

Chapter 4

Scatterplots and Correlation

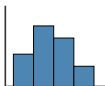


Intended Learning Outcomes

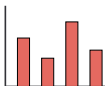
- Create and interpret scatterplots
- Identify explanatory and response variables
- Describe form, direction, and strength
- Recognise outliers
- Calculate and interpret correlation
- Apply properties of correlation
- Distinguish correlation from causation

Why Study Relationships Between Variables?

- So far, we have focused on **univariate** data analysis (one variable at a time).
- To investigate variables, we use tools like **histograms**, **bar graphs**, **stem plots**, and **time series**.



Histogram



Bar Graph



Time Series

- Many real-world questions involve more than one variable.

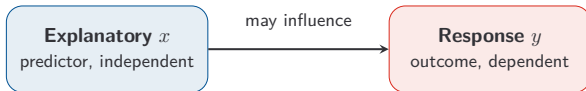
Examples of Bivariate Questions

1. Does **study time** affect **exam scores**?
2. Does **height** predict **earnings**?
3. Does **temperature** influence **ice cream sales**?
4. Does **exercise** relate to **resting heart rate**?

Response and Explanatory Variables

Response and Explanatory Variables

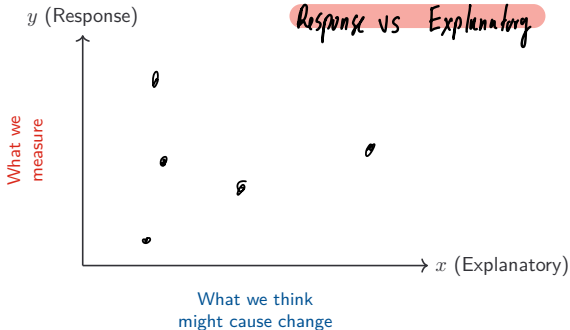
- The **response variable** y is the outcome of interest.
- The **explanatory variable** x may help explain or predict y .



Axis Placement Rule

Key Point: Place the **explanatory variable** on the x -axis and the **response variable** on the y -axis.

Convention for titles:



Example 1.13: Identifying Variables

For each scenario, identify the explanatory and response variables:

1. (Study time, Exam scores)

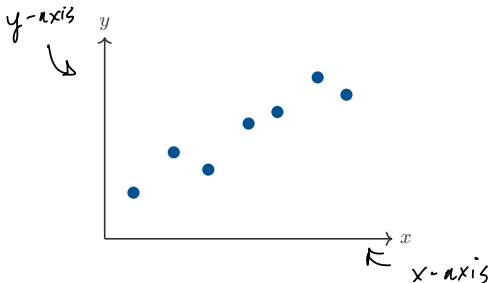
2. (Ice cream sales, Temperature)

3. (Education level, Income)

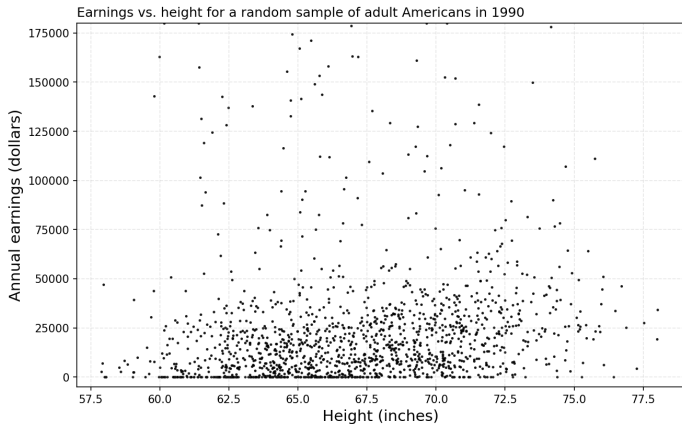
Scatterplot: Definition

Scatterplot

A **scatterplot** displays the relationship between two quantitative variables. Each individual appears as a **point** at coordinates (x_i, y_i) .

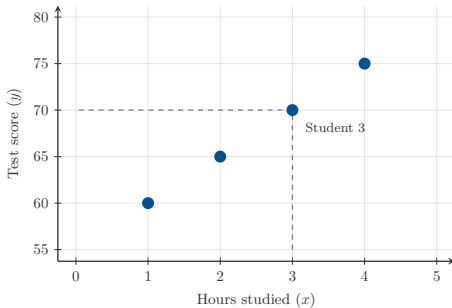


Heights and earnings



Study Time and Test Score

Hours studied (x)	1	2	3	4
Test score (y)	60	65	70	75



Constructing a Scatterplot

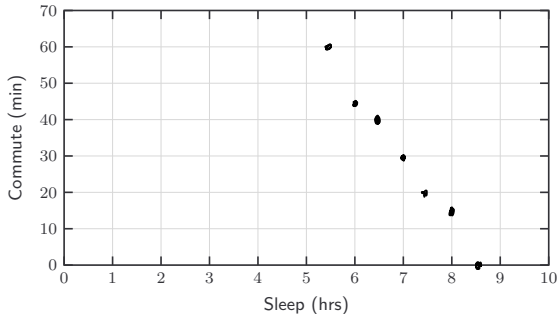
How to Create a Scatterplot

1. Identify the two quantitative variables
2. Decide which variable goes on each axis
3. Label the x -axis (horizontal)
4. Label the y -axis (vertical)
5. Plot each individual as a point at (x_i, y_i)
6. Add a descriptive title

Example 1.6: Plotting Data by Hand

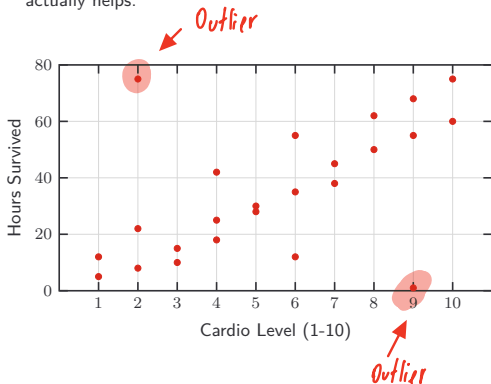
Plot the scatterplot for the following data:

y →	Commute (min)	60	45	30	15	0	20	40
x →	Sleep (hrs)	5.5	6.0	7.0	8.0	8.5	7.5	6.5



Example 1.7: Reading a Scatterplot

Context: A researcher analyzes data from a Zombie Outbreak Simulation to see if “Rule #1: Cardio” actually helps.



1. What is the **longest time** anyone survived?

~75 hours

2. Look at the person with **Cardio Level 10**. What was their lowest survival time?

60 hours

3. Are there any outlying points? If so, circle one and explain why it is an outlier.

(2, 75)

4. How many people with **Cardio Level 4** survived longer than 40 hours?

One individual

Plotting several variables simultaneously

We can encode a third (categorical) variable using:

- **Color:** different groups shown in different colors
- **Shape:** different groups shown with different marker shapes
- **Size:** continuous third variable shown by point size

for categorical data

For quantitative data

This allows us to explore relationships between three or more variables simultaneously.

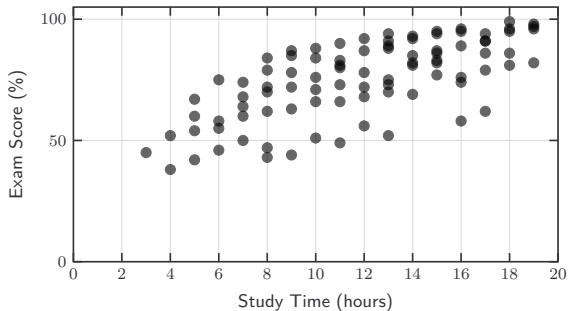
student_id	study_time	exam_score	attended_lectures	sleep_quantity
482910334	5	42	No	4
299401855	12	78	Yes	4.5
850223190	8	84	No	6
110594823	13	94	Yes	9
673820019	18	99	Yes	10
559201182	9	72	No	9

x

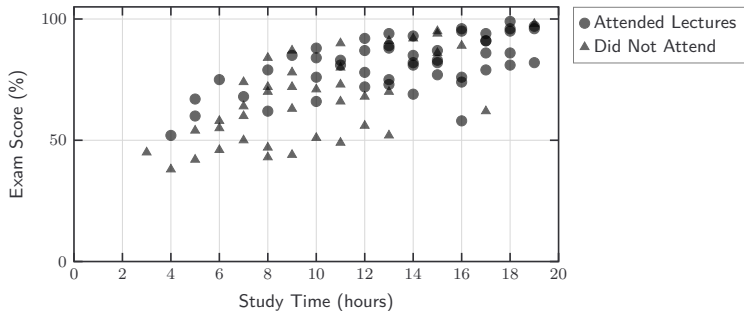
y

Example 1.8: Exam Scores vs. Study Time

Context: Study time vs. exam score, with different colors indicating sleep level (low, medium, high).

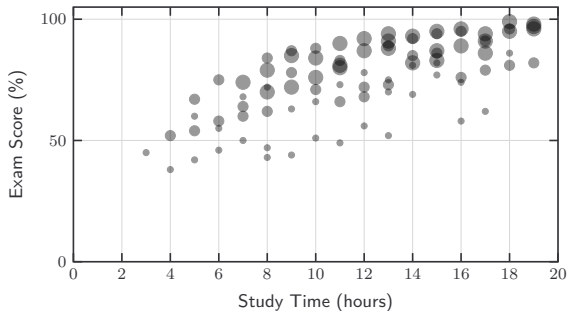


Example 1.10: Adding a Third Variable with Shape

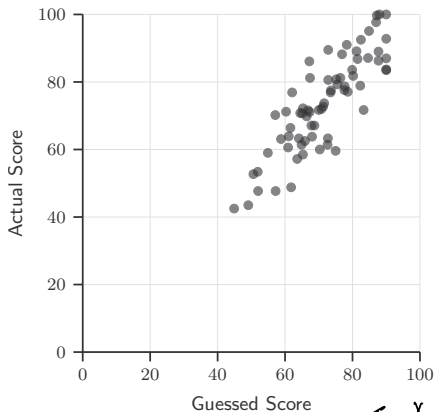


Example 1.11: Adding a Third Variable with Size

Context: Study time vs. exam score, with point size indicating hours of sleep (larger = more sleep).



How good were UWaterloo students at predicting their performance?

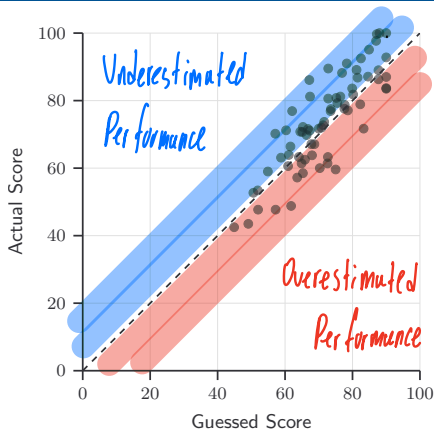


If $|x - y| \leq 5$, then
the student received +3

y: their actual performance
← x: Score guessed by students

How good were UWaterloo students at predicting their performance?

Adding a reference line

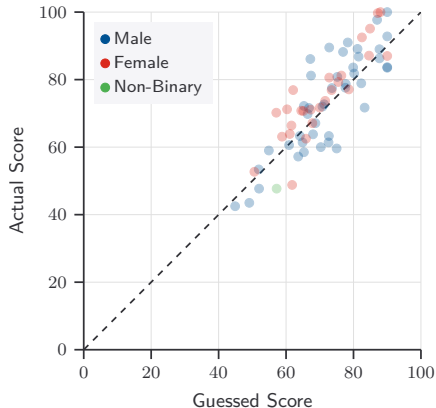


$y=x$ line represents perfect predictions

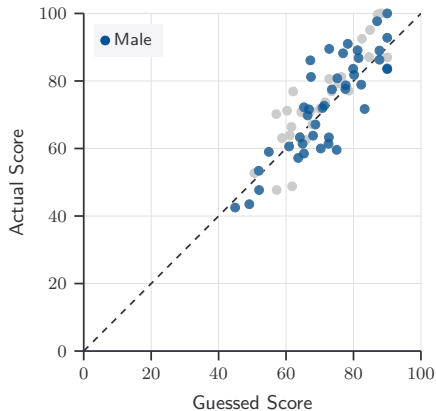
Points close to line ($|x-y| \leq 5$) got the bonus.

How good were UWaterloo students at predicting their performance?

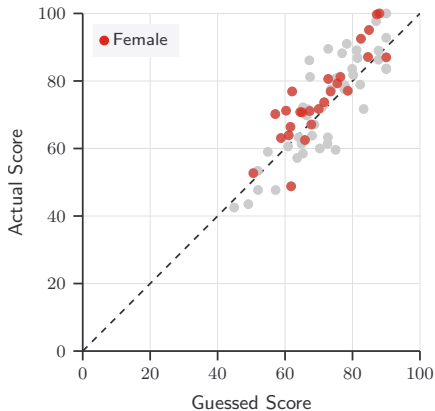
Broken down by gender



How good were men at predicting their performance?



How good were women at predicting their performance?





This is a bit harder than the level of interpretation we will focus on in this course, but it shows how scatterplots can be used to explore complex relationships between multiple variables.

Interpreting Scatterplots at the DS1000 level: Four Key Features

Outliers

Points far from
the pattern

Form

Linear or
nonlinear?

Direction

Positive or
negative?

Strength

Strong, moderate,
or weak?

Outliers in Scatterplots

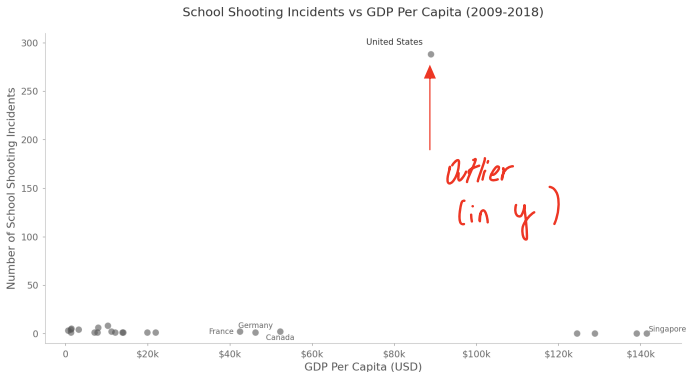
Outlier

A point that falls far from the overall pattern of the data.

Types of unusual points:

- **Vertically unusual:** far above or below the pattern
- **Horizontally unusual:** extreme x -value
- **Both:** unusual in both directions

Example 1.15: Outlier Example

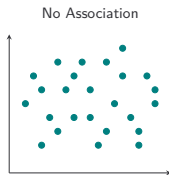
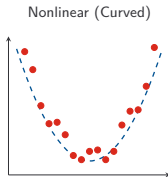
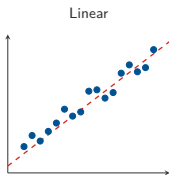


Form: Linear vs. Nonlinear

Form

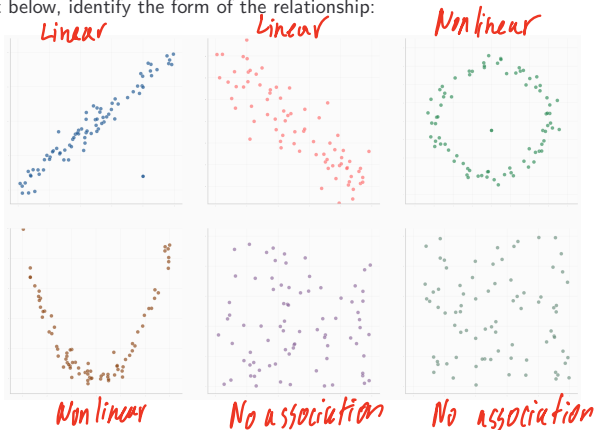
The **form** of a scatterplot describes the overall shape of the relationship

- **Linear:** Points cluster around a straight line.
- **Nonlinear:** Points follow a curved pattern (quadratic, exponential, logarithmic, etc.).
- **No association:** Points appear scattered with no discernible pattern.



Example 1.17: Finding Form

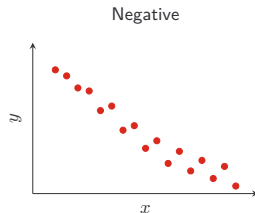
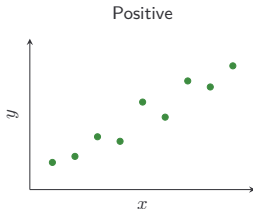
For each plot below, identify the form of the relationship:



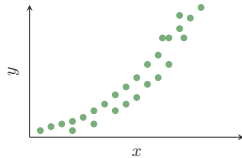
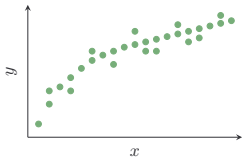
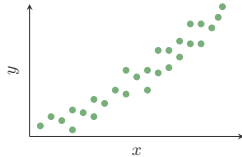
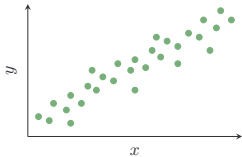
Direction

Direction

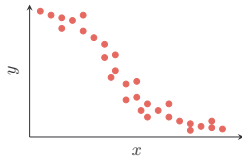
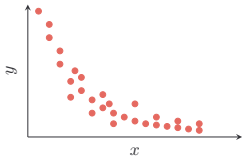
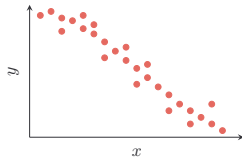
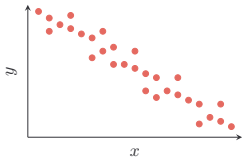
- **Positive:** as x increases, y tends to **increase**.
- **Negative:** as x increases, y tends to **decrease**.



Positive Association



Negative Association



Strength: Strong, Moderate, or Weak

Strength

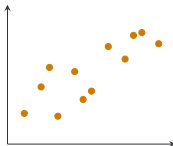
The **strength** of a relationship describes how closely the data points follow the overall pattern:

- **Strong:** Points lie close to the pattern with minimal scatter.
- **Moderate:** Points show noticeable scatter but the pattern remains clear.
- **Weak:** Points are widely scattered and the pattern is difficult to discern.

Strong



Moderate



Weak

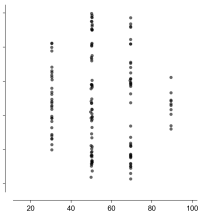
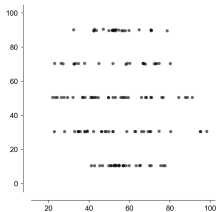
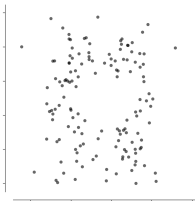
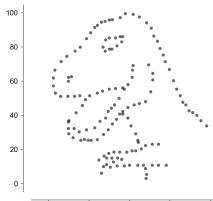


Describing a Scatterplot: Template

How to Describe a Scatterplot

1. Note any **outliers** or unusual points
2. State the **form**: linear, curved, or no clear form
3. State the **direction**: positive, negative, or none
4. State the **strength**: strong, moderate, or weak

Which plot below has the strongest linear relationship?



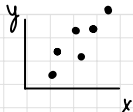
The Problem with “Eye-balling” Data

- **Subjectivity:** What looks “strong” to me might look “moderate” to you.
- **Scale:** Changing the axis scales can make the same data look steeper or flatter.
- **Precision:** We cannot compare relationships across different datasets (e.g., Height vs. Weight) using just adjectives.

We need an objective, numerical measure: **Correlation** (r).

Last time

- Scatterplots



- x : explanatory, y : response

- Describing scatter plots

- Outliers

- Shape/form: linear / non-linear / no association

- Strength: adherence to form

- Direction: + or -

} Not rigorous

Goal: Define a measure $\text{corr}(x, y)$ which quantifies the strength and direction present between the two variables.

Annotations:
- Arrow from x to "one variable"
- Arrow from y to "one variable"

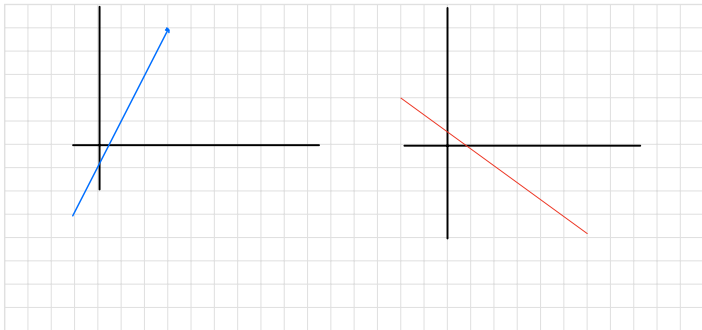
What do we want $\text{corr}(x, y)$ to satisfy?

- Symmetric: $\text{corr}(x, y) = \text{corr}(y, x)$

- Invariant under unit changes

e.g. corr of height (cm) and weight (kg) should be the same as height (in) and weight (lbs).
(same form and strength)

To achieve this, we need to be able to compare units.



What is a z -score?

z -score (Standardized Score)

The z -score of a value tells us how many standard deviations it is from the mean:

$$z_i = \frac{x_i - \bar{x}}{s_x}$$

- x_i : the value of interest
- \bar{x} : the mean of all values
- s_x : the standard deviation

How many std deviations the

- A positive z -score means the value is above the mean.
- A negative z -score means the value is below the mean.
- The larger the absolute value of z , the farther from the mean.

Example 1.23: Calculating z -scores

- Example 1:** Suppose the mean height in a class is 170 cm with a standard deviation of 8 cm. What is the z -score for a student who is 178 cm tall? $x = 178, \bar{x} = 170, s_x = 8$

$$z = \frac{x - \bar{x}}{s_x} = \frac{178 - 170}{8} = \frac{8}{8} = 1.$$

- Example 2:** If the mean test score is 75 with a standard deviation of 10, what is the z -score for a score of 60? $y = 60, \bar{y} = 75, s_y = 10$

$$z = \frac{60 - 75}{10} = \frac{-15}{10} = -1.5$$

- Example 3:** The average commute time is 30 minutes, with a standard deviation of 5 minutes. What is the z -score for a 35-minute commute? $x = 35, \bar{x} = 30, s_x = 5$

$$z = \frac{35 - 30}{5} = \frac{5}{5} = 1$$

Correlation Coefficient: Definition

Goal: ① find $\bar{x}, s_x, \bar{y}, s_y$
② find z_{x_i}, z_{y_i}
③ Multiply z-scores
④ "Average" them

Correlation Coefficient r

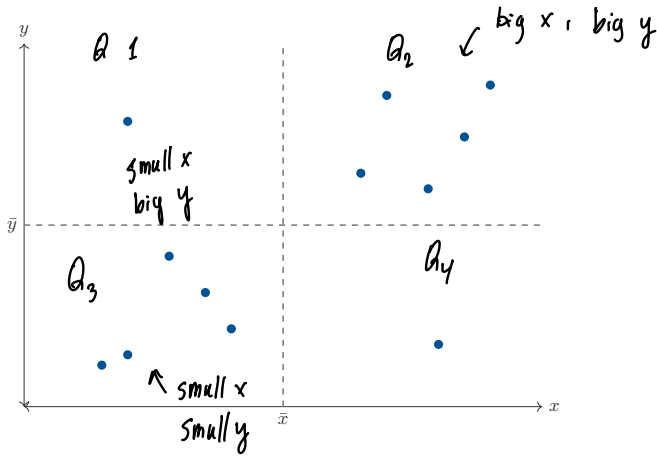
"average" \downarrow z_x z_y

$$r = \text{corr}(x, y) = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

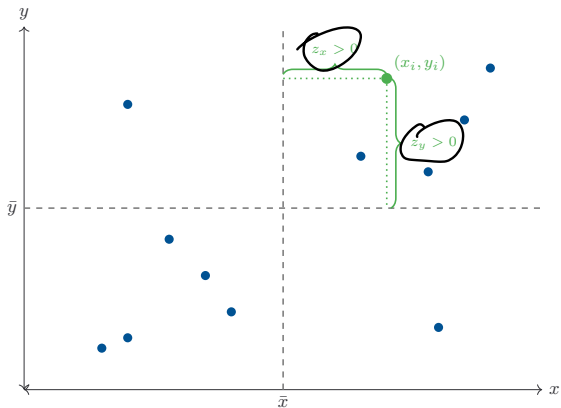
The correlation r is the "average" of the products of the **z-scores** for x and y .

When we need to emphasize the variables involved, we will use the notation $\text{corr}(x, y)$.

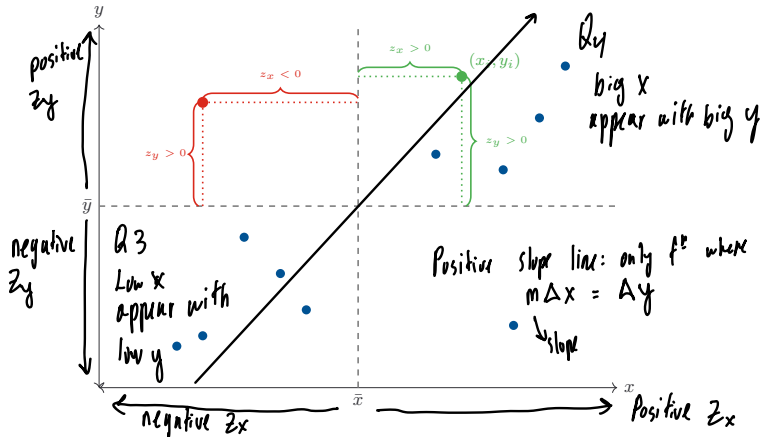
The Idea Behind Correlation



The Idea Behind Correlation

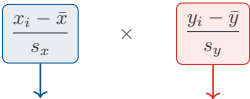


The Idea Behind Correlation



Correlation Formula Breakdown

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \times \left(\frac{y_i - \bar{y}}{s_y} \right)$$



z-score for x *z-score for y*

Example 1.25: Computing Correlation

The Data:

x	-1	0	1
y	1	3	2

$$n=3$$

Step 1: Compute Means and Standard Deviations

$$\bar{x} = \frac{-1+0+1}{3} = 0$$

$$s_x^2 = \frac{1}{2} \left((-1-0)^2 + (0-0)^2 + (1-0)^2 \right) \\ = \frac{1}{2} (1+0+1) = \frac{2}{2} = 1$$

$$\bar{y} = 2, \quad s_y = 1$$

$$\Rightarrow s_x = \sqrt{s_x^2} = \sqrt{1} = 1$$

$$\bar{x} = 0, \bar{y} = 2, s_x = 1, s_y = 1.$$

$$\begin{aligned}\bar{x} &= \frac{-1+0+1}{3} = 0 & s_x^2 &= \frac{1}{2}((-1-0)^2 + (0-0)^2 + (1-0)^2) \\ & & &= \frac{1}{2}(1+0+1) = \frac{2}{2} = 1 \\ \bar{y} &= 2, & s_y &= 1\end{aligned}$$

Step 2: Compute the z -scores and their products for each data point.

x_i	y_i	$z_x = \frac{x-0}{1} = x$	$z_y = \frac{y-2}{1}$	$z_x z_y$
-1	1	$\frac{-1-0}{1} = \frac{-1}{1} = -1$	$\frac{1-2}{1} = \frac{-1}{1} = -1$	$(-1)(-1) = 1$
0	3	0	1	$0 \cdot 1 = 0$
1	2	1	0	$1 \cdot 0 = 0$
Sum of Products:				1

$$\begin{aligned}\text{Corr}(x, y) &= r \\ &= \frac{1}{n-1} \sum_{i=1}^n z_{x_i} z_{y_i} \\ &\quad \downarrow \quad \downarrow \\ &\quad \frac{x_i - \bar{x}}{s_x} \quad \frac{y_i - \bar{y}}{s_y}\end{aligned}$$

From the table, we calculated $\sum z_x z_y = 1$.

Step 3: Divide the sum of the product of the z -scores by $n - 1$

$$\text{Corr}(x, y) = \frac{1}{3-1} \cdot 1 = \frac{1}{2} \cdot 1 = \boxed{\frac{1}{2}}$$

Example 1.26: Computing Correlation

Find the correlation

x	8	5	5	6
y	3	3	6	8

Step 1: Compute Means and Standard Deviations

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

$$s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$s_y = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\bar{x} = 6, \bar{y} = 5, s_x = \sqrt{2}, s_y = \sqrt{6}$$

Step 2: Compute the z -scores and their products.

Aside: $\sqrt{a} \cdot \sqrt{b}$
 $= \sqrt{a \cdot b}$

x_i	y_i	$z_x = \frac{x-6}{\sqrt{2}}$	$z_y = \frac{y-5}{\sqrt{6}}$	$z_x z_y$
8	3	$\frac{8-6}{\sqrt{2}} = \frac{2}{\sqrt{2}}$	$-\frac{2}{\sqrt{6}}$	$-\frac{4}{\sqrt{12}}$
5	3	$\frac{5-6}{\sqrt{2}} = -\frac{1}{\sqrt{2}}$	$-\frac{2}{\sqrt{6}}$	$\frac{2}{\sqrt{12}}$
5	6	$-\frac{1}{\sqrt{2}}$	$\frac{1}{\sqrt{6}}$	$-\frac{1}{\sqrt{12}}$
6	8	$\frac{0}{\sqrt{2}}$	$\frac{3}{\sqrt{6}}$	$\frac{0}{\sqrt{12}}$
				$-\frac{3}{\sqrt{12}}$

$\leftarrow \frac{(-4 + 2 - 1 + 0)}{\sqrt{12}}$

$$\text{Sum of products} = \frac{-3}{\sqrt{12}}$$

Step 3: Divide by $n - 1$ to get r .

$$\begin{aligned}\text{Corr}(x, y) &= \frac{1}{n-1} \sum z_x z_y = \frac{1}{4-1} \cdot \frac{-3}{\sqrt{12}} \\ &= \frac{1}{3} \cdot \frac{(-3)}{\sqrt{12}} = \frac{-1}{\sqrt{12}} \approx -0.289\end{aligned}$$

Alternative Formulas for Correlation

Alternative Correlation Formulas

$$r = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i) (\sum_{i=1}^n y_i)}{\sqrt{\left[n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 \right] \left[n \sum_{i=1}^n y_i^2 - (\sum_{i=1}^n y_i)^2 \right]}}$$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Using alternative correlation formulas

Problem: Calculate r using all three formulas for: $(2, 3), (4, 7), (6, 5)$

Step 1: Setting up the table

	Data points			Sum
x_i	2	4	6	$\sum x_i = 12$
y_i	3	7	5	$\sum y_i = 15$
x_i^2	4	16	36	$\sum x_i^2 = 56$
y_i^2	9	49	25	$\sum y_i^2 = 83$
$x_i y_i$	6	28	30	$\sum x_i y_i = 64$

Summary:

- $n = 3$
- $\bar{x} = 4, \bar{y} = 5$
- $s_x = 2, s_y = 2$

Using alternative correlation formulas (continued)

	Data points			Sum
x_i	2	4	6	$\sum x_i = 12$
y_i	3	7	5	$\sum y_i = 15$
x_i^2	4	16	36	$\sum x_i^2 = 56$
y_i^2	9	49	25	$\sum y_i^2 = 83$
$x_i y_i$	6	28	30	$\sum x_i y_i = 64$

$$r = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \cdot \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

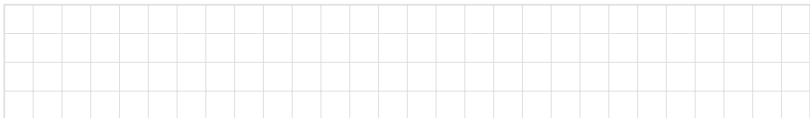
$$r = \frac{5(64) - (12)(15)}{\sqrt{5(56) - (12)^2} \sqrt{5(83) - (15)^2}}$$

Another approach

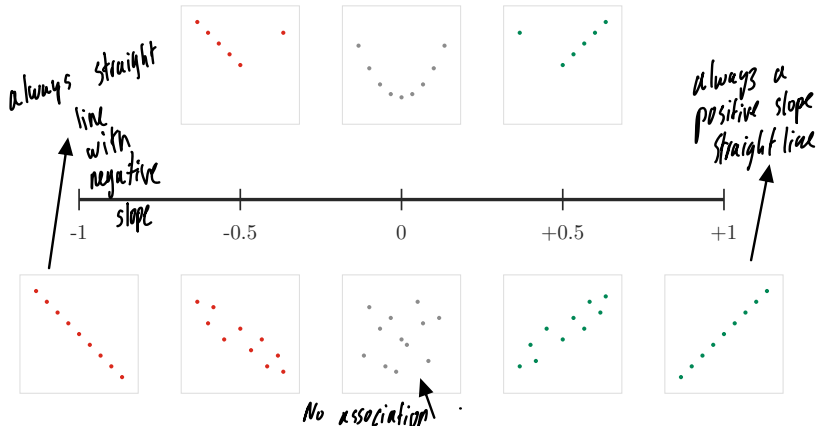
Method 2: Deviation Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}}$$

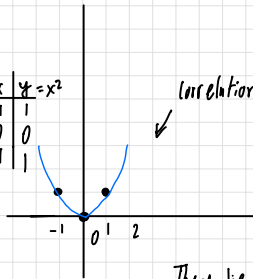
x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
2	3			
4	7			
6	5			
Sum:				



Visualizing Different Values of r

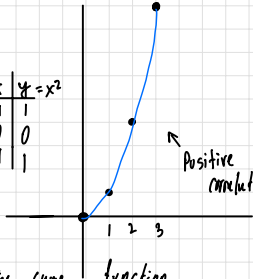


x	y = x ²
-1	1
0	0
1	1



Correlation = 0

x	y = x ²
-1	1
0	0
1	1



Positive correlation

They lie on the same function

Last time

- Z-score of an observation x_i of the variable x is \swarrow quantitative
$$z_i = \frac{x_i - \bar{x}}{s_x}$$

\nwarrow How many stds away an observation is from the mean

- For any x, y quantitative variables

$$\text{Corr}(x, y) = \frac{1}{n-1} \sum_{i=1}^n z_{x_i} z_{y_i} = r.$$

\uparrow
A measure of the strength (and also direction) of the linear relationship b/w x and y .

Property 1: Boundedness

Bounded Range

The correlation coefficient is always between -1 and $+1$:

$$-1 \leq r \leq +1$$

- $r = +1$: Perfect positive linear relationship
- $r = -1$: Perfect negative linear relationship
- $r = 0$: No linear relationship



You will **never** compute a correlation outside this range.

Property 2: Symmetry

Symmetry

The correlation between x and y equals the correlation between y and x :

$$\text{corr}(x, y) = \text{corr}(y, x)$$

- There is no distinction between explanatory and response variables when computing correlation.
- The correlation between “Height” and “Weight” is the same as between “Weight” and “Height.”

Example 1.35: Switching Variables

Problem: A researcher reports that the correlation between “hours of exercise per week” and “resting heart rate” is $r = -0.45$.

Find the correlation between “resting heart rate” and “hours of exercise per week.”

$$\begin{aligned} \text{Corr}(y, x) &= \text{Corr}(x, y) && \text{by symmetry} \\ &= -0.45 \end{aligned}$$

Property 3: Unit Invariance

Scale Invariance

Multiplying either variable by a positive constant does **not** change correlation:

$$\text{corr}(Ax, By) = \text{corr}(x, y) \quad \text{for } A, B > 0$$

↑ ↑
change of units.

Why? Correlation is based on **standardized scores** (z -scores), which have no units.

Example:

- Convert height from inches to centimeters (multiply by 2.54)
- Convert weight from lbs to kilograms (divide by 2.2)
- The correlation remains exactly the same!

Example 1.38: Currency Conversion

Problem: The correlation between “advertising spending (in USD)” and “sales revenue (in USD)” is $r = 0.72$.

A European branch wants to analyze the same data but in Euros (1 USD = 0.85 EUR).

What is $\text{corr}(\text{spending in EUR}, \text{revenue in EUR})$?

$100 \text{ USD} = 0.85 \text{ EUR} (100 \text{ USD})$

Goal: $\text{Corr}(w, z)$

Note that $w = 0.85x$, $z = 0.85y$

$$\begin{aligned}\text{Corr}(w, z) &= \text{Corr}(0.85x, 0.85y) \\ &= \text{Corr}(x, y) \quad \text{by unit in variance}\end{aligned}$$

$\overset{\text{USD}}{=} 0.85 \text{ EUR} \cdot 100$
 $= 85 \text{ EUR}$

Property 4: Translation Invariance

Translation Invariance

Adding a constant to either variable does **not** change correlation:

$$\text{corr}(x + a, y + b) = \text{corr}(x, y)$$

Why? Shifting data up/down or left/right doesn't change the **shape** of the point cloud.

Example: Midterm and final exam scores

Problem: The correlation between midterm scores x and final exam scores y is $r = 0.68$. The professor decides to add 5 bonus points to everyone's midterm and 10 bonus points to everyone's final exam. What is the new correlation?

$$\text{Corr}(x + 5, y + 10) = \text{Corr}(x, y) \quad \text{by translation invariance} \\ = 0.68$$

Combining Invariance Properties

 **Key Point:** Correlation is **invariant under linear transformations**:

$$\text{corr}(Ax + a, By + b) = \text{corr}(x, y) \quad \text{for } A, B > 0$$

This means you can:

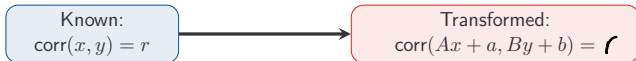
- Change units (Fahrenheit to Celsius, miles to kilometers)
- Shift baselines (add/subtract constants)
- Scale data (multiply by positive constants)

...and correlation **will not change**.

Using Properties to Simplify Calculations

The properties of correlation let us solve problems **without computation**.

Strategy: If you know $\text{corr}(x, y)$, you can immediately deduce correlations for any linear transformation (function) of x or y .



Example 1.39: Temperature Conversion

Problem: A researcher finds that the correlation between daily temperature (in Celsius) and ice cream sales is $r = 0.82$.

What is the correlation if temperature is measured in Fahrenheit instead?

Recall: $F = \frac{9}{5}C + 32$

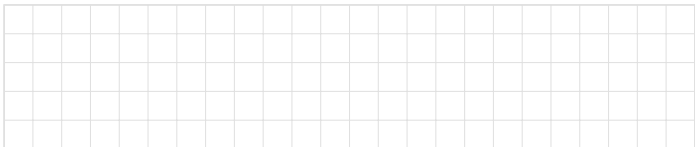
$$\begin{aligned}\text{Corr} \left(\frac{9}{5}x + 32, y \right) &= \text{Corr} \left(\frac{9}{5}x, y \right) && \text{by trans. invariance} \\ &= \text{Corr} (x, y) && \text{by unit invariance}\end{aligned}$$

Example 1.40: Standardized Test Scores

Problem: The correlation between SAT Math scores and first-year GPA is $r = 0.65$.

A new “adjusted SAT” is created: $\text{SAT}_{\text{adj}} = 2 \cdot \text{SAT} - 800$

What is $\text{corr}(\text{SAT}_{\text{adj}}, \text{GPA})$?



Example 1.36: Symmetry in Action

Problem: A study reports that the correlation between years of education and annual income is $r = 0.58$.

A journalist writes: "The correlation between income and education is different from the correlation between education and income."

Is this correct?

Sign Changes Under Negative Scaling

Effect of Negative Multipliers

For $A, B \neq 0$:

$$\text{corr}(Ax, By) = \text{sign}(A) \cdot \text{sign}(B) \cdot \text{corr}(x, y)$$

$\text{sign}(A)$	$\text{sign}(B)$	Effect on r	Example
+	+	Same sign	$\text{corr}(2x, 3y) = r$
+	-	Flips sign	$\text{corr}(2x, -3y) = -r$
-	+	Flips sign	$\text{corr}(-2x, 3y) = -r$
-	-	Same sign	$\text{corr}(-2x, -3y) = r$

Example 1.41: Negative Scaling Application

Problem: The correlation between “hours of sleep” x and “number of errors on a task” y is $r = -0.60$.

A researcher redefines the variables as:

- “Sleep deficit” = $8 - \text{hours of sleep}$ $w = 8 - x$

- “Accuracy” = $100 - \text{errors}$ $z = 100 - y$

Find $\text{corr}(\text{sleep deficit}, \text{accuracy})$.

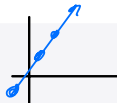
$$-x + 8, -y + 100$$

Goal: Find $\text{corr}(w, z)$

$$\begin{aligned}\text{corr}(w, z) &= \text{corr}(8 - x, 100 - y) \\ &= \text{corr}(-x, -y) \\ &= (-1)(-1) \text{corr}(x, y)\end{aligned}$$

Property 5: Perfect Correlation


Perfect Correlation



- $r = +1$: All points fall **exactly** on a line with **positive slope**.
- $r = -1$: All points fall **exactly** on a line with **negative slope**.



Example: If $y = 3x + 7$, then $r = 1$ because every point lies exactly on the line.


 $\text{sgn}(3)$

Example 1.43: Does Slope Affect Correlation?

Problem: Three datasets have the following exact relationships:

- Dataset A: $y = 10x$ (steep line)
- Dataset B: $y = 2x$ (moderate line)
- Dataset C: $y = 0.5x$ (shallow line)

} All have correlation $r=1$

Which dataset has the highest correlation?

$$\begin{aligned} \text{Corr}(\underbrace{ax+b}_y, x) &= \text{Corr}(ax, x) && \text{(Property 4)} \\ &= \text{sgn}(a) \text{Corr}(x, x) && \text{by trans. invariance} \\ \text{general form of line} &&& \downarrow \\ \text{sgn}(a) &= \begin{cases} 1, & \text{if } a > 0 \\ -1, & \text{if } a < 0 \end{cases} \\ &&& \frac{a}{|a|} \end{aligned}$$

Example 1.44: What About Zero Slope?

Problem: Consider the data points $(1, 5), (2, 5), (3, 5), (4, 5), (5, 5)$.

All points lie exactly on the horizontal line $y = 5$. What is the correlation?

The correlation is undefined as
the z-scores are undefined for y :

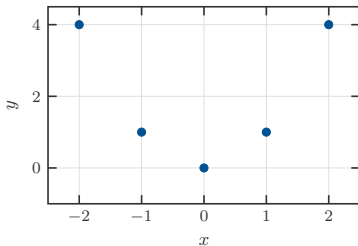
$$s_y = 0 \quad \text{hence} \quad z_{y_i} = \frac{y_i - \bar{y}}{s_y} \\ = \frac{y_i - \bar{y}}{0}$$



Property 6: Correlation Measures Linear Relationships Only

⚠ Caution: Correlation (r) only measures the strength of **linear** relationships.

Check that the correlation of the following points is $r = 0$:



Example 1.65: Test Score Transformations

Problem: The correlation between raw test scores (x) and project grades (y) is $r = 0.75$. The instructor applies these transformations:

- Curved scores: $x_{\text{curved}} = 1.2x + 10$
- Weighted projects: $y_{\text{weighted}} = 0.8y$

(a) What is $\text{corr}(x_{\text{curved}}, y_{\text{weighted}})$?

(b) What is $\text{corr}(y, x)$?

$$\begin{aligned} \text{a) } \text{Corr}(x_c, y_w) &= \text{Corr}(1.2x + 10, 0.8y) \\ &= \text{Corr}(1.2x, 0.8y) \\ &= \text{Corr}(x, y) = 0.75 \end{aligned}$$

translation invariance
change of units

$$\text{b) } \text{Corr}(y, x) = \text{Corr}(x, y) = 0.75$$

Example 1.65: Test Score Transformations

Problem: The correlation between raw test scores (x) and project grades (y) is $r = 0.75$. The instructor applies these transformations:

- Curved scores: $x_{\text{curved}} = 1.2x + 10$
- Weighted projects: $y_{\text{weighted}} = 0.8y$

(a) What is $\text{corr}(x_{\text{curved}}, y_{\text{weighted}})$?

(b) What is $\text{corr}(y, x)$?

$$\begin{aligned} \text{a) } \text{Corr}(x_c, y_w) &= \text{Corr}(1.2x + 10, 0.8y) \\ &= \text{Corr}(1.2x, 0.8y) \\ &= \text{Corr}(x, y) = 0.75 \end{aligned}$$

translation invariance
change of units

$$\text{b) } \text{Corr}(y, x) = \text{Corr}(x, y) = 0.75$$

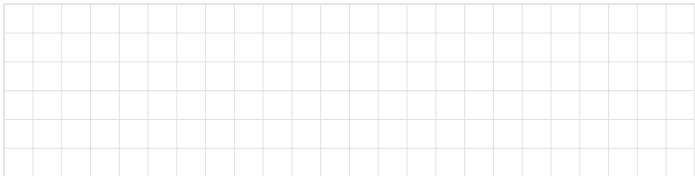
Example 1.66: Health Metrics

Problem: A health study finds $\text{corr}(\text{BMI}, \text{blood pressure}) = 0.42$.

Define new variables:

- “BMI deficit” = $25 - \text{BMI}$ (how far below “healthy” BMI)
- “Blood pressure in kPa” = $\text{BP in mmHg} \times 0.133$

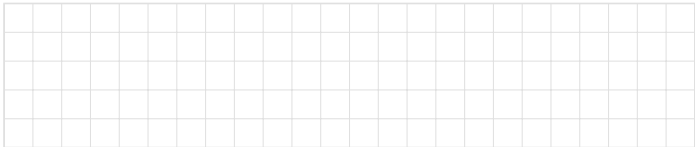
Find $\text{corr}(\text{BMI deficit}, \text{BP in kPa})$.



Example 1.67: True or False?

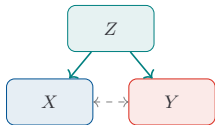
Problem: Identify which claims are **correct** or **incorrect**.

1. "If $r = 0$, then x and y are unrelated."
2. "The correlation between age and height is different if I measure height in feet vs. meters."
3. "If all points lie on the line $y = -5x + 100$, then $r = -1$."
4. "A steeper regression line means a higher correlation."



Correlation Does Not Imply Causation

⚠ Caution: Just because two variables are correlated does **not** mean one causes the other.



Confounding



Reverse



Coincidence

Example 1.47: Confounding Example

Observation: Ice cream sales and drowning deaths are positively correlated.

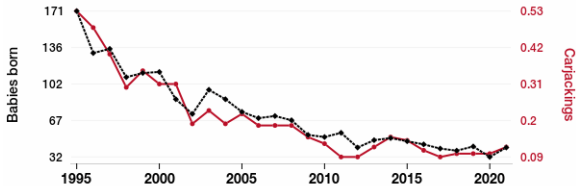
Does ice cream cause drowning?

Example 1.50: Spurious Correlations

Popularity of the first name Alix

correlates with

Carjackings in the US



◆--- Babies of all sexes born in the US named Alix · Source: US Social Security Administration

●— Rate of nonfatal carjacking victimization per 1,000 persons age 16 or older (3-year moving averages) · Source: Bureau of Justice Statistics

1995-2021, $r=0.968$, $r^2=0.937$, $p<0.01$ · tylervigen.com/spurious/correlation/5912

Properties of Correlation: Summary

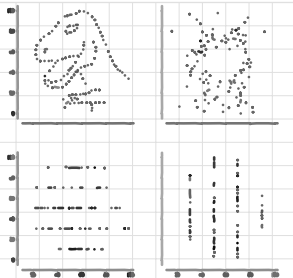
Property	What It Means
1. Boundedness	$-1 \leq r \leq +1$
2. Symmetry	$\text{corr}(x, y) = \text{corr}(y, x)$
3. Unit Invariance	$\text{corr}(Ax, By) = A B \text{corr}(x, y)$ for $A, B > 0$
4. Translation Invariance	Adding constants doesn't change r
5. Perfect Correlation	$r = \pm 1$ means points lie exactly on a line
6. Linearity Only	r only detects linear patterns

Pitfalls of Summary Statistics

All of these scatterplots have the same correlation.

Correlation only tells us the strength (and direction) of the **linear** relationship.

Takeaway: whenever possible, visualize the data



```
Group 1:  
dino: mean_x=54.26, mean_y=47.83, sd_x=16.77, sd_y=26.94, corr=-0.06  
away: mean_x=54.27, mean_y=47.83, sd_x=16.77, sd_y=26.94, corr=-0.06  
h_lines: mean_x=54.26, mean_y=47.83, sd_x=16.77, sd_y=26.94, corr=-0.06  
v_lines: mean_x=54.27, mean_y=47.84, sd_x=16.77, sd_y=26.94, corr=-0.07
```

Chapter 4 Summary

Scatterplots

- Plot two or more quantitative variables on a scatterplot
- Explanatory on x ; Response on y
- Identify: outliers, form, direction, strength

Correlation Formula

z -score:

$$z = \frac{x - \bar{x}}{s_x}$$

$$r = \frac{1}{n-1} \sum \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

Properties of r

- Bounded: $-1 \leq r \leq +1$
- Symmetric: order of variables doesn't matter
- Unit-free: invariant to (positive) scaling and shifting
- Measures **linear** relationships only

Caution

- Correlation \neq Causation