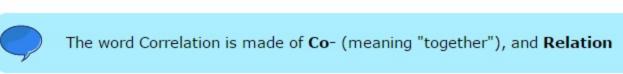
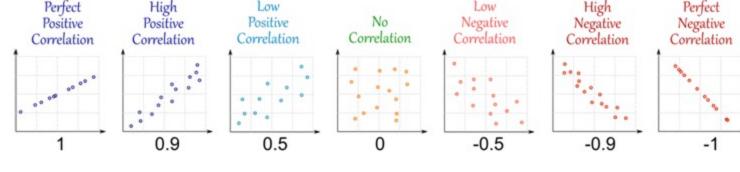
Correlation

When two sets of data are strongly linked together we say they have a High Correlation.



- Correlation is Positive when the values increase together, and
- Correlation is Negative when one value decreases as the other increases

Here we look at linear correlations (correlations that follow a line).



1 is a perfect positive correlation

negative.

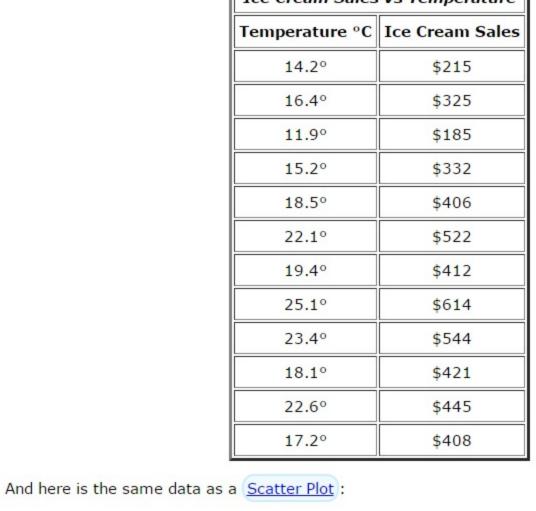
Correlation can have a value:

- 0 is no correlation (the values don't seem linked at all) -1 is a perfect negative correlation
- The value shows **how good the correlation is** (not how steep the line is), and if it is positive or

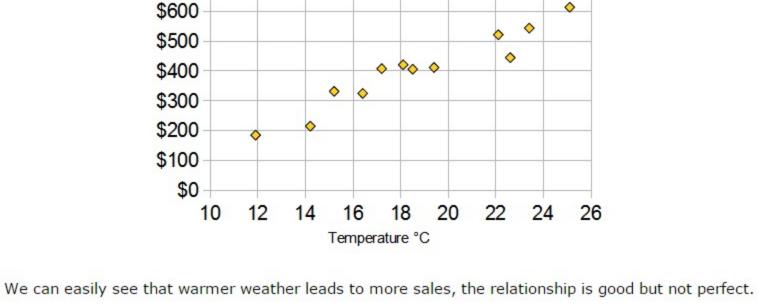
Example: Ice Cream Sales

The local ice cream shop keeps track of how much ice cream they sell versus the temperature on that day, here are their figures for the last 12 days:

Ice Cream Sales vs Temperature



\$100



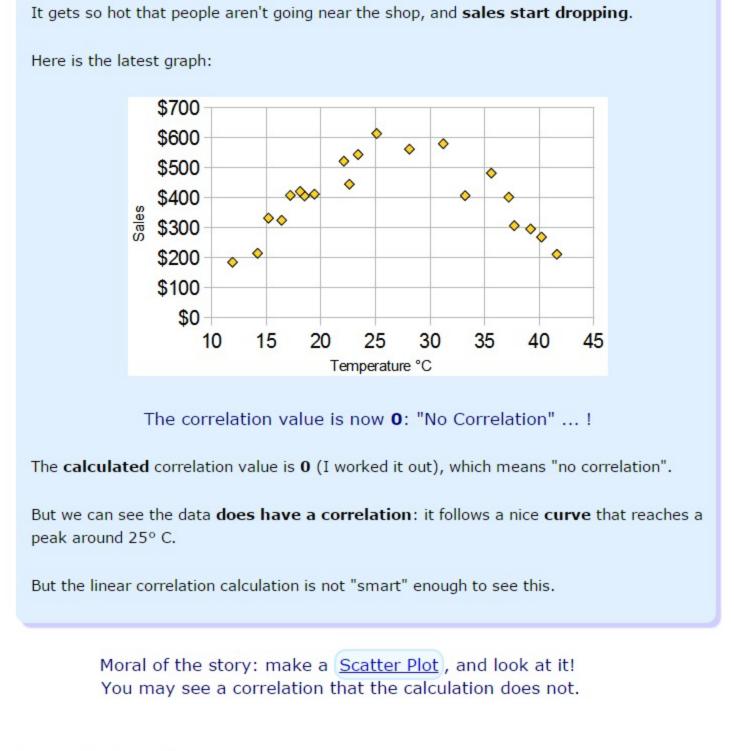
Correlation Is Not Good at Curves

In fact the correlation is **0.9575** ... see at the end how I calculated it.

\$700

Our Ice Cream Example: there has been a heat wave!

The correlation calculation only works well for relationships that follow a straight line.



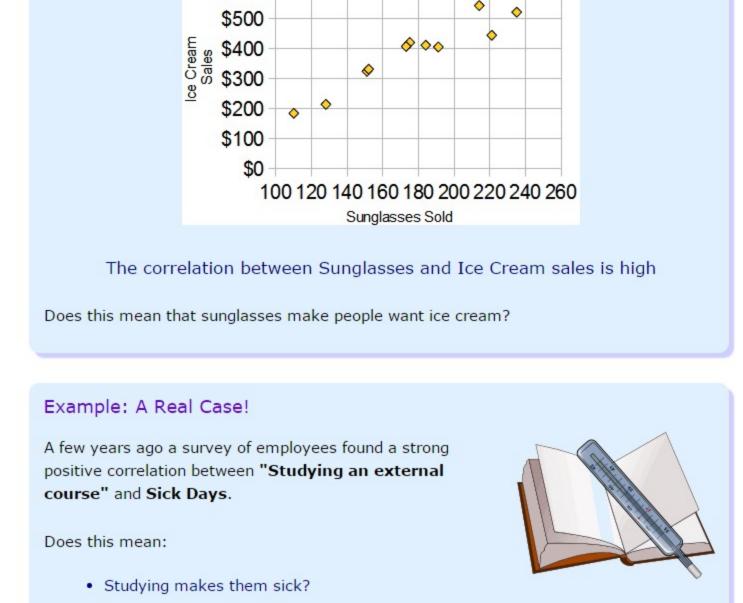
"Correlation Is Not Causation" ... which says that a correlation does not mean that one thing causes the other (there could be other reasons the data has a good correlation).

Example: Sunglasses vs Ice Cream Our Ice Cream shop finds how many sunglasses were sold by a big store for each day and

Correlation Is Not Causation

compares them to their ice cream sales: \$700

\$600



in Excel or LibreOffice Calc but here is how to calculate it yourself:

them "b")

How To Calculate

Let us call the two sets of data "x" and "y" (in our case Temperature is x and Ice Cream Sales is y):

How did I calculate the value **0.9575** at the top?

Sick people study a lot?

Without further research we can't be sure why.

· Or did they lie about being sick to study more?

• Step 3: Calculate: a × b, a² and b² for every value • Step 4: Sum up $\mathbf{a} \times \mathbf{b}$, sum up $\mathbf{a^2}$ and sum up $\mathbf{b^2}$

-\$77

\$212

\$142

Calculate ab, a² and b²

20.3

46.2

12.3

11.6

41.0

22.1

0.5

0.0

5.3

34,969

5,929

47,089 4,900

14,400

44,944

20,164

16

100

a×b

842

177

245

408

1,357

667

-1

7

1,476

Step 5: Divide the sum of a × b by the square root of [(sum of a²) × (sum of b²)]

Here is how I calculated the first Ice Cream example (values rounded to 1 or 0 decimal places):

• Step 2: Subtract the mean of x from every x value (call them "a"), do the same for y (call

I used "Pearson's Correlation". There is software that can calculate it, such as the CORREL() function

Subtract Mean "a" "b" Temp °C Sales -4.514.2 \$215 -\$187

\$325

\$614

\$544

16.4

25.1

23.4

. Step 1: Find the mean of x, and the mean of y

11.9 \$185 -\$217 15.2 \$332 -3.5-\$70 18.5 \$406 -0.2\$4 3.4 22.1 \$522 \$120 19.4 \$412 0.7 \$10

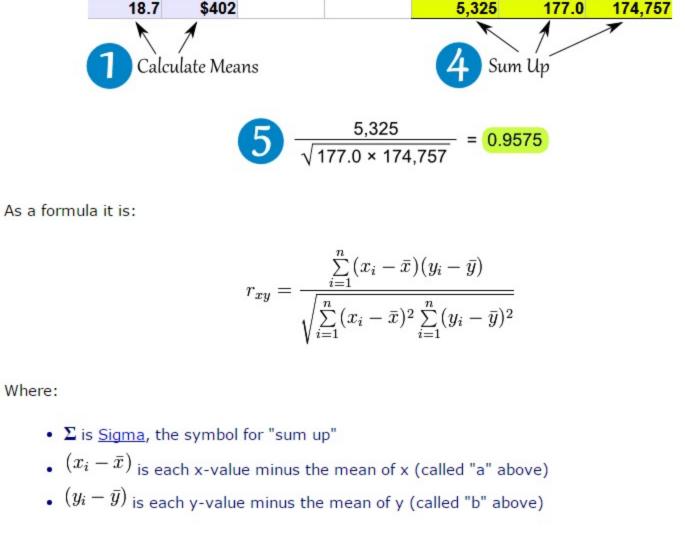
-2.3

-6.8

6.4

4.7

-0.6\$19 -11 361 18.1 \$421 0.4 22.6 3.9 168 15.2 \$445 \$43 1,849 17.2 -9 2.3 \$408 -1.5



You probably won't have to calculate it like that, but at least you know it is not "magic", but simply

You can calculate it in one pass through the data. Just sum up x, y, x^2 , y^2 and xy (no need for a or b calculations above) then use the formula:

Note for Programmers

a routine set of calculations.

Your turn:

Where:

 $r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$ Other Methods

There are other ways to calculate a correlation coefficient, such as "Spearman's rank correlation

Question 1 Question 2 Question 3 Question 4 Question 5

Question 6 Question 7 Question 8 Question 9 Question 10

coefficient", but I prefer using a spreadsheet like above.

How do you get from

now to next?

STORAGE FOR WHAT'S NEXT