

CLASSICAL PROBABILITY

→ states that possible outcome of any event
it predicts the likeliness of an event happening.

$$P(E) = \frac{\text{no. of favourable outcomes}}{\text{tot no. of exhaustive events}}$$

→ if a random exp. or trial results of n exhaustive mutually exclusive & equally likely outcomes ; out of which m are favourable to the occurrence of an event E . $P(E) = m/n$

- LIMITATIONS :
- various outcomes of the exp. are not likely. (die roll)
 - only works well in simple well defined experiments
 - doesn't use actual data or obs., thus less reliable

STATISTICAL PROBABILITY

→ statistical representation of any random event involves data collection, analysis and interpretation
if a trial is repeated n no. of times under essentially identical conditions then $P(E) = \lim_{n \rightarrow \infty} \frac{m}{n}$.

- LIMITATIONS :
- requires large sample size, a no. of trials is done
 - ~~less~~^{no} unique events, repeated trials may occur
 - relies heavily on past data, which may not reflect future outcomes

Additive law of probability

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(B \cap C) - P(C \cap A) + P(A \cap B \cap C)$$

Multiplicative law of probability

$$P(A \cap B) = P(A) \cdot P(B|A) = P(B) \cdot P(A|B)$$

Total probability theorem

$$P(A) = \sum P(E_j) \cdot P(A|E_j)$$

Conditional probability

$$P(A|B) = P(A \cap B) / P(B)$$

Bayes' theorem

$$P(E_i|A) = \frac{P(E_i) \cdot P(A|E_i)}{\sum P(E_i) \cdot P(A|E_i)}$$

$$= P(E_i) \cdot P(A|E_i) / P(A)$$

Measures of Central Tendency

MEAN

$$\text{Sample mean: } \bar{x} = \frac{\sum x_i}{N} \quad (\text{subset of a population})$$

avg. of
set of
values

$$\text{Population mean: } \mu = \frac{1}{N} \sum x_i \quad (\text{for complete population})$$

Grouped data : $\frac{\sum f_i x_i}{\sum f_i}$

MEDIAN

middle

value

when in
Sorted data

Ungrouped data :

$$Md = \begin{cases} \frac{(N/2)^{\text{th}} \text{ obs.} + (N/2 + 1)^{\text{th}} \text{ obs.}}{2}, & N \text{ even} \\ \frac{N+1}{2}^{\text{th}} \text{ obs.}, & N \text{ odd} \end{cases}$$

Grouped data :

$$Md = l + \left(\frac{N/2 - Cf}{f} \right) \times h$$

freq.

► Median class : Class in which $N/2$ lies before cumulative freq.

l : lower limit, f : median class freq., h : width of class

Cf : cumulative freq. of preceding class.

MODE

most

frequently
occ. value

ungrouped : (counting)

$$\text{Grouped : } M_o = l + \frac{(f_1 - f_0)}{(2f_1 - f_0 - f_2)} \times h$$

Modal class : class with highest frequency

l: lower limit

f₁: freq. of modal class

h: width of class

f₀: freq. of ~~moda~~ preceding class

f₂: freq. of succeeding class

$$\text{MODE} = 3 \text{ MEDIAN} - 2 \text{ MEAN}$$

$$\text{COMBINED MEAN} \quad \bar{\bar{x}} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

Absolute dispersion

Measures of dispersion

→ expressed in terms
of original unit

RANGE : diff. b/w max and minimum values

$$\text{RANGE} = (\text{max value}) - (\text{min value})$$

VARIANCE dispersion : $\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$: population variance

of a set of values : $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$: sample variance
→ (estimate of the populatn) → BESSEL'S CORRECTION

$$\text{STANDARD DEVIATION} \quad S.D = \sqrt{\sigma^2}$$

measure of spread of data : $S.D \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$: population variance

of data : $\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$: sample variance

MEAN \rightarrow PARTITION $M.D = \frac{\sum f_i |x_i - A|}{N}$: grouped data

DEVIATION

avg. dist
blw each
data point
and mean
of dataset.

 $M.D = \frac{\sum |x_i - \bar{x}|}{N} : \text{sample}$
 $M.D = \frac{\sum |x_i - A|}{N} : \text{population}$

COMBINED VARIANCE $= \sigma_{12}^2 = \frac{\eta_1(\sigma_1^2 + d_1^2) + \eta_2(\sigma_2^2 + d_2^2)}{(\eta_1 + \eta_2)}$

$d_i = \bar{x}_i - \bar{x}$ \rightarrow combined mean
individual mean

COEFFICIENT

OF VARIAT^N $CV = \frac{\sigma}{M} \times 100 (\%) : \text{population}$

standr

measure of
dispersⁿ $CV = \frac{s}{\bar{x}} \times 100 (\%) : \text{sample}$

relative to
mean. (finding variability)

- PARTITION VALUES : divide the series into no. of equal parts

Quartile $Q_i = iN/4$ th

Decile $D_i = iN/10$ th

Percentile $P_i = iN/100$ th

for grouped data : $Q_i = l + \frac{iN/4 - Cf}{f} \times h$

$D_i = l + \frac{iN/10 - Cf}{f} \times h$

$P_i = l + \frac{iP_i/100 - Cf}{f} \times h$

ABSOLUTE DISPERSION

(disp. is exp. in terms of original units)

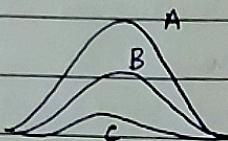
distance measure

- Range
- Interquartile Range
 $(Q_3 - Q_1)$

spread of dispersion

- mean deviation
- Standard deviation

Kurtosis : measure of peakness of data distribution.



A: LEPTOKURTIC

$$\beta_2 > 3 \quad \gamma_2 > 0$$

B: MESOKURTIC

$$\beta_2 = 3 \quad \gamma_2 = 0$$

C: PLATYKURTIC

$$\beta_2 < 3 \quad \gamma_2 < 0$$

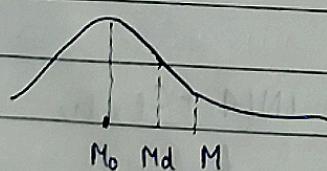
$$\beta_1 = \frac{M_3^2}{M_2^3}, \quad \gamma_1 = \sqrt{\beta_1} \quad \beta_2 = \frac{M_4}{M_2^2}, \quad \gamma_2 = \beta_2 - 3$$

Skewness : measure of symmetry $S_k = 3(M - Md) \text{ or } 3(M - Mo)$



Symmetric $\beta_1 = 0 \quad \gamma_1 = 0$

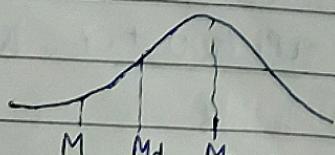
$$M = Md = Mo$$



Right skewed $\beta_1 > 0$

+ve skewed

$$M > Md > Mo$$



Left skewed

-ve skewed

$$M < Md < Mo$$

$$\beta_1 < 0$$

KARL PEARSON'S SKEWNESS COEFF = mean-mode or mean-median

BOULEY'S SKEWNESS COEFF = $\frac{Q_3 + Q_1 - 2\text{median}}{Q_3 - Q_1}$

| | |
|----------|-------|
| Page No. | _____ |
| Date | _____ |

Moments

These moments are quantitative measures related to the shape of dataset's distribution / spread and central tendency.

RAW MOMENTS (Moments about origin)

→ measures values of random variable relative to the origin. used to describe overall shape of distribution.

$$M'_r = E[x-A]^r = \frac{1}{N} \cdot \sum_{i=1}^N (x_i - A)^r \quad (A: \text{arbitrary point})$$

$$M'_r (\text{about origin}) = E[X^r]$$

$$M'_1 = (\text{first raw moment}) = \bar{x} \quad (\text{mean})$$

$$M'_2 = E(X^2) = E[(x-A)^2]$$

$$M'_3 = E(X^3) = E[(x-A)^3] \dots$$

CENTRAL MOMENTS

→ calculated with respect to the mean.
used to understand the shape and variability of dataset.

$$M_r = E[(x-\bar{x})^r] = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^r$$

$$M_1 \quad (\text{first central moment}) = 0 \quad (\text{zero})$$

$$M_2 \quad (\text{second " " }) = E[(x-\bar{x})^2] = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \sigma^2 \quad (\text{variance})$$

RELATION B/W RAW MOMENTS & CENTRAL MOMENTS

$$M_2 = M'_2 - (M'_1)^2$$

$$M_3 = M'_3 - 3M'_2 M'_1 + 2M'_1^3$$

$$M_4 = M'_4 - 4M'_3 M'_1 + 6M'_2 (M'_1)^2 - 3(M'_1)^4$$

First (now or central moments) : mean

Second (central moment) : variance

Third central moment : skewness

Fourth central moment : kurtosis

SKEWNESS

$$\beta_1 = \frac{M_3^2}{M_2^3}, \gamma_1 = \sqrt{\beta_1}$$

KURTOSIS

$$\beta_2 = \frac{M_4}{M_2^2}, \gamma_2 = \beta_2 - 3$$

Random Variable

- variable whose values depend on the outcome of a random phenomenon.

Discrete random variable: have finite distinct possible values. (PMF probability functn)

continuous random variable: can take any value within a certain range. (PDF probability functn)

Probability Mass function: it is the prob. of a discrete random variable, it must satisfy the foll. condn : 1. $P(X) \geq 0$ 2. $\sum P(X) = 1$

Probability Density function: it is the prob. functn of a contin' random variable, it must satisfy the foll. condn

$$1. f(x_i) \geq 0 \quad 2. \int_{-\infty}^{\infty} f(x).dx = 1$$

Cumulative Distribution function: it describes the prob. that a random variable takes on a value less than or equal to a given point. $F(x) = P(X \leq x)$

→ both discrete and ~~random~~ contn random vars.

$$1. 0 \leq F(x) \leq 1 \quad 2. F(x) \leq F(y) \text{ if } x \leq y$$

Expectation (expected value) of random var

measure of its avg. or mean value, based on its probability distn. It provides the weighted avg. of all possible values of random var can take

$$(\text{DISCRETE RANDOM VAR}) \quad E(X) = \sum_n x_i \cdot f(x_i)$$

$$(\text{CONTINUOUS RANDOM VAR}) \quad E(X) = \int_{-\infty}^{\infty} x \cdot f(x) \cdot dx$$

Properties

1. $E(X+Y) = E(X) + E(Y)$
2. $E(XY) = E(X) \cdot E(Y)$ X, Y are independent
3. $E[a \cdot f(x)] = a \cdot E[f(x)]$
4. $E[af(n)+b] = aE[f(n)] + b$
 $\hookrightarrow E(b) = b$

Variance

- measures the spread or dispersion of possible values around the expected value.

$$\text{Var}(X) = E(X^2) - (E(X))^2$$

$$= \sum n^2 \cdot f(n) - (\sum n \cdot f(n))^2 : \text{discrete}$$

or

$$= \int n^2 \cdot f(n) - (\int n \cdot f(n))^2 : \text{continuous}$$

Standard Deviation

- measure of the avg. devn from the mean and is expressed in same units as random var.

$$\sigma = \sqrt{\text{Var}(X)}$$

Properties

1. $V(ax) = a^2 V(X)$
2. $V(X+b) = V(X) + b$
3. $V(aX+b) = a^2 V(X) + b$

Bivariate Random Variable

involves 2 random variable X and Y , defined on the same probability space.

Measures of Association

Joint Probability distribution: joint distⁿ of 2 random var X and Y gives the probability of both variables taking specific values simultaneously.

Joint probability mass function: it is the joint probability function of discrete random var X and Y , for their simultaneous occurrence.

$$P(X=x, Y=y) = f(x, y)$$

it satisfies the condⁿ: 1. $f(x, y) \geq 0$ 2. $\sum_n \sum_y f(x, y) = 1$

Joint probability density function: it is the joint probability function of continuous random var X & Y , for their simultaneous occurrence.

$$P(X=x, Y=y) = f(x, y)$$

it satisfies the condⁿ: 1. $f(x, y) \geq 0$ 2. $\int_x \int_y f(x, y) = 1$

Marginal distribution: gives the distribution of one var X or Y alone regardless of other.

Marginal PDF: for discrete random variable

$$f_X(x) = \sum_y f(x, y) \quad \text{or} \quad f_X(x) = \sum_y f(x, y)$$

Marginal PDF: for continuous random variable

$$f_X(x) = \int_y f(x, y), dy \quad f_Y(y) = \int_x f(x, y), dx$$

| | | | | | | |
|-----------------|----------|----------|----------|----------|----------|----------|
| $x \setminus y$ | y_1 | y_2 | y_3 | \dots | y_m | |
| x_1 | $f(1,1)$ | $f(1,2)$ | \dots | $f(1,m)$ | $g_1(x)$ | |
| x_2 | $f(2,1)$ | $f(2,2)$ | \dots | \vdots | $g_2(x)$ | |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| x_m | $f(m,1)$ | $f(m,2)$ | \dots | $f(m,m)$ | $g_m(x)$ | |
| | $h_1(y)$ | $h_2(y)$ | \dots | $h_m(y)$ | 1 | |

conditional distribution.

$$\text{PMF discrete case: } P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

$$\text{PDF continuous case: } f_{x|y}(x|y) = \frac{f(x,y)}{f(y)} \quad f(y|x) = \frac{f(x,y)}{f(x)}$$

Measure of Association

COVARIANCE: a measure of how 2 variables are linearly related.

$$\text{cov}(X,Y) = E[(X - E(X))(Y - E(Y))] = E[XY] - E(X)E(Y)$$

$$\text{cov}(X,Y) = E[XY] - E(X)E(Y)$$

if X and Y are independent
 $\text{cov}(X,Y) = 0$

CORRELATION: measure of linear relationship b/w 2 variables.
 (value b/w -1 and 1)

$$\text{cor}(X,Y) = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$$

$$\text{VARIANCE: population: } \sigma^2 = \frac{1}{N} \sum_i (x_i - \mu)^2 \quad \text{sample: } s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

| | | |
|--|------------------|------------------|
| $\text{cor} = -1$ | $\text{cor} = 1$ | $\text{cor} = 0$ |
| $\text{cov} < 0$ | $\text{cov} > 0$ | $\text{cov} = 0$ |
| $\text{cor} = 0 \dots$ | $\text{cor} = 0$ | $\text{cor} = 0$ |
| $\text{cov} > 0$ | | $\text{cov} = 0$ |
| $-1 \leq \text{cor} \leq 1$ | | |
| $r = -1$: strongly negatively related | | |
| $r = 0$: no relation | | |
| $r = 1$: strongly positively related. | | |

Moment generating function

it is used to generate moments ~~of distribution~~
of the distribution.

$$M_x(t) = E[e^{tx}]$$

$$M_x(t) = E[e^{tx}] = \begin{cases} \sum_n e^{xt} \cdot f(n) & : \text{discrete} \\ \int_n e^{tx} \cdot f(x) & : \text{continuous} \end{cases}$$

$$\begin{aligned} M_x(t) &= E[e^{tx}] = E \left[1 + (xt) + \frac{(xt)^2}{2!} + \dots + \frac{(xt)^n}{n!} \right] \xrightarrow{\text{e } x \text{ series}} \\ &= E(1) + tE(x) + \frac{E(x^2)}{2!} \cdot t^2 + \dots + \frac{t^n E(x^n)}{n!} \\ &= 1 + \frac{t}{1!} M'_1 + \frac{t^2}{2!} M'_2 + \dots + \frac{t^n}{n!} M'_n \end{aligned}$$

$$M_x(t) = \sum_{r=0}^{\infty} t^r M'_r \quad \text{where } M'_r = E[x^r]$$

(Alternate method to find RAW MOMENTS)

diffn ①

$$M_n'(t) = \frac{d M_n(t)}{dt} = \left. \frac{d(0 + E(x) + 2t^2/2! + \dots)}{dt} \right|_{t=0} = \frac{E(x)}{1!}$$

$$M_1' = E(x)$$

$$M_2' = M_n''(t)$$

$$M_3' = M_n'''(t) \text{ and so on.}$$

$$M_r' = \left. \frac{d^r M_n(t)}{dt^r} \right|_{t=0}$$

Properties

1. if 2 random var have the same MGF, they have the same distribution.
2. MGF of the sum of a no. of independent random var is equal to the product of their respective MGF.

$$M_{x_1+x_2+\dots+x_n}(t) = M_{x_1}(t) \times M_{x_2}(t) \times \dots \times M_{x_n}(t)$$

$$3. M_{cx}(t) = M_x(ct)$$

4. effect of change of origin and scale.

$$U = u - a \rightarrow \text{shifting of origin}$$

\downarrow
shifting scale.

$$M_U(t) = E[e^{Ut}] = E\left[e^{(u-a)t}\right] = e^{-at} M_x(t)$$

$$M_U(t) = e^{-at} M_x(t/h)$$

DISCRETE DISTRIBUTIONS

BINOMIAL

DISTRIBUTN A random experiment whose outcome is either 'success' or 'failure' with probability p and q resp.

$$\text{PMF: } P(X=x) = {}^n C_x (p)^x (q)^{n-x}, \quad x=0, 1, \dots, n$$

- PROPERTIES:

$$\text{Mean: } E(X) = np$$

$$\text{Variance: } V(X) = npq$$

if n is large, it approx
the normal distr for
 $p \approx 0.5$

POISSON

DISTRIBUTN

Models the no. of ~~trials~~ events occurring in a fixed interval of time or space, given a constant rate λ and independence of events.

→ no. of trials is infinitely large ($n \rightarrow \infty$)

→ probability of success is very small ($p \rightarrow 0$)

$$\text{PMF: } P(X=x) = \frac{\lambda^x}{x!} e^{-\lambda} \quad x=0, 1, 2, \dots, \infty, \lambda \geq 0$$

- PROPERTIES

$$\text{Mean: } E(X) = \lambda$$

$$\text{Variance: } V(X) = \lambda$$

it is often used as an approx for binomial dist
where $n \rightarrow \infty, p \rightarrow 0$.

GEOMETRIC
DISTRIBUTN

no. of trials required to get first success in series of independent Bernoulli trials.

$$\text{PMF: } P(X=x) = (1-p)^{x-1} p = pq^{x-1} \quad x=0, 1, 2, \dots, \infty$$

- PROPERTIES

$$\text{Mean: } E(X) = q/p$$

~~Interpretation~~

$$\text{Variance: } V(X) = \frac{q}{p^2}$$

CONTINUOUS DISTRIBUTIONS

NORMAL DISTRIBUTN

A bell shaped distribution characterized by its mean μ and s.d σ .

$$\text{PDF: } f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- Mean = μ
- Variance = σ^2

→ total area under curve is 1.

UNIFORM DISTRIBUTN

All values within a specified interval $[a, b]$ are equally likely.

$$\text{PDF: } f(x) = \begin{cases} 0 & , x < a \text{ or } x > b \\ \frac{1}{b-a} & , a \leq x \leq b \end{cases}$$

- Mean = $\frac{a+b}{2}$
- Variance = $\frac{(b-a)^2}{12}$

EXPONENTIAL DISTRIBUTN

POISSON PROCESS

models time until next event occurs in a \sim

$$\text{PDF: } f(x) = \lambda e^{-\lambda x} , x \geq 0$$

- Mean = $1/\lambda$
- Variance = $1/\lambda^2$

CORRELATION & REGRESSION

LINEAR

CORRELATN

measures the strength and direction of a linear relationship b/w 2 variables X and Y .

$$\text{correlation coefficient } (r) : r = \frac{\text{cov}(X,Y)}{\sigma_x \sigma_y}$$

RANK

CORRELATION

COEFFICIENT

$$\text{spearman's rank } (P) : P = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

$$\text{correlation coefficient}$$

REGRESSION

describes the relationship b/w a dependent variable Y and one more independent variable X .

$$\text{simple linear regression} : Y = a + bx$$

REGRESSN

COEFFICIENT

$$R: \text{of } Y \text{ on } X \Rightarrow b_{yx} = r \frac{\sigma_y}{\sigma_x} \quad \} \text{ if } r=0$$

$$R: \text{of } X \text{ on } Y \Rightarrow b_{xy} = r \frac{\sigma_x}{\sigma_y} \quad \} \quad b_{yx} = b_{xy} = 0$$

• Properties

$$1. \text{ correlation coeff } r = \sqrt{b_{yx} b_{xy}}$$

$b_{yx} < 1, b_{xy} > 1$

2. if $b_{yx} < 1, b_{xy} > 1$

$$3. b_{yx} + b_{xy} > r$$

r

$$4. b_{yx} = r \frac{\sigma_y}{\sigma_x} = r \frac{\sigma_v}{\sigma_u} \quad \left(M = \frac{n-a}{n}, V = \frac{y-b}{k} \right)$$

\rightarrow independent of origin shifting, but depends on scale.

CENTRAL LIMIT THEOREM

central limit theorem states that the sum (or average) of a large no. of independent and identically distributed random variables approaches a NORMAL DISTRIBUTION, regardless of original distribution.

- APPLICABILITY:
- works for independent random variables with same distribution.
 - sample size should be large ($n \geq 30$)

STATEMENT: if $X_1, X_2, X_3, \dots, X_n$ are (independent & identically distributed) random variables with mean M & variance σ^2 , the standardized sum:

$$Z = \frac{X - M}{\sigma} \quad X \sim N(M, \sigma^2)$$

IMPLICATIONS

- allows approximation of probabilities using normal distribution.
- Justifies use of normality of large-sample hypothesis testing.

HYPOTHESIS TESTING

FORMATION OF HYPOTHESIS

A hypothesis is a claim or statement about a population parameter (e.g. mean proportion)

null hypothesis (H_0): assumes no effect or no difference. It is the statement being tested

ex: $H_0: \mu = \mu_0$ (population mean equal to some value)
alternate hypothesis (H_1): Represents what we want to prove or detect.

ex: $H_1: \mu \neq \mu_0$ (mean is not equal to specific value)

LARGE SAMPLE TESTS

For large samples ($n \geq 30$), the z-test is used.
It assumes the sampling distribution of the test statistics follows a normal distribution.

a) test for single proportion: $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$

b) test for difference of proportions: $Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_1 + n_2}}}$

c) test for single mean: $Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$

d) test for difference of means: $Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$

e) test for standard deviations: f-test

SMALL
SAMPLE
TESTS

For small samples ($n < 30$), the t-test is used, which accounts for additional uncertainty due to small sample size.

a) t-test for single mean: $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$

①

b) t-test for difference of means (equal variances):

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$$

c) t-test for correlation coefficient

$$t = r \sqrt{\frac{n-2}{1-r^2}}$$

F-TEST FOR
RATIO OF
VARIANCES

F-test is used to compare variances of 2 samples

$$F = \frac{s_1^2}{s_2^2} \quad (F \geq 1)$$

CHI-SQUARE
TEST

Chi-square test (χ^2) is used for categorical data

a) Goodness of Fit: tests whether observed frequencies (O_i) match expected frequencies (E_i)

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

b) Test for Independence: tests whether 2 attributes are independent.

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

TESTING FOR STATISTICAL HYPOTHESIS

Step 1: Null hypothesis (H_0)

~~Step~~

Step 2: Alternate hypothesis (H_1)

$$\left\{ \begin{array}{ll} H_1 : M \neq M_0 & (2\text{-tailed}) \\ H_1 : M > M_0 & (\text{right-tailed}) \\ H_1 : M < M_0 & (\text{left-tailed}) \end{array} \right.$$

Step 3: Level of significance

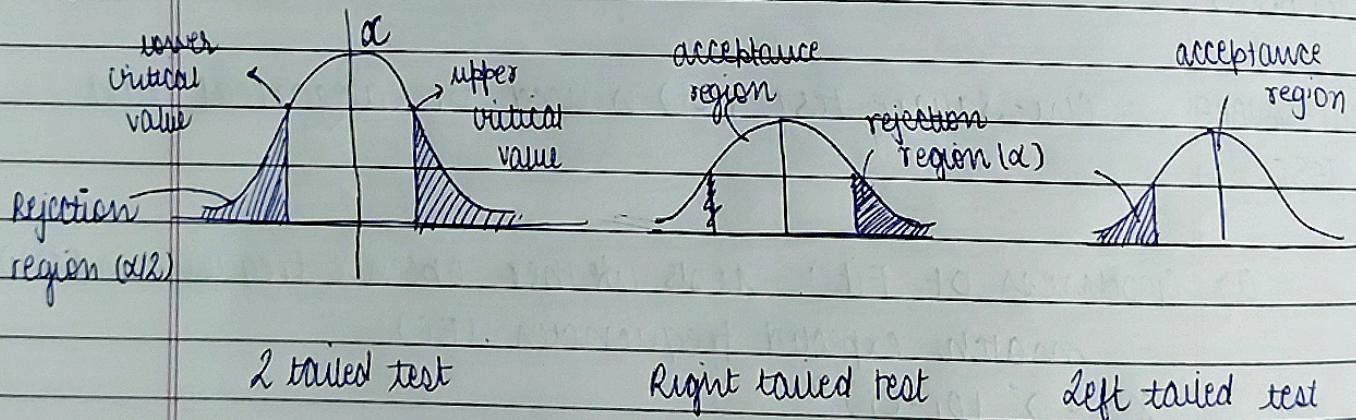
→ the probability of the variate falling in the critical region.

Step 4: Test statistic

compute test statistic Z under null hypothesis.

$$Z = \frac{T - E(T)}{S_E(T)}$$

Step 5: conclusion



$$|Z| > z_\alpha : \text{Reject } H_0 \quad |Z| < z_\alpha : \text{accept } H_0$$

| | 1% (0.01) | 5% (0.05) | 10% (0.1) |
|--------------|-------------------------|---------------------|----------------------|
| Two tailed | $ z_{\alpha/2} = 2.58$ | $ z_\alpha = 1.96$ | $ z_\alpha = 1.645$ |
| Right tailed | $z_\alpha = 1.33$ | $z_\alpha = 1.615$ | $z_\alpha = 1.28$ |
| Left tailed | $z_\alpha = -1.33$ | $z_\alpha = -1.615$ | $z_\alpha = -1.28$ |

Z-test

TEST OF SIGNIFICANCE FOR LARGE SAMPLES

$n > 30$

single proportion

$$S.E.(p) = \sqrt{\frac{pq}{n}}$$

$$E(p) = E\left(\frac{x}{n}\right) = \frac{np}{n} = p$$

$\hat{p} = \frac{x}{n}$: observed proportion of success

p : probability of success

q : probability of failure

difference of proportion

$$Z = \frac{p_1 - p_2}{\sqrt{pq \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \quad q = 1 - p$$

single mean

$$Z = \frac{\bar{X} - M}{\sigma / \sqrt{n}}$$

\bar{X} : mean sample M : population mean

σ : S.D n : sample size

$$\text{limits of population mean } \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < M < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

diff'n of means for 2 large samples

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

σ : S.D \bar{X} : sample mean

n : sample size

$$\sigma_1^2 = \sigma^2 ; \sigma^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2}$$

diff'n of standard deviation

$$Z = \frac{s_1 - s_2}{\sqrt{\frac{\sigma_1^2}{2n_1} + \frac{\sigma_2^2}{2n_2}}}$$

σ : population standard devn

s : sample standard devn

n : sample size

t-test

TEST OF SIGNIFICANCE FOR SMALL SAMPLE

$n < 30$

Mean
of
random
sample

$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

\bar{X} : sample mean

S : sample S.D.

μ : population mean

n : sample size

Mean of

$$t = \frac{(\bar{X} - \bar{Y})}{S\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

\bar{X}, \bar{Y} : sample mean

n : sample size

small
sample

from

Normal

population

degree of freedom = $n_1 + n_2 - 2$

$$S^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2} \quad \text{or} \quad S^2 = \frac{\sum (x_1 - \bar{x}_1)^2 + \sum (x_2 - \bar{x}_2)^2}{n_1 + n_2 - 2}$$