# 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 Covid-19 deaths 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  library(lubridate) # to deal with dates, you will need to install
3
4  # read in data
5  d <- read_csv("deaths_fatality_type.csv")
6
7  # select columns
8  d |>
9    select(date, death_covid)
10
11 # arrange by deaths in descending order
12 d |>
13   arrange(-death_covid)
14
15 # make a new column which is true if reported deaths are negative
16 d |>
17   mutate(deaths_negative = deaths_total<0)
18
19 # filter out negtaive deaths
20 d |>
21   filter(deaths_total>0)
22
23 # summarize the total number of deaths over all days
24 d |>
25   summarize(total_covid_deaths = sum(death_covid))
26
```

# 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!

```
d_with_year <- d |>
  mutate(year = year(date))

d_with_year |>
  group_by(year) |>
  summarize(total_deaths = sum(death_covid))
```

Assigning the dataset with new year column to a new dataset

Here is the output:

```
> d_with_year <- d |>
+   mutate(year = year(date))
> d_with_year |>
+   group_by(year) |>
+   summarize(total_deaths = sum(death_covid))
# A tibble: 3 × 2
   year total_deaths
  <dbl>        <dbl>
1  2020         1193
2  2021         5617
3  2022         2887
```

Total deaths by year!

# 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