

Tidy Data

Fall 2022

September 28 2022

Last time: Combining data sets

COMBINE CASES

The diagram illustrates the combination of two data frames, X and y, using a plus sign (+) as the operator. Data frame X is represented by a gray grid:

A	B	C
a	t	1
b	u	2
c	v	3

Data frame y is represented by an orange grid:

A	B	C
C	v	3
d	w	4

The resulting combined data frame is shown below the addition operation.

COMBINE VARIABLES

The diagram illustrates the combination of two data frames, X and y, into a wider frame using a plus sign (+) and an equals sign (=) as the operators. Data frame X is represented by a gray grid:

A	B	C
a	t	1
b	u	2
c	v	3

Data frame y is represented by an orange grid:

A	B	D
a	t	3
b	u	2
d	w	1

The resulting combined data frame is shown below the addition operation.

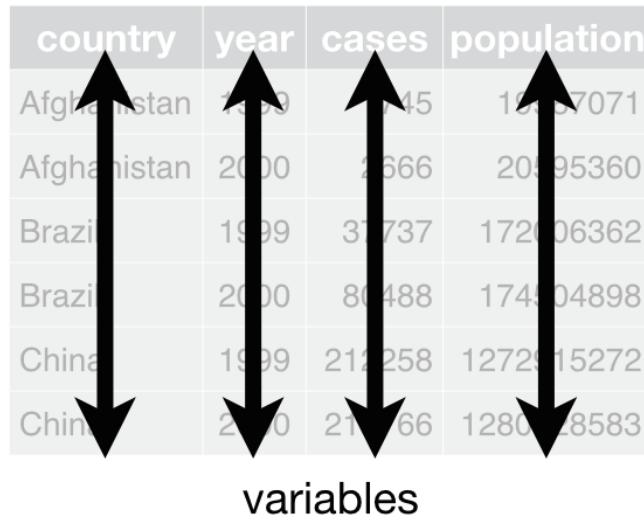
dplyr data transformation cheatsheet

What are tidy data?

1. Each **variable** forms a column
2. Each **observation** forms a row
3. Each **value** has its own cell

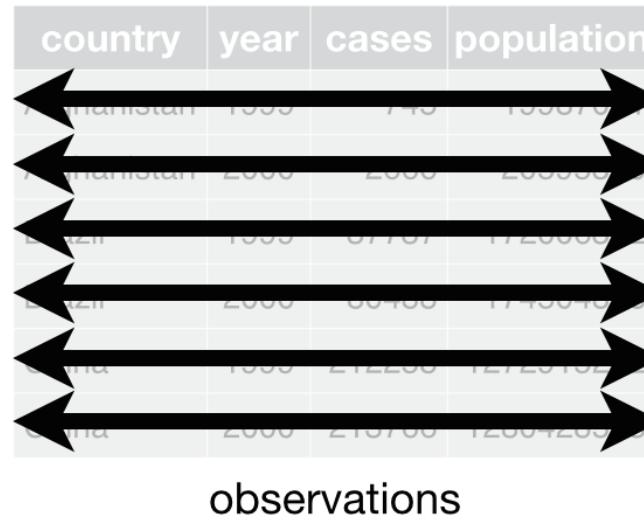
country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

variables



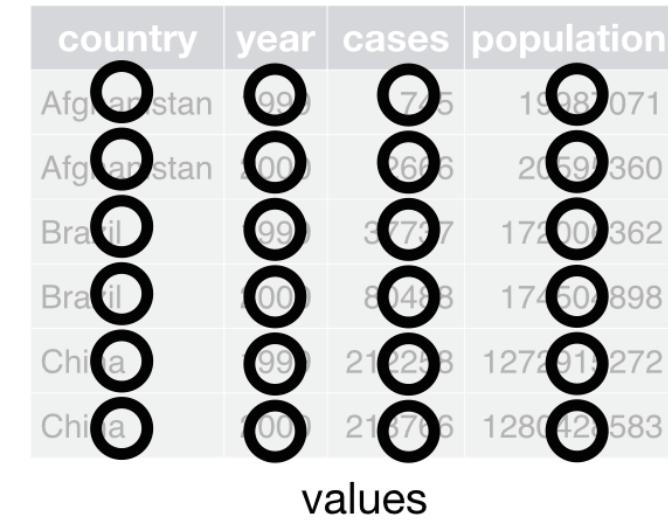
country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

observations



country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

values



Untidy data: example 1

```
untidy_data <- tibble(  
  name = c("Ana", "Bob", "Cara"),  
  wt_07_01_2021 = c(100, 150, 140),  
  wt_08_01_2021 = c(104, 155, 138),  
  wt_09_01_2021 = c(NA, 160, 142)  
)  
untidy_data  
# A tibble: 3 × 4  
  name  wt_07_01_2021 wt_08_01_2021 wt_09_01_2021  
  <chr>     <dbl>        <dbl>        <dbl>  
1 Ana        100         104         NA  
2 Bob        150         155         160  
3 Cara       140         138         142
```

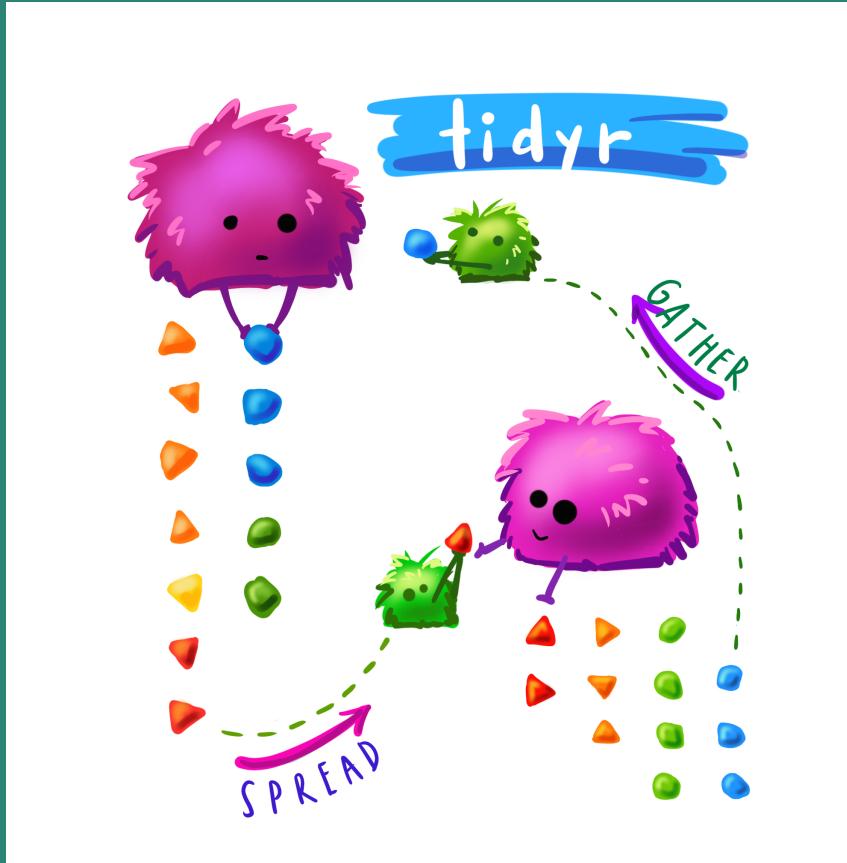
Tidy data: example 1

```
library(lubridate)
library(stringr)
untidy_data %>%
  pivot_longer(names_to = "date",
              values_to = "weight",
              cols = -name) %>%
  mutate(date = str_remove(date, "wt_"),
         date = dmy(date))      # dmy() is a function in the lubridate package
# A tibble: 9 × 3
  name   date       weight
  <chr> <date>     <dbl>
1 Ana    2021-01-07  100
2 Ana    2021-01-08  104
3 Ana    2021-01-09  NA
4 Bob    2021-01-07  150
5 Bob    2021-01-08  155
6 Bob    2021-01-09  160
7 Cara   2021-01-07  140
8 Cara   2021-01-08  138
9 Cara   2021-01-09  142
```

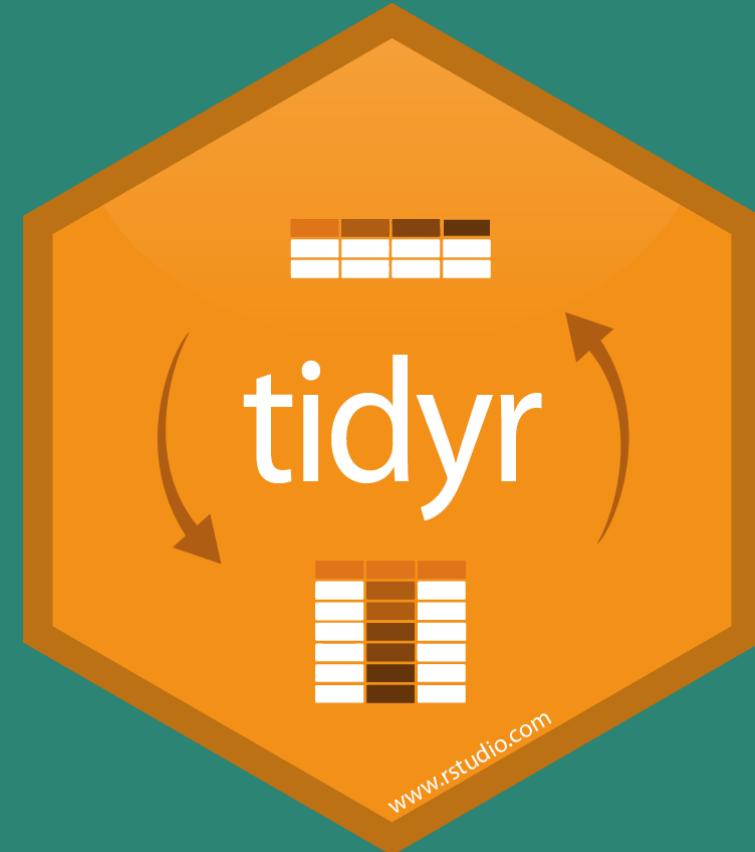
We will learn how to do this!

Reshaping data

wide vs. long



Allison Horst



tidyverse

Wide vs. long data

- Wide data has one row per subject, with multiple columns for their repeated measurements
- Long data has multiple rows per subject, with one column for the measurement variable and another indicating from when/where the repeated measures are from

id	SBP_visit1	SBP_visit2	SBP_visit3
a	130	110	112
b	120	116	122
c	130	136	138
d	119	106	118

wide

id	visit	SBP
a	1	130
b	1	120
c	1	130
d	1	119
a	2	110
b	2	116
c	2	136
d	2	106
a	3	112
b	3	122
c	3	138
d	3	118

long

Wide to long: pivot_longer()

```
BP_wide
# A tibble: 4 × 5
  id   sex   SBP_v1 SBP_v2 SBP_v3
  <chr><chr>  <dbl>  <dbl>  <dbl>
1 a     F        130    110    112
2 b     M        120    116    122
3 c     M        130    136    138
4 d     F        119    106    118
```

```
BP_long <- BP_wide %>%
  pivot_longer(names_to = "visit",
               values_to = "SBP",
               cols = SBP_v1:SBP_v3)

BP_long
# A tibble: 12 × 4
  id   sex   visit   SBP
  <chr><chr> <chr>  <dbl>
1 a     F     SBP_v1  130
2 a     F     SBP_v2  110
3 a     F     SBP_v3  112
4 b     M     SBP_v1  120
5 b     M     SBP_v2  116
6 b     M     SBP_v3  122
7 c     M     SBP_v1  130
8 c     M     SBP_v2  136
9 c     M     SBP_v3  138
10 d    F     SBP_v1  119
11 d    F     SBP_v2  106
12 d    F     SBP_v3  118
```

Wide to long: `pivot_longer()`

`pivot_longer` lengthens data, increasing the number of rows and decreasing the number of columns.

Need to **specify**:

- *new column names*
 - `names_to`: stores row names of wide data's columns
 - `values_to`: stores data values
- *which columns to pivot*

Long to wide: pivot_wider()

```
BP_long
# A tibble: 12 × 4
  id   sex   visit     SBP
  <chr> <chr> <chr>    <dbl>
1 a     F     SBP_v1    130
2 a     F     SBP_v2    110
3 a     F     SBP_v3    112
4 b     M     SBP_v1    120
5 b     M     SBP_v2    116
6 b     M     SBP_v3    122
7 c     M     SBP_v1    130
8 c     M     SBP_v2    136
9 c     M     SBP_v3    138
10 d    F     SBP_v1    119
11 d    F     SBP_v2    106
12 d    F     SBP_v3    118
```

```
BP_wide1 <- BP_long %>%
  pivot_wider(names_from = "visit",
              values_from = "SBP")
BP_wide1
# A tibble: 4 × 5
  id   sex   SBP_v1  SBP_v2  SBP_v3
  <chr> <chr>    <dbl>    <dbl>    <dbl>
1 a     F        130     110     112
2 b     M        120     116     122
3 c     M        130     136     138
4 d     F        119     106     118
```

Long to wide: `pivot_wider()`

`pivot_wider` increases number of columns and decreases the number of rows. Need to specify:

- *new column names*
 - `names_from`: get the name of the column
 - `values_from`: get the cell values from

Separate Info

```
BP_long
# A tibble: 12 × 4
  id   sex  visit    SBP
  <chr><chr><chr>  <dbl>
1 a     F    SBP_v1  130
2 a     F    SBP_v2  110
3 a     F    SBP_v3  112
4 b     M    SBP_v1  120
5 b     M    SBP_v2  116
6 b     M    SBP_v3  122
7 c     M    SBP_v1  130
8 c     M    SBP_v2  136
9 c     M    SBP_v3  138
10 d    F    SBP_v1  119
11 d    F    SBP_v2  106
12 d    F    SBP_v3  118
```

```
BP_long1 <- BP_long %>%
  separate(visit, c("acrnym", "visit"))
BP_long1
# A tibble: 12 × 5
  id   sex  acrnym visit    SBP
  <chr><chr><chr>  <chr>  <dbl>
1 a     F    SBP    v1      130
2 a     F    SBP    v2      110
3 a     F    SBP    v3      112
4 b     M    SBP    v1      120
5 b     M    SBP    v2      116
6 b     M    SBP    v3      122
7 c     M    SBP    v1      130
8 c     M    SBP    v2      136
9 c     M    SBP    v3      138
10 d    F    SBP    v1      119
11 d    F    SBP    v2      106
12 d    F    SBP    v3      118
```

Remove and Clean up Columns

```
BP_long1
# A tibble: 12 × 5
  id   sex acrnym visit    SBP
  <chr> <chr> <chr> <chr> <dbl>
1 a     F     SBP   v1     130
2 a     F     SBP   v2     110
3 a     F     SBP   v3     112
4 b     M     SBP   v1     120
5 b     M     SBP   v2     116
6 b     M     SBP   v3     122
7 c     M     SBP   v1     130
8 c     M     SBP   v2     136
9 c     M     SBP   v3     138
10 d    F     SBP   v1     119
11 d    F     SBP   v2     106
12 d    F     SBP   v3     118
```

Goal: Get rid of `acrnym` column and remove the string `v` from the `visit` variable's values.

```
library(stringr)
BP_long2 <- BP_long1 %>%
  select(-acrnym) %>%
  mutate(visit = str_replace(visit, "v", ""))
  mutate_at(vars(visit), as.double)

BP_long2
# A tibble: 12 × 4
  id   sex   visit    SBP
  <chr> <chr> <dbl> <dbl>
1 a     F       1     130
2 a     F       2     110
3 a     F       3     112
4 b     M       1     120
5 b     M       2     116
6 b     M       3     122
7 c     M       1     130
8 c     M       2     136
9 c     M       3     138
10 d    F       1     119
11 d    F       2     106
12 d    F       3     118
```

Make cleaned-up long data wide

```
head(BP_long2, 4)
# A tibble: 4 × 4
  id    sex    visit    SBP
  <chr> <chr> <dbl>   <dbl>
1 a      F        1     130
2 a      F        2     110
3 a      F        3     112
4 b      M        1     120
```

```
BP_wide2 <- BP_long2 %>%
  pivot_wider(names_from = "visit",
              values_from = "SBP")
BP_wide2
# A tibble: 4 × 5
  id    sex    `1`    `2`    `3`
  <chr> <chr> <dbl>   <dbl>   <dbl>
1 a      F      130    110    112
2 b      M      120    116    122
3 c      M      130    136    138
4 d      F      119    106    118
```

Problem: have numbers as column names

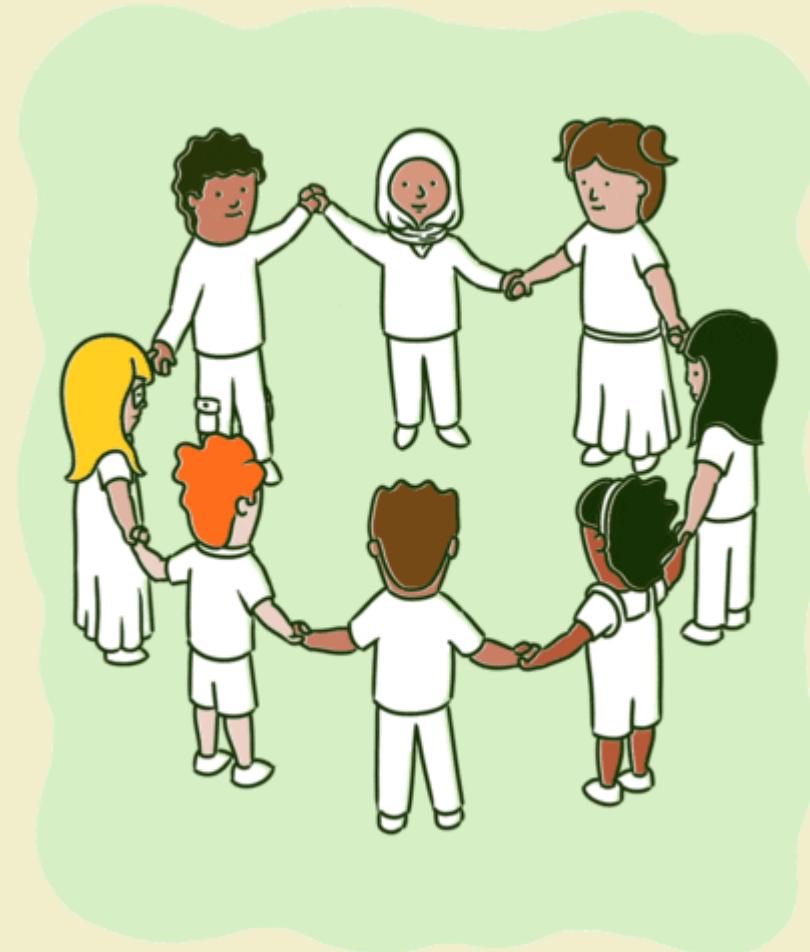
Solution: have row names start with the key column's name separated by a character

```
BP_wide3 <- BP_long2 %>%
  pivot_wider(names_from = "visit",
              values_from = "SBP",
              names_prefix = "value_")
```

```
BP_wide3
# A tibble: 4 × 5
  id    sex    value_1  value_2  value_3
  <chr> <chr>   <dbl>    <dbl>    <dbl>
1 a      F       130     110     112
2 b      M       120     116     122
3 c      M       130     136     138
4 d      F       119     106     118
```

Group Activity 1

10:00



- Let's go over to maize server/ local Rstudio and our class [moodle](#)
- Get the class activity 8.Rmd file
- Work on problem 1
- Ask me questions

janitor::clean_names

convert all column names to `*-case!`

(not so awesome
column names)

`fun(units)`
`sampled`
`suchness`
`000.00F`

MESS



1M,00P

clean_names()
(your data frame +
case choice here)

CHOOSE:
✓ snake
lower_camel
big_camel
screaming_snake
+ more!

WAY MORE
DEALWITHABLE
COLUMN
NAMES!

hope-restr
-wc
-than
yay-yay
cool
good-feel
so-much-nicer
Teddy
snake

HORST 2019

Allison Horst

Clean messy column names with janitor::clean_names()

```
bar <- tibble("First Name"= c("Yi","DJ"), "last init" = c("C","R"), "% in" = c(0.1, 0.5),
             "~~~~~"= 1:2, " " " =3:2," hi"=c("a","b"), "null"=c(NA,NA))
bar
# A tibble: 2 × 7
`First Name` `last init` `% in` `~~~~~` `` `` ` hi` null
<chr>         <chr>       <dbl>   <int> <int> <chr> <lgl>
1 Yi            C           0.1      1     3 a     NA
2 DJ            R           0.5      2     2 b     NA
```

```
library(janitor)
bar %>% clean_names() %>%          # in the janitor package
  remove_empty(c("rows","cols"))      # also useful
# A tibble: 2 × 6
  first_name last_init percent_in      n      x hi
  <chr>       <chr>       <dbl> <int> <int> <chr>
1 Yi          C           0.1      1     3 a
2 DJ          R           0.5      2     2 b
```

Removing missing data: `drop_na()`

A simple example:

```
foo <- tibble(id = 3:5,
               name = c("Al",
                        "Wu",
                        "Flo"),
               height = c(3, NA, 2.8),
               years = c(50, 33, NA))

foo
# A tibble: 3 × 4
  id   name  height years
  <int> <chr>  <dbl> <dbl>
1     3 Al        3     50
2     4 Wu       NA     33
3     5 Flo      2.8    NA
```

Remove *all* rows with any missing data

```
foo %>% drop_na()
# A tibble: 1 × 4
  id   name  height years
  <int> <chr>  <dbl> <dbl>
1     3 Al        3     50
```

Remove rows with `NA` in selected columns

```
foo %>% drop_na(height)
# A tibble: 2 × 4
  id   name  height years
  <int> <chr>  <dbl> <dbl>
1     3 Al        3     50
2     5 Flo      2.8    NA
```

Replace NAs with another value: `replace_na()`

Use with `mutate()`

```
foo  
# A tibble: 3 × 4  
  id name  height years  
  <int> <chr>   <dbl> <dbl>  
1    3 Al        3     50  
2    4 Wu       NA     33  
3    5 Flo      2.8    NA
```

```
foo %>% mutate(height = replace_na(height , 0),  
                  years = replace_na(years, 0))  
# A tibble: 3 × 4  
  id name  height years  
  <int> <chr>   <dbl> <dbl>  
1    3 Al        3     50  
2    4 Wu        0     33  
3    5 Flo      2.8     0
```

Date manipulation

	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

lubridate

Dates with `lubridate`

- Convert characters to special "Date" type
- Easy date magic examples:
 - add and subtract dates
 - convert to minutes/years/etc
 - change timezones
 - add 1 month to a date...



Allison Horst

What kind of date do you have?

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 wide variety of input formats.

2017-11-28T14:02:00	ymd_hms() , ymd_hm() , ymd_h() . ydm_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")
2:01	hms::hms() Also lubridate::hms() , hm() and ms() , which return periods.* hms::hms(sec = 0, min = 1, hours = 2)

```
timedata <- tibble(name = c("Yi", "Mo", "Dee"),  
                    dob=c("10/31/1952",  
                          "1/12/1984",  
                          "2/02/2002"),  
                    age= c("69", "38", "19"))
```

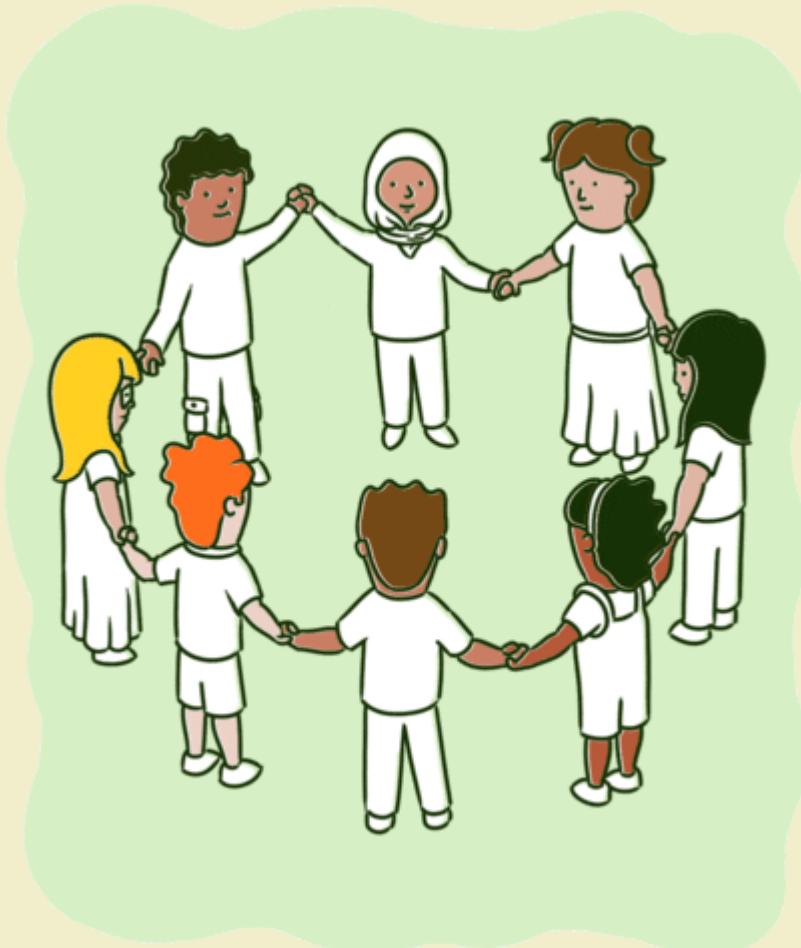
```
timedata %>%  
  mutate(dob_date = mdy(dob),  
         age_num = as.numeric(age))  
# A tibble: 3 × 5  
  name    dob        age   dob_date  age_num  
  <chr>   <chr>     <dbl>   <date>      <dbl>  
1 Yi      10/31/1952 69     1952-10-31  69  
2 Mo      1/12/1984  38     1984-01-12  38  
3 Dee     2/02/2002  19     2002-02-02  19
```

Math with dates

```
timedata %>%
  mutate(dob = mdy(dob),
         dob_year = year(dob),                                # convert to a date
         time_since_birth = dob %--% today(),                # extract the year
         age = time_since_birth %/% years(1),                 # create an "interval"
         dobplus = dob + days(30)                            # modulus on "years"
  )
# A tibble: 3 × 6
  name    dob        age  dob_year time_since_birth      dobplus
  <chr>   <date>    <dbl>   <dbl>   <Interval>       <date>
1 Yi     1952-10-31    69    1952  1952-10-31 UTC--2022-09-28 UTC 1952-11-30
2 Mo     1984-01-12    38    1984  1984-01-12 UTC--2022-09-28 UTC 1984-02-11
3 Dee    2002-02-02    20    2002  2002-02-02 UTC--2022-09-28 UTC 2002-03-04
```

Group Activity 2

10:00



- Work on problems 2
- Ask me questions
- Any hw-related questions?

Acknowledgements

Parts of the slides are adapted from Jessica Minnier's OCTRI BERD R Courses instruction materials.