



## M6-02: Correlation

Part of the "Towards Machine Learning" Learning Badge

Video Walkthrough: <https://discovery.cs.illinois.edu/m6-02/>

**The correlation coefficient has no units since it is calculated from standard units (z-scores) and is not affected by a change of scale:**

- Adding a constant to all  $X$  or  $Y$  values does NOT change  $r$ .
- Multiplying all  $X$  or  $Y$  values by a positive constant does NOT change  $r$ . What does multiplying by a negative constant do?
- Interchanging all  $X$  and  $Y$  values does not change  $r$ . (The correlation between height and weight is the same as the correlation between weight and height).
- Changing units does NOT change  $r$ . (So the correlation between height in inches and weight in pounds is the same as between height in meters and weight in kilograms.)

### Correlation is NOT Causation

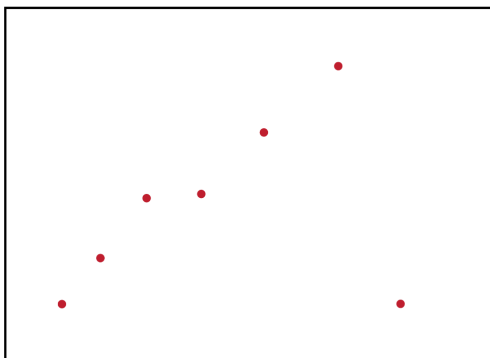
Correlation only measures the association between two variables. Correlation does not imply causation. Though the correlation between the weight and the math ability of children in a school district may be positive, that does not mean that doing math makes children heavier or that putting on weight improves the children's math skills. Age is a confounding variable: older children are both heavier and better at math than younger children, on average.

**Outliers can have a strong effect on the correlation coefficient,  $r$ .**

Outliers should only be excluded for good reason. The correlation coefficient should be used with caution when there are outliers. Outliers can be typos, lies, real data, etc.

**Puzzle:** For which plot below does the outlier raise the correlation coefficient and for which plot does it lower the correlation coefficient?

Plot A



Plot B

