

# What Is Statistics?

Statistics is a branch of mathematics that deals with collecting, analyzing, interpreting, and visualizing empirical data. Descriptive statistics and inferential statistics are the two major areas of statistics. Descriptive statistics are for describing the properties of sample and population data (what has happened). Inferential statistics use those properties to test hypotheses, reach conclusions, and make predictions (what can you expect).

- ## Use of Statistics in Machine Learning

1. Asking questions about the data
2. Cleaning and preprocessing the data
3. Selecting the right features
4. Model evaluation
5. Model prediction

- ## Population and Sample

- ### Population: ( $N$ )

In statistics, the population comprises all observations (data points) about the subject under study.

An example of a population is studying the voters in an election. In the 2019 Lok Sabha elections, nearly 900 million voters were eligible to vote in 543 constituencies.

- ### Sample: ( $n$ )

In statistics, a sample is a subset of the population. It is a small portion of the total observed population.

An example of a sample is analyzing the first-time voters for an opinion poll.

## • **Measures of Central Tendency**

Measures of central tendency are the measures that are used to describe the distribution of data using a single value. Mean, Median and Mode are the three measures of central tendency

### Mean:

The arithmetic mean is the average of all the data points.

If there are n number of observations and  $x_i$  is the ith observation, then mean is:

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

Example :-

Consider the data frame below that has the names of seven employees and their salaries.

	Name	Salary
0	Jane	50000
1	Michael	54000
2	Willian	50000
3	Rosy	189000
4	Hana	55000
5	Ferdie	40000
6	Graeme	59000

To find the mean or the average salary of the employees, you can use the mean() functions in Python.

```
print(df['Salary'].mean())
71000.0
```

Solution :-

$$50000 + 54000 + 50000 + 189000 + 55000 + 40000 + 59000 \\ = 497000$$

There are 7 employees.

$$\text{Mean Salary} = \frac{497000}{7} = 71000$$

Ans:

## Median:

Median is the middle value that divides the data into two equal parts once it sorts the data in ascending order.

If the total number of data points (n) is odd, the median is the value at position  $(n+1)/2$ .

When the total number of observations (n) is even, the median is the average value of observations at  $n/2$  and  $(n+2)/2$  positions.

The median() function in Python can help you find the median value of a column. From the above data frame, you can find the median salary as:

```
print(df['Salary'].median())
54000.0
```

Question :-

Consider the data frame below that has the names of seven employees and their salaries.

	Name	Salary
0	Jane	50000
1	Michael	54000
2	Willian	50000
3	Rosy	189000
4	Hana	55000
5	Ferdie	40000
6	Graeme	59000

Solution

**Step 1: Sort salaries**

[40000, 50000, 50000, 54000, 55000, 59000, 189000]

**Step 2: Find the middle value (since there are 7 values)**

The 4th value (middle one) = 54000

**Median Salary = 54000**

$$\Rightarrow \left( \frac{7+1}{2} \right) \Rightarrow 4^{\text{th}} \text{ value} =$$

## Mode:

The mode is the observation (value) that occurs most frequently in the data set. There can be over one mode in a dataset.

The mode salary from the data frame can be calculated as:

```
print(df['Salary'].mode())
0    50000
dtype: int64
```

Question

Given below are the heights of students (in cm) in a class:

155, 157, 160, 159, 162, 160, 161, 165, 160, 158

Mode = 160 cm.

## Outliers

Outliers are values "outside" the other values:

```
99,86,87,88,111,86,103,87,94,78,300,85,86
```

Outliers can change the mean a lot. Sometimes we don't use them (they might be an error), or we use the median or the mode instead.

# # Descriptive v/s Inferential Statistics

## • Descriptive Statistics

**Descriptive Statistics** summarizes (describes) observations from a set of data.

Since we register every newborn baby, we can tell that 51 out of 100 are boys.

From these collected numbers, we can predict a 51% chance that a new baby will be a boy.

It is a mystery that the ratio is not 50%, like basic biology would predict. We only know that we have had this tilted sex ratio since the 17th century.

### Note

Raw observations are only data. They are not real knowledge.

You use **Descriptive Statistics** to transform raw observations into data that you can understand.

## Descriptive Statistics Measurements

Descriptive statistics are broken down into different measures:

### Tendency (Measures of the Center)

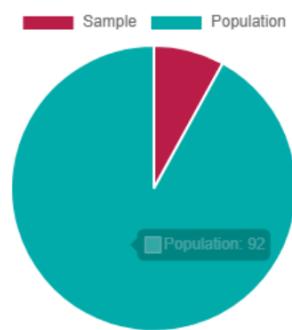
- The Mean (the average value)
- The Median (the mid point value)
- The Mode (the most common value)

### Spread (Measures of Variability)

- Min and Max
- Standard Deviation
- Variance
- Skewness
- Kurtosis

## • Inferential Statistics

**Inferential statistics** are methods for quantifying properties of a population from a small **Sample**:



You take data from a sample and make a prediction about the whole population.

For example, you can stand in a shop and ask a **sample of 100 people** if they like chocolate.

From your research, using inferential statistics, you could predict that 91% of **all shoppers** like chocolate.

## Incredible Chocolate Facts

Nine out of ten people love chocolate.

50% of the US population cannot live without chocolate every day.

You use **Inferential Statistics** to predict whole domains from small samples of data.

→ **Descriptive Statistics** is broken down into **Tendency** and **Variability**.

**Tendency** is about **Center Measures**:

- The **Mean** (the average value)
- The **Median** (the mid point value)
- The **Mode** (the most common value)

**Variability** uses these measures:

- Min and Max
- Variance
- Deviation
- Distribution
- Skewness
- Kurtosis

## The Variance

In statistics, the **Variance** is the average of the squared differences from the **Mean Value**.

In other words, the variance describes how far a set of numbers is **Spread Out** from the mean (average) value.

Mean value is described in the previous chapter.

*Question*

This table contains 11 values:

7	8	8	9	9	9	10	11	14	14	15
---	---	---	---	---	---	----	----	----	----	----

Calculate the Variance:

```
// Calculate the Mean (m)
let m = (7+8+8+9+9+10+11+14+14+15)/11;

// Calculate the Sum of Squares (ss)
let ss = (7-m)**2 + (8-m)**2 + (8-m)**2 + (9-m)**2 + (9-m)**2 + (9-m)**2 + (9-m)**2 + (10-m)**2 + (11-m)**2 + (14-m)**2 + (15-m)**2;

// Calculate the Variance
let variance = ss / 11;
```

To calculate the variance of the Grade, use the following:

```
print(df['Grade'].var())
685.6190476190476
```

The **variance** is a statistical measure that describes how spread out the values in a dataset are. There are **two types of variance** depending on the context:

### 1. Population Variance ( $\sigma^2$ )

Used when you have data for the **entire population**.

**Formula:**

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

- $\sigma^2$ : population variance
- $N$ : total number of data points
- $x_i$ : each individual data point
- $\mu$ : population mean

### 2. Sample Variance ( $s^2$ )

Used when you have a **sample** (subset) of the population.

**Formula:**

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

- $s^2$ : sample variance
- $n$ : number of sample data points
- $x_i$ : each sample data point
- $\bar{x}$ : sample mean

#### Why divide by $(n - 1)$ instead of $n$ ?

This is called **Bessel's correction**. It corrects the bias in estimating the population variance from a sample.

Metric	Formula	Denominator	Used For
Population Variance	$\frac{1}{N} \sum (x_i - \mu)^2$	$N$	Entire population
Sample Variance	$\frac{1}{n-1} \sum (x_i - \bar{x})^2$	$n - 1$	Subset of population

Why divide by  $(n - 1)$  instead of  $n$ ? In sample variance

#### Bessel's Correction (Dividing by $n-1$ ):

- By dividing by  $(n - 1)$  instead of  $n$ , we compensate for this **underestimation**.
- This makes the sample variance an **unbiased estimator** of the population variance.

#### Think of it like this:

- Using the **sample mean** makes the data look **less spread out** than they truly are.
- So we **inflate** the result slightly by dividing by  $(n - 1)$  rather than  $n$ .

# Standard Deviation

**Standard Deviation** is a measure of how spread out numbers are.

The symbol is  $\sigma$  (Greek letter sigma).

The formula is the  $\sqrt{\text{variance}}$  (the square root of the variance).

The Standard Deviation is (in JavaScript):

```
// Calculate the Mean (m)
let m = (7+8+8+9+9+10+11+14+15)/11;

// Calculate the Sum of Squares (ss)
let ss = (7-m)**2 + (8-m)**2 + (8-m)**2 + (9-m)**2 + (9-m)**2 + (9-m)**2 + (10-m)**2 + (11-m)**2 + (14-m)**2 + (15-m)**2;

// Calculate the Variance
let variance = ss / 11;

// Calculate the Standard Deviation
let std = Math.sqrt(variance);
```

**Deviation** is a measure of **Distance**.

**How far** (on average), all values are from **the Mean** (the Middle).

You can find the standard deviation using the `std()` function in Python.

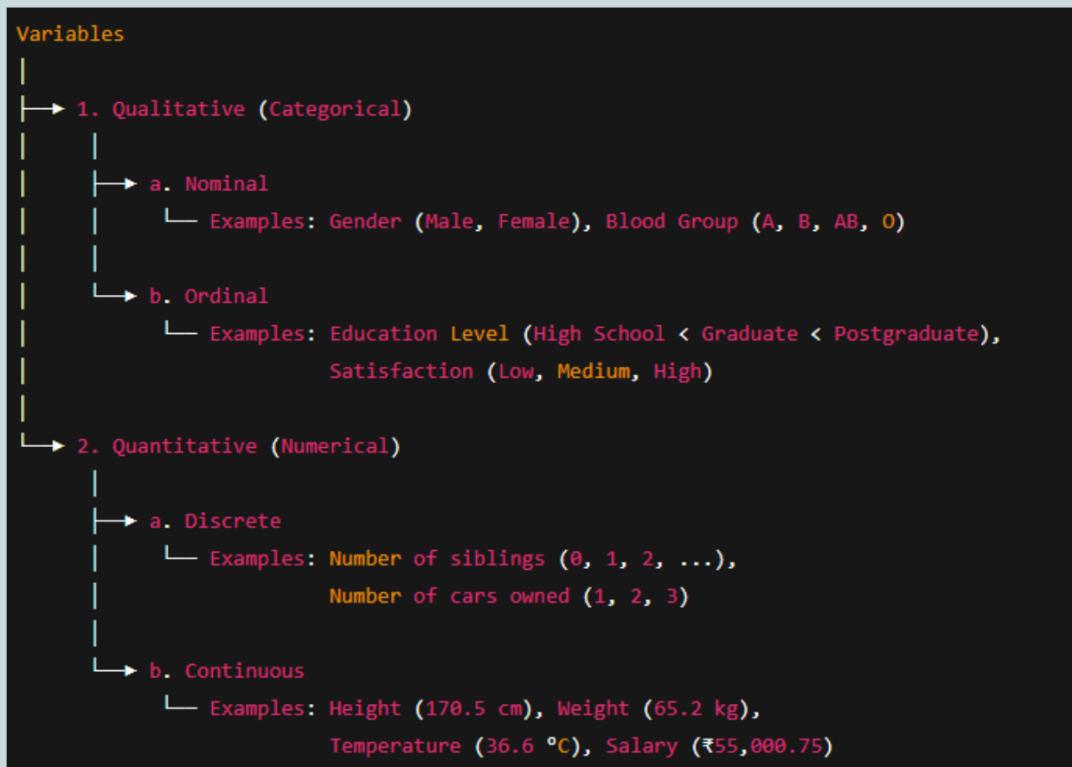
```
print(df['Grade'].std())
26.184328282754315
```

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$$

$$\sigma = \sqrt{\sigma^2} = \sqrt{\text{Variance}}$$

## • Variables

Any characteristic that can be measured or categorized.



## 1. Qualitative Variables (Categorical)

Represent categories or labels; non-numeric.

### a. Nominal

Categories with **no natural order**.

**Example:** Gender (Male, Female), Blood Type (A, B, O)

### b. Ordinal

Categories with a **meaningful order**, but **no fixed difference** between levels.

**Example:** Satisfaction (Low, Medium, High), Education Level (UG < PG < PhD)

## 2. Quantitative Variables (Numerical)

Represent **numeric values**; can be measured or counted.

### a. Discrete

**Countable** values with gaps; usually whole numbers.

**Example:** Number of siblings (0, 1, 2), Number of books

## b. Continuous

Can take **any value within a range**, including decimals.

**Example:** Height (165.3 cm), Temperature (36.6°C), Salary (₹50,000.75)

### • Random variables

A **random variable** is a variable that takes on **numerical values** determined by the **outcome of a random experiment**.

In simple terms:

A random variable maps the outcomes of a random process (like rolling a die or flipping a coin) to numbers.

## Types of Random Variables

### 1. Discrete Random Variable

Takes a **finite or countable** number of values.

**Example:**

Let X be the number shown on a rolled die.

Possible values:  $X=\{1,2,3,4,5,6\}$

### 2. Continuous Random Variable

Can take **any value within a range** (infinite possibilities).

**Example:**

Let Y be the height of a person in cm.

Possible values:  $Y \in [140, 200]$  (like 170.3, 165.8 etc.)

## Example:

**Experiment:** Toss a coin twice

- Sample space: {HH, HT, TH, TT}

- Define a random variable  $X$ : number of heads

Then:

- $HH \rightarrow X = 2$
- $HT \rightarrow X = 1$
- $TH \rightarrow X = 1$
- $TT \rightarrow X = 0$

So,  $X \in \{0, 1, 2\} \rightarrow$  Discrete Random Variable

## • Histogram

A histogram is a graphical representation used to visualize the distribution of a numerical dataset. It shows how frequently data falls within specific ranges or intervals (called bins).

It looks like a bar chart, but it shows the distribution (not categories).

## Key Features of a Histogram:

- X-axis: Intervals (bins) of data values
- Y-axis: Frequency (how many data points fall in each bin)
- Bars touch each other (unlike bar charts)

Example:

Suppose we have test scores of 15 students:

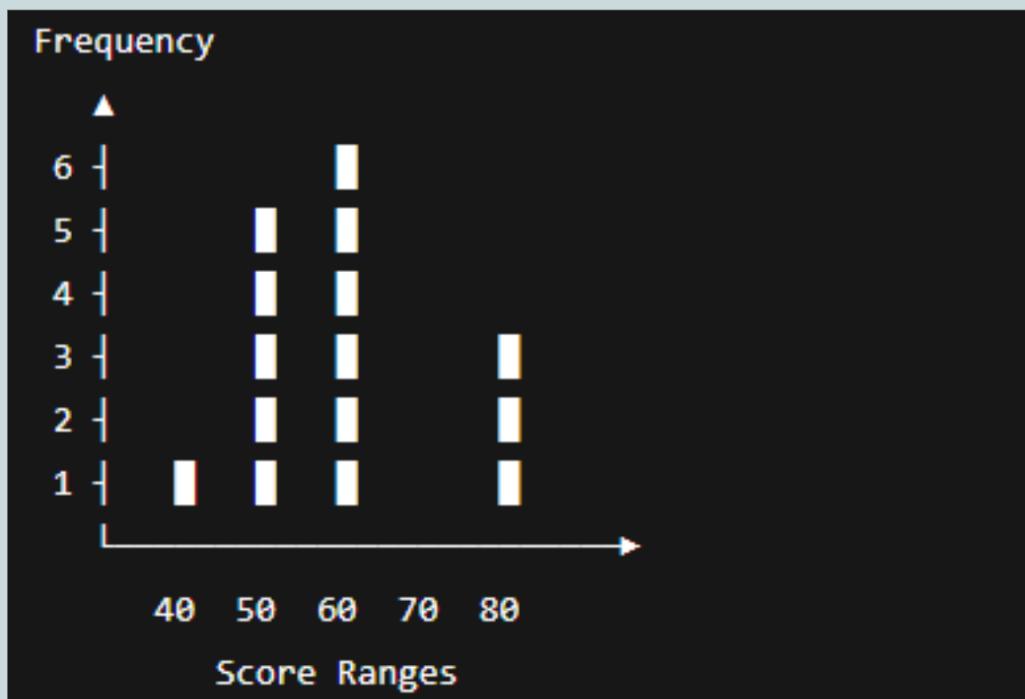
makefile

Scores: 45, 50, 52, 53, 55, 57, 59, 61, 62, 65, 67, 68, 70, 72, 75

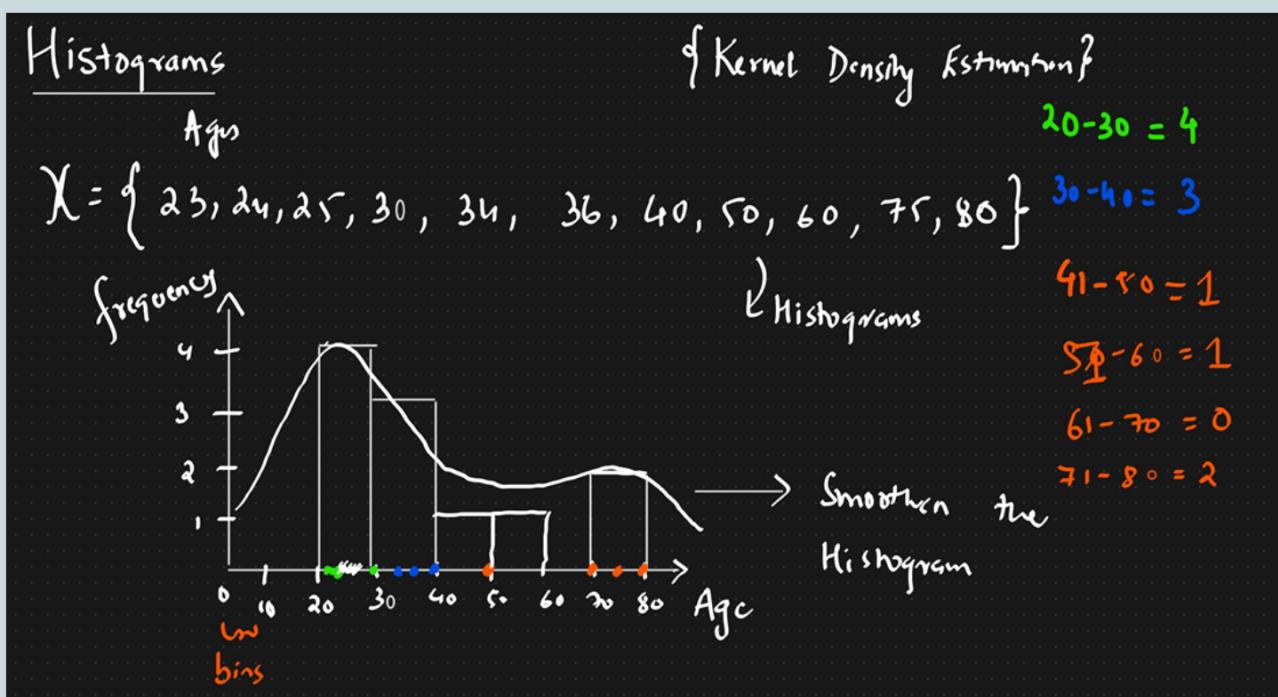
We can create bins:

Bin (Score Range)	Frequency (Count)
40–49	1
50–59	5
60–69	6
70–79	3

Solution →



Ques: →



## • Percentile and Quartiles

Percentiles And Quartiles

$$\text{Percentage} = \{1, 2, 3, 4, 5, 6\}$$

# No. of odd numbers = 3

$$\text{Percentage of odd numbers in this group} = \frac{3}{6} \times 100 = 50\%$$

Percentiles: A percentile is a value below which certain percentage of observations lie.

$$\{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, \boxed{9}, 9, 10\} \quad n=9 \quad \frac{3+4}{2} = \underline{\underline{3.5}}$$

$$\begin{aligned}\text{Percentile of Value } x &= \frac{\# \text{ of values below } x \times 100}{n} \\ &= \frac{11}{14} \times 100 \\ &= 78.57\% \text{ of value } 9\end{aligned}$$

←  $\underline{\underline{3.75}}$

Percentile  
Ranking

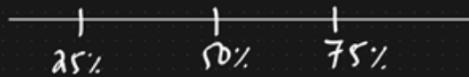
$78.57\% \text{ is } \underline{\underline{3.75}}$

$$\begin{aligned}\Rightarrow \text{Value} &= \frac{\text{Percentile}}{100} \times (n+1) \\ &= \frac{78.5}{100} \times (15) \\ &= \underline{\underline{3.75}} \approx \underline{\underline{3.5}}\end{aligned}$$

78.57 percentage of entire distribution is less than 9

## ② Quartiles

$$\left. \begin{array}{l} 25\% = 1^{\text{st}} \text{ Quartile} \\ 50\% = 2^{\text{nd}} \text{ Quartile} \\ 75\% = 3^{\text{rd}} \text{ Quartile} \end{array} \right\}$$



## • 5 Number Summary

1). Minimum

2). 1<sup>st</sup> Quartile (25%) (Q1)

3). Median

4). 3<sup>rd</sup> Quartile (75%) (Q3)

5). Maximum

$$Q1 = \frac{\text{Percentile}}{100} \times (n+1)$$

Eg: 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

[Lower fence  $\rightarrow$  Higher fence]

$\downarrow$   
Outlier

Anything outside this  $\uparrow$  is considered as outliers

$$\text{Lower fence} = Q1 - 1.5(IQR)$$

$$\text{Higher fence} = Q3 + 1.5(IQR)$$

$$IQR = Q3 - Q1$$

Inter Quartile Range

(25%)

$$Q1 = \frac{\text{Percentile}}{100} \times (n+1) = \frac{25}{100} \times (9+1) = 5^{\text{th}} \text{ position}$$

$\boxed{\text{value} \Rightarrow 3}$

75%

$$Q_3 = \frac{75}{100} \times 20 = 15^{\text{th}} \text{ position}$$

I value  $\Rightarrow 7$

$$\Rightarrow IQR = Q_3 - Q_1 \\ = 7 - 3 = 4$$

$$\Rightarrow \text{Lower fence} = Q_1 - 1.5(IQR) \\ = 3 - 1.5(4) \\ = 3 - 6 = -3 \\ =$$

$$\Rightarrow \text{Higher fence} = Q_3 + 1.5(IQR) \\ = 7 + 1.5 \times 4 \\ = 13$$

$$[-3, 13]$$

So, anything below -3 or greater than 13 is considered as outlier

So, 27 is outlier

So, After removing outlier

e.g.  $\rightarrow 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9$

So, Minimum = 1

1<sup>st</sup> Quartile = 3

Median = 5

3<sup>rd</sup> Quartile = 7

Maximum = 9

Box - Plot  $\Rightarrow$  This plot is used to visualize the outliers

# Covariance And Correlation

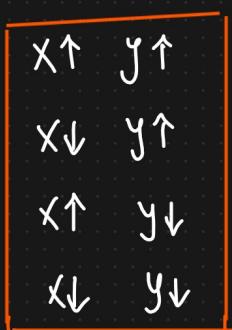
Covariance and correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with changes in another variable.

## Covariance

**Definition:** Covariance is a measure of how much two random variables change together. If the variables tend to increase and decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.

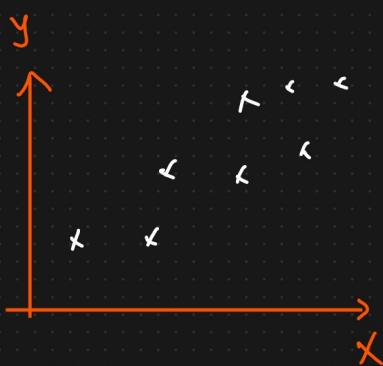
[Quantify the Relationship between X and Y]

X	Y
2	3
4	5
6	7
8	9

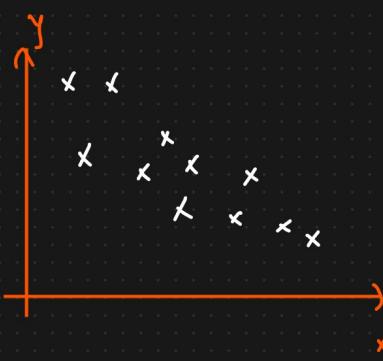


Dataset

↓ ↑ Size of house	Price ↑ ↓
1200	45 lakhs
1300	50 lakh
1500	75 lakh



⇒ +ve Covariance ⇒ +ve value



X	Y
7	10
6	12
5	14
4	16

⇒ -ve Covariance ⇒ -ve value

## Covariance

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

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

$$\boxed{\text{Cov}(X, X) = \text{Var}(X)} \quad \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

$x_i \rightarrow$  Data point of random variable  $X$

$\bar{x} \rightarrow$  Sample mean of  $n$

$y_i \rightarrow$  Data points of random variable  $Y$

$\bar{y} \rightarrow$  Sample mean of  $Y$

## Students

Hour Studied ( $X$ )

2

3

4

5

6

Exam Score ( $Y$ )

50

60

70

80

90

$x \uparrow y \uparrow \Rightarrow +ve$   
 $x \downarrow y \downarrow$  covariance

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

$$\textcircled{1} \quad \bar{x} = \frac{2+3+4+5+6}{5} = 4 //$$

$$\textcircled{2} \quad \bar{y} = \frac{50+60+70+80+90}{5} = 70 //$$

$$\text{Cov}(X, Y) = (2-4)(50-70) + (3-4)(60-70) + (4-4)(70-70) + (5-4)(80-70) + (6-4)(90-70)$$


---

4

$$\text{Cov}(X, Y) = \underline{\underline{20}}.$$

$\Rightarrow$  The positive covariance indicates the no. of hours studied increased the exam score also.

$$\left\{ \begin{array}{l} X \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{l} Y \\ 50 \\ 60 \\ 70 \\ 80 \end{array} \right\} \Rightarrow \underline{\underline{-ve}}$$

$x \uparrow y \downarrow$   
 $x \downarrow y \uparrow$

0.96

$\text{Cov}(A, B) \quad \text{Cov}(B, C)$

0.88

$$\begin{array}{ll} -200 & -300 \\ +100 & +300 \\ \hline 20 & 30 \end{array}$$

$\text{Cov}(A, B)$

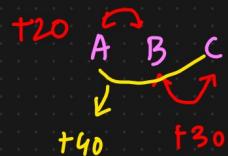
$\text{Cov}(B, C)$

$\text{Cov}(A, C)$

Advantages

[-1 to 1]

Disadvantage



- ① Quantify the Relationship between X and Y

- ① Covariance does not have a Specific limit value.

$$\text{Cov}(X, Y) \Leftarrow -\infty \text{ to } \infty$$

- ② Correlation

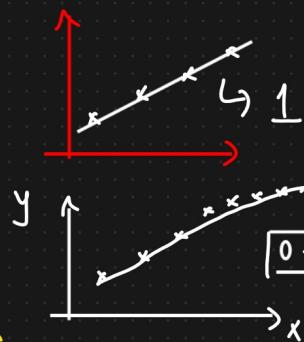
→ Pearson Correlation Coefficient  
→ Spearman Rank Correlation

- ① Pearson Correlation Coefficient  $\Rightarrow [-1 \text{ to } 1]$

$$\rho_{x,y} = \frac{\text{Cov}(X, Y)}{\sigma_x \cdot \sigma_y} = \frac{20}{\sigma_x \cdot \sigma_y} \Rightarrow 0 \text{ to } 1$$

- ① The more the value towards +1 the more +ve correlated X & Y is.  
② The more the value towards -1 the more -ve correlated it is (X, Y)

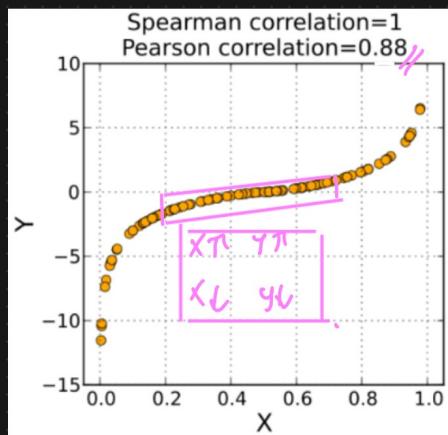
- ③ Spearman Rank Correlation



Pearson Correlation

Correlation for

non linear data



A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data points with greater  $x$  values than that of a given data point will have greater  $y$  values as well. In contrast, this does not give a perfect Pearson

$\Downarrow$   
 $\Rightarrow X \uparrow Y \uparrow$   
 $\Rightarrow X \downarrow Y \downarrow$   
 $\Downarrow$

Pearson Correlation  
 $= 0.88$



x	y	$R(x)$	$R(y)$
1	2	2	1
3	4	3	2
5	6	4	3
7	8	5	5
0	7	1	4

$$\gamma_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))} \Leftrightarrow$$

## Feature Selection

Size of house  $\uparrow$       No. of Room  $\uparrow$       location  $\uparrow$       ~~No. of people stay in the house~~

are correlated

