# Data Summarization

# Recap

- `select()`: subset and/or reorder columns

- `filter()`: remove rows

- `arrange()`: reorder rows

- `mutate()`: create new columns or modify them

- `select()` and `filter()` can be combined together

- remove a column: `select()` with `!` mark (`!col_name`)

- you can do sequential steps: especially using pipes **%>%**

 Cheatsheet

# Another Cheatsheet

https://raw.githubusercontent.com/rstudio/cheatsheets/main/data-transformation.pdf
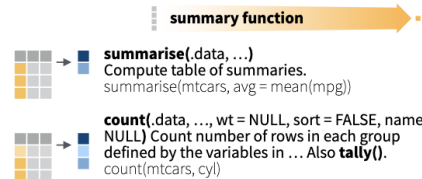
# Data Summarization

- Basic statistical summarization

  - `mean(x)`: takes the mean of x

  - `sd(x)`: takes the standard deviation of x

  - `median(x)`: takes the median of x

  - `quantile(x)`: displays sample quantiles of x. Default is min, IQR, max

  - `range(x)`: displays the range. Same as `c(min(x), max(x))`

  - `sum(x)`: sum of x

  - `max(x)`: maximum value in x

  - `min(x)`: minimum value in x

# Some examples

We can use the `CO_heat_ER` object from the `dasehr` package to explore different ways of summarizing data. (This dataset contains information about the number and rate of visits for heat-related illness to ERs in Colorado from 2011-2022, adjusted for age.) The `head` command displays the first rows of an object:

```
library(dasehr)
head(CO_heat_ER)

# A tibble: 6 × 7
  county      rate lower95cl upper95cl visits  year gender
  <chr>      <dbl>    <dbl>     <dbl>  <dbl> <dbl> <chr>
1 Statewide   5.64     4.70      6.59    140  2011 Female
2 Statewide   7.39     6.30      8.47    183  2011 Male
3 Statewide   6.51     5.80      7.23    323  2011 Both genders
4 Statewide   5.64     4.72      6.57    146  2012 Female
5 Statewide   7.56     6.48      8.65    193  2012 Male
6 Statewide   6.58     5.88      7.29    339  2012 Both genders
```

# Behavior of `pull()` function

`pull()` converts a single data column into a vector. This allows you to run summary functions.

```
CO_heat_ER %>% pull(visits)
```

# Statistical summarization the "tidy" way

**Add the `na.rm` = argument for missing data**

```
CO_heat_ER %>% pull(visits) %>% mean()

[1] NA

CO_heat_ER %>% pull(visits) %>% mean(na.rm=T)

[1] 9.791114
```

# Summarization on tibbles (data frames)

# Summarize the data: `dplyr summarize()` function

`summarize` creates a summary table.

Multiple summary statistics can be calculated at once (unlike `pull()` which can only do a single calculation on one column).

```
# General format - Not the code!
{data to use} %>%
    summarize({summary column name} = {function(source column)},
             {summary column name} = {function(source column)})
```

# Summarize the data: `dplyr summarize()` function

```
CO_heat_ER %>%
  summarize(mean_visits = mean(visits))

# A tibble: 1 × 1
  mean_visits
        <dbl>
1          NA

CO_heat_ER %>%
  summarize(mean_visits = mean(visits, na.rm = TRUE))

# A tibble: 1 × 1
  mean_visits
        <dbl>
1        9.79
```

# Summarize the data: `dplyr summarize()` function

`summarize()` can do multiple operations at once. Just separate by a comma.

```
CO_heat_ER %>%
  summarize(mean_visits = mean(visits, na.rm = TRUE),
            median_visits = median(visits, na.rm = TRUE),
            mean_rate = mean(rate, na.rm = TRUE))

# A tibble: 1 × 3
  mean_visits median_visits mean_rate
        <dbl>         <dbl>     <dbl>
1        9.79             0      1.87
```

# Summarize the data: `dplyr summarize()` function

Note that `summarize()` creates a separate tibble from the original data.

If you want to save a summary statistic in the original data, use `mutate()` instead to create a new column for the summary statistic.

# summary() Function

Using `summary()` can give you rough snapshots of each numeric column
(character columns are skipped):

```
summary(CO_heat_ER)
```

```
    county                rate              lower95cl            upper95cl
 Length:2340         Min.   : 0.000    Min.   : 0.000     Min.   :  0.000
 Class :character    1st Qu.: 0.000    1st Qu.: 0.000     1st Qu.:  0.000
 Mode  :character    Median : 0.000    Median : 0.000     Median :  0.000
                     Mean   : 1.869    Mean   : 1.119     Mean   :  2.755
                     3rd Qu.: 0.000    3rd Qu.: 0.000     3rd Qu.:  0.000
                     Max.   :89.275    Max.   :43.398     Max.   :151.420
                     NA's   :832       NA's   :832        NA's   :832

    visits              year            gender
 Min.   :  0.000   Min.   :2011    Length:2340
 1st Qu.:  0.000   1st Qu.:2014    Class :character
 Median :  0.000   Median :2016    Mode  :character
 Mean   :  9.791   Mean   :2016
 3rd Qu.:  0.000   3rd Qu.:2019
 Max.   :494.000   Max.   :2022
 NA's   :832
```

# Summary & Lab Part 1

- summary stats (`mean()`) work with `pull()`

- don't forget the `na.rm = TRUE` argument!

- `summary(x)`: quantile information

- `summarize`: creates a summary table of columns of interest

 Class Website

 Lab

# **distinct()** values

`distinct(x)` will return the unique elements of column `x`.

```
CO_heat_ER %>%
  distinct(gender)

# A tibble: 3 × 1
  gender
  <chr>
1 Female
2 Male
3 Both genders
```

# How many `distinct()` values?

`n_distinct()` tells you the number of unique elements. *Must pull the column first!*

```
CO_heat_ER %>%
  pull(gender) %>%
  n_distinct()
```

```
[1] 3
```

# dplyr: count

Use `count` to return row count by category.

```
CO_heat_ER %>% count(gender)

# A tibble: 3 × 2
  gender            n
  <chr>         <int>
1 Both genders    780
2 Female          780
3 Male            780
```

# dplyr:count

Multiple columns listed further subdivides the count.

```
CO_heat_ER %>% count(county, gender)

# A tibble: 195 × 3
   county    gender            n
   <chr>     <chr>         <int>
 1 Adams     Both genders    12
 2 Adams     Female          12
 3 Adams     Male            12
 4 Alamosa   Both genders    12
 5 Alamosa   Female          12
 6 Alamosa   Male            12
 7 Arapahoe  Both genders    12
 8 Arapahoe  Female          12
 9 Arapahoe  Male            12
10 Archuleta Both genders    12
# 	 185 more rows
```

# Grouping

# Perform Operations By Groups: dplyr

`group_by` allows you group the data set by variables/columns you specify:

```
CO_heat_ER_grouped <- CO_heat_ER %>% group_by(gender)
CO_heat_ER_grouped

# A tibble: 2,340 × 7
# Groups:   gender [3]
   county      rate lower95cl upper95cl visits  year gender
   <chr>      <dbl>     <dbl>     <dbl>  <dbl> <dbl> <chr>
 1 Statewide   5.64      4.70      6.59    140  2011 Female
 2 Statewide   7.39      6.30      8.47    183  2011 Male
 3 Statewide   6.51      5.80      7.23    323  2011 Both genders
 4 Statewide   5.64      4.72      6.57    146  2012 Female
 5 Statewide   7.56      6.48      8.65    193  2012 Male
 6 Statewide   6.58      5.88      7.29    339  2012 Both genders
 7 Statewide   4.94      4.06      5.82    124  2013 Female
 8 Statewide   6.72      5.72      7.72    178  2013 Male
 9 Statewide   5.82      5.16      6.49    302  2013 Both genders
10 Statewide   3.52      2.80      4.25     92  2014 Female
# ⏳ 2,330 more rows
```

# Summarize the grouped data

It's grouped! Grouping doesn't change the data in any way, but how **functions operate on it**. Now we can summarize `visits` by group:

```
CO_heat_ER_grouped %>%
  summarize(avg_visits = mean(visits, na.rm = TRUE))

# A tibble: 3 × 2
  gender         avg_visits
  <chr>              <dbl>
1 Both genders       16.3
2 Female              4.77
3 Male                9.00
```

# Use the **pipe** to string these together!

Pipe `CO_heat_ER` into `group_by`, then pipe that into `summarize`:

```
CO_heat_ER %>%
  group_by(gender) %>%
  summarize(avg_visits = mean(visits, na.rm = TRUE))

# A tibble: 3 × 2
  gender         avg_visits
  <chr>               <dbl>
1 Both genders         16.3
2 Female                4.77
3 Male                  9.00
```

# Group by as many variables as you want

`group_by` gender and year:

```
CO_heat_ER %>%
  group_by(year, gender) %>%
  summarize(avg_visits = mean(visits, na.rm = TRUE))

# A tibble: 36 × 3
# Groups:   year [12]
     year gender        avg_visits
    <dbl> <chr>              <dbl>
 1   2011 Both genders       11.3
 2   2011 Female              4.32
 3   2011 Male                6.06
 4   2012 Both genders       12.8
 5   2012 Female              4.76
 6   2012 Male                6.71
 7   2013 Both genders       12.4
 8   2013 Female              3.72
 9   2013 Male                6.11
10   2014 Both genders        9.67
#  26 more rows
```

# Counting

There are other functions, such as `n()` count the number of observations (NAs included).

```
CO_heat_ER %>%
  group_by(gender) %>%
  summarize(n = n(),
            mean = mean(visits, na.rm = TRUE))

# A tibble: 3 × 3
  gender           n  mean
  <chr>        <int> <dbl>
1 Both genders   780 16.3
2 Female         780  4.77
3 Male           780  9.00
```

# Counting

`count()` and `n()` can give very similar information.

```
CO_heat_ER %>% count(gender)


# A tibble: 3 × 2
  gender            n
  <chr>         <int>
1 Both genders    780
2 Female          780
3 Male            780


CO_heat_ER %>% group_by(gender) %>% summarize(n()) # n() typically used with summarize


# A tibble: 3 × 2
  gender         `n()`
  <chr>         <int>
1 Both genders    780
2 Female          780
3 Male            780
```

A few miscellaneous topics ..

# Base R functions you might see: `length` and `unique`

These functions require a column as a vector using `pull()`.

```
CO_heat_ER_gen <- CO_heat_ER %>% pull(gender) # pull() to make a vector
CO_heat_ER_gen %>% unique() # similar to distinct()

[1] "Female"        "Male"          "Both genders"
```

# Base R functions you might see: **length** and **unique**

These functions require a column as a vector using `pull()`.

```
CO_heat_ER_gen %>% unique() %>% length() # similar to n_distinct()

[1] 3
```

# \* New! \* Many dplyr functions now have a **.by=** argument

Pipe `CO_heat_ER` into `group_by`, then pipe that into `summarize`:

```
CO_heat_ER %>%
  group_by(gender) %>%
  summarize(avg_visits = mean(visits, na.rm = TRUE),
            max_visits = max(visits, na.rm = TRUE))
```

is the same as..

```
CO_heat_ER %>%
  summarize(avg_visits = mean(visits, na.rm = TRUE),
            max_visits = max(visits, na.rm = TRUE),
            .by = county)
```

# `summary()` vs. `summarize()`

- `summary()` (base R) gives statistics table on a dataset.

- `summarize()` (dplyr) creates a more customized summary tibble/dataframe.

# Summary & Lab Part 2

- `count(x)`: what unique values do you have?

    - `distinct()`: what are the distinct values?

    - `n_distinct()` with `pull()`: how many distinct values?

- `group_by()`: changes all subsequent functions

    - combine with `summarize()` to get statistics per group

    - combine with `mutate()` to add column

- `summarize()` with `n()` gives the count (NAs included)

 Class Website

 Lab

# Extra Slides: More advanced summarization

# Data Summarization on data frames

- Statistical summarization across the data frame

    - `rowMeans(x)`: takes the means of each row of x

    - `colMeans(x)`: takes the means of each column of x

    - `rowSums(x)`: takes the sum of each row of x

    - `colSums(x)`: takes the sum of each column of x

```
yearly_co2 <- yearly_co2_emissions
```

# rowMeans() example

Get means for each row.

Let's see what the mean CO2 emissions is across years for each row (country):

```
yearly_co2 %>%
  select(starts_with("201")) %>%
  rowMeans(na.rm = TRUE) %>%
  head(n = 5)

[1]  10254   5106 129800    487  32040

yearly_co2 %>%
  group_by(country) %>%
  summarize(mean = rowMeans(across(starts_with("201")), na.rm = TRUE)) %>%
  head(n = 5)

# A tibble: 5 × 2
  country      mean
  <chr>       <dbl>
1 Afghanistan 10254
2 Albania      5106
3 Algeria    129800
4 Andorra       487
5 Angola      32040
```

# `colMeans()` example

Get means for each column.

Let's see what the mean is across each column (year):

```r
yearly_co2 %>%
  select(starts_with("201")) %>%
  colMeans(na.rm = TRUE) %>%
  head(n = 5)
```

```
    2010      2011      2012      2013      2014
165334.1 171764.9 174033.4 174856.2 175992.5
```

```r
yearly_co2 %>%
  summarize(across(starts_with("201"), ~mean(.x, na.rm = TRUE)))
```

```
# A tibble: 1 × 5
  `2010`  `2011`  `2012`  `2013`  `2014`
   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 165334. 171765. 174033. 174856. 175993.
```