

Numerical Summaries

Summarizing data in R 1/2

- Have seen `summary` (5-number summary of each column). But what if we want:
 - a summary or two of just one column
 - a count of observations in each category of a categorical variable
 - summaries by group
 - a different summary of all columns (eg. SD)
- To do this, meet pipe operator `%>%`. This takes input data frame, does something to it, and outputs result. (Learn: **Ctrl-Shift-M**.)

Summarizing data in R 2/2

- Output from a pipe can be used as input to something else, so can have a sequence of pipes.
- Summaries include: `mean`, `median`, `min`, `max`, `sd`, `IQR`, `quantile` (for obtaining quartiles or any percentile), `n` (for counting observations).
- Use our Australian athletes data again.

Packages for this section

```
library(tidyverse)
```

Summarizing one column

- Mean height:

```
athletes %>% summarize(m=mean(Ht))
```

```
# A tibble: 1 x 1
      m
  <dbl>
1  180.
```

or to get mean and SD of BMI:

```
athletes %>% summarize(m=mean(BMI),s=sd(BMI))
```

```
# A tibble: 1 x 2
      m      s
  <dbl> <dbl>
1  23.0  2.86
```

Quartiles

- `quantile` calculates percentiles (“fractiles”), so we want the 25th and 75th percentiles:

```
athletes %>% summarize( Q1=quantile(Wt, 0.25),
                        Q3=quantile(Wt, 0.75))
```

```
# A tibble: 1 x 2
      Q1    Q3
  <dbl> <dbl>
1  66.5  84.1
```

Creating new columns

- These weights are in kilograms. Maybe we want to summarize the weights in pounds.
- Convert kg to lb by multiplying by 2.2.
- Create new column and summarize that:

```
athletes %>% mutate(wt_lb=Wt*2.2) %>%
  summarize(Q1_lb=quantile(wt_lb, 0.25),
            Q3_lb=quantile(wt_lb, 0.75))
```

```
# A tibble: 1 x 2
  Q1_lb Q3_lb
  <dbl> <dbl>
1  146.  185.
```

Counting how many

for example, number of athletes in each sport:

```
athletes %>% count(Sport)
```

```
# A tibble: 10 x 2
```

	Sport	n
	<chr>	<int>
1	BBall	25
2	Field	19
3	Gym	4
4	Netball	23
5	Row	37
6	Swim	22
7	T400m	29
8	TSprnt	15
9	Tennis	11
10	WPolo	17

Counting how many, variation 2:

Another way (which will make sense in a moment):

```
athletes %>% group_by(Sport) %>%  
  summarize(count=n())
```

```
# A tibble: 10 x 2
```

	Sport	count
	<chr>	<int>
1	BBall	25
2	Field	19
3	Gym	4
4	Netball	23
5	Row	37
6	Swim	22
7	T400m	29
8	TSprnt	15
9	Tennis	11
10	WPolo	17

Summaries by group

- Might want separate summaries for each “group”, eg. mean and SD of height for males and females. Strategy is `group_by` (to define the groups) and then `summarize`:

```
athletes %>% group_by(Sex) %>%  
  summarize(mean_Ht = mean(Ht), sd_Ht = sd(Ht))
```

```
# A tibble: 2 x 3  
  Sex    mean_Ht sd_Ht  
  <chr>    <dbl> <dbl>  
1 female    175.  8.24  
2 male     186.  7.90
```

Count plus stats

- If you want number of observations per group plus some stats, you need to go the `n()` way:

```
athletes %>% group_by(Sex) %>%  
  summarize(n = n(), mean_Ht = mean(Ht), sd_Ht = sd(Ht))
```

```
# A tibble: 2 x 4  
  Sex      n mean_Ht sd_Ht  
  <chr> <int>    <dbl> <dbl>  
1 female   100    175.  8.24  
2 male    102    186.  7.90
```

- This explains second variation on counting within group: “within each sport/Sex, how many athletes were there?”

Summarizing several columns

- Standard deviation of each (numeric) column:

```
athletes %>% summarize(across(where(is.numeric), \ (x) sd(x)))
```

```
# A tibble: 1 x 11
  RCC    WCC    Hc    Hg  Ferr  BMI  SSF  `"%Bfat"`  LBM  Ht  Wt
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 0.458  1.80  3.66  1.36  47.5  2.86  32.6    6.19  13.1  9.73  13.9
```

- Median and IQR of all columns whose name starts with H:

```
athletes %>% summarize(across(starts_with("H"),
                               list(
med = \(x) median(x),
iqr = \(x) IQR(x))))
```

```
# A tibble: 1 x 6
  Hc_med Hc_iqr Hg_med Hg_iqr Ht_med Ht_iqr
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1  43.5  4.98  14.7  2.07  180.  12.2
```

Same thing by group

```
athletes %>%
  group_by(Sex) %>%
  summarize(across(starts_with("H"),
                    list(
med = \(h) median(h),
iqr = \(h) IQR(h))))
```

```
# A tibble: 2 x 7
  Sex    Hc_med Hc_iqr Hg_med Hg_iqr Ht_med Ht_iqr
<chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 female  40.6  4.03  13.5  1.60  175  8.68
2 male   45.5  2.57  15.5  0.975 186. 11.3
```