

# Oh the Summaries . . .

Directions: Follow along with the slides and answer the questions in **BOLDED** font in your journal.

## Just the beginning

- Means, medians, MAD & SD are just a few examples of **numerical summaries**.
- **Numerical summaries** are numbers that describe characteristics of the data.
  - Means & medians describe the *center* of the data.
  - MAD & SD describe the *spread* of the data.
- In this lab, we will look at what other numbers are useful for describing data.

## To start

- Load your *Personality Color* data again and name it: **colors**.
- In the lines of code that appear on the following slides:
  - Replace any `~x` with the name of the variable for your *predominant* color score.

## Extreme values

- Besides looking at *typical* values, sometimes we want to see *extreme* values, like the smallest and largest values.
- To find these values, we can calculate the **min** and **max**.

```
min(~x, data=colors)
```

```
max(~x, data=colors)
```

- Or use the **range** function to compute both

```
range(~x, data=colors)
```

## Range

- While the **range** function will compute the smallest and largest values, the **range** is also a term we use to describe the spread of our data.
- Calculate it by taking: **max - min**.
- The range is often much less informative than other measures of *spread*.
- **What makes the range less informative than the MAD?**
  - Can you think of examples where the MAD will give you a better idea of the *variability* than the **range**?

## Quartiles (Q1 & Q3)

- We often use the **median** to describe the *center* of our data because half of the data is smaller than the median and the other half is larger.
- If instead we found a value that was larger than just 25% of our data, we would have computed the *1st quartile*.
- If we found the value that was larger than 75% of our data, we call that the *3rd quartile*.
- **Why do you think we use the names ‘1st and 3rd quartiles’?**

## The Inter-Quartile-Range (IQR)

- Just like we used the **min** and **max** to compute the **range**, we can also use the *1st* and *3rd* quartiles to compute the IQR.
- The IQR is another way to describe *spread*.
  - It describes how *wide* or *narrow* the middle 50% of our data are.
  - If the IQR is a small number, then the middle 50% of our data is close to the **median**.
  - Otherwise, the middle 50% of our data is further away from the median.

## Finding the IQR

- Make a histogram of your *predominant* color’s scores.
- Visually (Don’t worry about being super-precise):
  - Cut the distribution into quarters so the *number of data points* is equal for each piece. (Each piece should contain 25% of the data.)
  - **Write down the numbers that split the data up into these 4 pieces.**
  - **How long is the interval of the middle two pieces?**
  - This length is the *IQR*.

## Calculating the IQR

- Calculate the IQR by using either of the following

```
IQR(~x, data = colors)
```

```
iqr(~x, data = colors)
```

- **Compare your visual estimate of the IQR to the actual IQR.**

## Other quantiles

- The median, 1st and 3rd *quartiles* can also be called the 50th, 25th and 75th *quantiles*.
  - They are called *quantiles* because they describe the *quantity* of data that is smaller than that value.

- The 25th quantile is the value that is larger than 25% of the data.
- We can compute quantiles too!

```
qdata(~x, data = colors, p = 0.35)
```

- Where *p* stands for the *percentage* of data you'd like our value to be larger than.

## Boxplots

- By using the medians, quartiles, and min/max, we can construct a new single variable plot called the **box and whisker** plot, often shortened to just a **boxplot**.
- Try making one of your predominant color.

```
bwplot(~x, data=colors)
```

- Sketch your boxplot in your journal. Label the min, max, Q1, Q3, and the median.
- How would you interpret your boxplot? Where is the bulk of your data? Where is it centered? Can you say anything about its shape?

## Our favorite summaries

- Numerical summaries are brief ways to describe our data, using numbers.
- However, computing lots of different summaries can be tedious.
- Use the following command to compute some of our *favorite* summaries.

```
favstats(~x, data=colors)
```

- Which summaries are displayed?
- What do you think *n* stands for?