#### Outline

- Workflow assignment
- Tidy data + tutorial
- Go over wrangling homework
- Automation lecture + tutorial
- Break-out groups

# Workflow assignment

# Pick a project and critique the current workflow

- Something you're working on or something similar to what you want/plan to work on
- Ideally, contains raw data that needs to be cleaned
- If you're not sure about the project you've picked, discuss it with me

### Two parts

- Written report (2-4 pages, 2x spaced)
  - Turn in on canvas by 1:30pm on 2/10
- Class presentation (7-8 minutes)
  - In class on 2/10 or 2/17
- Each part will
  - Describe workflow (1/3 of length/time)
  - Critique the efficiency, fidelity, & reproducibility (remaining 2/3 of length/time)

# PSYC 259: Principles of Data Science

Week 4: Part 1

Tidy Data

# Why do data need to be wrangled and tidied?

- Data files are often created without the foresight of how they might be used by computers
  - Instead, they are made to look nice for humans
- We need to tidy data to make it convenient for programming
  - Format is dictated by analyses
  - You may need multiple formats

#### Tidy data (Wickham, 2014)

- Each variable is a column
  - Variable: "All values that measure the same underlying attribute"
- Each observation is a row
  - Observation: Unit for which all variables were measured (person, session, trial, etc.)
- Each type of observational unit is a table
  - e.g., participant questionnaire vs trial-by-trial data

### Common types of "un-tidy" data: Column names are values

```
relig_income
#> # A tibble: 18 x 11
      religion `<$10k` `$10-20k` `$20-30k` `$30-40k` `$40-50k` `$50-75k` `$75-100k`
                                     <dbl>
      <chr>
                 <dbl>
                           <dbl>
                                               <dbl>
                                                          <dbl>
                                                                    <dbl>
                                                                               <dbl>
   1 Agnostic
                    27
                              34
                                        60
                                                  81
                                                                                 122
                                                            76
                                                                      137
   2 Atheist
                                        37
                                                  52
                    12
                              27
                                                             35
                                                                       70
                                                                                  73
   3 Buddhist
                    27
                              21
                                        30
                                                  34
                                                            33
                                                                       58
                                                                                  62
#> 4 Catholic
                             617
                                       732
                                                            638
                   418
                                                  670
                                                                     1116
                                                                                 949
```

```
relig_income
#> # A tibble: 18 x 11
      religion `<$10k` `$10-20k` `$20-30k` `$30-40k` `$40-50k` `$50-75k` `$75-100k`
                           <dbl>
                                     <dbl>
                                                <dbl>
                                                          <dbl>
                                                                    <dbl>
      <chr>
                 <dbl>
                                                                                <dbl>
#>
   1 Agnostic
                    27
                              34
                                         60
                                                   81
                                                             76
                                                                      137
                                                                                  122
   2 Atheist
                    12
                              27
                                         37
                                                   52
                                                             35
                                                                       70
                                                                                   73
   3 Buddhist
                    27
                              21
                                         30
                                                   34
                                                             33
                                                                       58
                                                                                   62
#> 4 Catholic
                   418
                             617
                                        732
                                                  670
                                                            638
                                                                     1116
                                                                                  949
```

```
relig_income
#> # A tibble: 18 x 11
      religion `<$10k` `$10-20k` `$20-30k` `$30-40k` `$40-50k` `$50-75k` `$75-100k`
                 <dbl>
                           <dbl>
                                      <dbl>
                                                <dbl>
                                                           <dbl>
                                                                     <dbl>
      <chr>
                                                                                 <dbl>
   1 Agnostic
                    27
                              34
                                         60
                                                              76
                                                                                  122
                                                   81
                                                                       137
   2 Atheist
                              27
                                         37
                                                              35
                                                                        70
                                                                                   73
   3 Buddhist
                              21
                                         30
                                                              33
                                                                        58
                                                                                   62
#> 4 Catholic
                   418
                              617
                                        732
                                                  670
                                                             638
                                                                      1116
                                                                                   949
```

- Income level can't be made a factor to use in a model
- Challenging to summarize, each income level needs to be treated as a separate variable
- Need to use different verbs to subset observations (filter for religion, select for income)

# Solution: pivot\_longer() turns columns into rows

```
relig_income
#> # A tibble: 18 x 11
     religion `<$10k` `$10-20k` `$20-30k` `$30-40k` `$40-50k` `$50-75k` `$75-100k`
     <chr>
                <dbl>
                          <dbl>
                                   <dbl>
                                             <dbl>
                                                       <dbl>
                                                                 <dbl>
                                                                           <dbl>
#>
#> 1 Agnostic
                   27
                                      60
                                                81
                            34
                                                          76
                                                                  137
                                                                             122
#> 2 Atheist 12
                                      37
                                                52
                            27
                                                          35
                                                                   70
                                                                              73
  3 Buddhist
                27
                            21
                                      30
                                                34
                                                          33
                                                                   58
                                                                              62
#> 4 Catholic
                 418
                            617
                                     732
                                               670
                                                         638
                                                                 1116
                                                                             949
```

```
relig_income %>%
  pivot_longer(!religion, names_to = "income", values_to = "count")
#> # A tibble: 180 x 3
      religion income
#>
                                  count
      <chr>
               <chr>
                                  <dbl>
#>
   1 Agnostic <$10k
                                      27
    2 Agnostic $10-20k
                                      34
   3 Agnostic $20-30k
                                      60
    4 Agnostic $30-40k
                                      81
    5 Agnostic $40-50k
                                      76
```

## Common types of "un-tidy" data: Multiple variables in one column

- Mixing types in the same column
- No meaningful summary of "value"
- Impossible to use values in a model (i.e., outcome ~ order\*sex\*block)

# Solution: pivot\_wider() turns rows into columns

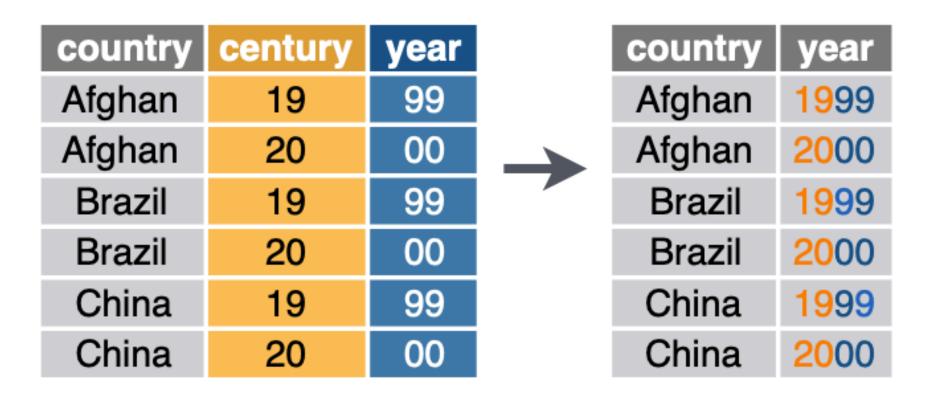
# Common types of "un-tidy" data: Variables glued into one column

| country | year | rate     |               | country | year | cases | pop |
|---------|------|----------|---------------|---------|------|-------|-----|
| Α       | 1999 | 0.7K/19M |               | Α       | 1999 | 0.7K  | 19M |
| Α       | 2000 | 2K/20M   | $\rightarrow$ | Α       | 2000 | 2K    | 20M |
| В       | 1999 | 37K/172M |               | В       | 1999 | 37K   | 172 |
| В       | 2000 | 80K/174M |               | В       | 2000 | 80K   | 174 |
| С       | 1999 | 212K/1T  |               | С       | 1999 | 212K  | 1T  |
| С       | 2000 | 213K/1T  |               | С       | 2000 | 213K  | 1T  |

```
separate(table3, rate, sep = "/",
into = c("cases", "pop"))
```

- Solution: use separate()

## Common types of "un-tidy" data: Variables split across columns



unite(table5, century, year, col = "year", sep = "")

- Solution: use unite()

# Common types of "un-tidy" data: Observations split across tables

```
    DF
    A
    B
    C

    x
    a
    t
    1

    x
    b
    u
    2

    x
    c
    v
    3

    z
    c
    v
    3

    z
    d
    w
    4
```

**bind\_rows(...**, .id = NULL)
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured)

- Solution: use bind\_rows()
- Solution: use vroom() if data are split across files

## Common types of "un-tidy" data: Variables split across tables

left\_join(x, y, by = NULL, copy=FALSE, suffix=c(".x",".y"),...)

Join matching values from y to x. **left\_join(**x, y, by = NULL,

- Solution: use left\_join() or other joins
- Why not bind\_cols?

# Data tidying tutorial

Example scripts 3 and 4 in "259-data-wrangling"