# Sociology Quant Camp

## Introduction to R
## Module 2: piping and tidyverse

**Monica Alexander, Statistical Sciences and Sociology**

# Using the tidyverse to manipulate real data

- In the previous module, we saw some functions and loaded in the `tidyverse` package

- Tidyverse has a range of functions that make it easier to manipulate real data

- Things like: adding columns, selecting columns, filtering out rows based on certain values…

- These functions have been specifically designed to work with datasets with lots of variables of different types

# A first example

- Let's read in the shelter data and select some columns

- Note that `colnames()` is a useful function to see what the columns are called

# Demo: selecting columns

# The pipe |>

- An alternative way of writing code

- Makes the code read more like a sentence

- Read the pipe as "and then"

- So here we are taking the data AND THEN selecting columns

# Core tidyverse functions

- `select`: select columns

- `arrange`: sort/arrange by value

- `mutate`: make a new column

- `filter`: filter out certain rows

- `summarize`: produce summaries of data

- `group_by`: group the data by certain variable(s)

```r
1   library(tidyverse)
2
3   # read in data
4   d <- read_csv("shelter.csv")
5
6   # select
7   d |>
8     select(occupancy_date, organization_name,
9            occupancy_rate_beds, occupancy_rate_rooms)
10
11  # assign this to an object
12
13  dr <- d |>
14    select(occupancy_date, organization_name,
15           occupancy_rate_beds, occupancy_rate_rooms)
16
17  # arrange
18  dr |>
19    arrange(occupancy_rate_beds)
20
21  # mutate
22
23  dr |>
24    mutate(less_than_50pc_beds = occupancy_rate_beds<50)
25
26  # filter
27
28  dr |>
29    filter(!is.na(occupancy_rate_beds))
30
31  # summarize
32
33  dr |>
34    summarize(min_rate = min(occupancy_rate_beds, na.rm = TRUE))
```

# Demo: tidyverse functions

# group_by

- The group_by function is extremely powerful when used in conjunction with summarize to get summaries by groups

- Note that we can thread together multiple pipes!

```
dr |>
  group_by(occupancy_date) |>
  summarize(min_rate = min(occupancy_rate_beds, na.rm = TRUE))
```

Here is the output:

```
> dr |>
+   group_by(occupancy_date) |>
+   summarize(min_rate = min(occupancy_rate_beds, na.rm = TRUE))
# A tibble: 246 × 2
   occupancy_date min_rate
   <date>            <dbl>
 1 2024-01-01         43.8
 2 2024-01-02         50
 3 2024-01-03         50
 4 2024-01-04         50
 5 2024-01-05         50
 6 2024-01-06         63.6
 7 2024-01-07         63.6
 8 2024-01-08         66.7
 9 2024-01-09         52.5
10 2024-01-10         50
# i 236 more rows
# i Use `print(n = ...)` to see more rows
```

# Demo: more complicated tidyverse functions

# Where to get help

- Lots of good, free online sources

    - R for Data Science: https://www.tidyverse.org/learn/

    - Telling stories with data: https://tellingstorieswithdata.com/

    - Tidyverse skills for data science: https://jhudatascience.org/tidyversecourse/intro.html

- Google/Stack Overflow

- Email

- Practice, practice, practice; don't be afraid of mistakes