model can be measured with the *effective degrees of freedom* (EDF)

- Linear model: maximally smooth
  - EDF=2 (intercept + slope parameter)
- Interpolating mode: best fit
  - EDF=n (one parameter for every observation)

## Outline

Categorization and its discontents

Join the dots

Smoothing splines

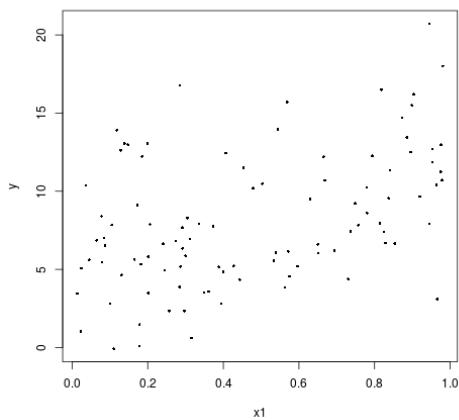
Splines in R

## Spline models in R

- Do not use the `splines` package.
- Use the `gam` function from the `mgcv` package to fit your spline models.
- The `gam` function chooses number and placement of knots for you and estimates the size of the tuning parameter  $\lambda$  automatically.
- You can use the `gam.check` function to see if you have enough knots. Also re-fit the model explicitly setting a larger number of knots (e.g. `double`) to see if the fit changes.

## Penalized spline

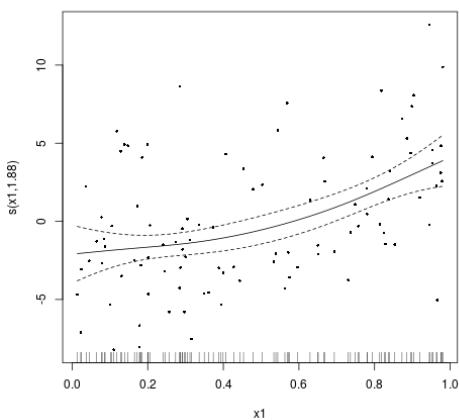
Some simulated data



- A `gam` fit to some simulated data
- Model has 9 degrees of freedom
- Smoothing reduces this to 2.88 effective degrees of freedom

## Penalized spline

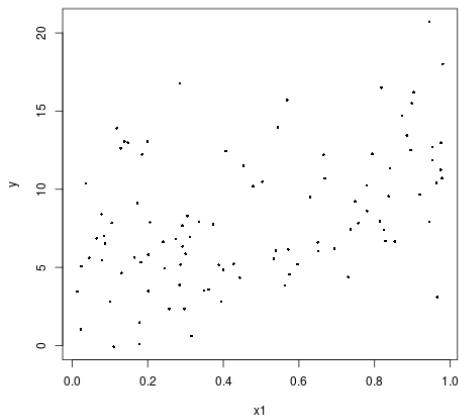
A gam fit with default options



- A `gam` fit to some simulated data
- Model has 9 degrees of freedom
- Smoothing reduces this to 2.88 effective degrees of freedom

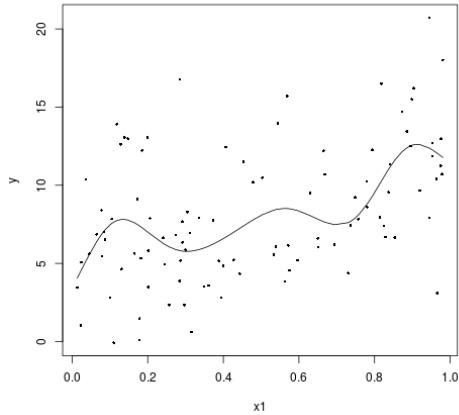
## Unpenalized spline

Some simulated data



- An unpenalized spline using the same spline basis as the gam fit.
- Model has 9 degrees of freedom

## Unpenalized spline



- An unpenalized spline using the same spline basis as the gam fit.
- Model has 9 degrees of freedom

## Conclusions

- Epidemiologists like to turn continuous variables into categories.
- Statisticians do not like categorization because it loses information.
- Splines are a flexible class of models that avoid categorization but also avoid making strong assumptions about the shape of a dose-response relationship.
- Penalized regression splines are based on compromise between goodness-of-fit and smoothness.
- Most of the decisions in fitting a penalized regression spline can be made for you
  - Degree of smoothing
  - Number of knots
  - Placement of knots

# More Advanced Graphics in R

Martyn Plummer

University of Warwick, UK

SPE 2023, Tartu



Overview of graphics systems  
○○○○○

Device handling  
○○○

Base graphics  
○○○○○○○○○○

Grid graphics  
○○○

## Outline

Overview of graphics systems

Device handling

Base graphics

Grid graphics



Overview of graphics systems  
●○○○○

Device handling  
○○○

Base graphics  
○○○○○○○○○○

Grid graphics  
○○○

## Graphics Systems in R

R has several different graphics systems:

- ▶ Base graphics (the `graphics` package)
- ▶ Lattice graphics (the `lattice` package)
- ▶ Grid graphics (the `grid` package)
- ▶ Grammar of graphics (the `ggplot2` package)

Why so many? Which one to use?



## Base Graphics

- ▶ The oldest graphics system in R.
- ▶ Based on S graphics (Becker, Chambers and Wilks, *The New S Language*, 1988)
- ▶ Implemented in the base package `graphics`
  - ▶ Loaded automatically so always available
- ▶ Ink on paper model; once something is drawn “the ink is dry” and it cannot be erased or modified.

## Grid Graphics

- ▶ A complete rewrite of the graphics system of R, independent of base graphics.
- ▶ Programming with graphics:
  - ▶ Grid graphics commands create graphical objects (Grobs)
  - ▶ Printing a Grob displays it on a graphics device
  - ▶ Functions can act on grobs to modify or combine them
- ▶ Implemented in the base package `grid`, and extended by CRAN packages `gridExtra`, `gridDebug`, ...
- ▶ Described by the package author Paul Murrell in the book *R Graphics (2nd edition)*, 2011.

## Grammar of Graphics

- ▶ Originally described by Leland Wilkinson in the book *The Grammar of Graphics*, 1999 and implemented in the statistical software nViZn (part of SPSS)
- ▶ Statistical graphics, like natural languages, can be broken down into components that must be combined according to certain rules.
- ▶ Provides a *pattern language* for graphics:
  - ▶ geometries, statistics, scales, coordinate systems, aesthetics, themes, ...
- ▶ Implemented in R in the CRAN package `ggplot2`
- ▶ Described more fully by the `ggplot2` package author Hadley Wickham in the book *ggplot2: Elegant Graphics for Data Analysis*, 2009.

## Putting It All Together

- ▶ Base graphics are the default, and are used almost exclusively in this course
- ▶ `grid` provides alternate low-level graphics functions.
  - ▶ Experts only
- ▶ `ggplot2` is an alternate, high-level graphics package built on `grid`.
- ▶ All graphics packages take time to learn...

## Graphics Devices

Graphics devices are used by all graphics systems.

- ▶ Plotting commands will draw on the current *graphics device*
- ▶ The *default* graphics device is a window on your screen:
  - In RStudio `RStudioGD()`
  - On Windows `windows()`
  - On Unix/Linux `x11()`
  - On Mac OS X `quartz()`
- It normally opens up automatically when you need it.
- ▶ You can have several graphics devices open at the same time (but only one is current)

## Graphics Device in RStudio

RStudio has its own graphics device `RStudioGD` built into the graphical user interface

- ▶ You can see the contents in a temporary, larger window by clicking the zoom button.
- ▶ You can write the contents directly to a file with the export menu
- ▶ Sometimes the small size of the `RStudioGD` device causes problems. Open up a new device calling `RStudioGD()`. This will appear in its own window, free from the GUI.

## Writing Graphs to Files

There are also non-interactive graphics devices that write to a file instead of the screen.

`pdf` produces Portable Document Format files

`win.metafile` produces Windows metafiles that can be included in Microsoft Office documents (windows only)

`postscript` produces postscript files

`png`,  `bmp`,  `jpeg` all produce bitmap graphics files

- ▶ Turn off a graphics device with `dev.off()`. Particularly important for non-interactive devices.
- ▶ Plots may look different in different devices

## Types of Plotting Functions

- ▶ High level
  - ▶ Create a new page of plots with reasonable default appearance.
- ▶ Low level
  - ▶ Draw elements of a plot on an existing page:
    - ▶ Draw title, subtitle, axes, legend ...
    - ▶ Add points, lines, text, math expressions ...
- ▶ Interactive
  - ▶ Querying mouse position (`locator`), highlighting points (`identify`)

## Base x-y Plots

- ▶ The `plot` function with one or two numeric arguments
- ▶ Scatterplot or line plot (or both) depending on `type` argument: "`l`" for `lines`, "`p`" for `points` (the default), "`b`" for `both`, plus quite a few more
- ▶ Also: formula interface, `plot(y~x)`, with arguments similar to the modeling functions like `lm`

## Customizing Plots in Base

- ▶ Most plotting functions take optional parameters to change the appearance of the plot
  - ▶ e.g., `xlab`, `ylab` to add informative axis labels
- ▶ Most of these parameters can be supplied to the `par()` function, which changes the default behaviour of subsequent plotting functions
- ▶ Look them up via `help(par)`! Here are some of the more commonly used:
  - ▶ Point and line characteristics: `pch`, `col`, `lty`, `lwd`
  - ▶ Multiframe layout: `mfrow`, `mfcol`
  - ▶ Axes: `xlim`, `ylim`, `xaxt`, `yaxt`, `log`

## Adding to Plots in Base

- ▶ `title()` add a title above the plot
- ▶ `points()`, `lines()` adds points and (poly-)lines
- ▶ `text()` text strings at given coordinates
- ▶ `abline()` line given by coefficients (*a* and *b*) or by fitted linear model
- ▶ `axis()` adds an axis to one edge of the plot region.  
Allows some options not otherwise available.

## Strategy for Customization of Base Graphics

- ▶ Start with default plots
- ▶ Modify parameters (using `par()` settings or plotting arguments)
- ▶ Add more graphics elements. Notice that there are graphics parameters that turn things *off*, e.g. `plot(x, y, xaxt="n")` so that you can add completely customized axes with the `axis` function.
- ▶ Put all your plotting commands in a script or inside a function so you can start again

## Demo 1

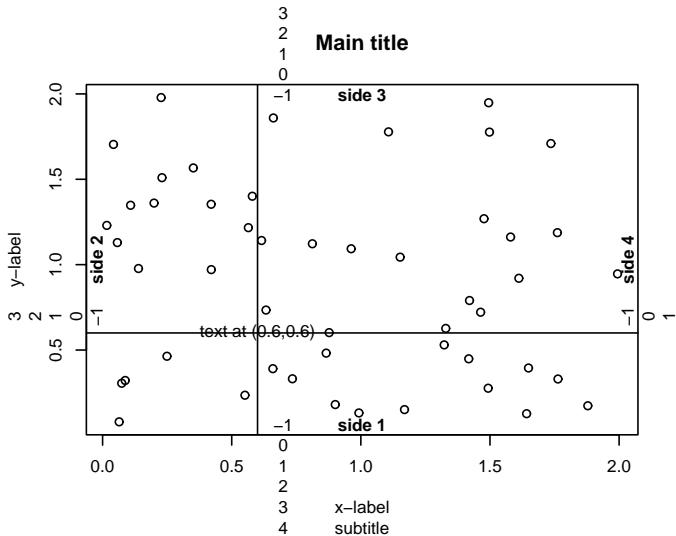
```
library(ISwR)
par(mfrow=c(2,2))
matplot(intake)
matplot(t(intake))
matplot(t(intake),type="b")
matplot(t(intake),type="b",pch=1:11,col="black",
       lty="solid", xaxt="n")
axis(1,at=1:2,labels=names(intake))
```

## Margins

- ▶ R sometimes seems to leave too much empty space around plots (especially in multi-frame layouts).
- ▶ There is a good reason for it: You might want to put something there (titles, axes).
- ▶ This is controlled by the `mar` parameter. By default, it is `c(5, 4, 4, 2)+0.1`
  - ▶ The units are *lines of text*, so depend on the setting of `pointsize` and `cex`
  - ▶ The sides are indexed in clockwise order, starting at the bottom (1=bottom, 2=left, 3=top, 4=right)
- ▶ The `mtext` function is designed to write in the margins of the plot
- ▶ There is also an *outer margin* settable via the `oma` parameter. Useful for adding overall titles etc. to multiframe plots

## Demo 2

```
x <- runif(50,0,2)
y <- runif(50,0,2)
plot(x, y, main="Main title", sub="subtitle",
      xlab="x-label", ylab="y-label")
text(0.6,0.6,"text at (0.6,0.6)")
abline(h=.6,v=.6)
for (side in 1:4)
  mtext(-1:4,side=side,at=.7,line=-1:4)
mtext(paste("side",1:4), side=1:4, line=-1,font=2)
```



A set of small, light-blue navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and table of contents.

## Overview of graphics systems

## Device handling

## Base graphics



## Grid graphics

## A Few Words on Grid Graphics

- ▶ Experts only, but ...
  - ▶ Recall that `ggplot2` uses `grid`
  - ▶ The key concepts you need are *grobs* and *viewports*

A set of small, semi-transparent navigation icons typically found in presentation software like Beamer. They include symbols for back, forward, search, and other document-related functions.

20/22

## Overview of graphics systems

## Device handling

## Base graphics



## Grid graphics

# Grobs: Graphical Objects

- ▶ Grobs are created by plotting functions in `grid`, and `ggplot2`
  - ▶ Grobs are only displayed when they are printed
  - ▶ Grobs can be modified or combined before being displayed
  - ▶ The `ggplot2` package uses the `+` operator to combine grobs representing different elements of the plot

A set of small, light-blue navigation icons typically found in presentation software like Beamer. They include symbols for back, forward, search, and table of contents.

## Viewports

- ▶ The plotting region is divided into viewports
- ▶ Grobs are displayed inside a viewport
  - ▶ Viewports can be different sizes (inches, centimetres, lines of text, or relative units)
  - ▶ Each viewport may have its own coordinate systems

## Survival analysis with competing risks

Janne Pitkäniemi

Finnish Cancer Registry  
Tampere University  
Statistical Practice in Epidemiology (2019, Tartu)

1 / 29

## Points to be covered

1. Survival or time to event data & censoring.
  2. Competing risks: event-specific cumulative incidences & hazards.
  3. Kaplan–Meier and Aalen–Johansen estimators.
  4. Regression modelling of hazards: Cox model.
  5. Packages `survival`, `mstate`, (`cmprisk`).
  6. Functions `Surv()`, `survfit()`, `plot.survfit()`, `coxph()`.

## Survival time – time to event

**Time** spent (`lex.dur`) in a given **state** (`lex.Cst`) from its beginning till a certain *endpoint* or *outcome event* (`lex.Xst`) or *transition* occurs, changing the state to another.

Examples of such times and outcome events:

- ▶ lifetime: birth → death,
  - ▶ duration of marriage: wedding → divorce,
  - ▶ healthy exposure time:  
start of exposure → onset of disease,
  - ▶ clinical survival time:  
diagnosis of a disease → death.

3 / 29

## Ex. Survival of 338 oral cancer patients

Important variables:

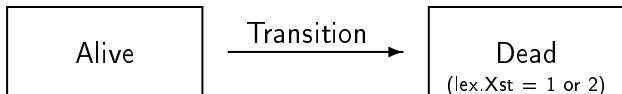
- ▶ time = duration of patientship from diagnosis (**entry**) till death (death) or censoring (Alive), (lex.Cst is (Alive))
- ▶ event = indicator for the outcome and its observation at the end of follow-up (**exit**):  
0 = censoring,  
1 = death from oral cancer

Special features:

- ▶ Two possible endpoints
- ▶ Censoring – incomplete observation of the survival time.

## Set-up of classical survival analysis

- ▶ **Two-state model:** only one type of event changes the initial state.
- ▶ Major applications: analysis of lifetimes since birth and of survival times since diagnosis of a disease until death from any cause.



- ▶ **Censoring:** Death and final lifetime not observed for some subjects due to emigration or closing the follow-up while they are still alive

## Distribution concepts: hazard function

The **hazard rate** or **intensity** function  $\lambda(t)$

$$\lambda(t) = P(t < T \leq t + \Delta | T > t) / \Delta, \text{ for small } \Delta$$

- ≈ the conditional probability that the event occurs in a short interval  $(t, t + \Delta]$ , given that it does not occur before  $t$ , divided by interval length.

In other words, during a short interval

$$\text{risk of event} \approx \text{hazard} \times \text{interval length}$$

## Distribution concepts: survival and cumulative hazard functions

### Survival function

$$S(t) = P(T > t),$$

= probability of avoiding the event at least up to  $t$  (the event occurs only after  $t$ ).

The **cumulative hazard** (or integrated intensity):

$$\Lambda(t) = \int_0^t \lambda(u)du$$

Connections between the functions:

$$S(t) = \exp\{-\Lambda(t)\}$$

## Observed data on survival times

For individuals  $i = 1, \dots, n$  let

$T_i$  = time to outcome event,

$U_i$  = time to censoring.

Censoring is assumed **noninformative**, i.e.

independent from occurrence of events.

We observe

$y_i = \min\{T_i, U_i\}$ , i.e. the exit time, and

$\delta_i = 1_{\{T_i < U_i\}}$ , indicator (1/0) for the outcome event occurring first, before censoring.

Censoring must properly be taken into account in the statistical analysis.

## Approaches for analysing survival time

- ▶ **Parametric model** (like Weibull, gamma, etc.) on hazard rate  $\lambda(t) \rightarrow$  Likelihood:

$$L = \prod_{i=1}^n \lambda(y_i)^{\delta_i} S(y_i)$$

- ▶ **Piecewise constant rate** model on  $\lambda(t)$   
– see Bendix's lecture on time-splitting (Poisson likelihood).
- ▶ **Non-parametric** methods, like  
Kaplan–Meier (KM) estimator of survival curve  $S(t)$  and Cox proportional hazards model on  $\lambda(t)$ .

# R package survival

Tools for analysis with one outcome event.

- ▶ `Surv(time, event) -> sobj`  
creates a **survival object** `sobj` assuming that all start at 0, containing pairs  $(y_i, \delta_i)$ ,
  - ▶ `Surv(entry, exit, event) -> sobj2`  
creates a survival object from entry and exit times,
  - ▶ `survfit(sobj ~ x) -> sfo`  
creates a **survfit** object `sfo` containing KM or other non-parametric estimates (also from a fitted Cox model),
  - ▶ `plot(sfo)`  
plot method for survival curves and related graphs,
  - ▶ `coxph(sobj ~ x1 + x2)`  
fits a Cox model with covariates `x1` and `x2`.
  - ▶ `survreg()` – parametric survival models.

A set of small, semi-transparent navigation icons typically found in presentation software like Beamer. They include symbols for back, forward, search, and table of contents.

10 / 29

10 / 29

## Ex. Oral cancer data (cont'd)

```

> orca$suob <- Surv(orca$time, 1*(orca$event > 0) )
> orca$suob[1:7] # + indicates censored observations
[1] 5.081+ 0.419 7.915 2.480 2.500 0.167 5.925+
> km1 <- survfit( suob ~ 1, data = orca)
> km1
# brief summary
Call: survfit(formula = suob ~ 1, data = orca)

```

```
[1,] 338    229   5.42    4.33   6.92
> summary(km1)      # detailed KM-estimate
Call: survfit(formula = suob ~ 1, data = orca)

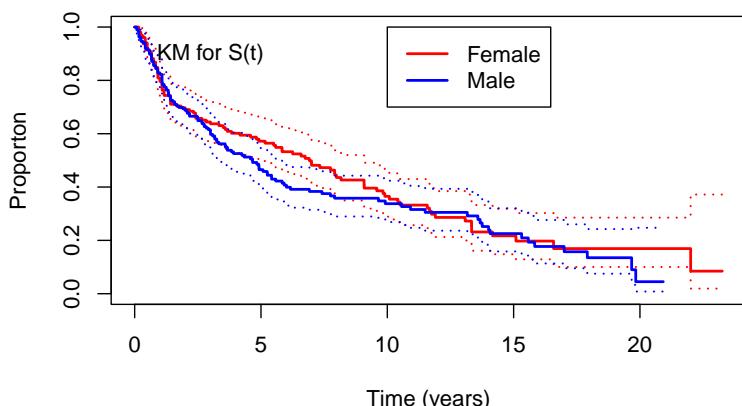
  time n.risk n.event survival std.err lower 95% CI upper 95% CI
  0.085  338      2  0.9941  0.00417   0.9859  1.0000
  0.162  336      2  0.9882  0.00588   0.9767  1.0000
  0.167  334      4  0.9763  0.00827   0.9603  0.9993
  0.170  330      2  0.9704  0.00922   0.9525  0.9899
  0.246  328      1  0.9675  0.00965   0.9487  0.9871
  0.249  327      1  0.9645  0.01007   0.9450  0.9841
  0.252  326      3  0.9556  0.01120   0.9339  0.9781
  0.329  323      1  0.9527  0.01155   0.9303  0.9761
  0.334  322      1  0.9497  0.01189   0.9267  0.9731
  0.413  321      1  0.9467  0.01221   0.9231  0.9711
```

11 / 29

11 / 29

## Oral cancer: Kaplan-Meier estimates

### **Estimated survival (95% CI) and CDF**

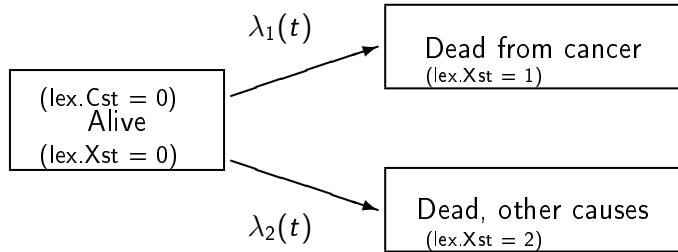


A set of small, light-blue navigation icons typically found in presentation software like Beamer. They include symbols for back, forward, search, and table of contents.

12 / 39

## Competing risks model: causes of death

- ▶ Often the interest is focused on the risk or hazard of dying from one specific cause.
  - ▶ That cause may eventually not be realized, because a **competing cause** of death hits first.



- ▶ Generalizes to several competing causes.

13 / 29

13 / 29

## Competing events & competing risks

In many epidemiological and clinical contexts there are competing events that may occur before the target event and remove the person from the population at risk for the event, e.g.

- ▶ *target event*: occurrence of endometrial cancer, *competing events*: hysterectomy or death.
  - ▶ *target event*: relapse of a disease (ending the state of remission),  
*competing event*: death while still in remission.
  - ▶ *target event*: divorce,  
*competing event*: death of either spouse.

14 / 29

## Event-specific quantities

## Cumulative incidence function (CIF) or

$$F_c(t) = P(T \leq t \text{ and } C = c), \quad c = 1, 2,$$

From these one can recover

- ▶  $F(t) = \sum_c F_c(t)$ , CDF of event-free survival time  $T$ , i.e. cumulative risk of any event by  $t$ .
  - ▶  $S(t) = 1 - F(t)$ , **event-free survival function**, i.e. probability of avoiding all events by  $t$ , but  $S(t) \neq F_1(t) + F_2(t)$

15 / 29

15 / 29

## Event-specific quantities (cont'd)

### Event- or cause-specific hazard function

$$\begin{aligned}\lambda_c(t) &= \lim_{\Delta \rightarrow 0} \frac{P(t < T \leq t + \Delta \text{ and } C = c \mid T > t)}{\Delta} \\ &= \frac{f_c(t)}{1 - F(t)}\end{aligned}$$

CIF = risk of event  $c$  over risk period  $[0, t]$  in the presence of competing risks,  
also obtained

$$F_c(t) = \int_0^t \lambda_c(v) S(v) dv, \quad c = 1, 2,$$

More on the technical definitions of relevant quantities:  
<http://bendixcarstensen.com/AdvCoh/papers/fundamentals.pdf>

## Warning of “net risk” and “cause-specific survival”

- ▶ The “**net risk**” of outcome  $c$  by time  $t$ , assuming hypothetical elimination of competing risks, is often defined as
$$F_1^*(t) = 1 - S_1^*(t) = 1 - \exp\{-\Lambda_1(t)\} \neq S(t)$$
- ▶ In clinical survival studies, function  $S_1^*(t)$  is often called “**cause-specific survival**”, or “**net survival**”
- ▶ Yet, these \*-functions,  $F_1^*(t)$  and  $S_1^*(t)$ , lack proper probability interpretation when competing risks exist.
- ▶ Hence, their use should be viewed critically (Andersen & Keiding, *Stat Med*, 2012)

## Analysis with competing events

Let  $U_i$  = censoring time,  $T_i$  = time to first event, and  $C_i$  = variable for event 1 or 2. We observe

- ▶  $y_i = \min\{T_i, U_i\}$ , i.e. the exit time, and
- ▶  $\delta_{ic} = 1_{\{T_i < U_i \text{ & } C_i=c\}}$ , indicator (1/0) for event  $c$  being first observed,  $c = 1, 2$ .

Non-parametric estimation of CIF

- ▶ Let  $t_1 < t_2 < \dots < t_K$  be the  $K$  distinct time points at which any outcome event was observed,  
Let also  $\tilde{S}(t)$  be KM estimator for overall  $S(t)$ .
- ▶ **Aalen-Johansen estimator (AJ)** for the cumulative incidence function  $F(t)$  should be used

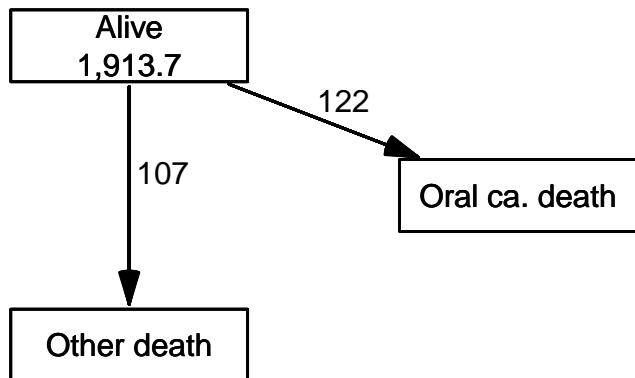
# R tools for competing risks analysis

- ▶ `survfit( Surv(...,type="mstate") )` in Survival-package can be fitted for any transition of a multistate model and to obtain A-J estimates.
  - ▶ Package `cmprsk` – `cuminc(ftime, fstatus, ...)` computes CIF-estimates, and can be compared in more than two samples.  
`plot.cuminc()` plots them.
  - ▶ Package Epi – Lexis tools for multistate analyses  
Will be advertised by Bendix!

19 / 29

## Box diagram for transitions

NOTE: entry.status has been set to "Alive" for all.  
NOTE: entry is assumed to be 0 on the stime timescale.

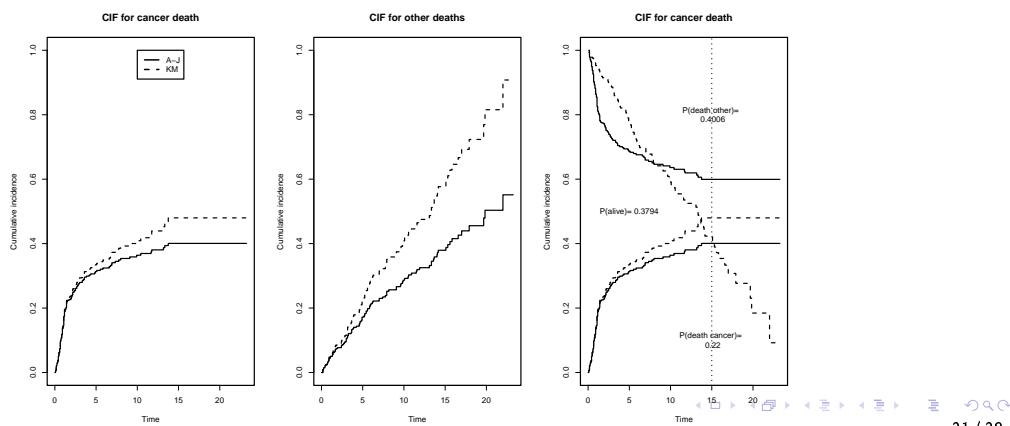


◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ 20 / 29

22 / 22

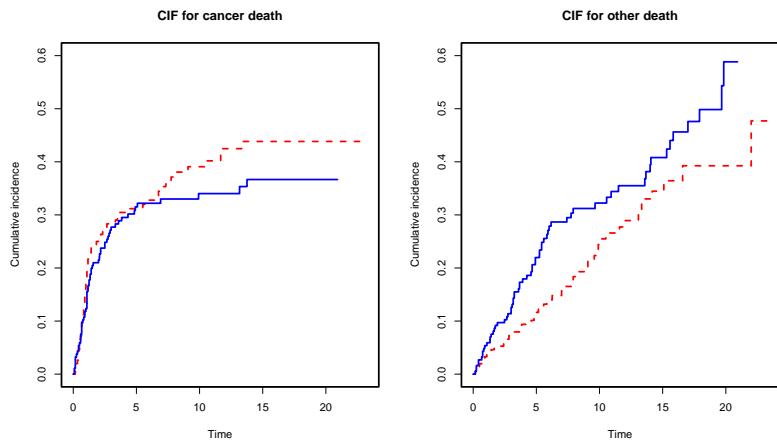
## Ex. Survival from oral cancer

- ▶ AJ-estimates of CIFs (solid) for both causes.
  - ▶ Naive KM-estimates of CIF (dashed) > AJ-estimates
  - ▶ CIF curves may also be stacked (right).



21 / 29

Ex. CIFs by cause in men and women



CIF for cancer higher in women (chance?) but for other causes higher in men (no surprise).

22 / 29

22 / 29

# Regression models for time-to-event data

Regression models for hazards can be defined e.g. for

(a) hazards, multiplicatively:

$$\lambda_i(t) = \lambda_0(t; \alpha) r(\eta_i), \quad \text{where}$$

$\lambda_0(t; \alpha)$  = baseline hazard and

$r(\eta_i)$  = relative rate function, typically  $\exp(\eta_i)$

(b) hazards, additively:

$$\lambda_i(t) = \lambda_0(t; \alpha) + \eta_i.$$

23 / 29

## Relative hazards model or Cox model

In model (b), the baseline hazard  $\lambda_0(t, \alpha)$  may be given a parametric form (e.g. Weibull) or a piecewise constant rate (exponential) structure.

Often a parameter-free form  $\lambda_0(t)$  is assumed. Then

$$\lambda_i(t) = \lambda_0(t) \exp(\eta_1),$$

specifies the **Cox model** or the **semiparametric proportional hazards model**.

$\eta_i = \beta_1 x_{i1} + \cdots + \beta_p x_{ip}$  not depending on time.

Generalizations: **time-dependent** covariates  $x_{ij}(t)$

24 / 29

## PH model: interpretation of parameters

Present the model explicitly in terms of  $x$ 's and  $\beta$ 's.

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_1 x_{i1} + \cdots + \beta_p x_{ip})$$

Consider two individuals,  $i$  and  $i'$ , having the same values of all other covariates except the  $j^{\text{th}}$  one.

The ratio of hazards is constant:

$$\frac{\lambda_i(t)}{\lambda_{i'}(t)} = \frac{\exp(\eta_i)}{\exp(\eta_{i'})} = \exp\{\beta_j(x_{ij} - x_{i'j})\}.$$

Thus  $e^{\beta_j} = \text{HR}_j = \text{hazard ratio}$  or relative rate associated with a unit change in covariate  $X_j$ .

## Ex. Total mortality of oral ca. patients

Fitting Cox models with sex and sex + age.

```
> cm0 <- coxph(suob ~ sex, data = orca)
> summary(cm0)
    coef exp(coef) se(coef)   z Pr(>|z|)
sexMale 0.126     1.134    0.134 0.94    0.35
          exp(coef) exp(-coef) lower .95 upper .95
sexMale     1.13      0.882     0.872      1.47

> cm1 <- coxph(suob ~ sex + age, data = orca)
> summary(cm1)
    exp(coef) exp(-coef) lower .95 upper .95
sexMale     1.49      0.669     1.14     1.96
age         1.04      0.960     1.03     1.05
```

The M/F contrast visible only after age-adjustment.

## Predictions from the Cox model

- ▶ Individual survival *times* cannot be predicted but ind'l survival *curves* can.  
PH model implies:

$$S_i(t) = [S_0(t)]^{\exp(\beta_1 x_{i1} + \cdots + \beta_p x_{ip})}$$

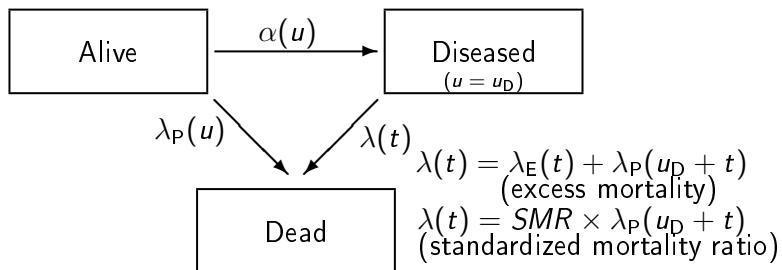
- ▶ Having estimated  $\beta$  by partial likelihood, the baseline  $S_0(t)$  is estimated by Breslow method
- ▶ From these, a survival curve for an individual with given covariate values is predicted.
- ▶ In R: `pred <- survfit(mod, newdata=...)` and `plot(pred)`, where `mod` is the fitted `coxph` object, and `newdata` specifies the covariate values. `newdata` is always needed for predictions.

## Modelling with competing risks

Main options, providing answers to different questions.

- (a) Cox model for event-specific hazards  $\lambda_c(t) = f_c(t)/[1 - F(t)]$ , when e.g.  
the interest is in the biological effect of the prognostic factors on the fatality  
of the very disease that often leads to the relevant outcome.
- (b) **Fine–Gray model** for the hazard of the subdistribution  
 $\gamma_c(t) = f_c(t)/[1 - F_c(t)]$  when we want to assess the impact of the factors  
on the overall cumulative incidence of event  $c$ .  
– Function `crr()` in package `cmprsk`.

## Competing risks model: excess hazard of death



where

- $\lambda_P(u)$  is the hazard of dying from any cause among disease-free members
- $\lambda_E(t)$  is the excess hazard of dying from the disease among diseased cohort members

# Representation of follow-up

**Bendix Carstensen** Steno Diabetes Center Copenhagen  
Herlev, Denmark  
<http://BendixCarstensen.com>

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

From C:\Bendix\teach\SPE\git\lectures\time-rep\time-rep.tex

Tuesday 30 May, 2023, 08:43

1 / 40

# Representation of follow-up

**Bendix Carstensen**

Representation of follow-up

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

time-split

- ▶ In follow-up studies we estimate rates from:
  - ▶  $D$  — events, deaths
  - ▶  $Y$  — person-years
  - ▶  $\hat{\lambda} = D/Y$  rates
  - ▶ ... empirical counterpart of intensity — an **estimate**
- ▶ Rates differ between persons.
- ▶ Rates differ **within** persons:
  - ▶ by age
  - ▶ by calendar time
  - ▶ by disease duration
  - ▶ ...
- ▶ Multiple timescales.
- ▶ Multiple states (little boxes — later)

## Representation of follow-up data

A cohort or follow-up study records **events** and **risk time**

The outcome is thus **bivariate**:  $(d, y)$

Follow-up **data** for each individual must therefore have (at least) three pieces of information recorded:

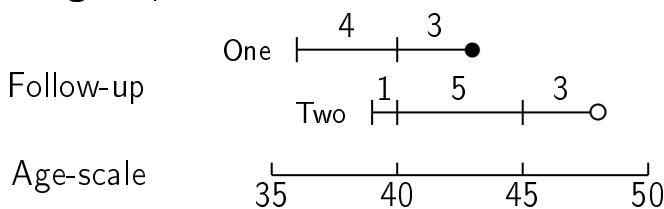
Date of entry	entry	date variable
Date of exit	exit	date variable
Status at exit	fail	indicator (mostly 0/1)

These are specific for each **type** of outcome.

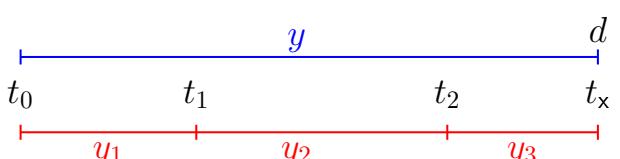
## Stratification by age

If follow-up is rather short, age at entry is OK for age-stratification.

If follow-up is long, stratification by categories of **current age** is preferable.



- allowing rates to vary across age-bands
- how do we do the split and why is it OK?



Probability

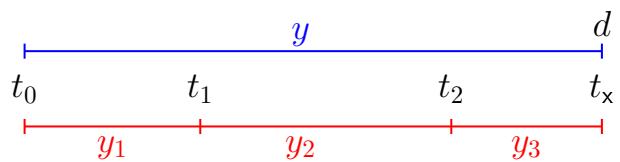
$$P(d \text{ at } t_x \mid \text{entry } t_0)$$

$$\begin{aligned} &= P(\text{surv } t_0 \rightarrow t_1 \mid \text{entry } t_0) \\ &\times P(\text{surv } t_1 \rightarrow t_2 \mid \text{entry } t_1) \\ &\times P(d \text{ at } t_x \mid \text{entry } t_2) \end{aligned}$$

log-Likelihood

$$d \log(\lambda) - \lambda y$$

$$\begin{aligned} &= 0 \log(\lambda) - \lambda y_1 \\ &+ 0 \log(\lambda) - \lambda y_2 \\ &+ d \log(\lambda) - \lambda y_3 \end{aligned}$$



Probability log-Likelihood

$$\begin{aligned}
 & P(d \text{ at } t_x \mid \text{entry } t_0) && d \log(\lambda) - \lambda y \\
 & = P(\text{surv } t_0 \rightarrow t_1 \mid \text{entry } t_0) && = 0 \log(\lambda_1) - \lambda_1 y_1 \\
 & \times P(\text{surv } t_1 \rightarrow t_2 \mid \text{entry } t_1) && + 0 \log(\lambda_2) - \lambda_2 y_2 \\
 & \times P(d \text{ at } t_x \mid \text{entry } t_2) && + d \log(\lambda_3) - \lambda_3 y_3
 \end{aligned}$$

— allows different rates ( $\lambda_i$ ) in each interval

Representation of follow-up (time-split)

6 / 40

## Dividing time into bands requires:

**Origin:** The date where the time scale is 0:

- ▶ Age — 0 at date of birth
- ▶ Disease duration — 0 at date of diagnosis
- ▶ Occupation exposure — 0 at date of hire

**Intervals:** How should it be subdivided:

- ▶ 1-year classes? 5-year classes?
- ▶ Equal length?

**Aim:** Separate rate in each interval, mimicking continuous time by using small intervals:

— time at the beginning of interval as quantitative variable.

Representation of follow-up (time-split)

7 / 40

## Example: cohort with 3 persons:

Id	Bdate	Entry	Exit	St
1	14/07/1952	04/08/1965	27/06/1997	1
2	01/04/1954	08/09/1972	23/05/1995	0
3	10/06/1987	23/12/1991	24/07/1998	1

- ▶ Age bands: 10-years intervals of current age.
- ▶ Split  $Y$  for every subject accordingly
- ▶ Treat each segment as a separate unit of observation.
- ▶ Keep track of exit status ( $D$ ) in each interval.

Representation of follow-up (time-split)

8 / 40

## Splitting the follow-up

	subj. 1	subj. 2	subj. 3
Age at Entry:	13.06	18.44	4.54
Age at eXit:	44.95	41.14	11.12
Status at exit:	Dead	Alive	Dead
<i>Y</i>	31.89	22.70	6.58
<i>D</i>	1	0	1

Representation of follow-up (time-split)

9 / 40

Age	subj. 1		subj. 2		subj. 3		$\sum$	
	<i>Y</i>	<i>D</i>	<i>Y</i>	<i>D</i>	<i>Y</i>	<i>D</i>	<i>Y</i>	<i>D</i>
0–	0.00	0	0.00	0	5.46	0	5.46	0
10–	6.94	0	1.56	0	1.12	1	8.62	1
20–	10.00	0	10.00	0	0.00	0	20.00	0
30–	10.00	0	10.00	0	0.00	0	20.00	0
40–	4.95	1	1.14	0	0.00	0	6.09	1
$\sum$	31.89	1	22.70	0	6.58	1	60.17	2

Representation of follow-up (time-split)

10 / 40

## Splitting the follow-up

id	Bdate	Entry	Exit	St	risk	int
1	14/07/1952	03/08/1965	14/07/1972	0	6.9432	10
1	14/07/1952	14/07/1972	14/07/1982	0	10.0000	20
1	14/07/1952	14/07/1982	14/07/1992	0	10.0000	30
1	14/07/1952	14/07/1992	27/06/1997	1	4.9528	40
2	01/04/1954	08/09/1972	01/04/1974	0	1.5606	10
2	01/04/1954	01/04/1974	31/03/1984	0	10.0000	20
2	01/04/1954	31/03/1984	01/04/1994	0	10.0000	30
2	01/04/1954	01/04/1994	23/05/1995	0	1.1417	40
3	10/06/1987	23/12/1991	09/06/1997	0	5.4634	0
3	10/06/1987	09/06/1997	24/07/1998	1	1.1211	10

Keeping track of calendar time too?

Representation of follow-up (time-split)

11 / 40

## Follow-up intervals on several timescales

- ▶ The risk-time is the same on all timescales
- ▶ Only need the entry point on each time scale:
  - ▶ Age at entry.
  - ▶ Date of entry.
  - ▶ Time since treatment at entry.
    - if time of treatment is the entry, this is 0 for all.
- ▶ Response variable in analysis of rates:

$$(d, y) \quad (\text{event}, \text{duration})$$

- ▶ Covariates in analysis of rates:
  - ▶ timescales
  - ▶ other (fixed) measurements
- ▶ ... do not confuse **duration** and **timescale** !

Representation of follow-up (time-split)

12 / 40

## Follow-up data in Epi — Lexis objects I

```
> thoro[1:6,1:8]
   id sex birthdat contrast injecdat volume exitdat exitstat
1  1   2 1916.609        1 1938.791     22 1976.787      1
2  2   2 1927.843        1 1943.906     80 1966.030      1
3  3   1 1902.778        1 1935.629     10 1959.719      1
4  4   1 1918.359        1 1936.396     10 1977.307      1
5  5   1 1902.931        1 1937.387     10 1945.387      1
6  6   2 1903.714        1 1937.316     20 1944.738      1

> thL <- Lexis(entry = list(age = injecdat - birthdat,
+                           dte = injecdat,
+                           tfi = 0 ),
+                  exit = list(dte = exitdat),
+                  exit.status = as.numeric(exitstat == 1),
+                  data = thoro)
```

Representation of follow-up (time-split)

13 / 40

## Follow-up data in Epi — Lexis objects II

```
NOTE: entry.status has been set to 0 for all.
NOTE: Dropping 2 rows with duration of follow up < tol
> summary(thL, timeScales = TRUE)
Transitions:
  To
From 0   1 Records: Events: Risk time: Persons:
  0 504 1964    2468    1964  51934.08    2468

Timescales:
age dte tfi
" " " "
```

Representation of follow-up (time-split)

14 / 40

## Definition of Lexis object

```
thL <- Lexis(entry = list(age = injecdat-birthdat,
                           dte = injecdat,
                           tfi = 0),
               exit = list(dte = exitdat),
               exit.status = as.numeric(exitstat==1),
               data = thoro)
```

`entry` is defined on **three** timescales,  
but `exit` is only needed on **one** timescale (or vice versa):

Follow-up time is the same on all timescales: `exitdat - injecdat`

One element of entry and exit must have same name (`dte`).

Representation of follow-up (time-split)

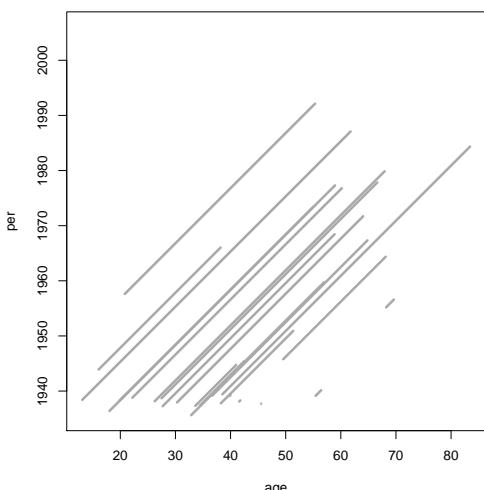
15 / 40

## The looks of a Lexis object

```
> thL[1:4,1:9]
   age      dte tfi lex.dur lex.Cst lex.Xst lex.id
1 22.18 1938.79  0   37.99      0       1       1
2 49.54 1945.77  0   18.59      0       1       2
3 68.20 1955.18  0    1.40      0       1       3
4 20.80 1957.61  0   34.52      0       0       4
...
> summary(thL)
Transitions:
  To
From  0    1 Records: Events: Risk time: Persons:
  0 504 1964     2468    1964    51934.08     2468
```

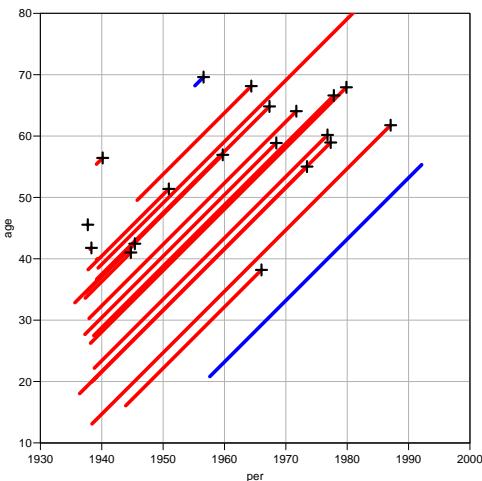
Representation of follow-up (time-split)

16 / 40



```
> plot( thL, lwd=3 )
Representation of follow-up (time-split)
```

17 / 40



## Lexis diagram

```
> plot( thL, 2:1, lwd=5, col=c("red","blue")[thL$contrast],
+       grid=TRUE, lty.grid=1, col.grid=gray(0.7),
+       xlim=1930+c(0,70), xaxs="i", ylim= 10+c(0,70), yaxs="i", las=1 )
> points( thL, 2:1, pch=c(NA,3)[thL$lex.Xst+1],lwd=3, cex=1.5 )
```

Representation of follow-up (time-split)

18 / 40

### EINLEITUNG

IN DIE

### THEORIE

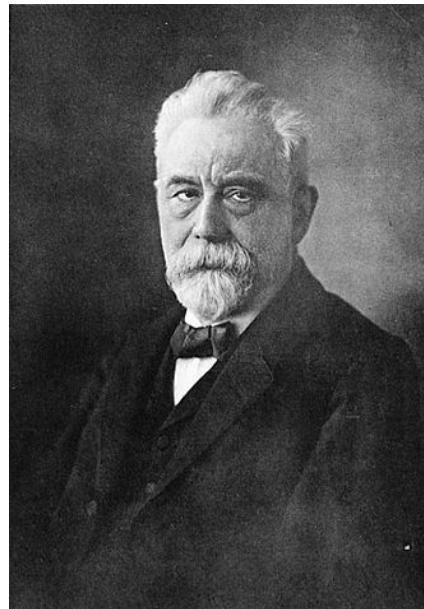
DER

### BEVÖLKERUNGSSTATISTIK

VON

W. LEXIS

DR. DER STAATSWISSENSCHAFTEN UND DER PHILOSOPHIE,  
O. PROFESSOR DER STATISTIK IN DORPAT.



18 / 40

Representation of follow-up (time-split)<sup>1875</sup>

## Splitting follow-up time

```
> spl1 <- splitLexis( thL, time.scale="age", breaks=seq(0,100,20) )
> round(spl1,1)
   age     dte    tfi lex.dur lex.Cst lex.Xst    id sex birthdat contrast injecdat vo
1 22.2 1938.8  0.0   17.8      0     0  1   2 1916.6      1 1938.8
2 40.0 1956.6 17.8   20.0      0     0  1   2 1916.6      1 1938.8
3 60.0 1976.6 37.8    0.2      0     1  1   2 1916.6      1 1938.8
4 49.5 1945.8  0.0   10.5      0     0  640 2 1896.2      1 1945.8
5 60.0 1956.2 10.5    8.1      0     1  640 2 1896.2      1 1945.8
6 68.2 1955.2  0.0    1.4      0     1 3425 1 1887.0      2 1955.2
7 20.8 1957.6  0.0   19.2      0     0 4017 2 1936.8      2 1957.6
8 40.0 1976.8 19.2   15.3      0     0 4017 2 1936.8      2 1957.6
...
...
```

Representation of follow-up (time-split)

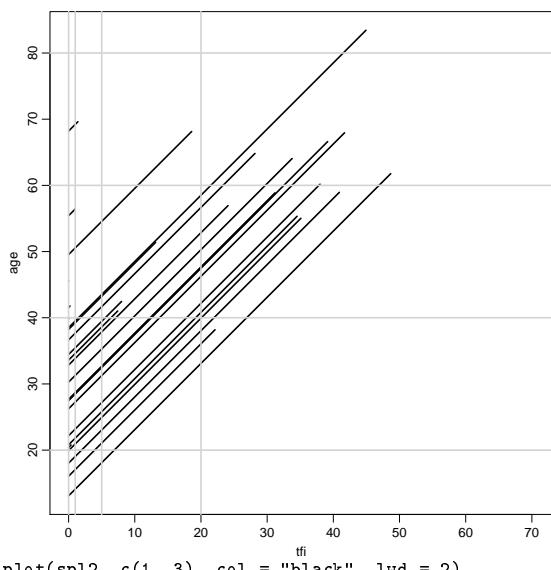
19 / 40

## Split on another timescale

```
> spl2 <- splitLexis( spl1, time.scale="tfi", breaks=c(0,1,5,20,100) )
> round( spl2, 1 )
   lex.id age  dte  tfi lex.dur lex.Cst lex.Xst    id sex birthdat contrast inj
 1      1 22.2 1938.8  0.0     1.0     0     0     1   2 1916.6     1 19
 2      1 23.2 1939.8  1.0     4.0     0     0     1   2 1916.6     1 19
 3      1 27.2 1943.8  5.0    12.8     0     0     1   2 1916.6     1 19
 4      1 40.0 1956.6 17.8     2.2     0     0     1   2 1916.6     1 19
 5      1 42.2 1958.8 20.0    17.8     0     0     1   2 1916.6     1 19
 6      1 60.0 1976.6 37.8     0.2     0     0     1   1   2 1916.6     1 19
 7      2 49.5 1945.8  0.0     1.0     0     0     0  640   2 1896.2     1 19
 8      2 50.5 1946.8  1.0     4.0     0     0     0  640   2 1896.2     1 19
 9      2 54.5 1950.8  5.0     5.5     0     0     0  640   2 1896.2     1 19
10     2 60.0 1956.2 10.5     8.1     0     0     1  640   2 1896.2     1 19
11     3 68.2 1955.2  0.0     1.0     0     0     0 3425   1 1887.0     2 19
12     3 69.2 1956.2  1.0     0.4     0     0     1 3425   1 1887.0     2 19
13     4 20.8 1957.6  0.0     1.0     0     0     0 4017   2 1936.8     2 19
14     4 21.8 1958.6  1.0     4.0     0     0     0 4017   2 1936.8     2 19
15     4 25.8 1962.6  5.0    14.2     0     0     0 4017   2 1936.8     2 19
16     4 40.0 1976.8 19.2     0.8     0     0     0 4017   2 1936.8     2 19
17     4 40.8 1977.6 20.0    14.5     0     0     0 4017   2 1936.8     2 19
```

Representation of follow-up (time-split)

20 / 40



```
plot(spl2, c(1, 3), col = "black", lwd = 2)
```

Representation of follow-up (time-split)

21 / 40

age	tfi	lex.dur	lex.Cst	lex.Xst
22.2	0.0	1.0	0	0
23.2	1.0	4.0	0	0
27.2	5.0	12.8	0	0
40.0	17.8	2.2	0	0
42.2	20.0	17.8	0	0
60.0	37.8	0.2	0	1

## Splitting on several timescales

```
> spl1 <- splitLexis(thL, time.scale = "age", breaks = seq(0, 100, 20))
> spl2 <- splitLexis(spl1, time.scale = "tfi", breaks = c(0, 1, 5, 20, 100))
> summary(spl2)
```

Transitions:

To	From	0	1	Records:	Events:	Risk time:	Persons:
	0	8250	1964	10214	1964	51934.08	2468

```
> library(popEpi)
> splx <- splitMulti(thL, age = seq(0, 100, 20), tfi = c(0, 1, 5, 20, 100))
> summary(splx)
```

Transitions:

To	From	0	1	Records:	Events:	Risk time:	Persons:
	0	8248	1964	10212	1964	51916.98	2468

```
> # NOTE: splitMulti excludes follow-up outside range of breaks
```

Representation of follow-up (time-split)

22 / 40

## Likelihood for time-split data

- ▶ We assume that rates are constant in each (small) intervals
- ▶ Each observation in the dataset represents an interval, contributing a term to the (log-)likelihood for the rate
- ▶ Each **term** looks like a contribution from a Poisson variate (albeit with values only 0 or 1)
- ▶ So the likelihood from a single **person** looks like the likelihood from several independent Poisson variates
- ▶ ... but the data are neither independent nor Poisson

Representation of follow-up (time-split)

23 / 40

## Analysis of time-split data

Observations (records) classified by  $p$ —person and  $i$ —interval

- ▶  $d_{pi}$  — events in the variable: `lex.Xst & lex.Xst!=lex.Cst`
- ▶  $y_{pi}$  — risk time: `lex.dur` (duration)
- ▶ Covariates are:
  - ▶ timescales (age, period, time in study)
  - ▶ other variables for this person (constant in each interval).
- ▶ Likelihood for rates for one person is identical to a Poisson likelihood for many independent Poisson variates
- ▶ Modeling rates using `glm` or `gam`:  
time-scales and other covariates are treated alike

Representation of follow-up (time-split)

24 / 40

## Fitting a simple model—data:

```
> stat.table(contrast,
+             list(D = sum(lex.Xst),
+                  Y = sum(lex.dur),
+                  Rate = ratio(lex.Xst, lex.dur, 100)),
+             margin = TRUE,
+             data = spl2)

-----
contrast      D      Y    Rate
-----
1        928.00 20094.74   4.62
2       1036.00 31839.35   3.25
Total     1964.00 51934.08   3.78
-----
```

Representation of follow-up (time-split)

25 / 40

## Fitting a simple model

```
-----  
contrast      D      Y    Rate  
-----  
1           928.00 20094.74   4.62  
2          1036.00 31839.35   3.25  
-----
```

```
> m0 <- glm((lex.Xst==1) ~ factor(contrast) - 1,  
+             offset = log(lex.dur / 100),  
+             family = poisson,  
+             data = spl2)  
> round(ci.exp(m0), 2)  
              exp(Est.) 2.5% 97.5%  
factor(contrast)1      4.62 4.33 4.93  
factor(contrast)2      3.25 3.06 3.46
```

... a Poisson model for mortality using log-person-years as offset

Representation of follow-up (time-split)

26 / 40

## Fitting a simple model

```
-----  
contrast      D      Y    Rate  
-----  
1           928.00 20094.74   4.62  
2          1036.00 31839.35   3.25  
-----
```

```
> m0 <- glm(cbind(lex.Xst, lex.dur / 100) ~ factor(contrast) - 1,  
+             family = poisreg,  
+             data = spl2)  
> round(ci.exp(m0), 2)  
              exp(Est.) 2.5% 97.5%  
factor(contrast)1      4.62 4.33 4.93  
factor(contrast)2      3.25 3.06 3.46
```

... a Poisson model for mortality rates based on deaths and person-years

Representation of follow-up (time-split)

27 / 40

## Fitting a simple model

The wrapper `glm.Lexis` requires that `lex.Cst` and `lex.Xst` are factors  
—use `factorize` to make them:

```
> splf <- factorize(spl2)  
> m0 <- glm.Lexis(splf, ~ factor(contrast) - 1, scale = 100)  
stats::glm Poisson analysis of Lexis object splf with log link:  
Rates for the transition:  
0->1  
, lex.dur (person-time) scaled by 100  
> round(ci.exp(m0), 2)  
              exp(Est.) 2.5% 97.5%  
factor(contrast)1      4.62 4.33 4.93  
factor(contrast)2      3.25 3.06 3.46
```

... a Poisson model for mortality rates based on deaths and person-years in a  
`Lexis` object

Representation of follow-up (time-split)

28 / 40

## Fitting a simple model — aggregate data

contrast	D	Y	Rate
1	928.00	20094.74	4.62
2	1036.00	31839.35	3.25

As long as we only use covariates that take only a few values, we can model the aggregate data directly:

```
> mx <- glm(cbind(c(928, 1036), c(20094.74, 31839.35) / 100) ~ factor(1:2) - 1,  
+           family = poisreg )  
> round(ci.exp(mx), 2)  
  
exp(Est.) 2.5% 97.5%  
factor(1:2)1    4.62 4.33 4.93  
factor(1:2)2    3.25 3.06 3.46
```

Representation of follow-up (time-split)

29 / 40

## SMR

### Bendix Carstensen

Representation of follow-up

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

SMR

## Cohorts where all are exposed

When there is no comparison group we may ask:

Do mortality rates in cohort differ from those of an **external** population, for example:

Rates from:

- ▶ Occupational cohorts
- ▶ Patient cohorts

compared with reference rates obtained from:

- ▶ Population statistics (mortality rates)
- ▶ Hospital registers (disease rates)

SMR (SMR)

30 / 40

## Cohort rates vs. population rates: RSR

- ▶ **Additive:**  $\lambda(a) = \delta(a) + \lambda_p(a)$
- ▶ Note that the survival (since  $a = a_0$ , say) is:

$$\begin{aligned} S(a) &= \exp\left(-\int_{a_0}^a \delta(a) + \lambda_p(a) da\right) \\ &= \exp\left(-\int_{a_0}^a \delta(a) da\right) \times S_p(a) \\ \Rightarrow r(a) &= S(a)/S_p(a) = \exp\left(-\int_{a_0}^a \delta(a) da\right) \end{aligned}$$

- ▶ **Additive** model for **rates**  $\Leftrightarrow$  **Relative survival** model.

## Cohort rates vs. population rates: SMR

- ▶ **Multiplicative:**  $\lambda(a) = \theta \times \lambda_p(a)$
- ▶ Cohort rates proportional to reference rates,  $\lambda_p$ :  
 $\lambda(a) = \theta \times \lambda_p(a)$  —  $\theta$  the same in all age-bands.
- ▶  $D_a$  deaths during  $Y_a$  person-years in age-band  $a$  gives the likelihood:

$$\begin{aligned} D_a \log(\lambda(a)) - \lambda(a) Y_a &= D_a \log(\theta \lambda_p(a)) - \theta \lambda_p(a) Y_a \\ &= D_a \log(\theta) + D_a \log(\lambda_p(a)) - \theta (\lambda_p(a) Y_a) \end{aligned}$$

- ▶ The constant  $D_a \log(\lambda_p(a))$  does not involve  $\theta$ , and so can be dropped.

- ▶  $\lambda_p(a) Y_a = E_a$  is the “expected” number of cases in age  $a$ , so the log-likelihood contribution from age  $a$  is:

$$D_a \log(\theta) - \theta (\lambda_p(a) Y_a) = D_a \log(\theta) - \theta (E_a)$$

- ▶ The log-likelihood is similar to the log-likelihood for a rate, so:

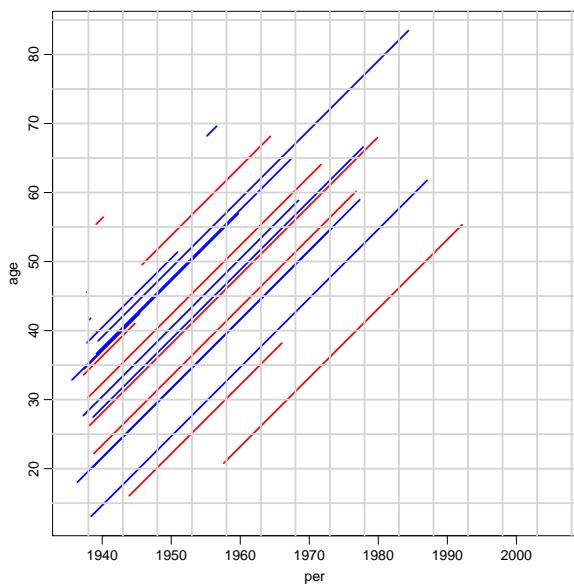
$$\hat{\theta} = \sum_a D_a / \sum_a E_a = \text{Observed/Expected} = \text{SMR}$$

## Modeling the SMR in practice

- ▶ As for the rates, the SMR can be modelled using individual data.
- ▶ Response is  $d_i$ , the event indicator (`lex.Xst`).
- ▶ log-offset is the expected value for each piece of follow-up,  
 $e_i = y_i \times \lambda_p$  (`lex.dur * rate`)
- ▶  $\lambda_p$  is the population rate corresponding to the age, period and sex of the follow-up period  $y_i$ .

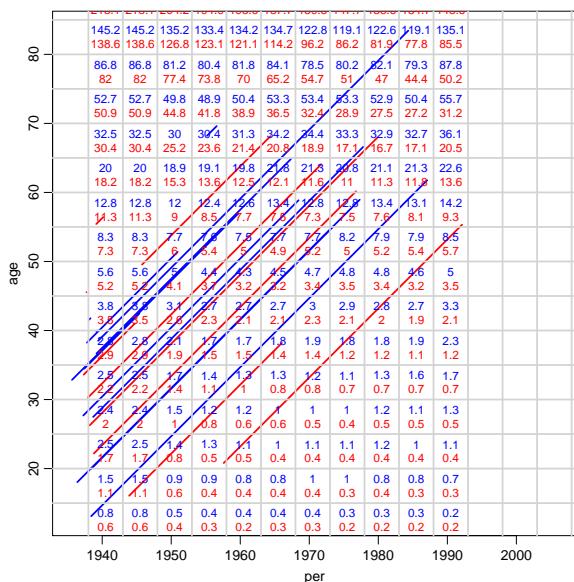
SMR (SMR)

34 / 40



SMR (SMR)

35 / 40



SMR (SMR)

36 / 40

## Split the data to fit with population data

```
> thad <- splitMulti(thL, age=seq(0,90,5), dte=seq(1938,2038,5) )
> summary( thad )
```

Transitions:

From	To	Records:	Events:	Risk time:	Persons:
0	1	21059	1939	22998	51787.96
					2463

Create variables to fit with the population data

```
> thad$agr <- timeBand( thad, "age", "left" )
> thad$per <- timeBand( thad, "dte", "left" )
> round( thad[1:5,c("lex.id","age","agr","dte","per","lex.dur","lex.Xst","sex")], 2
lex.id    age      dte lex.dur lex.Xst agr per sex
  1 22.18 1938.79    2.82      0  20 1938   2
  1 25.00 1941.61    1.39      0  25 1938   2
  1 26.39 1943.00    3.61      0  25 1943   2
  1 30.00 1946.61    1.39      0  30 1943   2
  1 31.39 1948.00    3.61      0  30 1948   2
```

SMR (SMR)

37 / 40

```
> data( gmortDK )
> dim( gmortDK )
[1] 418  21
> gmortDK[1:6,1:6]
  agr per sex risk dt rt
1  0  38    1 996019 14079 14.135
2  5  38    1 802334  726  0.905
3 10  38    1 753017  600  0.797
4 15  38    1 773393 1167  1.509
5 20  38    1 813882 2031  2.495
6 25  38    1 789990 1862  2.357

> gmortDK$per <- gmortDK$per+1900
> #
> thadx <- merge( thad, gmortDK[,c("agr","per","sex","rt")] )
> #
> thadx$E <- thadx$lex.dur * thadx$rt / 1000
```

SMR (SMR)

38 / 40

```
> stat.table(contrast,
+             list( D = sum(lex.Xst),
+                   Y = sum(lex.dur),
+                   E = sum(E),
+                   SMR = ratio(lex.Xst, E)),
+             margin = TRUE,
+             data = thadx)

-----  
contrast      D      Y      E      SMR  
-----  
1          917.00 20045.46 214.66  4.27  
2         1022.00 31742.51 447.21  2.29  
Total     1939.00 51787.96 661.87  2.93  
-----
```

SMR (SMR)

39 / 40

```

-----
contrast      D      Y      E      SMR
-----
1           917.00 20045.46 214.66   4.27
2          1022.00 31742.51 447.21   2.29
-----
> m.SMR <- glm(cbind(lex.Xst, E) ~ factor(contrast) - 1,
+                 family = poisreg,
+                 data = thadx)
> round(ci.exp(m.SMR), 2)
              exp(Est.) 2.5% 97.5%
factor(contrast)1      4.27 4.00  4.56
factor(contrast)2      2.29 2.15  2.43

```

- ▶ Analysis of SMR is like analysis of rates:
- ▶ Replace  $Y$  with  $E$  — that's all! (`glm.Lexis` not usable)
- ▶ ... it's the calculation of  $E$  that is difficult

# Nested case-control and case-cohort studies

Tuesday, 27 August, 2019

Esa Läärä & Martyn Plummer

Statistical Practice in Epidemiology with R

University of Tartu, Estonia

23 to 28 August, 2019

## Points to be covered

- ▶ Outcome-dependent sampling designs a.k.a.  
**case-control** studies vs. **full cohort** design.
- ▶ **Nested case-control** study (NCC): sampling of controls from risk-sets during follow-up of study population.
- ▶ **Matching** in selection of control subjects in NCC.
- ▶ R tools for NCC: function `ccwc()` in `Epi` for sampling controls, and `clogit()` in `survival` for model fitting.
- ▶ **Case-cohort** study (CC): sampling a subcohort from the whole cohort as it is at the start of follow-up.
- ▶ R tools for CC model fitting: function `cch()` in `survival`

## Example: Smoking and cervix cancer

Study population, measurements, follow-up, and sampling design

- ▶ Joint cohort of  $N \approx 500\,000$  women from 3 Nordic biobanks.
- ▶ Follow-up: From variable entry times since 1970s till 2000.
- ▶ For each of 200 cases, 3 controls were sampled; matched for biobank, age ( $\pm 2$  y), and time of entry ( $\pm 2$  mo).
- ▶ Frozen sera of cases and controls analyzed for cotinine etc.

Main result: Adjusted OR = 1.5 (95% CI 1.1 to 2.3) for high ( $>242.6$  ng/ml) vs. low ( $<3.0$  ng/ml) cotinine levels.

Simen Kapeu et al. (2009) *Am J Epidemiol*

## Example: USF1 gene and CVD

Study population, measurements, follow-up, and sampling design

- ▶ Two FINRISK cohorts, total  $N \approx 14000$  M & F, 25-64 y.
- ▶ Baseline health exam, questionnaire & blood specimens at recruitment in the 1990s – Follow-up until the end of 2003.
- ▶ Subcohort of 786 subjects sampled.
- ▶ 528 incident cases of CVD; 72 of them in the subcohort.
- ▶ Frozen blood from cases and subcohort members genotyped.

Main result: Female carriers of a high risk haplotype had a 2-fold hazard of getting CVD [95% CI: 1.2 to 3.5]

Komulainen *et al.* (2006) *PLoS Genetics*

Nested case-control and case-cohort studies

3 / 31

## Full cohort design & its simple analysis

- ▶ **Full cohort design:** Data on exposure variables obtained for all subjects in a large study population.
- ▶ Summary data for crude comparison:

	Exposed	Unexposed	Total
Cases	$D_1$	$D_0$	$D$
Non-cases	$B_1$	$B_0$	$B$
Group size at start	$N_1$	$N_0$	$N$
Follow-up times	$Y_1$	$Y_0$	$Y$

- ▶ Crude estimation of **hazard ratio**  $\rho = \lambda_1/\lambda_0$ :  
**incidence rate ratio** IR, with standard error of  $\log(\text{IR})$ :

$$\hat{\rho} = \text{IR} = \frac{D_1/Y_1}{D_0/Y_0} \quad \text{SE}[\log(\text{IR})] = \sqrt{\frac{1}{D_1} + \frac{1}{D_0}}.$$

- ▶ More refined analyses: Poisson or Cox regression.

Nested case-control and case-cohort studies

4 / 31

## Problems with full cohort design

Obtaining exposure and covariate data

- ▶ Slow and expensive in a big cohort.
- ▶ Easier with questionnaire and register data,
- ▶ Extremely costly and laborious for e.g.
  - measurements from biological specimens, like genotyping, antibody assays, etc.
  - dietary diaries & other manual records

*Can we obtain equally valid estimates of hazard ratios etc. with nearly as good precision by some other strategies?*

Yes – we can!

Nested case-control and case-cohort studies

5 / 31

## Estimation of hazard ratio

The incidence rate ratio can be expressed:

$$\begin{aligned} \text{IR} &= \frac{D_1/D_0}{Y_1/Y_0} = \frac{\text{cases: exposed / unexposed}}{\text{person-times: exposed / unexposed}} \\ &= \frac{\text{exp're odds in cases}}{\text{exp're odds in p-times}} = \mathbf{\text{exposure odds ratio (EOR)}} \end{aligned}$$

= Exposure distribution in cases vs. that in cohort!

Implication for more efficient design:

- ▶ *Numerator*: Collect exposure data on all cases.
- ▶ *Denominator*: Estimate the ratio of person-times  $Y_1/Y_0$  of the exposure groups in the cohort by **sampling** “control” subjects, on whom exposure is measured.

## Case-control designs

General principle: Sampling of subjects from a given study population is *outcome-dependent*.

Data on risk factors are collected separately from

- (I) **Case group**: All (or high % of) the  $D$  subjects in the study population (total  $N$ ) encountering the outcome event during the follow-up.
- (II) **Control group**:
  - ▶ Random **sample** (simple or stratified) of  $C$  subjects ( $C \ll N$ ) from the population.
  - ▶ Eligible controls must be bf risk (alive, under follow-up & free of outcome) at given time(s).

## Study population in a case-control study?

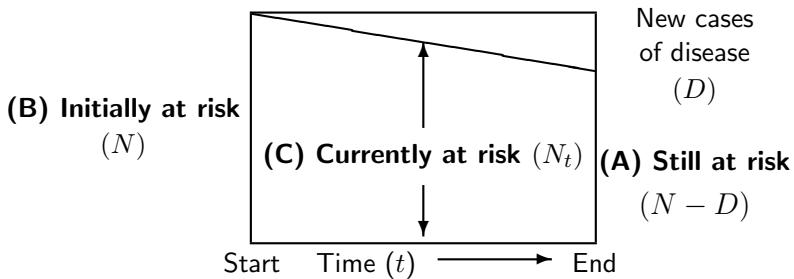
Ideally: The study population comprises subjects who *would be* included as cases, if they got the outcome in the study

- ▶ *Cohort-based studies*: **cohort** or **closed** population of well-identified subjects under intensive follow-up for outcomes (e.g. biobank cohorts).
- ▶ *Register-based studies*: **open** or **dynamic** population in a region covered by a disease register.
- ▶ *Hospital-based studies*: dynamic **catchment** population of cases – may be hard to identify (e.g. hospitals in US).

In general, the role of control subjects is to represent the distribution of person-times by exposure variables in the underlying population from which the cases emerge.

## Sampling of controls – alternative frames

Illustrated in a simple longitudinal setting:  
Follow-up of a cohort over a fixed risk period & no censoring.



Rodrigues, L. & Kirkwood, B.R. (1990). Case-control designs of common diseases ... *Int J Epidemiol* **19**: 205-13.

Nested case-control and case-cohort studies

9 / 31

## Sampling schemes or designs for controls

### (A) Exclusive or traditional, “case-noncase” sampling

- ▶ Controls chosen from those  $N - D$  subjects still at risk (healthy) *at the end* of the risk period (follow-up).

### (B) Inclusive sampling or case-cohort design (CC)

- ▶ The control group – *subcohort* – is a random sample of the cohort ( $N$ ) *at start*.

### (C) Concurrent sampling or density sampling

- ▶ Controls drawn *during the follow-up*
- ▶ **Risk-set or time-matched sampling:**  
A set of controls is sampled from the *risk set at each time  $t$  of diagnosis* of a new case – a.k.a. **nested case-control design** (NCC)

Nested case-control and case-cohort studies

10 / 31

## Nested case-control – two meanings

- ▶ In some epidemiologic books, the term “nested case-control study” (NCC) covers jointly all variants of sampling: **(A)**, **(B)**, and **(C)**, from a cohort.

Rothman *et al.* (2008): *Modern Epidemiology*, 3rd Ed.  
Dos Santos Silva (1999): *Cancer Epidemiology*. Ch 8-9

- ▶ In biostatistical texts NCC typically refers only to the variant of concurrent or density sampling **(C)**, in which *risk-set* or *time-matched* sampling is employed.

Borgan & Samuelsen (2003) in *Norsk Epidemiologi*  
Langholz (2005) in *Encyclopedia of Biostatistics*.

- ▶ We shall follow the biostatisticians!

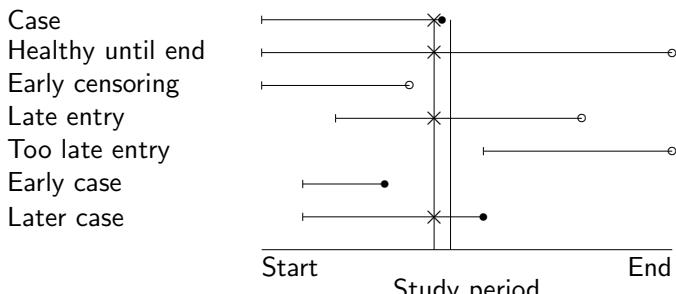
Nested case-control and case-cohort studies

11 / 31

## NCC: Risk-set sampling with staggered entry

Sampling frame to select controls for a given case:

Members ( $\times$ ) of the **risk set** at  $t_k$ , i.e. the population at risk at the time of diagnosis  $t_k$  of case  $k$ .



**Sampled risk set** contains the case and the control subjects randomly sampled from the non-cases in the risk set at  $t_k$ .

## Use of different sampling schemes

### (A) Exclusive sampling, or “textbook” case-control design

- ▶ Almost exclusively(!) used in studies of epidemics.
- ▶ (Studies on birth defects with *prevalent* cases.)

### (B) Inclusive sampling or case-cohort design

- ▶ Good esp. for multiple outcomes, if measurements of risk factors from stored material remain stable.

### (C) Concurrent or density sampling (without or with time-matching)

- ▶ The only logical design in an open population.
- ▶ Most popular in chronic diseases (Knol *et al.* 2008).

Designs (B) and (C) allow valid estimation of hazard ratios  $\rho$  without any “rare disease” assumption.

## Case-control studies: Textbooks vs. real life

- ▶ Many epi texts focus on the traditional design: **exclusive sampling** of controls, ignoring other designs.
- ▶ Claim: “*Odds ratio is the only estimable parameter.*”
- ▶ Yet, over 60% of published case-control studies apply **concurrent sampling** or **density sampling** of controls from an **open** or **dynamic** population.
- ▶ Thus, the parameter most often estimated is the **hazard ratio (HR)** or **rate ratio  $\rho$** .
- ▶ Still, 90% of authors really estimating HR, reported as having estimated an OR (e.g. Simen Kapeu *et al.* 2009)

Knol *et al.* (2008). What do case-control studies estimate?

*Am J Epidemiol* **168**: 1073-81.

## Exposure odds ratio – estimate of what?

- ▶ Crude summary of case-control data

	exposed	unexposed	total
cases	$D_1$	$D_0$	$D$
controls	$C_1$	$C_0$	$C$

- ▶ Depending on study base & sampling strategy, the **exposure odds ratio**

$$\text{EOR} = \frac{D_1/D_0}{C_1/C_0} = \frac{\text{cases: exposed / unexposed}}{\text{controls: exposed / unexposed}}$$

is a consistent estimator of

- (a) hazard ratio, (b) risk ratio, (c) risk odds ratio,
- (d) prevalence ratio, or (e) prevalence odds ratio

- ▶ **NB.** In case-cohort studies with variable follow-up times  $C_1/C_0$  is substituted by  $\widehat{Y}_1/\widehat{Y}_0$ , from estimated p-years.

Nested case-control and case-cohort studies

15 / 31

## Precision and efficiency

With exclusive (**A**) or concurrent (**C**) sampling of controls (unmatched), the estimated variance of  $\log(\text{EOR})$  is

$$\begin{aligned}\widehat{\text{var}}[\log(\text{EOR})] &= \frac{1}{D_1} + \frac{1}{D_0} + \frac{1}{C_1} + \frac{1}{C_0} \\ &= \text{cohort variance} + \text{sampling variance}\end{aligned}$$

- ▶ Depends basically on the numbers of cases, with  $\geq 4$  controls per case.
- ▶ Is not much bigger than  $1/D_1 + 1/D_0$  = variance in a full cohort study with same numbers of cases.
- ⇒ Usually  $< 5$  controls per case is enough.
- ⇒ *These designs are very cost-efficient!*

Nested case-control and case-cohort studies

16 / 31

## Estimation in concurrent or density sampling

- ▶ Assume a simple situation: Prevalence of exposure in the study population stable over time.
- ⇒ The exposure odds  $C_1/C_0$  among controls  
= a consistent estimator of exposure odds  $Y_1/Y_0$  of person-times.
- ▶ Therefore, the crude EOR =  $(D_1/D_0)/(C_1/C_0)$   
= a consistent estimator of hazard ratio  $\rho = \lambda_1/\lambda_0$ .
- ▶ Variance of  $\log(\text{EOR})$  estimated as above.
- ▶ Yet, stability of exposure distribution may be unrealistic, especially in a closed study population or cohort.
- ▶ Solution: **Time-matched** sampling of controls from **risk sets**, i.e. NCC, & matched EOR to estimate HR.

Prentice & Breslow (1978), Greenland & Thomas (1982).

Nested case-control and case-cohort studies

17 / 31

## Matching in case-control studies

- = **Stratified sampling** of controls, e.g. from the same region, sex, and age group as a given case
- ▶ **Frequency matching or group matching:**  
For cases in a specific stratum (e.g. same sex and 5-year age-group), a set of controls from a similar subgroup.
- ▶ **Individual matching** (1:1 or 1:m matching):  
For each case, choose 1 or more (rarely > 5) closely similar controls (e.g. same sex, age within  $\pm 1$  year).
- ▶ **NCC:** Sampling from risk-sets implies time-matching at least. Additional matching for other factors possible.
- ▶ **CC:** Subcohort selection involves no matching with cases.

## Virtues of matching

- ▶ Increases *efficiency*, if the matching factors are both
  - (i) strong *risk factors* of the disease, and
  - (ii) *correlated* with the main exposure.
- Major reason for matching.
- ▶ *Confounding* due to poorly quantified factors (sibship, neighbourhood, etc.) may be removed by close matching – only if properly analyzed.
- ▶ Biobank studies: Matching for storage time, freeze-thaw cycle & analytic batch improves **comparability of measurements** from frozen specimens
  - Match on the time of baseline measurements within the case's risk set.

## Warnings for overmatching

Matching a case with a control subject is a different issue than matching an unexposed subject to an exposed one in a cohort study – much trickier!

- ▶ Matching on an *intermediate* variable between exposure and outcome.
  - ⇒ *Bias!*
- ▶ Matching on a *surrogate* or *correlate* of exposure, which is not a true risk factor.
  - ⇒ *Loss of efficiency.*
- **Counter-matching:** Choose a control which is not similar to the case w.r.t a correlate of exposure.
  - ⇒ Increases efficiency!
  - Requires appropriate weighting in the analysis.

## Sampling matched controls for NCC using R

- ▶ Suppose key follow-up items are recorded for all subjects in a cohort, in which a NCC study is planned.
- ▶ Function `ccwc()` in package `Epi` can be used for risk-set sampling of controls.
  - Arguments:

```
entry : Time of entry to follow-up  
exit : Time of exit from follow-up  
fail : Status on exit (1 for case, 0 for censored)  
origin : Origin of analysis time scale (e.g. time of birth)  
controls : Number of controls to be selected for each case  
match : List of matching factors  
data : Cohort data frame containing input variables
```

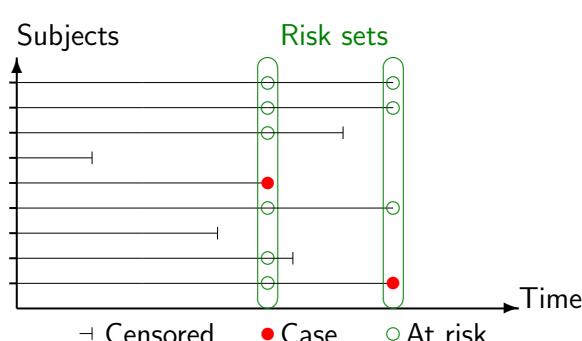
- ▶ Creates a data frame for a NCC study, containing the desired number of matched controls for each case.

## Analysis of matched studies

- ▶ Close matching induces a new parameter for each matched case-control set or stratum.
  - ⇒ **unconditional logistic regression** breaks down.
- ▶ Matching on well-defined variables (like age, sex)
  - include these factors as covariates.
- ▶ Matching on “soft” variables (like sibship) can be dealt with **conditional logistic regression**.
- ▶ Same method in matched designs **(A)**, exclusive, and **(C)**, concurrent, but interpretation of  $\beta_j$ s differs:
  - (A)**  $\beta_j = \log$  of risk odds ratio (ROR),
  - (C)**  $\beta_j = \log$  of hazard ratio (HR).

## Full cohort design: Follow-up & risk sets

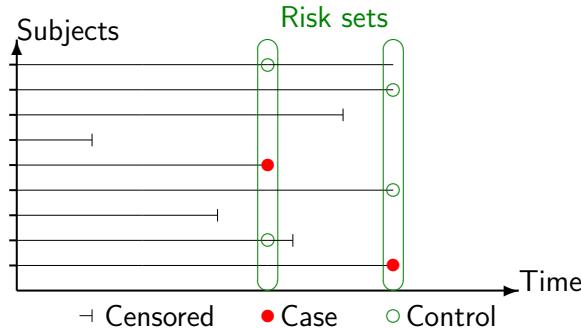
Each member of the cohort provides exposure data for all cases, as long as this member is at risk, i.e. (i) alive, (ii) not censored & (iii) free from outcome.



Times of new cases define the **risk-sets**.

## Nested case-control (NCC) design

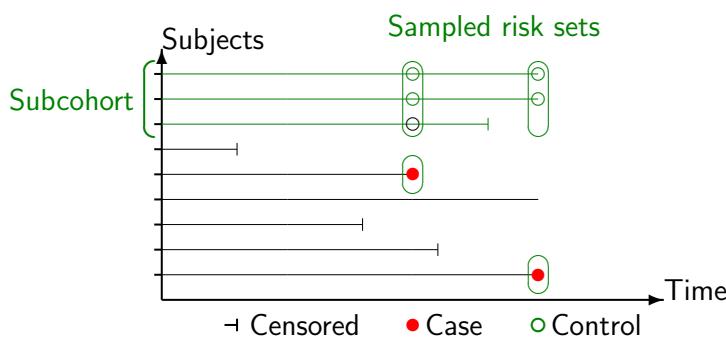
Whenever a new case occurs, a set of controls (here 2/case) are sampled from its risk set.



**NB.** A control once selected for some case can be selected as a control for another case, and can later on become a case, too.

## Case-cohort (CC) design

**Subcohort:** Sample of the whole cohort randomly selected at the outset.  
– Serves as a reference group for all cases.



**NB.** A subcohort member can become a case, too.

## Modelling in NCC and other matched studies

Cox proportional hazards model:

$$\lambda_i(t, x_i; \beta) = \lambda_0(t) \exp(x_{i1}\beta_1 + \cdots + x_{ip}\beta_p),$$

Estimation: partial likelihood  $L^P = \prod_k L_k^P$ :

$$L_k^P = \exp(\eta_{i_k}) / \sum_{i \in \tilde{R}(t_k)} \exp(\eta_i),$$

where  $\tilde{R}(t_k)$  = **sampled risk set** at observed event time  $t_k$ , containing the case + sampled controls ( $t_1 < \cdots < t_D$ )

⇒ Fit stratified Cox model, with  $\tilde{R}(t_k)$ 's as the strata.

⟲ **Conditional logistic regression**

– function `clogit()` in `survival`, wrapper of `coxph()`.

## Modelling case-cohort data

Cox's PH model  $\lambda_i(t) = \lambda_0(t) \exp(\eta_i)$  again, but ...

- ▶ Analysis of survival data relies on the theoretical principle that *you can't know the future*.
- ▶ Case-cohort sampling breaks this principle:  
cases are sampled based on what *is known* to be happening to them during follow-up.
- ▶ The union of cases and subcohort is a mixture
  1. random sample of the population, and
  2. "high risk" subjects who are *certain* to become cases.
- ⇒ Ordinary Cox partial likelihood is wrong.
- ▶ Overrepresentation of cases must be corrected for, by  
(I) **weighting**, or (II) **late entry method**.

## Correction method I – weighting

The method of **weighted partial likelihood** borrows some basic ideas from survey sampling theory.

- ▶ Sampled risk sets  
 $\tilde{R}(t_k) = \{\text{cases}\} \cup \{\text{subcohort members}\}$  at risk at  $t_k$ .
- ▶ Weights:
  - $w = 1$  for all cases (within and outside the subcohort),
  - $w = N_{\text{non-cases}}/n_{\text{non-cases}} = \text{inverse of sampling-fraction } f$  for selecting a non-case to the subcohort.
- ▶ Function `coxph()` with option `weights = w` would provide consistent estimation of  $\beta$  parameters.
- ▶ However, the SEs must be corrected!
- ▶ R solution: Function `cch()` – a wrapper of `coxph()` – in package `survival`, with `method = "LinYing"`.

## Comparison of NCC and CC designs

- ▶ Statistical efficiency

Broadly similar in NCC and CC with similar numbers of cases and controls.
- ▶ Statistical modelling and valid inference

Straightforward for both designs with appropriate software, now widely available for CC, too
- ▶ Analysis of outcome rates on several time scales?

**NCC:** Only the time scale used in risk set definition can be the time variable  $t$  in the baseline hazard of PH model.

**CC:** Different choices for the basic time in PH model possible, because subcohort members are not time-matched to cases.

## Comparison of designs (cont'd)

### ► Missing data

NCC: With close 1:1 matching, a case-control pair is lost, if either of the two has data missing on key exposure(s).

CC: Missingness of few data items is less serious.

### ► Quality and comparability of biological measurements

NCC: Allows each case and its controls to be matched also for analytic batch, storage time, freeze-thaw cycle, → better comparability.

CC: Measurements for subcohort performed at different times than for cases → differential quality & misclassification.

### ► Possibility for studying many diseases with same controls

NCC: Complicated, but possible if matching is not too refined.

CC: Easy, as no subcohort member is "tied" with any case.

## Conclusion

- "Case-controlling" is very cost-effective.
- Case-cohort design is useful especially when several outcomes are of interest, given that the measurements on stored materials remain stable during the study.
- Nested case-control design is better suited e.g. for studies involving biomarkers that can be influenced by analytic batch, long-term storage, and freeze-thaw cycles.
- Matching helps in improving efficiency and in reducing bias – but only if properly done.
- Handy R tools are available for all designs.

## Some topics on causal inference

Krista Fischer

Institute of Mathematics and Statistics, University of Tartu  
Estonian Genome Center, Institute of Genomics, University of Tartu

Statistical Practice in Epidemiology, Tartu 2019

### How to define a causal effect?

### Causal graphs, confounding and adjustment

### Causal models for observational data

Instrumental variables estimation and Mendelian randomization

### Summary and references

### References

## Statistical associations vs causal effects in epidemiology

Does the exposure (smoking level, obesity, etc) have a causal effect on the outcome (cancer diagnosis, mortality, etc)?

is not the same question as

Is the exposure associated with the outcome?

Conventional statistical analysis will answer the second one, but not necessarily the first.

# Statistical associations vs causal effects in epidemiology

Does the exposure (smoking level, obesity, etc) have a causal effect on the outcome (cancer diagnosis, mortality, etc)?

is not the same question as

Is the exposure associated with the outcome?

Conventional statistical analysis will answer the second one, but not necessarily the first.

# Statistical associations vs causal effects in epidemiology

Does the exposure (smoking level, obesity, etc) have a causal effect on the outcome (cancer diagnosis, mortality, etc)?

is not the same question as

Is the exposure **associated** with the outcome?

Conventional statistical analysis will answer the second one, but not necessarily the first.

## What is a causal effect?

*There is more than just one way to define it.*

A causal effect may be defined:

- ▶ At the individual level:  
*Would my cancer risk be different if I were a (non-)smoker?*
  - ▶ At the population level:  
*Would the population cancer incidence be different if the prevalence of smoking were different?*
  - ▶ At the exposed subpopulation level:  
Would the cancer incidence in smokers be different if they were nonsmokers?

None of these questions is “mathematical” enough to provide a mathematically correct definition of causal effect.

## What is a causal effect?

*There is more than just one way to define it.*

A causal effect may be defined:

- ▶ At the individual level:  
*Would my cancer risk be different if I were a (non-)smoker?*
  - ▶ At the population level:  
*Would the population cancer incidence be different if the prevalence of smoking were different?*
  - ▶ At the exposed subpopulation level:  
Would the cancer incidence in smokers be different if they were nonsmokers?

None of these questions is “mathematical” enough to provide a mathematically correct definition of causal effect

## What is a causal effect?

*There is more than just one way to define it.*

A causal effect may be defined:

- ▶ At the individual level:  
*Would my cancer risk be different if I were a (non-)smoker?*
  - ▶ At the population level:  
*Would the population cancer incidence be different if the prevalence of smoking were different?*
  - ▶ At the **exposed subpopulation level**:  
*Would the cancer incidence in smokers be different if they were nonsmokers?*

None of these questions is “mathematical” enough to provide a mathematically correct definition of causal effect.

## What is a causal effect?

*There is more than just one way to define it.*

A causal effect may be defined:

- ▶ At the individual level:  
*Would my cancer risk be different if I were a (non-)smoker?*
  - ▶ At the population level:  
*Would the population cancer incidence be different if the prevalence of smoking were different?*
  - ▶ At the exposed subpopulation level:  
Would the cancer incidence in smokers be different if they were nonsmokers?

None of these questions is “mathematical” enough to provide a mathematically correct definition of causal effect

## Causal effects and counterfactuals

- ▶ Defining the causal effect of an observed exposure always involves some **counterfactual** (what-if) thinking.
  - ▶ The individual causal effect can be defined as the difference

$$Y(X=1) - Y(X=0)$$

- . where  $Y(1) = Y(X = 1)$  and  $Y(0) = Y(X = 0)$  are defined as individual's potential (counterfactual) outcomes if this individual's exposure level  $X$  were set to 1 or 0, respectively.

- ▶ Sometimes people (e.g J. Pearl) use the “do” notation to distinguish counterfactual variables from the observed ones:  $Y(\text{do}(X = 1))$  and  $Y(\text{do}(X = 0))$ .

## Causal effects and counterfactuals

- ▶ Defining the causal effect of an observed exposure always involves some counterfactual (what-if) thinking.
  - ▶ The individual causal effect can be defined as the difference

$$Y(X=1) - Y(X=0)$$

. where  $Y(1) = Y(X = 1)$  and  $Y(0) = Y(X = 0)$  are defined as individual's potential (counterfactual) outcomes if this individual's exposure level  $X$  were set to 1 or 0, respectively,

- Sometimes people (e.g J. Pearl) use the “do” notation to distinguish counterfactual variables from the observed ones:  $Y(\text{do}(X = 1))$  and  $Y(\text{do}(X = 0))$ .

## Causal effects and counterfactuals

- ▶ Defining the causal effect of an observed exposure always involves some counterfactual (what-if) thinking.
  - ▶ The individual causal effect can be defined as the difference

$$Y(X = 1) - Y(X = 0)$$

. where  $Y(1) = Y(X = 1)$  and  $Y(0) = Y(X = 0)$  are defined as individual's potential (counterfactual) outcomes if this individual's exposure level  $X$  were set to 1 or 0, respectively.

- Sometimes people (e.g J. Pearl) use the “*do*” notation to distinguish counterfactual variables from the observed ones:  $Y(\text{do}(X = 1))$  and  $Y(\text{do}(X = 0))$ .

## The “naïve” association analysis

- ▶ With a binary exposure  $X$ , compare average outcomes in exposed and unexposed populations:

$$E(Y|X=1) - E(Y|X=0)$$

*Is cancer incidence different in smokers and nonsmokers?*

- But mostly:

$$E(Y|X=1) \neq E(Y(1))$$

*Cancer risk in smokers is not the same as the potential cancer risk in the population if everyone were smoking*

- Similarly:

$$E(Y|X=0) \neq E(Y(0))$$

- ▶ In most cases there is always some **unobserved confounding** present and therefore the naïve analysis does not provide causal effect estimates.

## Counterfactual outcomes in different settings

- ▶ **Randomized trials:** probably the easiest setting to imagine  $Y(X)$  for different  $X$
  - ▶ “Actionable” exposures: smoking level, vegetable consumption, . . . – potential interventions may alter exposure levels in future.
  - ▶ Non-actionable exposures: e.g. genotypes. It is difficult to ask *“What if I had different genes?”*. Still useful concept to formalize genetic effects (heritability, attributable risk).
  - ▶ Combinations: With  $X$ – a behavioral intervention level,  $Z$ –smoking level and  $Y$ –a disease outcome, one could formalize the effect of intervention on outcome by using  $Y(X, Z(X))$

## Counterfactual outcomes in different settings

- ▶ Randomized trials: probably the easiest setting to imagine  $Y(X)$  for different  $X$
  - ▶ “Actionable” exposures: smoking level, vegetable consumption, . . . – potential interventions may alter exposure levels in future.
  - ▶ Non-actionable exposures: e.g. genotypes. It is difficult to ask “*What if I had different genes?*”. Still useful concept to formalize genetic effects (heritability, attributable risk).
  - ▶ Combinations: With  $X$  – a behavioral intervention level,  $Z$  – smoking level and  $Y$  – a disease outcome, one could formalize the effect of intervention on outcome by using  $Y(X, Z(X))$

## Counterfactual outcomes in different settings

- ▶ Randomized trials: probably the easiest setting to imagine  $Y(X)$  for different  $X$
  - ▶ “Actionable” exposures: smoking level, vegetable consumption, . . . – potential interventions may alter exposure levels in future.
  - ▶ **Non-actionable exposures:** e.g. genotypes. It is difficult to ask *“What if I had different genes?”*. Still useful concept to formalize genetic effects (heritability, attributable risk).
  - ▶ Combinations: With  $X$  – a behavioral intervention level,  $Z$  – smoking level and  $Y$  – a disease outcome, one could formalize the effect of intervention on outcome by using  $Y(X, Z(X))$

## Counterfactual outcomes in different settings

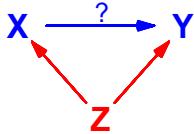
- ▶ Randomized trials: probably the easiest setting to imagine  $Y(X)$  for different  $X$
  - ▶ “Actionable” exposures: smoking level, vegetable consumption, . . . – potential interventions may alter exposure levels in future.
  - ▶ Non-actionable exposures: e.g. genotypes. It is difficult to ask *“What if I had different genes?”*. Still useful concept to formalize genetic effects (heritability, attributable risk).
  - ▶ **Combinations:** With  $X$ – a behavioral intervention level,  $Z$ –smoking level and  $Y$ –a disease outcome, one could formalize the effect of intervention on outcome by using  $Y(X, Z(X))$

## Classical/generalized regression estimates vs causal effects?

- ▶ In the presence of confounding, regression analysis provides a biased estimate for the true causal effect
  - ▶ To reduce such bias, one needs to collect data on most important confounders and adjust for them
  - ▶ However, too much adjustment may actually introduce more biases
  - ▶ Causal graphs (Directed Acyclic Graphs, DAGs) may be extremely helpful in identifying the optimal set of adjustment variables

## Adjustment for confounders I

“Classical” confounding: situation where third factors Z influence both, X and Y



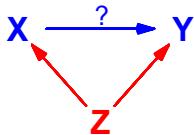
For instance, one can assume:  $X = Z + U$  and  $Y = Z + V$ , where  $U$  and  $V$  are independent of  $Z$ .

$X$  and  $Y$  are independent, conditional on  $Z$ , but marginally dependent.

One should adjust the analysis for  $Z$ , by fitting a regression model for  $Y$  with covariates  $X$  and  $Z$ . There is a causal effect between  $X$  and  $Y$ , if the effect of  $X$  is present in such model.

## Adjustment for confounders I

“Classical” confounding: situation where third factors Z influence both, X and Y



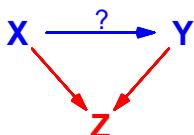
For instance, one can assume:  $X = Z + U$  and  $Y = Z + V$ , where  $U$  and  $V$  are independent of  $Z$ .

$X$  and  $Y$  are independent, conditional on  $Z$ , but marginally dependent.

One should adjust the analysis for  $Z$ , by fitting a regression model for  $Y$  with covariates  $X$  and  $Z$ . There is a causal effect between  $X$  and  $Y$ , if the effect of  $X$  is present in such model.

## Adjustment may sometimes make things worse

Example: the effect of X and Y on Z:



A simple model may hold:  $Z = X + Y + U$ , where  $U$  is independent of  $X$  and  $Y$ .

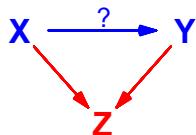
Hence  $Y = Z - X - U$ .

We see the association between  $X$  and  $Y$  only when the “effect” of  $Z$  has been taken into account. But this is not the causal effect of  $X$  on  $Y$ .

One should NOT adjust the analysis for  $Z$ !

Adjustment may sometimes make things worse

Example: the effect of X and Y on Z:



A simple model may hold:  $Z = X + Y + U$ ,

where  $U$  is independent of  $X$  and  $Y$ .

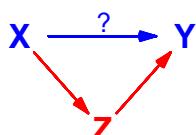
Hence  $Y = Z - X - U$ .

We see the association between  $X$  and  $Y$  only when the “effect” of  $Z$  has been taken into account. But this is not the causal effect of  $X$  on  $Y$ .

**One should NOT adjust the analysis for  $Z$ !**

### More possibilities: mediation

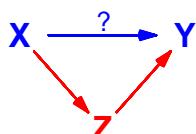
Example: the effect of X on Y is (partly) **mediated** by Z:



$$Y \equiv X + Z + U.$$

## More possibilities: mediation

Example: the effect of X on Y is (partly) **mediated** by Z:

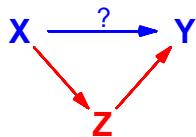


$$Y \equiv X + Z + U.$$

If you are interested in the **total effect** of  $X$  on  $Y$  – don't adjust for  $Z$ !

## More possibilities: mediation

Example: the effect of X on Y is (partly) mediated by Z:

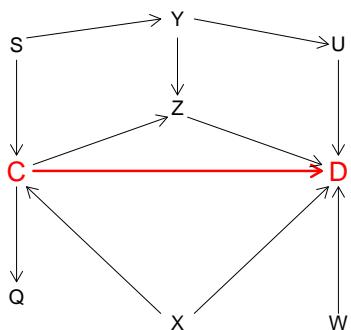


$$Y = X + Z + U,$$

If you are interested in the **total effect** of  $X$  on  $Y$  – don't adjust for  $Z$ !

If you are interested in the **direct effect** of  $X$  on  $Y$  – adjust for  $Z$ .  
(Only if the  $Z-Y$  association is unconfounded)

Actually there might be a complicated system of causal effects:



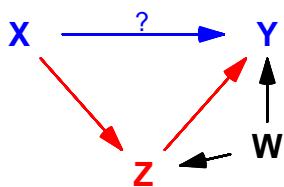
C-smoking; D-cancer

Q, S, U, W, X, Y, Z - other factors that influence cancer risks and/or smoking (genes, social background, nutrition, environment, personality, ...)

*To check for confounding,*

1. Sketch a causal graph
  2. Remove all arrows corresponding to the causal effect of interest (thus, create a graph where the causal null-hypothesis would hold).
  3. Remove all nodes (and corresponding edges) except those contained in the exposure ( $C$ ) and outcome ( $D$ ) variables and their (direct or indirect) ancestors.
  4. Connect by an undirected edge every pair of nodes that both share a common child and are not already connected by a directed edge.
    - ▶ If now  $C$  and  $D$  are still associated, we say that the  $C - D$  association is confounded
    - ▶ Identify the set of nodes that need to be deleted to separate  $C$  and  $D$  – inferences conditional on these variables give unconfounded estimates of the causal effects.

## Example: mediation with confounding



Follow the algorithm to show that one should adjust the analysis for  $W$ . If  $W$  is an unobserved confounder, no valid causal inference is possible in general. However, the total effect of  $X$  on  $Y$  is estimable.

## “Mendelian randomization” – genes as Instrumental Variables

- Most of the exposures of interest in chronic disease epidemiology cannot be randomized.
  - Sometimes, however, nature will randomize for us: there is a SNP (Single nucleotide polymorphism, a DNA marker) that affects the exposure of interest, but not directly the outcome.
  - Example: a SNP that is associated with the enzyme involved in alcohol metabolism, genetic lactose intolerance, etc.

However, the crucial assumption that the SNP cannot affect outcome in any other way than throughout the exposure, cannot be tested statistically!

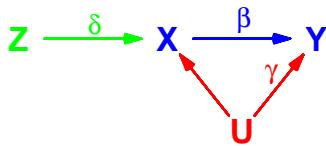
“Mendelian randomization” – genes as Instrumental Variables

- ▶ Most of the exposures of interest in chronic disease epidemiology cannot be randomized.
  - ▶ Sometimes, however, nature will randomize for us: there is a SNP (Single nucleotide polymorphism, a DNA marker) that affects the exposure of interest, but not directly the outcome.
  - ▶ Example: a SNP that is associated with the enzyme involved in alcohol metabolism, genetic lactose intolerance, etc.

However, the crucial assumption that the SNP cannot affect outcome in any other way than throughout the exposure, cannot be tested statistically!

## General instrumental variables estimation

A causal graph with exposure  $X$ , outcome  $Y$ , confounder  $U$  and an *instrument*  $Z$ :



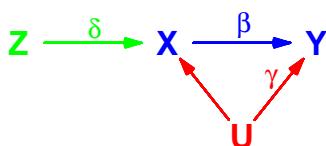
Simple regression will yield a biased estimate of the causal effect of  $X$  on  $Y$ , as the graph implies:

$$Y = \alpha_y + \beta X + \gamma U + \epsilon, E(\epsilon|X, U) = 0$$

so  $E(Y|X) = \alpha_y + \beta X + \gamma E(U|X)$ .

Thus the coefficient of  $X$  will also depend on  $\gamma$  and the association between  $X$  and  $U$ .

## General instrumental variables estimation



$$Y = \alpha_y + \beta X + \gamma U + \epsilon, E(\epsilon|X, U) = 0$$

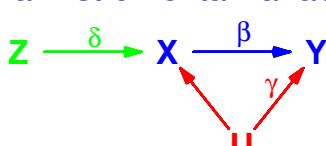
**How can  $Z$  help?**

If  $E(X|Z) = \alpha_x + \delta Z$ , we get

$$E(Y|Z) = \alpha_y + \beta E(X|Z) + \gamma E(U|Z) = \alpha_y + \beta(\alpha_x + \delta Z) = \alpha_y^* + \beta \delta Z.$$

As  $\delta$  and  $\beta\delta$  are estimable, also  $\beta$  becomes estimable.

## General instrumental variables estimation



$$Y = \alpha_y + \beta X + \gamma U + \epsilon, E(\epsilon|X, U) = 0$$

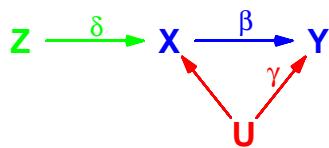
**How can  $Z$  help?**

If  $E(X|Z) = \alpha_x + \delta Z$ , we get

$$E(Y|Z) = \alpha_y + \beta E(X|Z) + \gamma E(U|Z) = \alpha_y + \beta(\alpha_x + \delta Z) = \alpha_y^* + \beta \delta Z.$$

As  $\delta$  and  $\beta\delta$  are estimable, also  $\beta$  becomes estimable.

## General instrumental variables estimation



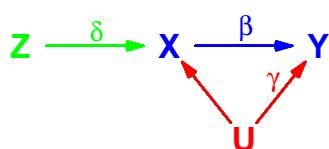
$$Y = \alpha_y + \beta X + \gamma U + \epsilon, \quad E(\epsilon|X, U) = 0$$

## How can Z help?

If  $E(X|Z) = \alpha_x + \delta Z$ , we get

$$E(Y|Z) = \alpha_y + \beta E(X|Z) + \gamma E(U|Z) = \alpha_y + \beta(\alpha_x + \delta Z) = \alpha_y^* + \beta\delta Z.$$

## General instrumental variables estimation



$$Y = \alpha_Y + \beta X + \gamma U + \epsilon, \quad E(\epsilon|X, U) = 0$$

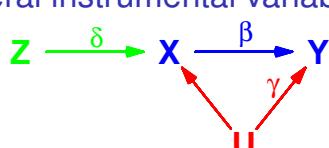
## How can Z help?

If  $E(X|Z) = \alpha_x + \delta Z$ , we get

$$E(Y|Z) = \alpha_Y + \beta E(X|Z) + \gamma E(U|Z) = \alpha_Y + \beta(\alpha_X + \delta Z) = \alpha_Y^* + \beta\delta Z.$$

As  $\delta$  and  $\beta\delta$  are estimable, also  $\beta$  becomes estimable.

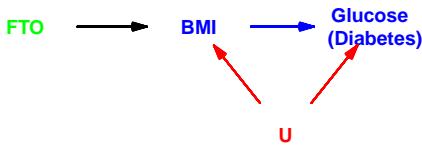
## General instrumental variables estimation



1. Regress  $X$  on  $Z$ , obtain an estimate  $\hat{\delta}$
  2. Regress  $Y$  on  $Z$ , obtain an estimate  $\hat{\delta}\beta$
  3. Obtain  $\hat{\beta} = \frac{\hat{\delta}\beta}{\hat{\delta}}$
  4. Valid, if  $Z$  is not associated with  $U$  and does not have any effect on  $Y$  (other than mediated by  $X$ )
  5. Standard error estimation is more tricky – use for instance library(sem).function tsLS().

## Mendelian randomization example

FTO genotype, BMI and Blood Glucose level (related to Type 2 Diabetes risk; Estonian Biobank, n=3635, aged 45+)



- ▶ Average difference in Blood Glucose level (Glc, mmol/L) per BMI unit is estimated as 0.085 (SE=0.005)
  - ▶ Average BMI difference per FTO risk allele is estimated as 0.50 (SE=0.09)
  - ▶ Average difference in Glc level per FTO risk allele is estimated as 0.13 (SE=0.04)
  - ▶ **Instrumental variable estimate of the mean Glc difference per BMI unit is 0.209 (se=0.078)**

## IV estimation in R (using library (sem)):

## IV estimation: can untestable assumptions be tested?

```

> summary(lm(Glc~bmi+fto, data=fen))
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 1.985     0.106   18.75 <2e-16 ***
bmi         0.088     0.004   23.36 <2e-16 ***
fto         0.049     0.030    1.66    0.097 .
For Type 2 Diabetes:
> summary(glm(t2d~bmi+fto, data=fen, family=binomial))
Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -7.515     0.187  -40.18 <2e-16 ***
bmi         0.185     0.006   31.66 <2e-16 ***
fto         0.095     0.047    2.01    0.044 *

```

Does FTO have a direct effect on Glc or T2D?

A significant FTO effect would not be a proof here (nor does non-significance prove the opposite)! (WHY?)

## Can we test pleiotropy?

A naïve approach would be to fit a linear regression model for  $Y$ , with both  $X$  and  $G$  as covariates.

But in this case we estimate:

$$E(Y|X, G) = \text{const} + \beta_{pl}G + \beta X + \gamma E(U|X, G).$$

It is possible to show that  $U$  is not independent of either  $X$  nor  $G$  – therefore, the coefficient of  $G$  in the resulting model would be nonzero even if  $\beta_{pl} = 0$ .

Therefore there is no formal test for pleiotropy possible in the case of one genetic instrument – only biological arguments could help to decide, whether assumptions are likely to be fulfilled

In the case of *multiple genetic instruments* and *meta-analysis*, sometimes the approach of *Egger regression* can be used (Bowden et al, 2015). But even that is not an assumption-free method!

## Summary

- ▶ There is no unique definition of “the causal effect”
  - ▶ The validity of any causal effect estimates depends on the validity of the underlying assumptions.
  - ▶ Adjustment for other available variables may remove (some) confounding, but it may also create more confounding. **Do not adjust for variables that may themselves be affected by the outcome.**
  - ▶ Instrumental variables approaches can be helpful, but beware of assumptions!

## Some references

- ▶ A webpage by Miguel Hernan and Jamie Robins:  
<http://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>
  - ▶ *An overview of Mendelian randomization:* Sheehan, N., Didelez, V., et al., Mendelian Randomization and Causal Inference in Observational Epidemiology, PLoS Med. 2008 August; 5(8).
  - ▶ *A way to correct for pleiotropy bias:* Bowden J, Davey Smith G, Burgess S, Mendelian randomization with invalid instruments. Int J Epidemiol. 2015 Apr;44(2):512-25.
  - ▶ ... and how to interpret the findings (warning against overuse): Burgess, S., Thompson, S.G., Interpreting findings from Mendelian randomization using the MR-Egger method, Eur J Epidemiol (2017).
  - ▶ A lot of ongoing research on Mendelian randomization ... (see recent papers by Jack Bowden, Stephen Burgess and others on methodological advances, follow e.g. [@MR\\_lit](#) on Twitter for case studies) and causal inference in general

# Causal Inference 2: Model-based estimation of causal contrasts

Tuesday, 6 June, 2023

Esa Läärä

Statistical Practice in Epidemiology with R

2 to 7 June, 2023

University of Tartu, Estonia

## Outline

- ▶ Causal questions
- ▶ Associational contrasts of risks and hazards
- ▶ Causal estimands: contrasts of counterfactual quantities
- ▶ Outcome regression models
- ▶ Standardization or g-formula
- ▶ Exposure modelling, propensity scores and weighting
- ▶ Double robust estimators and machine learning algorithms

## Some literature

- ▶ Austin & Stuart (2015) *Stat Med* 34(28):3661-3679.
- ▶ Funk et al. (2011) *Am J Epidemiol* 173(7):761767
- ▶ Hernan & Robins (2020). *Causal Inference: What if?*. CRC Press.
- ▶ Luque Fernandez et al. (2018) *Stat Med* 2018;37(16):2530-2546
- ▶ Sjölander (2016) *Eur J Epidemiol* 31:563-574
- ▶ Smith et al. (2022) *Stat Med* 2022;41(2):407-432.
- ▶ Zhou et al. (2022) PSweight vignette.

## Causal question in PECOT format & Example

P **Population:** 2900 women with breast cancer (Rotterdam study)

E **Exposure:** Hormonal treatment (HT)

C **Comparator:** Placebo, no HT

O **Outcome:** Recurrence or death

T **Time frame:** 10 y from surgery to outcome

Causal questions of interest – comparisons of counterfactuals:

- What is the 10-year risk  $\pi^1$  of the outcome, if everybody in P were exposed to HT, as compared with  $\pi^0$ , the risk if nobody were exposed?
- What is the 10-year risk  $\pi_1^1$  of the outcome, among those in P, who are factually exposed to HT, as compared with the risk  $\pi_1^0$ , if they were not exposed?

## Risks by factual exposure and their associational contrasts

- ▶ Let  $T = \text{time to outcome event}$ ,  $Y(t) = \mathbf{1}_{\{T \leq t\}}$  = binary indicator (1/0) for the outcome to occur by time  $t$ , and  $X$  be a risk factor,
- ▶ Let  $\pi_x(t) = E[Y(t)|X = x]$  = risk of the outcome to occur by time  $t$  in the subset of the target population factually exposed to level  $X = x$ :

$$\pi_x(t) = P\{Y(t) = 1 | X = x\} = P\{T \leq t | X = x\}.$$

- ▶ For simplicity, let  $X$  be dichotomous & suppress dependence on time, too. Assume also absence of competing events.
- ▶ Common **associational contrasts** of risks between exposure groups:

- **Risk difference**  $\tau = \pi_1 - \pi_0 = E(Y|X = 1) - E(Y|X = 0)$ ,
- **Risk ratio**  $\phi = \pi_1/\pi_0$ ,
- **Risk odds ratio**  $\psi = \frac{\omega_1}{\omega_0} = \frac{\pi_1/(1-\pi_1)}{\pi_0/(1-\pi_0)}$ .

## Hazards by factual exposure and their associational contrasts

- ▶ Let the hazard of outcome  $Y$  at  $t$  among those exposed to level  $X = x$  be
$$\lambda_x(t) = \lambda(t | X = x) = \lim_{h \rightarrow 0} P\{Y(t+h) = 1 | Y(t) = 0, X = x\}/h.$$
- ▶ Common associational contrasts:
  - **Hazard difference**  $\delta(t) = \lambda_1(t) - \lambda_0(t)$ ,
  - **Hazard ratio**  $\rho(t) = \lambda_1(t)/\lambda_0(t)$ .  
This is often assumed constant  $\rho$  – as in **Cox regression**.

- ▶ Other quantities of interest: contrasts of **restricted mean survival time**

$$\mathbf{RMST}(\tau) = \int_0^\tau [1 - \pi_x(t)] dt, \quad x = 0, 1.$$

- ▶ So far all these associational contrasts are **marginal** or **unconditional**, meaning that they are not conditioned on or stratified by any covariate  $Z$ .

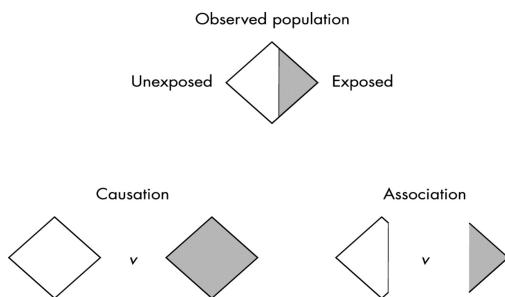
## Conditional associational contrasts

- ▶ Let  $Z$  be a covariate (can be multivariable) and
 
$$\pi_{xz} = P\{Y = 1 \mid X = x, Z = z\} = E(Y|X = x, Z = z)$$
 be the risk of outcome by  $t$  in a population group where  
 $X = x$  and  $Z = z$ ,  $x = 0, 1$ .
- ▶ Let  $\theta$  be any type of marginal associational contrast, and  $\theta_z$  be the analogous **conditional associational contrast** between exposed and unexposed among those with  $Z = z$ .
- ▶ For instance,  $\tau_z = \pi_{1z} - \pi_{0z}$  is the risk difference conditional on  $Z = z$ , i.e.  $z$ -specific risk difference.
- ▶ Similarly  $\phi_z = \pi_{1z}/\pi_{0z}$  and  $\psi_z = \pi_{1z}(1 - \pi_{1z})/[\pi_{0z}(1 - \pi_{0z})]$  are the  $z$ -specific risk ratio and odds ratio, respectively.
- ▶ Conditional contrasts of hazards:  $\delta_z = \lambda_{1z} - \lambda_{0z}$  and  $\rho_z = \lambda_{1z}/\lambda_{0z}$ , where  $\lambda_{xz} = \lambda_{xz}(t) = \lambda(t|X = x, Z = z)$ .

Causal Inference 2: Model-based estimation of causal contrasts

6 / 29

## Associational and causal contrasts



- ▶ **Associational:** Contrast of risks between the **subsets** of the population determined by the subjects' **factual** exposure value.
- ▶ **Causal:** Contrast of risks in the **entire population** under the alternative **potential** or **counterfactual** exposure values;  
see Hernan (2004), Hernan & Robins (2006), H&R (2020)

Causal Inference 2: Model-based estimation of causal contrasts

7 / 29

## Causal estimands: contrasts of counterfactual risks

- ▶ Let  $T^{X=x}$ , or in short  $T^x$ , be time to event, and  
 $Y^{X=x}(t) = Y^x(t) = \mathbf{1}_{\{T^x \leq t\}}$  indicate (1/0) event to occur by  $t$ ,  
 if exposure  $X$  were – **counterfactually** – forced to value  $x$  in the  
whole target population.
- ▶ The **counterfactual** risk for exposure value  $X = x$   

$$\pi^x(t) = P\{Y^{X=x}(t) = 1\} = E[Y^{X=x}(t)] = P\{T^x \leq t\}$$
- ▶ **Marginal causal contrasts** of risk, suppressing dependence on time.
  - risk difference (RD)  $\tau^c = \pi^1 - \pi^0 = P\{Y^{X=1} = 1\} - P\{Y^{X=0} = 1\}$ ,
  - risk ratio (RR)  $\phi^c = \pi^1/\pi^0$ ,
  - risk odds ratio (OR)  $\psi^c = [\pi^1/(1 - \pi^1)]/[\pi^0/(1 - \pi^0)]$ ,
- ▶ Alternative notation: Judea Pearl's (2010) **do-operator**  

$$P\{Y = 1|\text{do}(X = x)\} = P\{Y^{X=x} = 1\}$$
.

Causal Inference 2: Model-based estimation of causal contrasts

8 / 29

## Identifiability of causal contrasts of risks

- If in the pertinent causal diagram, there are **open non-causal paths** between  $X$  and  $Y$  – e.g. due to unmeasured confounders  $U$  – the causal contrasts of interest are not identified  $\Leftrightarrow$  **residual confounding**.
- Let  $Z'$  be a set of observed covariates that are **non-descendants** of  $X$ . If  $Z \subset Z'$  (discrete) were sufficient to **block** all non-causal paths, then counterfactual risks are identified by **standardization** – or **g-formula**:

$$\begin{aligned}\pi^x &= E(Y^{X=x}) = E_Z[E_Y(Y|X=x, Z)] \\ &= \sum_z P\{Y=1 | X=x, Z=z\}P\{Z=z\}.\end{aligned}$$

- Causal contrasts  $\tau^c, \phi^c, \psi^c$  are obtained from  $\pi^1$  and  $\pi^0$  thus derived.

- If  $X$  is randomized, then  $X \perp\!\!\!\perp Z \cup U$ , and it holds simply

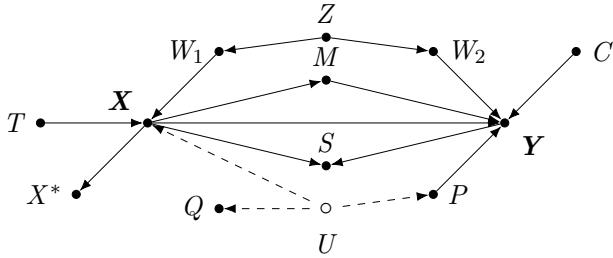
$$\pi^x = P\{Y^{X=x}=1\} = P\{Y=1 | X=x\} = \pi_x, \text{ and e.g.}$$

$$\tau^c = \pi^1 - \pi^0 = \pi_1 - \pi_0 = \tau.$$

Causal Inference 2: Model-based estimation of causal contrasts

9 / 29

## Example: Identifying causal contrast from DAG



- **Causal paths**  $X \rightarrow Y$  and  $X \rightarrow M \rightarrow Y$ : Don't block!
- **Non-causal paths** between  $X$  and  $Y$ : Block!
  - If already blocked, don't open (e.g. by conditioning on  $S$ ).
- **Backdoor paths**  $X \leftarrow W_1 \leftarrow Z \rightarrow W_2 \rightarrow Y$  and  $X \leftarrow U \rightarrow P \rightarrow Y$ : Block with minimal effort. – **Sufficient sets**:  $P$  plus one from  $\{Z, W_1, W_2\}$ . – If  $P$  unobserved, substitute by  $Q$ , proxy of  $U$ .
- No need to adjust for  $T$ . – Adjusting for  $C$  can improve precision.

Causal Inference 2: Model-based estimation of causal contrasts

10 / 29

## Conditional causal contrasts

- Let  $Z$  be a covariate (can be vector-valued). **Conditional causal effect** of  $X$  given  $Z = z$  is defined by counterfactual  $z$ -conditional risks  $\pi_z^x = P\{Y^{X=x}=1 | Z=z\}$ .
- These have their own identifiability conditions.
- For instance,  $Z = \text{sex}$ , or any other variable that divides the target population to interesting subsets or strata.
- **Conditional or  $z$ -specific causal contrasts** of risks are, for instance

$$\begin{aligned}\tau_z^c &= \pi_z^1 - \pi_z^0 = P\{Y^{X=1}=1 | Z=z\} - P\{Y^{X=0}=1 | Z=z\}, \\ \phi_z^c &= \pi_z^1/\pi_z^0 = P\{Y^{X=1}=1 | Z=z\}/P\{Y^{X=0}=1 | Z=z\}\end{aligned}$$

- If  $\tau_z^c$  has the same value for all  $z$ , RD is **homogenous**. Otherwise it is **heterogenous** or **modified** by  $Z$ . – Similarly defined for RR and OR.

Causal Inference 2: Model-based estimation of causal contrasts

11 / 29

## Causal contrasts in factual exposure groups

- ▶ Causal risk difference **among exposed** is defined  
 $\tau_1^c = P\{Y^{X=1} = 1 | X = 1\} - P\{Y^{X=0} = 1 | X = 1\}$ ,  
also known as **average treatment effect among treated (ATT)**.  
– The contrast **among unexposed (ATU)** is analogously defined.
- ▶ The effect often heterogenous, and groups noncomparable.
- ▶ If  $Z$  is a sufficient set, g-formulas for these contrasts are

$$\text{ATT} = \pi_1 - \sum \pi_{0z} P\{Z = z | X = 1\} = \text{"observed - expected"},$$

$$\text{ATU} = \sum \pi_{1z} P\{Z = z | X = 0\} - \pi_0 = \text{"expected - observed"}.$$

- ▶ Different standard populations for ATT, ATU, and for marginal contrast ATE, i.e. **average treatment effect in the whole population**:

$$\text{ATE} = \tau^c = \pi^{X=1} - \pi^{X=0} = \sum_z \pi_{1z} P\{Z = z\} - \sum_z \pi_{0z} P\{Z = z\}.$$

## Example: Causal contrasts in exposure groups

- ▶ Suppose a population in which 45 % are truly exposed.
- ▶ Let a binary  $Z$  be sufficient to remove confounding, and  
 $P\{X = 1 | Z = 1\} = 0.75$  and  $P\{X = 1 | Z = 0\} = 0.25$ .
- ▶ Let also factual risks  $\pi_{xz} = P\{Y = 1 | X = x, Z = z\}$   
by  $X$  and  $Z$  be as shown in the cells of the table below ( $x, z = 0, 1$ )
- ⇒ Values for marginal risks,  $\pi_1, \pi_0$ , crude contrast  $\tau = \pi_1 - \pi_0$ , and  
conditional contrasts  $\tau_z = \pi_{1z} - \pi_{0z}$  are shown in margins:

	$Z = 1$	$Z = 0$	$\pi_x$
$X = 1$	0.50	0.20	$\pi_1 = 0.40$
$X = 0$	0.25	0.10	$\pi_0 = 0.13$
Contrasts	$\tau_1 = 0.25$	$\tau_0 = 0.10$	$\tau = 0.27$

## Example: Causal contrasts in exposure groups (cont'd)

- ⇒ Observed risks (bolded), counterfactual risks and their causal contrasts are

True exposure group		Marginal	
	$x = 1$	$x = 0$	Marginal
Risk $\pi_x^{X=1}$ if exposed	<b>0.40</b>	0.26	$\pi^{X=1} = 0.32$
Risk $\pi_x^{X=0}$ if unexposed	0.20	<b>0.13</b>	$\pi^{X=0} = 0.16$
Contrasts $\tau_x^c$	ATT = 0.20	ATU = 0.13	$\tau^c = \text{ATE} = 0.16$

- ▶ Here, the causal risk difference is bigger among exposed. Thus, being exposed seems to be a modifier of the effect of exposure on this scale!
- ▶ Interestingly, the causal risk ratio is homogenous.

**NB** Popular design for straightforward estimation of ATT:  
**matched cohort study**.

## Causal contrasts of hazards

- The counterfactual hazard, when everybody were exposed to  $X = x$ :

$$\lambda^x(t) = \lim_{h \rightarrow 0} \frac{1}{h} P\{t < T^x \leq t + h \mid T^x > t\} = \lim_{h \rightarrow 0} \frac{1}{h} \frac{P\{Y^x(t+h) = 1\}}{P\{Y^x(t) = 0\}}$$

- Marginal causal contrasts of hazards:

– hazard difference (HD)  $\delta^c(t) = \lambda^1(t) - \lambda^0(t)$ ,

– hazard ratio (HR)

$$\rho^c(t) = \frac{\lambda^1(t)}{\lambda^0(t)} = \frac{\lim P\{Y^1(t+h) = 1 \mid Y^1(t) = 0\}/h}{\lim P\{Y^0(t+h) = 1 \mid Y^0(t) = 0\}/h}$$

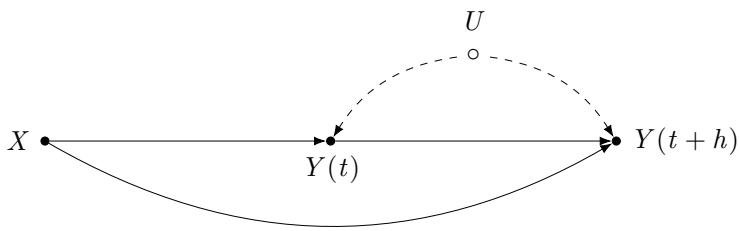
- Identified by corresponding associational contrasts, if  $X$  is randomized.

- However, if  $X$  has any effect, conditions  $Y^1(t) = 0$  and  $Y^0(t) = 0$  imply that when  $t > 0$ , the counterfactual populations at risk will be different.

- Even if exposure groups were comparable at  $t = 0$ , after that they are not.

## Hazard of hazard ratios

- The hazard at any time  $t > 0$  is affected by known and unknown causes of the outcome  $\Rightarrow$  individual **frailty**  $U$  varies in the population.



- $Y(t)$  is a **collider** on the path from  $X$  to  $Y(t+h)$  via  $U$ .
- Conditioning on  $Y(t) = 0$  opens this non-causal path, thus inducing bias when evaluating the causal effect of  $X$  by  $\rho^c(t)$  for  $t > 0$ .

See e.g. [Aalen et al. \(2015\)](#)

## Hazard of hazard ratios (cont'd)

- If exposure to  $X = 1$  increases the risk vs.  $X = 0$ , among the more frail subjects, those exposed tend to get the outcome earlier than the unexposed. Later on, the remaining exposed are less frail than the unexposed.
- $\Rightarrow$  Interpretation of hazard ratios in causal terms is problematic, even if  $\rho(t)$  were perfectly constant over time (like in Cox regression).
- This problem is less pronounced with rare outcomes than common ones. Yet, causal contrasts of risks are generally preferable to those of hazards
- Moreover, due to frailty,  $\rho(t)$  may be far from constant. It can be quite high early on, but then decrease and even dive below 1 at some  $t^*$ . Then, it may be hazardous to conclude that the exposure is preventive for  $t > t^*$ .
- Thus, letting  $\rho(t)$  to vary by  $t$  in modelling is often advised.

See e.g. [Hernan \(2010\)](#), [Stensrud et al. \(2019\)](#), [Martinussen \(2021\)](#)

## Example: WHI Trial on MHT and CHD (Manson et al. 2003)

Table 2. Estrogen plus Progestin and the Risk of CHD, According to Year of Follow-up. <sup>a</sup>			
Year of Follow-up	CHD	Hazard Ratio for CHD (95% CI)	
	Estrogen-plus-Progestin Group	Placebo Group	
no. of cases (annualized percentage)			
1	42 (0.50)	23 (0.29)	1.81 (1.09–3.01)
2	38 (0.45)	28 (0.35)	1.34 (0.82–2.18)
3	19 (0.23)	15 (0.19)	1.27 (0.64–2.50)
4	32 (0.39)	25 (0.32)	1.25 (0.74–2.12)
5	29 (0.41)	19 (0.28)	1.45 (0.81–2.59)
≥6	28 (0.37)	37 (0.56)	0.70 (0.42–1.14)

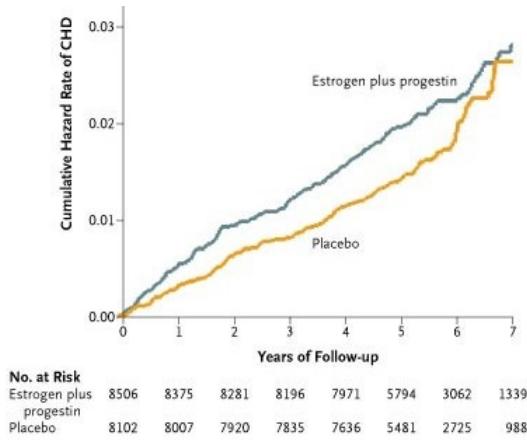
\* CHD includes acute myocardial infarction (MI) necessitating hospitalization, silent myocardial infarction as determined by serial electrocardiography, and death due to CHD. There were nine silent myocardial infarctions (four in the estrogen-plus-progestin group and five in the placebo group). Hazard ratios are stratified according to age, presence or absence of a previous coronary event, and randomly assigned diet modification group and are adjusted for previous coronary-artery bypass grafting or percutaneous transluminal cor-

Causal Inference 2: Model-based estimation of causal contrasts

18 / 29

- ▶ Women, 50–79 y, MHT:  $N_1 = 8506$ , placebo:  $N_0 = 8102$
- ▶ Followed-up for max 8.6 y, mean 5.6 y.
- ▶ Cases & rates/ $10^4$  y  $D_1 = 188$ ,  $I_1 = 39$ ,  $D_0 = 147$ ,  $I_0 = 33$ .
- ▶ Crude IR = 1.20, adjusted 1.24 (1.00–1.54)

## Example: WHI Trial (cont'd)



- ▶ Curves of cumulative hazard approximate the development of cumulative risks  $\pi_x(t)$  over time.
- ▶ In early years, the curve of MHT runs on top, reflecting higher hazard in that period.
- ▶ By 6–7 year, cumulative risks appear to have reached same level.

Causal Inference 2: Model-based estimation of causal contrasts

19 / 29

## From hazards to causal contrasts of risk

- ▶ With censored data and especially when the exposure is not randomized, causal contrasts of risk require a valid model for hazards as function of time, exposure  $X$ , and relevant covariates  $Z$ .
- ▶ Suppose  $Z$  blocks all non-causal paths. Then **counterfactual conditional hazards**  $\lambda^x(t|Z = z)$  are identified by observable hazards  $\lambda_x(t|Z = z)$ :

$$\lambda_x(t|Z = z) = \lim_{h \rightarrow 0} P\{Y(t+h) = 1 | Y(t) = 0, X = x, Z = z\}/h$$

- ▶ When  $Z$  is continuous or vector-valued, a reasonable model for  $\lambda_x(t|Z = z)$  can often be built on e.g. Cox regression or its time-dep extensions.
- ▶ Discrete-time hazards  $P\{Y(t_{k+1}) = 1 | Y(t_k) = 0, X = x, Z = z\}$  and binary models on these also used (e.g. by Hernan & Robins).

Causal Inference 2: Model-based estimation of causal contrasts

20 / 29

## From hazards to causal contrasts of risk (cont'd)

- ▶ When no competing events exist, counterfactual  $z$ -conditional risks  $\pi^x(t|Z = z)$  are identified from factual  $z$ -conditional hazards and risks:  
$$\pi^x(t|Z = z) = \pi_x(t|Z = z) = 1 - \exp \left\{ - \int_0^t \lambda_x(u|Z = z) du \right\}.$$
- ▶ Counterfactual marginal risks are then obtained using the g-formula.  
If  $Z$  is discrete-valued, the formula is

$$\pi^x(t) = \sum_z \pi_x(t|Z = z) P\{Z = z\}.$$

- ▶ Thus, marginal risks are again weighted averages of the conditional ones, and use of g-formula corresponds to **direct standardization**.

NB. For those more mathematically oriented: If  $Z$  contains discrete and/or continuous variables, the sum is substituted by a **Stieltjes-integral**, and the g-formula is technically expressed as  $\pi^x(t) = E_Z[\pi_x(t|Z)] = \int_z \pi_x(t|Z = z) dF(z)$ , where  $F(z)$  is the joint distribution function of  $Z$  in the population.

## Causal modelling: Outcome regression

Modelling how expected values, risks, hazards, etc. depend on exposure  $X$  and covariates  $Z$  (modifiers, and/or confounders). – Common elements:

- ▶ Each subject  $i$  ( $i = 1, \dots, N$ ) has an own **profile**, i.e. vector  $(x_i, z_i^T)$  of values of  $X$  and covariates  $Z$ .
- ▶ In the spirit of **generalized linear models**, let vector  $(\alpha, \beta, \gamma^T)$  contain regression coefficients, and specify the **linear predictor**  
– assuming so far no **interactions**, nor **effect modifications**

$$\eta_i = \alpha + \beta x_i + \gamma^T z_i$$

- ▶ **Product terms** can be added for interactions and modifications if needed, and **splines** may be used for continuous covariates.
- ▶ Further model specification depends on the type of outcome variable, causal contrasts of interest, and importance and choice of time scale(s).

## Binary regression and classical interpretations of coefficients

- ▶ Basic model for risks  $\pi$  in fixed risk periods with complete follow-up without censoring and competing events:

$$g\{\pi(x_i)\} = \alpha + \beta x_i + \gamma^T z_i, \quad i = 1, \dots, N.$$

- ▶ Link  $g(\cdot)$  and causal interpretation of  $\beta$ , assuming the validity of model (including homogeneity or non-modification of the contrast in question) and that  $Z$  blocks all backdoor paths:
  - id  $\Rightarrow \beta =$  risk difference (RD)  $\tau^c$  for  $X = 1$  vs.  $X = 0$ , adjusted for  $Z$
  - log  $\Rightarrow \beta =$  log of risk ratio (RR)  $\phi^c$  – “ –
  - logit  $\Rightarrow \beta =$  log of conditional risk odds ratio (OR),  $\psi_z^c$ , – “ –

**NB.** This is different from marginal OR due to **non-collapsibility**.
- ▶ Random component: Binomial family – Fitting: some GLM program.
- ▶ Issues with id & log in keeping predicted  $\hat{\pi}(\cdot)$  between 0 and 1.

## Binary regression and causal contrasts by g-formula

- ▶ Assuming that  $Z$  is sufficient to block non-causal paths, a logistic model is fitted, which may even contain product terms allowing modification

$$\text{logit}(\pi_i) = \log[\pi_i/(1 - \pi_i)] = \alpha + \beta x_i + \gamma^T z_i + \delta^T(x_i z_i), \quad i = 1, \dots, n.$$

- ▶ For each individual  $i$ , predicted risks are computed for both  $X = 1$  and  $X = 0$ , keeping  $Z = z_i$  as it is

$$\tilde{\pi}_i^{X_i=x} = \text{expit}\{\hat{\alpha} + \hat{\beta}x + \hat{\gamma}^T z_i + \hat{\delta}^T(xz_i)\}, \quad x = 0, 1.$$

- ▶ Marginal counterfactual risks for  $x = 1, 0$  are estimated applying **g-formula**:

$$\hat{\pi}^{X=x} = \hat{E}_Z[E(Y|X=x, Z)] = \frac{1}{n} \sum \tilde{E}(Y_i|X_i=x, Z=z_i) = \frac{1}{n} \sum \tilde{\pi}_i^{X_i=x}$$

as the data provide a non-parametric estimate of the joint distribution of  $Z$

- ▶ Estimators of **marginal causal contrasts** of risks are now, e.g.

$$\hat{\tau}^c = \hat{\pi}^{X=1} - \hat{\pi}^{X=0}, \quad \hat{\psi}^c = [\hat{\pi}^1/(1 - \hat{\pi}^1)]/[\hat{\pi}^0/(1 - \hat{\pi}^0)]$$

## Exposure modelling, propensity scores and weighting

Let  $X$  be a binary exposure variable. Assume again that  $Z$  is a sufficient set

- ▶ **Exposure model** predicting individual  $X_i$ :s by confounders is fitted

$$\text{logit}[P\{X_i = 1|Z = z_i\}] = \alpha^* + z_i^T \gamma^*, \quad i = 1, \dots, N.$$

- ▶ **Propensity scores**  $\text{PS}_i$ , or fitted probabilities of exposure:

$$\text{PS}_i = \hat{P}\{X_i = 1|Z = z_i\} = \text{expit}(\hat{\alpha}^* + z_i^T \hat{\gamma}^*).$$

- ▶ Individual **weights**  $W_i = w(\text{PS}_i, X_i)$  are computed (see next slide).
- ▶ Counterfactual risks are estimated as weighted averages of the outcome in the two exposure groups

$$\hat{\pi}^{X=x} = \frac{\sum_{i=1}^n \mathbf{1}_{\{X_i=x\}} W_i Y_i}{\sum_{i=1}^n \mathbf{1}_{\{X_i=x\}} W_i} = \frac{\sum_{X_i=x} W_i Y_i}{\sum_{X_i=x} W_i}, \quad x = 0, 1$$

- ▶ From these, marginal causal contrasts are estimated as before.

## Exposure modelling, propensity scores and weighting (cont'd)

- ▶ **Inverse probability weights** (IPW) are used to estimate marginal causal contrasts (like ATE) in the whole target population. They are based on inverses of the fitted probabilities of belonging to the realized exposure group:

$$W_i = w(\text{PS}_i, X_i) = \frac{\mathbf{1}_{\{X_i=1\}}}{\text{PS}_i} + \frac{\mathbf{1}_{\{X_i=0\}}}{1 - \text{PS}_i}, \quad i = 1, \dots, n.$$

- ▶ If the interest is on causal contrasts among the treated (like ATT), the **treated weights** are used:  $W_i = 1$  for  $X_i = 1$ , and  $W_i = \text{PS}_i/(1 - \text{PS}_i)$  for  $X_i = 0$ .
- ▶ Other: **overlap weights**, **matching weights**, **entropy weights**.
- ▶ The goodness-of-fit of the exposure model needs to be assessed. For that purpose, various measures of **covariate balance** are developed.

## Double robust (DR) estimators and machine learning methods

- ▶ The validity of either g-formula or PS-based weighting depends on, how accurately the outcome model or exposure model is specified.
- ▶ **Double robust (DR)** estimation of causal contrasts:  
Combination of g-formula and IPW. – Alternatives
  - **Augmented IPW** (AIPW); see [Jonsson Funk et al. \(2011\)](#),
  - **Targeted maximum likelihood estimation** (TMLE);  
see [Schuler & Rose \(2015\)](#), [Luque-Fernandez et al. \(2018\)](#)

Validity of a DR estimator requires that either the exposure model or the outcome model is correctly specified.

- ▶ Algorithms developed for **supervised learning** increase flexibility in modelling both outcome and exposure (see [Bi et al. 2019](#), [Blakely et al. 2020](#)).

## Causal contrasts from censored time-to-event data

- ▶ Various methods to estimate counterfactual risks  $\pi^{X=x}(t)$  and their contrasts (see [Denz et al. 2023](#)) – For instance
- (a) Fit a Cox model  $\lambda(t|x_i, z_i) = \lambda_0(t) \exp(\beta x_i + \gamma^T z_i)$ , take estimates of coefficients and baseline cumulative hazard  $\hat{\Lambda}_0(t)$  from which:
$$\tilde{\pi}_i^{X=x}(t) = 1 - \exp\{-\hat{\Lambda}_0(t) \exp(\hat{\beta}x_i + \hat{\gamma}^T z_i)\}.$$
Counterfactuals  $\pi^{X=x}(t)$  and contrasts are then estimated by g-formula.
- (b) Get weights  $W_i$  from an exposure model, fit Cox with “intercept only” and  $X$  as a strata() variable and  $W_i$ :s as weights, and estimate  $\hat{\pi}^{X=x}(t)$  using `survfit()`, etc.
- ▶ Other: IPW Kaplan-Meier, use of pseudo-values, DR methods, ...
- ▶ Competing event settings: additional complexities in defining and analysing causal contrasts (see [Rudolph et al. \(2020\)](#), [Young et al. \(2020\)](#)).

## Conclusion

- ▶ Careful specification of causal question and estimands needed.
- ▶ Selection of confounders: efficient blocking of backdoor paths.
- ▶ Basic estimation methods: outcome regression & g-formula, exposure modelling & PS-weighting, double robust estimators.
- ▶ Sufficiently flexible models desirable to reduce misspecification bias.
- ▶ Statistical inference: robust covariance matrix & delta method, bootstrapping, efficient influence curve, etc. – ignored here.
- ▶ We also limited to time-fixed exposure (binary) and confounders.
- ▶ Extensions are available to deal with polytomous exposure as well as time-varying exposure and confounding.

# Multistate models

**Bendix Carstensen** Steno Diabetes Center Copenhagen  
Herlev, Denmark  
<http://BendixCarstensen.com>

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

From C:\Bendix\teach\SPE\git\lectures\multistate\multistate.tex

Monday 29 May, 2023, 13:01

1 / 27

# Multistate models

**Bendix Carstensen, Martyn Plummer**

Multistate models

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

ms-Markov

## Common assumptions in survival analysis

1. Subjects are **either** “healthy” **or** “diseased”, with no intermediate state.
2. The disease is **irreversible**, or requires intervention to be cured.
3. The time of disease incidence is known **exactly**.
4. The disease is **accurately** diagnosed.

These assumptions are true for **death** and many **chronic diseases**.

A question of definition:

– consider occurrence of **recording of** a given disease

## A model for cervical cancer

Invasive squamous cell cancer of the cervix is preceded by cervical intraepithelial neoplasia (CIN)



Purpose of a screening programme is to detect and treat CIN — status of persons obtained at screening dates

Aim of the modeling the **transition rates** between **states**, is to be able predict how population moves between **states**

- ▶ Transition rates between states
- ▶ Probability of state occupancy

## Markov models for multistate processes

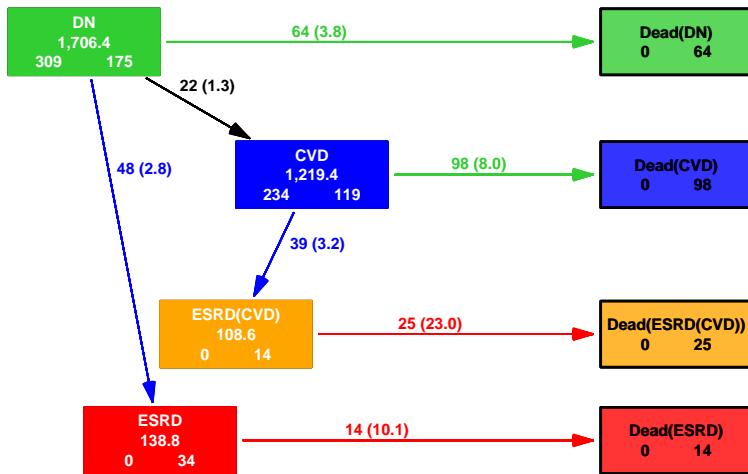
The natural generalization of Poisson regression to multiple disease states:

- ▶ transition between states depends **only** on current state
- ▶ — this is the **Markov** property
- ▶ ⇒ transition rates are constant over time
- ▶ (time-fixed) covariates may influence transition rates
- ▶ the formal Markov property is **very** restrictive
- ▶ in the clinical litterature “Markov model” is often used about any type of multistate model

## Components of a multistate (Markov) model

- ▶ Define the disease states
- ▶ Define which transitions between states are allowed
- ▶ Select covariates influencing transition rates (may be different between transitions)
- ▶ Not a trivial task — do we want e.g.
  - ▶ cause of death (CVD, Cancer, Other)
  - ▶ disease status at death (prev.CVD, prev.Can, neither)

## A more complicated multistate model



Multistate models (ms-Markov)

6 / 27

## Likelihood for a multistate model

- ▶ The likelihood of the model depends on the probability of being in state  $j$  at time  $t_1$ , given that you were in state  $i$  at time  $t_0$ .
- ▶ Assume transition rates constant in small time intervals
- ▶ ⇒ each interval for a person contributes term(s) to the likelihood
- ▶ one term for each possible transition between states
- ▶ the total likelihood for person  $p$  in intervals  $i$  is a product of these terms,  $d_{pi} \log(\lambda_{pi}) - \lambda_{pi} y_{pi}$
- ▶ ⇒ each term has the form of the likelihood for a Poisson variate  $d$  with mean  $\lambda y$

Multistate models (ms-Markov)

7 / 27

## Likelihood for a multistate model

- ▶ each term has the form of the likelihood for a Poisson variate  $d$  with mean  $\lambda y$
- ▶ terms are **not** independent, but the total likelihood is a product; hence of the same form as the likelihood from independent Poisson variates
- ▶ but observations from intervals from one person are neither Poisson nor independent

Multistate models (ms-Markov)

8 / 27

## Realms of multistate modeling

- ▶ intensities — dimension  $\text{time}^{-1}$
- ▶ state probabilities — dimensionless,  $\text{time}^0$   
 $\text{integral}$  of intensities w.r.t. to time
- ▶ sojourn times — dimension  $\text{time}^1$   
 $\text{integral}$  of state probabilities w.r.t. to time

## Classes of multistate models

- ▶ Markov model: transition between states depends **only** on current state  $\Rightarrow$  transition rates are constant  
**time-homogeneous Markov model**
- ▶ If transition rates depend on the **same timescale** only we have a **time-inhomogeneous Markov model**
- ▶ If transition rates depend on the time since entry to the current state we have a **semi-Markov model**
- ▶ If transition rates depend on several timescales we have a **general multistate model** (there is no name for this)

## Computing state probabilities from intensities in multistate models

- ▶ time-homogeneous Markov model:  
closed-form formulae exist
- ▶ time-inhomogeneous Markov model:  
closed-form formulae exist (a bit more complicated)
- ▶ semi-Markov model:  
no closed form formulae exist
- ▶ general multistate model:  
no closed form formulae exist

No formulae means that any inference on state probabilities and sojourn times must be based on **simulation** from the model.

# Multistate models with Lexis

Bendix Carstensen

Multistate models

SPE, Tartu, Estonia,

June 2023

<http://BendixCarstensen.com/SPE>

ms-Lexis

## Example: Renal failure data from Steno

Hovind P, Tarnow L, Rossing P, Carstensen B, and Parving H-H: Improved survival in patients obtaining remission of nephrotic range albuminuria in diabetic nephropathy. *Kidney Int.*, 66(3):1180–1186, 2004.

- ▶ Endpoint of interest: Death or end stage renal disease (ESRD), i.e. dialysis or kidney transplant.
- ▶ 96 patients entering at nephrotic range albuminuria (NRA), i.e. U-alb > 300mg/day.
- ▶ Is remission from this condition (i.e return to U-alb < 300mg/day) predictive of the prognosis?

Multistate models with Lexis (ms-Lexis)

12 / 27

	Total	Remission	
		Yes	No
No. patients	125	32	93
No. events	77	8	69
Follow-up time (years)	1084.7	259.9	824.8

Cox-model:

Timescale: Time since nephrotic range albuminuria (NRA)

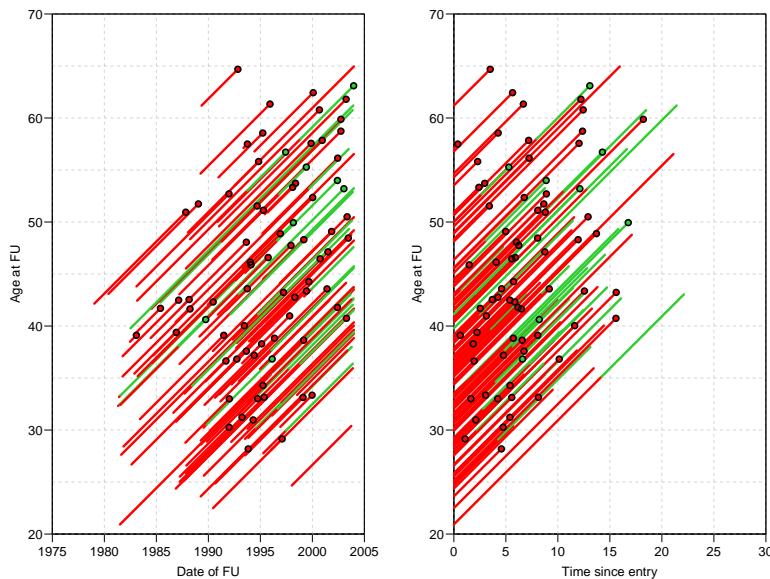
Entry: 2.5 years of GFR-measurements after NRA

Outcome: ESRD or Death

Estimates:	RR	95% c.i.	p
Fixed covariates:			
Sex (F vs. M):	0.92	(0.53,1.57)	0.740
Age at NRA (per 10 years):	1.42	(1.08,1.87)	0.011
Time-dependent covariate:			
Obtained remission:	0.28	(0.13,0.59)	0.001

Multistate models with Lexis (ms-Lexis)

13 / 27



Multistate models with Lexis (ms-Lexis)

14 / 27

## Features of the analysis

- ▶ Remission is included as a time-dependent variable.
- ▶ Age at entry is included as a fixed variable.

```
renal[1:5,]
id      dob      doe      dor      dox event
17 1967.944 1996.013      NA 1997.094    2
26 1959.306 1989.535 1989.814 1996.136    1
27 1962.014 1987.846      NA 1993.239    3
33 1950.747 1995.243 1995.717 2003.993    0
42 1961.296 1987.884 1996.650 2003.955    0
```

Note patient 26, 33 and 42 obtain remission.

Multistate models with Lexis (ms-Lexis)

15 / 27

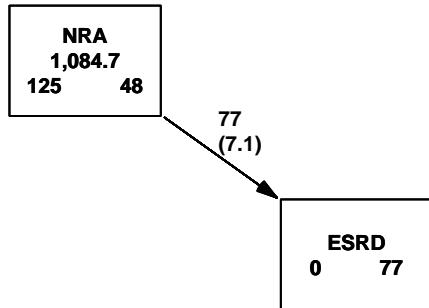
```
> Lr <- Lexis(entry = list(per = doe,
+                         age = doe-dob,
+                         tfi = 0),
+               exit = list(per = dox),
+               exit.status = event>0,
+               states = c("NRA", "ESRD"),
+               data = renal)
> summary(Lr)

Transitions:
  To
From  NRA  ESRD  Records:  Events: Risk time: Persons:
   NRA    48     77      125       77    1084.67      125
```

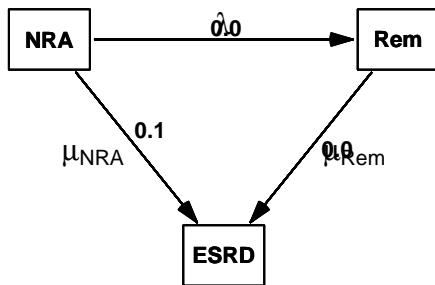
Multistate models with Lexis (ms-Lexis)

16 / 27

```
> boxes(Lr, boxpos = list(x = c(25, 75),
+                           y = c(75, 25)),
+       scale.R = 100, show.BE = TRUE )
```



## Illness-death model



$\lambda$ : remission rate.

$\mu_{NRA}$ : mortality/ESRD rate **before** remission.

$\mu_{rem}$ : mortality/ESRD rate **after** remission.

## Cutting follow-up at remission: cutLexis

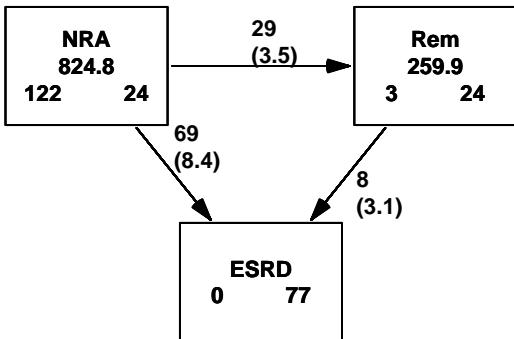
```
> Lc <- cutLexis(Lr, cut = Lr$dor,
+                   timescale = "per",
+                   new.state = "Rem",
+                   precursor.states = "NRA")
> summary(Lc)
```

Transitions:

To	From	NRA	Rem	ESRD	Records:	Events:	Risk time:	Persons:
NRA	NRA	24	29	69	122	98	824.77	122
Rem	Rem	0	24	8	32	8	259.90	32
	Sum	24	53	77	154	106	1084.67	125

## Showing states and FU: boxes.Lexis

```
> boxes(Lc, boxpos = list(x = c(15, 85, 50),
+                           y = c(85, 85, 20)),
+       scale.R = 100, show.BE = TRUE)
```



Multistate models with Lexis (ms-Lexis)

20 / 27

## Cutting follow up at events: cutLexis

```
> Lc <- cutLexis( Lr, cut = Lr$dor,
+                   timescale = "per",
+                   new.state = "Rem",
+                   precursor.states = "NRA",
+                   split.states = TRUE )
> summary( Lc )

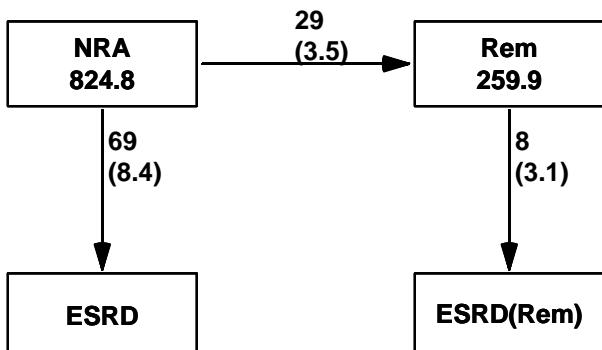
Transitions:
  To
From  NRA  Rem  ESRD  ESRD(Rem)  Records:  Events:  Risk time:  Persons:
  NRA   24    29    69      0        122     98   824.77    122
  Rem    0    24     0      8        32      8   259.90     32
  Sum   24    53    69      8       154    106  1084.67    125
```

Multistate models with Lexis (ms-Lexis)

21 / 27

## Showing states and FU: boxes.Lexis

```
> boxes(Lc, boxpos = list(x = c(15, 85, 15, 85),
+                           y = c(85, 85, 20, 20)),
+       scale.R = 100)
```

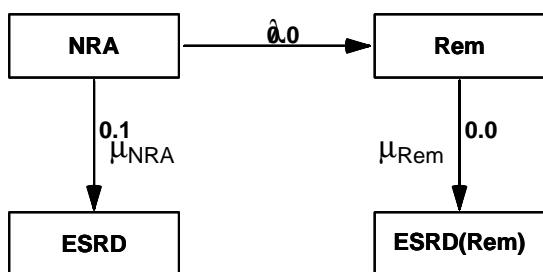


Multistate models with Lexis (ms-Lexis)

22 / 27

## Likelihood for a general MS-model

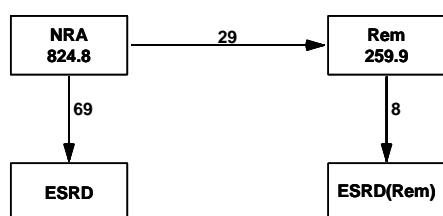
- ▶ Product of likelihoods for each transition
  - each one as for a survival model
- ▶ **Risk time** is the risk time in the “From” state
- ▶ **Events** are transitions to the “To” state
- ▶ All other transitions out of “From” are treated as **censorings**
- ▶ Possible to fit models separately for each transition



Cox-analysis with remission as time-dependent covariate:

- ▶ Ignores  $\lambda$ , the remission rate.
- ▶ Assumes  $\mu_{NRA}$  and  $\mu_{rem}$  use the same timescale.

## Model for all transitions



### Cox-model:

- ▶ Different timescales for transitions possible
- ▶ ... only one per transition
- ▶ No explicit representation of estimated rates.

### Poisson-model:

- ▶ Timescales can be different
- ▶ Multiple timescales can be accommodated simultaneously
- ▶ Explicit representation of all transition rates

## Calculating state probabilities

$$\begin{aligned} & P \{ \text{Remission before time } t \} \\ &= \int_0^t \lambda(u) \exp \left( - \int_0^u \lambda(s) + \mu_{\text{NRA}} ds \right) du \end{aligned}$$

$$\begin{aligned} & P \{ \text{Being in remission at time } t \} \\ &= \int_0^t \lambda(u) \exp \left( - \int_0^u \lambda(s) + \mu_{\text{NRA}}(s) ds \right) \times \\ & \quad \exp \left( - \int_u^t \mu_{\text{rem}}(s) ds \right) du \end{aligned}$$

Note  $\mu_{\text{rem}}$  could also depend on  $u$ , time since obtained remission.

Sketch of programming, assuming that  $\lambda$  (`lambda`),  $\mu_{\text{NRA}}$  (`mu.nra`) and  $\mu_{\text{rem}}$  (`mu.rem`) are known at any age (stored in vectors)

```
c.rem      <- cumsum(lambda)
c.mort.nra <- cumsum(mu.nra)
c.mort.rem <- cumsum(mu.rem)
pr1 <- cumsum(lambda * exp(-(c.rem + c.mort.nra)))

intgr(t,s) <-
function(t,s){
lambda[s] * exp(-(c.rem[s] + c.mort.nra[s])) *
exp(-(c.mort.rem[t] - c.mort.rem[s]))
}
for(t in 1:100) p2[t] <- sum(intgr(t,1:t))
```

If  $\mu_{\text{rem}}$  also depends on time since remission, then `c.mort.rem` should have an extra argument—technically very complicated