

Statistics

• intro to statistics

Descriptive

- measure of central tendency
- measure of dispersion
- measure of symmetry

Inferential

- Probability Distribution
- PMF and PDF & CDF
- Central Limit Theorem
- Statistical test

Definition → Statistics is a mathematical science of including methods of collecting, organizing and analysing data in such a way that meaningful conclusion can be drawn from them.

Why stat in data science →

In DS/ML → We try to learn pattern in data

→ it is raw fact and pieces of information that can be stored, measured and re-accessed.

~~x^{xx}~~ data is used to bring insights to increase the company's revenue ~~x^{xx}~~

Collecting

Source of data

organizing

to reaccess it efficiently

analyzing

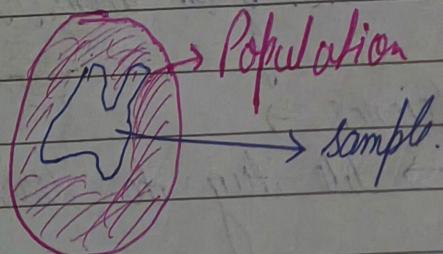
finding insights from data

Types of statistics

1) Descriptive statistics → it consist of organizing and summarising the complete data / population. → Population summarisation

2) Inferrential statistics → you can not count no. of trees in a jungle if someone ask you this. you will use your inference.

→ It consist of using ~~data~~ sample data that has been measured to form conclusion about a population.
With given sample we conclude something about a population.

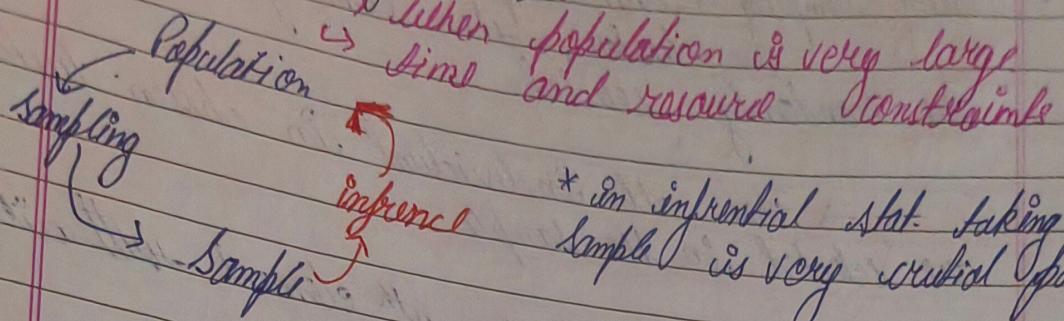


NOTE → When we want exact statistics to take any decision we use Descriptive statistics

Techniques of descriptive stat

- 1) Measures of central tendency (mean, median, mode)
- 2) Measures of asymmetry (skewness, kurtosis)
- 3) Measures of dispersion (std deviation, Variance)

Techniques of inferential stat



* in inferential stat. taking sample is very crucial part



Types of sampling →

- 1) Simple Random Sampling
- 2) Stratified Sampling
- 3) Cluster Sampling
- 4) Systematic Sampling

- Simple Random Sampling : every member of population(s) has an equal chance of being selected in the sample
- Disadvantages
 - A possibility of members not being part of sample from a certain group

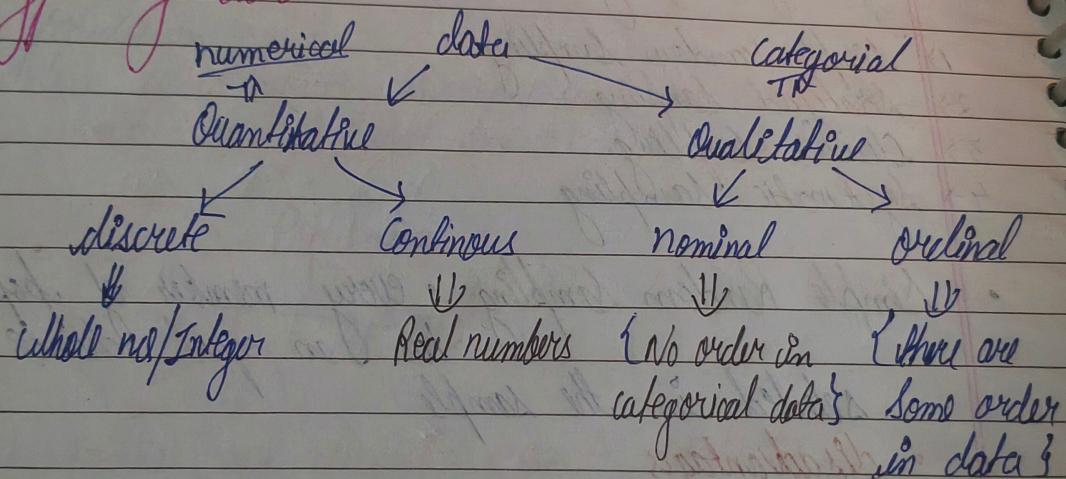
- Stratified Sampling → strata → layers/groups.
 - different distinct categories are there
 - A simple random sampling could be chosen from each strata or layer.

Cluster sampling → divides population into group or clusters, then some of these clusters are randomly selected.

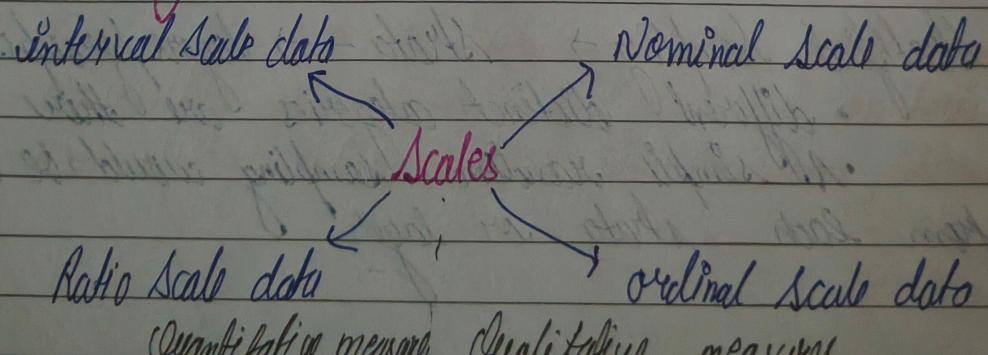
→ then all individual in chosen cluster are selected in the sample.

Systematic sampling → every n^{th} element will be selected

* Types of data



* Scales of measurement →



Nominal Scale Data

- Qualitative or categorical data
 - No order in data
 - Gender, color, loc etc.
 - Charts → pie, bar, ~~etc.~~
- May
→ Count
→ Show percentage
→ Use pie chart
to represent.

Ordinal Scaled Data

- Here order or rank has a meaning
- diff. can not be measured

plots → pie, bar, freq. or count

Interval Scaled Data

- The rank and order has a meaning
- diff. can be measured (excluding ratio)
- It does not have zero starting value.

plots → Histogram, Scatter plot, line chart

Ratio Scaled Data

- Order and rank has a meaning
- Difference and Ratio are measured
- It does have a 0 starting point Compulsory

Measure of Central tendency

Descriptive stat

→ summarization of data without removing or subtracting any value or instance

- Central tendency tries to represent the center point of a data set.
- mean
 - median
 - mode

mean → symbolic representation is \bar{x} for mean.

$$\bar{x} = \frac{1+2+3+4+5}{5}$$

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

→ summing up all the values of observation and dividing them with the no. of observation.

median → physical mid point of data

1. Sort the data

2. if count is even → median = avg of two middle values

if count is odd → median = $\left(\frac{n}{2}\right)^{\text{th}}$ value

Note → mean is affected by outliers whereas median is not affected at all.

mode → maximum frequency (the element that repeated highest no. of time)

xx with mean, median, mode we are trying to represent the central most value in our data.

if data is numerical data (continuous) → mean, median
if data is categorical data → mode.

Measure of dispersion

before understanding the range we will understand about some topic.

1) Range → difference between the minimum and the maximum value.

$$\text{Range} \Rightarrow \text{Max} - \text{Min}$$

Note → outlier affect Range that why we come upto another topic which is Quartiles.

• Percentage → {1, 2, 3, 4, 5}

$$\% \text{ of odd No} = \frac{3}{5} \times 100 = 60\%$$

60% of the no. numbers are odd

* percentage = $\frac{\text{favourable}}{\text{Total}} \times 100$

DAILY

Percentile → A percentile is a value below which a certain percentage of observations lie.

11, 2, 3, 4, 5, 6, 7, 8, 9, 10

* What is percentile rank of 5.

No below 5 = 2, Total no = 10

$$\frac{2}{10} \times 100 = 20\%$$

∴ 20 percentile of data lie below 20.

* What value exist at 75th percentile

$$\text{Value} = \frac{\text{percentile}}{100} \times (\text{Total value} + 1)$$

$$\text{Value} = \frac{75}{100} \times 11 = \frac{33}{4}$$

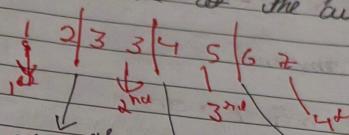
$$\text{Value} = 8.25$$

Here 8.25th value in data set is value below which 75% of data exist

Quartile \rightarrow

Quartiles are the values that divide a list of numbers into quarters.

- Put the no. in order
- Then cut the no. into 4 equal parts
- The quartiles are at the cut points, 1st, 2nd, 3rd, 4th



$$Q_1 \rightarrow 2$$

$$Q_2 \rightarrow 4$$

$$Q_3 \rightarrow 6$$

If total no. is odd

$$Q_1 = \left(\frac{n+1}{4}\right)^{\text{th}} \text{ no.}$$

$$Q_3 = 3\left(\frac{n+1}{4}\right)^{\text{th}} \text{ no.}$$

$$Q_2 = \left(\frac{n+1}{2}\right)^{\text{th}} \text{ no.}$$

median

If total no. is even

$$Q_1 = \frac{n^{\text{th}}}{4} \text{ no.}$$

$$Q_3 = \frac{3n^{\text{th}}}{4} \text{ no.}$$

$$Q_2 = \left(\frac{n^{\text{th}}}{2}\right) + \left(\frac{n+1}{2}\right)^{\text{th}}$$

2

Quartile 1 | Quartile 2 | Quartile 3 | Quartile 4

Q_1 Q_2 Q_3 Q_4

↓ divide into 2 equal parts

Five point summary
 Q₀ (min) → 0 → 0% percentile
 Q₁ → 1 → 25% percentile
 Q₂ → 2 → 50% percentile
 Q₃ → 3 → 75% percentile
 Q₄ (max) → 4 → 100% percentile

Transaction amount
 1000 → min 4000
 2000 → Q₁ = 5000
 This means 25% of transaction amount is equals to or below 5000 in data

Similarly for Q₂ = 10000
 Q₃ = 18000

Disadvantages of range → outlier affects the range

outlier	outlier	outlier	outlier	
Q ₀	Q ₁	Q ₂	Q ₃	Q ₄

These are the outlier ranges
 min < Q₁
 max > Q₄

1 | 2 | 3 | 3 | 3 |

percentile
percentile
percentile
percentile
percentile
percentile

12
Hence answer
3000

range

Interquartile Range →

PAGE NO.
DATE

IQR $\rightarrow Q_3 - Q_1$ { Since IQR deals with Q_3 and Q_1 , it is not affected by outliers. Outlier detection becomes easy with IQR.

$Q_1 = \text{Median of } Q_1$
 $Q_3 = \text{Median of } Q_3$

Range of interquartile range for symmetric data is zero.

$Q_1 - Q_3 = 0$
 $(Q_1 - Q_3) = 0$
Defined range
Range

maximum
minimum

Mean Deviation

by this we can
interpret that
on an average each
of the data is
1.2 unit away from mean value.

$$\text{formulae} = \frac{\sum_{i=1}^n |U_i - \bar{U}|}{n} + \frac{\sum_{i=1}^n |K_i - M|}{n}$$

- Variance \rightarrow The average of squared differences from the mean.

Population Variance

$$\sigma^2 = \frac{\sum_{i=1}^n |U_i - M|^2}{n}$$

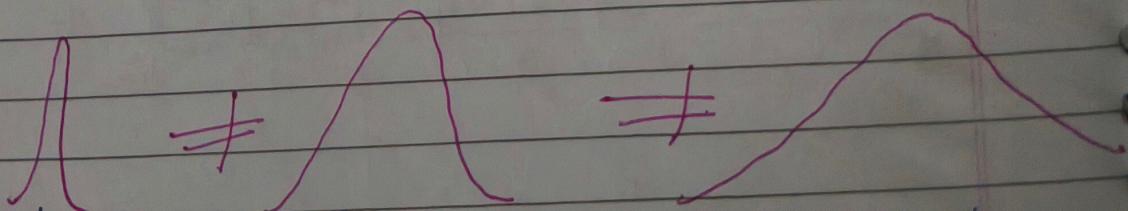
population
mean

Sample Variance

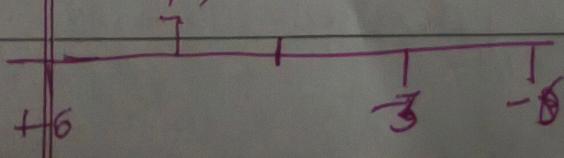
$$S^2 = \frac{\sum_{i=1}^n |U_i - U_s|^2}{n-1}$$

sample
mean

- Variance talks about spread at an overall level. & Variance increase as spread increases



Why squaring $\rightarrow \frac{+3 + 0 - 3 - 0}{4} = 0$



- steps to calculate variance
 - calculate mean
 - for each no. in data, calculate diff b/w the mean and no.
 - squared the difference
 - calculate average of squared differences.

Standard Deviation
 it is a measure of how spread out numbers are.
 square root of variance

$$\sigma = \sqrt{\text{Var}}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

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

why std. deviation?

1. Variance can be huge no. because it talks about spread of an overall level. Comparison of each no. w.r.t variance becomes difficult.

$$2. \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}} \rightarrow \text{dimensions are squared}$$

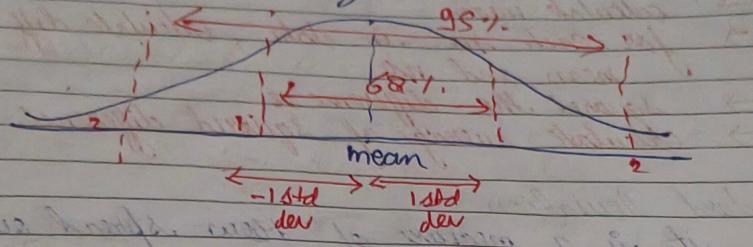
in variance dimension changes?

$$\text{ex} \rightarrow \frac{(1m - 2m)^2 + (3m - 2m)^2}{2} = \frac{1m^2 + 2m^2}{2} = m^2$$

Here dimension changed from m to m²

→ std dev. → std way of knowing where your data lies.

if data is normalize



Why we using squared diff in variance.

i) it negates the values if not sq.

$$\text{eg} \rightarrow \text{all } = 0 \quad +2 \quad +1 \quad -1 \quad -2 \\ = +2 + 1 - 1 - 2 = \frac{0}{4} = 0$$

→ Here is the spread = 0, rd.
if it is zero, then all values = mean

ii) if we make absolute instead of sq.

$$\text{SC-1} \quad -1 \quad | \quad 1 \quad | \quad 3 \quad | \quad 3 \quad | \quad 6 \quad = \frac{3+3+6+6}{4} \\ = 4.5$$

$$\text{SC-2} \quad +4 \quad | \quad 1 \quad | \quad -1 \quad | \quad 1 \quad | \quad -5 \quad = \frac{1+2+4+5}{4} \\ = 4.5$$

∴ Here if we see we get same value 4.5

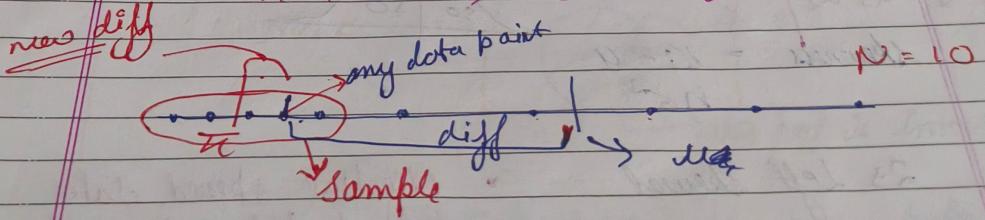
But do both dataset have same variance?
→ NO

→ That's why we use squared difference

Why in population variance it divides by N but
in sample variance it is $n-1$
when we calculate any sample stat.

↳ when we do not have access to complete population
we use this because sample variance will be
biased estimator

$$S^2 = \frac{1}{n-1} \sum (u_i - \bar{u})^2$$

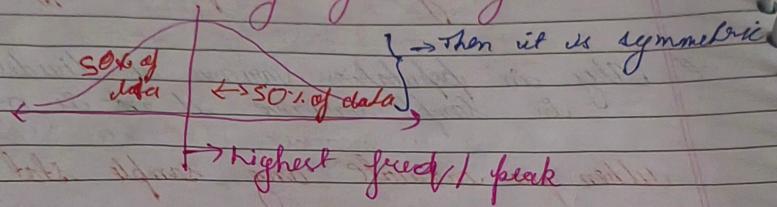


$$(u_i - \bar{u})^2 < \frac{(u_i - \bar{u})^2}{n}$$

So to get ans near to σ^2 we reduce
 $n = n+1$ as it is denominator
(denominator ↓, ans ↑)

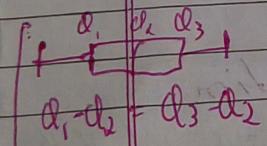
* Measure of symmetry → talks about shape of data

skewness → measure of symmetry

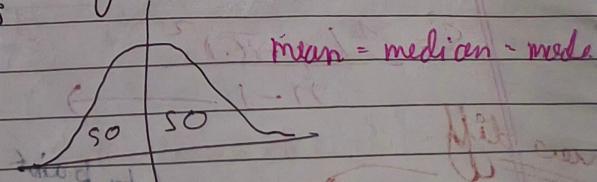


skewness → inclined to some side

1) No skewness $\Rightarrow \{ \text{symmetric distribution}$
 $\{ \text{skewness} = 0 \}$

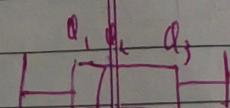


$$\text{skewness} = \frac{\bar{x} - M}{n^{\frac{3}{2}}}$$

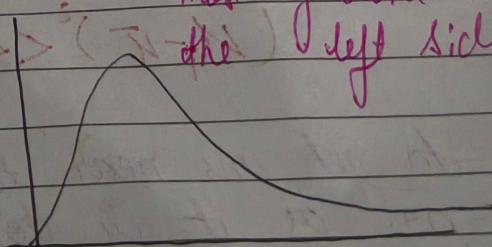


2) Left skewed

- Negative skewed data
- most of data lies at the left side.



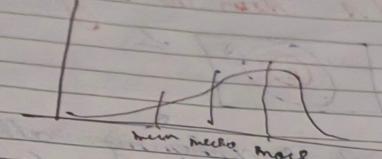
$$d_3 - d_2 > d_4 - d_1$$



3) Right Skewed

$$\{ \text{mean} > \text{median} > \text{mode}$$

3) Right skewed data



mean > median > mode

Set → collection of unordered unique elements
Properties of set

1) intersection → common element

$$A \cap B = \boxed{\textcircled{1} \textcircled{2}}$$

→ This part is intersection

2) Union → all distinct elements from both sample

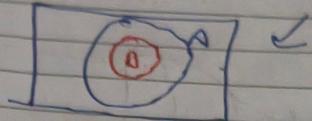
$$A \cup B = \boxed{\textcircled{1} \textcircled{2}}$$

3) difference → elements which are present only in

$$A - B = \boxed{\textcircled{1}}$$

$$B - A = \boxed{\textcircled{2}}$$

4 → Subset \rightarrow if all elements of B present in A then we say B is subset of A $\rightarrow A \supset B$ or $B \subset A$

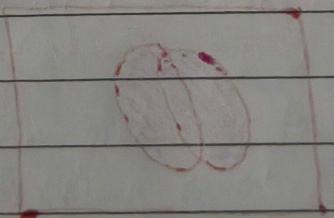
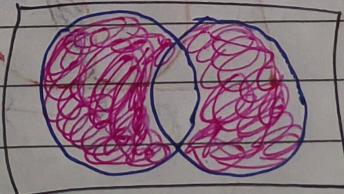


5 → Superset \rightarrow A is containing all element of B - then A is superset of B $\rightarrow A \supset B$

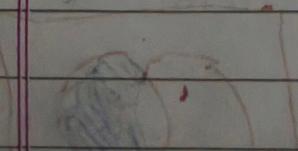
6 → Symmetric Diff \rightarrow off of intersection

* The element of that are distinct in both excluding intersection

$$A \Delta B \text{ or } \{A \cup B - A \cap B\}$$



$$= A \cup B$$



$$= A \cap B$$

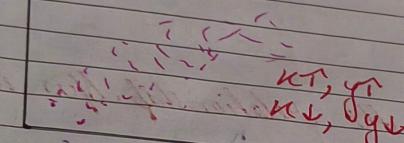


$$= A - B$$

Covariance And Correlation

Here we say as x increase
 y also increase

+ understanding the relationship



direct relationship
 \Rightarrow price of house based on area / locality.

	Total transaction count	Total transaction amount
5	10K	
4	8K	
3	3K	

indirect relationship
 \Rightarrow No. of years and alcohol consumption

Quantifying / measuring relationship

$$\text{Covariance} \rightarrow \text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

x, y are features

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

Co-Variance \rightarrow spread
 exist negative

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

So we can say that variance is the relationship of feature with itself.

So we can say it is telling relationship with itself.

P Similarly Covariance is nothing but we are trying to calculate relationship of any two features with respect of other.

T) x Advantages of Covariance

→ We now know the relationship b/w x & y.

disadvantage $\{ \text{Cov}(u, y) = \frac{\sum_{i=1}^n (u_i - \bar{u})(y_i - \bar{y})}{n-1} \}$

$$\text{Cov}(u, y) = 50 \quad \text{Cov}(A, B) = 100.$$

Can we say $\text{Cov}(u, y) = \frac{1}{2} \text{Cov}(A, B)$?

No.

Why $\rightarrow \text{Cov}(2x, y) \rightarrow$ range $\rightarrow (-\infty, \infty)$

The we can not compare strength as both may have different starting point. Or 'O'.

2 → No comparison of strength of relation, NO any standardised scale to interpret the strength.

⇒ Covariance has dimension.

$$\text{Cov}(u, y) = \frac{\sum_{i=1}^n (u_i - \bar{u})(y_i - \bar{y})}{n-1}$$

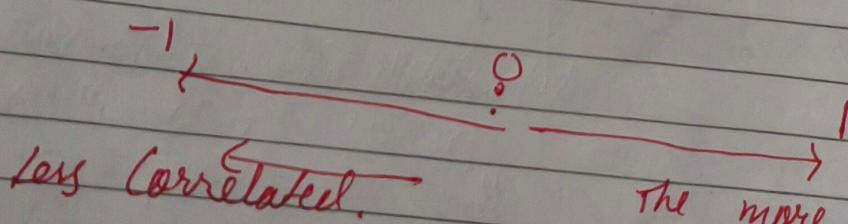
$$\text{Cov}(\text{transcition}, \text{height}) = \text{Rs. ft} - 450 \text{ Rs. ft.}$$

∴ we can not compare two different dimensions. ⇒

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

$$f(u, y) = \frac{\text{Cov}(u, y)}{\sigma_u \sigma_y} = [-1 \text{ to } 1]$$

→ it solve dimension problem as it cancel out and become dimensionless



less correlated.

The more positively correlated features are

- Note →
- Pearson Correlation Coefficient always measures the linear relationship.
 - Correlation is never a slope
 - Non linear data have pearson correlation = 0

for non linear data relation ship we have

→ Spearman Rank Correlation

$$\rho_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

$R(x) \rightarrow$ Rank of x
 $R(y) \rightarrow$ " " y

Rank →
1 → Sort value
2 → Highest will
3st etc will
increase
3 → use the Rank
in position of data
points.

PROBAB

• Rank

Set of Rank