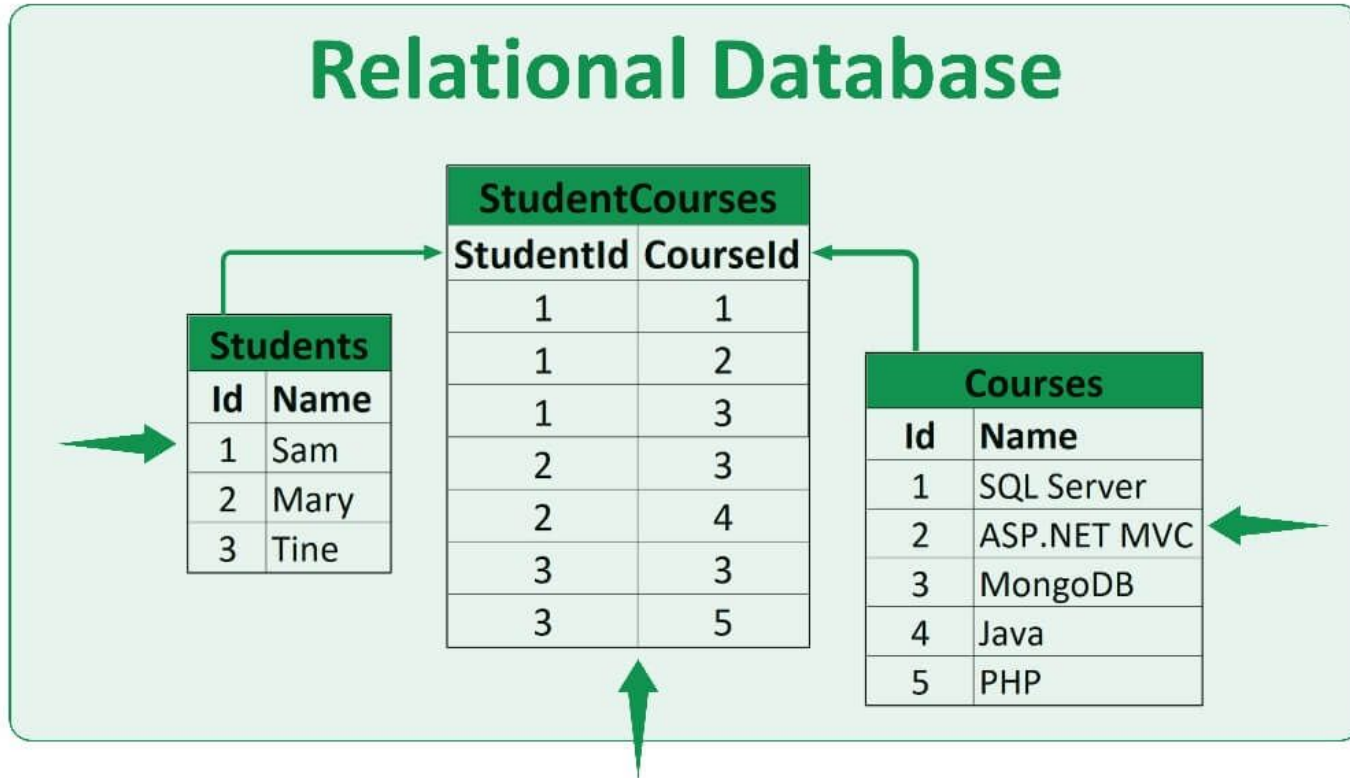


Afternoon stretch: How are these tables related?



Time, Dates, and Joins

Dates and Time

- tidyverse data types:
 - <date> : date, can take many formats
 - <time> : time, can include time zone
 - <dtm> : date plus time, a unique instant in time
 - Base R calls it POSIXct, but it is the same thing
- always use the simplest possible data type



```
flights |>
  select(year, month, day, hour, minute) |>
  mutate(departure = make_datetime(year, month, day, hour, minute))
#> # A tibble: 336,776 × 6
#>   year month   day hour minute departure
#>   <int> <int> <int> <dbl> <dbl> <dtm>
#> 1  2013     1     1     5     15 2013-01-01 05:15:00
#> 2  2013     1     1     5     29 2013-01-01 05:29:00
#> 3  2013     1     1     5     40 2013-01-01 05:40:00
```

Math with Date-times — Lubridate provides three classes of timespans to facilitate math with dates and date-times.

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

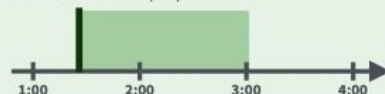
A normal day

```
nor <- ymd_hms("2018-01-01 01:30:00", tz="US/Eastern")
```



Periods track changes in clock times, which ignore time line irregularities.

`nor + minutes(90)`



Durations track the passage of physical time, which deviates from clock time when irregularities occur.

`nor + dminutes(90)`



Intervals represent specific intervals of the timeline, bounded by start and end date-times.

`interval(nor, nor + minutes(90))`



Math with Date-times — Lubridate provides three classes of timespans to facilitate math with dates and date-times.

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

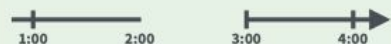
A normal day

```
nor <- ymd_hms("2018-01-01 01:30:00",tz="US/Eastern")
```



The start of daylight savings (spring forward)

```
gap <- ymd_hms("2018-03-11 01:30:00",tz="US/Eastern")
```



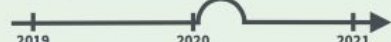
The end of daylight savings (fall back)

```
lap <- ymd_hms("2018-11-04 00:30:00",tz="US/Eastern")
```



Leap years and leap seconds

```
leap <- ymd("2019-03-01")
```



Periods track changes in clock times, which ignore time line irregularities.

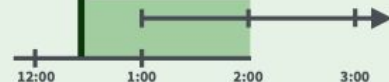
`nor + minutes(90)`



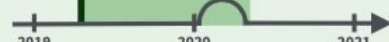
`gap + minutes(90)`



`lap + minutes(90)`



`leap + years(1)`

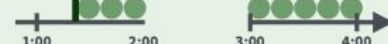


Durations track the passage of physical time, which deviates from clock time when irregularities occur.

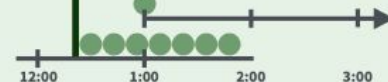
`nor + dminutes(90)`



`gap + dminutes(90)`



`lap + dminutes(90)`



`leap + dyears(1)`



Intervals represent specific intervals of the timeline, bounded by start and end date-times.

`interval(nor, nor + minutes(90))`



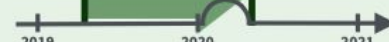
`interval(gap, gap + minutes(90))`



`interval(lap, lap + minutes(90))`



`interval(leap, leap + years(1))`



Unix Epoch: January 1, 1970

The way R is able to add or subtract dates is by converting it to time since the Unix epoch.

Returned when converting date class to numeric class.

Give it a try, what do these calls return?

Call	Returns
<code>Sys.time()</code>	
<code>today()</code>	
<code>as.numeric(Sys.time())</code>	
<code>as.numeric(today())</code>	

Unix Epoch: January 1, 1970

The way R is able to add or subtract dates is by converting it to time since the Unix epoch.

Returned when converting date class to numeric class.

Give it a try, what do these calls return?

Call	Returns
<code>Sys.time()</code>	"YYYY-MM-DD HH:MM:SS TZ"
<code>today()</code>	"YYYY-MM-DD"
<code>as.numeric(Sys.time())</code>	Seconds since epoch
<code>as.numeric(today())</code>	Days since epoch

ISO8601: international standard for writing dates from biggest to smallest separated by "-"

Standardizing Non- ISO8601 Formats

- Lot's of different ways to express dates and time
 - MM/DD/YY
 - DD - MM - YY
 - DD Month YY
 - AM/PM
 - 24 HR
 - Time zones
 - Daylight savings
 - ...
- R can manipulate all of these types of dates by parsing the strings to create a uniform format ready for analysis

2017-11-28T14:02:00

2017-22-12 10:00:00

11/28/2017 1:02:03

1 Jan 2017 23:59:59

20170131

July 4th, 2000

4th of July '99

2001: Q3

07-2020

2:01

Non-ISO8601 Formats

- Lot's of different ways to express dates and time
 - MM/DD/YY
 - DD - MM - YY
 - DD Month YY
 - AM/PM
 - 24 HR
 - Time zones
 - Daylight savings
 - ...
- R can manipulate all of these types of dates by parsing a string to create a uniform format ready for analysis

Date-times



2017-11-28 12:00:00

A **date-time** is a point on the timeline stored as the number of seconds since 1970-01-01 00:00:00 UTC

dt < **as_datetime**(1511870400)
"2017-11-28 12:00:00 UTC"

PARSE DATE-TIMES (Convert strings or numbers to date-times)

1. Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
2. Use the function below whose name replicates the order. Each accepts a *tz* argument to set the time zone, e.g. `ymd(x, tz = "UTC")`.

2017-11-28T14:02:00 **ymd_hms()**, **ymd_hm()**, **ymd_h()**.
`ymd_hms("2017-11-28T14:02:00")`

2017-22-12 10:00:00 **ydm_hms()**, **ydm_hm()**, **ydm_h()**.
`ydm_hms("2017-22-12 10:00:00")`

11/28/2017 1:02:03 **mdy_hms()**, **mdy_hm()**, **mdy_h()**.
`mdy_hms("11/28/2017 1:02:03")`

1 Jan 2017 23:59:59 **dmy_hms()**, **dmy_hm()**, **dmy_h()**.
`dmy_hms("1 Jan 2017 23:59:59")`

20170131 **ymd()**, **ydm()**. `ymd(20170131)`

July 4th, 2000 **mdy()**, **myd()**. `mdy("July 4th, 2000")`

4th of July '99 **dmy()**, **dym()**. `dmy("4th of July '99")`

2001: Q3 **yq()** Q for quarter. `yq("2001: Q3")`

07-2020 **my()**, **ym()**. `my("07-2020")`

2:01 **hms::hms()** Also **lubridate::hms()**, **hm()** and **ms()**, which return periods.* `hms::hms(seconds = 0, minutes = 1, hours = 2)`

Defining Custom Formats

Similar to regex, there is a syntax to express the format of your data.

How would you express:

1. May 30, 1985
2. 05/30/85
3. 1985-05-30

Type	Code	Meaning	Example
Year	%Y	4 digit year	2021
	%y	2 digit year	21
Month	%m	Number	2
	%b	Abbreviated name	Feb
	%B	Full name	February
Day	%d	One or two digits	2
	%e	Two digits	02
Time	%H	24-hour hour	13
	%I	12-hour hour	1
	%p	AM/PM	pm
	%M	Minutes	35
	%S	Seconds	45
	%OS	Seconds with decimal component	45.35
	%Z	Time zone name	America/Chicago
	%Z	Offset from UTC	+0800

Defining Custom Formats

Similar to regex, there is a syntax to express the format of your data.

How would you express:

May 30, 1985	%B %e %Y
05/30/85	%m/%d/%y
1985-05-30	%Y-%m-%d

<https://r4ds.hadley.nz/datetimes#tbl-date-formats>

Type	Code	Meaning	Example
Year	%Y	4 digit year	2021
	%y	2 digit year	21
Month	%m	Number	2
	%b	Abbreviated name	Feb
	%B	Full name	February
Day	%d	One or two digits	2
	%e	Two digits	02
	%H	24-hour hour	13
Time	%I	12-hour hour	1
	%p	AM/PM	pm
	%M	Minutes	35
	%S	Seconds	45
	%OS	Seconds with decimal component	45.35
	%Z	Time zone name	America/Chicago
	%Z	Offset from UTC	+0800

Questions and switch to M6.R

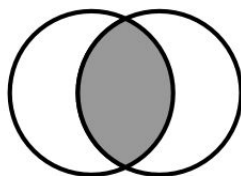
Joins: Because data is never in just one place



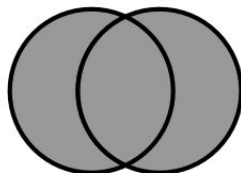
Joins are a super useful diagnostic tool if you have a strong conceptual understanding of how they work

Conceptual explanations often use diagrams like:

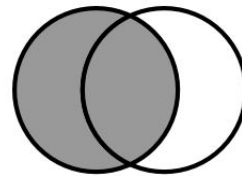
x		y	
1	x1	1	y1
2	x2	2	y2
3	x3	4	y3



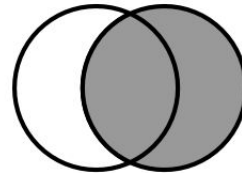
`inner_join(x, y)`



`full_join(x, y)`



`left_join(x, y)`



`right_join(x, y)`

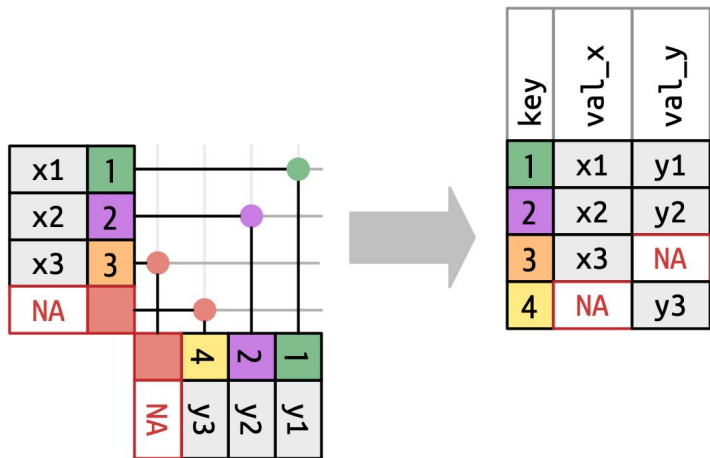
Joins: Because data is never in just one place



Combining data sets

x		y	
1	x1	1	y1
2	x2	2	y2
3	x3	4	y3

Mutation Joins: add new variables to one data frame by matching observations in another.



The key(s) is used to decide which observations (rows) go where.

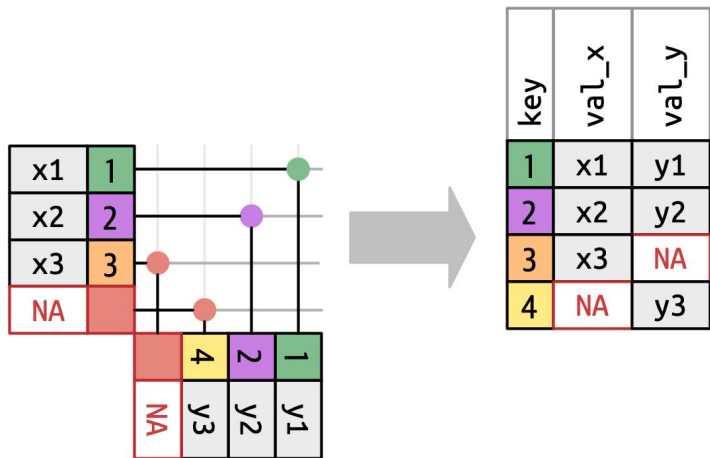
Joins: Because data is never in just one place



Combining data sets

x		y	
1	x1	1	y1
2	x2	2	y2
3	x3	4	y3

Mutation Joins: add new variables to one data frame by matching observations in another.



The key(s) is used to decide which observations (rows) go where.

What do you think would happen if:

- you had a typo in a key?
- a key matched multiple observations?

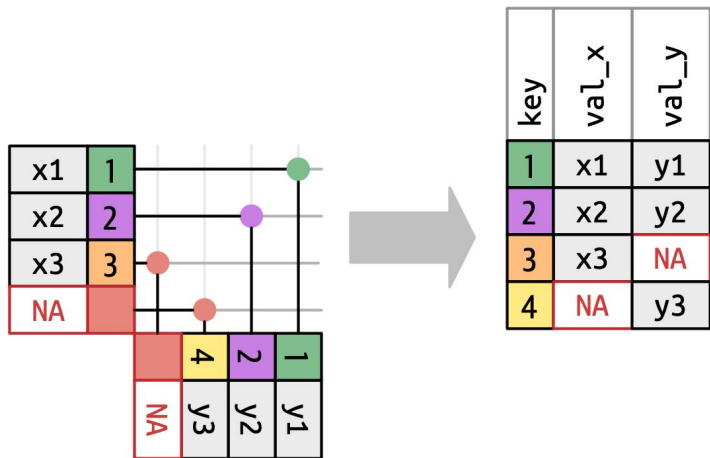
Joins: Because data is never in just one place



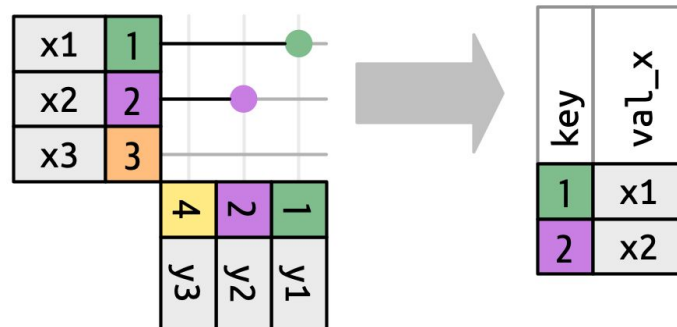
Combining data sets

x		y	
1	x1	1	y1
2	x2	2	y2
3	x3	4	y3

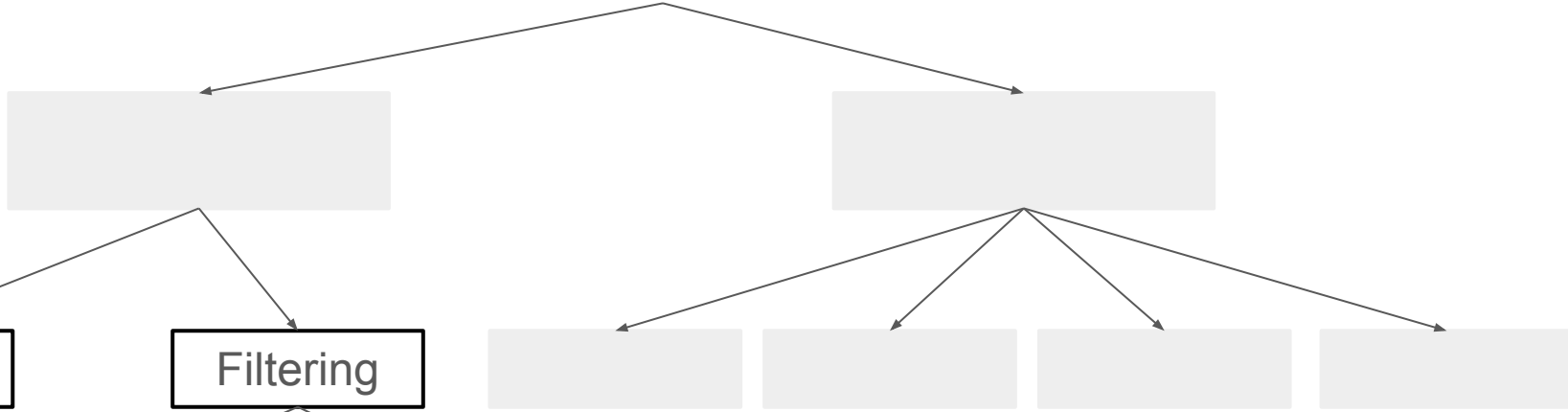
Mutation Joins: add new variables to one data frame by matching observations in another.



Filtering Joins: filter observations from one data frame based on whether or not they match an observation in another.



Types of Joins



Can filter on those that match or those that don't match

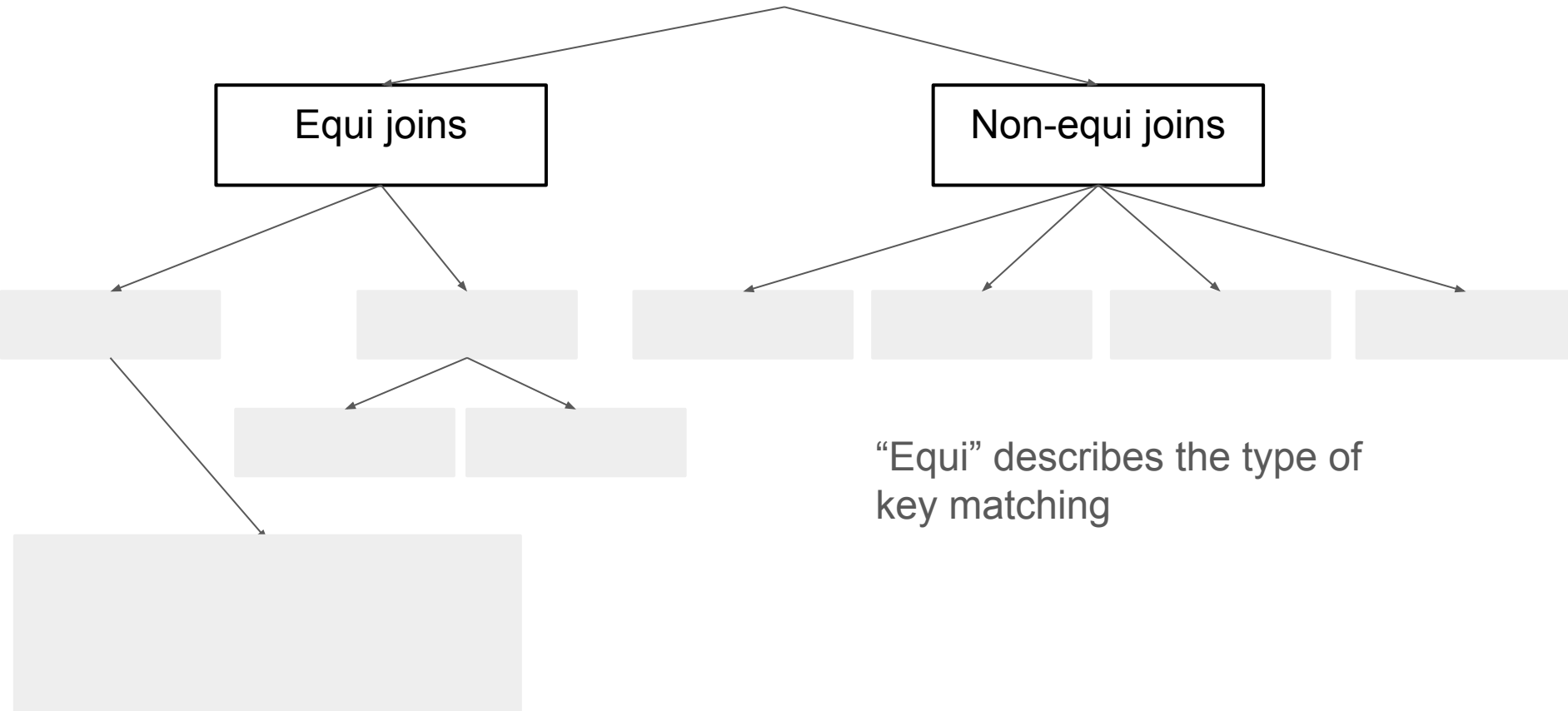
Super useful when mutating join is giving unexpected results

Types of Joins

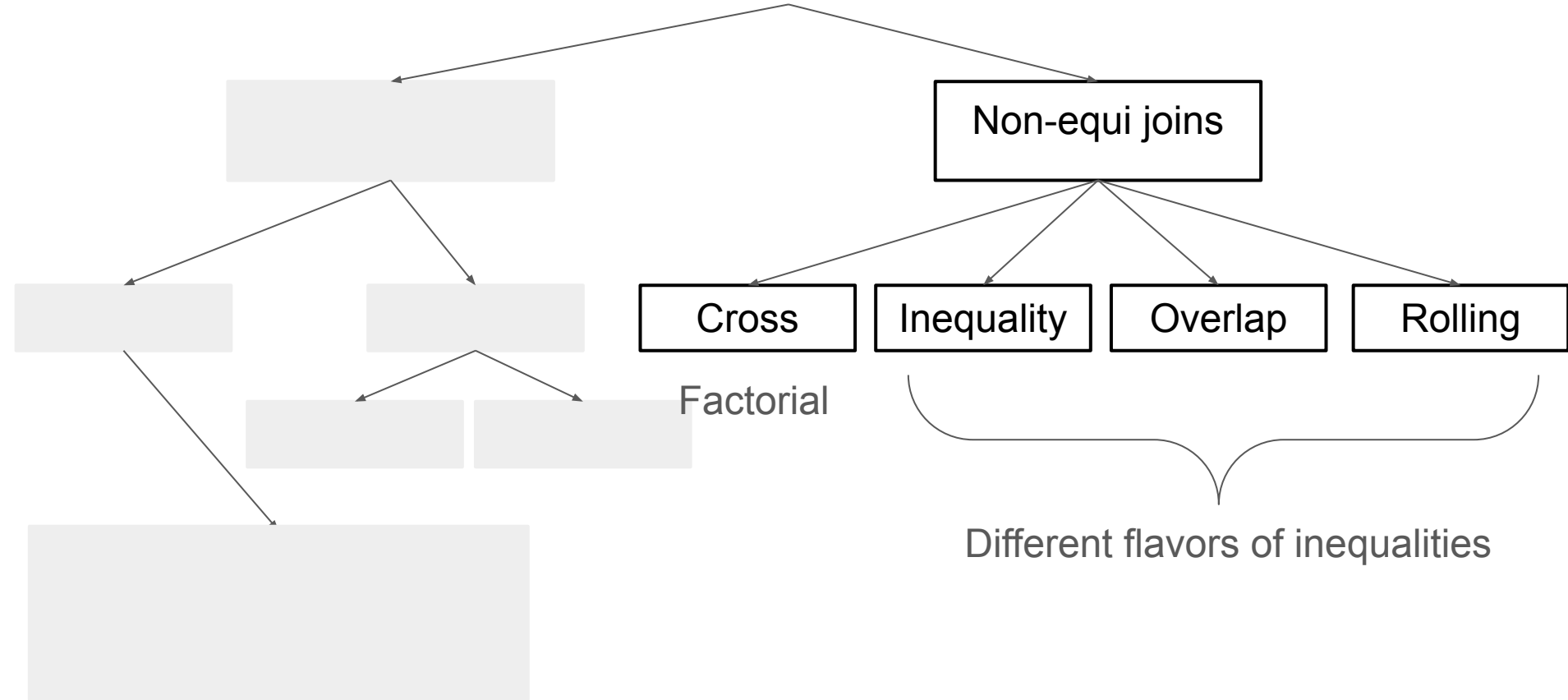
Equi joins

Non-equi joins

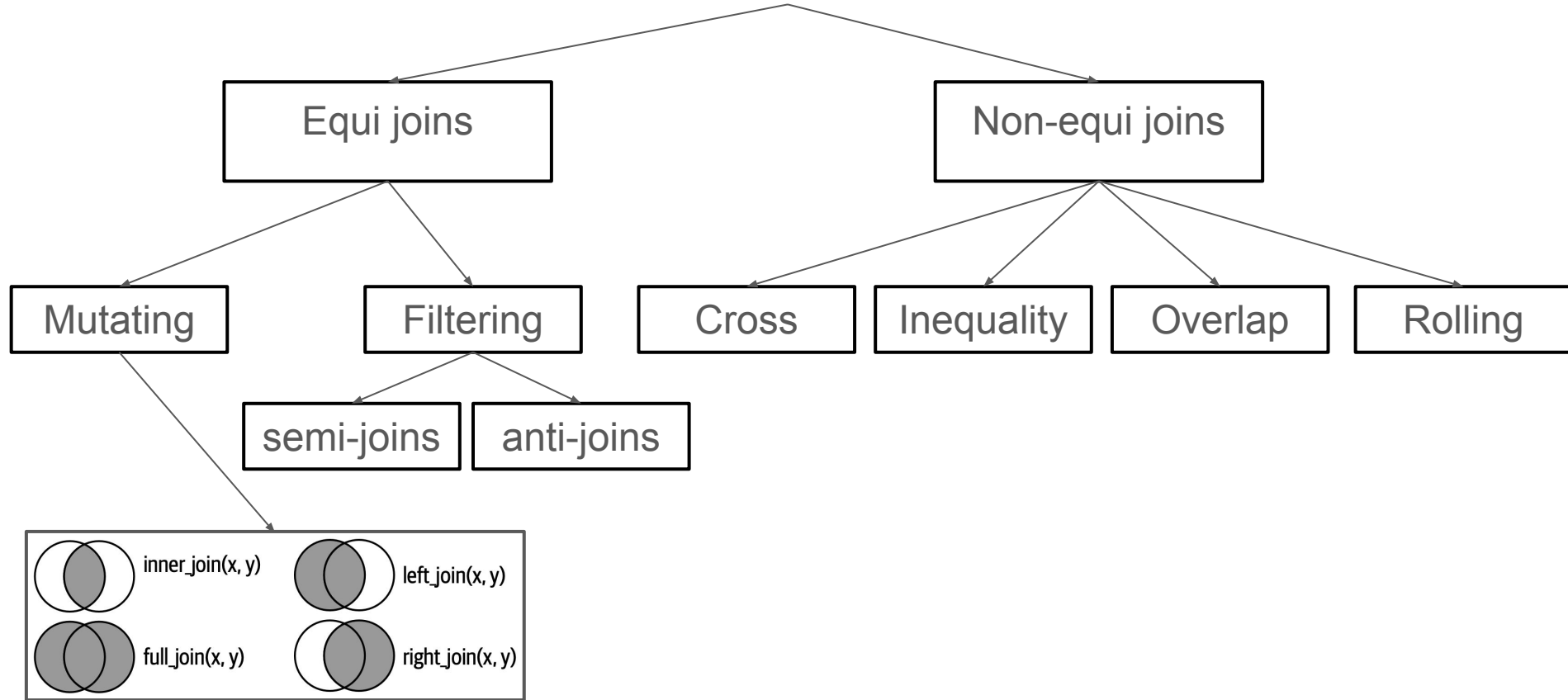
“Equi” describes the type of
key matching



Types of Joins



Types of Joins



Let's Practice

X

Patient ID	Diag. Test
12	Y
25	N
72	Y

Y

Patient ID	Test Result
12	Positive
72	Negative
83	Negative

What is the key?

What is the output of an

`inner_join(x,y)`

`anti_join(x,y)`

`anti_join(y,x)`

`full_join(x,y)`

Questions and switch to M6.R