

# Correlations

# Recap

## Population Variability

### Sums of squares

$$SS = \sum (X_i - \mu_x)^2$$

### Variance

$$\sigma^2 = \frac{\sum (X_i - \mu_x)^2}{N} = \frac{SS}{N}$$

### Standard deviation

$$\sigma = \sqrt{\frac{\sum (X_i - \mu_x)^2}{N}} = \sqrt{\frac{SS}{N}} = \sqrt{\sigma^2}$$

## Sample variability

### Sums of squares

$$SS = \sum (X_i - \bar{X})^2$$

### Variance

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{N - 1} = \frac{SS}{N - 1}$$

### Standard deviation

$$s = \sqrt{\frac{\sum (X_i - \bar{X})^2}{N - 1}} = \sqrt{\frac{SS}{N - 1}} = \sqrt{s^2}$$

# Bi-variate descriptives

## Covariation

"Sum of the cross-products"

## Population

$$SP_{XY} = \Sigma(X_i - \mu_X)(Y_i - \mu_Y)$$

## Sample

$$SP_{XY} = \Sigma(X_i - \bar{X})(Y_i - \bar{Y})$$

# Covariance

Sort of like the variance of two variables

## Population

$$\sigma_{XY} = \frac{\Sigma(X_i - \mu_X)(Y_i - \mu_Y)}{N}$$

## Sample

$$s_{XY} = cov_{XY} = \frac{\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{N - 1}$$

## Covariance table

$$\mathbf{K}_{\mathbf{XX}} = \begin{bmatrix} \sigma_X^2 & cov_{XY} & cov_{XZ} \\ cov_{YX} & \sigma_Y^2 & cov_{YZ} \\ cov_{ZX} & cov_{ZY} & \sigma_Z^2 \end{bmatrix}$$

# Correlation

- Measure of association
- How much two variables are *linearly* related
- -1 to 1
- Sign indicates direction of relationship
- Invariant to changes in mean or scaling

# Correlation

Pearson product moment correlation

## Population

$$\rho_{XY} = \frac{\Sigma z_X z_Y}{N} = \frac{SP}{\sqrt{SS_X} \sqrt{SS_Y}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

## Sample

$$r_{XY} = \frac{\Sigma z_X z_Y}{n - 1} = \frac{SP}{\sqrt{SS_X} \sqrt{SS_Y}} = \frac{s_{XY}}{s_X s_Y}$$

```
data %>% ggplot(aes(x = x, y = y)) + geom_point(size = 3) + theme_bw()
```



What is the correlation between these two variables?



```
data %>% ggplot(aes(x = x, y = y)) + geom_point(size = 3) + theme_bw()
```



What is the correlation between these two variables?

```
data %>% ggplot(aes(x = x, y = y)) + geom_point(size = 3) + theme_bw()
```



What is the correlation between these two variables?

# Effect size

- Recall that z-scores allow us to compare across units of measure; the products of standardized scores are themselves standardized.
- The correlation coefficient is a **standardized effect size** which can be used to communicate the strength of a relationship.
- Correlations can be compared across studies, measures, constructs, time.
- Example: the correlation between age and height among children is  $r = .70$ . The correlation between self- and other-ratings of extraversion is  $r = .25$ .

# What is a large correlation?

- [Cohen \(1988\)](#): .1 (small), .3 (medium), .5 (large)
  - Often forgot: Cohen said only to use them when you had nothing else to go on, and has since regretted even suggesting benchmarks to begin with.
- $r^2$ : Proportion of variance "explained"
  - as [Ozer & Funder \(2019\)](#) discuss, we're not really explaining anything and the change in scale can mess up our interpretations if we're not careful.

# What are good benchmarks?

From Ozer & Funder (2019)

- Classic social psych studies:  $r = .36 - .42$
- Scarcity increases the perceived value of a commodity  $r = .12$
- People attribute failures to bad luck  $r = .10$
- Communicators perceived as more credible are more persuasive  $r = .10$
- People in a bad mood are more aggressive  $r = .41$
- Antihistamine and symptom relief  $r = .11$
- Ibuprofen and pain relief  $r = .14$
- Height and weight  $r = .44$

# What affects correlations?

It's not enough to calculate a correlation between two variables. You should always look at a figure of the data to make sure the number accurately describes the relationship. Correlations can be easily fooled by qualities of your data, like:

- Skewed distributions
- Outliers
- Restriction of range
- Nonlinearity

# Skewed distributions

```
p = data %>% ggplot(aes(x=x, y=y)) + geom_point()  
ggMarginal(p, type = "density")
```



# Outliers

```
data %>% ggplot(aes(x=x, y=y)) + geom_point()
```





# Outliers

```
data %>% ggplot(aes(x=x, y=y)) +  
  geom_point() +  
  geom_smooth(method = "lm",  
              se = FALSE,  
              color = "red") +  
  geom_smooth(data = data[-51,],  
              method = "lm",  
              se = FALSE)
```



# Restriction of range

```
data %>%  
  ggplot(aes(x=x, y=y)) +  
    geom_point() +  
    geom_smooth(method = "lm",  
                se = FALSE,  
                color = "red")
```



# Restriction of range

```
data %>%  
  ggplot(aes(x=x, y=y)) +  
    geom_point() +  
    geom_smooth(method = "lm",  
                se = FALSE,  
                color = "red") +  
    geom_point(data = real_data) +  
    geom_smooth(method = "lm",  
                se = FALSE,  
                data = real_data,  
                color = "blue")
```



# Nonlinearity

```
data %>% ggplot(aes(x=x, y=y)) + geom_point() + geom_smooth(method = "lr
```



# It's not always apparent

Sometimes issues that affect correlations won't appear in your graph, but you still need to know how to look for them.

- Low reliability
- Content overlap
- Multiple groups

# Multiple groups

```
data %>% ggplot(aes(x=x, y=y)) + geom_point() + geom_smooth(method = "lr
```



# Multiple groups

```
data %>% ggplot(aes(x=x, y=y, color = gender)) + geom_point() + geom_smooth()
```



# Special cases of the Pearson correlation

- **Spearman correlation coefficient**

- Applies when both X and Y are ranks (ordinal data) instead of continuous
- Denoted  $\rho$  by your textbook, although I prefer to save Greek letters for population parameters.

- **Point-biserial correlation coefficient**

- Applies when Y is binary.
  - NOTE: This is not an appropriate statistic when you artificially dichotomize data.

- **Phi (  $\phi$  ) coefficient**

- Both X and Y are dichotomous.