

Lecture 13: Data wrangling

Reminder:

- Exam 1 on Oct. 9 (next Monday)
- Review questions on course website
- Review day on Friday in class

So far

- `select`: choose certain columns
- `filter`: choose certain rows
- `summarize`: calculate summary statistics
- `group_by`: group rows together
- `mutate`: create new columns
- `count`: count the number of rows
- `arrange`: re-order the rows

Do dogs help exam stress?

- Data collected on 284 students at a mid-size Canadian university
- Students randomly assigned to one of three treatment groups: handler-only contact, indirect contact, and direct contact
- Well-being and ill-being measures recorded before and after treatment for each student
- Approach: compare pre/post measures of well-being and ill-being

Recording well-being and ill-being measures

- Likert items for each well-being / ill-being measure
- Average the likert items to get a score for each measure
- E.g.:
 - Positive affect score is the average of 5 Likert items
 - Social connectedness is the average of 20 Likert items

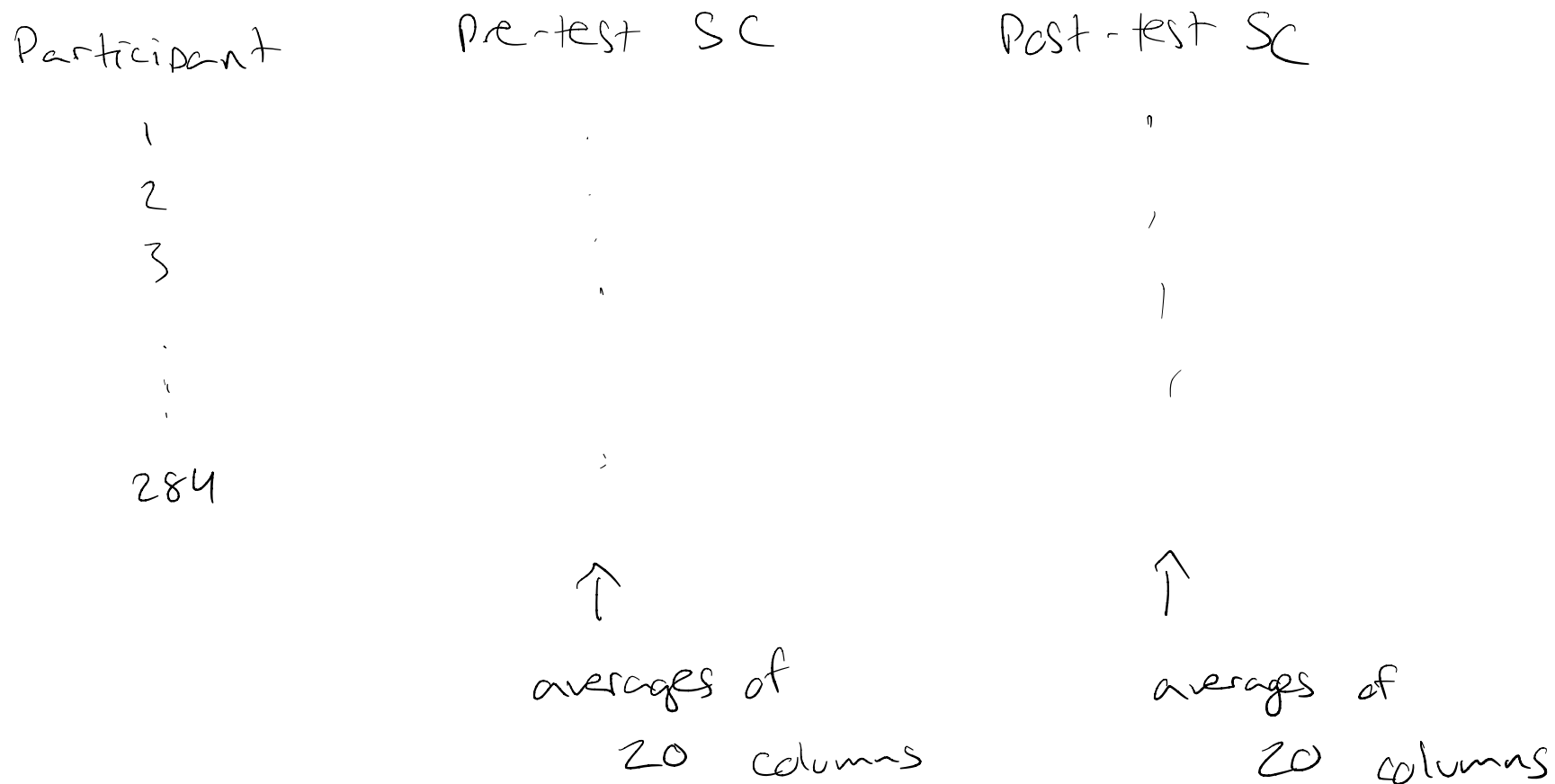
Example Likert item for social connectedness

“I am able to relate to my peers.”

- Strongly disagree (1)
- Disagree
- Somewhat disagree
- Somewhat agree
- Agree
- Strongly agree (6)

Our goal for today

- Calculate the pre- and post-treatment social connectedness scores for each participant
- **Question:** What do we want the final data to look like?



Initial data processing

- Social connectedness is the average of 20 Likert items
- These items should take values between 1 and 6
- However:

```
1 raw_data |>  
2   select(starts_with("SC")) |>  
3   max(na.rm=T)
```

[1] 66

↑ mistake!

Let's code these as missing (NA)

Handling errors

```
1 cleaned_data <- raw_data |>
2   mutate(SC1_1 = ifelse(SC1_1 > 6, NA, SC1_1),
3          SC1_2 = ifelse(SC1_2 > 6, NA, SC1_2),
4          ...)
```

Are there any issues with this approach?

- tedious
- prone to errors

across

```
1 example_df
```

```
  x1 x2 x3 y1 y2
1  5  7  8  5 NA
2  2  4  4  2  4
3  4  8  7  5  5
```

```
1 example_df |>
```

```
2   summarize(across(c(x1, x2, x3, y1, y2), mean))
```

```
  x1      x2      x3 y1 y2
1 3.666667 6.333333 6.333333 4 NA
```

need a function that
calculates the mean,
ignoring missing values, and
only takes one argument

← calculate the mean

for each column
specified

Question: What if I want to ignore NAs when computing the mean?

$na.rm = T$

(x_1, x_2, \dots)

across

```
1 example_df
```

```
  x1 x2 x3 y1 y2
1  5  7  8  5 NA
2  2  4  4  2  4
3  4  8  7  5  5
```

```
1 example_df |>
2   summarize(across(c(x1, x2, x3, y1, y2),
3                   function(x) {mean(x, na.rm=T)} ))
```

```
  x1      x2      x3 y1  y2
1 3.666667 6.333333 6.333333 4 4.5
```

anonymous function

computes mean, ignores NAs

across

```
1 example_df
```

```
  x1 x2 x3 y1 y2
1  5  7  8  5 NA
2  2  4  4  2  4
3  4  8  7  5  5
```

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {x + 1} ))
```

```
  x1 x2 x3 y1 y2
1  6  8  9  6 NA
2  3  5  5  3  5
3  5  9  8  6  6
```

change this to replace $x > 6$ with NA

Question: What if I want to replace values > 6 with NA?

across

```
1 example_df
```

	x1	x2	x3	y1	y2
1	5	7	8	5	NA
2	2	4	4	2	4
3	4	8	7	5	5

```
1 example_df |>  
2   mutate(across(c(x1, x2, x3, y1, y2),  
3             function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

Handling errors

still manually specifying each
column (tedious!)

```
1 cleaned_data <- raw_data |>  
2   mutate(across(c(SC1_1, SC1_2, ... ),  
3                 function(x) {ifelse(x > 6, NA, x)}))
```

Question: Are there any issues with this approach?

Handling errors

all of the SC columns

```
1 cleaned_data <- raw_data |>
2   mutate(across(starts_with("SC"),
3                 function(x) {ifelse(x > 6, NA, x)}))
4
5 cleaned_data |>
6   select(starts_with("SC")) |>
7   max(na.rm=T)
```

↑ replace values > 6
with NA

[1] 6

More data cleaning

- For some Social Connectedness items, “6” means “more connected”
 - e.g.: “I find myself actively involved in people’s lives.”
- For some Social Connectedness items, “6” means “less connected”
 - e.g.: “I feel like an outsider.” $\leftarrow \text{new response} = 7 - \text{response}$
- We want higher scores to always mean “more connected”

We need to reverse the scores for some Social Connectedness items!

More data cleaning

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

Suppose we want to reverse the scores for x1 and x3

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} )) |>
4   select(num_range("x", c(1, 3)))
```

	x1	x3
1	5	NA
2	2	4
3	4	NA

give me columns start with "x" and
have 1 or 3
=> x1 x3

More data cleaning

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

Suppose we want to reverse the scores for x1 and x3

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ),
4         across(num_range("x", c(1, 3)),
5                 function(x) {7 - x}))
```

	x1	x2	x3	y1	y2
1	2	NA	NA	5	NA
2	5	4	3	2	4
3	3	NA	NA	5	5

$$2 = 7 - 5$$

$$3 = 7 - 4$$

With the dog data

```
1 cleaned_data <- raw_data |>
2   mutate(across(starts_with("SC"),
3                 function(x) {ifelse(x > 6, NA, x)}),
4           across(num_range("SC1_",
5                             c(3, 6, 7, 9, 11, 13, 15, 17, 18, 20)),
6                 function(x) {7 - x}),
7           across(num_range("SC2_",
8                             c(3, 6, 7, 9, 11, 13, 15, 17, 18, 20)),
9                 function(x) {7 - x}))
```

"SC", c(3, 6, 7, ...)

SC3 SC6

Averaging columns

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)}),
4           across(num_range("x", c(1, 3)),
5                 function(x) {7 - x}))
```

	x1	x2	x3	y1	y2
1	2	NA	NA	5	NA
2	5	4	3	2	4
3	3	NA	NA	5	5

Question: What if I want to calculate the average of the x columns for each row?

Averaging columns

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ),
4           x_mean = (x1 + x2 + x3)/3)
```

	x1	x2	x3	y1	y2	x_mean
1	5	NA	NA	5	NA	NA
2	2	4	4	2	4	3.333333
3	4	NA	NA	5	5	NA

← new column to contain row averages of x1, x2, x3

Averaging columns

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ),
4           x_mean = mean(c(x1, x2, x3), na.rm=T))
```

	x1	x2	x3	y1	y2	x_mean
1	5	NA	NA	5	NA	3.8
2	2	4	4	2	4	3.8
3	4	NA	NA	5	5	3.8

calculate the mean across all rows in x1, x2, x3
want: 5
3.33
4

(not being done row wise)

Averaging columns

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} ))
```

	x1	x2	x3	y1	y2
1	5	NA	NA	5	NA
2	2	4	4	2	4
3	4	NA	NA	5	5

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3                 function(x) {ifelse(x > 6, NA, x)} )) |>
4   rowwise() |> ← calculate for each row
5   mutate(x_mean = mean(c(x1, x2, x3), na.rm=T))
```

A tibble: 3 × 6

Rowwise:

	x1	x2	x3	y1	y2	x_mean
	<int>	<int>	<int>	<int>	<int>	<dbl>
1	5	NA	NA	5	NA	5
2	2	4	4	2	4	3.33
3	4	NA	NA	5	5	4

Averaging columns

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3               function(x) {ifelse(x > 6, NA, x)} ))
```

```
  x1 x2 x3 y1 y2
1   5 NA NA  5 NA
2   2  4  4  2  4
3   4 NA NA  5  5
```

```
1 example_df |>
2   mutate(across(c(x1, x2, x3, y1, y2),
3               function(x) {ifelse(x > 6, NA, x)} )) |>
4   rowwise() |>
5   mutate(x_mean = mean(c_across(starts_with("x")), na.rm=T))
```

```
# A tibble: 3 × 6
```

```
# Rowwise:
```

```
  x1    x2    x3    y1    y2  x_mean
<int> <int> <int> <int> <int> <dbl>
1     5    NA    NA     5    NA     5
2     2     4     4     2     4    3.33
3     4    NA    NA     5     5     4
```

↑ across columns, within each row

With the dog data

```
1 cleaned_data <- raw_data |>
2   mutate(across(starts_with("SC"),
3     function(x) {ifelse(x > 6, NA, x)}),
4     across(num_range("SC1_",
5       c(3, 6, 7, 9, 11, 13, 15, 17, 18, 20)),
6     function(x) {7 - x}),
7     across(num_range("SC2_",
8       c(3, 6, 7, 9, 11, 13, 15, 17, 18, 20)),
9     function(x) {7 - x})) |>
10 rowwise() |>
11 mutate(sc_pre = mean(c_across(starts_with("SC1_")), na.rm=T),
12       sc_post = mean(c_across(starts_with("SC2_")), na.rm=T))
```

```
1 cleaned_data |>
2   select(sc_pre, sc_post)
```

A tibble: 284 × 2

Rowwise:

	sc_pre	sc_post
	<dbl>	<dbl>
1	3.9	3.8
2	5.15	5.26
3	4.1	4.15
4	4.65	5.1
5	3.65	3.6
6	4.35	4.65

7	4.75	4.4
8	4.6	4.65
9	4.2	4.15
10	5.8	5.75

Class activity

<https://sta279->

[f23.github.io/class_activities/ca_lecture_13.html](https://sta279-f23.github.io/class_activities/ca_lecture_13.html)

