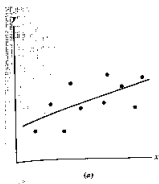## Correlation

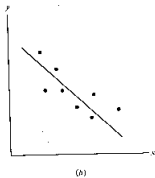*Goal:* Find cause and effect links between variables.

What can we conclude when two variables are highly **correlated**?



| **Positive Correlation** | **Negative Correlation** |
|:---:|:---:|
| High values of $x$ | High values of $x$ |
| are associated with | are associated with |
| high values of $y$. | low values of $y$. |

The **correlation coefficient,** $R^2$ is a number between 0 and 1. Values near 1 show strong correlation; values near 0 show weak correlation.
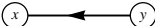
## Causation

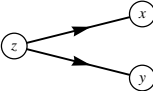If we have high correlation, we'd like to determine causation.

To visually represent the direction of causality between variables, use arrows. For example, if $x$ causes $y$, we draw an arrow from $x$ to $y$.

The ways in which two variables may have strong correlation are:

  I. Simple Causality    $x \longrightarrow y$

 II. Reverse Causality    $x \longleftarrow y$

III. Mutual Causality    $x \rightleftarrows y$

IV. Hidden/Confounding Variable    $z \begin{smallmatrix} \nearrow x \\ \searrow y \end{smallmatrix}$

 V. Complete Accident/Coincidence    $x$       $y$

# Simple Causality

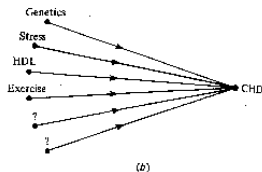I. Simple Causality $\quad x \longrightarrow y$

We say that variables $x$ and $y$ are related by **simple causality** if the level of $x$ *determines* the level of $y$.

Example 2 (pp. 171–173) deals with high blood pressure. After plotting blood pressure ($x$) with deaths from heart disease ($y$), there is high correlation.

A chain of causation can be deduced that makes the argument for simple causality:

high blood pressure $\rightarrow$ arteries clog $\rightarrow$ lack of oxygen in heart $\rightarrow$ heart disease

Many factors have been determined that increase the chance for heart disease.



(b)

# Reverse Causality

II. Reverse Causality    $x$ ⟵ $y$

We say that variables $x$ and $y$ are related by **reverse causality** if the level of $x$ *is determined by* the level of $y$.

*Example.* Islanders in South Pacific deter-
mined that healthy people had body lice and
sick people didn't. The islanders concluded
that more body lice means better health.
However, everyone had lice and lice prefer
healthy hosts.



Figure 7

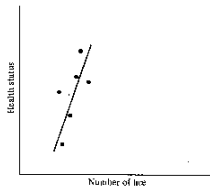*Example.* Human birth rate and
stork population: "storks bring babies".

## Mutual Causality / Feedback

III. Mutual Causality   $x \rightleftarrows y$

We say that variables $x$ and $y$ are related by **mutual causality** if changes in $x$ produce changes in $y$ and vice versa.
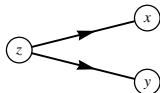
*Example.*  Car dealers:
If you plot car sales and advertising budget
for a large set of car dealers, you will likely
find a strong correlation.

Do car sales pay for advertising
or does advertising drive sales?

They are mutually reinforcing,
so this is an example of mutual causality.

## Hidden Variable Causes Both

IV. Hidden/Confounding Variable 

We say that $x$ and $y$ are in a **spurious relationship** if the levels of both $x$ and $y$ are determined by the level of a **confounding variable** $z$.

*Example.* In a city, the number of churches there are is highly correlated with the number of liquor stores.

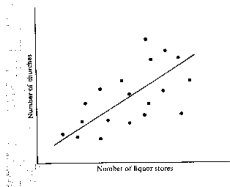▶ Simple causation would imply:

▶ Reverse causation would imply:



Figure 10

In this instance, there is a confounding variable: _____ .

## Complete Accident

V. Complete Accident/Coincidence  $(x)$  $(y)$

If none of the above four cases apply, $x$ and $y$ are unrelated.

Take two dice. Roll each five times. Plot the value of one die versus the value of the other die for the five rolls. Often there will be no correlation.
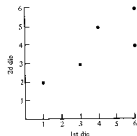
One instance of correlation occurred, with an $R^2$ of 0.672 (relatively high!)



| Roll | Number showing | |
| | 1st die | 2d die |
| #1 | 1 | 2 |
| #2 | 3 | 3 |
| #3 | 4 | 5 |
| #4 | 5 | 4 |
| #5 | 6 | 6 |

An example of a correlation by coincidence.

*Example.* Perhaps with students and SSN's?

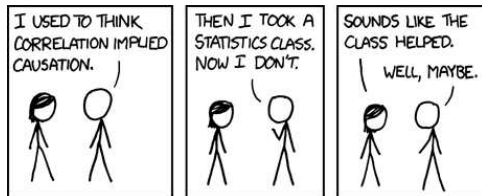▶ The chance of this occurring decreases as more observations are taken.

## Correlation does not imply causation!

*Groupwork:* Justify the correlations between the following variables:

▶ As ice cream sales increase, the rate of drowning deaths increase.

▶ The more firemen fighting the fire, the larger the fire grows.

▶ With fewer pirates on the open seas, global warming has increased.

▶ The more people in my Facebook group, the faster it grows.

*What is the joke below?*



Source: http://xkcd.com/552/

# Calculating the $R^2$ Statistic

The **correlation coefficient** $R^2$ ("**R-Squared**") is a value between 0 and 1 which helps measure the goodness of fit of a *linear regression*.

To calculate $R^2$, you need to calculate:

▶ The **error sum of squares**: $SSE = \sum_i \left[ y_i - f(x_i) \right]^2$.

⋆ *SSE* is the variation between the data and the regression line. ⋆

▶ The **total corrected sum of squares**: $SST = \sum_i \left[ y_i - \bar{y} \right]^2$, where $\bar{y}$ is the average $y_i$ value.

⋆ *SST* is the variation solely due to the data. ⋆

▶ Now calculate $R^2 = 1 - \frac{SSE}{SST}$.

⋆ $R^2$ is the proportion of variation explained by the line. ⋆

$R^2$ near 0 $\Rightarrow$ low correlation.          $R^2$ near 1 $\Rightarrow$ high correlation.

# Calculating the $R^2$ Statistic

*Example.* (cont. from notes p. 24) What is the correlation coefficient of the data set: $\{(1.0, 3.6), (2.1, 2.9), (3.5, 2.2), (4.0, 1.7)\}$?

Recall that the regression line is $f(x) = -0.605027x + 4.20332$.

▶ The **error sum of squares**: $SSE = \sum_i \left[y_i - f(x_i)\right]^2$.

$SSE = (3.6 - f(1.0))^2 + (2.9 - f(2.1))^2 + (2.2 - f(3.5))^2 + (1.7 - f(4.0))^2$
$\quad = (.0017)^2 + (-0.033)^2 + (0.114)^2 + (-0.083)^2 = 0.0210$

▶ The **total corrected sum of squares**: $SST = \sum_i \left[y_i - \bar{y}\right]^2$.

First calculate $\bar{y} = (3.6 + 2.9 + 2.2 + 1.7)/4 = 2.6$

$SST = (3.6 - 2.6)^2 + (2.9 - 2.6)^2 + (2.2 - 2.6)^2 + (1.7 - 2.6)^2$
$\quad = (1)^2 + (0.3)^2 + (-0.4)^2 + (-0.9)^2 = 2.06$

▶ Now calculate $R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{0.0210}{2.06} = 1 - .01 = 0.99$.