# Maths final notes

**Statistics:**  a branch of mathematics dealing with the collection, analysis, interpretation, and presentation of masses of numerical data.

**Random variable** is a variable whose value is unknown or a function that assigns values to each of an experiment's outcomes.

**Joint PD:** is a statistical measure that is used to calculate the probability of 2 events occurring together at the same time

**Marginal PD:** gives the probability of various values of the variable in the subset without reference to the value of the other variable

**Probability mass function (PMF)** is a function that gives the probability that a discrete random variable is exactly equal to some value. Sometimes it is also known as the discrete density function

**probability density function (PDF)**, or density of a continuous random variable,  is used to specify the probability of the random variable falling within a particular range of values, This probability is given by the integral of this variable's PDF over that range

| Difference between Mutually exclusive and independent events | |
| --- | --- |
| **Mutually exclusive events** | **Independent events** |
| When the occurrence is not simultaneous for two events then they are termed as Mutually exclusive events. | When the occurence of one event does not control the happening of the other event then it is termed as an independent event. |
| The non-occurrence of an event will end up in the occurrence of an event. | There is no influence of an occurrence with another and they are independent of each other. |

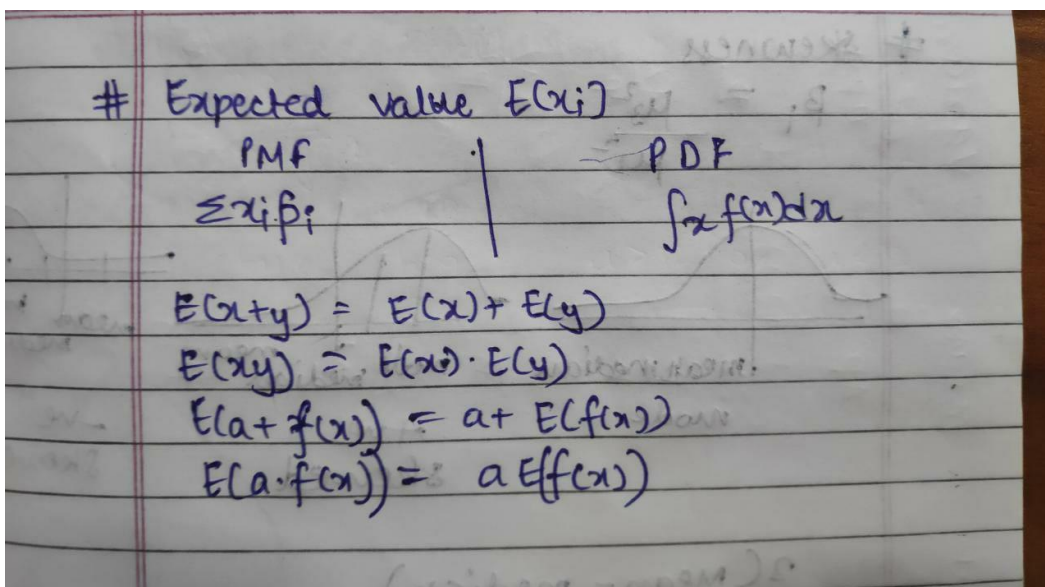| The mathematical formula for mutually exclusive events can be represented as P(X and Y) = 0 | The mathematical formula for independent events can be defined as P(X and Y) = P(X) P(Y) |
|---|---|
| The sets will not overlap in the case of mutually exclusive events. | The sets will overlap in the case of independent events. |

**Exhaustive events :** In probability, a set of events is collectively exhaustive if they cover all of the probability space:

**Discrete random variable:** A discrete random variable has a countable number of possible values. The probability of each value of a discrete random variable is between 0 and 1, and the sum of all the probabilities is equal to 1. It has PMF(probability mass function)

**Continous random variable:** Continuous random variables, on the other hand, take on values that vary continuously within one or more real intervals, and have a cumulative distribution function (CDF) that is absolutely continuous. It has pdf (probabilty distribution function)

**Mathematical expectation:** Mathematical expectation, also known as the expected value, which is the summation of all possible values from a random variable.
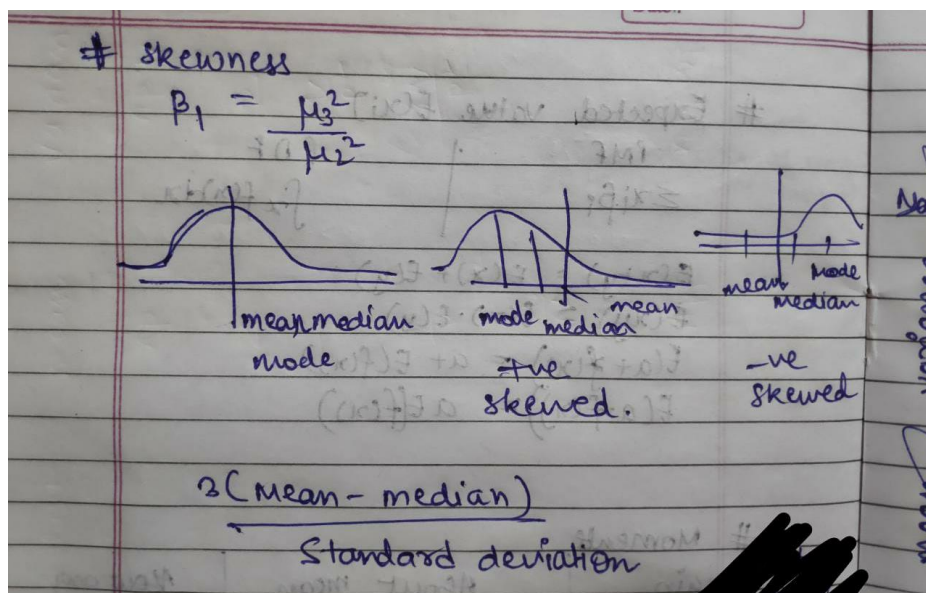
It is also known as the product of the probability of an event occurring, denoted by P(x), and the value corresponding with the actually observed occurrence of the event.
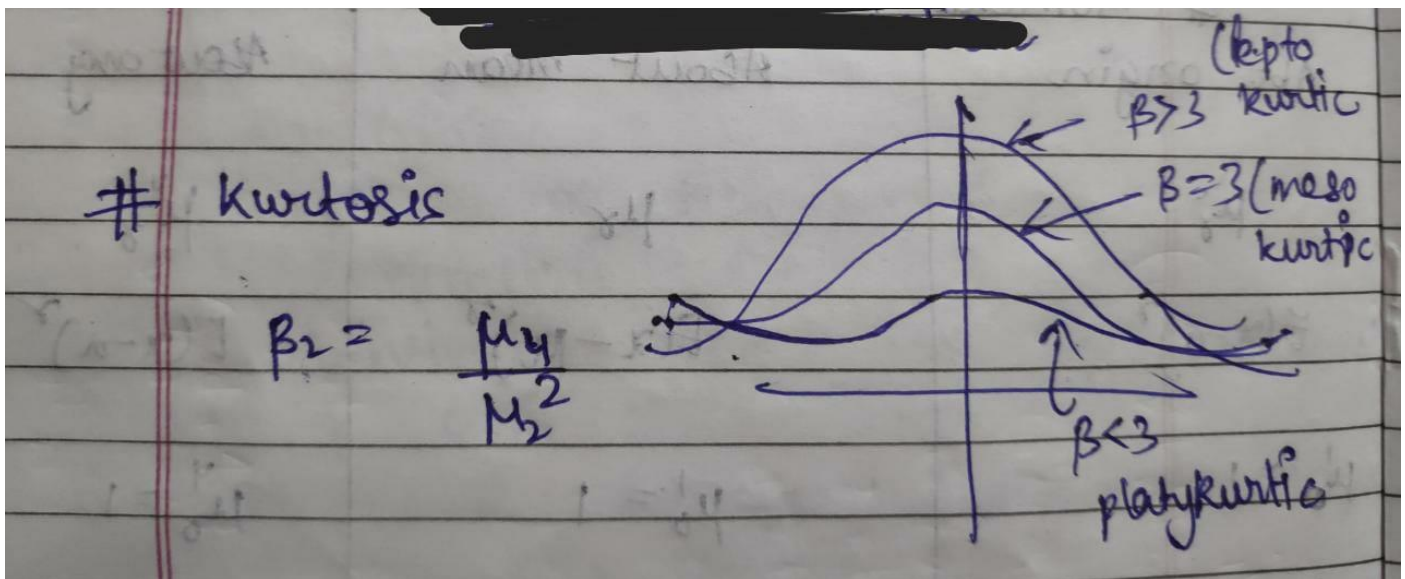
**Moments**

# Moments

| Abt origin | About mean | About any |
|---|---|---|
| $\mu_r'$ | $\mu_r$ | $\mu_r'$ |
| $E(x-0)^r$ | $E(x-\mu)^r$ | $E(x-a)^r$ |
| $\mu_0' = 1$ | $\mu_0' = 1$ | $\mu_0' = 1$ |
| $\mu_1' = E(x) = $ mean | $\mu_1 = 0$ | |
| $\mu_2' = E(x^2)$ | $\mu_2 = $ variance <br> $\sigma^2 = E(x^2) - (E(x))^2$ | |

**Skewness:** Skewness is a measure of symmetry, or more precisely, the lack of symmetry. A distribution, or data set, is symmetric if it looks the same to the left and right of the center point.

# skewness

$$\beta_1 = \frac{\mu_3^2}{\mu_2^3}$$

mean median mode

mode median +ve skewed.

mean mode median −ve skewed

$$\frac{3(\text{mean} - \text{median})}{\text{Standard deviation}}$$

Kurtosis:

Kurtosis is a measure of whether the data are heavy-tailed or light-tailed relative to a normal distribution. That is, data sets with high kurtosis tend to have heavy tails, or outliers. Data sets with low kurtosis tend to have light tails, or lack of outliers. A uniform distribution would be the extreme case.



## Change of origin in mgf

Changes the mean but variance remains same

## Change in scale in mgf

Changes the mean and variance

## Weak law of large numbers:

**Central limit theorem:**

The central limit theorem states that if you have a population with mean μ and standard deviation σ and take sufficiently large random samples from the population **with replacement**, then the distribution of the sample means will be approximately normally distributed.

**Population**: it is the collection of all items of interest to our study and is usually denoted by an uppercase N, the numbers we have obtained when using a population are called parameters (u and sigma)

**Sample:**

Sample :
          The selection of some section from the whole population is known as Sample.

**Sampling terminology**

Sampling Terminology

- Total number of items / population,
  Population Size $= N$.
- Mean of the population,
  Population Mean $(\mu) = \dfrac{\Sigma X}{N}$
- Variance of the Population,
  Population Variance $(\sigma^2) = \Sigma(X_i - \mu)^2 / N$
- Number of items / Population,
  Sample Size $= n$.
- Mean of the sample,
  Sample mean $(\bar{x}) = \Sigma x / n$.
- Variance of the Sample,
  Sample Variance $(s^2) = \Sigma(x_i - \bar{x})^2 / n$

## Sampling:

sampling : process of drawing sample
Fundamental assumption.
↳ random sampling

## Simple sampling:

simple sampling : A special case of
random sampling in which each event
has the same probability of success
and the chance of success for diff
events are independent.

## Parameters and staistics:

Statistical constants of the $pop^n$ suchas
$\mu$ & $\sigma$ are parameters.

↳ constants of sample → $\bar{x}$, & $S(\sigma)$
are called statistics
(Greek → population , Roman → statistics)

## Objective of sampling:

Objectives of sampling
↳ max info abt $pop^n$
↳ min effort, cost, time
↳ find best possible value of the paramet-
er under specific $cond^n$.
↳ Sampling determines the reliability of
these estimates.

**Sampling error:**

Sampling error is the difference between the sample measure, and the corresponding population measure due to the fact that the sample is not a perfect representation of the population.

**Properties os sampling means and variance**

Properties of the Distribution of Sample Means.

- The mean of the sample means will be the same as the population mean.

The standard deviation of the sample means will be smaller than the standard deviation of the population, and will be equal to the population standard deviation divided by the square root of the sample size.

**Satistical interference**

# Statistical interference is the logic of sampling theory is the logic of induction in which we pass from a particular (sample) to general (population). Such generalis- ation from sample to population is statistical interference.

## sampling distribution

A sampling distribution is the distribution of frequencies of mean of all possible size. Samples of size $n$ is sampling distribution of means (similar for $s$)
( after drawing each sample we put back the sample so that pop$^n$ remain)

## Standard error

\# Standard error (SE)
↳ Standard deviation of sampling distribution.
↳ used to assess the difference b/w expected and observed values.
↳ $\frac{1}{SE}$ = precision. (reciprocal)

↳ $n \geqslant 30$ (large sample)
$n < 30$ (small sample)
↳ sampling distribution of large sample ≃ normal (assumption)

**Statistical hypothesis**

(SH)

# Statistical hypothesis : assumption about the population which may or may not be true are called SH.

**Testing hypothesis**

# Testing hypothesis : The methode consists in assuming the hypothesis is correct and then computing the p of getting the observed sample.
↳ If this p is less than a preassigned value then hypothesis is wrong.

**Errors type1 and type2**

#

Errors

| Type I | Type II |
|---|---|
| rejected while it should be accepted | accepted to while it should be rejected |
| aims at limiting to preassigned value | aims to minimise |

best way to minimise both is to increase
sample size (n)

# Null hypothesis

# Null hypothesis: The hypothesis formulated for the sake of rejecting it under the assumption that it is true.
↳ denoted by $H_0$.

↳ By accepting null hypothesis, we mean "that on the basis of statistic calculated from the sample we do not reject it".

↳ acceptance ↛ true
  rejection ↛ false.

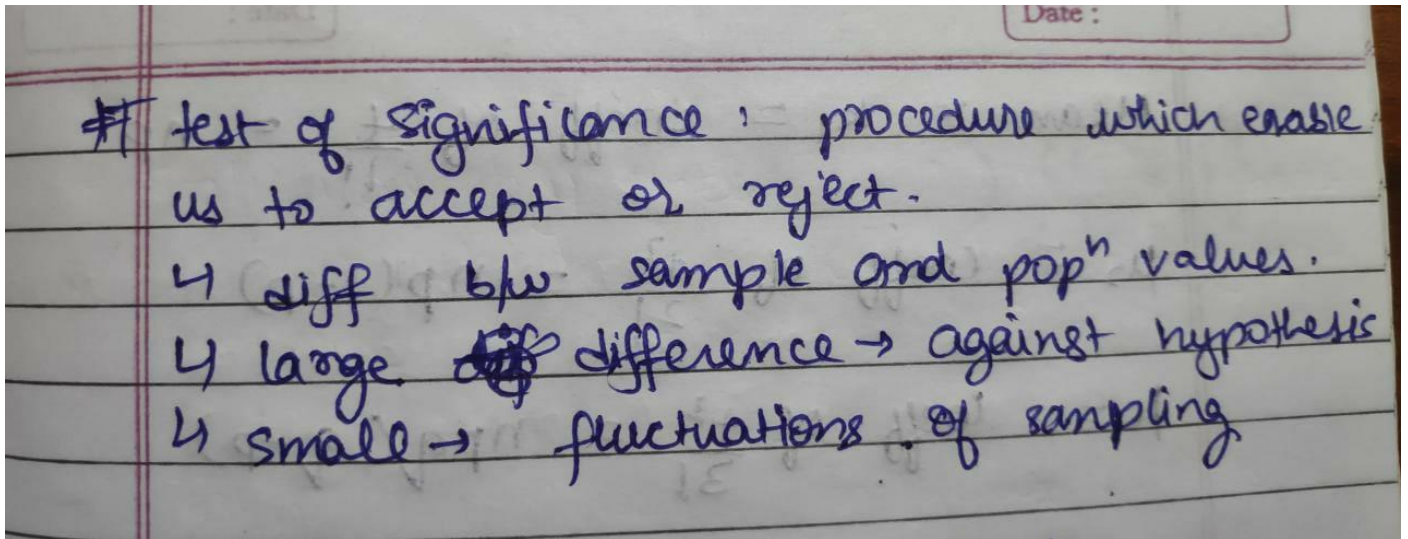## Level of significance and critical region

# level of significance:
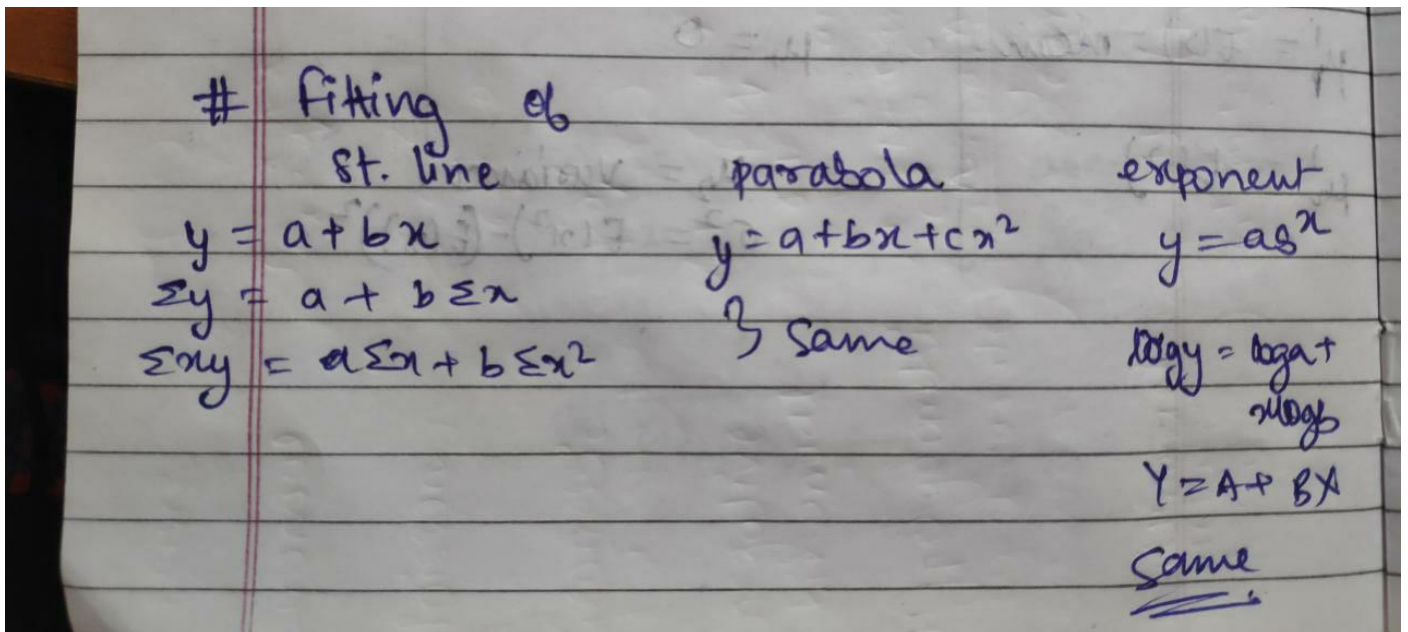↳ the p level below which we reject hypothesis

# critical region: the region in which sample value falling is rejected
generally take 2 regions (5% and 1% of area of normal curve)

area of both sides are considered → double tail. eg: coin toss
area of 1 side (right) → single tail.

## Test of significance

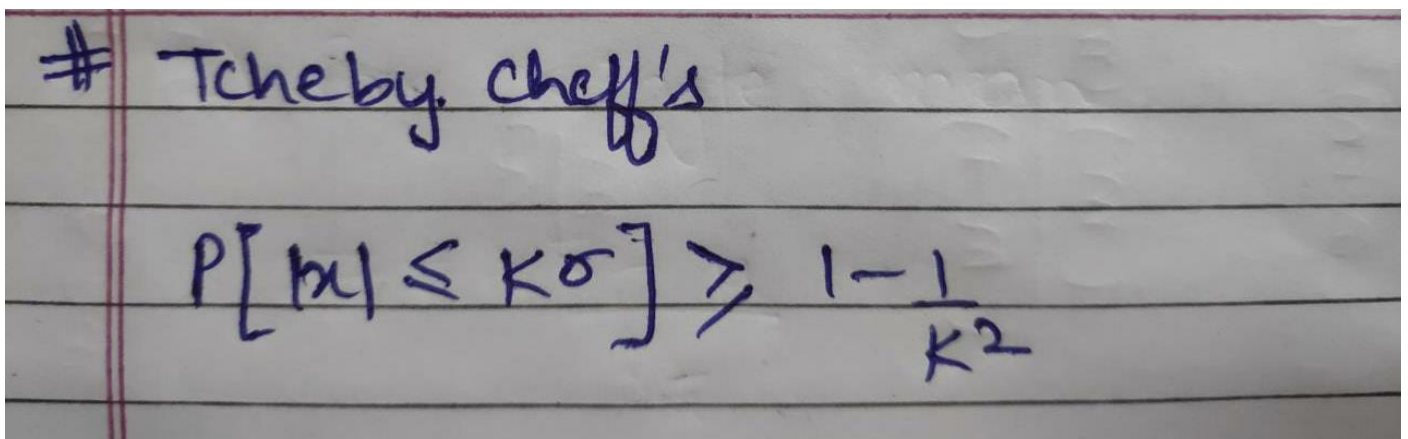#1 test of Significance : procedure which enable us to accept or reject.
↳ diff b/w sample and pop$^n$ values.
↳ large difference → against hypothesis
↳ small → fluctuations of sampling

## Fitting to straight line, parabola, exponential

# Fitting of

st. line

$y = a + bx$
$\Sigma y = a + b\Sigma x$
$\Sigma xy = a\Sigma x + b\Sigma x^2$

parabola

$y = a + bx + cx^2$
3 same

exponent

$y = ag^x$

$\log y = \log a + x\log g$

$Y = A + BX$

Same

## Tchebycheffs

# Tcheby. cheff's

$$P\left[|x| \leq K\sigma\right] \geq 1 - \frac{1}{K^2}$$

**Markov**

$$P(x \geq a) \leq \frac{E(x)}{a}$$

Markov's inequality

**Attribute sampling** refers to a statistical sampling tool used by the auditors to analyze the features of a particular population

**Confidence limits** for the mean are an interval estimate for the mean. Confidence limits tell you how accurate your estimate of the mean is likely to be.

 **The chi-square** compares the size any discrepancies between the expected results and the actual results, given the size of the sample and the number of variables in the relationship.

**Degrees of Freedom** refers to the maximum number of logically independent values, which are values that have the freedom to vary, in the data sample.

 **Goodness of fit**   test helps you see if your sample data is accurate or somehow skewed

**The significance level**, also denoted as alpha, is a measure of the strength of the evidence that must be present in your sample before rejecting the null.

In a **test of independence,** we state the null and alternative hypotheses in words. Since the contingency table consists of two factors, the null hypothesis states that the factors are independent and the alternative hypothesis states that they are not independent (dependent).

The **contingency coefficient** is a coefficient of association that tells whether two variables or data sets are independent or dependent of each other

**Yates' correction** is to prevent the overestimation of statistical significance for small data when 'zero cells' are present in a 2 × 2 contingency table.

**Point estimation**, in statistics , the process of finding an approximate value of some parameter—such as the mean (average)—of a population from random samples of the population

**interval estimation** is the use of sample data to calculate an interval of possible values of an unknown population parameter;

**Maximum likelihood** estimation is a method that determines values for the parameters of a model. The parameter values are found such that they maximise the likelihood that the process described by the model produced the data that were actually observed.

| Test statistic | Null and alternative hypotheses | Statistical tests that use it |
|---|---|---|
| t-value | Null: The means of two groups are equal<br><br>Alternative: The means of two groups are not equal | T-test<br><br>Regression tests |
| z-value | Null: The means of two groups are equal<br><br>Alternative:The means of two groups are not equal | Z-test |
| F-value | Null: The variation among two or more groups is greater than or equal to the variation between the groups<br><br>Alternative: The variation among two or more groups is smaller than the variation between the groups | ANOVA<br><br>ANCOVA<br><br>MANOVA |
| X2-value | Null: Two samples are independent<br><br>Alternative: Two samples are not independent (i.e. they are correlated) | Chi-squared test<br><br>Non-parametric correlation tests |

**Chi sqaure variate**

# Chi square

$$\chi^2 = \sum_{i=0}^{n} \left( \frac{(O_i - E_i)^2}{E_i} \right)$$

d.f of Binomial = $n-1$

df of poisson = $n-2$

df of normal = $n-3$

**Test of significance for large samples:**

1) For a signle proportion

# Test of significance for large samples

(1) for single proportion.

$\mu = np$

$\sigma^2 = npq$

$z = \dfrac{x - \mu}{\sigma}$

$z < 1.96$ w/c (5% significant)

$z < 2.58$ (1% significant)

2) For difference between proportion

(2) for diff b/w proportional

$$z = \frac{P_1 - P_2}{\sqrt{PQ\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2}$$

$$Q = 1 - P$$

1.645 → 5%.

2.33 → 1%.

## 3) For single mean



Page No.

Date :

3) for single mean

$$z = \frac{\bar{x} - \mu}{?}$$

$\mu \rightarrow$ mean of pop$^n$
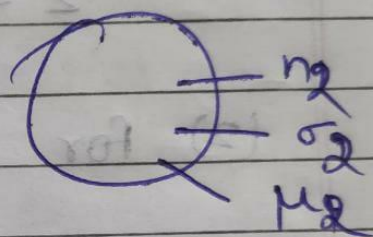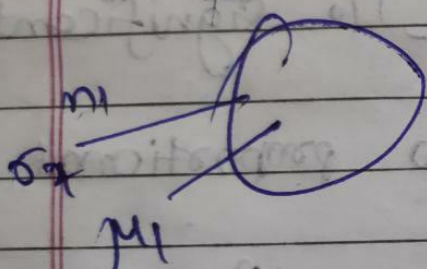
$\sigma_{\bar{x}} \rightarrow$ standard error of means

$\sigma_{\bar{x}} \rightarrow \dfrac{\sigma}{\sqrt{n}}$     $\sigma \rightarrow$ var of pop$^n$

$\rightarrow \dfrac{S}{\sqrt{n}}$     $S \rightarrow$ var of sample

## 4) Difference between 2 proportionate

4) diff b/w 2 proportionate



$$z = \frac{\mu_1 - \mu_2}{\sigma_{\bar{x}_1 - \bar{x}_2}}$$

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad \text{2nd} \left.\begin{array}{l} \sigma_1, \sigma_2 \\ pop^n \end{array}\right\}$$

**Gamma function properties**

# Gamma function properties

$$\frac{\overline{|a}}{\lambda^a} = \int_0^\infty x^{a-1} e^{-\lambda x} \, dx$$

$$\overline{|a+1} = a\overline{|a}$$

$$\overline{|\tfrac{1}{2}} = \sqrt{\pi}$$

$$\overline{|n} = n!$$

# Distributions

| Name | $f(x)$ | MGF | cumulant | mean | variance |
|---|---|---|---|---|---|
| Binomial | $^nC_x p^x q^{n-x}$ | $(q+pe^t)^n$ | | $np$ | $npq$ |
| Poisson | $\dfrac{z^r e^{-z}}{r!}$ | | | $z$ | $z$ |
| Geometric | $pq^n$ | | | $\dfrac{1}{p}$ | $\dfrac{q}{p^2}$ |
| negative binomial | $^{x-1}C_{r-1}\, p^r q^{x-r}$ | | | $\dfrac{r}{p}$ | $\dfrac{rq}{p^2}$ |
| Hypergeometric (Multivariate) (extension) | $\dfrac{^aC_n\,^{N-a}C_{n-x}}{^NC_m}$ | | | $\dfrac{na}{N}$ | |
| Multinomial $P(X_1=x_1, X_2=x_2 \cdots)$ | $\dfrac{n!}{x_1!\cdot x_2!\cdots x_n!}\left(p_1^{x_1} p_2^{x_2} p_3^{x_3} \cdots p_n^{x_n}\right)$ | | | $np_i$ | $np_i\, q_i$ |

| Name | $f(x)$ | MGF | | mean | variance |
|---|---|---|---|---|---|
| Normal distribution | $\dfrac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$ | $e^{\mu t + \frac{t^2 \sigma^2}{2}}$ | | $\mu$ | $\sigma^2$ |

$r^{th}$ moment

| Name | $f(x)$ | | mean | variance |
|---|---|---|---|---|
| $B_1$ | $\begin{cases}\dfrac{x^{m-1}(1-x)^{n-1}}{B(m,n)} & x\leqslant 0 \\ & m,n>0 \\ 0 & \text{otherwise}\end{cases}$ | $\dfrac{\left(\overline{m+r}\right)\left(\overline{m+n}\right)}{\left(\overline{m+r+n}\right)\left(\overline{m}\right)}$ | $\dfrac{m}{m+n}$ | $\dfrac{mn}{(m+n)^2(m+n+1)}$ |
| $B_2$ | $\begin{cases}\dfrac{x^{m-1}}{B(m,n)(1+x)^{m+n}} & x\geqslant 0 \\ & m,n>0 \\ 0 & \text{otherwise}\end{cases}$ | $\dfrac{\left(\overline{m+r}\right)\left(\overline{n-r}\right)}{\overline{m}\;\overline{n}}$ | $\dfrac{m}{n-1}$ | $\dfrac{(m)(m+n-1)}{(n-1)^2(n-2)}$ |
| $\gamma$ | $\begin{cases}\dfrac{a^m x^{m-1} e^{-ax}}{\overline{m}} & x\geqslant 0 \\ 0 & \text{otherwise}\end{cases}$ | | $\dfrac{m}{a}$ | $\dfrac{m}{a^2}$ |

Skewness $\to \dfrac{4}{m}$

kurtosis $\to \dfrac{3}{(m)\overline{16}}$ . $\dfrac{3m+6}{m}$

$m \to$ shape parameter
$a \to$ rate parameter
$\dfrac{1}{a} \to$ scale parameter

Exponential $\theta \cdot e^{-\theta x}$    $\dfrac{(\theta)}{(\theta - t)}$    $\dfrac{t}{\theta} + \dfrac{t^2}{2\theta^2} + \dfrac{t^3}{3\theta^3} + \dots$    $\dfrac{1}{\theta}$    $\dfrac{1}{\theta^2}$

uniform $f(x) = \dfrac{1}{b-a}$    $\mu_r' = \dfrac{b^{r+1} - a^{r+1}}{(r+1)(b-a)}$    $\dfrac{1}{2}(b+a)$    $\dfrac{(a-b)^2}{12}$

f distribution $F(\nu_1, \nu_2)$
$$= \frac{S_1^2}{S_2^2}$$
$\nu_1 \to$ df of sample 1
$\nu_2 \to$ df of sample 2
$S_1^2 \to$ variance of S1
$S_2^2 \to$ variance of S2

     $\dfrac{\nu_2}{\nu_2 - 2}$    $\dfrac{2\nu_2^2(\nu_1 + \nu_2 - 2)}{\nu_1(\nu_2 - 2)^2(\nu_2 - 4)}$

Cauchy $\left( \dfrac{1}{\pi} \dfrac{1}{1+x^2} \right) -\infty < x < \infty$    $\mu$

0

Replace $x \to \dfrac{x - \mu}{\lambda}$