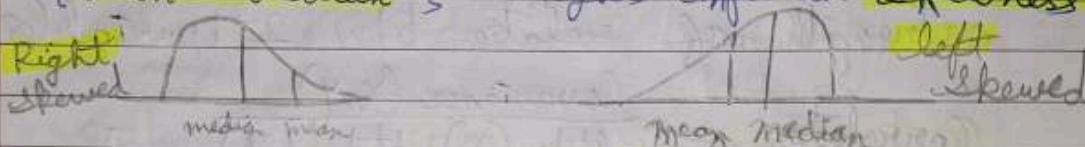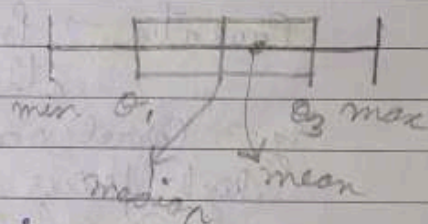<u>Descriptive Statistics</u> → to describe data

a) Measures of <u>central tendency</u> → mean, median, mode
median more useful in data with outliers
mode not significant for continuous data (high PDF)
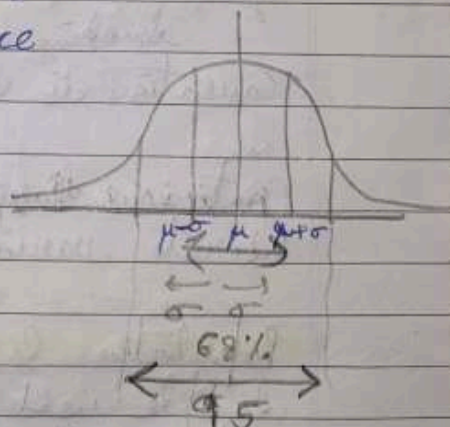
a) { Mean − Median } → gives info on <u>skewness</u> (point)

<u>Right</u> skewed

<u>left</u> skewed

median mean                mean median

a) Measures of <u>dispersion / spread</u> :
Range, I Q R (Inter Quartile Range) $= Q_3 - Q_1$,
↳ max − min

a) <u>Box Plot / Whiskers</u>
<u>Box and Whiskers Plot</u>

min. $Q_1$        $Q_3$ max
median    mean

a) More measures of <u>dispersion</u> :
Standard Deviation, variance

a) <u>68 − 95 − 99.7 rule</u> →
Valid on bell shaped data
(Perfectly apply on normal
distribution)

$\mu - \sigma$   $\mu$   $\mu + \sigma$
$\sigma$   $\sigma$
68%
95

a) <u>Chebyshev's Theorem</u>
Atleast $\left(1 - \frac{1}{k^2}\right)^{th}$ of data lies within
$\pm K$ standard deviations regardless of shape
of distribution
(75% data lies b/w $\mu - 2\sigma$ & $\mu + 2\sigma$)

Measures of association → Covariance, correlation, Causation

$$\text{Covariance} = \frac{1}{(n-1)} \sum_{i=1}^{n} (x_i - \bar{x})(Y_i - \bar{Y})$$

Covariance changes with unit of data hence we use correlation. Correlation measures strength of relation b/w 2 variables

$$\text{Correlation} = \frac{\text{Covariance }(X,Y)}{\text{Stdev}(X)\;\text{Stdev}(Y)}$$

Correlation $\in [-1, 1]$

Causation → proving that 1 variable causes other. Not the same as correlation

eg → Correl (occupancy, hotel prices) = + high but high prices → ⊘ high occupancy

eg → Correl (Smoking, Cancer) = +1 & Smoking causes cancer

Causation is huge topic out of scope

Continuous variables → PDF Probability distribution
Discrete variables → PMF Prob Mass fn

Population & Sample
Sample used as usually we can't access whole population, its cost effective & feasible

Sampling techniques such that sample is good representative of population → whole field in marketing
Random Sampling

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mu$ = Population mean
$\sigma$ = Population Std Dev
$\bar{x}$ = Sample Mean
$s$ = Sample Std Deviation

## Central Limit Theorem

Sample mean is normally distributed with mean equal to population mean, irrespective of distribution type of population (be it unimodal, multimodal, symetric, skewed, discrete, continuous)

$$\bar{x} \sim Normal\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

Bernoulli Process / Trial → 2 outcomes only → win or lose
Binomial distribution → n trials (independent)

success = $p$    failure = $1-p$

Random variable $X$ = no of successes in $n$ trials

$$P(X = x) = \frac{n!}{x!\,(n-x)!} p^x (1-p)^{1-x} = {}^nc_x \, p^x (1-p)^{1-x}$$

Binomial Distribution    mean = $np$    Variance = $npq$

Eg → No of fraud reports among $n$ tax reports, no of students passing exam
here $x \in [0, n]$    (Unlike Poisson)

Poisson Distribution    $$P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

Mean = $\lambda$ = Variance    $x \in [0, 1, 2 \dots \infty]$

$T$ | Student T distribution → symetric, centered at 0, only 1 parameter → doF (Degrees of freedom)

As DOF → $\infty$, T distribution → Standard Normal Distribution

Confidence interval

$z$ statistic $= \dfrac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim$ Normal $(0, 1)$

$T$ statistic $= \dfrac{\bar{X} - \mu}{s / \sqrt{n}} \sim t_{n-1}$

$\sigma \to$ population std dev   $s \Rightarrow$ Sample std dev

Constructing confidence interval $\to$

$\bar{X} - |z_{\alpha/2}| \dfrac{\sigma}{\sqrt{n}} < \mu < \bar{X} + |z_{\alpha/2}| \dfrac{\sigma}{\sqrt{n}}$   $z_{\alpha/2}$

$\to [(1-\alpha)$ confidence interval for population mean $]$

$[$ for $\alpha$ confidence interval, $\left| z_{\frac{1}{2} - \frac{\alpha}{2}} \right| \dfrac{\sigma}{\sqrt{n}}$ margin

which makes intuitive sense $]$

If $\sigma$ not known use $s$ &

$\pm |t_{\alpha/2}| \dfrac{s}{\sqrt{n}} =$ margin   $z_{\frac{1}{2} - \frac{\alpha}{2}}$

for Confidence Interval for population proportion

$\hat{P} - |z_{\alpha/2}| \sqrt{\dfrac{\hat{P}(1-\hat{P})}{n}} < P < \hat{P} + |z_{\alpha/2}| \sqrt{\dfrac{\hat{P}(1-\hat{P})}{n}}$

(we don't use T distribution for population proportion)
$p \to$ population proportion
$\hat{P} \to$ Sample proportion

Sample size, how big to take? $\to$
Use that industry rule of thumb
or get % confidence needed, % tolerance (margin of error
put in above equ get $n$

==Hypothesis testing== → ==Null Hypothesis H_0==, Alternate
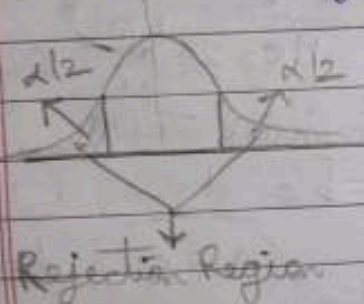                                               Hypothesis H_A

==Steps→==

i) Formulate Hypothesis $H_0$ & $H_A$
ii) Calculate T statistic $= (\bar{x} - \mu)/(s/\sqrt{n})$
iii) Cutoff value for T-statistic ($\alpha$ = significance level)
iv) Check whether T-stat falls in the rejection region

Then conclusion → accept or reject $H_0$

2) Types of $H_0$

$\mu = \ldots$                                    $\mu \geq \ldots$              $\mu \leq \ldots$

==Two Tailed Hypothesis test==        ↓                          ↓
                                          Rejection region        Rejection region
                                                                          on RHS

$\alpha/2$                $\alpha/2$       Rejection region              ↓
                                          on LHS →
                                                   ==One Tailed==
                                                   Hypothesis tests

Rejection Region

Null hypothesis cannot have $\{<, >, \neq\}$

==Type I errors :== False +ves (Rejecting $H_0$ when it is true)
      $\alpha$ = Probablity of type I error
==Type II errors :== False -ve (Not rejecting $H_0$ when its false)
      $\beta$ = probablity of type II error

Other typical types of hypothesis tests
==Difference in means test==
  ↳ 2 versions → with or without equal variance assumption
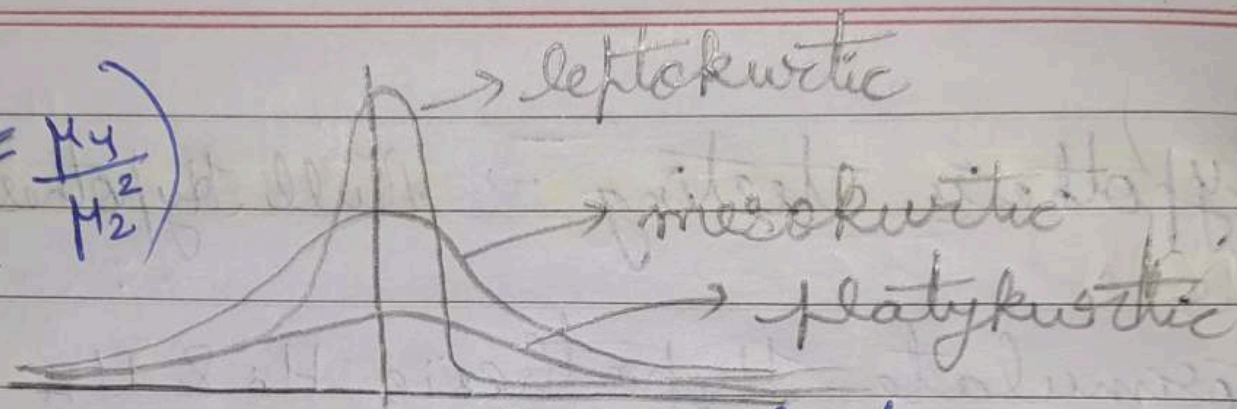==Paired T test== for diff in means
Diff in means tests → only for 2 populations not multiple

Kurtosis $\left(= \dfrac{\mu_4}{\mu_2^2}\right)$

→ leptokurtic

→ mesokurtic

→ platykurtic

Mesokurtic → Normal distribution → Kurtosis =

$\mu_4 = 4^{th}$ central moment