

1

N

D

E

NAME: Arpan Sapkota

STD.: MDS

SEC.: 1 | L

ROLL NO.: 07

SUB.: for

Statistical Methods for Data Science (MDS553)

3 Generalized Power Series Distribution

Variable

1. Qualitative : Categorical, Response, Attributes.

2. Quantitative : Numerical.

→ 1. Discrete → Counting process.

→ 2. Continuous → Measuring process. → Interval.

Random Variable → Unbiased, fair, $X = \text{no. of passed student.}$

The variable which takes different values generated by random experiment is known as random variable. Random Variable is denoted by x, y and $z.$

For Example:

If we toss a coin two times, then total no. of possible cases = $2^2 = 4.$ (sample space)
 $(\text{events})^{(\text{no. of trials})} \Rightarrow$ Mutually Exclusive, collectively exhaustive cases.

$$\{(HH, HT, TH, TT)\} \quad X \quad P(X)$$

Let $x = \text{no. of heads}$	0	$\frac{1}{4}$
	1	$\frac{2}{4}$

If implies x is a random variable.

$$\sum P(x) = 1.$$

$f(x) \rightarrow$ Probability density function (P.d.f) $\int_{-\infty}^{\infty} f(x) dx = 1, 0 \leq f(x) \leq 1, C$

$P(x) \rightarrow$ Probability Mass function (P.m.f) $\sum P(x) = 1, 0 \leq P(x) \leq 1, D$

Mathematical Expectation (Avg)

let x be a discrete random variable which takes the values x_1, x_2, \dots, x_n with respective probabilities $P(x_1), P(x_2), \dots, P(x_n)$ then the mathematical expectation of random variable x is given by,

$$E(x) = x_1 P(x_1) + \dots + x_n P(x_n) = \sum x_i P(x_i) \quad (\text{Discrete})$$

where,

$P(x) =$ Probability Mass function.

By, $E(x) = \int x f(x) dx$ where $f(x) =$ probability Density function
 (continuous)

1 2 3 4 5 6 }

$$E(X) = \bar{X} = \sum x P(x)$$

x = values on the faces of the dice.

$$= \int x f(x) dx$$

X	P(x)	x P(x)
1	1/6	1/6
2	1/6	2/6
3	1/6	3/6
4	1/6	4/6
5	1/6	5/6
6	1/6	6/6

$$E(X) = \sum x P(x) = \frac{1}{6} [1+2+3+4+5+6]$$

6

6

6

6

6

Variance in terms of Mathematical Expectation

$$V(X) = \text{Variance of } X = E[(X - E(X))^2]$$

$$\therefore V(X) = E(X^2) - [E(X)]^2 \quad \text{where, } E(X^2) = \sum x^2 P(x)$$

Hence, for continuous,

$$E(X^2) = \int_{-\infty}^{\infty} x^2 f(x) dx$$

Moments (Shape of distribution? \Rightarrow Moments)

The variance, power of deviations taken from chosen value is known as moments. The chosen value may be any arbitrary value (assumed mean).

If deviations are taken from mean then they are called central moments.

If any arbitrary value then called raw moments.

Central Moments (μ_r).The r th order central moment is denoted by μ_r and is given by,

$$\text{Definition: } \mu_r = \frac{1}{n} \sum (x - \bar{x})^r \quad [\text{Individual Series}]$$

First Order Central Moment (μ_1) = $\frac{1}{n} \sum (x - \bar{x})^1 = \bar{x}$ Second Order Central Moment (μ_2) = $\frac{1}{n} \sum (x - \bar{x})^2 = \text{Variance}$

(contd.)

Raw Moments (μ_r')

The r th order raw moment about any arbitrary value (assumed mean)

'a' is given by,

$$\mu_r' = \frac{1}{n} \sum (x - a)^r$$

r th raw moment about origin i.e. $a=0$

$$\mu_r' = \frac{1}{n} \sum (x - 0)^r = \frac{1}{n} \sum x^r$$

$$\Rightarrow \mu_r' = E(x^r)$$

$$r=1, \mu_1' = E(x) = \text{Mean}$$

$$r=2, \mu_2' = E(x^2)$$

* Note

$$\mu_2 = \text{Variance}$$

$$= E(x^2) - [E(x)]^2$$

$$= \mu_2' - (\mu_1')^2$$

$$* \text{ Coeff. of Skewness } (\beta_1) = \frac{\mu_3^2}{\mu_2^2} \Rightarrow r_1 = \pm \sqrt{\beta_1}$$

$$* \text{ Coeff. of Kurtosis } (\beta_2) = \frac{\mu_4}{\mu_2^2}$$

Moment Generating function.

The moment generating function of random variable x is given by,

$$M_x(t) = E(e^{tx}) = \sum_x e^{tx} p(x) \quad (\text{Discrete})$$

$$= \int_x e^{tx} f(x) dx \quad (\text{Continuous})$$

* Note

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

$$= E \left[1 + \frac{tx}{1!} + \frac{t^2 x^2}{2!} + \dots + \frac{t^r x^r}{r!} + \dots \right]$$

$$e^a = 1 + a + \frac{a^2}{2!} + \dots$$

$$= 1 + t E(x) + \frac{t^2}{2!} E(x^2) + \dots + \frac{t^r}{r!} E(x^r) + \dots$$

$$M_x(t) = 1 + \frac{t}{1!} \mu_1' + \frac{t^2}{2!} \mu_2' + \dots + \frac{t^r}{r!} \mu_r' + \dots$$

* Note:

$$\mu_1' = \left. \frac{d}{dt} M_x(t) \right|_{t=0} \quad \text{and, } \mu_2' = \left. \frac{d^2}{dt^2} M_x(t) \right|_{t=0}$$

$$\mu_1' = \left. \frac{d}{dt} M_x(t) \right|_{t=0}$$

Properties

1. $M_{ax}(t) = M_x(at)$

2. If X and Y are two independent variables then the moment generating function of their sum ($X+Y$) is

$$M_{x+y}(t) = M_x(t) M_y(t) \quad | E(XY) = E(X) E(Y)$$

If X and Y are independent

Characteristic function

$$\Phi_x(t) = E[e^{itx}]$$

$$= \sum_x e^{itx} p(x) \quad (\text{Discrete})$$

$$= \int_x e^{itx} f(x) dx \quad (\text{Continuous})$$

Covariance between X and Y + Note: If X and Y are independent

Covariance between X and Y . $E(XY) = E(X) E(Y)$

$$\text{Cov}(X, Y) = \frac{1}{n} \sum (x - \bar{x})(y - \bar{y}) \Rightarrow \text{Cov}(X, Y) = 0$$

$$= \frac{1}{n} \sum xy - \bar{x}\bar{y}$$

$$\text{Cov}(X, Y) = E(XY) - E(X) E(Y)$$

Discrete Probability Distribution

1. Binomial

2. Poisson

3. Negative Binomial

1. Binomial Distribution.

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

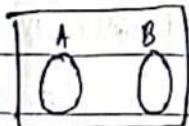
n = no. of trials (fixed)

p = prob. of success

$q = 1-p$ = prob. of failure

X = no. of successes

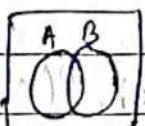
ME Events \Rightarrow Disjoint Set



\rightarrow coin Toss

\rightarrow Card: R & B

NA ME Events \Rightarrow Joint Set



\rightarrow Simultaneous occur

\rightarrow cards: R & Queen

Generalized Power Series Distribution (GPSD)

A discrete random variable X is said to follow a general power series distribution (GPSD) if its probability mass function is given by,

$$P(X=x) = P(x) = \sum_{\theta} a_x \theta^x ; \quad x = 0, 1, 2, \dots ; \quad a_x > 0$$

$f(\theta)$

$0 ; \text{ otherwise}$

$$\text{where, } f(\theta) = \sum_{x \in S} a_x \theta^x ; \quad \theta > 0 \quad (\text{power series})$$

So that $f(\theta)$ is positive, finite and differentiable and it is a non-empty countable subset of non-negative integers.

Special Case of GPSD

1. Binomial Distribution.

In GPSD, take $\theta = p$, $f(\theta) = (1+\theta)^n$ and $x = 0, 1, 2, \dots, n$

Then,

$$f(\theta) = (1+\theta)^n = \sum_{x=0}^n a_x \theta^x$$

$$\text{Now, the binomial expansion of } (1+\theta)^n \text{ is } (1+x)[a+b]^n = \sum_{x=0}^n {}^n C_x a^x b^{n-x}$$

$$(1+\theta)^n = \sum_{x=0}^n {}^n C_x \theta^x (1)^{n-x}$$

$$\therefore P(X=x) = \sum_{x=0}^n {}^n C_x \theta^x$$

$$= \sum_{x=0}^n {}^n C_x \theta^x$$

$$\therefore a_x = {}^n C_x$$

$$\text{Since, } P(X=x) = a_x \theta^x = {}^n C_x \left(\frac{p}{1-p}\right)^x f(\theta) \quad (1+p)^n$$

$$= {}^n C_x p^x (1-p)^{n-x} = {}^n C_x p^x (1-p)^{n-x}$$

$$\left(\frac{1-p+p}{1-p}\right)^n$$

$$\therefore P(X=x) = {}^n C_x p^x q^{n-x} ; x=0,1,2,\dots,n$$

This is the probability mass function of Binomial distribution with parameters 'n' and 'p'.

Hence, the binomial distribution is the special case of GPSD.

Under the following condition Binomial distribution is applied.

1. No. of trials 'n' is fixed.
2. Events are mutually exclusive i.e. only two possible outcomes like head or tail, defective or non-defective, pass or fail, etc
3. Events are statistically independent.
4. Probability of success 'p' remains constant in each trials.
5. $p+q=1$; q = probability of failure.

Definition:

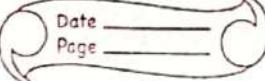
If a discrete random variable X follow binomial distribution with parameter 'n' and 'p' then the probability mass function of the distribution is given by

$$P(X=x) = {}^n C_x p^x q^{n-x} ; x=0,1,2,\dots,n$$

$$p+q=1$$

at least (\geq) : at least 2 $\Rightarrow x \geq 2$

at most (\leq) : at most 2 $\Rightarrow x \leq 2$



It gives the probability of getting 'x' success out of 'n' independent trials.
where,

p = prob. of success.

x = no. of successes

Note:

If $X \sim B(n, p)$ \Rightarrow i.e. X follows Binomial distribution with parameters n & p.
Then,

$$\text{Mean} = np \quad \text{and}$$

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

Problems on Binomial Distribution.

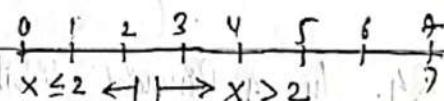
Q2. for a binomial distribution with $n=7$ and $p=0.2$, find.

- (a) $P(X=5)$ (b) $P(X>2)$ (c) $P(X>3)$ (d) $P(X \geq 4)$

$$n=7, \quad p=0.2$$

$$q=1-p=0.8$$

If $X \sim B(n, p)$, then $P(X=x) = {}^n C_x p^x q^{n-x}$; $x=0, 1, 2, \dots, n$,



$$(b) P(X>2) = 1 - P(X \leq 2)$$

$$= 1 - [P(X=0) + P(X=1) + P(X=2)]$$

$$= 1 - [{}^7 C_0 (0.2)^0 (0.8)^{7-0} + {}^7 C_1 (0.2)^1 (0.8)^{7-1} + {}^7 C_2 (0.2)^2 (0.8)^{7-2}]$$

$$= 1 - [0.2097 + 0.367 + 0.271]$$

$$= 0.1480$$

$$(d) P(X \geq 4) = P(X=4) + P(X=5) + P(X=6) + P(X=7)$$

$$= {}^7 C_4 (0.2)^4 (0.8)^{7-4} + {}^7 C_5 (0.2)^5 (0.8)^{7-5} + {}^7 C_6 (0.2)^6 (0.8)^{7-6}$$

$$+ {}^7 C_7 (0.2)^7 (0.8)^{7-7}$$

$$= 0.0333$$

Q5. At a particular university it has been found that 20% of the student withdraw without completing the Business statistics course. Assume that 18 student have registered for the course this semester.

- (a) What is the probability that none will withdraw?
- (b) What is the probability that at least one will withdraw?
- (c) What is the probability that at most 2 will withdraw?

$X = \text{no. of withdraw}$

$$p = \text{prob. of withdraw} = 0.20 \Rightarrow q = 0.80$$

$n = 18$

Here, $X \sim B(n, p)$ so, $P(X=x) = {}^n C_x p^x q^{n-x}$; $x = 0, 1, 2, \dots, n$

$$(a) P(X=0) = {}^{18} C_0 (0.20)^0 (0.80)^{18-0}$$

$$(b) P(X \geq 1) = 1 - P(X < 1) = 1 - P(X=0) = 1 - 0.018 = 0.982$$

$$(c) P(X \leq 2) = P(X=0) + P(X=1) + P(X=2)$$

$$= 0.018 + {}^{18} C_1 (0.2)^1 (0.80)^{18-1} + {}^{18} C_2 (0.2)^2 (0.80)^{18-2}$$

Q4. In a Binomial distribution with 6 independent trials the probabilities of 3 and 4 successes are found to be 0.2457 and 0.089. find the parameter p of the distribution. ($p = 0.326$)

$n = 6$

1 2 3 4 (5) 6

$$p = 0.326$$

$$q = 1-p = 1-0.326$$

$$P(X=3) = 0.2457, P(X=x) = {}^n C_x p^x q^{n-x}$$

$$P(X=4) = 0.089$$

Now,

$$\frac{P(X=3)}{P(X=4)} = \frac{0.2457}{0.089} = \frac{{}^6 C_3 p^3 q^{6-3}}{2.76} = \frac{2.76}{(0.326)^3 (0.674)^3}$$

$$\text{or, } \frac{209}{152} = 2.76 \Rightarrow 1-p = \frac{2.76 \times 3}{4}$$

$$\text{or, } \frac{1-1}{p} = \frac{2.07}{1} \Rightarrow \frac{1-1}{p} = 2.07 \Rightarrow q = \frac{1}{3.07}$$

Q13. A fair coin is tossed ten times. Find the probability of obtaining

(a) Exactly 4 heads.

(b) No heads

$$n = 10$$

(c) At least one head

$$p = 0.5 \quad (\because \text{since coin is fair})$$

(d) At most three heads

$$x = \text{no. of heads}$$

(e) More than 8 heads

$$q = 1-p = 0.5$$

(f) 3 heads and 7 tails.

$$\text{NoH} : 10, 9, 8, \dots, 0$$

Here,

$$\text{NoT} : 0, 1, 2, \dots, 10$$

$$X \sim B(n, p) \text{ so, } P(X=x) = {}^n C_x p^x q^{n-x}, x=0, 1, 2, \dots, n$$

$$(a) P(X=4) = {}^{10} C_4 (0.5)^4 (0.5)^{10-4} = 0.2051$$

$$(b) P(X=0) = {}^{10} C_0 (0.5)^0 (0.5)^{10-0} = 0.000977$$

$$(c) P(X \geq 1) = 1 - P(X < 1) = 1 - P(X=0) = 1 - 0.000977$$

$$(d) P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3)$$

$$(e) P(X \geq 8) = P(X=9) + P(X=10)$$

$$(f) \text{Prob. [8 heads and 7 tails]} = P(X=3)$$

Q15. Prob. of being liberal = 0.30

prob. of being conservative = 0.55

prob. of middle of the road = 0.15

$$\text{If } X \sim B(n, p) \text{ then } P(X=x) = {}^n C_x p^x q^{n-x}; x=0, 1, 2, \dots, n$$

$$n = 10$$

(a) $X = \text{no. of liberal}$

(b) $X = \text{no. of middle of the road.}$

$$p = 0.30 \Rightarrow q = 0.70$$

$$p = 0.15$$

$$P(X=4) = {}^{10} C_4 (0.30)^4 (0.70)^{10-4}$$

$$P(X=2)$$

(b) $X = \text{no. of conservative}$

$$p = 0.55$$

$$P(X=0) =$$

(c) $X = \text{no. of liberal if } p = 0.30$

$$P(X \geq 8) = P(X=8) + P(X=9) + P(X=10)$$

Q17

$$n=6$$

$$P = \text{prob. of illegal} = 0.45$$

$$q = 1-p = 0.55$$

$$P(X=x) = {}^n C_x p^x q^{n-x}; x=0,1,\dots,n.$$

$$p+q=1.$$

$X = \text{no. of illegal}$

$$\text{(i). } P(X=2) =$$

$$\text{(ii). } P(X=0) =$$

$$\text{(iii). } P = \text{prob. of legal} = 0.55$$

$$q = 1-p = 0.45$$

$X = \text{no. of legal}$

Q18.

mean = 0.4 and s.d. = 0.6

If $X \sim B(n,p)$ then mean = np (and) variance ($= npq$)

$$\text{and } P(X=x) = {}^n C_x p^x q^{n-x}; (x=0,1,2,\dots,n)$$

$$\therefore np = 0.4 \text{ and } npq = 0.6^2$$

$$0.4(q) = 0.6^2 \Rightarrow q = 0.6^2 / 0.4 = 0.9$$

$$\Rightarrow q = 0.6^2 / 0.4 = 0.9$$

$$\therefore p = 1-q = 0.1$$

$$\text{further, } np = 0.4 \Rightarrow n = 0.4 / 0.1 = 4$$

Prob. at least one success

$$P(X \geq 1) = 1 - P(X < 1) = 1 - P(X=0) = 1 - {}^4 C_0 (0.1)^0 (0.9)^4 = 0.3439$$

Q18.

$$P(X \geq 2) = ?$$

Q24.

$$n=5$$

$$P = \text{prob. of correct answer} = 1/4 = 0.25$$

$$(q = 1-p = 0.75)$$

$X = \text{no. of correct answers}$

Here, $X \sim B(n,p)$, so $P(X=x) = {}^n C_x p^x q^{n-x}; x=0,1,2,\dots,n$

$$\text{(i). } P(X=5) = {}^5 C_5 (0.25)^5 (0.75)^{5-5} = 0.25^5 = 0.00098$$

$$\text{(ii). } P(X \geq 4) = P(X=4) + P(X=5) = {}^4 C_4 (0.25)^4 (0.75)^{4-4} + 0.00098$$

$$= 0.0156.$$

(Events) no. of trials.

Date _____
Page _____

Q25

7 → (6, 1), (1, 6), (5, 2), (2, 5), (4, 3), (3, 4)

11 → (6, 5), (5, 6)

12 → (6, 6) $n = 5$.

p = prob. that a player gets audited

$$= \frac{9}{36} = \frac{1}{4}$$

Favorable no. of cases for sum

$$7 \text{ or } 11 \text{ or } 12 = 9$$

$$q = 1 - p = 0.75$$

Total cases = 86

$$n = 5$$

$\lambda = \text{no. of times a player gets audited.}$

Here, X follows $B(n, p)$ so, $P(X=x) = {}^n C_x p^x q^{n-x}$, $x = 0, 1, 2, \dots, n$.

$$P(X \geq 1) = 1 - P(X=0)$$

$$= 1 - [{}^5 C_0 (0.25)^0 (0.75)^{5-0}] = 0.9627.$$

Mean, Variance and Moment generating function of Binomial Distribution with parameters 'n' and 'p'.

Mean

If $X \sim B(n, p)$ then $P(X=x) = {}^n C_x p^x q^{n-x}$; $x = 0, 1, \dots, n$.

$$p+q=1$$

$$E(X) = \sum_{x=0}^n x P(x)$$

$$= \sum_{x=0}^n x \{ {}^n C_x p^x q^{n-x} \}$$

$${}^n C_r = n!$$

$$= \sum_{x=0}^n x \cdot \frac{n!}{(n-x)! x!} p^x q^{n-x} / (n-x)! x!$$

$$= \sum_{x=0}^n x \cdot \frac{n(n-1)!}{(n-x)! x!(x-1)!} p^{x-1} p q^{n-x} / [(a+b) = \sum_{x=0}^n {}^n C_x p^x q^{n-x}]$$

$$= np \sum_{x=1}^n \frac{(n-1)!}{(n-x)! (x-1)!} p^{x-1} q^{n-x}$$

$$= np \sum_{x=1}^n {}^{n-1} C_{x-1} p^{x-1} q^{(n-x)-(x-1)}$$

$$E(X) = np [q+p]^{n-1}$$

$$\therefore E(X) = np \quad [p+q=1]$$

Variance

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

$$= E(X^2) - (np)^2 \quad \text{--- (1)}$$

$$\text{Now, } E(X^2) = \sum_{x=0}^n x^2 p(x)$$

$$= \sum_{x=0}^n [x(x-1)(+x)] p(x)$$

$$= \sum_{x=0}^n x(x-1)p(x) + (\sum_{x=0}^n x p(x))$$

$$= \sum_{x=0}^n x(x-1) \frac{n!}{(n-x)!} p^x q^{n-x} + E(X)$$

$$= \sum_{x=0}^n x(x-1) \frac{n!}{(n-x)!} \frac{p^x q^{n-x}}{(n-x)! x!} + np$$

$$= \sum_{x=0}^n x(x-1) \frac{n(n-1)(n-2)!}{(n-x)!, 2(x-1)(x-2)!} p^2 p^{x-2} q^{n-x} + np$$

$$= n(n-1)p^2 \sum_{x=2}^{n-2} {}_{x-2}^n C p^{x-2} q^{n-x} + np$$

$$= n(n-1)p^2 [p^{n-2} + \dots + p^{n-2}] + np$$

$$= n(n-1)p^2 [q+p]^{n-2} + np$$

$$= n(n-1)p^2 + np \quad [p+q=1]$$

Now, eqn (1). becomes.

$$\text{Var}(X) = n(n-1)p^2 + np - n^2 q^2 = n^2 p^2 - np^2 + np - n^2 p^2 \\ = np(1-p)$$

$$\therefore \text{Var}(X) = npq$$

Moment Generating function (mgf)

The mgf of binomial r.v. X is given by,

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

$$= \sum_{x=0}^n e^{tx} P(x)$$

$$= \sum_{x=0}^n e^{tx} {}^n C_x p^x q^{n-x}$$

$$= \sum_{x=0}^n {}^n C_x (pe^{tx}) q^{n-x}$$

$$\boxed{M_X(t) = [q + pe^t]^n} \quad \because [q+p]^n = \sum_{x=0}^n {}^n C_x p^x q^{n-x}$$

Mean and Variance using mgf.

The r^{th} raw moment about origin is given by

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

$$\mu_1' = E(X-0) = E(X)$$

$$\mu_2' = E(X-0)^2 = E(X^2)$$

We know,

$$M_X(t) = (q + pe^t)^n$$

$$\frac{d M_X(t)}{dt} = n(q + pe^t)^{n-1} pe^t$$

$$\text{Now, } \mu_1' = E(X) = \left. \frac{d M_X(t)}{dt} \right|_{t=0} = n(q + pe^t)^{n-1} pe^0$$

$$= n(q + p)^{n-1} p$$

$$= np \quad [\because q+p=1]$$

Further,

$$\frac{d^2 M_X(t)}{dt^2} = np \left[e^{t(n-1)} (q + pe^t)^{n-2} pe^t + (q + pe^t)^{n-1} e^t \right]$$

$$= np \left[e^{t(n-1)} \frac{d}{dt} (q + pe^t)^{n-1} + (q + pe^t)^{n-1} \frac{d}{dt} e^t \right]$$

$$\therefore \frac{d}{dt} (q + pe^t)^{n-1} = (q + pe^t)^{n-1} \cdot \frac{d}{dt} (q + pe^t)$$

$$= (n-1)(q + pe^t)^{n-1-1} \cdot (0+pe^t) = (n-1)(q + pe^t)^{n-2} pe^t$$

$$= np [e^{(n-1)}(q+pe^t)^{n-2} pe^t + (q+pe^t)^{n-1} e^t]$$

$$\therefore \mu_2' = \frac{d^2 M_x(t)}{dt^2} \Big|_{t=0} = np [e^0(n-1)(q+pe^0)^{n-2} pe^0 + (q+pe^0)^{n-1} e^0]$$

$$= np [(n-1)p+1]$$

$$= n^2 p^2 - np^2 + np$$

$$\therefore \text{Variance} = E(x^2) - [E(x)]^2$$

$$> n^2 p^2 - np^2 + np - n^2 p^2$$

$$= np(1-p)$$

$$\therefore \mu_2 = npq.$$

0 — x — n Binomial

(v) 0 — x — n Poisson

(vi) 0 — x — ∞ Normal

Ex: Ques. I

Ans. 1. If $X \sim N(\mu, \sigma^2)$ then $E(X^2) = \mu^2 + \sigma^2$

(Ans. 2)

Ans. 3. $E(X^2) = (EX)^2 + E(X^2 - EX)^2$

Poisson Distribution

If a discrete random variable x follows poisson distribution with λ , then the probability mass function of the distribution is given by,

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} ; x = 0, 1, 2, \dots$$

Note:

If $X \sim P(\lambda)$ then mean = λ and variance = λ .
where λ = average.

Example:

X = No. of calls in a day

λ = Average no. of calls in a day.

Poisson Approximation to Binomial distribution or relation between binomial and poisson distribution.

The binomial distribution reduces to poisson distribution under following conditions:

i) If no. of trials 'n' is very large i.e. $n \rightarrow \infty$ ($n > 20$)

ii) If probability of success p is very small i.e. $p \rightarrow 0$ ($p \leq 0.05$)
 $\Rightarrow \lambda = np$

$$P(X=x) = \frac{e^{-np} (np)^x}{x!}$$

Poisson distribution from GPSD.

$$P(x) = \frac{a_x \theta^x}{f(\theta)}$$

In GPSD, take $\theta = \lambda$, $f(\theta) = e^\theta$ and

$S = \{0, 1, 2, \dots, \infty\}$ then, where, $f(\theta) = \sum_{x=0}^{\infty} a_x \theta^x$

$$f(\theta) = e^\theta = \sum_{x=0}^{\infty} a_x \theta^x$$

$$e^x = 1 + x + \frac{x^2}{2!} + \dots$$

$$\text{Now, } e^\theta = 1 + \theta + \frac{\theta^2}{2!} + \dots = \sum_{x=0}^{\infty} \frac{1}{x!} \theta^x$$

Now

$$\sum_{x=0}^{\infty} \frac{\theta^x}{x!} = \sum_{x=0}^{\infty} a_x \theta^x$$

$$\Rightarrow a_x = \frac{1}{x!}$$

The proof of GPSD is given by,

$$P(x) = \frac{a_x \theta^x}{f(\theta)} = \frac{\left(\frac{1}{x!}\right) \theta^x}{e^\theta}$$

$$\therefore P(x) = \frac{e^{-\lambda} \lambda^x}{x!}; x = 0, 1, 2, \dots$$

This is the proof of poisson distribution with parameter λ .
Hence, poisson distribution is the special case of GPSD.

Problems:

Q. 1. $\lambda = 4.2$

(a). $P(X \geq 5) = 1 - P(X < 5)$

$$(X \leq 4) = 1 - [P(X=0) + P(X=1) + P(X=2) + P(X=3) + P(X=4)]$$

If $X \sim P(\lambda)$ then, $P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}; x = 0, 1, 2, \dots$

$$\begin{aligned} P(X \geq 5) &= 1 - \left[\frac{e^{-4.2} 4.2^0}{0!} + \frac{e^{-4.2} 4.2^1}{1!} + \frac{e^{-4.2} 4.2^2}{2!} + \frac{e^{-4.2} 4.2^3}{3!} + \frac{e^{-4.2} 4.2^4}{4!} \right] \\ &= 1 - e^{-4.2} [1 + 4.2 + 8.82 + 12.348 + 12.965] \\ &= 0.41 \end{aligned}$$

Q. 3. X = no. of price takers in every 3 years.

λ = average no. of price takers in every 3 years.
 $= 4$.

Here, $X \sim P(\lambda)$ so, $P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}; x = 0, 1, 2, \dots$

Q. 7 no. price hikes in a randomly selected period of 3 years is 12

$$P(X=0) = \frac{e^{-4} 4^0}{0!} = 0.0183$$

$$P(X=1) = \frac{e^{-4} 4^1}{1!} = 0.1464$$

$$P(X=2) = \frac{e^{-4} 4^2}{2!} = 0.1952$$

$$\begin{aligned} P(X \geq 5) &= 1 - [P(X < 5)] \\ &= 1 - [P(X=0) + P(X=1) + P(X=2) + P(X=3) + P(X=4)] \\ &= 1 - [0.0183 + \frac{e^{-4} 4^1}{1!} + 0.1464 + \frac{e^{-4} 4^2}{2!} + 0.1952] \end{aligned}$$

Q4. Given, Average no. of sales in week-days = 1.6

(a) X = no. of sales in two days

λ = average no. of sales in two days.

$$= \frac{1.6}{5} \times 2$$

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}; x=0, 1, 2, \dots$$

$$= 0.64$$

$$P(X=4) = \frac{e^{-0.64} 0.64^4}{4!} = 0.0037$$

(b) X = no. of sales in a day

λ = average no. of sales in a day.

$$= 1.6$$

$$= \frac{1.6}{5}$$

$$= 0.32$$

$$P(X \geq 3) = 1 - P(X < 3)$$

$$= 1 - [P(X=0) + P(X=1) + P(X=2)]$$

Q8. $X = \text{no. of arrivals in 1 minute}$
 $\lambda = \text{average no. of arrival in a minute} = 6.$

$$\text{(i)} \quad P(X > 3) = 1 - P(X \leq 3) \\ = 1 - [P(X=0) + P(X=1) + P(X=2) + P(X=3)]$$

$$\text{(ii)} \quad P(X = 3 \text{ or } 4) = P(X=3) + P(X=4) \quad P(A \text{ or } B) = P(A) + P(B) - P(A \cap B)$$

$$\text{(iii)} \quad P(X=1)$$

$$\text{(iv).} \quad P(X > 2) = 1 - P(X \leq 2) \\ = 1 - [P(X=0) + P(X=1)]$$

(c) (i) $X = \text{no. of arrivals in two minutes.}$

$$\lambda = \text{average no. of arrivals in two minutes.} \\ = 6 \times 2 = 12$$

$$\text{(ii).} \quad P(X=1) =$$

$$\text{(iii).} \quad P(X > 2) = 1 - P(X \leq 2)$$

$$= 1 - [P(X=0) + P(X=1)].$$

10. $A = \text{average no. of vehicles per cycle (45 seconds)}$

$$= \frac{10}{60} \times 45 \\ = 7.5$$

$X = \text{no. of vehicles per cycle.}$

$$45 - \frac{10}{60} \times 45 \\ = 60$$

20. $n = \text{no. of disputes} = 80$

$$p = \text{prob. of non-manual disputes.} \\ = 1 - 0.96$$

$$p \leq 0.05 \\ n > 20$$

$$p = 0.04$$

Now, using Poisson approximation to binomial we get,

$$\eta = np = 60 (0.04) = 3.2$$

$$\therefore P(X=x) = \frac{e^{-np} (np)^x}{x!}$$

$$\textcircled{a} \quad P(X=7) = \frac{e^{-3.2} 3.2^7}{7!} = 0.027$$

$$\textcircled{b} \quad P(X=0) = \frac{e^{-3.2} 3.2^0}{0!} = e^{-3.2} = 0.0407$$

$$\textcircled{c} \quad P(X \geq 2) = 1 - P(X < 2)$$

$$= 1 - [P(X=0) + P(X=1)]$$

$$= 1 - [0.0407 + \frac{e^{-3.2} 3.2^1}{1!}]$$

$$25. \quad p = \text{prob. of defective} = \frac{1}{400} = 0.0025$$

$$\eta = 100$$

X = no. of defectives.

Using poisson approximation to binomial.

We get,

$$\lambda = np = 100 (0.0025) = 0.25$$

$$\textcircled{a} \quad P(X=0)$$

$$\textcircled{b} \quad P(X \leq 2)$$

$$\textcircled{c} \quad P(X \geq 1) = 1 - P(X=0)$$

$$\textcircled{d} \quad P(X \geq 3)$$

Mean, Variance and Moment Generating Function of Poisson Distribution with parameter λ .

If $X \sim P(\lambda)$ then the probability of the distribution is,

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad x=0,1,2,\dots \quad [e^{\lambda} = 1 + \lambda + \frac{\lambda^2}{2!} + \dots]$$

Mean:

$$E(X) = \sum x P(x) = \sum_{x=0}^{\infty} x \cdot \frac{e^{-\lambda} \lambda^x}{x!}$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} \frac{x \cdot \lambda^x}{x!(x-1)!} = \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}$$

$$= \lambda e^{-\lambda} \left[1 + \frac{\lambda}{1!} + \frac{\lambda^2}{2!} + \dots \right]$$

$$= \lambda e^{-\lambda} \cdot e^{\lambda}$$

$$\therefore E(X) = \lambda$$

Variance:

$$\text{Var}(X) = E(X^2) - [E(X)]^2 = E(X^2) - (\lambda)^2 \quad \text{--- (1)}$$

$$\text{where, } E(X^2) = \sum_{x=0}^{\infty} x^2 P(x) = \sum_{x=0}^{\infty} [x(x-1) + x] P(x),$$

$$= \sum_{x=0}^{\infty} x(x-1) P(x) + \sum_{x=0}^{\infty} x P(x) = \sum_{x=0}^{\infty} x(x-1) \frac{e^{-\lambda} \lambda^x}{x!} + E(X)$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} \frac{x(x-1)\lambda^x}{x(x-1)(x-2)!} + \lambda = \lambda^2 e^{-\lambda} \left[\sum_{x=2}^{\infty} \frac{\lambda^{x-2}}{(x-2)!} \right] + \lambda$$

$$= \lambda^2 e^{-\lambda} \left[1 + \frac{\lambda}{1!} + \frac{\lambda^2}{2!} + \dots \right] + \lambda$$

$$= \lambda^2 e^{-\lambda} e^{\lambda} + \lambda$$

$$\therefore E(X^2) = \lambda^2 + \lambda$$

Now, (1) becomes,

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

$$\boxed{\text{Var}(X) = \lambda}$$

Moment Generating Function (mgf)

The mgf of discrete random variable X is

$$M_X(t) = E(e^{tX}) = \sum_{x=0}^{\infty} e^{tx} p(x) = \left(\sum_{x=0}^{\infty} e^{tx} \frac{e^{-\lambda} \lambda^x}{x!} \right).$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} (\lambda e^t)^x = e^{-\lambda} \left[1 + \lambda e^t + \frac{(\lambda e^t)^2}{2!} + \dots \right]$$

$$= e^{-\lambda} e^{\lambda e^t} = e^{\lambda e^t - \lambda} = e^{\lambda(e^t - 1)}$$

$$\therefore M_X(t) = e^{\lambda(e^t - 1)}$$

Mean and Variance using mgf

The r th order raw moment about origin is,

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

$$M'_r = E(X^r) = \left. \frac{d M_X(t)}{dt} \right|_{t=0}$$

$$\begin{aligned} \text{Now, } M'_r &= \left. \frac{d M_X(t)}{dt} \right|_{t=0} = \left. \frac{d}{dt} e^{\lambda(e^t - 1)} \right|_{t=0} \\ &= \left. \frac{d}{dt} e^{\lambda(e^t - 1)} \right|_{t=0} = e^{\lambda(e^t - 1)} \cdot \lambda e^t \end{aligned}$$

$$\therefore E(X) = e^{\lambda(e^0 - 1)} \cdot \lambda e^0 = \lambda$$

$$M'_2 = E(X^2) = \left. \frac{d^2 M_X(t)}{dt^2} \right|_{t=0}$$

$$\begin{aligned} \text{Now, } M'_2 &= \left. \frac{d^2 M_X(t)}{dt^2} \right|_{t=0} = \left. \frac{d}{dt} \left[e^{\lambda(e^t - 1)} \cdot \lambda e^t \right] \right|_{t=0} \\ &= \lambda \left[e^t \left\{ e^{\lambda(e^t - 1)} \cdot \lambda e^t \right\} + e^{\lambda(e^t - 1)} \cdot \lambda e^t \right] \end{aligned}$$

$$\begin{aligned} E(X^2) &= \lambda \left[e^0 \left\{ e^{\lambda(e^0 - 1)} \cdot \lambda e^0 \right\} + e^0 \cdot e^{\lambda(e^0 - 1)} \cdot \lambda e^0 \right] \\ &= \lambda[\lambda + 1] = \lambda^2 + \lambda \end{aligned}$$

$$\text{Variance} = E(X^2) - [E(X)]^2 = \lambda^2 + \lambda - \lambda^2 = \lambda$$

$$\text{Var}(X) = \lambda$$

\therefore If $X \sim P(\lambda)$ then mean of X = Variance of $X = \lambda$

Negative Binomial Distribution (NBD)

$p = \text{Prob. of success in each trial}$

NA (ME)

$q = 1 - p = \text{Prob. of failure in ...}$

$k = \text{fixed no. of successes}$

Dep.

Indep.

$x = \text{no. of failures (r.v.)}$

$P(A \cap B)$

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

$x+k = \text{total no. of trials}$

If a trial is repeated independently till fixed no. of success occurs with prob. 'p' which is constant for each trial, then such experiment follows negative Binomial Distribution.

Derivation of NBD

Consider a binomial experiment consists of 'n' trials, and the prob. of success 'p' is constant for each trial. Let there be 'x' failures preceding the k^{th} success in $(x+k)$ trials. Here, $n = x+k$ trials are required to produce 'k' success (which is fixed) and x is a random variable.

Now, we need to find the probability that k^{th} success occurs in $(x+k)$ trials, which is the probability that exactly 'x' failures preceding the k^{th} success in $x+k$ trials.

Here, the last trials must be success with prob. 'p' and in the remaining $x+k-1$ trials, there are $(k-1)$ success and whose prob. is given by,

$${}_{x+k-1} C_{k-1} p^{k-1} q^x, \quad p+q=1$$

Now,

By Multiplication Theorem of Probability, the prob. of x failures preceding the k^{th} success in $x+k$ trials is given by,

$$({}_{x+k-1} C_{k-1} p^{k-1} q^x) \cdot p$$

$$= {}_{x+k-1} C_{k-1} p^k q^x$$

This is the prob. (prob. mass function) of NBD with parameter 'p' & 'k'

Note: If $X \sim NB(k, p)$ then,

$$P(X=x) = \frac{x+k-1}{k} C_{k-1} p^k q^x ; \quad x=0, 1, 2, \dots$$

$$= \frac{x+k-1}{k} C_x p^k q^x$$

$$C_{k-1} = \frac{x+k-1}{k} C_x$$

$$(1-p)^{-k} \sum_{x=0}^{\infty} C_x p^x q^{x-k}$$

Negative Binomial Distribution from GPSD.

In GPSD, take, $\theta = q/p$, and $f(\theta) = (1-\theta)^{-k}$

and, $S = \{0, 1, 2, \dots, \infty\}$

$$\text{Now, } f(\theta) = \sum_{x=0}^{\infty} a_x \theta^x = (1-\theta)^{-k} \quad \text{--- (i)}$$

$$\text{Now, } (1-\theta)^{-k} = \sum_{x=0}^{\infty} \left(\frac{x+k-1}{k} C_x \right) \theta^x \quad \text{--- (ii)}$$

From (i) & (ii), we get.

$$a_x = \frac{x+k-1}{k} C_x$$

From GPSD,

$$P(x) = \frac{a_x \theta^x}{f(\theta)} = \frac{\left(\frac{x+k-1}{k} C_x \right) \left(\frac{q/p}{1+q/p} \right)^x}{\left(1 - \frac{q/p}{1+q/p} \right)^{-k}}$$

$$P(x) = \frac{x+k-1}{k} C_x q^x p^x ; \quad x=0, 1, 2, \dots$$

Mean, Variance and Moment Generating of NBD.

If a d.r.v $X \sim NB(k, p)$ then the prob. of the distribution is given by.

$$P(X=x) = \frac{x+k-1}{k} C_{k-1} p^k q^x ; \quad x=0, 1, 2, \dots$$

$k > 0$

$$\begin{aligned}
 &= x+k-1 C_x p^k q^{x-k} \cdot (1-q) \cdot \frac{(1-q)}{1-p} \sum_{n=0}^{\infty} \frac{x+n-1}{n!} C_n q^n \\
 \text{Mean: } E(X) &= \sum_{x=0}^{\infty} x P(X=x) = \sum_{x=0}^{\infty} x \cdot x+k-1 C_x p^k q^{x-k} \\
 &= p^k \sum_{x=0}^{\infty} \frac{(x+k-1)!}{(x+k-1-x)! x!} q^x \\
 &= p^k \sum_{x=0}^{\infty} x \cdot (x+k-1)! \cdot q^{x-1} \cdot \frac{q}{(k+1)! x(x-1)!} \\
 &= q p^k k \sum_{x=1}^{\infty} \frac{(x+k-1)!}{x! k! (x-1)!} q^{x-1} \\
 &= q p^k k \sum_{x=1}^{\infty} x+k-1 C_{x-1} q^{x-1} \\
 &= q p^k k \sum_{x=1}^{\infty} (x-1)+(k+1)-1 C_{x-1} q^{x-1} \\
 &= k q p^k (1-q)^{-k-1} \\
 &= k q p^k p^{-k} \quad \because p+q=1 \\
 \therefore E(X) &= \boxed{kq}
 \end{aligned}$$

Variance:

$$\text{Var}(X) = E(X^2) - [E(X)]^2 = E(X^2) - \left(\frac{kq}{p}\right)^2 - ①$$

$$\begin{aligned}
 \text{Now, } E(X^2) &= \sum_{x=0}^{\infty} x^2 p(x) = \sum_{x=0}^{\infty} [x(x-1) + x] p(x) \\
 &= \sum_{x=0}^{\infty} x(x-1) p(x) + \sum_{x=0}^{\infty} x p(x) \\
 &= \sum_{x=0}^{\infty} x(x-1) \frac{x+k-1}{x!} C_x p^k q^{x-k} + E(X) \\
 &= p^k \sum_{x=0}^{\infty} x(x-1) \frac{(x+k-1)!}{(x+k-1-x)! x!} q^x + E(X) \\
 &= p^k q^2 \sum_{x=1}^{\infty} x(x-1) \frac{(x+k-1)!}{(k-1)! x(x-1)(x-2)!} q^{x-2} + E(X)
 \end{aligned}$$

$$\begin{aligned}
 &= p^k q^2 k(k+1) \sum_{x=2}^{\infty} \frac{(x+k-1)!}{(k+1)k(x-1)!(x-2)!} q^{x-2} + E(x) \\
 &= p^k q^2 k(k+1) \sum_{x=2}^{\infty} \frac{(x+k-1)!}{(k+1)!(x-2)!} q^{x-2} + E(x) \\
 &= p^k q^2 k(k+1) \sum_{x=2}^{\infty} \frac{x+k-1}{C_{x-2}} C_{x-2} q^{x-2} + E(x). \quad C_{x-2} = \frac{(x+k-1)!}{(x+k-1-x+2)!} \\
 &= p^k q^2 k(k+1) \sum_{x=2}^{\infty} \frac{(x-2)+(k+2)-1}{C_{x-2}} C_{x-2} q^{x-2} + E(x) \quad = \frac{(x+k-1)!}{(k+1)!(x-2)!} \\
 &= p^k q^2 k(k+1) \left(\frac{1}{p} \right)^{-k-2} + \frac{kq}{p} \quad (1-q)^{-n} = \sum_{x=0}^{\infty} C_x q^x \\
 &= p^k q^2 k(k+1) \frac{1}{p}^{-k-2} + \frac{kq}{p} \\
 &= \frac{q^2 k^2}{p^2} + \frac{q^2 k}{p} + \frac{kq}{p} = \frac{q^2 k^2}{p^2} + \frac{q^2 k}{p^2} + \frac{kq}{p}.
 \end{aligned}$$

Now, Eq ① becomes,

$$\begin{aligned}
 \text{Var}(X) &= \frac{q^2 k^2}{p^2} + \frac{q^2 k}{p} + \frac{kq}{p} - \frac{k^2 q^2 k}{p^2} \\
 &= \frac{kq}{p} \left(\frac{q}{p} + 1 \right) \\
 &= \frac{kq}{p} \left(q+p \right)
 \end{aligned}$$

$$\boxed{\text{Var}(X) = \frac{kq}{p^2}}$$

Moment Generating function

$$\begin{aligned}
 \text{Now, the mgf } M_X(t) &= E(e^{tx}) = \sum_{x=0}^{\infty} e^{tx} p(x) = \sum_{x=0}^{\infty} e^{tx} x^{x+k-1} C_x p^k q^x \\
 &= p^k \sum_{x=0}^{\infty} x^{x+k-1} C_x (qe^t)^x = p^k (1-qe^t)^{-k}
 \end{aligned}$$

$$\boxed{M_X(t) = \frac{p^k}{(1-qe^t)^k}}$$

Mean and Variance using MGF

If $X \sim NB(k, p)$, then,

$$M_X(t) = \frac{p^k}{(1-qe^t)^k}$$

$$(1-qe^t)^k$$

The r th raw moment about origin,

$$M'_r = d^r M_X(t)$$

$$d^{tr}$$

$$\Big|_{t=0}$$

Now, Mean:

$$M'_1 = E(X) = \frac{d M_X(t)}{dt} \Big|_{t=0}$$

$$\Big|_{t=0}$$

$$\begin{aligned} \text{Now } \frac{d M_X(t)}{dt} &= \frac{d \{p^k (1-qe^t)^{-k}\}}{dt} = p^k \frac{d}{dt} (1-qe^t)^{-k} \cdot \frac{d(1-qe^t)}{dt} \\ &= p^k (-k) (1-qe^t)^{-k-1} (-q)e^t \end{aligned}$$

$$= p^k k q e^t (1-qe^t)^{-k-1}$$

$$\therefore M'_1 = p^k k q e^0 (1-qe^0)^{-k-1} = p^k k q (1-q)^{-k-1} = p^k k q p^{-k-1}$$

$$\boxed{E(X) = \frac{kq}{p}}$$

$$M'_2 = \frac{d^2 M_X(t)}{dt^2} \Big|_{t=0}$$

$$\frac{d^2 M_X(t)}{dt^2} = p^k k q \frac{d}{dt} \{e^t (1-qe^t)^{-k-1}\}$$

$$= p^k k q [(1-qe^t)^{-k-1} e^t + e^t (-k-1) (1-qe^t)^{-k-2} (-qe^t)]$$

$$\therefore M'_2 = p^k k q [(1-qe^0)^{-k-1} e^0 + e^0 (k+1) (1-qe^0)^{-k-2} q e^0]$$

$$= p^k k q [(1-q)^{-k-1} + (k+1)(1-q)^{-k-2} q]$$

$$= p^k k q [p^{-k-1} + (k+1) p^{-k-2} q]$$

$$= \frac{kq}{p} \left[1 + \frac{(k+1)q}{p} \right]$$

$$E(X) = \frac{kq}{p} + \frac{k^2 q^2}{p^2} + \frac{kq^2}{p^2}$$

Variance:

$$\begin{aligned} \text{Var}(X) &= E(X^2) - [E(X)]^2 \\ &= \frac{kq}{p} + \frac{k^2 p^2}{p^2} + \frac{kq^2}{p^2} - \frac{k^2 q^2}{p^2} \\ &= \frac{kq}{p} (1 + \frac{1}{p}) = \frac{kq}{p^2} (p+q) \\ \therefore \text{Var}(X) &= \frac{kq}{p^2} \end{aligned}$$

- Q. If the prob. that an applicant will pass the driving license test on any given trial is 0.75. Find the prob. that an applicant will finally pass the test in the fourth try.

$$x+k=4$$

$$P(X=3) = ?$$

$$x=3$$

$$k=1$$

$$p = 0.75$$

$$q = 0.25$$

$$x+k-1$$

$$C_x p^x q^{x+k-1}$$

p = prob. if an applicant will pass the test = 0.75

$$q = 1-p = 0.25$$

Here, X = no. of failures.

Here $X \sim NB(k, p)$ so, $P(X=x) = {}^{x+k-1}C_x p^k q^x$; $x=0, 1, \dots$

$$\text{we have, } x+k=4 ; k=1$$

$$\therefore x=3$$

Now, the prob. that applicant will pass the test on fourth try is.

$$P(X=3) = {}^{3+1-1}C_3 (0.75)^1 (0.25)^3 = 0.01$$

- Q. The probability of hitting the target at any trial is 0.2. If a shooter aims at a target, find the prob. that the fifth fire is second hit.

$$p = 0.2, q = 0.8$$

$$x+k=5, k=2$$

Q3. A boy is throwing stones at a target. If the prob. of hitting the target at any point is $\frac{1}{2}$, what is the prob. that his 10th throw is 5th hit?

Soln,

$p = \text{prob. of hitting a target at any trial} (= 0.5)$

$$q = 0.5$$

$$k = 5$$

$$x+k = 10$$

$$P(X=5) = {}^{2+k-1}C_k p^k q^{k-1}$$

Now, the prob. that his 10th throw is 5th hit is given by,

$$P(X=5) = {}^{5+5-1}C_5 (0.5)^5 (0.5)^5 = 0.123$$

The prob. that a strike will make a goal in a football match is 0.2. What is the prob. that he will do fourth goal in fifth shot (0.0051)

Multinomial Distribution

An experiment is said to be multinomial experiment, if:

- (i) The experiment consists of n fixed trials.
- (ii) For each trial there are $k \geq 2$ i.e. 3 or more possible outcomes.
- (iii) The trials are independent.
- (iv) The prob. for each outcome remains the same from trial to trial.

The examples of the multinomial experiments are,

- (i) No. of throw by a fair dice in which each throw can result six different outcomes.
- (ii) A number of selection or drawing's of balls at random with replacement from a box containing 30 balls of which 10 are white, 15 are black and 5 are red.

Definition:

The random variable $X = X_1, X_2, \dots, X_k$ denoting the outcome of n trials, where $X_i = \text{freq}^n$ or number of outcome E_i with respective probability

P_i ($i = 1, 2, \dots, k$), is said to have multinomial distribution with parameters

$(n, P_1, P_2, \dots, P_k)$ if its prob. is given by,

$$P(X_1=x_1, X_2=x_2, \dots, X_k=x_k) = \frac{n!}{x_1! x_2! \dots x_k!} P_1^{x_1} P_2^{x_2} \dots P_k^{x_k}$$

Where, $\sum_{i=1}^k x_i = n$ and $\sum_{i=1}^k P_i = 1$.

Q. An urn contains 8 red balls, 3 yellow balls and 9 white balls.

6 balls are drawn at random with replacement. what is the prob. that 2 are red, 1 is yellow and 3 are white?

Soln:

Since, 6 balls are randomly selected with replacement so there are 6 independent trials of a multinomial experiment.

$P_1 = \text{Prob. of getting a red ball in any trial} = 8/20$

$P_2 = \text{prob. of getting a yellow ball in any trial} = 3/20$

$P_3 = \text{prob. of getting a white ball in any trial} = 9/20$

The prob. of multinomial dist. with parameters $(n, p_1, p_2, \dots, p_k)$ is

$$P(X_1=x_1, X_2=x_2, \dots, X_k=x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

The prob. of getting 2R, 1Y and 3W is given by,

$$P(X_1=2, X_2=1, X_3=3) = \frac{6!}{2! 1! 3!} \left(\frac{8}{20}\right)^2 \left(\frac{3}{20}\right)^1 \left(\frac{1}{20}\right)^3$$

- Q. The painted light balls produced by a company are 50% red, 30% blue and 20% green. In a sample of 5 balls; find the prob that 2 red, 1 is green and 2 are blue.

Soln,

$$P_1 = \text{Red} = \left(\frac{1}{2}\right)$$

$$P_2 = \text{Blue} = \left(\frac{3}{5}\right)$$

$$P_3 = \text{Green} = \left(\frac{2}{5}\right)$$

$$P(X_1=2, X_2=1, X_3=2) = \frac{5!}{2! 1! 2!} (0.5)^2 (0.2)^1 (0.3)^2$$

$$= 0.135$$

Reduction of Multinomial distribution into binomial distribution.

When $k=2$, then the multinomial distribution has prob. mass function

$$P(X_1=x_1, X_2=x_2) = \frac{n!}{x_1! x_2!} p_1^{x_1} p_2^{x_2} \quad \text{if } x_1+x_2=n$$

$$P(x_1, n-x_1) = \frac{n!}{x_1! (n-x_1)!} p_1^{x_1} (1-p_1)^{n-x_1} \quad \text{if } x_1 = n \quad \text{and} \quad p_2 = 1-p_1$$

$$\text{If } p_1=p \text{ and } x_1=x \text{ then } P(x) = \frac{n!}{x! (n-x)!} p^x (1-p)^{n-x}$$

$$\Rightarrow P(x) = {}^n C_x p^x q^{n-x}$$

Moment Generating function of Multinomial distribution.

If $X = (X_1, X_2, \dots, X_k) \sim MD(n, P_1, P_2, \dots, P_k)$
 Then,

$$P(X) = \frac{n!}{x_1! x_2! \dots x_k!} P_1^{x_1} P_2^{x_2} \dots P_k^{x_k}$$

Now, the moment generating function
 of X is given by,

$$M_X(t) = M_{X_1, X_2, \dots, X_k}(t_1 + t_2 + \dots + t_k)$$

$$P(X) = \frac{n!}{x_1! x_2! \dots x_k!} P_1^{x_1} P_2^{x_2} \dots P_k^{x_k}$$

$$\sum P(X) = \sum \left[\frac{n!}{x_1! x_2! \dots x_k!} P_1^{x_1} P_2^{x_2} \dots P_k^{x_k} \right]$$

$$= (P_1 + P_2 + \dots + P_k)^n$$

$$= 1$$

$$e^{t_1 x_1} e^{t_2 x_2} \dots e^{t_k x_k}$$

$$= E[e^{t_1 X_1 + t_2 X_2 + \dots + t_k X_k}]$$

$$= \sum_n \left[e^{t_1 x_1 + t_2 x_2 + \dots + t_k x_k} \right] \frac{n!}{x_1! x_2! \dots x_k!} P_1^{x_1} P_2^{x_2} \dots P_k^{x_k}$$

$$= \sum_n \frac{n!}{x_1! x_2! \dots x_k!} (P_1 e^{t_1})^{x_1} (P_2 e^{t_2})^{x_2} \dots (P_k e^{t_k})^{x_k}$$

$$\therefore M_X(t) = (P_1 e^{t_1} + P_2 e^{t_2} + \dots + P_k e^{t_k})^n$$

This is the moment generating function of multinomial distribution.

If $t_i \neq 0$ and $t_i = 0$ for all $i = 2, 3, \dots, k$

$$M_{X_1}(t_1) = M_{X_1, X_2, X_3, \dots, X_k}(t_1, 0, \dots, 0)$$

$$= (P_1 e^{t_1} + P_2 e^0 + \dots + P_k e^0)^n$$

$$= (P_1 e^{t_1} + P_2 + \dots + P_k)^n \quad \because P_1 + P_2 + \dots + P_k = 1$$

$$= [P_1 e^{t_1} + (1 - P_1)]^n \quad P_2 + \dots + P_k = 1 - P_1$$

$$\therefore M_{X_1}(t_1) = (P_1 e^{t_1} + q_1)^n$$

This is the moment generating function of binomial distribution with parameters n and P_1 .

Hence, by uniqueness property of mgf $X_1 \sim B(n, P_1)$

\Rightarrow for all $i = 1, 2, \dots, k$,

$$X_i \sim B(n, P_i)$$

$$\therefore E(X_i) = n P_i \quad \text{and} \quad V(X_i) = n P_i q_i$$

Characteristic function of Multinomial Distribution

$$\phi_X(t) = E[e^{itX}]$$

If $X = (X_1, X_2, \dots, X_k) \sim M(n, p_1, p_2, \dots, p_k)$

then the prob. of the distribution is given by.

$$= \sum e^{itx} p(x) \quad (D)$$

$$= \int_{-\infty}^{\infty} e^{itx} f(x) dx \quad (C)$$

$-\infty < x < \infty$

$$P(X_1=x_1, X_2=x_2, \dots, X_k=x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

Now, the characteristic function of X is given by,

$$\phi_X(t) = \phi_{X_1, X_2, \dots, X_k}(t_1, t_2, \dots, t_k)$$

$$= E[e^{it_1 x_1} e^{it_2 x_2} \dots e^{it_k x_k}]$$

$$= \sum_x [e^{it_1 x_1} e^{it_2 x_2} \dots e^{it_k x_k}] \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

$$= \sum_x \frac{n!}{x_1! x_2! \dots x_k!} (p_1 e^{it_1})^{x_1} (p_2 e^{it_2})^{x_2} \dots (p_k e^{it_k})^{x_k}$$

$$\phi_X(t) = (p_1 e^{it_1} + p_2 e^{it_2} + \dots + p_k e^{it_k})^n.$$

Covariance between x_i and x_j in Multinomial distribution.

$$(i+j=1, 2, \dots, k)$$

$$\text{cov}(x_i, x_j) = \frac{E(x_i x_j) - E(x_i) E(x_j)}{n}$$

Covariance between x_i and x_j .

$$\text{cov}(x_i, x_j) = E(x_i x_j) - E(x_i) E(x_j)$$

$$= E[\{x_i - E(x_i)\}\{x_j - E(x_j)\}]$$

$$= E(x_i x_j) - (np_i)(np_j) \quad (1)$$

$$= \frac{1}{n} \sum xy - \bar{x}\bar{y}$$

$$\text{Now, } E(x_i x_j) = \left. \frac{d^2 M_X(t)}{dt_i dt_j} \right|_{t=0} = E(xy) - E(x) E(y)$$

$$= \left. \frac{d^2}{dt_i dt_j} [p_1 e^{it_1} + p_2 e^{it_2} + \dots + p_k e^{it_k}]^n \right|_{t=0}$$

$$= \left. \frac{d}{dt_i} \left[\frac{d}{dt_j} (p_1 e^{it_1} + \dots + p_k e^{it_k})^n \right] \right|_{t=0}$$

$$= \left. \frac{d}{dt_i} \left[n (p_1 e^{it_1} + \dots + p_k e^{it_k})^{n-1} p_j e^{it_j} \right] \right|_{t=0}$$

$$= n p_j e^{it_j} \left. \frac{d}{dt_i} \{ (p_1 e^{it_1} + \dots + p_k e^{it_k})^{n-1} \} \right|_{t=0}$$

$$\begin{aligned}
 E(X_i X_j) &= n p_i e^{t_i} (n-1) (p_1 e^{t_1} + \dots + p_k e^{t_k})^{n-2} p_j e^{t_j} \Big|_{t=0} \\
 &= n p_j (n-1) (p_1 + \dots + p_k)^{n-2} p_i e^{t_i} \Big|_{t=0} \\
 &= n p_i p_j (n-1) (1)^{n-2} \left[\because \sum_{i=1}^k p_i = 1 \right] \\
 &\quad [p_1 e^{t_1} + \dots + p_i e^{t_i} + \dots + p_j e^{t_j} + \dots + p_k e^{t_k}]^n
 \end{aligned}$$

Now, eqn ① becomes,

$$\text{Cov}(X_i, X_j) = n p_i p_j (n-1) - (n p_i) (n p_j).$$

$$= n^2 p_i p_j - n p_i p_j - n^2 p_i p_j.$$

$$\therefore \text{Cov}(X_i, X_j) = -n p_i p_j$$

It means, the multinomial distribution is negatively correlated

Correlation between X_i & X_j :

$$\begin{aligned}
 r &= \frac{\text{Cov}(X_i, X_j)}{\text{SD}(X_i) \cdot \text{SD}(X_j)} & V(X_i) &= n p_i q_i \\
 &= \frac{-n p_i p_j}{\sqrt{n p_i (1-p_i)} \sqrt{n p_j (1-p_j)}} & &= n p_i (1-p_i) \\
 &= -\frac{\sqrt{p_i} \sqrt{p_j}}{\sqrt{(1-p_i)} \sqrt{(1-p_j)}}
 \end{aligned}$$

$$\therefore r = -\left(\frac{p_i p_j}{(1-p_i)(1-p_j)} \right)^{1/2}$$

A fair dice is thrown five times,

- ① What is the prob of getting one 6, one 5, two 4, one 3 and no others.
- ② What is the prob. that the sum of the scores is exactly 8?

$$\begin{aligned}
 P(\text{Sum } 8) &= P(1, 1, 2, 2, 2) + P(1, 1, 1, 2, 3) + \\
 &\quad P(1, 1, 1, 1, 4)
 \end{aligned}$$

Trails:	1	2	3	4	5	Sum
Faces:	1	1	2	2	2	8
	1	1	1	2	3	8
	1	1	1	1	4	8

Soh:

Since, a fair dice is thrown five times so we have 5 independent multinomial experiments and at each trial we have 6 possible outcomes i.e. 1, 2, 3, 4, 5 and 6 each with prob. $1/6$ & 1
 $\therefore n = 5, X_i = 1, 2, 3, 4, 5, 6, f(x_1) = 1/6$

① Now, the prob. of getting one 6, one 5, two 4, one 3 and no other is given by,

$$P(X_1=0, X_2=0, X_3=1, X_4=2, X_5=1, X_6=1)$$

$$= \frac{5!}{0! 0! 1! 2! 1! 1!} \left(\frac{1}{6} \right)^0 \left(\frac{1}{6} \right)^0 \left(\frac{1}{6} \right)^2 \left(\frac{1}{6} \right)^1 \left(\frac{1}{6} \right)^1 \left[\frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k} \right]$$

$$= 0.0077.$$

② A sum of 8 can be obtained from any three mutually exclusive ways (1, 1, 2, 2, 2) or (1, 1, 1, 2, 3) or (1, 1, 1, 1, 4) [In any Order]

$$\therefore P(\text{sum } 8) = P(X_1=2, X_2=3) + P(X_1=3, X_2=1, X_3=1) + P(X_1=4, X_4=1)$$

$$= \frac{5!}{2! 3!} \left(\frac{1}{6} \right)^2 \left(\frac{1}{6} \right)^3 + \frac{5!}{3! 1! 1!} \left(\frac{1}{6} \right)^3 \left(\frac{1}{6} \right)^1 \left(\frac{1}{6} \right)^1 + \frac{5!}{4! 1!} \left(\frac{1}{6} \right)^4 \left(\frac{1}{6} \right)^1$$

$$= 0.0045$$

H.W.

Type	D	A	B	AB
Prob.	0.44	0.42	0.10	0.04

In a random sample of 10 Americans what is the prob. that 6 have D, 2 have A, 1 has B & 1 has AB?

① The prob. that player A will win any game is 30%, the prob. that player B will win is 20% and the prob. that player C will win is 50%. If they play 7 games, what is the prob. that A will win 2 games, player B will win 2 games and player C will win 3 games?

③ A fair dice is thrown 15 times. What is the prob. of getting five 6s, four 5s, three 4s, two 3s, one 2 and no 1?

Non-Parametric Tests \rightarrow No Distribution, No Assumption on Data.

Hypothesis Testing

\rightarrow parametric

\rightarrow Non-parametric

Parameter \rightarrow population $\rightarrow \mu, \rho, \sigma^2$

Statistics \rightarrow Sample $\rightarrow \bar{x}, r, s^2$

Parameter

Statistical Measures which are obtained from population (constant)

Statistics

Statistical Measures which are obtained from Sample. (variable)

Types of Hypothesis

1. Null Hypothesis (H_0) : $\mu = \mu_0$.

\rightarrow No difference between sample & population parameters.

\rightarrow True - Expected = 0

2. Alternative hypothesis (H_1)

\rightarrow Complementary of Null Hypothesis.

\rightarrow Mutually Exclusive

$H_1: \mu \neq \mu_0$ (TTT) \rightarrow No direction

, $H_1: <, >, =$ (OTT)

Types of Error:

\rightarrow Reject H_0 when it is true

\rightarrow Accept H_0 when it is false.

Type-I : $\alpha = \text{prob} = \text{level of significance}$

Type-II : β

Type-I Error $\rightarrow \downarrow$ Risk

α

producer's Risk

Situation Decision	H_0 True	H_0 False
Reject H_0	WD: Type-I	Correct Decision
Accept H_0	Correct Decision	WD: Type-II

Type-II Error $\rightarrow \uparrow$ Risk

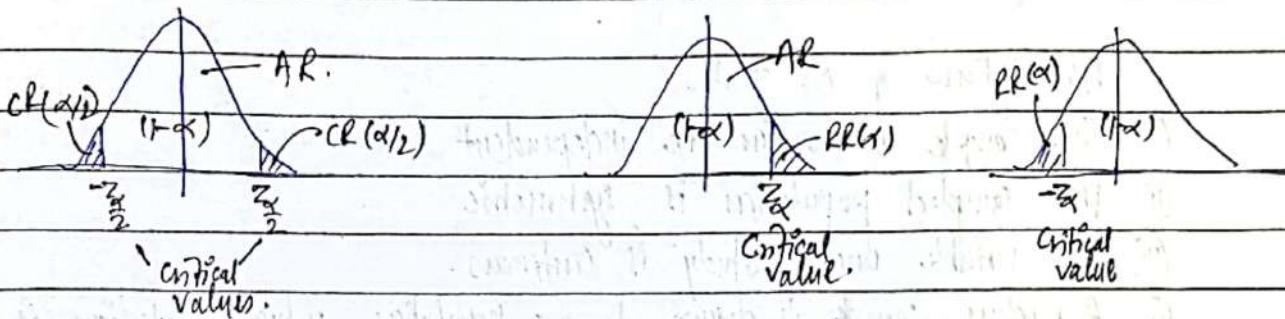
β

consumer's Risk.

Critical Value (Significant Value)

- value that separates the acceptance and rejection region.
- It depends on:
 - ① Size of the level of significance used.
 - ② Type of the used. (OTT or TTT)

TTT UTT (RTT) LTT



$p\text{-value} > \alpha$ insignificant $\Rightarrow H_0$ is true (Do not reject H_0)

$p\text{-value} \leq \alpha$ significant $\Rightarrow H_0$ is False (H_1 is True) (Reject H_0)

Degree of freedom (df)

- ① One Sample $\rightarrow df = n - 1$

Identification of One-tailed test and two-tailed test.

One Tailed Test (OTT)

If the direction of difference like at least, at most, increase, decrease, minimum, maximum, superior, inferior, less than, more than, etc.

are involved in the statement of hypothesis then we use OTT

Otherwise Two Tailed Test (TTT)

Non-Parametric Test

Non-parametric (NP) tests are defined as those statistical tests of hypothesis in which no parameter is involved and which are based on other statistic.

In other words, those test of hypothesis are said to be non-parametric tests if the hypothesis does not involve any parameter of the population and if the measurements are on the nominal or ordinal scale.

Assumptions of NP Tests.

- (i) The sample observations are independent
- (ii) The sampled population is symmetric
- (iii) The variable under study is continuous.
- (iv) A random sample is drawn from a population whose median is unknown.

Median Test

$n > 20 \rightarrow$ Large Sample.

Individual Series

$$\text{Median } (Md) = \text{Value of } \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term.}$$

Median Test is a non-parametric test used to test the difference in medians of two independent distributions. In other words this test is applied to test whether two independent random samples have been drawn from two populations with same median or not. Also it is used to test whether the two treatments applied in an experiment are equally effective or not.

Procedure

$H_0 : Md_1 = Md_2$ i.e. There is no significant difference between two medians.

$H_1 : Md_1 \neq Md_2$ i.e. There is significant difference between two medians (TGT)

or, $H_1 : Md_1 > Md_2$ (Right) } OTT

or, $H_1 : Md_1 < Md_2$ (Left)

Small Sample Size

$n_1 \leq 10$ and $n_2 \leq 10$

Test - statistic

The test - statistic is obtained from following procedure :-

- (i) Combine all the observations of both samples in ascending order of magnitude and find the value of median.
- (ii) Count no. of observations less or equal to Median in first sample and denote it by 'a'.
- (iii) The test - statistic is 'a'.

Critical Region

We can obtain the p-value (P_0), the probability associated with the value as extreme as the observed 'a' for n_1 and n_2 .

$$\text{i.e. } P_0 = P(A \geq a)$$

$$\text{Where, } P_0 = \frac{n_1 C_a + n_2 C_{k-a}}{n_1 + n_2 C_k}; \quad k = \frac{n_1 + n_2}{2}$$

$a = 0, 1, 2, \dots, \min(n_1, k)$

Decision

For TTT:

If p-value = $2P_0 > \alpha$, then we do not reject H_0 .
otherwise reject H_0 .

For OTT

If p-value = $P_0 > \alpha$, then we do not reject H_0 .
otherwise reject H_0 .

The following are the yield data of 10 plots under two treatment x & y.

x : 46 45 32 42 39 48 49 30 51 34

y : 44 40 59 47 55 47 50 71 43 55

Use Median test to test the effectiveness of two treatments.

$\alpha = 0.05$, if not given.

Date _____
Page _____

Soln,

H_0 : Both the treatments X and Y are equally effective.

H_1 : Both the treatments X and Y are not equally effective (TTT)

Test statistic

Under H_0 , the test statistic is 'a'

Calculation

Combining both samples in ascending order of magnitude

30, 32, 34, 39, 40, 42, 43, 44, 45, (46, 47) 47,
48, 49, 50, 51, 55, 55, 59, 71

$M_d = \text{Value of } \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}, n=20$

= value of 10.5th term.

$$\approx \frac{46+47}{2} = 46.5$$

Now, no. of observations term or equal to Median in the first sample = 7
 $\therefore a = 7$

Calculation of p-value.

$$n_1 = 10$$

$$n_2 = 10$$

$$k = \frac{n_1+n_2}{2} = 10$$

$$= \sum_{q=7}^{10} \frac{n_1 C_q}{n_1+n_2} \frac{n_2 C_{k-q}}{C_k}$$

$$= \sum_{q=7}^{10} \frac{^{10}C_q}{20 C_{10}} \frac{^{10}C_{10-q}}{^{10}C_k}$$

$$= \frac{1}{20 C_{10}} \left[^{10}C_7 \cdot ^{10}C_{10-7} + ^{10}C_8 \cdot ^{10}C_{10-8} + ^{10}C_9 \cdot ^{10}C_{10-9} + ^{10}C_{10} \cdot ^{10}C_{10-10} \right]$$

$$= 0.089$$

Decision:

$$P\text{-value} = 2P_0 \quad [\because P_0 \text{ is TTT}]$$

$$\approx 2(0.089)$$

=

$$\alpha = 0.05$$

Since $P\text{-value} > \alpha$ so we do not reject H_0 .

Hence, both the treatment x and y are equally effective.

Median Test Small Sample Size $n_1, n_2 \leq 10$

- Q. A quality controller wishes to determine whether there is a difference in outcome between two different tools of software I and II. The following data shows the outcome of two different tools. Can the controller conclude that a difference exists? Use median test at 5% level of significance.

Software I 24.0 16.7 22.8 19.8 18.9

Software II 23.2 19.8 18.1 17.6 20.2 17.8

Soln,

H_0 : There is no significant difference between the outcomes of two different tools.

H_1 : There is significant difference between the outcomes of two different tools.

(TTT)

Combining both samples in ascending order of magnitude.

16.7 17.6 17.8 18.1 18.6 19.8 19.8 20.2 22.8 23.2 24.0

$$Md = \text{Value of } \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}; \quad n=11 \quad [\because n_1=5 \text{ and } n_2=6]$$

= value of 6th term.

$$= 19.8$$

No. of observations sum or equal to median is first sample (a) = 9.

Test statistic: Under H_0 , $a=3$

Calculation of p-value:

We have, $n_1 = 5, n_2 = 6$

$$\therefore k = \frac{n_1+n_2}{2} = 5.5 \approx 6$$

$$P_0 = P(A \geq a)$$

$$= \sum_{a=3}^{n_1} \frac{\binom{n_1}{a} \binom{n_1+n_2}{k}}{\binom{n_1+n_2}{k}}$$

$$= \frac{5!}{3!} \frac{6!}{6-3!} + \frac{5!}{4!} \frac{6!}{6-4!} + \frac{5!}{5!} \frac{6!}{6-5!}$$

$$= \frac{5!}{6!}$$

$$= 0.608$$

Since, the problem is two-tailed test so,

$$P\text{-value} = 2P_0 = 2(0.608) \approx 1$$

Decision:

We have, $\alpha = 0.05$ and p-value 1.

Since p-value > α so we do not reject H_0 .

Hence, there is no significant difference between the outcomes of two different tools.

Large Sample Size i.e. $n_1 > 10$ or $n_2 > 10$. [Chi-Square Test]

[χ^2 -test]

	No. of obs $\leq M_d$	No. of obs $> M_d$	Total
First	a	c	a+c
Second	b	d	b+d
Total	(a+b)	(c+d)	$n = a+b+c+d$

Test Statistic. Under H_0 :

$$\chi^2 = \frac{n(ad - bc)^2}{(a+c)(b+d)(a+b)(c+d)}$$

$$n \chi^2_{(1)}$$

$$df$$

$$df = (r-1)(c-1) \\ = (2-1)(2-1) = 1$$

Decision ..

AR

RR

$$\chi^2_{1,1} = 3.84$$

If test statistics value falls in Rejection Region(RR) then we reject H_0 otherwise don't reject H_0 .

If $\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$, \Rightarrow Reject H_0 else don't reject H_0 .

- Q. The same C programming papers were marked by two A and B. The final marks were recorded as follows:

Teacher A 73 89 82 43 80 73 66 45 93 36 77 60

Teacher B 88 78 91 48 85 74 77 31 38 62 76 77

Using median test at 5% level of significance to determine if the marks distribution of two teachers differ significantly.

Stem (10)	leaf (1)
3	6 1
4	3 5 8
5	
6	6 0 2
7	3 3 7 8 4 7 8 6 7
8	9 2 0 8 5
9	3 1

Soln,

H_0 : The marks distribution of two teachers do not differ significantly.

H_1 : The marks distribution of two teachers differ significantly.

Calculation.

Arranging given data in ascending order of magnitude.

31, 36, 43, 45, 48, 60, 62, 66, 73, 74, 76, 77, 77, 77, 78, 78, 80, 82, 85, 88, 89
91, 93.

$$Md = \text{Value of } \left(\frac{n+1}{2} \right)^{\text{th}} \text{ item}, n=24$$

= Value of 12.5th term

$$= \frac{76+77}{2} = 76.5$$

Test-statistic : For large sample, the test-statistic is given by,

$$\chi^2 = \frac{n(ad - bc)^2}{(a+c)(b+d)(a+b)(c+d)}$$

$$= 24(49-45)^2$$

12.12.12.12

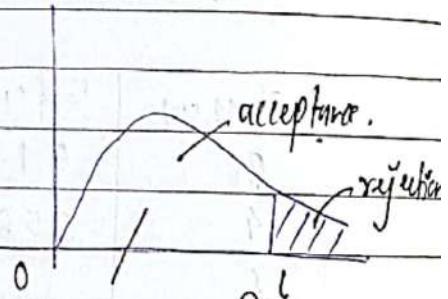
$$= 0.667$$

	$X_i \leq Md$	$X_i > Md$	Total
First	a=7	c=5	12
Second	b=5	d=7	12
Total	12	12	n=24

Critical Value :

We have, $\alpha = 0.05$, and $df = 1$

$$\chi_{\text{tab}}^2 = \chi_{1, 0.05}^2 = 3.841$$



Decision:

Since, the test-statistic value falls in the acceptance region, so do not reject H_0 .

Since $\chi_{\text{cal}}^2 < \chi_{\text{tab}}^2$ so we do not reject H_0 .

Hence, the marks distribution of two teachers do not differ significantly.

Mann Whitney 'U' test.

This test is most powerful non-parametric test for testing the hypothesis of difference between two independent locations of two independent random samples. This is non-parametric alternative test of t-test.

Procedure

$H_0: M_{d1} = M_{d2}$ i.e. There is no significant difference between two location or two Medians.

$H_1: M_{d1} \neq M_{d2}$ i.e. There is significant difference between two location or two Med.

or $H_1: M_{d1} > M_{d2}$ (Right) } OTT

or $H_1: M_{d1} < M_{d2}$ (Left)

Small Sample Size, i.e. n_1 or $n_2 \leq 10$

Test statistic: Under H_0 .

$$U_0 = \text{minimum } \{ U_1 \text{ and } U_2 \}$$

$$\text{Where, } U_1 = \frac{n_1 \cdot n_2 + \frac{n_1(n_1+1)}{2} - R_1}{2}$$

$$U_2 = \frac{n_1 \cdot n_2 + \frac{n_2(n_2+1)}{2} - R_2}{2}$$

Where, $R_1 = \text{sum of the combined ranks given to the first sample.}$

$R_2 = \text{sum of the combined ranks given to the second sample.}$

$$U_1 + U_2 = n_1 \cdot n_2$$

Critical Value

for two tailed test, $U_{tab} = U_{\alpha/2}, (n_1, n_2)$

for One tailed test, $U_{tab} = U_{\alpha}, (n_1, n_2)$

Decision

If $U_0 > U_{tab}$ then we do not reject H_0 . Otherwise reject H_0 .

Q. The nicotine contents of two brands of cigarettes, measured in milligram was found to be as follows.

Brand A : 2.1 4.0 6.3 5.4 4.8 3.7 6.1 3.3

Brand B : 4.1 0.6 3.1 2.5 3.0 6.2 1.6 2.2 1.9 5.4

Is there any significant difference between two brands of cigarettes.
Use Mann-Whitney U-test at $\alpha = 0.10$?

Soln

H_0 : There is no significant difference between two brands of cigarettes.

H_1 : There is significant " " " " " "

Test statistic: Under H_0 .

$$U_0 = \text{minimum } \{ U_1 \text{ and } U_2 \}$$

Check

$$\left\{ \begin{array}{l} U_1 + U_2 = n_1 n_2 \\ 23 + 57 = 8(10) \end{array} \right\}$$

Calculation:

A	B	Combined Rank		
		A	B	
2.1	4.1	4	12	n_{10}
4.0	0.6	10.5	1	$U_1 = n_1 n_2 + n_1(n_1+1) - R_1$
4.0 $\rightarrow 10$	6.3	3.1	18	2
4.0 $\rightarrow 11$	5.4	2.5	14.5	$= 8(10) + 8(8+1) - 93$
21/2	4.8	4.0	13	2
$= 10.5$	3.7	6.2	9	$= 0.3$
6.1	1.6	16	2	and
3.1	2.2	8	5	$U_2 = n_1 n_2 + n_2(n_2+1) - R_2$
1.9			3.	2
5.4		14.5		$= 8(10) + 10(10+1) - 78$
$R_1 = 93$	$R_2 = 78$			2.
				$= 57$

Critical Value:

We have $\alpha = 0.10$ and $(n_1, n_2) = (8, 10)$

$$\therefore U_{0.10, (8, 10)} = U_{0.10, (8, 10)} \quad (\text{TTT})$$

$$= 20$$

Decision :

Since, $U_0 = 28 > U_{0.10, (8,10)} = 20$, we do not reject H_0 .

Hence, there is no significant difference between.....

Large Sample Size i.e. $n_1 > 10$ or $n_2 > 10$

For Large Sample size the distribution of U_0 is normal with mean

$$= \frac{n_1 n_2}{2} \quad \text{and S.d.} = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

$$\therefore Z = \frac{U_0 - \mu_{U_0}}{\sigma_{U_0}}$$

$$= \frac{U_0 - \frac{n_1 n_2}{2}}{\sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}}$$

$$X \sim N(\mu, \sigma^2)$$

i.e. X follows ND.
with mean μ & SD σ .

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

std. Normal.

Critical Value:

The critical value of the test-statistic is obtained from z-table at α level of significance for one tailed and two tailed test.

$$\therefore Z_{\text{tab}} = Z_\alpha$$

$$\alpha = 0.05$$

One Two

1.645	1.96
-------	------

Decision:

If $|Z_{\text{cal}}| > Z_{\text{tab}}$ then we reject H_0 otherwise do not reject H_0 .

Note:

In case of tied rank

$$\sigma_{U_0} = \sqrt{\frac{n_1 n_2}{n(n-1)}} \left[\frac{n^3 - n}{12} - \sum T_i^3 \right]$$

$$\text{Where } T_i = \sum_{j=1}^{12} (t_{ij}^3 - t_{ij})$$

$t_{ij} = \text{no. of times } i^{\text{th}} \text{ rank is repeated.}$

Q. A farmer wishes to determine whether there is a difference in yields between two different varieties of wheat I and II. The following data shows the production of wheat per unit area using the two varieties. Can the farmer conclude at significance level 0.01 that a difference exists?

Wheat I : 15.9 15.3 16.4 14.9 15.3 16.0 14.6 15.3 14.7 16.6 16.0

Wheat II : 16.3 16.8 17.1 16.9 18.0 15.6 18.1 17.2 15.4

$n_1 = 11$ and $n_2 = 9$ i.e. Large Sample Size (Z-test)

H_0 : There is no significant difference in the yields of wheat I and II or the difference does not exist.

H_1 : There is significant diff. " " " " or the difference exists.

Calculation:

Wheat I	Wheat II	Combined Rank	Naw
15.9	16.3	9	12.5
15.3	16.8	5	15
16.4	17.1	12.5	17
14.9	16.9	3	16
15.3	16.0	5	11
16.0	15.6	10.5	8
14.6	18.1	2	20
15.3	17.2	5	18
14.5	15.4	1	7
16.6		14	
16.0		10.5	

$$R_1 = 77.5 \quad R_2 = 132.5$$

$$t_1 = 3, \quad t_2 = 2, \quad t_3 = 2$$

$$[U_1 + U_2 = n_1 n_2]$$

$$\Rightarrow 87.5 + 11.5 = 11(9)$$

Now,

$$E(U_0) = \frac{n_1 n_2}{2} = \frac{99}{2} = 49.5$$

$$S.d(U_0) = \sqrt{\frac{n_1 n_2}{n(n-1)} \left\{ \frac{n^3 - n}{12} - \sum T_i^2 \right\}}$$

Where, $E\bar{U}_0 = \frac{\sum (U_i^3 - \bar{U}_i)}{12} = \frac{3^3 - 3}{12} + \frac{2^3 - 2}{12} + \frac{2^3 - 2}{12} = 3.$

$$\therefore Sd(U_0) = \sigma_{U_0} = \sqrt{\frac{11(9)}{20(20-1)} \left\{ \frac{20^3 - 20}{12} - 3 \right\}} \\ = 13.13.$$

Test Statistic : Under H_0 .

$$Z = \frac{U_0 - E(U_0)}{Sd(U_0)} \quad [U_0 = \min \{U_1, U_2\}] \\ = \frac{11.5 - 9.5}{13.13} \\ = -2.89 \\ \therefore |Z_{cal}| = 2.89$$

Calculated > tabulated.

Critical Value.

We have, $\alpha = 0.01$

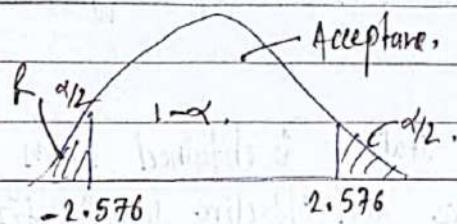
Then we reject H_0 otherwise do not reject H_0 .

$$\therefore Z_{tab} = Z_{0.01} (TIT) \\ = 2.576$$

Decision:

Since $|Z_{cal}| > Z_{tab}$ so we reject H_0 .

Hence, there is significant difference in the yields of wheat I and II.



Paired Test (for Dependent Samples)

Wilcoxon Matched Pairs Signed Rank Test. [small sample size $n \leq 25$]

This test is applied to test the hypothesis concerning the difference between two treatments used on the two random samples.
It can also be applied to test the hypothesis concerning the effectiveness of some activity (Treatment).

Procedure:

$H_0 : Md_1 = Md_2$ i.e. no. sig. diff before after | Md_1 = Median score before Treatment.

$H_1 : Md_1 \neq Md_2$ i.e. sig. diff before after | Md_2 = Median score after Treatment

or, $H_1 : Md_1 > Md_2$ (right) | DTT.

or $H_1 : Md_1 < Md_2$ (left)

Test Statistic: Under H_0 .

$$T = \text{minimum } \{ s(+), s(-) \}$$

$s(+)$ = sum of the ranks with '+' sign

$s(-)$ = sum of the ranks with '-' sign.

Ranks are given to the absolute value of difference i.e. $|d|$

Critical Value.

The critical value of the test statistic is obtained from Wilcoxon table at α level of significance and effective sample size $n_e(n)$

$$\text{i.e. } T_{\text{tab}} = T_{\alpha, n_e}$$

$$n_e = n - \text{no. of } d_i = 0$$

Decision

If $T \leq T_{\alpha, n_e}$ then we reject H_0 . Otherwise do not reject H_0 .

Seven prospective graduate student took a test twice with the following scores:

First Attempt : 4.70 5.30 6.10 4.90 6.00 5.90 5.80

Second Attempt : 5.10 5.50 6.80 4.90 5.85 6.20 5.18

Test whether there is significant difference between scores in first attempt and second attempt using Wilcoxon Matched Pair Signed rank test.

$n=7$

H_0 : There is no significant difference between the scores of first & second attempts.

H_1 : There is significant " " " " " (TTT)

Test statistic : Under H_0 .

$$T = \text{minimum of } \{ S(+), S(-) \}$$

$$\therefore T = 3$$

Calculation.

First attempt (x)	Second attempt (y)	$d = x - y$	Rank of d	Rank with sign
4.70	5.10	-40	6	6
5.30	5.50	-20	4	4
6.10	6.80	10	1	1
4.90	4.90	0		
6.00	5.85	15	2	2
5.90	6.20	30	5	5
5.80	5.18	-18	3	3

$S(+)=3$ $S(-)=25$

Critical value

We have, $\alpha = 0.05$ and $n = n_e = 7$,

$$\therefore T_{7, 0.05} (\text{TTT}) = 2.$$

Decision

Since $T(=3) > T_{7, 0.05} (=2)$, we do not reject H_0 . Hence there is no significant difference between the scores of first and second attempts.

ignore '0'

Date _____
Page _____

Q. To evaluate a speech & reading course, a group of 10 subjects were asked to read two comparable articles one before the course and one after the course. Their scores on reading test are as follows.

Before course (X) : 57 40 64 70 90 59 76 98 70 83

After course (Y) : 60 90 62 70 95 58 80 99 75 94

Test whether the course is beneficial using the Wilcoxon Matched Pairs Signed rank test at 5% level of significance.

H_0 : $M_{dX} = M_{dY}$ i.e. the course is not beneficial

H_1 : $M_{dX} < M_{dY}$ i.e. " " " beneficial (OTT)

Test Statistic : Under H_0 ,

$$T = \text{minimum } \{ S(+), S(-) \}$$

Calculation:

X	Y	d	Rank of d	Rank with + sign	- sign
57	60	-3	4		4
40	90	-50	9		9
64	62	2	3		3
70	70	0	—		
90	95	-5	6.5		6.5
59	58	1	1.5		1.5
76	80	-4	5		5
98	99	-1	1.5		1.5
70	95	-25	6.5		6.5
83	94	-11	8		8

$$S(+)=4.5 \quad S(-)=40.5$$

Critical Value:

We have, $\alpha = 0.05$ & $n_e = n-1 = 10-1 = 9 \quad \therefore T_{9, 0.05} (\text{OTT}) = 8$

Decision: Since, $T(=4.5) < T_{9, 0.05} (\text{OTT}) = 8$. Hence the course is beneficial.
We reject H_0 .

Kolmogorov-Smirnov Test (K-S Test) - One Sample.

It is used to test the difference between observed and expected frequencies.

$H_0: F(x) = F(y)$ i.e. There is no significant difference between observed and expected frequencies or frequencies are uniformly distributed.

$H_1: F(x) \neq F(y)$ i.e. There is significant " " " are not (TTT)

or $H_1: F(x) > F(y)$ by OTT

or $H_1: F(x) < F(y)$

Test statistic : Under H_0 ,

$$D_o = \text{maximum } f |f_e - f_o|$$

where,

$$f_e = \text{Relative Cumulative frequency} = \frac{c_{fe}}{n}$$

$$f_o = \frac{c_{fo}}{n}$$

$$c_{fe} = \text{Cumulative frequency}$$

$$o = \text{Observed frequency}$$

$$e = \text{expected frequency.}$$

Critical value or tabulated value :

The critical value of the test statistic is obtained from Kolmogorov-Smirnov test table at α level of significance and n no. of observations for one tailed and two tailed tests.

Decision:

If $D_o \geq D_{n,\alpha}$ then we reject H_0 otherwise do not reject H_0 .

Q. The number of laptop in 10 different departments are given below. Test whether the laptop are uniformly distributed over the entire office use KS test.

Department No.: 1 2 3 4 5 6 7 8 9 10

No. of Laptop: 8 10 9 12 15 7 5 12 13 9

H_0 : The laptops are uniformly distributed over the entire office.
 H_1 : The laptops are not " "

Test-statistic : Under H_0 :

$$D_0 = \text{maximum } |f_e - f_o| : \text{expected freq} n = NP$$

Calculation :

Dept.	Observed freq. (f_e)	Expected freq. (f_o)	CFO	$F_o = \frac{f_o}{n}$	Cfe	$f_e = \frac{f_e}{n}$	$ f_e - f_o $
1	8	10	8	8/100	10	11/100	1/100
2	10	10	18	18/100	20	20/100	2/100
3	9	10	27	27/100	30	30/100	3/100
4	12	10	39	39/100	40	40/100	1/100
5	15	10	54	54/100	50	50/100	4/100
6	7	10	61	61/100	60	60/100	1/100
7	5	10	66	66/100	90	90/100	4/100
8	12	10	78	78/100	80	80/100	2/100
9	13	10	91	91/100	90	91/100	1/100
10	9	10	100	100/100	100	100/100	0

$$\sum f_e = 100$$

$$\therefore D_0 = 4/100 = 0.04.$$

Critical value,

we have, $\alpha = 0.05$, if $n=100$.

$$\therefore D_{n,\alpha} = D_{100, 0.05} (\text{TTT}) = \frac{1.22}{\sqrt{100}} = 0.122.$$

Decision:

Since $D_0 < D_{n,\alpha}$, so we do not reject H_0 .

Hence, laptops are uniformly distributed over the entire offices.

Q. In a certain computer hardware manufacturing industry six different types of machines are working to cut pieces of wires. The number of wires of unequal length recorded in a day is as follows:

Machine	1	2	3	4	5	6
---------	---	---	---	---	---	---

No. of wires	2	0	4	8	5	11
--------------	---	---	---	---	---	----

Do these data provide sufficient evidence that the six machines equally cut the wires of unequal length? Apply Kolmogorov-Smirnov test at 5% level of significance.

So,

H_0 : The data provide sufficient evidence that the six machines equally cut the wires of unequal length.

H_1 : The data do not provide sufficient evidence (TTT)

Test-Statistic: Under H_0 , $\rightarrow e = \sum_{i=1}^n \frac{e_i}{n}$ where $e_i = \frac{|O_i - E_i|}{\sqrt{E_i}}$

Machine Observed(O_i) Expected(E_i) c_{f0} $F_0 = c_{f0}/n$ c_{fe} $F_e = c_{fe}/n$ $|F_e - F_0|$

1	2	5	2	$2/30$	9	$5/30$	$3/30$
---	---	---	---	--------	---	--------	--------

2	0	5	2	$2/30$	10	$10/30$	$8/30$
---	---	---	---	--------	----	---------	--------

3	4	5	6	$6/30$	15	$15/30$	$9/30$
---	---	---	---	--------	----	---------	--------

4	8	5	14	$14/30$	120	$20/30$	$6/30$
---	---	---	----	---------	-----	---------	--------

5	5	5	19	$19/30$	25	$25/30$	$0/30$
---	---	---	----	---------	----	---------	--------

6	11	5	30	$80/30$	30	$30/30$	0
---	----	---	----	---------	----	---------	---

$$\therefore D_0 = 9/30 = 0.3$$

Critical Value: or Tabulated Value,

$$D_{n,\alpha} = D_{30, 0.05} (\text{TTT}) = 0.242$$

Decision: Since, $D_0 > D_{n,\alpha}$ so we reject H_0 .

Hence, the data do not provide sufficient evidence that the six machines...

Kolmogorov-Smirnov Test for two sample.

It is non-parametric test used to test whether two independent samples are from same population or not.

Procedure:

$H_0: f(x) = f(y)$ i.e. the two independent samples are from same population or same distribution.

$H_1: F(x) \neq F(y)$ i.e. the two independent samples are not from ... (TTT).

or $H_1: F(x) > F(y) \quad ? \text{ OTT}$

or $H_1: F(x) < F(y)$

Small Sample Size

(i) $n_1 = n_2 \leq 40$

(ii) if $n_1 \neq n_2 \leq 20$

Test-statistic: Under H_0

$$D_0 = \text{maximum } \{ |F(x_i) - F(y_j)| \}$$

where, $f(x) = \frac{f_x}{n_1}$ and, $f(y) = \frac{f_y}{n_2}$

Critical Value:

$$D_{tab} = D_{(n_1, n_2), \alpha}$$

Decision:

If $D_0 \geq D_{(n_1, n_2), \alpha}$ then we reject H_0 , otherwise do not reject H_0 .

Large Sample size

$$n_1 = n_2 > 40 \quad \text{and} \quad n_1 \neq n_2 > 20$$

Test-statistic

For two tailed test,

$$D_0 = \text{maximum } \{ |F(x) - f(y)| \}$$

For one tailed test,

$$\chi^2 = \frac{1}{D_0} \left(\frac{n_1 n_2}{n_1 + n_2} \right) \quad \text{OTT, } \chi^2 \Rightarrow df = 2$$

Critical Value

$$\chi^2_{tab} = \chi^2_{2,\alpha} \rightarrow df = 2$$

Decision

If calculated \geq tabulated then we reject H₀, otherwise do not reject H₀.

Q. Life in years of two types cells used in laptop are given below:

Cell X: 4 6 5 6 3 4 5 3 | 5

Cell Y: 2 5 4 3 4 2 4 3 | 5.

Test whether, life of two brand of cells are same or not? Use ks test at $\alpha = 0.05$.

Soln

H₀: The life of two brands of cells are same.

H₁: " " " " are not same (TTT).

$$D_0 = \text{maximum } \sum |f(x_i) - f(y_j)|$$

Calculation:

Life in years	Freq. of X	Freq. of Y	$f(x)$	$f(y)$	$f(x)$	$f(y)$	$ f(x) - f(y) $
2	0	2	0	2	0/9	2/9	2/9
3	2	2	2	4	2/9	4/9	2/9
4	2	3	4	7	4/9	7/9	3/9
5	3	2	7	9	7/9	9/9	2/9
6	2	0	9	9	9/9	9/9	0

$$n_1 = 9 \quad n_2 = 9.$$

$$\therefore D_0 = 3/9.$$

Critical Value: We have $n_1 = n_2 = 9$ and $\alpha = 0.05$

$$\therefore D_{9,0.05} (\text{TTT}) = 5/4.$$

Decision :

Since $D_0 < D_{0.05}$, so we do not reject H_0 .
 Hence the life of two brands of cells are same.

- Q. Given below represent monthly income distn of employees in a hardware company and software company.

Income (000 Rs)	No. of employees in Hw Company	No. of emp. in Sw Company
20-30	6	12
30-40	10	18
40-50	11	16
50-60	13	12
60-70	25	10
70-80	15	12
80-90	10	10
	$n_1 = 90$	$n_2 = 90$

Do the income distribution support that income of employees in hardware company is more than income of employees in software company? Use Wilks test.

$$Md_1 = Md_2$$

$$Md_1 \neq Md_2$$

H_0 : The income of employees of Hardware company is not more than that of Software company.

H_1 : The income " " " " " " is more " "

Test Statistics: Under H_0 ,

$$\chi^2 = 4D_0^2 \left(\frac{n_1 n_2}{n_1 + n_2} \right) \sim \chi^2_{(1,2)}$$

Munie (000 Rs)	H.W Comp. fx	S.W. Company fy	$f(x) = \frac{fx}{90}$	$f(y) = \frac{fy}{90}$	$ f(x) - f(y) $
20-30	6	12	6/90	12/90	
30-40	10	18	10/90	30/90	
40-50	11	16	11/90	16/90	
50-60	13	12	13/90	12/90	
60-70	25	10	25/90	10/90	
70-80	15	12	15/90	12/90	
80-90	10	10	10/90	10/90	
	90	90			

$$n_1 = 90 \text{ & } n_2 = 90$$

$$D_0 = 19/90$$

$$\chi^2 = 4 \left(\frac{19}{90} \right)^2 \left[\frac{90(90)}{90+90} \right]$$

$$\chi^2_{\text{cal}} = 8.022$$

Critical Value, we have $\alpha = 0.05$ & d.f. = 2.

$$\chi^2_{\text{tab}} = \chi^2_{2,0.05} = 5.991$$

Decision

Since, $\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$ so, we reject H_0 .

Kruskal-Wallis H Test $H \Rightarrow p\text{-value}$ (At least 3 sample)
 OR Kruskal's Wallis One way ANOVA Test.

If it is used to test whether the three or more populations differ in location (median) or not.

Notation:

$k = \text{no. of samples} (> 3)$

$n_1 = \text{size of first sample}$

$n_2 = \text{size of second sample}$

$n_k = \text{size of } k^{\text{th}} \text{ sample.}$

Procedure:

$H_0: M_{d1} = M_{d2} = \dots = M_{dk}$ i.e. there is no significant difference between k -populations medians or k -sample medians.

$H_1: M_{d1} \neq M_{d2} \neq \dots \neq M_{dk}$ i.e. there is significant ... or at least one of the median is different.

Small Sample Size i.e. $k=3$ and $n_i \leq 5$; $i=1, 2, \dots, k$

Test Statistic: Under H_0

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1)$$

$R_i = \text{Sum of the combined ranks given to the first sample.}$

$\bar{R}_k = \text{Sum of the combined ranks given to the } k^{\text{th}} \text{ sample.}$

$$n = n_1 + \dots + n_k$$

If tied occurs then corrected test statistic is obtained as follows.

$$\text{Corrected } H = \frac{H}{CF}$$

where, $CF = \text{correction factor}$

$$= 1 - \frac{\sum (f_i^3 - f_i)}{n^3 - n}$$

$f_i = \text{no. of times } i^{\text{th}} \text{ rank is repeated}$

Critical Value.

For different values of n_i , ($i=1, 2, \dots, k$), the probability P_0 associated with the value, an extreme as observed H_0 (Hact) is obtained from the K-W table i.e. $P_0 = P(H > H_{act})$.

Decision

If p-value (P_0) $\leq \alpha$ then we reject H_0 . Otherwise do not reject H_0 .

Large Sample Size i.e. $k \geq 3$ and $n_i > 5$

In this case, the sampling distribution of ' H ' can be approximated by a chi-square distribution with $(k-1)$ degrees of freedom.

$$X^2_{tab} = X^2_{(k-1), \alpha}$$

Decision:

If $H \geq X^2_{(k-1), \alpha}$ then we reject H_0 . Otherwise don't reject H_0 .

Q. For the following scores of 3 groups, apply Kruskal Wallis H test to test the hypothesis that the three groups are not significantly different.

Group

Scores.

A 96 128 83 61 101

B 82 124 132 135 104

C 115 149 166 147 -

Solution.

$$k=3, n_1=5, n_2=5 \text{ and } n_3=4 \therefore n=14$$

H_0 : Three groups are not significantly different.

H_1 : Three groups are significantly different.

Calculation: combined rank

	Scores.					R_i	R_i^2/n_i
A	96 (14)	128 (9)	83 (3)	61 (1)	101 (5)	$R_1 = 22$	$R_1^2/n_1 = 96.8$
B	82 (12)	124 (8)	132 (10)	135 (11)	104 (6)	$R_2 = 37$	$R_2^2/n_2 = 27.38$
C	115 (7)	149 (13)	166 (14)	147 (12)	-	$R_3 = 46$	$R_3^2/n_3 = 529$

$$\leq \frac{R_i^2}{n} = 89.406$$

Test Statistic : Under H_0 ,

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1) = \frac{12}{14(14+1)} (899.6) - 3(14+1)$$

$$\therefore H = 6.405.$$

Critical Value :

(look near to 6.405)

For, $H = 6.405$, $n_1 = 5$, $n_2 = 5$, $n_3 = 7$ and $k = 3$ we have

$$p\text{-value } (P_0) = 0.049.$$

Decision.

We have, $\alpha = 0.05$ and $p\text{-value} = 0.049$.

Since $-p\text{-value} < \alpha$ so we reject H_0 . Hence, there is significant difference between three groups.

- Q. The following data represent the operating times in hours for 3 types of scientific pocket calculators before a change is required.

	Calculator A	4.9	6.1	4.8	4.6	5.3		Calculator B	5.5	5.4	6.2	5.8	5.5	5.2	4.8		Calculator C	6.4	5.8	5.6	6.5	6.3	6.6		

Use Kruskal Wallis test, at the 0.01 level of significance test the hypothesis that the operating times for all three calculators are equal.

$$K = 3, n_1 = 5, n_2 = 7 \text{ and } n_3 = 6.$$

H_0 : The operating time of all three calculators are same.

H_1 : The operating time of all three calculators are different.

Calculations:

Combined Rank

Calculator A	4.9 (4)	6.1 (12)	4.8 (1)	4.6 (2)	5.3 (11)														
B	5.5 (8.5)	5.4 (7)	6.2 (13)	5.8 (11)	5.5 (8.5)	5.2 (5)	4.8 (3)												
C	6.4 (15)	5.8 (18)	5.6 (10)	6.5 (16)	6.5 (16)	6.6 (12)													

$$t_1 = 2, R_1 = 25, R_2 = 56 \text{ and } R_3 = 90.$$

Since 5.5 is repeated two times.

Correction factor,

$$C.F = 1 - \frac{\sum (t_i^3 - t_i)}{n^3 - n} = 1 - \frac{12^3 - 12}{18^3 - 18} = 0.999$$

Test Statistic : Under H_0

$$H_{\text{corrected}} = \frac{12 - \sum r_i^2 - 3(n+1)}{n(n+1) - 1 - n}$$

$$= \frac{12 - \sum r_i^2 - 3(18+1)}{18(18+1)} = \frac{12 - [25 + 56 + 90] - 3(18+1)}{18(18+1)} = \frac{12 - 171 - 57}{324} = \frac{-114}{324} = -0.35$$

Critical Value

We have, $k=3$ and $\alpha = 0.05$

$$\chi^2_{\text{tab}} = \chi^2_{(k-1), \alpha} = \chi^2_{2, 0.05} = 5.991$$

Decision : Since, $H_{\text{corrected}} > \chi^2_{\text{tab}}$ so we reject H_0 . Hence, the operating time of all three calculators are different.

Friedman Test or Friedman two way ANOVA test ($k \geq 3$)

The Friedman test is designed to test whether the population of $k (\geq 3)$ treatments are identical or not i.e. to test whether all the treatments have equal effect or not, or to test whether all treatments are equally effective or not.

Procedure :

H_0 : $Md_1 = Md_2 = \dots = Md_k$ i.e. there is no significant difference between k -population or k -sample medians.

H_1 : $Md_1 \neq Md_2 \neq \dots \neq Md_k$ i.e. there is significant or at least one median is different.

Test Statistic : Under H_0 ,

$$Fr = \frac{12}{n k(k+1)} \sum_{i=1}^k r_i^2 - 3n(k+1) ; n = \text{size of each sample},$$

[Note that, here ranks are given according to block or column wise]

If there is tie in the observation then the corrected value of the test-statistic is calculated as,

$$\text{Friedman} = F_r \cdot \frac{1}{k}, \quad (f = 1 - \frac{\sum (t_i^3 - t_i)}{n(k^3 - k)})$$

Where, t_i = no. of times i th rank is repeated.

Case-I :- Small Sample Size. P.e.

(i) $k=3$ and $2 \leq n \leq 9$.

(ii) $k=4$ and $2 \leq n \leq 5$.

Critical Value.

For given sample size n , no. of samples k and the calculated value of the test-statistic i.e. F_r , we obtain the p-value P_0 from the Friedman probability table.

Decision: If $P\text{-value } (P_0) > \alpha$ then we do not reject H_0 .

Otherwise reject H_0 .

Case-II: Except above two cases, we use large sample size case.

Here, the dist'n of the test-statistic follows χ^2 -distribution with $(n-1)$ d.f.
i.e. $F_r \sim \chi^2_{(k-1)}$.

Decision: If $F_r > \chi^2_{(k-1), \alpha}$ then we reject H_0 . Otherwise don't reject H_0 .

Q. The Scores of 3 matched groups under the six conditions are given below.

Group. Condition.

	I	II	III	IV	V	VI
A	9	5	2	5	6	7
B	6	4	3	4	6	5
C	5	1	3	3	6	5

Apply the Friedman two way ANOVA test to identify if there is significant difference in the variation between matched groups. Use 5% level of significance.

Soln, We have, $k=3$ and $n=6$

H_0 : There is no significant difference in the variation between three matched groups.

H_1 : There is significant " "

Calulations:

Conditions.

Group	I	II	III	IV	V	VI	R_i	R_i^2
A	9(3)	5(3)	2(1)	5(3)	6(2)	7(3)	15	225
B	6(2)	4(2)	9(1.5)	4(1)	6(2)	5(1.5)	12	144
C	5(1)	1(1)	3(2.5)	3(1)	8(2)	5(1.5)	9	81

$$t_1 = 2, t_2 = 3, t_3 = 2$$

Calulation:

$$cf = 1 - \frac{\sum (t_i^3 - t_i)}{n(k^3 - k)} = 1 - \frac{(2^3 - 2)(3^3 - 3)(2^3 - 2)}{6(3^3 - 3)} \approx 0.75.$$

Test-statistic: Under H_0 ,

$$\begin{aligned} F_{\text{current}} &= \frac{12}{nk(k+1)} \sum_{i=1}^k R_i^2 - 3n(k+1) / cf \\ &= \frac{12}{6(3)(4)} \left[\frac{225 + 144 + 81}{3(6)(4)} \right] = 4. \\ &\quad 0.75. \end{aligned}$$

Critical Value: We have $k=3$, $n=6$ and $F_{\text{crit}} = 4$. We hence

$$P\text{-value} (P_0) = 0.184.$$

Decision: we have $\alpha = 0.05$, $P_0 = 0.184$.

Since: $P\text{-value} > \alpha$ so we do not reject H_0 . There is no significant difference in the variation between three matched group.

- a) Three different advertising media TV, Radio and News paper are being compared to study their effectiveness in promoting sales of WaiWai noodles. Each advertising media is exposed for specified period of time and sales (per day) from 10 stores located at different areas are recorded.

Advertising Media.

Stores.

	A	B	C	D	E	F	G	H	I	J
TV	20	21	15	12	14	17	21	16	20	18
Radio	7	9	11	12	10	10	14	12	8	7
News Paper	8	6	11	12	9	6	8	10	8	6

Are three advertising media equally effective, use Friedman two way ANOVA test.

$$S_{ij}^2, k=3, n=10$$

H_0 : There is no significant difference....

H_1 : There is significant difference....

Adv. Media.

Stores.

	A	B	C	D	E	F	G	H	I	J	f_1	f_2
TV	20(3)	21(3)	15(3)	12(4)	14(3)	17(3)	21(3)	16(3)	20(3)	18(3)	21	84
Radio	7(1)	9(2)	11(1.5)	12(4)	10(2)	10(2)	14(2)	12(2)	8(1.5)	7(1)	18	324
News Paper	8(2)	6(1)	11(1.0)	12(2)	9(1)	6(1)	8(1)	10(1)	8(1.5)	6(1)	13	169

$$f_1 = 2$$

$$f_2 = 2$$

$$f_3 = 2$$

Calculation.

$$(F = 1 - \frac{\sum f_i^2 + \bar{t}}{n(k-1)} = 1 - \left[\frac{(2^2-2) + (12^2-3) + (10^2-2)}{10(3^2-3)} \right] = 0.825$$

 Test statistic: Under H_0

$$\text{Fractuted} = \left(\frac{12}{n(k-1)} \sum_{i=1}^k \frac{f_i^2 - 3(k+1)}{cf} \right) = 15.76$$

Critical Value:

We have $k=3$, $n=10$, Fractuted = 15.76.

$$(M-1) = (3-1) = 2 = \text{d.f.}, \chi^2_{0.01, 2} = 5.99$$

$$f_2 > \chi^2_{0.01, 2}$$

 Decision: Reject H_0 .

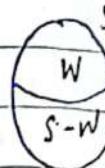
Hypothesis Testing.

Critical Region.

Let us define a sample space 'S' and we divide the whole sample space into two regions or two disjoint subsets say 'W' and 'S-W' such that.

$$W \cup S-W = S.$$

$$\text{and } W \cap S-W = \emptyset \text{ (null set)}$$



If the sample points fall into the region 'W', then we can reject the null hypothesis (H_0) and the region is called region of rejection. or if it is called critical region and if the sample points fall into the region ($S-W$), then we can accept the null hypothesis and the region is called region of acceptance.

Type - I error

$$\alpha = \text{Prob. (Type I error)}$$

$$= \text{Prob. (reject } H_0 \text{ when } H_0 \text{ is true)}$$

$$= \text{Prob. (x } \in W | H_0)$$

$$= \int_{x \in W} f(x | H_0) dx$$

Type - II error

$$\beta = \text{Prob. (Type II error)}$$

$$= \text{Prob. (accept } H_0 \text{ when } H_0 \text{ is false)}$$

$$= \text{Prob. (x } \in S-W | H_1)$$

$$= \int_{x \in S-W} f(x | H_1) dx.$$

Power of the test.

$$1-\beta = 1 - \text{Prob. (Type II error)}$$

$$= 1 - \text{Prob. (accept } H_0 \text{ when } H_0 \text{ is false)}$$

$$= \text{Prob. (reject } H_0 \text{ when } H_0 \text{ is false)}$$

$$= \text{Prob. (accept } H_1 \text{ when } H_1 \text{ is true)}$$

$$= \text{Prob. (correct decision)}$$

- Q. Tom a coin 5 times, if 2 or 3 heads, turn up then the coin is accepted as an unbiased coin. Find α , β and the power of the test when $H_0: p = 1/2$ against $H_1: p = 3/4$.

Soln,

Since, 1 coin is tossed five times so, $S = \{0, 1, 2, 3, 4, 5\}$, $n=5$.

Let $X = \text{no. of heads}$ then $X \sim B(n, p)$ so

$$P(X=x) = {}^n C_x p^x q^{n-x}; x=0, 1, 2, \dots, n; p+q=1$$

Here,

$$S-W = \{2, 3\}, W = \{0, 1, 4, 5\}.$$

Now,

$$P(X|H_0) = {}^5 C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^{5-2} = {}^5 C_2 \left(\frac{1}{2}\right)^5.$$

and,

$$P(X|H_1) = {}^5 C_2 \left(\frac{3}{4}\right)^2 \left(1-\frac{3}{4}\right)^{5-2} = {}^5 C_2 \left(\frac{3}{4}\right)^2 \left(\frac{1}{4}\right)^{5-2}.$$

$$\alpha = \text{Prb.}(X \in \{0, 1, 4, 5|H_0\})$$

$$\text{Here, } \alpha = 1 - \text{Prb.}(X \in \{2, 3|H_0\}) \quad \beta = \text{Prb.}(X \in S-W|H_1)$$

$$= P(X=2|H_0) + P(X=3|H_0) \quad = \text{Prb.}(X \in W|H_1)$$

$$= 1 - [{}^5 C_2 \left(\frac{1}{2}\right)^5 + {}^5 C_3 \left(\frac{1}{2}\right)^5] \quad = \text{Prb.}(X=2|H_1) + P(X=3|H_1)$$

$$= 1 - 0.625 \quad = {}^5 C_2 \left(\frac{3}{4}\right)^2 \left(\frac{1}{4}\right)^{5-2} + {}^5 C_3 \left(\frac{3}{4}\right)^3 \left(\frac{1}{4}\right)^{5-3}$$

$$\therefore \alpha = 0.375.$$

Power of test, $1-\beta =$

Q Let P be the Prb. of getting a head in a single toss of a coin. The coin is tossed 5 times and it is desired to test $H_0: P = \frac{1}{2}$ against $H_1: P = \frac{2}{3}$.

The null hypothesis H_0 is rejected if more than 3 heads are obtained. Find the Prb. of type I, type II error also, find power of test

(Given, $n=5$)

$X = \text{no. of heads}$.

$$S = \{0, 1, 2, 3, 4, 5\}$$

$H_0: P = \frac{1}{2}$ against $H_1: P = \frac{2}{3}$.

Here, $X \sim B(n, p)$

$$\text{So, } P(X=x) = {}^n C_x p^x q^{n-x}; x=0, 1, \dots, n.$$

$$W = \{4, 5\}.$$

$$S-W = \bar{W} = \{0, 1, 2, 3\}.$$

$$\text{Now, } \alpha = \text{Prb.}(X \in W|H_0)$$

$$= P(X=4|H_0) + P(X=5|H_0)$$

$$\int e^{\theta x} dx = \frac{e^{\theta x}}{\theta}$$

Date _____
Page _____

$$\therefore \alpha = {}^5C_4 \left(\frac{1}{3}\right)^4 \left(\frac{2}{3}\right)^5 \cdot 4 + {}^5C_5 \left(\frac{1}{3}\right)^5 \left(\frac{2}{3}\right)^5 \cdot 5$$

and,

$$\beta = \text{Prb. } (x \in S-W | H_1)$$

$$= 1 - P(x \in W | H_1)$$

$$= 1 - [P(x = 4 | H_1) + P(x = 5 | H_1)]$$

$$= 1 - [{}^5C_4 \left(\frac{2}{3}\right)^4 \left(\frac{1}{3}\right)^5 + {}^5C_5 \left(\frac{2}{3}\right)^5 \left(\frac{1}{3}\right)^5]$$

$$= 0.54.$$

$$\text{Power of the test} = 1 - \beta = 1 - 0.54 = 0.46.$$

Q. If $x \geq 1$ is the critical region for testing $H_0: \theta = 2$ against alternative hypothesis $H_1: \theta = 1$ on the basis of single observation from the popn,

$$f(x, \theta) = \theta e^{-\theta x}, \quad 0 \leq x < \infty.$$

obtain the values of type I and type II errors. Also find the power of the test.

Solu, Given, $W = \{x : x \geq 1\}$.

$$\text{and } S-W = \{x : x < 1\}.$$

Further, $H_0: \theta = 2$ vs. $H_1: \theta = 1$

$$\text{and } f(x, \theta) = \theta e^{-\theta x}, \quad x > 0$$

$$\alpha = f(x: x \geq 1 | H_0) = \int_1^\infty f(x | H_0) dx = \int_1^\infty [\theta e^{-\theta x}]_{\theta=2} dx.$$

$$= \int_1^\infty 2e^{-2x} dx = 2 \left[\frac{e^{-2x}}{-2} \right]_1^\infty = \left[e^{-2x} \right]_1^\infty = e^{-2} - e^{-\infty}$$

$$= e^{-2} = 0.13.$$

$$\therefore \alpha = 0.13.$$

$$\beta = \text{Prb. } (x \in S-W | H_1)$$

$$= f(x: x < 1 | \theta = 1)$$

$$= \int_0^1 f(x | \theta = 1) dx$$

$$= \int_0^1 [\theta e^{-\theta x}]_{\theta=1} dx.$$

$$= \int_0^1 e^{-x} dx$$

$$= \left[\frac{e^{-x}}{-1} \right]_0^1 = [e^{-x}]_0^1 = e^0 - e^{-1}$$

$$\therefore \beta = 1 - e^{-1} = 0.632$$

$$\therefore \text{power of test} = 1 - \beta = 1 - 0.632 = 0.368.$$

Q. Given the frequency function, $f(x, \theta) = \begin{cases} 1/\theta & ; 0 \leq x \leq \theta \\ 0 & ; \text{otherwise} \end{cases}$

and that you are testing the null hypothesis $H_0: \theta=1$ against $H_1: \theta=2$, by means of a single observed value of x . What would be the sizes of the type I, type II error, if you choose the interval (i). $0.5 \leq x$ (ii). $1 \leq x \leq 1.5$ as the critical regions? Also obtain the power of test.

(Given,

$$f(x, \theta) = 1/\theta ; 0 \leq x \leq \theta$$

$$H_0: \theta=1 \quad v.s. \quad H_1: \theta=2$$

$$\text{i}. \quad W = \{x : x > 0.5\}$$

$$\text{and } S-W = \{x : x < 0.5\}$$

$$\text{Now, } \alpha = f(x \in W | H_0)$$

$$= \int_{0.5}^{\theta} \left[\frac{1}{\theta} \right] dx = \int_{0.5}^{1} dx = [x]_{0.5}^1 = 0.5.$$

$$\beta = f(x \in S-W | H_1)$$

$$= \int_{0}^{0.5} \left[\frac{1}{2} \right] dx = \frac{1}{2} [x]_{0}^{0.5} = \frac{1}{2} (0.5) = \frac{1}{4} = 0.25.$$

$$\therefore \beta = 0.25$$

$$\therefore \text{Power of the test } (1-\beta) = 0.75$$

$$\text{ii}. \quad W = \{x : 1 \leq x < 1.5\} \text{ and } S-W = \{x : x < 1 \text{ or } x > 1.5\}.$$

$$\alpha = f(x \in W | H_0)$$

$$= \int_{1}^{1.5} \left[\frac{1}{\theta} \right] dx$$

here, the prob. does not exist,

$$\text{Since, } 0 \leq x \leq \theta \text{ and } \theta=1 \quad \therefore \alpha = 0$$

$$\beta = f(x \in S-W | H_1)$$

$$= 1 - \int_{x \in W} f(x | H_1) dx = 1 - \int_{1}^{1.5} \left[\frac{1}{2} \right] dx = 1 - \frac{1}{2} [x]_{1}^{1.5} = 1 - \frac{1}{2} [1.5 - 1] = 0.75.$$

$$\text{Power of test} = 1 - \beta = 0.25$$

$$X \sim N(\mu, \sigma^2)$$

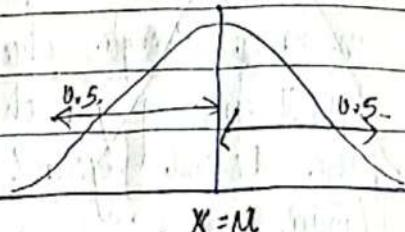
X follows Normal Distribution with mean μ and variance σ^2 .

$$Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$$

$$\downarrow \sigma$$

Standard Normal Distribution.

$$E[Z] = 0, \text{Var}[Z] = 1$$



$$\# \mu = \ell, \sigma = 2$$

$$\textcircled{1} \quad P(X > 10) = P\left(Z > \frac{10-\ell}{\sigma}\right)$$

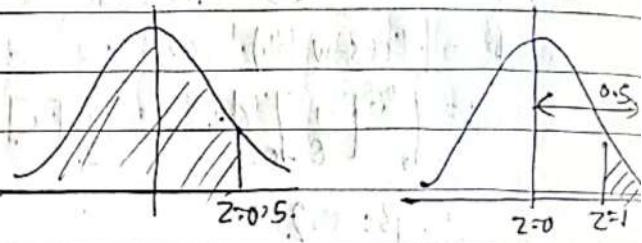
$$= P(Z > 1)$$

$$= 0.5 - P(0 < Z \leq 1)$$

$$\textcircled{11} \quad P(X < 9) = P\left(Z < \frac{9-\ell}{\sigma}\right)$$

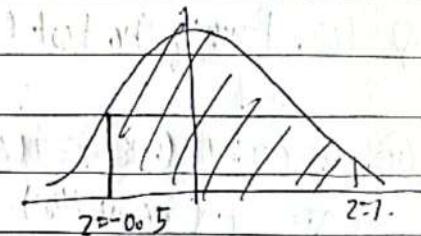
$$= P(Z < 0.5)$$

$$= 0.5 + P(0 < Z < 0.5)$$



$$\textcircled{11} \quad P(7 < X < 10) = P\left(\frac{7-\ell}{\sigma} < Z < \frac{10-\ell}{\sigma}\right)$$

$$= P(-0.5 < Z < 1)$$



- Q. Suppose it is desired to test the hypothesis $H_0: \mu = 35$ against the alternative $H_1: \mu \neq 35$ on the basis of a random sample of size 16 from a Normal population $N(\mu, 1)$. The decision rule is reject H_0 , if the sample mean $\bar{X} < 34.65$ or $\bar{X} > 35.35$

(1) Find the Prob. of type I error

(2) Find the Prob. of type II error (i) when $\mu = 36$ (ii) when $\mu = 36.1$

(3) Also find the power of the test.

Soln

$H_0: \mu = 35$ vs $H_1: \mu \neq 35$. $n=16$ & $\bar{X} \sim N(\mu, \sigma^2)$.

$W = \{\bar{X} : \bar{X} < 34.65 \text{ or } \bar{X} > 35.35\}$.

$S-W = \{\bar{X} : 34.65 < \bar{X} < 35.35\}$.

$\therefore \bar{X} \sim N(\mu, \sigma^2/n)$

Note: $\bar{X} \sim N(\mu, \sigma^2/n)$

$X \sim N(\mu, \sigma^2)$

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

 $n = \text{sample size}$

$\bar{X} \sim N(\mu, \sigma^2/n)$

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

$6/\sqrt{n}$

① Probability of type I error.

$\alpha = \text{probability } (\bar{X} \in W | H_0)$

$= P_{H_0}(\bar{X} \notin S-W : 34.65 < \bar{X} < 35.35)$

$= P_{H_0}(\bar{X} < 34.65) + P_{H_0}(\bar{X} > 35.35)$

$= P(z < 34.65 - 35) + P(z > 35.35 - 35)$

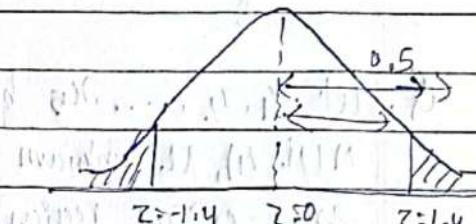
$= P(z < -1.4) + P(z > 1.4)$

(By Symmetric Law)

$= 2[0.5 - P(0 < z < 1.4)]$

$= 2[0.5 - 0.4192]$

$= 0.1616$

② $N=36$.Size or probability of type II error $\beta = P(\bar{X} \in S-W | H_1)$

$\beta = P_{H_1}(\bar{X} \notin S-W : 34.65 < \bar{X} < 35.35)$

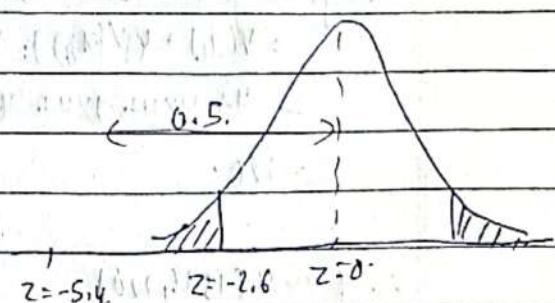
$= P(34.65 - 36 < z < 35.35 - 36)$

$= P(-5.4 < z < -2.6)$

$= 0.5 - P(-2.6 < z < 0)$

$= 0.5 - 0.4953$

$= 0.047$



Now,

$\text{Power of test} = 1 - \beta = 1 - 0.047 = 0.953$

$$\bar{X} \sim N(\mu, \sigma^2) \Rightarrow Z = \frac{\bar{X}-\mu}{\sigma}$$

$$\bar{X} \sim N(\mu, \sigma^2/n) \Rightarrow Z = \frac{\bar{X}-\mu}{\sigma/\sqrt{n}}$$

Q1. $\mu = 36.1$

$$\begin{aligned}\beta &= P(\bar{X} \in 34.65 < \bar{X} \leq 35.35 | H_1) \\ &= P(34.65 < \bar{X} \leq 35.35 | \mu = 36.1) \\ &= P\left(\frac{34.65 - 36.1}{\sigma/\sqrt{16}} < Z < \frac{35.35 - 36.1}{\sigma/\sqrt{16}}\right)\end{aligned}$$

$$= P(-5.8 < Z \leq -3)$$

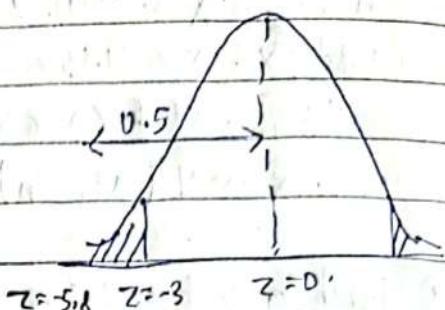
$$= 0.5 - P(-3 < Z \leq 0)$$

$$= 0.5 - 0.4987$$

$$= 0.0013$$

Now,

$$\text{Power of test} = 1 - \beta = 1 - 0.0013 = 0.9987$$



Q2. Let x_1, x_2, \dots, x_5 be a random sample of size 5 drawn from $N(\mu, 4)$, μ unknown. If you choose $x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5 \geq 0$ as a critical region for testing $H_0: \mu = 1$ against $H_1: \mu = 1$, find the size of the critical region and power of the test.

Soln,

$$n=5, X \sim N(\mu, 4)$$

$$\text{let } y = x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5$$

$$E(y) = E(x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5)$$

$$= E(x_1) + 2E(x_2) + 3E(x_3) + 4E(x_4) + 5E(x_5)$$

$$= \mu + 2\mu + 3\mu + 4\mu + 5\mu$$

$$= 15\mu$$

$$V(y) = V[x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5]$$

$$= V(x_1) + 4V(x_2) + 9V(x_3) + 16V(x_4) + 25V(x_5) \quad \because V(x) = \sigma^2 V(x) \text{ &}$$

$$= 4 + 4\mu + 9\mu + 16\mu + 25\mu$$

$$= 220$$

$$(\because x \sim N(\mu, 4))$$

x_i 's are independent

$$\therefore y \sim N(15\mu, 220)$$

$$W = \{y : y > 0\}$$

Critical Region

$$S-W = \{y : y \leq 0\}$$

Acceptable Region

Now,

$$\begin{aligned}
 \alpha &= P(y + w/H_0) && (\text{Type I error}) \\
 &= P(y > 0/\mu = -1) \\
 &= P\left(z \geq \frac{0 - (-1)}{\sqrt{220}}\right) && z = \frac{\bar{x} - \mu}{\sigma} \\
 &\Rightarrow P(z > 1.01) \\
 &= (0.5 - 0.3438) \\
 &= 0.1562
 \end{aligned}$$

$$\begin{aligned}
 \beta &= \text{prob. } (y + s - w/H_1) \\
 &= P(y \leq 0/\mu = 1) \\
 &= P\left(z \leq \frac{0 - (1)}{\sqrt{220}}\right) \\
 &= P(z < -1.01) \\
 &= 0.1562
 \end{aligned}$$

Cont..

Types of Hypothesis Testing.

- ① Simple Hypothesis.
- ② Composite Hypothesis

(i) Simple Hypothesis.

The hypothesis which completely specifies the distribution of random variable is known as simple hypothesis. For example, $X \sim N(\mu, \sigma^2)$ and the hypothesis to be tested as, $H_0: \mu = \mu_0, (\sigma^2 = \sigma_0^2)$ then under this hypothesis the form of the distribution is completely specified.

(ii) Composite Hypothesis.

A hypothesis which does not completely specifies the distributions then it is said to be composite hypothesis. For example $X \sim N(\mu, \sigma^2)$ and the hypothesis to be tested as $H_0: \mu = \mu_0$, then under this hypothesis σ^2 is not specified.

Likelihood function.

Let x_1, x_2, \dots, x_n be a random sample of size 'n' drawn from a population with pdf $f(x, \theta)$, then the joint pdf or pmf of x_1, x_2, \dots, x_n is said to be likelihood function of $x = x_1, x_2, \dots, x_n$ and it is denoted by L or $L(\theta)$ or $L(x, \theta)$. It is given by,

$$L = \prod_{i=1}^n f(x_i, \theta)$$

pdf = Prob. density fun.

$$= f(x_1, x_2, \dots, x_n, \theta)$$

pmf = Prob. mass fun.

$$= f(x_1, \theta) \cdot f(x_2, \theta) \cdots f(x_n, \theta)$$

iid = identically independent of

State the likelihood function of $x_i \sim P(\lambda)$; $i = 1, 2, \dots, n$

Soln, Here, $x_i \sim P(\lambda)$

$$\Rightarrow P(x_i) = e^{-\lambda} \lambda^{x_i}$$

$$x_i!$$

$$\therefore P(x_1) = \frac{e^{-\lambda} \lambda^{x_1}}{x_1!}, P(x_2) = \frac{e^{-\lambda} \lambda^{x_2}}{x_2!}, \dots, P(x_n) = \frac{e^{-\lambda} \lambda^{x_n}}{x_n!}$$

Now, the likelihood function of $X = x_1, x_2, \dots, x_n$ is given by,

$$L = \prod_{i=1}^n f(x_i; \lambda)$$

$$= P(x_1; \lambda), P(x_2; \lambda) \dots P(x_n; \lambda)$$

$$= \frac{e^{-\lambda} \lambda^{-x_1}}{x_1!}, \frac{e^{-\lambda} \lambda^{-x_2}}{x_2!}, \dots, \frac{e^{-\lambda} \lambda^{-x_n}}{x_n!}$$

$$= e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}$$

$$\prod_{j=1}^n (x_j!)$$

State the likelihood function of $x_i \sim N(\mu, \sigma^2)$, $i=1, 2, \dots, n$.

Soln: Here,

$x_i \sim N(\mu, \sigma^2)$ so

$$f(x_i) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \quad \text{if } i=1, 2, \dots, n.$$

If $x \sim N(\mu, \sigma^2)$

then

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

Now, the likelihood function of $x = x_1, x_2, \dots, x_n$ is given by

$$L = \prod_{i=1}^n f(x_i; \mu, \sigma^2)$$

$$= f(x_1; \mu, \sigma^2) f(x_2; \mu, \sigma^2) \dots f(x_n; \mu, \sigma^2)$$

$$= \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_1 - \mu}{\sigma} \right)^2} \dots \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_n - \mu}{\sigma} \right)^2}$$

$$\therefore L = \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n e^{-\frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2}$$

Best Critical Region and Best test.

A critical region of level size α is said to be a best critical region if it has minimum β or maximum power $(1-\beta)$ among all critical regions.

A test which minimizes β or maximizes the power $(1-\beta)$ for a desired level of significance α , is known as best test.

Most powerful critical region (MPCR) and most powerful test (MP test)

let x be an arbitrary sample point in a sample space's Then a critical region w is said to be a most powerful critical region (MPCR) or (BCR) of size α for testing a simple null hypothesis $H_0: \Theta = \Theta_0$ against a simple alternative hypothesis, $H_1: \Theta = \Theta_1$, if

$$\alpha = P(x \in w | H_0) \quad \text{--- (1)}$$

$$\text{and, } P(x \in w | H_1) \geq P(x \in w_i | H_1)$$

for every other critical region w_i satisfying eqn (1) and a test corresponding to this MPCR of size α is called most powerful (MP) test of size α .

Neyman Pearson Lemma (NP-Lemma)

NP Lemma provides the general method of finding a most powerful test of a simple null hypothesis against the simple alternative hypothesis.

Statement: Let $L(x|H_0)$ and $L(x|H_1)$ be the likelihood function under H_0 and H_1 respectively, the BCR to test the simple null hypothesis $H_0: \Theta = \Theta_0$ against the simple alternative hypothesis $H_1: \Theta = \Theta_1$ is given by,

$$\left(\frac{L(x|H_1)}{L(x|H_0)} \right) \geq k \quad \forall x \in w$$

$$\left(\frac{L(x|H_1)}{L(x|H_0)} \right) < k \quad \forall x \in S-w$$

Such that $\int_{x \in w} L(x|H_0) dx \leq \alpha$ fixed.

and k = non-negative constant which is derived from fixed α .

S = Sample space, w = critical region.

Then the test is most powerful test (MP-test) and critical region is BCR.

Proof:

Since, α is the prob. of type I error.

$$\alpha = \text{Prob(Type I error)}$$

$$= \text{Prob}(x \in w | H_0)$$

$$= \int_{x \in w} L(x|H_0) dx \quad \text{--- (1)}$$

and, $\beta = \text{Prob(Type II error)}$

$$= \text{Prob}(x \in S-w | H_1)$$

$$= \int_{x \in S-w} L(x|H_1) dx \quad \text{--- (II)}$$

let w^* be any other critical region with prob. type I error α_1 (which is less than or equal to α) and prob. of type II error β_1 .

$$\therefore \alpha_1 = P_{H_0}(x \in w^* | H_0)$$

$$\Rightarrow \alpha_1 = \int_{x \in w^*} L(x | H_0) dx$$

Similarly from B_1 ,

$$\beta_1 = \int_{x \in w^*} L(x | H_1) dx$$

Since, $\alpha \geq \alpha_1$. In fact, the latter is of size α .

$$\int_{x \in w} L(x | H_0) dx \geq \int_{x \in w^*} L(x | H_0) dx \quad \text{--- (1)}$$

further,

$$1 - \beta = \int_{x \in w} L(x | H_1) dx \quad \text{and} \quad 1 - \beta_1 = \int_{x \in w^*} L(x | H_1) dx$$

Now,

$$\beta_1 - \beta = \int_{x \in w} L(x | H_1) dx - \int_{x \in w^*} L(x | H_1) dx$$