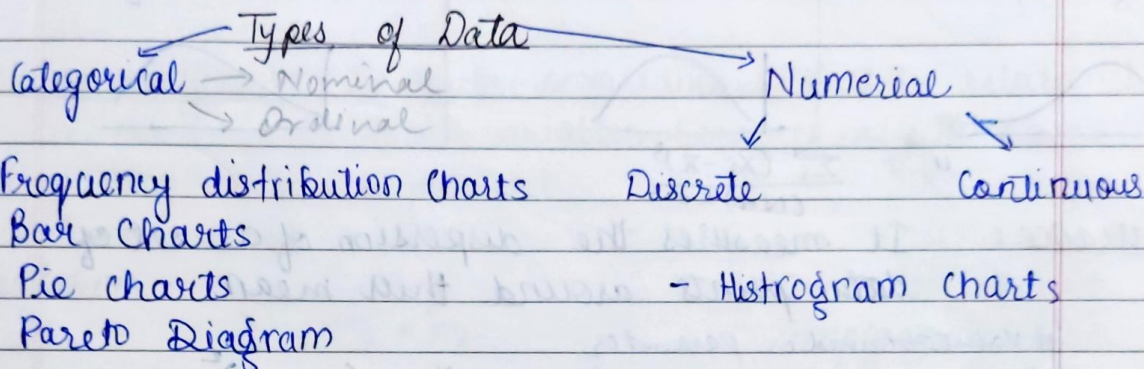# Statistics Essential for Data Science.

**Population:** - Collection of all items of interest to our study.
- $N$ (parameter)
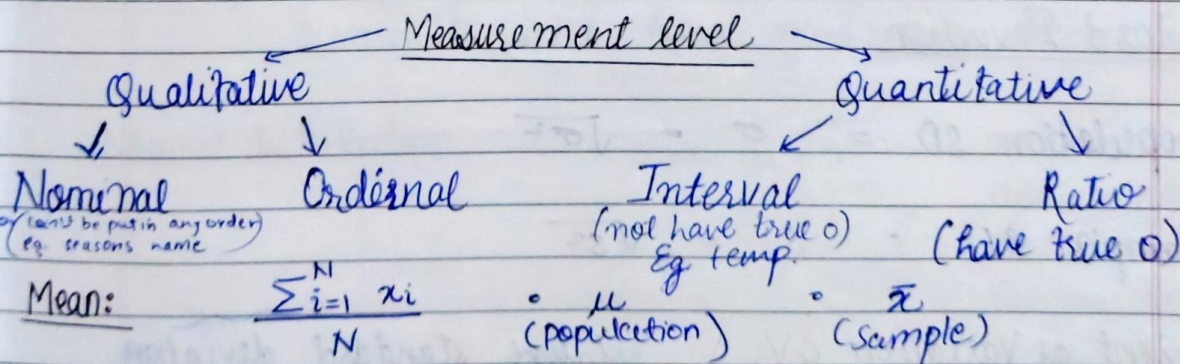
**Sample:** - A subset of a population
- $n$ (statistics)

## Types of Data

Categorical → Nominal
→ Ordinal

Numerical

↓ Discrete          ↘ Continuous

- Frequency distribution charts
- Bar Charts
- Pie charts
- Pareto Diagram

- Histrogram charts

' To show relationship b/w two numerical or Categorical
  identity → ✓ Cross Table (side by side bar charts)
          ✓ Scatter Plot
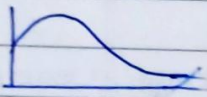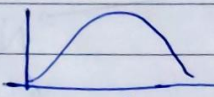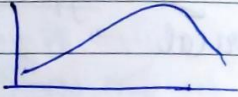
**Frequency:** Measure of occurence of variable

**Relative Frequency:** Measure the relative No. of occurance of variable. usually expressed in percentage.

## Measurement level

Qualitative                                    Quantitative

↓                    ↓                ↓                    ↓

**Nominal**          **Ordinal**      **Interval**         **Ratio**
or (can't be put in any order)                   (not have true 0)    (have true 0)
eg seasons name                                   Eg temp.

**Mean:** $\dfrac{\sum_{i=1}^{N} x_i}{N}$   • $\mu$ (population)   • $\bar{x}$ (sample)

easily affected by outliners

# measure of asymmetry

**skewness :** It indicates ~~wa~~ whether the data is concentrated on one side

| tve /right | 0 / symmetric | -ve / left |
|---|---|---|
| iv mean > median | - mean = median | - mean < median |
| - outliers are at right (tail) | - no c | - outliers are at to the left |

$$\tilde{u}_3 = \frac{\sum_i^N (X_i - \bar{x})^3}{(N-1) \times \sigma^3}$$

**Variance :** It measures the dispersion of a set of data points around their mean

# Value obtained is parameter

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

# value obtained is statistic

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

**? Why Squaring**
(i) Dispersion in non-negative
(ii) Non negative value don't cancel out
(iii) Amplifies the effect of large differences

## Standard Deviation :

population SD $= \sigma = \sqrt{\sigma^2}$

Sample SD $= s = \sqrt{s^2}$

## Coefficient of Variation CV.
relative standard deviation

population CV $= C_v = \frac{\sigma}{\mu}$

sample CV $= \hat{C}_v = \frac{s}{\bar{x}}$

Used to compare data sets in terms of variability

## covariance: $\quad -\infty$ to $+\infty$

- population $\quad \sigma_{xy} = \dfrac{\sum_{i=1}^{N}(x_i - \mu_x) \times (y_i - \mu_y)}{N}$

- Sample $\quad S_{xy} = \dfrac{\sum_{i=1}^{N}(x_i - \bar{x}) \times (y_i - \bar{y})}{n-1}$

Covariance may be: $\quad >0$, the 2 move together
$\qquad\qquad\qquad\qquad = 0$, the 2 are independent
$\qquad\qquad\qquad\qquad < 0$, the 2 move opposite

## Correlation Coefficient: It adjusts covariance, so that the relationship b/w 2 variables became easy and intuitive to interpret.

population $\qquad \dfrac{\sigma_{xy}}{\sigma_x \ast \sigma_y}$ $\qquad\qquad$ Correlation doesn't imply causation

Sample $\qquad \dfrac{S_{xy}}{S_x \ast S_y}$ $\qquad\qquad \left( \overset{}{-1} \longleftrightarrow 0 \longleftrightarrow +1 \right)$

$\qquad\qquad\qquad\qquad\qquad$ (min) $\qquad\qquad$ (max),
$\qquad\qquad\qquad\qquad$ move opposite $\mid$ move 2gether
$\qquad\qquad\qquad\qquad\qquad\qquad$ independent

Perfect positive correlation:
$$\dfrac{Cov.(x,y)}{Stdev(x) \times Stdev(y)} \quad = \quad \dfrac{Cov(y,x)}{Setdev(y) \ast Stdev(x)}$$

## Distribution: It is a function that shows the possible values for a variable, and how often they occur.

Leptokurtic ↘

### Normal distribution



$N \sim (\mu, \sigma^2)$
mean  variance
Mesokurtic
↓Platykurto

mean
mode
median

no skew
symmetrical
Bell curve
Gaussian distribution

# Standard normal distribution (z):

It is a special case of normal distribution

$$N \sim (\underset{\mu}{0}, \underset{\sigma}{1})$$

→ How to convert Normal Distribution to standard ND?

Z-score = $\dfrac{\text{Original value} - \text{Mean}}{\text{Standard deviation}}$ or $\dfrac{x - \mu}{\sigma}$

## Central limit theorem:

It states that if you have a population with mean $\mu$ and standard deviation $\sigma$ and take sufficiently large random samples from the population with replacement, then the distribution of the sample means will be approx normally distributed

## Standard Error:

The std. deviation of the distribution formed by sample means $\mu$

$$\underset{\text{(sample)}}{N} \sim \left(\mu, \left(\dfrac{\sigma^2}{n}\right)\right)$$

→ known population variance
→ sample size
mean   sample variance

Sample Variance = $\dfrac{\sigma^2}{n}$   ⇒   std. dev = $\sqrt{\dfrac{\sigma^2}{n}}$   ⇒   $\dfrac{\sigma}{\sqrt{n}}$

error
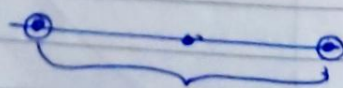
★ Std error decreases as sample size increases

## Estimator

It is an approximation depending solely on sample information

(i) point estimate

——————•————————

(ii) confidence interval estimate

⊙————•————⊙

interval

It is the range within which you expect the population parameter to be

Estimator of Parameter → Estimate
/how to estimate/ /what to estimate/   /concrete result/

Property $\longrightarrow$ Bias
$\downarrow$

Efficiency     An unbiased estimator has an
             expected value equal to the
             population parameter

The most efficient estimator is the unbiased estimator
with the smallest variance.

T score   $t_{n-1, \alpha} = \dfrac{\bar{x} - \mu_0}{s/\sqrt{n}}$

degree of freedom

matery

A t-score is one form of a standardized test statistic
Used when: (i) Sample size less than 30
           (ii) Has an unknown population standard variance.

$$\left[ \begin{array}{l} Z \sim N(0, 1) \\ Z \sim N(\bar{x} - \mu_0, 1) \end{array} \right] \begin{array}{l} \text{critical value} \\ \text{Z-score} \end{array}$$

mean          $\sigma$

it is standardized variable associated with the test.

                                    variance
confidence intervals, population unknown, t score

$$\left[ \bar{x} - t_{n-1, \alpha/2} \dfrac{s}{\sqrt{n}} , \quad \bar{x} + t_{n-1\alpha/2} \dfrac{s}{\sqrt{n}} \right]$$

                          variance
confidence intervals, population known, Z-score.

$$\left[ \bar{x} - Z_{\alpha/2} \dfrac{\sigma}{\sqrt{n}} , \quad \bar{x} + Z_{\alpha/2} \dfrac{\sigma}{\sqrt{n}} \right] \begin{array}{l} \text{reliability factor} \\ (t, z) \end{array}$$

                                    Margin of error

$\bar{x} \pm ME$

confidence level $= 1 + \alpha$

mean $\pm$ z score $\times$ std error
t-score

# Steps in data-driven decision-making

(i) Formulate a hypothesis

(ii) Find the right test

(iii) Execute the test

(iv) Make a decision

**Hypothesis** : "A idea that can be tested"

| Hypothesis | Notation | |
|---|---|---|
| Null | $H_0$ | (one to be tested) |
| Alternative | $H_1$ or $H_A$ | (everything else) |

**Significance level** : $\alpha$, It is probability of rejecting the null hypothesis, if it is true

**Error Types** : Types I
- False +ve
- Rejecting true null hypothesis.
- Probability of occuring $= \alpha$
- More dangerous

Type II
- false -ve
- Accepting false null hypothesis
- Less dangerous
- Probability of occuring : $\beta$

($\beta$ : depends mainly upon $\rightarrow$ sample size
$\rightarrow$ population variance)

# Power of test : $1 - \beta$

Rejecting false null hypothesis or,
Accepting true null hypothesis

**Testing :** It is done by standardizing the variable at hand and comparing it to the $z$ (critical value)

**Reject if :** Absolute values of $Z$-score $> z$

**p-value :** p-value is the smallest level of significance at which we can still reject the null hypothesis given the observed sample statistic.

Reject if $\alpha_p > p_\alpha$     Accept if $p > \alpha$     (evidence against null hypothesis)
(smallest value → strongest evidence)

(Previous)(Topic) Confidence Interval Estimator ; 2 sample.

## (i) Dependent sample

$$\bar{d} \pm t_{n-1, \alpha/2} \frac{S_d}{\sqrt{n}} \qquad \text{diff}^{\text{diff}} \text{ mean} + \underline{\text{std error} \times t\text{-score}}$$
(difference mean)

## (ii) Independent ⟶ Population Variance known

$$(\bar{x} - \bar{y}) \pm Z_{\alpha/2} \sqrt{\frac{\sigma^2_x}{n_x} + \frac{\sigma^2_y}{n_y}}$$

difference point
estimator
mean difference          variance of the difference

## (iii) Independent ⟶ PV not known ⟶ assumed equal

std dev

$$\text{pooled sample variance } (S^2p) = \frac{(n_x - 1) S^2_x + (n_y - 1) S^2_y}{n_x + n_y - 2}$$

$$(\bar{x} - \bar{y}) \pm t_{n_x + n_y - 2, \alpha/2} \sqrt{\frac{S^2_p}{n_x} + \frac{S^2_p}{n_y}}$$
mean difference

{ As standard normal distribution is symmetrical around 0, the 2 statement are equivalent:
$-4.67 < $ a negative $z$   $<=>$   $4.67 > $ a positive $z$ }

(iv) Independent $\longrightarrow$ $\sigma$ unknown $\longrightarrow$ assumed different

$$(\bar{x} - \bar{y}) \pm t_{v, \alpha/2} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

$$v = \frac{\left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}\right)^2}{\left[\frac{s_x^2}{n_x}\right]^2 \Big/ (n_x - 1) + \left[\frac{s_y^2}{n_y}\right]^2 \Big/ (n_y - 1)}$$

$\left\{ \begin{array}{l} \circ \text{ pvalue (one-sided)} = (1 - \text{value in table}) \\ \circ \text{ p-value (two-sided)} = (1 - \text{value in table}) \times 2 \end{array} \right\}$

• The closer to 0.000 the p-value, the better. (We rejects null hypothesis at all and uncommon)
• P-value is a normal universal concept that works with every distribution

## Linear Regression-
A linear regression is a linear approximati of a causal relationship b/w 2 or more variable.

Process :  Get sample data
$\downarrow$
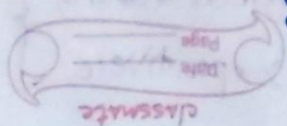
Design a model that works for that sample
$\downarrow$

Make prediction for whole population

population formula  $\left(\boxed{y = \beta_0 + \beta_1 x_1 + \varepsilon}\right) \rightarrow$ Error

Estimated or Predicted value $\quad$ constant $\quad$ Independent variable

$$\hat{y} = b_0 + b_1 x_1$$

| Corelation | Regression |
|---|---|
| * Reletionship | * One variable affects the other |
| * Not capture Causality | * Cause and effect |
| * $p(x,y) = p(y,x)$ | * One way |
| * single point | * line |



Sum of Square Total :  SST  $\sum_{i=1}^{n} (y_i - \bar{y})^2$

observed dependent variable / mean
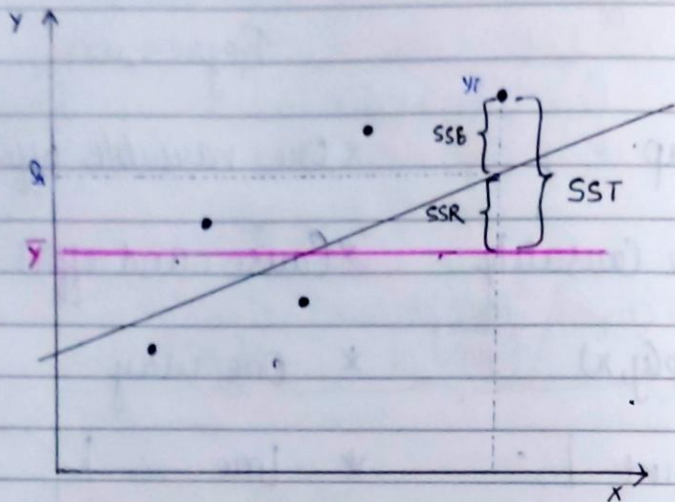
Sum of Square Regression : SSR  $\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$

predicted value

Sum of Square Error : SSE  $\sum_{i=1}^{n} e_i^2$

$$SST = SSE + SSR$$

$$\sum_{i=1}^{n} (y_i - \bar{y})^2 = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^{n} e_i^2$$

$$(SST) = (SSR) + (SSE)$$

| Total<br>variability | $=$ | Explained<br>variability | $+$ | Unexplained<br>variability |

$$R^2 = \frac{SSR}{SST}$$

shows how much of the total variability of the dataset is explained by our regression model

```
0 ─────────────────────────── 1
```

your regression explains NONE of the variability        your regression explains the entire variability

⭢ How to find Regression line?

Ordinary least squares (OLS)

min SSE

$S(b)$ is the OLS estimator of $\beta$ for a simple linear regression

$$S(b) = \sum_{i=1}^{n} (y_i - x_i^T b)^2 = \underline{(y - Xb)^T (y - Xb)}$$

linear algebra

$\boxed{\text{Other methods}}$ ⌐

- Bayesian method regression
- Kernal method regression
- Gaussian process regression

## Confidence Interval

**1 data set**

(1) Population Variance Known (z score)

(2) Population Variance unknown (t-score)

**2 data set**

(1) Dependent Set (t-score)

(2) Independent set
  (i) PV Known (z-score)
  (ii) PV unknown (t-score)

assumed equal (t-score)   assumed unequal (t score)