

## Introduction to fundamental statistical Concepts

\* Mean (the average) np.mean(value) @ arithmetic average

The mean (@average) is a measure of central tendency that represents the typical value in a dataset. It is calculated by, Mean =  $\frac{\sum \text{All values}}{\text{No. of values}} = \frac{\sum N}{N_{\text{num.}}}$

most common type of mean is arithmetic mean.

\* Median (middle value) np.median(value)

The median is the middle value of a dataset when arranged in ascending or descending order. It is a measure of central tendency that helps identify the center of the data.

How to calculate?

Sort the data.

$\frac{N+1}{2}$

a) Arrange the numbers in a order.

b) If the variable is odd, middle number is the median

If the variable is even, average of two middle number is median

$\frac{N}{2} \text{ or } \frac{N+1}{2}$

\* Correlation

Correlation is a statistical measure that describes the relationship between two variables. It tells us how strongly and in what direction the two variables move together.

Types of Correlation

a) Positive

one variable increases, other

b) Negative

one variable increases, other

c) No correlation

No relation between variables

\* Mode, is the elements with the common value in the data.

Unimodal - distribution of values with only one mode  
 bimodal - with two mode  
 multimodal - more than 2 mode  
 Mode can be both categorical/numerical

Python -

from statistics import multimode  
 $x = \text{multimode}(\text{value})$

Using R -

```
mode <- function(x) {
  Unique-values <- Unique(x)
  table <- tabulate(mode(n, unique-value))
  Unique-values [table == max(table)] }
```

## \* Correlation

Correlation is a statistical relationship between two variables and in what direction the two var

### Types of Correlation

#### a) Positive

one variable increases, other also increases.

#### b) Negative

one variable increases, other decreases.

#### c) No correlation

No relation between variable.

\* Mode, is the elements with the highest frequency or the most common value in the data.

**Shape of Distribution**

- Symmetry & skewness
  - Positive Skew (right-skewed) tail extends to right, mean > median.
  - Negative skew (left-skewed) tail extends to left, mean < median

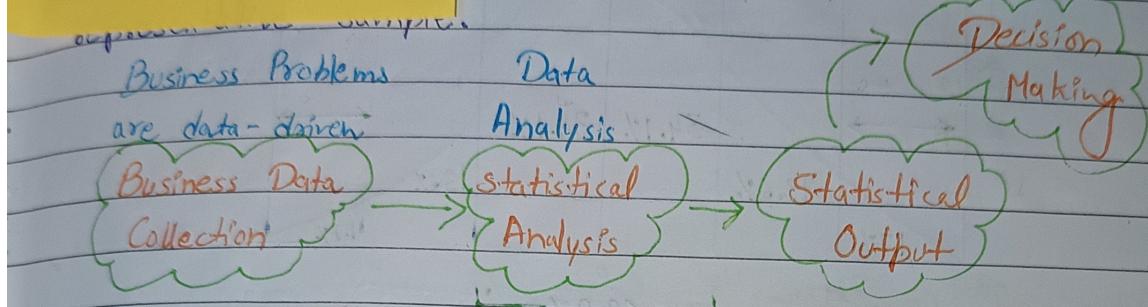
- Kurtosis
 

Describes the 'tailedness' of the distribution. It tells us how many outliers exists.

Date 27 March 2023  
DELTA Pg No.

Histograms → Statistical Analysis  
Histograms → Non-Statistical Analysis

Process of collecting and analysing the large quantities, especially for the purpose of making in a whole from those in a sample.



### Categories in Statistics (two major)

(a) Descriptive Statistics - organising/summarizing the data

Describe the basic features of the data in a study. They provide simple summaries about the sample and the measures. {Average/%} "Numerical Analysis/output"

(b) Inferential Statistics - for conclusions of the data,

Helps in generalising about the population by using various analytical tests and tools.

"Relative Analysis/output".

### Statistical Measures (Three major)

(a) Measure of Frequency

Describe basic features of data in a study. They provide simple summaries about the sample and the measures.

(b) Measure of Central Tendency (mean/median/mode.)

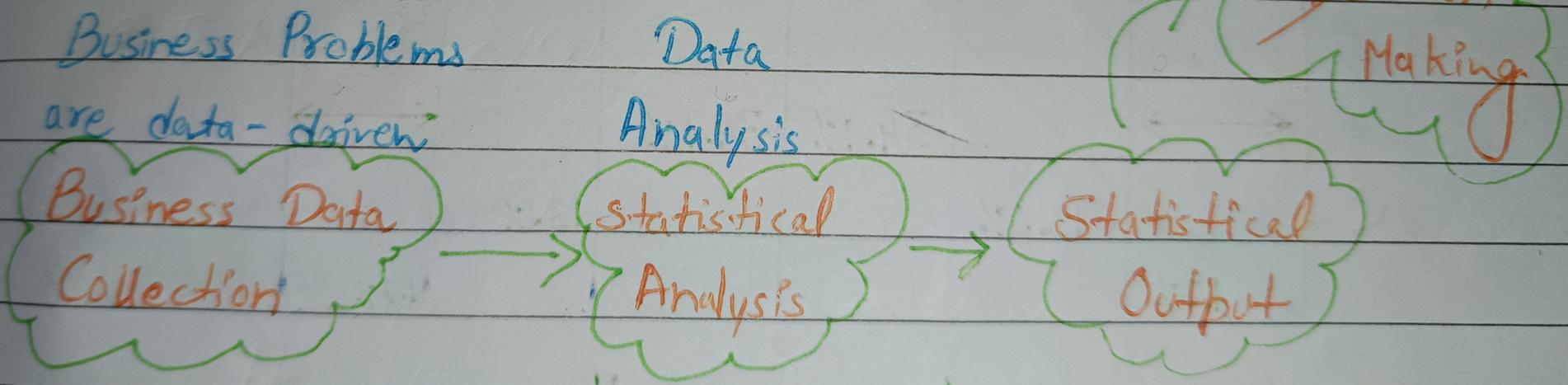
Helps in generalising about the population by using various analytical tests and tools. (center of distribution of Data)

(c) Measure of Spread (standard deviation/variance/IQR/Range)

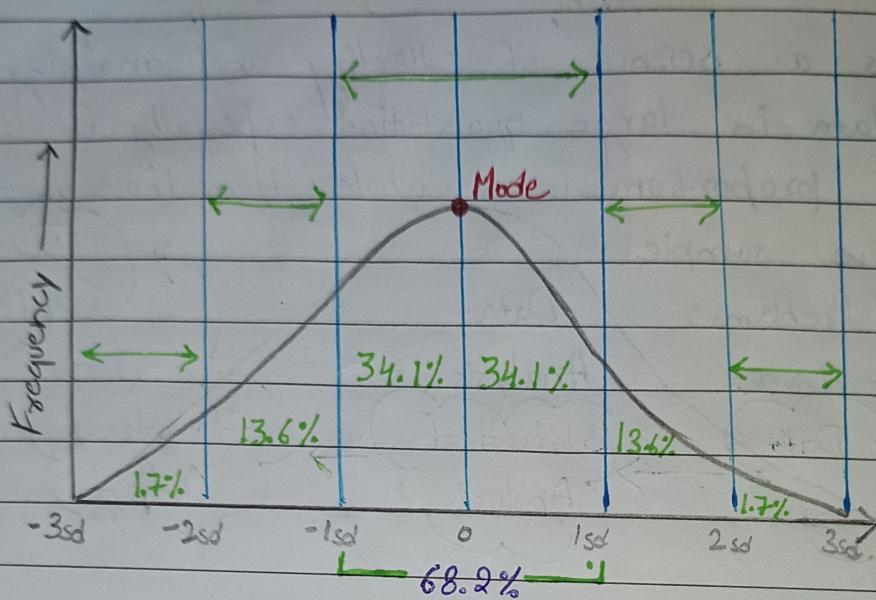
Helps in generalising about the population by using various analytical tests and tools.

Introduction to statistics → statistical Analysis  
→ Non-statistical Analysis

Statistics is a science of collecting and analysing the numerical data in large quantities, especially for the purpose of inferring proportions in a whole from those in a representative sample.



## • Distribution Curve (bell-curve)



↳ most of the data points are accumulated here.

**range**: distance between minimum and maximum values

**frequency**: No. of repetitions of values in x-axis.

**Central Tendency**: accumulation of data points towards the center.

### Null Hypothesis ( $H_0$ )

It is commonly accepted fact. In ML, it usually represents the conclusion that the model will not work. It is the opposite of the alternative hypothesis ( $H_A$ ). Researchers attempt to reject the null hypothesis.

$$z = \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \quad \text{and} \quad t = \frac{\bar{x}_n - \mu}{s/\sqrt{n}}$$

## Week 1 - (Thursday)

### (1.1) Fundamental Mathematics for Data Science. (1.1)

#### Key Mathematical Concepts for Data Science

- Linear Algebra and Probability

Understand vectors, matrices and distributions

- Essential Operations

Performs matrix manipulation and probability calculations

- Descriptive statistics

Summarise data and interpret visualisations

- Probability Applications

Apply distributions to predictive modelling.

### Linear Algebra: The backbone of Data Science. (1.2)

#### Essential Tools for complex Problem Solving.

- Focuses on vectors, matrices and linear transformations

- Powers data science, machine learning, physics and engineering

- Provides a structured approach to multi-variable problem-solving

- Enhances data analysis, algorithm optimisation and model development

#### Key Components in Linear Algebra.

##### The building blocks of Data magic

- Vectors

Lists of numbers that show direction or represent the data points

- Matrices

Grids of numbers used to store and organise data

- Linear Transformation

Ways to change or move data, like scaling or rotating.

## Applications of Linear Algebra in Data Science

Data Representation

Structures data with vectors and matrices

Dimensionality Reduction

Simplifies datasets with retaining insights

Machine Learning Algorithms

Powers algorithms like regression and SVM

Transformation & Projections

Converts data to reveal hidden patterns

Optimisation Problems

Solves tasks like minimising loss

Big Data Operations

Enables fast processing of large dataset

Neural Networks

Supports deep learning transformations

Feature Engineering

Create impactful features for models

Graph Theory

Models networks and connections

Broad Applications

Extends to AI, NLP and bioinformatics.

### Scalars Made Simple

(1.3).

Single Numerical value with no Direction

#### \* Types

- Real Scalars

Whole numbers, fractions, and irrational numbers

(e.g., 5, -3.2, 0.75, ...)

- Complex Scalars

Numbers in the form  $a+bi$ , where ' $i$ ' represents the imaginary unit

#### \* Operations

There are four operations that can be done with scalars :-

- Addition

- Subtraction

- Multiplication

- Division

#### \* Properties

Scalars are subjected to below properties :-

- Commutative

- Associative

- Distributive

What are vectors?

The Power of size and direction.

- A vector is a mathematical object that has both magnitude (size) and direction.
- A vector is typically represented as an arrow. The length shows the magnitude and the head shows the direction of the vector.
- Vectors are widely used in fields such as physics, engineering & data science for various applications.
- They are essential tools for accurate analysis, modelling and solving complex problems in these areas.

Magnitude

The length of the vector, showing the size or strength of the quantity.

Direction

The vector's orientation in space, indicated by an arrowhead.

### Notation Conventions.

Scalars	Vectors
• Represented by italicised letters (a, b, c)	• Represented by bold letter ( $\mathbf{v}$ ) or arrows ( $\vec{v}$ )
• Indicate single values without direction.	• Indicate quantities with magnitude and direction.
• Usually written in plain typeface.	• Can be expressed as column vector

$$\mathbf{v} = \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$

### Vector Magnitude

The distance from the vector's starting point to its endpoint in space. Used in distance measurements and in normalising vectors to unit length.  $\|\mathbf{v}\| = \sqrt{v_1^2 + v_2^2 + v_3^2}$

## Classification of vectors.

vectors types and their unique characteristics

- Position
  - represent a point in space relative to an origin
  - describe spatial relationships within a coordinate system

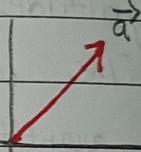
- Displacement
  - Indicate the change in position from one point to another
  - Capture both distance and direction of movement

- Zero
  - A vector with all components equal to zero
  - Represents no direction or magnitude

- Unit
  - vectors with a magnitude of one, used to show direction
  - Represent direction without considering magnitude.

## Visual Representation of vectors

- vectors are visually represented by arrows, showing both magnitude and direction.
- Its magnitude is represented by the arrow's length.
- Its direction is indicated by the arrow's orientation.



Same for 3-dimensional.

## Applications of vector Operations.

(1.4)

- Harnessing the power of vector operations
- Machine Learning Data point as vector

- Measure similarity
- Computer Graphics Group similar points
- Enable 3D rendering

Simulate lighting

- Life-like visuals
- Physical simulations Simulate motion
- Give direction + value
- Precise simulation

The role of matrices in Data Science. Rows - m Columns - n

### Data Representation

- Organises data into rows and columns
- Simplifies processing

### Transformation

- PCA reduces features, keeps variance
- Scaling standardises data.

### Machine Learning

- Regression finds co-efficients
- PCA extracts key components.

## Types of Probability

Theoretical	Based on logic and models	Dice
Empirical	From observations and data	weather forecasting
Subjective	Based on personal judgement	Estimating team history

## Types of measuring Data.

- Nominal Scale - categories without order
  - ordinal scale categories with order
  - Interval scale equal intervals, no true zero
  - Ratio scale Equal intervals, with a true zero
- median  
mean  
mode

## Types of Statistics.

Descriptive statistics →

- measures of Central Tendency
- Measures of dispersion (spread of data)
- Data visualization

- It summarises & describes the key characteristics of dataset, providing an overview of its main features

- Uses measures like mean, median and mode.

### Inferential statistics

- It makes prediction about a population based on sample data, allowing for generalisations and conclusions about larger groups.
- Involves hypothesis testing & estimation for informed decision making.

## Applications of Central Tendency and Dispersion.

### (A) Descriptive Analysis

Summarises key dataset characteristics highlighting central tendency and spread.

### (B) Hypothesis Testing

Determined if observed differences are statistically significant.

### (C) Risk Assessment

Assesses data variability and the potential for extreme values, crucial for predicting risk.

## Types of Matrix.

Square matrix

Identity matrix

Zero matrix

Diagonal matrix

Transpose

Determinant

## Principal Component Analysis (PCA)

It is a statistical technique used to simplify large datasets by reducing their dimensions while retaining most of the original information. It transforms correlated variables into a smaller set of uncorrelated variables called principal Components. This method is widely used in data preprocessing for machine learning, pattern recognition, and image processing.

In dimensionality reduction, techniques like principal components analysis (PCA) rely on matrices to reduce the number of features while preserving variance. Similarly, feature scaling uses matrix operations to standardise and normalise data, ensuring features are on a comparable scale.

In machine learning, matrices underpin key calculations. In linear regression, they help solve for model coefficients. PCA uses matrix decomposition to identify critical components, while neural networks depend on matrix operations to manage weights and activations across network layers, enabling fast computation of large datasets.

## LIVE SESSION

Math/Stats for DS

Week 1

\* Agenda

a) Prob. distribution → Descriptive

b) Basics of Statistics → Inferential

c) Simple descriptive analysis.

Houseprice → R (as a tool)

\* Descriptive Statistics.

Describing the data → Measure of numbers (frequency)

→ Measure of central Tendency

(Mean, median, mode.)

Avg

Central value

Highest frequency

\* Data (Prob. Distribution)

a) Quantitative

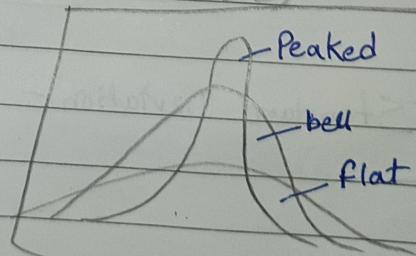
(Measured)

Discrete

Continuous

b) Qualitative.

(non-measurable)



ND 1) Bell shaped Curve

2) Mean & std dev

↳ the curve is symmetric about its mean

↳ std dev impacts the peakedness.

3) Area under the curve = prob. associated within the range.

4) Asymptotic → Tails extend to the infinity

5) Total Area Under Curve = 1

## Scales of Measurement.

### A) Nominal Scale.

Names, Labels, qualities.

No ranking / orders / Numbers

Mode / frequency / chi-square.

### B) Ordinal Scale

same as nominal but there is rank / order.

Median / Mode.

### C) Interval Scale.

- Equal interval

- No True Zero

- Equal difference b/w intervals.

Mean, SD, ANOVA, Regression.

### D) Ratio Scale

- true zero

- " + , - , X ,  $\div$ "

Mean, Median, Mode, SD, Geometric mean,

Harmean Mean.

## • Discrete Data. (Whole number)

(Bar Graph)

a) Bernoulli Distribution — a process with 2 possible outcomes

b) Uniform Distribution

c) Binomial Distribution —

d) Poisson Distribution — Models the prob. of events

Occurring in fixed intervals.

(Columns / variables / attributes)

Date /  
DELTA Pg No.

## Central Limit Theorem

### Random Variable

a random process, where each outcome is associated with some numbers.

Add  $N$  numbers of samples of this variable.

$$x_1 + x_2 + x_3 + \dots + x_n$$

The distribution of this sum looks more like a bell curve as  $N \rightarrow \infty$ .

$$\text{Mean, } \mu = E[x_i] = \sum_n P(X=x_i) \cdot x_i$$

$$\text{Variance, } \text{Var}(n) = E[(X-\mu)^2]$$

$$\text{Standard deviation, } \sigma = \sqrt{\text{Var}(n)}$$

### Normal Distribution

$$\text{pdf} = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

a) Bell Shaped

mean = distribution is symmetric b) depends on mean and std.

orbit mean

c) Area under curve  $\Rightarrow$  prob. associate with range

std = flatness or

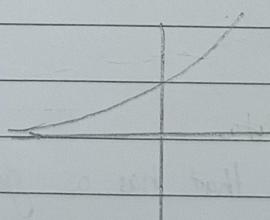
d) Asymptotic curve  $\Rightarrow$  tails extend till  $\infty$

peakness

e) Total area under the curve = 1

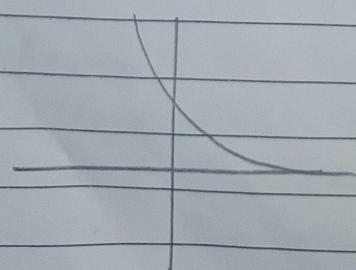
### Exponential Growth

$$e^n$$



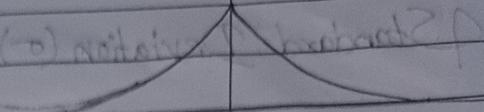
### Exponential decay

$$e^{-n}$$



Exponential distribution.

$$e^{-ln^2}$$



$$e^{-x^2}$$

bell curve.

normal distribution

### Variation

There are different measures commonly used are :

- Range (largest value - smallest value)
- Quartiles and Percentiles
- Interquartile Range.
- Standard Deviation.

### Quartiles.

$$\text{np.quantile(values, n)}$$

Between  $Q_0$  and  $Q_1$  are 25% lowest values in data.

b/w  $Q_1$  &  $Q_2$  are next 25%, and so on.

### Five Number Summary

$Q_0$  - smallest value.

$Q_1$  - separating value ( $P_{25\%}$ )

$Q_2$  - middle value (median) ( $P_{50\%}$ )

$Q_3$  - separating value ( $P_{75\%}$ )

$Q_4$  - largest value ( $P_{100\%}$ )

Outliers — extreme & unequal values in data.

Percentile, a value below which a certain percentage of observation lie.

$$\text{Percentile} = \frac{\text{value}}{n} \times 100$$

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

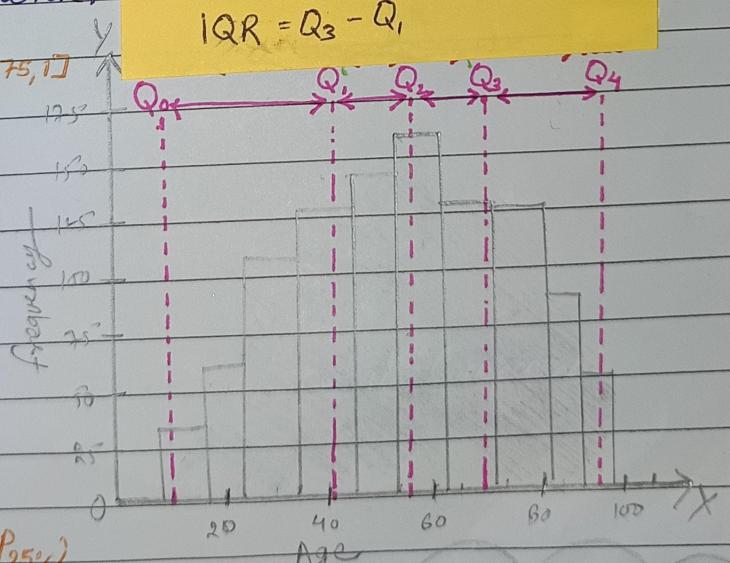
Removing Outliers (Generic)

Lowerfence  $\leftarrow$  Upperfence

$$LF = Q_1 - 1.5(IQR)$$

$$UF = Q_3 + 1.5(IQR)$$

$$IQR = Q_3 - Q_1$$



np.percentile(values, 65)

65th percentile.

## Variation

There are different measures of variation. The most commonly used are :

- Range (largest value - smallest value) np.ptp(values)

4 parts ← → 100 parts

- Quartiles and Percentiles

stats.iqr(values)

- Interquartile Range (IQR)  $[Q_1 - Q_3]$

- Standard Deviation.

Quartiles.

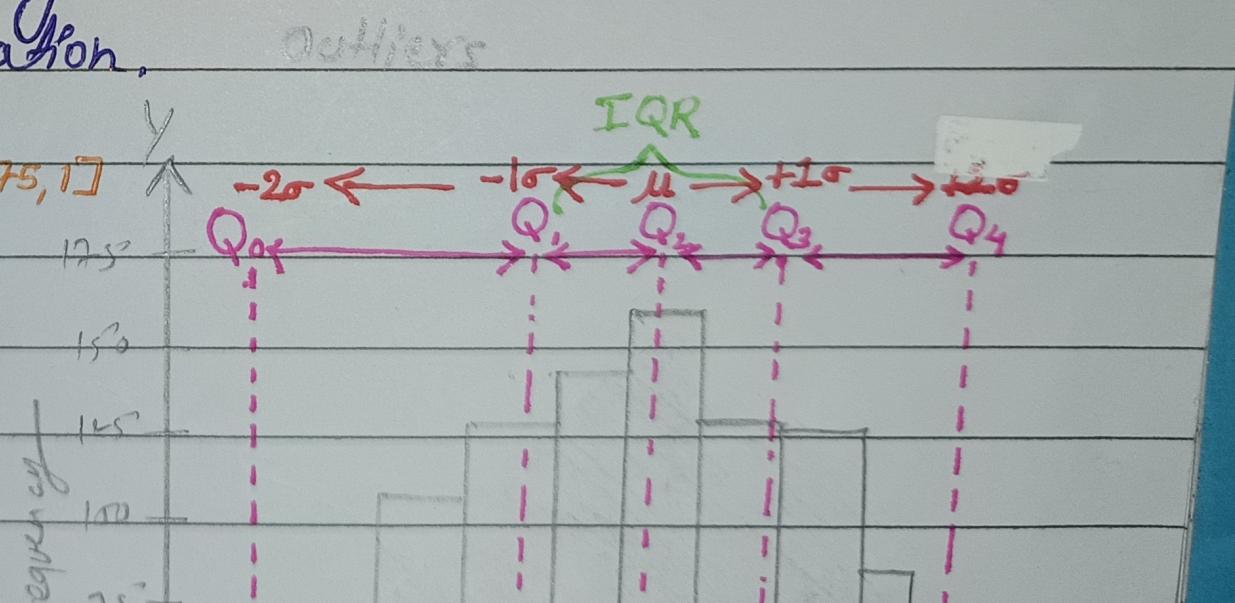
np.quantile(values, q)

Between  $Q_1$  and  $Q_2$ , are  
25% lowest values in data.

b/w  $Q_1$  &  $Q_2$  are next

$$Ad(25) = \frac{1}{6} \times 25 \times 15 \times (n+1)$$

$$Ad(25) = \frac{\mu}{100} \times 120$$



Standard Deviation ( $\sigma$ )

Variance

np.std(values)

Calculating Standard deviation :-

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

Calculating sample standard deviation :-

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

Why sample variance is  $\frac{1}{n-1}$ ?

random variable,  $X$   $\rightarrow$   $Z$  (statistics)

$Z$  = deviation of  $x$  from its mean in no. of its std.

$$Z = \frac{x - \mu}{\sigma} \quad \text{Always } \text{mean} : 0 \\ \text{std} : 1$$

Prob. density function

Date  
DELTA Pg No.

$$PDF = f(n) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{n^2}{2\sigma^2}}$$

Distribution.

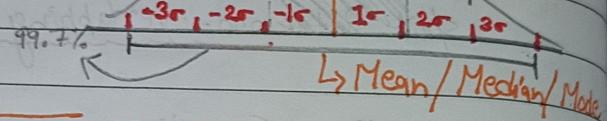
$$\text{(mean)} Z = 0 \\ \text{(std)} Z = 1$$

A) Gaussian / Normal Distribution

Empirical formula :-

68 - 95 - 99.7% Rule.

$$Z \text{ score} = \frac{x_i - \mu}{\sigma}$$



Standardization  $\leftarrow$  Z-score

$$\{1, 2, 3, 4, 5, 6, 7\} \rightarrow \text{Normal}$$



Z-score



$$\{-3, -2, -1, 0, 1, 2, 3\} \rightarrow \begin{array}{l} \text{Standard Normal} \\ \text{Distribution.} \\ \Sigma \mu = 0, \sigma = 1 \end{array}$$

$$* y \sim \text{std } (\mu = 0, \sigma = 1) *$$

It is the process of rescaling the data so that it has:

mean = 0

$$STD = 1, n_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Bounded data  
NN, KNN, CNN

Normalization  $\leftarrow$  Scaling

If refers to scaling the data to fit within a specific range — usually it's 0 to 1.

Minmax Scalar  $\rightarrow (0 \text{ to } 1)$

CNN  $\rightarrow$  Image classification.

Pixels

Normalisation

$0 - 255$

Minmax Scalar

$[0 - 1]$

## Elements of structured Data

Numerical  
numeric scale.

→ Continuous  
interval, float  
or numeric

→ Discrete  
only integers,  
count

Categorical  
enums, enumerated,  
factors, nominal.

→ Binary  
dichotomous, logical,  
indicator, boolean.

→ Ordinal  
ordered factor.

## Estimation of Location

- Mean (average)

the sum of all values divided by number of values.

- weighted mean (weighted average)

sum of all values times a weight divided by sum of weights.

- Median (50<sup>th</sup> percentile)

value such that one-half of data lies above & below.

- Percentile (quantile)

value such that p percent of data lies below.

- weighted median

value such that one-half of sum of weights lies above & below the sorted data.

- Trimmed mean (truncated mean)

average of all values after dropping a fixed number of / extreme values.

- Robust (resistant)

Not sensitive to extreme values.

- Outlier (extreme value)

a data value that is too differ from other values.