

## Variability in Statistics

**Range:** In statistics, the range is the smallest of all dispersion measures. It is the difference between the distribution's ~~two~~ extreme conclusions. In other words, the range is the difference between the distribution's maximum and minimum observations.

$$\text{Range} = X_{\max} - X_{\min}$$

where  $X_{\max}$  represent the largest observation and  $X_{\min}$  represent the smallest observation of the variable values.

## Percentiles, Quartiles and Interquartile Range (IQR)

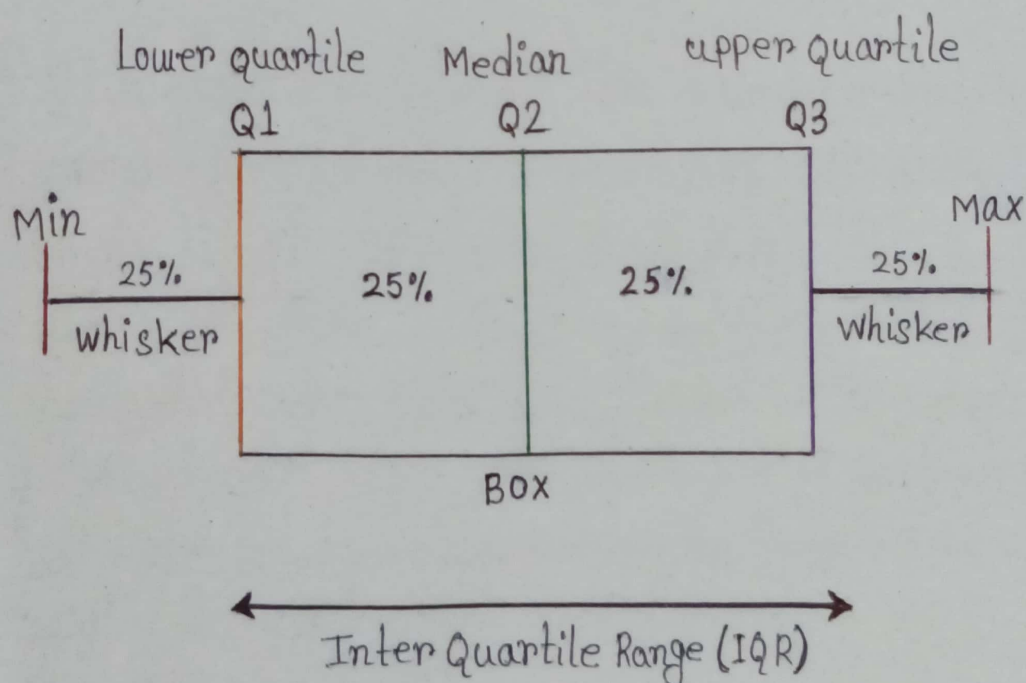
- **Percentiles** - It is a statistician's unit of measurement that indicates the value below which a given percentage of observations in a group of observation fall.

For instance, the value  $Q_x$  represents the 40th percentile of  $xx$  (0.40)

- **Quartiles** - values that divide the number of data points into four more or less equal parts, or quarters. Quartiles are the 0th, 25th, 50th, 75th, and 100th percentile values or the 0th, 25th, 50th, 75th, and 100th percentile values.

• Interquartile Range (IQR) - The difference between the third and first quartiles is defined by the interquartile range. The partitioned values that divide the entire series into four equal parts are known as quartiles. So, there are three quartiles. The first quartile, known as the lower quartile, is denoted by  $Q_1$ , the second quartile by  $Q_2$ , and third quartile by  $Q_3$ , known as the upper quartile. As a result, the interquartile range equals the upper quartile minus the lower quartile.

$$\begin{aligned} \text{IQR} &= \text{Upper Quartile} - \text{Lower Quartile} \\ &= Q_3 - Q_1 \end{aligned}$$





- **Variance** - The dispersion of a data collection is measured by variance. It is defined technically as the average of squared deviations from the mean.

Population Variance	Sample Variance
$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$	$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$
$\sigma^2$ = Population Variance $x_i$ = value of $i^{\text{th}}$ element $\mu$ = population mean $N$ = Population Size	$s^2$ = Sample variance $x_i$ = value of $i^{\text{th}}$ element $\bar{x}$ = Sample mean $n$ = Sample Size

- **Standard Deviation** - The standard deviation is a measure of data dispersion WITHIN a single sample selected from the study population. The square root of the variance is used to compute it. It simply indicates how distant the individual values in a sample are from the mean. To put it another way, how dispersed is the data from the sample? As a result, it is a sample statistic.

### Standard Deviation Formula

Population	Sample
$\sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}}$ <p>X - The value in the data distribution <math>\mu</math> - The population mean N - Total Number of Observations</p>	$s = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$ <p>X - The value in the data distribution <math>\bar{X}</math> - The Sample Mean n - Total Number of Observations</p>

- **Standard Error (SE)** - The standard error indicates how close the mean of any given sample from that population is to the true population mean. When the standard error rises, implying that the means are more dispersed, it becomes more likely that any given mean is an inaccurate representation of the true population mean. When the sample size is increased, the standard error decreases - as the sample size approaches the true population size, the sample means cluster more around the true population mean.

### Standard Error Formula

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$



## Relationship Between variables

- **Causality:** The term "causation" refers to a relationship between two events in which one is influenced by the other. There is causality in statistics when the value of one event, or variable, grows or decreases as a result of other events. Each of the events we just observed may be thought of as a variable, and as the number of hours worked grows, so does the amount of money earned. On the other hand, if you work fewer hours, you will earn less money.

Statistic assesses how much - and how far - the variables change in tandem. To put it another way, it's a measure of the variance between two variables. The metric, on the other hand, does not consider the interdependence of factors. Any positive or negative value can be used for the variance.

### The following is how the values are interpreted:

- **Positive covariance:** When two variables move in the same direction, this is called positive covariance.
- **Negative Covariance:** indicates that two variables are moving in opposite directions.

Single observed value of dependent variable

mean of all values of independent variable

Single observed value of independent variable

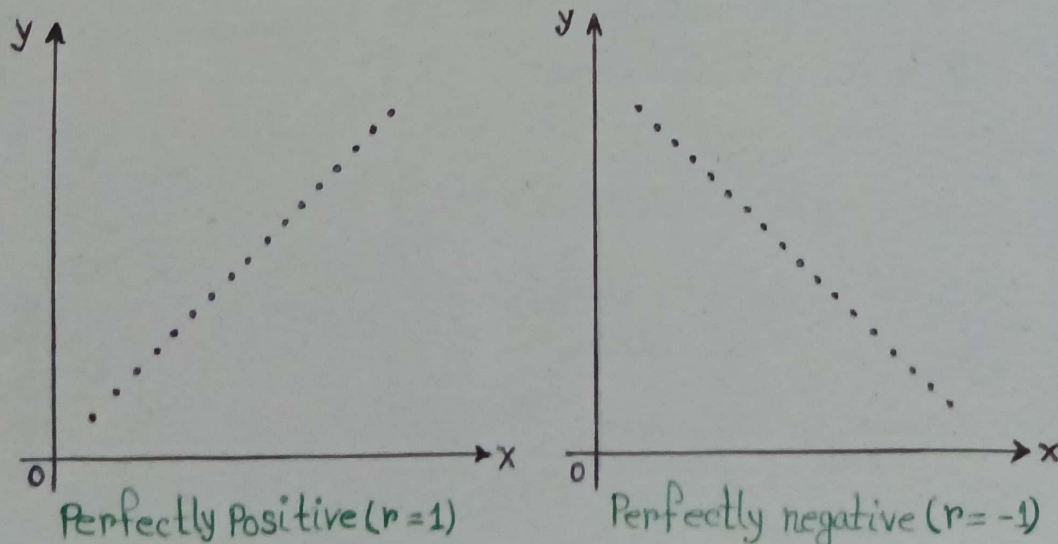
total count of sample values

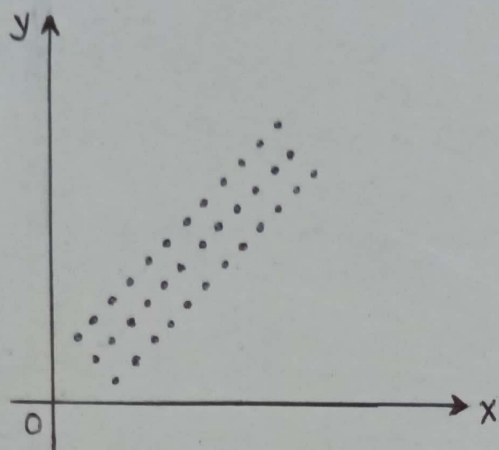
$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

mean of all values of independent variable

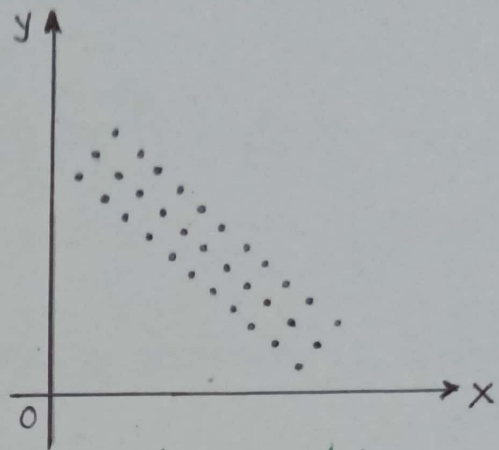
Population count minus one (Bessel's Correction)

- **Correlation**: Correlation is a statistical method for determining whether or not two quantitative or categorical variables are related. To put it another way, it's a measure of how things are connected. Correlation analysis is the study of how variables are connected.

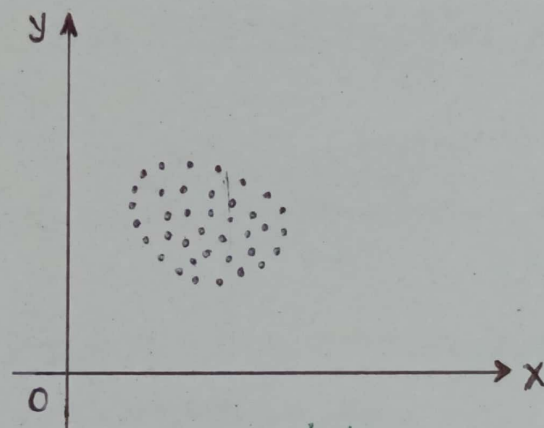




Positive correlation



Negative correlation



Uncorrelation

- Pearson's product moment correlation coefficient ( $r$ ) :

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}, \text{ where } \text{Cov}(x, y) = \frac{1}{n} \sum (x - \bar{x})(y - \bar{y})$$

$\sigma_x$  = standard Deviation of  $x$

$\sigma_y$  = standard Deviation of  $y$



## • Rank Correlation :

- Spearman's rank correlation coefficient (R)

$$R = 1 - \frac{6 \sum D^2}{N^3 - N}$$

where  $D$  = difference of the ranks of an individual in the two characters and

$N$  = number of individuals in the group.

- Coefficient of rank correlation (R) in case of 'tied rank's'

$$R = 1 - \frac{6 \left[ \sum D^2 + \sum \frac{t^3 - t}{12} \right]}{N^3 - N}$$

where  $t$  = number of individuals involved in a tie in any of the two series.

Correlations are useful because they allow you to forecast future behaviour by determining what relationship variables exist. In the practical field, such as government and healthcare, knowing what the future holds is critical. Budgets and company plans are also based on these facts.