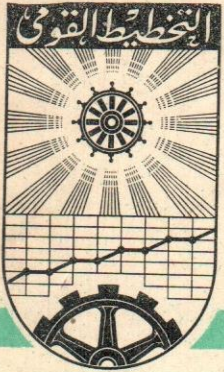


الجمهورية العربية المتحدة



مَعْمَدُ التَّخْطِيطِ الْقَوْمِىِّ نَسَبَةٌ فَقَطْ

Memo. No. 590.

Note On Statistical Methods

By

Dr. Moharram W. Mahmoud

B.Sc., M.Sc., M.A., Ph.D.

Operations Research Centre.

10/8/1965.

القاهرة

٣ شارع محمد منظر - بالزمالك

<u>Chapter</u>		<u>Page</u>
(I)	<u>Moments</u> Moments around the origin (m'_r) - Moments around the mean (m_r) - Expressing m'_r in terms of m_r Measure of skewness - Measure of kurtosis Numerical example	1
(II)	<u>Fundamental Concepts of Probability</u> Introduction & definitions - Basic theorems - Conditional probability - Independence	5
(III)	<u>Random Variables (Discontinuous & Continuous)</u> Binomial probability distribution - Its relationship to the incomplete Beta function Fitting a binomial distribution - Normal distribution - Fitting a normal curve.	11
(IV)	<u>Expected Values & Characteristic function</u> Theorems on expected values - Linear functions of a random variable - Quadratic function. Properties of Characteristic functions - Characteristic functions as moment generating functions - Characteristic functions used for calculating cumulants. One-to-one transformation Transformation $y = x^2$ χ^2 - distribution	26
(V)	<u>Some Exact Distributions</u> Poisson - Negative binomial - Sequential binomial - Hypergeometric - Gamma - Beta	48
(VI)	<u>Bivariate Normal & Binomial Distributions</u> Bivariate normal: (contours - Marginal and conditional distributions - Moments) Bivariate binomial: (Characteristic function - approximation to the bivariate normal)	59
(VII)	<u>Sampling and sampling distributions</u> Basic definitions - Estimation - Testing statistical hypothesis - Sampling distribution of the mean \bar{x} - Sampling distribution of s^2 - Sampling distribution of t - Sampling distribution of F .	70

(VIII) Testing statistical hypothesis and confidence interval

Around the mean and the difference between two means for large and small samples - Around the proportion and the difference between two proportions. Around the variance

(IX) χ^2 and its Applications 101

Test of goodness of fit - Contingency tables - 2X2 tables

(X) Simple analysis of variance 112

Introduction - Estimates for σ^2 and measure of discrepancy - analysis of variance table - Expected values of sum of squares - Bartlett's test for homogeneity of variances

(XI) Simple linear regression 128

The regression line $y = a + b(x - \bar{x})$ - Testing the significance of "a" & "b" - Regression analysis table - Expected values of sum of squares due to regression and sum of squares of residuals.

(XII) Multiple Regression 140

Introduction - Case of two independent variates - Generalization to more than two independent variates.

This note has been prepared by the author for a course in Statistical methods which is especially designed for those of mathematical background as engineers and scientists.

It was given as lectures for

(i) the group of operations research of the 4th long term training course 1965.

&(ii) the group who attended the short term training period (Jan. - Feb. 1965) at the Operations Research Centre.

The author takes this opportunity to acknowledge the encouragement of Dr Salah Hamid; the general director of the Operations Research Centre.

He would also like to thank Mrs. Nadia Amer and Miss Ellen Zaki for the great care they took in writing this note.

Moharram W. Mahmoud

(1)

(1)

MOMENTS

(1-1) Moments around the origin:-

$$m'_r(x) = \frac{1}{n} \sum x_i^r \quad (\text{ungrouped data})$$

$$= \frac{1}{n} \sum f_i x_i^r \quad (\text{grouped data})$$

For grouped data x_i denote the mid points of the intervals and f_i the corresponding frequencies

(1-2) Moments around the mean:-

$$m_r(x) = \frac{1}{n} \sum (x_i - \bar{x})^r \quad (\text{ungrouped data})$$

$$= \frac{1}{n} \sum f_i (x_i - \bar{x})^r \quad (\text{grouped data})$$

Remarks (i) $m'_0(x) = m_0(x) = 1$

(ii) $m'_1(x) = \bar{x}$ = the mean

(iii) $m_1(x) = 0$

(iv) $m_2(x) = s^2$ = the variance.

(1.3) Expressing $m_r(x)$ in terms of $m'_r(x)$

$$m_r(x) = \sum_{t=0}^r (-1)^t \binom{r}{t} (m'_1)^t m'_{r-t}$$

In particular

$$m_2 = m'_2 - m_1'^2$$

$$m_3 = m'_3 - 3m'_2 m'_1 + 2m_1'^3$$

$$m_4 = m'_4 - 4m'_3 m'_1 + 6m'_2 m_1'^2 - 3m_1'^4$$

Remark In the case of frequency distributions with equidistant intervals we replace the mid-points x_i by d_i where

$$d_i = \frac{x_i - a}{l}$$

where "a" is an arbitrary origin

"l" is the length of the interval

Then

$$m_r(x) = l^r m_r(d)$$

(1.4) Measures of skewness

$$\beta_1 = \frac{m_3^2}{m_2^3}$$

$$\gamma_1 = \sqrt{\beta_1}$$

Measure of Kurtosis

$$\beta_2 = \frac{m_4}{m_2^2}$$

$$\gamma_2 = \beta_2 - 3$$

(1.5) Numerical Example (illustrative)

Intervals	f_i	d_i	d_i^2	d_i^3	d_i^4
10-20	10	-1	1	-1	1
20-30	18	0	0	0	0
30-40	14	1	1	1	1
40-50	8	2	4	8	16
	50	20	56	72	152
$m_r(d)$	$\sum f_i d_i$		$\sum f_i d_i^2$	$\sum f_i d_i^3$	$\sum f_i d_i^4$
	0.4		1.02	1.44	3.04

(4)

$$m_2(d) = 1.02 - (0.4)^2 = 0.86$$

$$m_3(d) = 1.44 - 3(1.02)(0.4) + 2(0.4)^3 = 0.344$$

$$m_4(d) = 3.04 - 4(1.44)(0.4) + 6(1.02)(0.4)^2 - 3(0.4)^4 = 1.6394$$

$$m_1(x) = 10 m_1'(d) + a = 10(0.4) + 25 = 29$$

$$m_2(x) = s^2 = 10^2 (0.86) = 86$$

$$m_3(x) = 10^3 (0.344) = 344$$

$$m_4(x) = 10^4 (1.6394) = 16394$$

$$\beta_1 = (344)^2 / (86)^3 = 0.1860$$

$$\gamma_1 = \sqrt{0.1860} = 0.43$$

$$\& \beta_2 = 16394 / 7396 = 2.22$$

$$\therefore \gamma_2 = \beta_2 - 3 = 2.22 - 3 = -0.78$$

Exercise given the following frequency table, calculate the mean, the variance and β_1 & β_2

Intervals	f	Intervals	f
55 - 75	3	155 - 175	209
75 - 95	21	175 - 195	81
95 - 115	78	195 - 215	21
115 - 135	182	215 - 235	5
135 - 155	305		905

(5)

(II) FUNDAMENTAL CONCEPTS
OF PROBABILITY

(2-1) The theory of probability is a branch of applied mathematics dealing with the effects of chance. If we throw a die upon a board we are certain that one of the six faces will turn up, but whether a particular face will show, depends on what we call chance. Also, if equal numbers of white and black balls are put in an urn and we draw one of them blindly, we are certain that its colour will be either white or black, but whether it will be black, that depends on chance.

The word event which we are going to use frequently is used to signify an observation satisfying some specified conditions.

Two events are said to be "equally likely" if after taking into consideration all relevant evidence one of them cannot be expected in preference to the other. e.g. in the case of the urn with equal number of white and black balls, if we draw a ball, it is equally likely to be either white or black.

In the field of statistical analysis there would seem to be two definitions:

(i) Mathematical theory of arrangements which is as old as gambling & playing cards. The probability (p) of an event is the ratio of the no. of ways in which the event may happen divided by the total no. of ways in which the event may or may not happen. This is under the condition that all the events are equally likely. e.g. in the case of an unbiased coin, the probability that the head appears uppermost is $p = \frac{1}{2}$. Also in throwing a die the probability a particular face will show is $p = \frac{1}{6}$.

(ii) The frequency theory: If in a series of n independent trials which are absolutely identical, the event E is found to occur in m trials, then the probability of E is $\frac{m}{n}$.

This gives us a way to estimate probabilities from experimental results in a simple way.

As n increases $\frac{m}{n}$ tends to p , i.e. $p = \lim_{n \rightarrow \infty} \frac{m}{n}$.

(2-2) Definition (1) Fundamental probability set (F.P.S.). is that set of individuals or units from which the probability is calculated.

In the case of die, the F.P.S. given by the mathematical theory of arrangements would be 6. If the die is biased in some way and it is necessary to estimate a probability from the observations, then the F.P.S. would be the total number of throws of a die.

Definition (2) Mutually exclusive: Two events E_1, E_2 are said to be mutually exclusive if no element of the P.P.S. may possess both E_1, E_2 . In other words the two events do not occur together.

Remarks (i) $\Pr \{E_1 + E_2\}$ means the probability of E_1 or E_2 .

(ii) $\Pr \{E_1 E_2\}$ means the probability of E_1 & E_2 .

(2.3) Basic theorems:-

In the following theorems we are going to assume that the fundamental probability set N , consists of

n_1	elements possessing	E_1
n_2	"	" E_2
n_{12}	"	" E_1 & E_2
n_0	"	" \bar{E}_1 & \bar{E}_2

(where \bar{E}_1 means the event E_1 does not occur)

In other words $N = n_1 + n_2 + n_{12} + n_0$

Theorem (1) If E_1, E_2 are mutually exclusive and the only possible events, then

$$\Pr \{ E_1 \} + \Pr \{ E_2 \} = 1$$

Proof:- Since the two events are mutually exclusive then $n_{12} = 0$. Also the two events are the only possible then $n_0 = 0$.

Therefore

$$n_1 + n_2 = N$$

$$\therefore \frac{n_1}{N} + \frac{n_2}{N} = 1$$

$$\therefore \Pr \{ E_1 \} + \Pr \{ E_2 \} = 1$$

Theorem (2) $\Pr \{ E_1 + E_2 \} = \Pr \{ E_1 \} + \Pr \{ E_2 \} - \Pr \{ E_1 E_2 \}$

$$\Pr \{ E_1 \} = \frac{n_1 + n_{12}}{N}, \quad \Pr \{ E_2 \} = \frac{n_2 + n_{12}}{N}$$

$$\Pr \{ E_1 E_2 \} = \frac{n_{12}}{N}$$

$$\begin{aligned} \therefore \Pr \{ E_1 \} + \Pr \{ E_2 \} - \Pr \{ E_1 E_2 \} &= \frac{n_1 + n_{12}}{N} + \frac{n_2 + n_{12}}{N} - \frac{n_{12}}{N} \\ &= \frac{n_1 + n_2 + n_{12}}{N} \\ &= \Pr \{ E_1 + E_2 \} \end{aligned}$$

Cor. If E_1, E_2 are mutually exclusive then

$$\Pr \{ E_1 + E_2 \} = \Pr \{ E_1 \} + \Pr \{ E_2 \}$$

In general for k mutually exclusive properties we have

$$\Pr \left\{ \sum_{i=1}^k E_i \right\} = \sum_{i=1}^k \Pr \{ E_i \}$$

Definition Conditional probability of an event E_2 given event E_1 is the probability of E_2 referred to the F.P.S. for E_1 and it is written $\Pr \{E_2 | E_1\}$

Theorem (3) $\Pr \{E_1 E_2\} = \Pr \{E_1\} \cdot \Pr \{E_2 | E_1\}$
 $= \Pr \{E_2\} \Pr \{E_1 | E_2\}$

Proof:- $\Pr \{E_1\} = \frac{n_1 + n_{12}}{N}$

$\Pr \{E_2\} = \frac{n_2 + n_{12}}{N}$

$\Pr \{E_2 | E_1\} = \frac{n_{12}}{n_1 + n_{12}}$

$\Pr \{E_1 | E_2\} = \frac{n_{12}}{n_2 + n_{12}}$

$\therefore \Pr \{E_1\} \cdot \Pr \{E_2 | E_1\} = \frac{n_1 + n_{12}}{N} \times \frac{n_{12}}{n_1 + n_{12}} = \frac{n_{12}}{N}$

$\Pr \{E_2\} \cdot \Pr \{E_1 | E_2\} = \frac{n_2 + n_{12}}{N} \times \frac{n_{12}}{n_2 + n_{12}} = \frac{n_{12}}{N}$

$\therefore \Pr \{E_1\} \cdot \Pr \{E_2 | E_1\} = \Pr \{E_2\} \Pr \{E_1 | E_2\}$
 $= \frac{n_{12}}{N} = \Pr \{E_1 E_2\}$

(2-4) Independence: E_1 is independent of E_2 if

$$\Pr \{E_1\} = \Pr \{E_1 | E_2\}.$$

Consequently this implies (theorem 3) that E_2 is independent of E_1 .

(i0)

Cor. If E_1, E_2 are mutually independent then

$$\Pr \{ E_1 E_2 \} = \Pr \{ E_1 \} \Pr \{ E_2 \}$$

In general

$$\Pr \left\{ \bigcap_{i=1}^k E_i \right\} = \prod_{i=1}^k \Pr \{ E_i \}$$

III Random Variables

Discontinuous random

Variables (Binomial distribution)

Continuous random variables (Normal distribution)

(3.1) Statistical data consist of values obtained by measurements of a certain character or characters of a number of individuals. The observed values are not usually all the same. As the number of observations is increased the proportion of observations less than a fixed number X say become more & more stable. From the viewpoint of the frequency theory of probability this defines the probability that the observed value of the character falls below the number X .

The theoretical concept which we use to correspond with "observed values" is the random variable. x is said to be a random variable, if for any number X , there is a probability $F(X)$ that x is less than or equal to X i.e.

$$\Pr\{x \leq X\} = F(X) \text{ exists.}$$

Replacing X by x in $F(X)$ we get what is called the cumulative distribution function of x .

The properties of $F(x)$ are

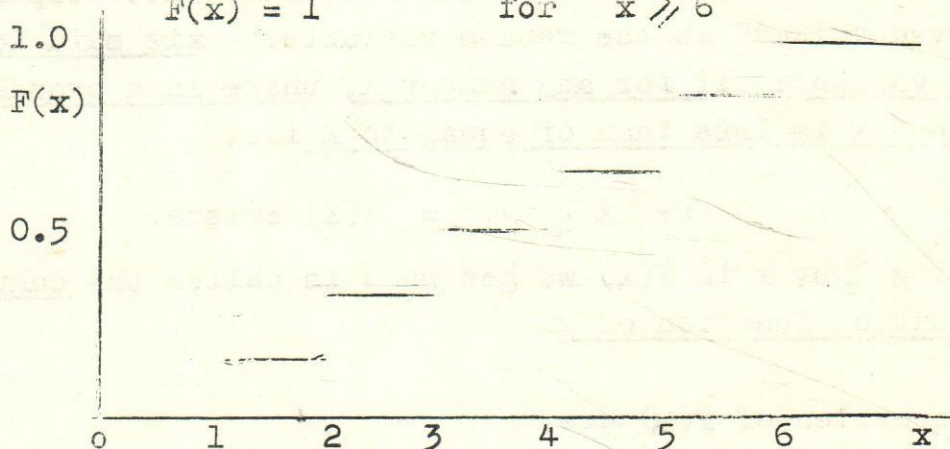
- (i) $0 \leq F(x) \leq 1$
- (ii) If $x_1 < x_2$ then $F(x_1) < F(x_2)$

i.e. $F(x)$ is non-decreasing function.

As an example, let us take an unbiased die.

If x denotes the number of dots on the upper-most face, then

$$\begin{aligned}
 F(x) &= 0 && \text{for } x < 1 \\
 F(x) &= \frac{1}{6} && \text{for } 1 \leq x < 2 \\
 F(x) &= \frac{1}{3} && \text{for } 2 \leq x < 3 \\
 F(x) &= \frac{1}{2} && \text{for } 3 \leq x < 4 \\
 F(x) &= \frac{2}{3} && \text{for } 4 \leq x < 5 \\
 F(x) &= \frac{5}{6} && \text{for } 5 \leq x < 6 \\
 F(x) &= 1 && \text{for } x \geq 6
 \end{aligned}$$



(3.2) Discontinuous random variables

Random variables that take isolated values as number of children in a family are called discontinuous. Their cumulative distribution functions are step-functions and the heights of the steps being equal to the corresponding probabilities. Suppose in general we have k definite values $x_1 < x_2 < x_3 \dots < x_k$ and the probabilities of x taking these values are $p_1, p_2, p_3, \dots, p_k$ respectively, i.e. $\Pr \{x = x_i\} = p_i$ $i=1, 2, \dots, k$

If $x_{i-1} \leq x < x_i$ then

$$F(x) = \Pr \{x \leq x_{i-1}\} = \sum_{j=1}^{i-1} p_j$$

Since the events $x = x_i$ are mutually exclusive then

$$\sum_{i=1}^k p_i = 1$$

A discontinuous random variable may have an infinity of possible discrete values

$$\dots x_{-2}, x_{-1}, x_0, x_1, x_2, \dots$$

If p_i be the probability of such a variable taking x_i then the conditions

$$p_i \geq 0$$

$$\sum_{i=-\infty}^{\infty} p_i = 1$$

must be satisfied.

Example Consider that we have ³ coins and we throw them on a board. The possible outcomes

(F.P.S.) are :-

				No. of cases	Probability
No head :	T	T	T	1	1/8
One head :	H	T	T	3	3/8
	T	H	T		
	T	T	H		
Two heads:	H	H	T	3	3/8
	H	T	H		
	T	H	H		
Three heads:	H	H	H	1	1/8

In this case we have $n=3$ independent trials. If x denotes the number of heads in tossing the three coins (or tossing one coin three times) then

$$\Pr \{x = 0 \mid n=3\} = 1/8$$

$$\Pr \{x = 1 \mid n=3\} = 3/8$$

$$\Pr \{x = 2 \mid n=3\} = 3/8$$

$$\Pr \{x = 3 \mid n=3\} = 1/8$$

It is obvious that $\sum p_i = 1/8 + 3/8 + 3/8 + 1/8 = 1$

Now since the probability that the head occurs in tossing a single coin is $\frac{1}{2}$, then we find that the probabilities that $x = 0, 1, 2, 3$, are the terms in the expansion $(\frac{1}{2} + \frac{1}{2})^3$ respectively. This of course is under the assumption that the two faces are equally likely. On the other hand if p denotes the probability that a head (H) occurs in tossing a single coin and $q = 1-p$ the probability that a tail (T) occurs i.e. p will be associated with (H) and q with (T), then

x	$P(x)$
0	q^3
1	$3 q^2 p$
2	$3 q p^2$
3	p^3

$$\sum_{x=0}^3 p(x) = q^3 + 3q^2p + 3qp^2 + p^3 = (q+p)^3 = 1$$

In general, suppose the probability of a success in a trial is "p" and the probability of a failure is $q = 1-p$. We can represent these probabilities in a functional form $f(\alpha)$ where $f(\alpha) = p$ for $\alpha = 1$ i.e. a success & $f(\alpha) = q$ for $\alpha = 0$ i.e. a failure. The probabilities associated with n trials which are mutually independent in the probability sense is

$$f(\alpha_1) f(\alpha_2) \dots f(\alpha_n)$$

The probability of x successes & $n-x$ failures in a specified order is

$$\{f(1)\}^x \{f(0)\}^{n-x} = p^x q^{n-x}$$

The no. of orders in which x successes & $n-x$ failures can occur is nC_x . These orders are mutually exclusive events, then

$$p(x) = {}^nC_x p^x q^{n-x} \quad (1)$$

This is the general term in the expansion $(q+p)^n$. $p(x)$ is called the binomial probability distribution

The mean of this distribution is given by

$$\begin{aligned} \mu'_1(x) &= \sum_{x=0}^n x {}^nC_x p^x q^{n-x} \\ &= np. \end{aligned} \quad (2)$$

Also

$$\begin{aligned} \mu'_2(x) &= \sum_{x=0}^n x^2 {}^nC_x p^x q^{n-x} \\ &= n(n-1)p^2 + np. \end{aligned}$$

$$\begin{aligned} \text{Hence } \mu_2(x) &= \mu'_2(x) - \left\{ \mu'_1(x) \right\}^2 \\ &= npq. \end{aligned}$$

i.e. the standard deviation = \sqrt{npq} .

Exercise

prove that

$$\beta_1 = \frac{1}{npq} - \frac{4}{n}$$

$$\beta_2 = 3 + \frac{1-6pq}{npq}$$

Relationship between the binomial distribution and the incomplete beta function.

The complete beta function

$$B(k, n-k+1) = \int_0^1 x^{k-1} (1-x)^{n-k} dx = \frac{\Gamma(k) \Gamma(n-k+1)}{\Gamma(n+1)}$$

$$= \frac{(k-1)! (n-k)!}{n!}$$

The incomplete beta function

$$B_p(k, n-k+1) = \int_0^p x^{k-1} (1-x)^{n-k} dx$$

$$\therefore I_p(k, n-k+1) = B_p(k, n-k+1) / B(k, n-k+1)$$

$$= \frac{n!}{(k-1)!(n-k)!} \int_0^p x^{k-1} (1-x)^{n-k} dx$$

$$= \frac{n!}{(k-1)!(n-k)!} \left[\frac{p^k q^{n-k}}{k} + \frac{p^{k+1} q^{n-k-1} (n-k)}{k(k+1)} + \dots \right.$$

$$\left. + \frac{p^n (n-k)!}{k(k+1) \dots (n-1)n} \right]$$

(18)

$$= \frac{n!}{k!(n-k)!} p^k q^{n-k} + \frac{n!}{(k+1)!(n-k-1)!} p^{k+1} q^{n-k-1} +$$

$$\dots\dots\dots + p^n$$

$$= P_{k,n} + P_{k+1,n} + \dots + P_{n,n}$$

where

$$P_{k,n} = \Pr \{x = k \mid n\}$$

$$\therefore I_p(k, n-k+1) = \sum_{x=k}^n P_{x,n} \quad (1)$$

Similarly it can be shown that

$$\begin{aligned} I_p(k+1, n-k) &= P_{k+1,n} + P_{k+2,n} + \dots + P_{n,n} \\ &= \sum_{x=k+1}^n P_{x,n} \end{aligned} \quad (2)$$

From (1) & (2) we get

$$P_{k,n} = I_p(k, n-k+1) - I_p(k+1, n-k) \quad (3)$$

Fitting a binomial distribution

Throwing 5 coins 100 times, the following table gives the frequency distribution of the number of coins (x) on which the head shows uppermost.

x:	0	1	2	3	4	5
f:	2	14	20	34	22	8

To fit a binomial distribution we 1st calculate the mean of the given frequency distribution

$$\text{i.e. } \bar{x} = 2.84$$

Then equating \bar{x} to the mean of the binomial

i.e. $\bar{x} = np$ where $n = 5$ we get

$$5p = 2.84.$$

$$\therefore p = 0.568$$

$$\therefore p(x) = \binom{5}{x} 0.568^x (0.432)^{5-x}$$

where $x = 0, 1, 2, 3, 4, 5$.

Values of $p(x)$ are calculated for the corresponding values of x as given in the table below:

x	$p(x)$	Expected frequencies	observed frequencies
0	0.015	1.5	2
1	0.099	9.9	14
2	0.260	26.0	20
3	0.342	34.2	34
4	0.225	22.5	22
5	0.059	5.9	8
		100.0	100.0

Note: The recurrent relation is

$$p(x+1) = \frac{n-x}{x+1} \cdot \frac{p}{q} \cdot p(x) \quad x=0,1,2,\dots,n-1$$

e.g. $p(1) = n \cdot \frac{p}{q} \cdot p(0)$

$$p(2) = \frac{n-1}{2} \cdot \frac{p}{q} \cdot p(1)$$

.....

$$p(n) = \frac{1}{n} \cdot \frac{p}{q} \cdot p(n-1)$$

(3-3) Continuous random variables

A continuous random variable (x) is a variable which may take any real value between certain limits such as age, height, weight etc. Its cumulative distribution function $F(x)$ is continuous and differentiable for all values of x . $F'(x)$ must exist and is called the probability density function of x and usually denoted by $p(x)$,

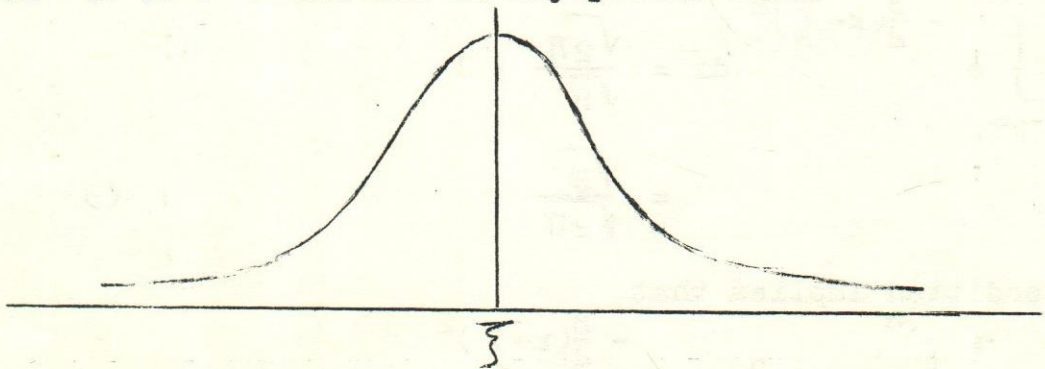
i.e. $p(x) = F'(x)$. Also

$$p(x) \geq 0$$

$$\& \int_{-\infty}^{\infty} p(x) dx = 1$$

As an example, let us consider the normal distribution. It was discovered by De Moivre in 1733, and at the same time but independently, Laplace & Gauss derived the formula for this distribution.

The normal curve is a bell shaped with its ends extending from $-\infty$ to $+\infty$ and has an asymptotic base.



To construct a mathematical function for this curve,

$$\frac{dy}{dx} = 0 \quad \text{at} \quad x = \xi$$

$$\& \frac{dy}{dx} = 0 \quad \text{at} \quad y = 0$$

$$\therefore \frac{dy}{dx} = -By(x - \xi) \quad (1) \quad B \text{ is } + \text{ve.}$$

It is clear that for $x < \bar{x}$, $\frac{dy}{dx}$ is +ve and for $x > \bar{x}$ $\frac{dy}{dx}$ is -ve. This means that the function attains its maximum at $x = \bar{x}$. Solving the differential equation (1) we get

$$\int \frac{dy}{y} = -B \int (x - \bar{x}) dx$$

$$\therefore \log y = -\frac{B}{2} (x - \bar{x})^2 + \log C$$

$$\therefore y = p(x) = C e^{-\frac{B}{2} (x - \bar{x})^2} \quad (2)$$

where C & B are Constants to be determined. We are going to impose two conditions:-

- (i) The area below the curve should be equal to unity.
- (ii) The variance is σ^2 i.e. $\mu_2 = \sigma^2$

The 1st condition implies that

$$C \int_{-\infty}^{\infty} e^{-\frac{B}{2}(x-\bar{x})^2} dx = 1$$

$$\text{But } C \int_{-\infty}^{\infty} e^{-\frac{B}{2}(x-\bar{x})^2} dx = \frac{\sqrt{2\pi}}{\sqrt{B}} C$$

$$\therefore C = \frac{\sqrt{B}}{\sqrt{2\pi}} \quad (3)$$

The 2nd Condition implies that

$$\mu_2 = C \int_{-\infty}^{\infty} (x - \bar{x})^2 e^{-\frac{B}{2}(x-\bar{x})^2} dx = \sigma^2$$

$$\text{But } C \int_{-\infty}^{\infty} (x - \bar{x})^2 e^{-\frac{B}{2}(x-\bar{x})^2} dx = \frac{\sqrt{2\pi}}{B^{3/2}} C = \frac{1}{B}$$

$$\text{i.e. } B = \frac{1}{\sigma^2} \quad (4)$$

$$\therefore C = \frac{1}{\sqrt{2\pi} \sigma} \quad (5)$$

$$\therefore p(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2\sigma^2} (x - \bar{x})^2} \quad (6)$$

\bar{x} & σ are the two parameters of the normal distribution.

If $Z = \frac{x - \bar{x}}{\sigma}$ i.e. the standardized normal variate

then

$$p(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} z^2} \quad (7)$$

which is the unit normal distribution

Remarks

(1) Now, if \bar{x} is known, formula (6) represents a family of normal curves having the same mean but differs in their variances. If σ is known, the same formula represents a family of normal curves having the same variance but differs in their means

$$(2) \quad \Pr \{x \leq x_1\} = \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^{x_1} e^{-\frac{1}{2} \left(\frac{x - \bar{x}}{\sigma}\right)^2} dx$$

$$\text{write } Z = \frac{x - \bar{x}}{\sigma}$$

$$\therefore dz = dx/\sigma$$

$$\begin{aligned} \therefore \Pr \{x \leq x_1\} &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{Z_1} e^{-\frac{1}{2} z^2} dz \\ &= \Pr \{Z \leq Z_1\} \end{aligned} \quad (8)$$

$\Pr \{x \leq x_1\}$ is called the cumulative distribution function denoted by $F(x)$ at $x = x_1$.

Also

$$\Pr \{Z \leq -z_1\} = \int_{-\infty}^{-z_1} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz$$

$$= \frac{1}{\sqrt{2\pi}} \int_{z_1}^{\infty} e^{-\frac{1}{2}z^2} dz = 1 - \Pr \{Z \leq z_1\}$$

The values of the cumulative distribution function at different values of Z are given in "Tables of the normal probability function"

Example: The mean weight of 500 students at a certain college is 151 lb. and the standard deviation is 15 lb. Assuming that the weights are normally distributed, find how many students weight

(i) more than 185 lb.

(ii) between 120 & 155 lb.

(i) standardized value of 185 is

$$Z = \frac{185 - 151}{15} = \frac{34}{15} = 2.30$$

$$\Pr \{Z > 2.3\} = 1 - \Pr \{Z \leq 2.3\} = 1 - .9893$$

$$= .0107$$

\therefore The number of students is $= 500 \times 0.0107 = 5$

(ii) standardized values

of 120 & 155 are

$$z_1 = \frac{120 - 151}{15} = -\frac{31}{15} = -2.0$$

$$z_2 = \frac{155 - 151}{15} = \frac{4}{15} = 0.3$$

$$\begin{aligned} \Pr \{ -2.0 \leq Z \leq 0.3 \} &= \Pr \{ Z \leq .3 \} - \Pr \{ Z \leq -2.0 \} \\ &= 0.6179 - .0228 = 0.5951 \end{aligned}$$

The number of students is = $500 \times 0.5951 = 298$

Fitting a normal curve

given a frequency distribution, it is required to fit a normal curve having the same mean and variance.

		(1)	(2)	(3)	(4)
Intervals	Observed frequencies	Standardized values	Areas $\Pr \{ Z \leq z \}$	Differences (Δ)	Expected freq.
60		-3.66	0.0001		
60- 65	3	-2.74	0.0031	0.0030	3.
65- 70	21	-1.83	0.0336	0.0305	30.5
70- 75	150	-0.91	0.1814	0.1478	147.8
75- 80	335	0.01	0.5040	0.3226	322.6
80- 85	326	0.93	0.8238	0.3198	319.8
85- 90	135	1.85	0.9678	0.1440	144.0
90- 95	26	2.77	0.9972	0.0294	29.4
95-100	4	3.68	0.9999	0.0027	2.7
	1000				

The following steps are followed:-

(i) Calculate the mean and the variance for the observed frequency distribution i.e.,

$$\bar{x} = 79.945$$

$$s = 5.445$$

- (ii) Calculate the standardized values for the upper limits of the intervals, e.g. the upper limit of the 1st interval becomes

$$\frac{65-79.945}{5.445} = -2.74$$

and so on as in column (1)

- (iii) Calculate the areas below these standardized values e.g.

$$\Pr \{x \leq -2.74\} = 0.0031$$

$$\Pr \{x \leq -1.83\} = 0.0336$$

and so on as in column (2)

- (iv) Calculate the areas corresponding to the intervals () by successive subtraction e.g. the area corresponding to the interval 60-65

$$= 0.0031 - 0.0001 = 0.0030.$$

and the area corresponding to the interval

$$65-70 = 0.0336 - 0.0031 = 0.0305$$

and so on as in column (3)

- (v) Calculate the expected frequencies N where N is the total number of frequencies.

e.g. the expected frequency corresponding to the interval

$$60-65 = 1000 \times 0.003 = 3$$

the expected frequency corresponding to the interval

$$65-70 = 1000 \times 0.0305 = 30.5 \text{ and so on as in column (4)}$$

(26)

(IV) Expected Values and Characteristic
Functions

IV Expected Values and Characteristic functions(4.1) Expected Values

Suppose N observations each of which may have any one of k finite number of values x_1, x_2, \dots, x_k . Suppose further that f_1 of these observations have the value x_1 , f_2 have the value x_2 , ... and in general f_i have the value x_i where $\sum f_i = N$.

$$\bar{x} = \frac{1}{N} \sum f_i x_i = \sum_i x_i \left(\frac{f_i}{N} \right)$$

as $n \rightarrow \infty$, $f_i / N \rightarrow p_i$

\therefore The quantity $\sum_{i=1}^k x_i p_i$ is called the expected value of the discontinuous random variable x which takes the possible values x_1, x_2, \dots, x_k with probabilities p_1, p_2, \dots, p_k respectively. The expected value is denoted by $\mathcal{E}(x)$ i.e.

$$\mathcal{E}(x) = \sum_{i=1}^k x_i p_i$$

The expected value is a theoretical parameter analogous to the observed sample mean.

The definition may be extended to discontinuous random values with an infinite number of possible values, i.e.

$$\mathcal{E}(x) = \sum_{-\infty}^{\infty} x_i p_i$$

For continuous random variables

$$\mathcal{E}(x) = \int_{-\infty}^{\infty} x p(x) dx$$

Any function of a random variable x , say $f(x)$, is also a random variable and has an expected value

$E\{f(x)\} = \sum f(x_i) p_i$ for discontinuous random variables and

$E\{f(x)\} = \int f(x) p(x) dx$ for continuous random variables.

Note The joint probability density function $p(x,y)$ has the following properties

- (i) $p(x,y) \geq 0$
- (ii) $Pr\{\alpha < x < \beta ; \alpha_1 < y < \beta_1\} = \int_{\alpha_1}^{\beta_1} \int_{\alpha}^{\beta} p(x,y) dx dy$
- (iii) $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x,y) dx dy = 1$
- (iv) $p_1(x) = \int_{-\infty}^{\infty} p(x,y) dy$
 $p_2(y) = \int_{-\infty}^{\infty} p(x,y) dx$
- (v) If x, y are independent
 $p(x,y) = p_1(x) \cdot p_2(y).$

Theorems on expected values

- (i) $E(A) = A$
- (ii) $E(Ax) = A E(x).$
- (iii) $E(x+y) = E(x) + E(y)$

Consider the case x, y are continuous random variables then

$$E(x+y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x+y) p(x,y) dx dy$$

(29)

$$\begin{aligned}
&= \int_{-\infty}^{\infty} x \int_{-\infty}^{\infty} p(x,y) dy dx + \int_{-\infty}^{\infty} y \int_{-\infty}^{\infty} p(x,y) dx dy \\
&= \int_{-\infty}^{\infty} x p_1(x) dx + \int_{-\infty}^{\infty} y p_2(y) dy \\
&= \bar{E}(x) + \bar{E}(y)
\end{aligned}$$

(iv) If x & y are independent then

$$\begin{aligned}
\bar{E}(xy) &= \bar{E}(x) \cdot \bar{E}(y) \\
\bar{E}(xy) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy p(x,y) dx dy \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy p_1(x) \cdot p_2(y) dx dy \\
&= \int_{-\infty}^{\infty} x p_1(x) dx \int_{-\infty}^{\infty} y p_2(y) dy \\
&= \bar{E}(x) \cdot \bar{E}(y)
\end{aligned}$$

Remarks

(i) The moments of a distribution are the expected values of the powers of the random variable which has the given distribution. i.e.

$$\begin{aligned}
\mu_r'(x) &= \bar{E}(x^r) = \int_{-\infty}^{\infty} x^r p(x) dx \\
\text{or} &= \sum_{i=-\infty}^{\infty} x_i^r p_i
\end{aligned}$$

$$\begin{aligned}
\mu_r(x) &= \bar{E}\{(x - \mu_1')^r\} = \int_{-\infty}^{\infty} (x - \mu_1')^r p(x) dx \\
\text{or} &= \sum_i (x_i - \mu_1')^r p_i
\end{aligned}$$

(ii) The variance

$$\sigma_x^2 = \mathcal{E}\{x - \mathcal{E}(x)\}^2$$

(iii) The coefficient of correlation

$$\rho = \frac{\mathcal{E}\{x - \mathcal{E}(x)\} \{y - \mathcal{E}(y)\}}{\sqrt{\mathcal{E}\{x - \mathcal{E}(x)\}^2 \cdot \mathcal{E}\{y - \mathcal{E}(y)\}^2}}$$

Linear functions of random variables

Given $\mathcal{E}(x_i) = \tau_i$, $\mathcal{E}(x_i - \mathcal{E}(x_i))^2 = \sigma_i^2$, then the standard error of any linear function

$$y = \sum_{i=1}^n \alpha_i x_i$$

of n random variables x_1, x_2, \dots, x_n is

$$\sigma_y^2 = \sum_{i=1}^n \alpha_i^2 \sigma_i^2 + 2 \sum_{i=1}^n \sum_{j=i+1}^n \alpha_i \alpha_j \sigma_i \sigma_j \rho_{ij}$$

α are constants, σ_i, σ_j are standard errors of x_i, x_j & ρ_{ij} their correlation coefficient.

$$\mathcal{E}(y) = \mathcal{E}\left\{\sum_{i=1}^n \alpha_i x_i\right\} = \sum_{i=1}^n \alpha_i \tau_i$$

$$\sigma_y^2 = \mathcal{E}\{y - \mathcal{E}(y)\}^2 = \mathcal{E}\left\{\sum_{i=1}^n \alpha_i x_i - \sum_{i=1}^n \alpha_i \tau_i\right\}^2$$

$$\begin{aligned}
&= E \left\{ \sum_{i=1}^n \alpha_i (x_i - \tau_i) \right\}^2 \\
&= E \left\{ \sum_{i=1}^n \alpha_i^2 (x_i - \tau_i)^2 + 2 \sum_i \sum_{j=i+1}^n \alpha_i \alpha_j (x_i - \tau_i)(x_j - \tau_j) \right\} \\
&= \sum_{i=1}^n \alpha_i^2 \sigma_i^2 + 2 \sum_i \sum_{j=i+1}^n \alpha_i \alpha_j \sigma_i \sigma_j p_{ij}
\end{aligned}$$

Examples

(1) $\bar{x} = \frac{1}{n} \sum x_i$

Here $\alpha = \frac{1}{n} \therefore E(\bar{x}) = \frac{1}{n} \sum_{i=1}^n \tau_i$

& $\sigma_{\bar{x}}^2 = \frac{1}{n^2} \sum \sigma_i^2 + \frac{2}{n^2} \sum_i \sum_{j=i+1}^n \sigma_i \sigma_j p_{ij}$

(2) If $y = x_1 \pm x_2$

$\therefore E(y) = \sum_{i=1}^2 \tau_i = \tau_1 \pm \tau_2$

$\sigma_y^2 = \sigma_1^2 + \sigma_2^2 \pm 2 \sigma_1 \sigma_2 p_{12}$

(3) If $y = x_1 \pm x_2$ and x_1, x_2 are independent then

$\sigma_y^2 = \sigma_1^2 + \sigma_2^2$

(4) If the x 's are independent and have the same standard error then

$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$

(5) If \bar{x}_1 is the mean of n_1 independent random variables, each with standard deviation σ_1 and \bar{x}_2 the mean of n_2 independent random variables each with standard deviation σ_2 , then

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

Quadratic functions of random variables

Given x_1, x_2, \dots, x_n are independent random variables,

$E(x_i) = \bar{x}_i$, $E(x_i - \bar{x}_i)^2 = \sigma_i^2$ for $i=1, 2, \dots, n$ then

$$E\left\{\sum_{i=1}^n (x_i - \bar{x})^2\right\} = \frac{n-1}{n} \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n (\bar{x}_i - \bar{x})^2$$

Let $y = \sum (x_i - \bar{x})^2$

$$\therefore y = \sum \left[(x_i - \bar{x}_i) + (\bar{x}_i - \bar{x}) - (\bar{x} - \bar{x}) \right]^2$$

$$= \sum_1 (x_i - \bar{x}_i)^2 + \sum_1 (\bar{x}_i - \bar{x})^2 + n(\bar{x} - \bar{x})^2$$

$$+ 2 \sum (x_i - \bar{x}_i)(\bar{x}_i - \bar{x}) - 2n(\bar{x} - \bar{x})^2$$

$$= \sum (x_i - \bar{x}_i)^2 + \sum_1 (\bar{x}_i - \bar{x})^2 - n(\bar{x} - \bar{x})^2 + 2 \sum (x_i - \bar{x}_i)(\bar{x}_i - \bar{x})$$

$$E(y) = \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n (\bar{x}_i - \bar{x})^2 - nE(\bar{x} - \bar{x})^2 + 0$$

$$= \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n (\bar{x}_i - \bar{x})^2 - \frac{1}{n} E\left\{\sum (x_i - \bar{x}_i)^2\right\}$$

$$\begin{aligned}
 &= \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n (\tau_i - \bar{\tau})^2 - \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \\
 &= \frac{n-1}{n} \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n (\tau_i - \bar{\tau})^2
 \end{aligned}$$

Lor:- (1) Suppose the x's have the same mean and standard deviation i.e.

$$\tau_i = \tau, \sigma_i = \sigma$$

$$\therefore \bar{\tau} = \frac{1}{n} \sum_i \tau_i = \tau$$

$$\mathcal{E}(y) = \mathcal{E} \sum (x_i - \bar{x})^2 = (n-1) \sigma^2$$

$$\frac{1}{n} \left\{ \sum_i (x_i - \bar{x})^2 \right\} = \frac{n-1}{n} \sigma^2$$

$$\& \mathcal{E} \left\{ \frac{1}{n-1} \sum_i (x_i - \bar{x})^2 \right\} = \sigma^2$$

Exercise:

Show that for two discontinuous random variables x & y

(i) $\mathcal{E}(x+y) = \mathcal{E}(x) + \mathcal{E}(y)$

(ii) $\mathcal{E}(xy) = \mathcal{E}(x) \mathcal{E}(y)$ if x & y are independent

$x_i \backslash y_j$	1	j	
1	p_{11}	p_{1j}	$p_{1.}$
.
.
.
i	p_{i1}	p_{ij}	$p_{i.}$
.
.
	$p_{.1}$	$p_{.j}$	1

$$\begin{aligned}
 \mathcal{E}(x_i + y_j) &= \sum_j \sum_i p_{ij} (x_i + y_j) \\
 &= \sum_j \sum_i p_{ij} x_i + \sum_j \sum_i p_{ij} y_j \\
 &= \sum_i p_{i.} x_i + \sum_j p_{.j} y_j \\
 &= \mathcal{E}(x) + \mathcal{E}(y)
 \end{aligned}$$

Remember that for x & y independent $p_{ij} = p_{i.} p_{.j}$

(4.2) Characteristic functions

$\varphi_x(t)$ is defined as the characteristic function of the random variable x if

$$\varphi_x(t) = \mathcal{E}(e^{itx})$$

This function will always exist since

$$|e^{itx}| = |\cos^2 tx + \sin^2 tx| = 1$$

$$\therefore \varphi_x(t) = \sum e^{itx} p(x) \text{ for discontinuous variables}$$

$$= \int e^{itx} p(x) dx \text{ for continuous variables}$$

Properties of characteristic functions

1- If "a" is constant then

$$\varphi_x(at) = \varphi_{ax}(t)$$

2- If x_1, x_2, \dots, x_n are independent

random variables then,

$$\varphi_{\sum x_j}(t) = \prod_{j=1}^n \varphi_{x_j}(t)$$

This is obvious because

$$\begin{aligned} \varphi_{\sum x}(t) &= \mathcal{E}(e^{it \sum x}) = \mathcal{E}\left\{e^{itx_1} e^{itx_2} \dots e^{itx_n}\right\} \\ &= \mathcal{E}(e^{itx_1}) \mathcal{E}(e^{itx_2}) \dots \mathcal{E}(e^{itx_n}) \\ &= \varphi_{x_1}(t) \cdot \varphi_{x_2}(t) \dots \varphi_{x_n}(t) \\ &= \sum_{j=1}^n \varphi_{x_j}(t) \end{aligned}$$

Lor (1) If the random variables have the same distribution then

$$\varphi_{\sum x}(t) = \left\{ \varphi_x(t) \right\}^n$$

Consequently it follows that

$$\varphi_{\bar{x}}(t) = \left\{ \varphi_x\left(\frac{t}{n}\right) \right\}^n$$

Characteristic functions as moment generating functions

$$\begin{aligned} \varphi_x(t) &= \int_{-\infty}^{\infty} e^{itx} p(x) dx \\ &= \int_{-\infty}^{\infty} p(x) \sum_{r=0}^{\infty} \frac{(itx)^r}{r!} dx \\ &= \sum_{r=0}^{\infty} \frac{(it)^r}{r!} \int_{-\infty}^{\infty} x^r p(x) dx \\ &= \sum_{r=0}^{\infty} \frac{(it)^r}{r!} \mu'_r(x). \end{aligned}$$

∴ Provided we can expand $\varphi_x(t)$ we can pick the moments around the origin. They are the coefficients of $\frac{(it)^r}{r!}$ in the expansion of $\varphi_x(t)$.

Also it can be shown that

$$i \mu'_1 = \varphi'_{x(0)}$$

$$i^2 \mu'_2 = \varphi''_{x(0)}$$

$$i^r \mu'_r = \varphi^{(r)}_{x(0)}.$$

Characteristic functions used for calculating cumulants

The cumulative function is defined as

$$\psi_x(t) = \log \{ \varphi_x(t) \}$$

$$\therefore \psi_x(t) = \log \left[1 + \frac{it}{1!} \mu'_1 + \frac{(it)^2}{2!} \mu'_2 + \frac{(it)^3}{3!} \mu'_3 + \dots \right]$$

$$= \frac{it}{1!} \mu'_1 + \frac{(it)^2}{2!} \mu'_2 + \frac{(it)^3}{3!} \mu'_3 + \dots$$

$$- \frac{1}{2} \left[\frac{it}{1!} \mu'_1 + \frac{(it)^2}{2!} \mu'_2 + \frac{(it)^3}{3!} \mu'_3 + \dots \right]^2$$

$$+ \frac{1}{3} \left[\frac{it}{1!} \mu'_1 + \frac{(it)^2}{2!} \mu'_2 + \frac{(it)^3}{3!} \mu'_3 + \dots \right]^3$$

$$- \dots \dots \dots$$

$$= itk_1 + \frac{(it)^2}{2!} k_2 + \frac{(it)^3}{3!} k_3 + \frac{(it)^4}{4!} k_4 + \dots$$

Equating the coefficients of the same power of t on both sides we get

$$k_1 = \mu'_1$$

$$k_2 = \mu'_2$$

$$k_3 = \mu'_3$$

$$k_4 = \mu'_4 - 3\mu'_2{}^2$$

$$k_5 = \mu'_5 - 10\mu'_3\mu'_2$$

Examples

(i) For the normal distribution we have

$$\begin{aligned} \varphi_x(t) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-\frac{1}{2}\left(\frac{x-\bar{x}}{\sigma}\right)^2} e^{itx} dx \\ &= e^{i\bar{x}t - \frac{1}{2}t^2\sigma^2} \end{aligned}$$

$$\therefore \psi_x(t) = \log \varphi_x(t) = i \xi t + \frac{(it)^2}{2} \sigma^2$$

$$\therefore k_1 = \xi$$

$$k_2 = \sigma^2$$

$$k_r = 0 \quad \text{for } r > 2$$

(ii) For the binomial distribution

$$\varphi_x(t) = \sum_{k=0}^n e^{itx} \binom{n}{x} p^x q^{n-x}$$

$$= \sum_{k=0}^n \binom{n}{x} (p e^{it})^x q^{n-x}$$

$$= (q + p e^{it})^n$$

$$\psi_x(t) = \log \varphi_x(t) = n \log (q + p e^{it})$$

$$= n \log [1 - p(1 - e^{it})]$$

$$= n \log \left[1 + p \left\{ \frac{it}{1!} + \frac{(it)^2}{2!} + \dots \right\} \right]$$

$$= \frac{(it)}{1!} np + \frac{(it)^2}{2!} n(p-p^2)$$

$$+ \frac{(it)^3}{3!} n(p - 3p^2 + 2p^3)$$

$$+ \frac{(it)^4}{4!} n(p - 7p^2 + 12p^3 - 6p^4)$$

$$+ \dots$$

$$\therefore k_1 = np$$

$$k_2 = n(p - p^2) = npq$$

$$\begin{aligned} k_3 &= n(p - 3p^2 + 2p^3) = np(1-p)(1-2p) \\ &= npq(q-p) \end{aligned}$$

$$\begin{aligned} k_4 &= np(1 - 7p + 12p^2 - 6p^3) \\ &= np \left[(1-p)(1-6p) + 6p^2(1-p) \right] \\ &= np(1-p)(1-6p+6p^2) \\ &= npq(1-6pq) \end{aligned}$$

Exercises (1) Given $p(x) = \frac{e^{-x} x^{l-1}}{(1)}$, the gamma distribut-

ion with parameter l , using the cumulant function prove that
 $2\beta_2 - 3\beta_1 - 6 = 0$

(2) Given $p_k = \frac{\lambda^k e^{-\lambda}}{k!}$ prove that

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

The Inverse theorem

This theorem states that the characteristic function uniquely determines the distribution function.

(i) If x is an absolutely continuous random variable, $p(x)$ is the p.d.f. and $\phi_x(t)$ is its characteristic function, then

$$p(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \varphi_x(t) dt$$

(ii) If p_k is the p.d.f. for a discontinuous random variable, then

$$p_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-itk} \varphi_k(t) dt$$

Example : Given

$$\varphi_x(t) = e^{-\frac{1}{2} t^2 \sigma^2}$$

then

$$p(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx - \frac{1}{2} t^2 \sigma^2} dt$$

$$p(x) = \frac{1}{2\pi} e^{-\frac{x^2}{2\sigma^2}} \int_{-\infty}^{\infty} e^{-\frac{1}{2} (t\sigma + \frac{ix}{\sigma})^2} dt$$

Taking $u = t + \frac{ix}{\sigma}$ then

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

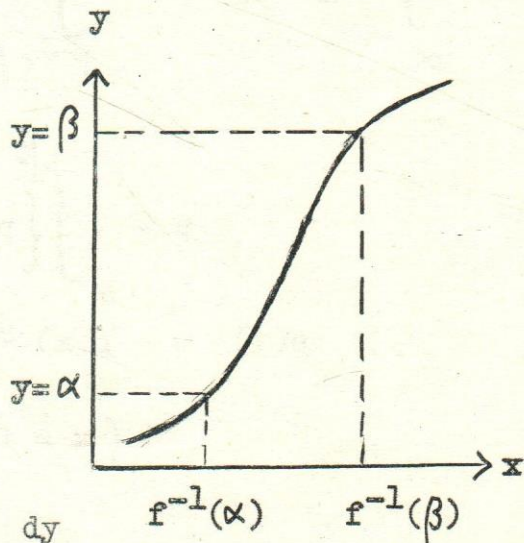
One-to-one transformation

Given $p(x)$ the probability density function of x & $y = f(x)$, we want $p(y)$ the probability density function of y .

In general, given $p(x_1, x_2, \dots, x_n)$, we want $p(y_1, y_2, \dots, y_n)$ we are going to deal with simple types of transformation for all of which $f(x)$ will be continuous function of x . In most cases the transformation will be also, one-to-one, i.e., to any value of x there corresponds one and only one value of, while to any value of y there corresponds one and only one value of x . Since $y = f(x)$ is continuous and differentiable, this implies that $f(x)$ must be either a decreasing or an increasing function of x . This also implies that $x = f^{-1}(y)$

(i) $f(x)$ is an increasing function of x

$$\begin{aligned}
 \Pr \{ \alpha < y < \beta \} \\
 &= \Pr \{ f^{-1}(\alpha) < x < f^{-1}(\beta) \} \\
 \text{i.e.} \quad \int_{\alpha}^{\beta} p(y) dy &= \int_{f^{-1}(\alpha)}^{f^{-1}(\beta)} p(x) dx \\
 &= \int_{\alpha}^{\beta} \left\{ p(x) \frac{dx}{dy} \right\} dy
 \end{aligned}$$



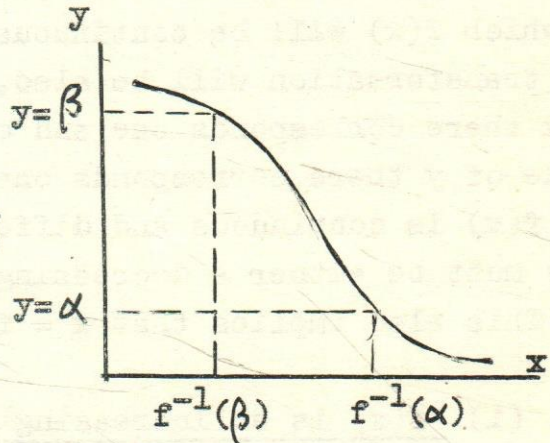
$$\therefore p(y) = p(x) \frac{dx}{dy} \quad \left(\frac{dx}{dy} \geq 0 \right) \quad (1)$$

$$(x = f^{-1}(y))$$

(ii) $f(x)$ is a decreasing function of x

$$\Pr \{ \alpha < y < \beta \}$$

$$= \Pr \{ f^{-1}(\beta) < x < f^{-1}(\alpha) \}$$



$$\text{i.e.} \quad \int_{\alpha}^{\beta} p(y) dy = \int_{f^{-1}(\beta)}^{f^{-1}(\alpha)} p(x) dx$$

$$= - \int_{\alpha}^{\beta} \left\{ p(x) \frac{dx}{dy} \right\} dy$$

$$\therefore p(y) = - p(x) \frac{dx}{dy} \quad \left(\frac{dx}{dy} \leq 0 \right) \quad (2)$$

$$(x = f^{-1}(y))$$

from (1) & (2) we have

$$p(y) = p(x) \left| \frac{dx}{dy} \right|$$

$$x = f^{-1}(y)$$

The transformation $y = x^2$

Pr y

$$= \Pr \{ \sqrt{\alpha} < x < \sqrt{\beta} \} + p \{ -\sqrt{\beta} < x < -\sqrt{\alpha} \}$$

$$\therefore \int_{\alpha}^{\beta} p(y) dy = \int_{\sqrt{\alpha}}^{\sqrt{\beta}} p(x) dx + \int_{-\sqrt{\beta}}^{-\sqrt{\alpha}} p(x) dx$$

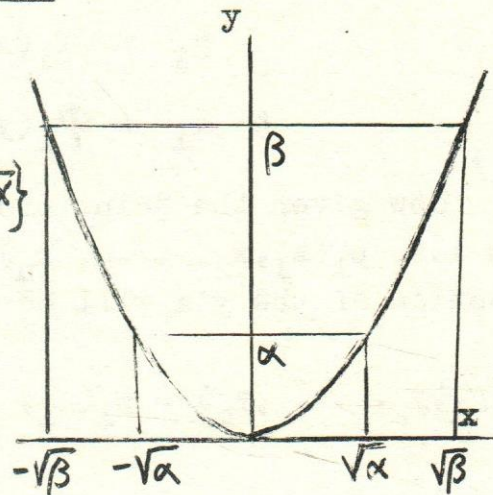
$$= \int_{\sqrt{\alpha}}^{\sqrt{\beta}} p(x) dx - \int_{-\sqrt{\alpha}}^{-\sqrt{\beta}} p(x) dx$$

$$= \int_{\alpha}^{\beta} p(x) \frac{dx}{dy} dy + \int_{\alpha}^{\beta} p(x) \frac{dx}{dy} dy$$

$x = \sqrt{y} \qquad \qquad \qquad x = -\sqrt{y}$

$$= \int_{\alpha}^{\beta} \frac{1}{2\sqrt{y}} \left\{ p(x)_{x=\sqrt{y}} + p(x)_{x=-\sqrt{y}} \right\} dy$$

$$\therefore p(y) = \frac{1}{2\sqrt{y}} \left\{ p(x)_{x=\sqrt{y}} + p(x)_{x=-\sqrt{y}} \right\} \quad (3)$$



One to one transformations of sets of n random variables

A set of n variables x_1, x_2, \dots, x_n and another set of the same number of variables y_1, y_2, \dots, y_n . A one-to-one transformation between these two sets means that each set of values of the x 's corresponds to a unique set of values of the y 's and conversely. Such a transformation may be written

$$y_i = f_i(x_1, x_2, \dots, x_n)$$

$$\& \ x_i = \phi_i(y_1, y_2, \dots, y_n)$$

Now given the joint probability density function of the x's i.e. $p_1(x_1, x_2, \dots, x_n)$, the joint probability density function of the y's will be

$$p_2(y_1, y_2, \dots, y_n) = p_1(1, 2, \dots, n) \left| \frac{\partial(x_1, x_2, \dots, x_n)}{\partial(y_1, y_2, \dots, y_n)} \right|$$

where $p_1(1, 2, \dots, n)$ is the joint probability density function of the x's expressed in terms of the y's.

Example If x_1, x_2 are two independent normally distributed variables with expected value zero and standard deviation , find the distribution of $\frac{1}{2}(x_1 + x_2)$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2}(x_1^2 + x_2^2)}$$

$$\text{Let } T = \frac{1}{2}(x_1 + x_2)$$

$$U = \frac{1}{2}(x_1 - x_2)$$

$$x_1 = T + u$$

$$x_2 = T - U$$

$$\frac{\partial(x_1, x_2)}{\partial(T, U)} = \begin{vmatrix} 1 & 1 \\ 1 & -1 \end{vmatrix} = -2$$

$$\begin{aligned}
 \therefore p(T, U) &= p(x_1, x_2) \left| \frac{\partial(x_1, x_2)}{\partial(T, U)} \right| \\
 &= \frac{1}{\pi \sigma^2} e^{-\frac{1}{\sigma^2} (T^2 + U^2)} \\
 \therefore p(T) &= \int_{-\infty}^{\infty} p(T, U) du = \frac{1}{\pi \sigma^2} \int_{-\infty}^{\infty} e^{-(T^2 + U^2)/\sigma^2} du \\
 &= \frac{1}{\sqrt{\pi} \sigma} e^{-T^2/\sigma^2}
 \end{aligned}$$

χ^2 - distribution

Let u_j be independent unit normal variates. χ^2 is defined as

$$\chi_y^2 = u_1^2 + u_2^2 + \dots + u_y^2 = \sum_{j=1}^y u_j^2 \quad (1)$$

$$\text{But } p(u_j) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} u_j^2} \quad -\infty \leq u_j \leq \infty \quad (2)$$

$$\text{write } v_j = u_j^2$$

$$\therefore p(v_j) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} v} v^{-\frac{1}{2}} \quad 0 \leq v \leq \infty \quad (3)$$

$$(i) \text{ To get } p(\chi_1^2) \quad \text{Let } v_1 = \chi_1^2$$

$$\therefore p(\chi_1^2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \chi_1^2} (\chi_1^2)^{-\frac{1}{2}} \quad 0 \leq \chi_1^2 \leq \infty$$

$$(ii) \text{ To get } p(\chi_2^2)$$

$$\text{Let } \chi_2^2 = v_1 + v_2 = \chi_1^2 + v_2, \quad y = \chi_1^2$$

$$\begin{aligned}\therefore p(x_1^2, v_2) &= p(x_1^2) \cdot p(v_2) \\ &= \frac{1}{2\pi} (x_1^2)^{-\frac{1}{2}} v_2^{-\frac{1}{2}} e^{-\frac{1}{2}(x_1^2 + v_2)}\end{aligned}$$

$$\begin{aligned}\therefore p(x_2^2, y) &= p(x_1^2, v_2) \left| \frac{\partial(x_1^2, v_2)}{\partial(x_2^2, y)} \right| \\ &= \frac{1}{2\pi} y^{-\frac{1}{2}} (x_2^2 - y)^{-\frac{1}{2}} e^{-\frac{1}{2}x_2^2} \quad 0 < y < x_2^2\end{aligned}$$

$$\begin{aligned}p(x_2^2) &= \int_0^{x_2^2} p(x_2^2, y) dy \\ &= \frac{1}{2} e^{-\frac{1}{2}x_2^2}\end{aligned}$$

(iii) To get $p(x_3^2)$

$$x_3^2 = x_2^2 + v_3, \quad y = x_2^2$$

$$\therefore p(x_3^2, v_3) = p(x_2^2) p(v_3)$$

$$= \frac{1}{2\sqrt{2\pi}} (x_3^2 - y)^{-\frac{1}{2}} e^{-\frac{1}{2}x_3^2}$$

$$\begin{aligned}p(x_3^2) &= \int_0^{x_3^2} p(x_3^2, y) dy \\ &= \frac{1}{\sqrt{2\pi}} (x_3^2)^{\frac{1}{2}} e^{-\frac{1}{2}x_3^2}\end{aligned}$$

\therefore In general

$$p(x^2) = c (x^2)^{\frac{1}{2}(v-2)} e^{-\frac{1}{2}x^2}$$

where $c = \frac{1}{2^{\gamma/2} \Gamma(\frac{\gamma}{2})}$

By mathematical induction, we proceed to find the p.d.f. of

$$x_{\gamma+1}^2 = x_{\gamma}^2 + v_{\gamma+1}$$

write $y = x^2$

$$\therefore p(x_{\gamma+1}^2, y) = \frac{c}{\sqrt{2\pi}} e^{-\frac{1}{2} x_{\gamma+1}^2} y^{(\gamma-2)/2} (x_{\gamma+1}^2 - y)^{-\frac{1}{2}}$$

$$\therefore p(x_{\gamma+1}^2) = \int_0^{x_{\gamma+1}^2} p(x_{\gamma+1}^2, y) dy$$

$$p(x_{\gamma+1}^2) = c' (x_{\gamma+1}^2)^{\frac{1}{2}(\gamma+1)-1} e^{-\frac{1}{2} x_{\gamma+1}^2}$$

This is the same as $p(x^2)$ but with $\gamma+1$ instead of γ which completes the proof.

Another proof

Let $v_j = u_j^2$

We know that $p(v_j) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} v} v^{-\frac{1}{2}}$

$$\begin{aligned} \text{and } \varphi_v(t) &= \frac{1}{\sqrt{2\pi}} \int_0^{\infty} e^{itv - \frac{1}{2}v} v^{-\frac{1}{2}} dv \\ &= (1 - 2it)^{-\frac{1}{2}} \end{aligned}$$

But v_j are independent where $j=1, 2, \dots, \gamma$

$$\begin{aligned} \therefore \varphi_{\sum v_j}(t) &= \{\varphi_v(t)\}^{\gamma} \\ &= (1 - 2it)^{-\frac{\gamma}{2}} \end{aligned}$$

i.e. $\varphi_{\chi^2}(t) = (1 - 2it)^{-\gamma/2}$

This is the characteristic function for Γ -function with parameter $\frac{1}{2} \gamma$

$$\therefore p(\chi_{\gamma}^2) = c e^{-\frac{1}{2} \chi_{\gamma}^2} (\chi_{\gamma}^2)^{\frac{1}{2} \gamma - 1}$$

Note : Given that u_j is a normal variate $N(0,1)$ & $v_j = u_j^2$ then

$$p(v_j) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} v_j} v_j^{-\frac{1}{2}}$$

By definition

$$\chi_{\gamma}^2 = \sum_{j=1}^{\gamma} u_j^2 = \sum_{j=1}^{\gamma} v_j$$

where v_j are independent variates

Let $y_j = \frac{1}{2} v_j$

$$\therefore p(y_j) = \frac{1}{\Gamma(\frac{1}{2})} e^{-y_j} y_j^{-\frac{1}{2}}$$

Hence $y_j = \frac{1}{2} v_j$ is distributed as Γ -variate with parameter $\frac{1}{2}$.

Consequently - applying the additive property of Γ -we have

$$\chi_{\gamma}^2 / 2 = \sum_{j=1}^{\gamma} y_j = \text{a } \Gamma\text{-variate with parameter } \gamma/2$$

$$\therefore p(\chi_{\gamma}^2 / 2) = \frac{1}{\Gamma(\gamma/2)} e^{-\frac{1}{2} \chi_{\gamma}^2} (\chi_{\gamma}^2 / 2)^{\frac{1}{2} \gamma - 1}$$

$$\therefore p(\chi_{\gamma}^2) = \frac{1}{2^{\frac{1}{2} \gamma} \Gamma(\gamma/2)} e^{-\frac{1}{2} \chi_{\gamma}^2} (\chi_{\gamma}^2)^{\frac{1}{2} \gamma - 1}$$

The Moments of χ^2 :

The cumulant function is $\varphi_{\chi^2}(t) = -\frac{1}{2} \gamma \log(1-2it)$

$$= \frac{it}{1!} \gamma + \frac{(it)^2}{2!} (2\gamma) + \frac{(it)^3}{3!} (8\gamma) + \frac{(it)^4}{4!} (48\gamma) + \dots$$

$$\therefore K_r = 2^{r-1} (r-1)! \gamma$$

Hence we have the following relationship $2\beta_2 - 3\beta_1 - 6 = 0$

Remark : It can be shown that

$$p(\chi_{\gamma}) = \frac{1}{2^{\frac{1}{2} \gamma - 1} \Gamma(\frac{1}{2} \gamma)} \chi_{\gamma}^{\gamma-1} e^{-\frac{1}{2} \chi_{\gamma}^2}$$

Also $\mu'_{2k}(\chi_{\gamma}) = \mu'_k(\chi_{\gamma}^2)$

A Discontinuous Distributions1- Poisson Distribution

It is in a sense a particular limiting form of the binomial distribution where $p \rightarrow 0$, $n \rightarrow \infty$ and np remains constant $= \lambda$ say.

$$\begin{aligned}
 P_{n,x} &= \frac{n!}{x! (n-x)!} p^x q^{n-x} \\
 &= \frac{n(n-1)\dots(n-x+1)}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\
 &= 1\left(1 - \frac{\lambda}{n}\right)\dots\left(1 - \frac{x-1}{n}\right) \left(\frac{\lambda}{n}\right)^x \frac{\left(1 - \frac{\lambda}{n}\right)^n}{\left(1 - \frac{\lambda}{n}\right)^x}
 \end{aligned}$$

when $n \rightarrow \infty$

$$P_{n,x} \rightarrow \frac{\lambda^x}{x!} e^{-\lambda}$$

The Poisson distribution is applicable to problems dealing with occurrence of events in a time interval of a given length such as:

- emission of rays from radio active substances
- certain traffic problems
- demands for telephone service
- bacteria count in cells
- distribution of bomb fragments in space.

The characteristic function

$$\begin{aligned}
 \varphi_x(t) &= \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} e^{-\lambda} e^{itx} \\
 &= e^{-\lambda(1-e^{it})}
 \end{aligned}$$

The cumulant function is

$$\psi_x(t) = \log \varphi_x(t) = -\lambda(1-e^{it})$$

all cumulants are constant and each $= \lambda$

$$\begin{aligned}
 \therefore \beta_1 &= 1/\lambda \\
 \beta_2 &= 3 + \frac{1}{\lambda}
 \end{aligned}$$

The relationship to Γ -function

$$\begin{aligned}\Gamma_x(n) &= \int_0^x v^{n-1} e^{-v} dv = \int_0^\infty v^{n-1} e^{-v} dv - \int_x^\infty v^{n-1} e^{-v} dv \\ &= \Gamma(n) + \int_x^\infty v^{n-1} d(e^{-v}) \\ &= \Gamma(n) - x^{n-1} e^{-x} - (n-1) x^{n-2} e^{-x} - (n-1)(n-2) \int_x^\infty e^{-v} v^{n-3} dv \\ \Gamma_x(n)/\Gamma(n) &= 1 - e^{-x} \left[1 + x + x^2/2! + \dots + x^{n-2}/(n-2)! + x^{n-1}/(n-1)! \right]\end{aligned}$$

$$I_x(n) = 1 - \sum_{t=0}^n e^{-x} x^t / t! = 1 - \Pr \{ r \leq n \}$$

$$\therefore \Pr \{ r \leq n \} = 1 - I_x(n)$$

2-The Negative Binomial

It is closely related to the Bernoulli binomial distribution. If we expand, according to the binomial theorem, $(q-p)^{-k}$ where $q=1+p$, $k > 0$, $p > 0$, we get the general term in

$$q^{-k} (1-p/q)^{-x} = q^{-k} \left[1 + k(p/q) + \frac{k(k+1)}{2!} \left(\frac{p}{q}\right)^2 + \dots \right]$$

$$\text{i.e.} \quad q^{-k} \frac{(k+x)}{x!} \left(\frac{p}{q}\right)^x$$

$$\text{Then} \quad p(x) = \frac{(k+x)}{x!} \left(\frac{p}{q}\right)^x \quad x=0 \text{ to } x=\infty$$

is called the negative binomial distribution.

(2.1) The moment generating function

$$\varphi_x(t) = \sum_{x=0}^{\infty} q^{-k} e^{tx} \frac{(k+x)}{x!} \left(\frac{p}{q}\right)^x = (q - pe^t)^{-k}$$

The mean is

$$\bar{G}(x) = \left. \frac{\partial \varphi}{\partial t} \right|_{t=0} = kp$$

$$\& \text{ Var}(x) = kpq$$

3- The Sequential Binomial

In this case we are going to fix the number of successes and vary the number of trials. We want k successes and we shall stop when they are achieved. The trials are independent and p is constant from trial to trial. The result will be the probability if $k-1$ achieved successes in the preceding $n-1$ trials multiplied by the probability of the n^{th} trial, i.e.

$$P_{n,k} = \binom{n-1}{k-1} p^{k-1} q^{n-k} p = \binom{n-1}{k-1} p^k q^{n-k}$$

(3.1) To prove that $\sum_{n=k}^{\infty} P_{n,k} = 1$

$$\begin{aligned} \sum_{n=k}^{\infty} \binom{n-1}{k-1} p^k q^{n-k} &= p^k + k p^k q + \frac{(k+1)k}{2!} p^k q^2 + \dots \\ &= p^k (1 - q)^{-k} = 1 \end{aligned}$$

(3.2) Moments

$$\mu'_1 = \sum_{n=k}^{\infty} n \binom{n-1}{k-1} p^k q^{n-k} = k/p$$

$$\mu'_2 = \sum_{n=k}^{\infty} n^2 \binom{n-1}{k-1} p^k q^{n-k} = k(k+1)/p^2 - k/p$$

$$\mu_2 = kq/p^2$$

Exercise Find μ_3 & μ_4

4- Hypergeometric Distribution

Finite population of balls of size N ; where there is Np black balls and Nq white balls. A sample of size n is drawn without replacement. What is the probability of k black balls in the sample of n ?

$$\Pr\{k=0\} = \frac{Nq}{N} \frac{Nq-1}{N-1} \frac{Nq-2}{N-2} \dots \frac{Nq-n+1}{N-n+1}$$

$$\Pr\{k=1\} = \frac{Np}{N} \frac{Nq}{N-1} \frac{Nq-1}{N-2} \dots \frac{Nq-n+2}{N-n+2} \binom{n}{1}$$

$$\Pr\{k=2\} = \frac{N_p}{N} \frac{N_p-1}{N-1} \frac{N_q}{N-2} \dots \frac{N_{q-n+3}}{N-n+3} \quad (2)$$

$$\Pr\{k=r\} = \frac{N_p}{N-1} \frac{N_p-1}{N-1} \dots \frac{N_p-r+1}{N-r+1} \frac{N_q}{N-r+2} \frac{N_q-1}{N-r+3} \dots \frac{N_{q-n+r+1}}{N-n+1} \quad \left(\frac{n}{r}\right)$$

$$= \frac{N_p!}{(N_p-r)!} \frac{N_q!}{(N_q-n+r)!} \frac{(N-n)!}{N!} \frac{n!}{r!(n-r)!}$$

$$= \frac{N_p!}{r!(N_p-r)!} \frac{N_q!}{(n-r)!(N_q-n+r)!} \frac{(N-n)!}{N!} \frac{n!}{n!}$$

$$= \binom{N_p}{r} \binom{N_q}{n-r} / \binom{N}{n}$$

(4.1) To prove that the sum of the probabilities = 1

$$\sum_{r=0}^n \frac{\binom{N_p}{r} \binom{N_q}{n-r}}{\binom{N}{n}} = \frac{1}{\binom{N}{n}} \sum_{r=0}^n \binom{N_p}{r} \binom{N_q}{n-r}$$

Now since $(1+x)^{N_p}(1+x)^{N_q} = (1+x)^N$

Equating the coefficients of x^r on both sides we get

$$\binom{N_q}{r} + \binom{N_p}{1} \binom{N_q}{r-1} + \dots + \binom{N_p}{r} = \binom{N}{r}$$

For $r=n$, we have

$$\binom{N_q}{n} + \binom{N_p}{1} \binom{N_q}{n-1} + \dots + \binom{N_p}{n} = \binom{N}{n}$$

$$\text{i.e. } \sum_{r=0}^n \binom{N_p}{r} \binom{N_q}{n-r} = \binom{N}{n}$$

Hence

$$\sum_{r=0}^n \frac{\binom{N_p}{r} \binom{N_q}{n-r}}{\binom{N}{n}} = \frac{1}{\binom{N}{n}} \sum_{r=0}^n \binom{N_p}{r} \binom{N_q}{n-r} = 1$$

(4.2) The moments

$$\mu_1' = \sum_{k=0}^n k \binom{Np}{k} \binom{Nq}{n-k} / \binom{N}{n}$$

$$= \sum_{k=0}^n \frac{Np!}{(k-1)!(Np-k)!} \frac{Nq!}{(n-k)!(Nq-n+k)!} \frac{n!(N-n)!}{N!}$$

$$= np \sum_{k=0}^n \binom{Np-1}{k-1} \binom{Nq}{n-k} / \binom{N-1}{n-1}$$

$$= np$$

$$\mu_2 = \sum_{k=0}^n k^2 \binom{Np}{k} \binom{Nq}{n-k} / \binom{N}{n} = \sum_{k=0}^n \{k(k+1) + k\} \binom{Np}{k} \binom{Nq}{n-k} / \binom{N}{n}$$

$$= np \left[1 + \frac{(n-1)(Np-1)}{N-1} \right]$$

$$\therefore \mu_2 = npq \left[1 - \frac{n-1}{N-1} \right]$$

As $N \rightarrow \infty$, $\mu_2 \rightarrow npq$

(4.3) Remark The problem can be represented in the following table

	A	\bar{A}	
Sample	r	n-r	n
Remainder	$Np-r$	$Nq-n+r$	$N-n$
Population	Np	Nq	N

As we mentioned before

$$\Pr\{k=r\} = \frac{Np!}{r! (Np-r)!} \frac{Nq!}{(n-r)!(Nq-n+r)!} \frac{(N-n)!n!}{N!}$$

In a 2×2 corresponding table, we have

a = r	c = n-r	n
b = $Np-r$	d = $Nq-n+r$	m = $N-n$
r = Np	s = Nq	N

Then $\Pr\{k=r\}$ becomes $\Pr\{k=r\} = \frac{r! s! m! n!}{a! b! c! d! N!}$

Also the mean $= np = n \frac{Np}{N} = \frac{nr}{N}$

$$\begin{aligned} \& \text{ the variance} &= npq \left(\frac{N-n}{N-1} \right) \\ &= np(1-p) \left(\frac{N-n}{N-1} \right) \\ &= \frac{nNp}{N} \frac{N-Np}{N} \frac{N-n}{N-1} \\ &= \frac{n r s m}{N^2(N-1)} \end{aligned}$$

(4.4) To test the significance of the difference between two treatments after being randomly assigned to a group of $N=n+m$ individuals

We are observing the absence or presence of reaction X. The 1st treatment T_1 is applied to n individuals & the 2nd treatment T_2 to m individuals. As a result $a/n, b/m$ show reaction X. The random process has been applied within the group of N individuals & its repetition would simply involve other random reassignments of the two treatments among the N . No assumption is made as to how the N individuals were selected from some large universe. On the null hypothesis there are $r=a+b$ individuals who will react & $s=c+d$ who do not react, whatever the arrangement of the treatments.

The chance that 'a' will react in n & 'b' in m , if the null hypothesis is true, will be

$$\Pr\{a|N, r, n\} = \frac{r!n!m!s!}{a!b!c!d!N!}$$

The mean of 'a' $= rn/N$ and

$$\text{the variance} = \frac{n r s m}{N^2(N-1)}$$

Example :

$$\begin{aligned} \text{Mean} &= rn/N = 6 \\ \text{Variance} &= 1.5 \\ \text{s.d.} &= 1.224 \end{aligned}$$

Using the approximation of the normal for the chance that

'a' = 2 or less we get $u = (2+0.5-6)/1.224 = -2.857$

	X	\bar{X}	
T_1	a	c	n
T_2	b	d	m
	r	s	N

	Survive	Die	
T_1	2	8	10
T_2	13	2	15
	15	10	25

From the table of the normal probability integral we get the probability of obtaining as large as a -ve deviation by chance is 0.00213. It is significantly small & it would be reasonable to suppose that T_2 is more effective than T_1 .

We can also use the hypergeometric

$$\Pr\{a=2\} = \frac{10!15!15!10!}{2!8!13! 2! 25!} = 0.001445$$

$$\Pr\{a=1\} = \frac{10!15!15!10!}{1! 9! 14! 1! 25!} = 0.0000458$$

$$\Pr\{a=0\} = \frac{10!15!15!10!}{0! 10! 15! 0! 25!} = 0.000003$$

$$\Pr\{a \leq 2\} = 0.001445 + 0.0000458 + 0.000003 = 0.0015$$

which is significant.

B Continuous Distributions

1- Gamma distribution

A continuous variable x is called a Gamma variate, if distributed with probability density function

$$p(x) = \frac{e^{-x} x^{l-1}}{\Gamma(l)} \quad (1)$$

throughout the range of values of x from $0 \rightarrow \infty$

The Characteristic function is given as

$$\varphi_x(t) = \int_0^{\infty} \frac{e^{-x} x^{l-1}}{\Gamma(l)} e^{itx} dx = (1 - it)^{-l} \quad (2)$$

$$\therefore \varphi_x(t) = -l \log(1-it) = l \left[it + \frac{(it)^2}{2} + \frac{(it)^3}{3} + \dots \right]$$

$$\therefore \mu'_1 = l, \mu_2 = l, \mu_3 = 2l, \mu_4 = 3l(l+2)$$

Theorem(1) If x is normally distributed with mean 'a' & s.d. = σ , then

$\frac{(x-a)^2}{2\sigma^2}$ is a gamma variate with parameter $\frac{1}{2}$

$$p(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp -\frac{1}{2} \left(\frac{x-a}{\sigma} \right)^2$$

$$\text{Let } y = \frac{1}{2} \left(\frac{x-a}{\sigma} \right)^2, \quad \therefore \frac{dy}{dx} = \frac{x-a}{\sigma^2}$$

$$\begin{aligned} \therefore p(y) &= \left[\frac{p(x)}{x-a=+\sigma\sqrt{2y}} + \frac{p(x)}{x-a=-\sigma\sqrt{2y}} \right] \left| \frac{dx}{dy} \right| \\ &= \frac{1}{\sqrt{\pi}} e^{-y} y^{-\frac{1}{2}} \end{aligned}$$

Hence y is a gamma variate with parameter $\frac{1}{2}$.

Theorem (2) The sum of two independent gamma variates with parameters l & m is a gamma variate with parameter $l+m$.

$$\text{Let } p(x) = \frac{e^{-x} x^{l-1}}{\Gamma(l)}, \quad p(y) = \frac{e^{-y} y^{m-1}}{\Gamma(m)}$$

$$\therefore p(x,y) = \frac{e^{-x-y} x^{l-1} y^{m-1}}{\Gamma(l) \Gamma(m)}$$

Write $w=x+y$, $v=x$; then $x=v$, $y=w-x$ & $|\partial(x,y)/\partial(w,v)| = 1$

$$p(w,v) = \frac{e^{-w} v^{l-1} (w-v)^{m-1}}{\Gamma(l) \Gamma(m)}$$

$$\therefore p(w) = \int_0^w p(w,v) dv = \frac{e^{-w} w^{l+m-1}}{\Gamma(l+m)}$$

Hence w is a gamma variate with parameter $l+m$

We can proceed in another way:-

$$\varphi_x(t) = (1-it)^{-l}, \quad \varphi_y(t) = (1-it)^{-m}$$

Since x, y are independent then

$$\varphi_{x+y}(t) = (1-it)^{-(l+m)}$$

which is the characteristic function of a gamma variate with parameter $l+m$

2- Beta Distribution

If x is a continuous variate distributed with probability density function

$$p(x) = \frac{x^{l-1} (1-x)^{m-1}}{B(l, m)}$$

throughout the range of values of $x=0$ to 1 ; then x is a Beta variate of the 1st kind with parameters ℓ, m

The characteristic function is given as

$$\begin{aligned}\varphi_x(t) &= \int_0^1 x^{\ell-1} (1-x)^{m-1} e^{itx} / B(\ell, m) dx \\ &= \frac{1}{B(\ell, m)} \left[B(\ell, m) + \frac{it}{1!} B(\ell+1, m) + \frac{(it)^2}{2!} B(\ell+2, m) \right. \\ &\quad \left. + \dots \dots \dots \right]\end{aligned}$$

The r^{th} moment is

$$\mu'_r = \frac{B(\ell+r, m)}{B(\ell, m)}$$

In particular

$$\mu'_1 = \frac{\ell}{\ell+m} \quad \& \quad \mu_2 = \frac{\ell m}{(\ell+m)^2(\ell+m+1)}$$

Theorem If x, y are independent gamma variates with parameters ℓ, m respectively, the quotient $x/(x+y)$ is a Beta variate with parameters ℓ, m of the 1st kind

Since x, y are independent, then

$$p(x, y) = \frac{e^{-(x+y)} x^{\ell-1} y^{m-1}}{\Gamma(\ell) \Gamma(m)}$$

Write $w = x/(x+y)$, $v = x$; then $x = v$, $y = v(1-w)/w$ and

$$| \partial(x, y) / \partial(w, v) | = v / w^2$$

$$\therefore p(w, v) = \frac{e^{-v/w} v^{\ell+m-1} (1-w)^{m-1} w^{-m-1}}{\Gamma(\ell) \Gamma(m)}$$

$$\therefore p(w) = \int_0^\infty p(w, v) dv = \frac{(1-w)^{m-1} w^{\ell-1}}{B(\ell, m)} = B_1(\ell, m)$$

Definition : A Beta variate of the 2nd type with +ve parameters ℓ, m is defined as a continuous variate x which is distributed with probability density function

$$p(x) = \frac{x^{\ell-1} (1+x)^{-(\ell+m)}}{B(\ell, m)}$$

throughout the range of values of $x=0$ to $x=\infty$ and denoted by $B_2(\ell, m)$.

$$E(x) = \int_0^1 x p(x) dx = \frac{1}{B(\ell, m)} \int_0^1 x (1+x)^{-(\ell+m)} dx$$

Write $y = x/(x+1)$

$$\therefore E(x) = \frac{1}{B(\ell, m)} \int_0^1 y^{\ell} (1-y)^{m-2} dy = \frac{\ell}{m-1}$$

$$\text{Also } \mu'_2 = \frac{\ell(\ell+1)}{(m-1)(m-2)}$$

$$\therefore \text{Var}(x) = \frac{\ell(\ell+m-1)}{(m-1)^2(m-2)}$$

Theorem : If x is a Beta variate of 2nd kind, its reciprocal is a variate of the same kind with parameters interchanged.

$$p(x) = \frac{x^{\ell-1}}{B(\ell, m)(1+x)^{\ell+m}} = B_2(\ell, m)$$

Write $y = 1/x$

$$\therefore p(y) = p(x) \left| \frac{\partial x}{\partial y} \right| = \frac{y^{m-1}}{B(\ell, m)(1+y)^{\ell+m}} = B_2(m, \ell)$$

EXERCISES

- 1- If $x_i (i=1, 2, \dots, n)$ are n independent variates normally distributed about a common mean μ , with s.d. σ_i and $\chi^2 = \sum x_i^2 / \sigma_i^2$ then $\frac{1}{2} \chi^2$ is a gamma variate with parameter $n/2$
- 2- Prove that a product of a $B_1(\ell, m)$ variate and an independent $\gamma(\ell+m)$ variate is a $\gamma(\ell)$ variate.
- 3- Prove that the quotient of two independent gamma variates with parameters ℓ, m is a $B_2(\ell, m)$ variate.

1- Bivariate Normal Distribution

(1-1) To build a mathematical equation of the surface representing the joint distribution of x, y . Take the origin at (x, y) and assume that:-

- (i) The marginal distribution of x is normal with s.d. = σ_x
- (ii) The regression of y on x is linear
- (iii) For given x , y varies normally
- (IV) Common array with s.d. = $\sigma_y \sqrt{1-\rho^2}$

From (i); if n_x is the frequency in the small interval δx

$$\therefore n_x = \frac{N}{\sqrt{2\pi} \sigma_x} e^{-\frac{1}{2} \frac{x^2}{\sigma_x^2}} \quad (1)$$

From (ii); $Y_x = \rho \frac{\sigma_y}{\sigma_x} x$ (2)

From (iii), (iv)

$$n_{xy} = \frac{n_x}{\sqrt{2\pi} \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left(\frac{y - Y_x}{\sigma_y}\right)^2} \delta y$$

$$= \frac{N}{2\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2} \left[\frac{x^2}{\sigma_x^2} + \frac{y - Y_x}{\sigma_y^2(1-\rho^2)} \right]} \delta x \delta y \quad (3)$$

Substituting from (2) in (3) we get

$$n_{xy} = \frac{N}{2\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[\frac{x^2}{\sigma_x^2} - 2\rho \frac{xy}{\sigma_x \sigma_y} + \frac{y^2}{\sigma_y^2} \right]} \delta x \delta y$$

If $\delta x, \delta y$ each $\rightarrow 0$, we can say

$$n_{xy} = Z \delta x \delta y$$

and if the volume under the surface is 1, then

$$Z = f(x,y) = \frac{1}{2\pi\sigma_x \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[\frac{x^2}{\sigma_x^2} - 2\rho \frac{xy}{\sigma_x \sigma_y} + \frac{y^2}{\sigma_y^2} \right]} \quad (4)$$

(1-2) Contours of the surface

$$Z = f(x,y) = \frac{1}{2\pi\sigma_1 \sigma_2 \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left(\frac{x^2}{\sigma_1^2} - 2\rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} \right)}$$

Z is constant when

$$\frac{x^2}{\sigma_1^2} - 2\rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} = k^2$$

which is the equation of an ellipse. If we change k, we get a family of coaxial similar ellipses. The regression line y on x is the diameter of the ellipse conjugate to y-axis also the regression line of x on y is conjugate to x-axis.

(i) To get the slope of major axis:-

The equation of the ellipse is

$$\frac{x^2}{\sigma_1^2} - 2\rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} = k^2$$

This can be written in the form

$$Ax^2 - 2Hxy + By^2 = k^2$$

If the angle of rotation of axis is θ , then θ is given by

$$\tan 2\theta = \frac{-2H}{A-B} = \frac{2\rho\sigma_1\sigma_2}{\sigma_1^2 - \sigma_2^2}$$

(62)

(ii) The area of the ellipse:-

Again the equation of the ellipse is

$$\frac{x^2}{\sigma_1^2} - 2\rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} = k^2 \quad (1)$$

Consider the concentric circle

$$x^2 + y^2 = r^2 \quad (2)$$

$$\therefore x^2 \left(\frac{1}{\sigma_1^2} - \frac{k^2}{r^2} \right) - 2\rho \frac{xy}{\sigma_1 \sigma_2} + y^2 \left(\frac{1}{\sigma_2^2} - \frac{k^2}{r^2} \right) = 0 \quad (3)$$

This equation has equal roots if

$$\rho^2 \frac{1}{\sigma_1^2 \sigma_2^2} = \left(\frac{1}{\sigma_1^2} - \frac{k^2}{r^2} \right) \left(\frac{1}{\sigma_2^2} - \frac{k^2}{r^2} \right)$$

$$\text{i.e. } \frac{1 - \rho^2}{\sigma_1^2 \sigma_2^2} r^4 - \left(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} \right) k^2 r^2 + k^4 = 0$$

 \therefore The area is $\pi r_1 r_2$ where

$$r_1^2 r_2^2 = \frac{k^4 \sigma_1^2 \sigma_2^2}{1 - \rho^2}$$

$$\therefore \text{The area is } \frac{\pi k^2 \sigma_1 \sigma_2}{\sqrt{1 - \rho^2}}$$

Also, for the ellipse

$$\frac{x^2}{\sigma_1^2} - 2\rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} = (k + \delta k)^2,$$

The area is

$$\frac{\pi (k + \delta k)^2 \sigma_1 \sigma_2}{\sqrt{1 - \rho^2}}$$

Hence, the area between the two ellipses is

$$dx dy = \frac{2k (\delta k) \sigma_1 \sigma_2}{\sqrt{1 - \rho^2}}$$

∴ The probability of a point to fall between the two ellipses is

$$\frac{dx \, dy}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}k^2}$$

$$= \frac{k(\delta k)}{1-\rho^2} e^{-\frac{1}{2(1-\rho^2)}k^2}$$

∴ The probability that the point will fall inside the ellipse k is

$$\frac{1}{1-\rho^2} \int_0^k k e^{-\frac{1}{2(1-\rho^2)}k^2} dk$$

$$= 1 - e^{-\frac{k^2}{2(1-\rho^2)}}$$

e.g. if the probability is 0.95, then

$$= 1 - e^{-\frac{1}{2(1-\rho^2)}k^2} = 0.95$$

$$\text{i.e. } e^{-\frac{1}{2(1-\rho^2)}k^2} = 0.05 \text{ from which we get } k.$$

(1-3) The Marginal & Conditional distributions

If the mean of the x 's is " a " and that of the y 's is b and if we write

$$x = \frac{x-a}{\sigma_1}, \quad y = \frac{y-b}{\sigma_2}$$

then the equation for the bivariate normal can be written in the form

$$f(x,y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}[x^2-2\rho xy + y^2]}$$

Integrating out y we get

$$f_1(x) = e \int_{-\infty}^{\infty} e^{-\frac{1}{2(1-\rho^2)} \left[\frac{x^2}{\sigma_1^2} - \rho \frac{xy}{\sigma_1 \sigma_2} + \frac{y^2}{\sigma_2^2} \right]} dy$$

$$= \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{x^2}{2\sigma_1^2}} = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-a)^2}{2\sigma_1^2}}$$

This is the marginal distribution of x .

Similarly we can get the marginal distribution of y

$$f_2(y) = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(y-b)^2}{2\sigma_2^2}}$$

Now the conditional probability function for y given x is

$$f(y|x) = f(x,y) / f_1(x)$$

$$\therefore f(y|x) = \frac{1}{\sqrt{2\pi}\sigma_2 \sqrt{1-\rho^2}} e^{-\frac{1}{2\sigma_2^2(1-\rho^2)} \left[y-b-\rho \frac{\sigma_2}{\sigma_1} (x-a) \right]^2}$$

This means that for fixed x , y is distributed normally with

$$(i) \quad E(y|x) = b + \rho \frac{\sigma_2}{\sigma_1} (x-a)$$

So the regression function of y on x is linear

$$(ii) \quad \text{Var}(y|x) = \sigma_2^2 (1-\rho^2)$$

Then the nearer ρ^2 is to 1, the smaller is the variance. If $\rho=0$, y does not depend on x and the two variates are independent. Similarly

$$f(x|y) = \frac{1}{\sqrt{2\pi}\sigma_1 \sqrt{1-\rho^2}} e^{-\frac{1}{2\sigma_1^2(1-\rho^2)} \left[x-a-\rho \frac{\sigma_1}{\sigma_2} (y-b) \right]^2}$$

and

$$(i) E(x|y) = a + \rho \frac{\sigma_1}{\sigma_2} (y-b)$$

$$(ii) \text{Var}(x|y) = \sigma_1^2 (1-\rho^2)$$

The Characteristic function

$$\begin{aligned} \varphi_{x,y}(t_1, t_2) &= \iint_{-\infty}^{\infty} e^{it_1 x_1 + it_2 x_2} f(x, y) dx dy \\ &= e^{-\frac{(it_1)^2}{2!} \sigma_1^2 - \frac{it_1}{1!} \frac{it_2}{2!} \sigma_1 \sigma_2 \rho - \frac{(it_2)^2}{2!} \sigma_2^2} \end{aligned}$$

$$\therefore \psi_{x,y}(t_1, t_2) = \frac{(it_1)^2}{2!} \sigma_1^2 + \frac{it_1}{1!} \frac{it_2}{1!} \rho \sigma_1 \sigma_2 + \frac{(it_2)^2}{2!} \sigma_2^2$$

$$\therefore K_{20} = \sigma_1^2, K_{11} = \sigma_1 \sigma_2 \rho, K_{02} = \sigma_2^2$$

and all other cumulants are = 0. In other words

$$\mu_{10} = 0, \mu_{20} = \sigma_1^2, \mu_{40} = 3\sigma_1^4$$

$$\mu_{01} = 0, \mu_{02} = \sigma_2^2, \mu_{04} = 3\sigma_2^4$$

$$\mu_{11} = \sigma_1 \sigma_2 \rho, \mu_{22} = (1+2\rho^2) \sigma_1^2 \sigma_2^2$$

$$\mu_{13} = 3\rho \sigma_1 \sigma_2^3$$

$$\mu_{31} = 3\rho \sigma_1^3 \sigma_2$$

Note: The bivariate normal distribution function contains five parameters. Four of them $(a, \sigma_1), (b, \sigma_2)$ characterize the two marginal distributions, while the fifth, ρ , characterizes the relationship between the two variables. It can be shown that the parameter ρ is equal to the mean value of the product of the standardized variables. As an estimate, r , of the parameter we compute $r = s_{12}/s_1 s_2$ where s_{12} denotes the estimate of the covariance. The parameter ρ is the correlation coefficient and r is its best estimate.

2- Bivariate Binomial Distribution

(2.1) In generalization to the binomial, consider the drawing of a sample of n from a population where the individuals may or may not have two attributes A & \bar{A} , B & \bar{B} . Suppose the proportions of the individuals with attributes AB , $\bar{A}B$, $A\bar{B}$, $\bar{A}\bar{B}$ are p_1, p_2, p_3, p_4 respectively where $p_1+p_2+p_3+p_4=1$

	A	\bar{A}	
B	p_1	p_3	p_1+p_3
\bar{B}	p_2	p_4	p_2+p_4
	p_1+p_2	p_3+p_4	1

	A	\bar{A}	
B	h	k	x_2
\bar{B}	j	l	$n-x_2$
	x_1	$n-x_1$	n

In exactly the same way as for the binomial, it is seen that the proportion of samples with h of the attribute (A B), j of the attribute ($\bar{A}B$), k of the attribute ($A\bar{B}$) & l of the attribute ($\bar{A}\bar{B}$), is

$$\frac{n!}{h! j! k! l!} p_1^h p_2^j p_3^k p_4^l$$

and the distribution of samples is given by the expansion of the multinomial form

$$(p_1+p_2+p_3+p_4)^n$$

The distribution given by this form is bivariate, one variate being the number of the A's i.e. (x_1) and the other being the number of the B's i.e. (x_2)

(2.2) The characteristic function

$$\varphi_{x_1, x_2} = E(e^{it_1 x_1 + it_2 x_2})$$

$$= \sum e^{it_1 x_1 + it_2 x_2} \frac{n!}{h! j! k! l!} p_1^h p_2^j p_3^k p_4^l$$

$$\begin{aligned}
&= \sum e^{it_1(h+j)+it_2(h+k)} \frac{n!}{h!j!k!l!} p_1^h p_2^j p_3^k p_4^l \\
&= \sum \frac{n!}{h!j!k!l!} \left\{ p_1 e^{i(t_1+t_2)} \right\}^h \left\{ p_2 e^{it_1} \right\}^j \left\{ p_3 e^{it_2} \right\}^k p_4^l \\
&= \left\{ p_1 e^{it_1+it_2} + p_2 e^{it_1} + p_3 e^{it_2} + p_4 \right\}^n
\end{aligned}$$

$$\begin{aligned}
\therefore \psi_{x_1, x_2}(t_1, t_2) &= n \log \left\{ p_1 e^{it_1+it_2} + p_2 e^{it_1} + p_3 e^{it_2} + p_4 \right\} \\
&= n \log \left[(p_1+p_2+p_3+p_4) + \frac{it_1}{1!} (p_1+p_2) \right. \\
&\quad + \frac{it_2}{1!} (p_1+p_3) + \frac{(it_1)^2}{2!} (p_1+p_2) + \frac{(it_2)^2}{2!} (p_1+p_3) \\
&\quad \left. + \frac{it_1}{1!} \frac{it_2}{1!} p_1 + \dots \right] \\
&= n \left[\frac{it_1}{1!} (p_1+p_2) + \frac{it_2}{1!} (p_1+p_3) \right. \\
&\quad + \frac{(it_1)^2}{2!} \left\{ p_1+p_2 - (p_1+p_2)^2 \right\} + \frac{(it_2)^2}{2!} \left\{ p_1+p_3 - (p_1+p_3)^2 \right\} \\
&\quad \left. + \frac{it_1}{1!} \frac{it_2}{1!} \left\{ p_1 - (p_1+p_2)(p_1+p_3) \right\} + \dots \right]
\end{aligned}$$

$$\therefore K_{10} = n(p_1+p_2), \quad K_{01} = n(p_1+p_3)$$

$$K_{20} = n(p_1+p_2) \left\{ 1 - (p_1+p_2) \right\}$$

$$K_{02} = n(p_1+p_3) \left\{ 1 - (p_1+p_3) \right\}$$

$$K_{11} = n \left\{ p_1 - (p_1+p_2)(p_1+p_3) \right\}$$

If we transfer the origin to the mean of the variates we have

$$\psi = \log \varphi$$

$$= -\frac{n}{2} \left[t_1^2(p_1+p_2) \{1-(p_1+p_2)\} + t_2^2(p_1+p_3) \{1-(p_1+p_3)\} + 2t_1t_2 \{p_1-(p_1+p_2)(p_1+p_3)\} \right] + O(n)$$

(2-3) The approximation of Biv. Binomial to the bivariate normal

When the distribution is expressed in standard measure we have

$$\varphi_{x_1-k_{10}, x_2-k_{01}}(t_1, t_2) = e^{it_1(x_1-k_{10})+it_2(x_2-k_{01})} \\ = e^{-(it_1K_{10}+it_2K_{01})} \varphi_{x_1, x_2}(t_1, t_2)$$

$$\therefore \varphi_{\frac{x_1-k_{10}}{\sigma_{20}}, \frac{x_2-k_{01}}{\sigma_{02}}}(t_1, t_2) \\ = e^{-\left(\frac{it_1K_{10}}{\sigma_{20}} + \frac{it_2K_{01}}{\sigma_{02}}\right)} \varphi_{x_1, x_2}(t_1, t_2)$$

$$\log \varphi_{\frac{x_1-k_{10}}{\sigma_{20}}, \frac{x_2-k_{01}}{\sigma_{02}}}(t_1, t_2) \\ = \frac{(it_1)^2}{2!} + \frac{(it_2)^2}{2!} + \frac{it_1}{1!} \frac{it_2}{1!} \frac{K_{11}}{\sigma_{20}\sigma_{02}} + O\left(\frac{1}{\sqrt{n}}\right) \\ \rightarrow -\frac{1}{2} (t_1^2 + t_2^2 + 2t_1t_2) \text{ as } n \rightarrow \infty \quad (1)$$

where

$$\rho = \frac{K_{11}}{\sigma_{20}\sigma_{02}} = \frac{p_1-(p_1+p_2)(p_1+p_3)}{\sqrt{[(p_1+p_2)(1-p_1-p_2)][(p_1+p_3)(1-p_1-p_3)]}}$$

$$\therefore \varphi_{\frac{x_1 - K_{10}}{\sigma_{20}}, \frac{x_2 - K_{01}}{\sigma_{02}}}(t_1, t_2) \rightarrow e^{-\frac{1}{2}(t_1^2 + 2\rho t_1 t_2 + t_2^2)}$$

which is the characteristic of the bivariate normal surface with zero mean, unit standard deviation and coefficient of correlation ρ i.e. of the form

$$p(x_1, x_2) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}(x_1^2 - 2\rho x_1 x_2 + x_2^2)} \quad (2)$$

In other words the standardized bivariate binomial distribution tends to the bivariate normal distribution (2) as $n \rightarrow \infty$.

(70)

(VII) Sampling and Sampling Distributions

Sampling and Sampling distributions

5-1 Basic definitions

Population is a totality of objects under Consideration It refers not alone to persons but to physical objects and operations as well.

A real population is one which actually exists in the form of specific population objects.

A hypothetical population is one which the statistician imagines to exist.

A finite population is one in which the number of population objects is countable and definitely limited in number.

An infinite population is one in which the number of objects or operations is unlimited or cannot be counted.

Sample is a part selected or drawn from a statistical population which is used as a basis of making estimates and inferences about the population.

If the selection takes place so that every unit has an equal chance of being selected, the sample is called random sample.

A small sample usually refers to one which has fewer than 30 sampling units. The value of the distinction between small and large samples lies in the fact that, although certain statistical operations can be applied to large random samples, they cannot be applied to small samples.

Statistic (S):- it refers to a value based on sample values and sometimes known as an estimate. It varies from one sample to another.

Parameter (P): it refers to a value calculated from the population. It is constant for this population.

Sampling error (E): it is the difference between the statistic and the corresponding parameter i.e. $E = S - P$.

If the expected value of the statistic (or the estimate) is equal to the corresponding parameter, the statistic is unbiased. Otherwise it is biased.

Sampling distribution: The frequency distribution of a statistic (an estimate) derived from an indefinitely large number of random samples each of size n is called a sampling distribution.

The standard deviation of this distribution is the standard error of the estimate. The less the variance of this distribution, the greater the precision of the estimate. The variance of a sampling distribution can be reduced by increasing the size of the sample or by using a more efficient sample design.

(5-2) Estimation: assume that the parent population is distributed in a form which would be completely determinate if we knew the value of some parameter θ . We are given a sample of values x_1, x_2, \dots, x_n .

We require to determine with the aid of the x 's, a number which can be taken to be the value of θ , or a range of numbers which can be taken to include that value. We cannot expect to find any method of estimation which can be generated to give a close estimate of θ on every occasion and for every sample. We must derive a rule, a method or a formula which will give good results "in the long run" or on the average. The method or rule of estimation is called an Estimator and the value to which it gives rise in particular cases, the Estimate. It is itself a random variable and has a distribution.

In general:-

A population with a density function $f(x; \theta_1, \theta_2, \dots, \theta_k)$ where x is the variate and $\theta_1, \theta_2, \dots, \theta_k$ are parameters. On the basis of a random sample of observations say x_1, x_2, \dots, x_n , wish to estimate one or more of the parameters $\theta_1, \theta_2, \dots, \theta_k$. The question is to find functions of the observations which may be represented by

$$\hat{\theta}_1(x_1, x_2, \dots, x_n), \hat{\theta}_2(x_1, x_2, \dots, x_n), \dots$$

such that the distribution of these of functions will be concentrated as closely as possible near the true values of the parameters. The properties of good estimators are :-

- (1) Consistent "t is a consistent estimator of θ if it converges to θ in probability"

$$\text{i.e. } P_r\{t \rightarrow \theta\} \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

- (2) unbiased: "t is said to be an unbiased estimator of θ if $E(t) = \theta$."

- (3) efficient: If for two consistent estimators t_1, t_2 , we have $\text{var}(t_1) < \text{var}(t_2)$ for all n , then t_1 is more efficient than t_2 for all sample sizes.

- (4) sufficient: If the joint density of the sample is expressible in the form

$$\begin{aligned} L(x_1, x_2, \dots, x_n, \theta) &= \prod_{i=1}^n f(x_i, \theta) \\ &= L_1(t, \theta) \cdot L_2(x_1, x_2, \dots, x_n) \end{aligned}$$

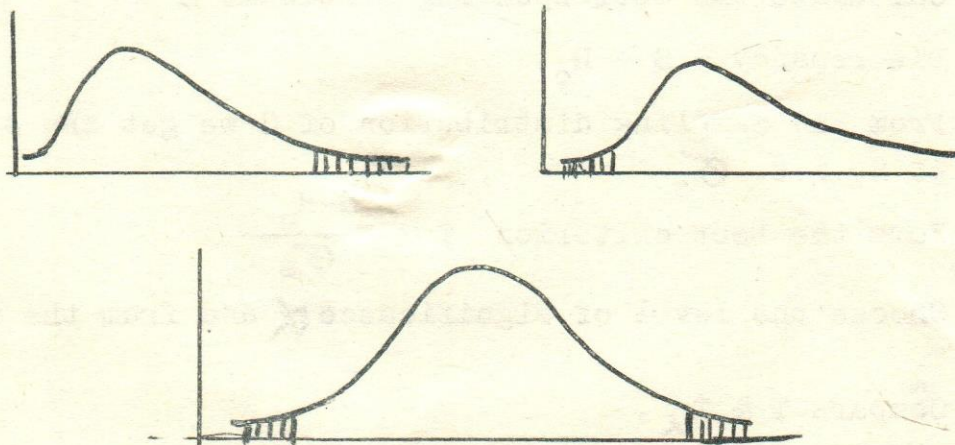
where L_1 does not contain the x 's and L_2 is independent of θ , then t is said to be a sufficient estimator of θ .

(5-3) Testing statistical hypotheses

An inference: is any conclusion made on the basis of experience.

Hypothesis: refers to the more systematic and formal statement of an inference usually in the form of a proposition

Statistical hypotheses: are usually classified into null and alternative hypotheses. A test of a such a hypothesis is consisting of testing the null hypothesis against the alternatives. In other words, the testing of a statistical hypothesis consists of formulating and applying an objective criterion for the purpose of accepting or rejecting the hypothesis. This criterion (test-statistic) will have a sampling distribution which indicated the probability of obtaining a value larger than the one obtained. If the obtained value of the test statistic is very likely value, then this is taken to be statistical evidence that the hypothesis is a plausible one. This process consists of selecting a critical region (a portion or portions of the sampling distribution of the statistic being used) and agreeing to reject the hypothesis under test whenever a random sample gives rise to a value of the statistic which falls in this critical region. Hence the sampling distribution of the test statistic must be known when the hypothesis is true



In some problems, the critical region will be only at one end either the upper or the lower and the test in such case is called a one-tailed test. In other problems, the critical region (regions of rejection) will be at both ends and the test is called a two-tailed test. This will depend upon whether we have one alternative hypothesis or two respectively.

The proportion of the probability distribution contained in the critical region is the level of significance of the test. In practice two levels are commonly used 1% or 5%

The size of the critical region is related to the risk one wants to accept in testing a statistical hypothesis. If the region is made too large, we reject too many true hypotheses.

Rejecting true hypothesis is called an error of the 1st kind or Type I error. If the region is too small, we accept too many false hypotheses.

Accepting false hypotheses is called an error of the 2nd kind or Type II error

Now for testing a statistical hypotheses $H = H_0$ the procedure is summerized as follows:-

(1) Null hypothesis $H = H_0$

Alternative hypotheses $H \neq H_0$

(2) Calculate the corresponding statistic S

(3) Discrepancy $= S - H_0$

(4) From the sampling distribution of S we get the standard error of S , i.e. σ_s

(5) Form the test criterion $T = \frac{S - H_0}{\sigma_s}$

(6) Choose the level of significance α and from the tables find

(7) Compare T & T_α .

If $T > T_\alpha$ reject the null hypothesis

If $T < T_\alpha$ accept the null hypothesis

(5.4) Sampling Distributions(1) Sampling distribution of the mean \bar{x}

A random sample of n independent values x_1, x_2, \dots, x_n is drawn from a normal population with mean \bar{x} and variance σ^2 .

The characteristic function for the normal distribution is

$$\varphi_x(t) = e^{it\bar{x} + \frac{(it)^2}{2} \sigma^2}$$

But

$$\begin{aligned} \varphi_{\bar{x}}(t) &= \left\{ \varphi_x\left(\frac{t}{n}\right) \right\}^n = \left\{ e^{\frac{it}{n}\bar{x} + \frac{(it)^2}{2} \frac{\sigma^2}{n}} \right\}^n \\ &= e^{it\bar{x} + \frac{(it)^2}{2} \frac{\sigma^2}{n}} \end{aligned}$$

This is the characteristic function of a normal distribution with mean \bar{x} and variance $\frac{\sigma^2}{n}$.

\therefore The sampling distribution of the means of ^{indep.} random samples of size n drawn from a normal population with mean \bar{x} and variance σ^2 is also a normal distribution with the same mean and variance $\frac{\sigma^2}{n}$.

(2) The sampling distribution of s^2

If x_i are n independent normal variates with mean \bar{x} and variance σ^2 then

$$u_i = \frac{x_i - \bar{x}}{\sigma} \quad (1)$$

is a unit normal variate.

Consequently $\chi_n^2 = \sum_{i=1}^n u_i^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2}$

If we replace \sum by \bar{x} then we get x_{n-1}^2

$$\begin{aligned} \text{i.e. } x_{n-1}^2 &= \frac{\sum (x_i - \bar{x})^2}{\sigma^2} \\ &= \frac{ns^2}{\sigma^2} \end{aligned} \quad (2)$$

where

$$s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$\therefore s^2 = x_{n-1}^2 \sigma^2 / n$$

$$\therefore p(s^2) = p(x_{n-1}^2) \left| \frac{d x_{n-1}^2}{d s^2} \right|$$

But

$$p(x_{n-1}^2) = \frac{1}{\Gamma\left(\frac{n-1}{2}\right) 2^{\frac{n-1}{2}}} e^{-\frac{1}{2} x_{n-1}^2} \left(x_{n-1}^2\right)^{\frac{n-3}{2}}$$

$$\therefore p(s^2) = \frac{1}{\Gamma\left(\frac{n-1}{2}\right) 2^{\frac{n-1}{2} \sigma^2}} e^{-\frac{ns^2}{2\sigma^2}} (s^2)^{\frac{n-3}{2}} \quad (3)$$

To get p(s)

$$\text{Let } y = s = \sqrt{s^2} \quad \text{i.e. } y^2 = s^2$$

$$\therefore p(s) = \frac{2}{\Gamma\left(\frac{n-1}{2}\right) 2^{\frac{n-1}{2} \sigma^2}} e^{-\frac{ns^2}{2\sigma^2}} s^{n-2} \quad (4)$$

To prove that \bar{x} & s^2 are independent

Normal variates x_1, x_2, \dots, x_n , each has expected value \bar{x} and s.d. σ .

$$\therefore p(x_1, x_2, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{1}{2\sigma^2} \sum (x_i - \bar{x})^2}$$

But

$$\sum_{i=1}^n (x_i - \bar{x})^2 = ns^2 + n(\bar{x} - \bar{x})^2$$

$$\bar{x} = \frac{\sum x}{n}, \quad s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$\therefore p(x_1, x_2, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{ns^2}{2\sigma^2}} e^{-\frac{n}{2\sigma^2} (\bar{x} - \bar{x})^2}$$

To use a transformation

$$x_i = \bar{x} + a_{i2}y_2 + a_{i3}y_3 + \dots + a_{in}y_n \quad (i=1, 2, \dots, n)$$

where $\sum a_{ij} = 0$, $\sum a_{ij}a_{ik} = 0$, $\sum a_{ik}^2 = 1$

$$j = 2, 3, \dots, n \text{ \& } j \neq k$$

$$x_1 = \bar{x} + \frac{1}{\sqrt{1.2}} y_2 + \frac{1}{\sqrt{2.3}} y_3 + \dots + \frac{1}{\sqrt{(n-1)n}} y_n$$

$$x_2 = \bar{x} - \frac{1}{\sqrt{1.2}} y_2 + \frac{1}{\sqrt{2.3}} y_3 + \dots + \frac{1}{\sqrt{(n-1)n}} y_n$$

$$x_3 = \bar{x} - \frac{2}{\sqrt{2.3}} y_3 + \dots + \frac{1}{\sqrt{(n-1)n}} y_n$$

$$\vdots$$

$$x_{n-1} = \bar{x} - \frac{n-2}{\sqrt{(n-2)(n-1)}} y_{n-1} + \frac{1}{\sqrt{(n-1)n}} y_n$$

$$x_n = \bar{x} - \frac{n-1}{\sqrt{(n-1)n}} y_n$$

$$\therefore \sum x_i = n \bar{x}$$

$$\sum (x_i - \bar{x})^2 = \sum_{j=2}^n y_j^2 \quad \text{i.e. } n s^2 = \sum_{j=2}^n y_j^2$$

$$\begin{aligned} \therefore p(\bar{x}, y_2, y_3, \dots, y_n) &= p(x_1, x_2, \dots, x_n) \quad |J| \\ &= c \cdot e^{-\frac{1}{2\sigma^2} \sum_{j=2}^n y_j^2} \cdot e^{-\frac{n}{2\sigma^2} (\bar{x} - \bar{\mu})^2} \end{aligned}$$

$$\begin{aligned} \therefore p(\bar{x}) &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} p(\bar{x}, y_2, y_3, \dots, y_n) dy_2 dy_3 \dots dy_n \\ &= c_1 \cdot e^{-\frac{n(\bar{x} - \bar{\mu})^2}{2\sigma^2}} \end{aligned}$$

$$p(\bar{x}) = \frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{n(\bar{x} - \bar{\mu})^2}{2\sigma^2}}$$

Also

$$p(y_2, y_1, \dots, y_n) = \int_{-\infty}^{\infty} p(\bar{x}, y_2, y_3, \dots, y_n) d\bar{x}$$

$$= c_2 \cdot e^{-\frac{1}{2\sigma^2} \sum_{j=2}^n y_j^2}$$

$$p(y_i) = c_3 \cdot e^{-\frac{1}{2\sigma^2} y_j^2}$$

$$\text{i.e. } y_j \sim N(0, \sigma)$$

$\therefore \bar{x}, y_2, y_3, \dots, y_n$ are mutually independent set of random variables

$\therefore \bar{x}, s^2 = \sum_{j=2}^n y_j^2$ are mutually independent

(3) The sampling distribution of "t"

Given a normal population with mean \bar{y} and variance σ^2 , the sampling distribution of the mean for independent random samples drawn from such a population is normal with mean \bar{y} & variance $\frac{\sigma^2}{n}$

$$\text{i.e. } \bar{x} \sim N\left(\bar{y}, \frac{\sigma^2}{n}\right)$$

(\sim is distributed as)

$$\therefore \frac{\bar{x} - \bar{y}}{\sigma/\sqrt{n}} \sim N(0, 1)$$

i.e. the test statistic $M = \frac{\bar{x} - \bar{y}}{\sigma/\sqrt{n}}$ is distributed as unit normal distribution. Usually σ for the population is unknown. An estimate for σ based on sample values can be derived:-

$$\therefore \frac{n s^2}{\sigma^2} = \chi_{n-1}^2$$

$$\therefore E\left(\frac{n s^2}{\sigma^2}\right) = n - 1$$

$$\therefore E\left(\frac{n}{n-1} s^2\right) = \sigma^2$$

$$\therefore \hat{\sigma}^2 = \frac{n}{n-1} s^2$$

$$\text{i.e. } \hat{\sigma} = \frac{\sqrt{n}}{\sqrt{n-1}} s$$

∴ The test statistic becomes

$$t = \frac{\bar{x} - \bar{y}}{\sigma / \sqrt{n}} = \frac{\bar{x} - \bar{y}}{s / \sqrt{n-1}} = \frac{\frac{\bar{x} - \bar{y}}{\sigma / \sqrt{n}}}{\chi_{n-1} / \sqrt{n-1}}$$

$$= \frac{\mu}{\chi_{n-1} / \sqrt{n-1}} = \frac{\mu}{\chi_y / \sqrt{y}} \quad (1)$$

Also

Given two normal populations $N(\bar{x}_1, \sigma_1^2)$ & $N(\bar{x}_2, \sigma_2^2)$ and two independent random samples of size n_1, n_2 are drawn. Then the sampling distribution of the difference between the two means i.e. $\bar{x}_1 - \bar{x}_2$ is normally distributed with mean $\bar{x}_1 - \bar{x}_2$ and variance $\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$.

Let the two populations have common variance $\sigma^2 = \sigma_1^2 = \sigma_2^2$ which is usually unknown. An estimate for σ^2 based on the two sample observations can be derived.

$$n_1 s_1^2 = \chi_{n_1-1}^2 \sigma^2$$

$$n_2 s_2^2 = \chi_{n_2-1}^2 \sigma^2$$

$$n_1 s_1^2 + n_2 s_2^2 = (\chi_{n_1-1}^2 + \chi_{n_2-1}^2) \sigma^2$$

$$= \chi_{n_1+n_2-2}^2 \sigma^2$$

$$\therefore E(n_1 s_1^2 + n_2 s_2^2) = (n_1 + n_2 - 2) \sigma^2$$

$$\therefore E\left(\frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2}\right) = \sigma^2$$

$$\therefore \hat{\sigma}^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2} = \frac{\chi_{n_1+n_2-2}^2 \sigma^2}{n_1+n_2-2}$$

But

$$u = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \text{ is distributed as unit normal.}$$

Since $\sigma^2 = \sigma_1^2 = \sigma_2^2$ and σ^2 is unknown, then the statistic becomes

$$\begin{aligned} t &= \frac{\bar{x}_1 - \bar{x}_2}{\hat{\sigma} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{\frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}}{\chi_{n_1+n_2-2} / \sqrt{n_1+n_2-2}} \\ &= \frac{u}{\chi_y / \sqrt{y}} \end{aligned} \quad (2)$$

To get $p(t)$

$$p(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} u^2}$$

&

$$p(\chi_y) = \frac{1}{2^{\frac{y}{2}-1} \Gamma(\frac{y}{2})} \chi^{y-1} e^{-\frac{1}{2} \chi^2}$$

$\therefore U$ & χ_y are independent

$$\therefore p(u, \chi_y) = p(u) p(\chi_y) = c \chi_y^{y-1} e^{-\frac{1}{2}(u^2 + \chi^2)}$$

Make the transformation

$$t = \frac{u}{\chi_y / \sqrt{y}}, \quad y = \chi_y$$

$$\text{i.e. } \chi_y = y, \quad u = ty / \sqrt{y}$$

$$\therefore \left| \frac{\partial(x, u)}{\partial(t, y)} \right| = y/\sqrt{y}$$

$$\therefore p(t, y) = c_1 y^\gamma e^{-\frac{1}{2} y^2 (1 + \frac{t^2}{y})}$$

$$\therefore p(t) = c_1 \int_0^\infty y^\gamma e^{-\frac{1}{2} y^2 (1 + \frac{t^2}{y})} dy \quad (y > 0)$$

write $z = \frac{1}{2} y^2 (1 + \frac{t^2}{y})$

$$\therefore p(t) = \frac{c_2}{(1 + \frac{t^2}{y})^{\frac{1}{2}(\gamma+1)}} \Gamma\left(\frac{\gamma+1}{2}\right)$$

where

$$c_2 = \frac{1}{\Gamma(\frac{1}{2}) \Gamma(\frac{\gamma}{2}) \sqrt{y}}$$

$$\therefore p(t) = \frac{1}{\sqrt{y} B(\frac{1}{2}, \frac{\gamma}{2})} (1 + \frac{t^2}{y})^{-\frac{1}{2}(\gamma+1)}$$

As $\gamma \rightarrow \infty$, $p(t, y) \rightarrow N(0, 1)$.

(4) The sampling distribution of F

Given two normal populations with variances σ_1^2, σ_2^2 respectively. (Usually σ_1^2, σ_2^2 are unknown)

Two independent random samples of size n_1, n_2 are drawn. If their variances are s_1^2, s_2^2 then

$$\hat{\sigma}_1^2 = \frac{n_1 s_1^2}{n_1 - 1}$$

$$\hat{\sigma}_2^2 = \frac{n_2 s_2^2}{n_2 - 1}$$

$$\therefore \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} = \frac{n_1 s_1^2 / (n_1 - 1)}{n_2 s_2^2 / (n_2 - 1)}$$

Also

$$n_1 s_1^2 = x_{n_1-1}^2 \sigma_1^2$$

$$n_2 s_2^2 = x_{n_2-1}^2 \sigma_2^2$$

and if $\sigma_1^2 = \sigma_2^2$, then

$$F = \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} = \frac{x_{n_1-1}^2 / (n_1-1)}{x_{n_2-1}^2 / (n_2-1)}$$

$$\therefore F_{\gamma_1, \gamma_2} = \frac{x_{\gamma_1}^2 / \gamma_1}{x_{\gamma_2}^2 / \gamma_2}$$

$\therefore x_{\gamma_1}^2, x_{\gamma_2}^2$ are independent

$$\therefore p(x_1^2, x_2^2) = p(x_1^2) p(x_2^2)$$

$$= c(x_{\gamma_1}^2)^{\frac{1}{2}(\gamma_1-2)} (x_{\gamma_2}^2)^{\frac{1}{2}(\gamma_2-2)} e^{-\frac{1}{2}(x_{\gamma_1}^2 + x_{\gamma_2}^2)}$$

where

$$c = \frac{1}{2^{\frac{1}{2}(\gamma_1+\gamma_2)} \Gamma(\frac{\gamma_1}{2}) \Gamma(\frac{\gamma_2}{2})}$$

Make the transformation

$$F = \frac{\gamma_2}{\gamma_1} \frac{x_{\gamma_1}^2}{x_{\gamma_2}^2}, \quad y = x_{\gamma_2}^2$$

$$\therefore \left| \frac{\partial (x_1^2, x_2^2)}{\partial (F, y)} \right| = \frac{\gamma_1}{\gamma_2} y \quad (85)$$

$$\therefore p(F, y) = C_1 F^{\frac{\gamma_1 - 2}{2}} y^{\frac{\gamma_1}{2} + \frac{\gamma_2}{2} - 1} e^{-y(1 + \frac{\gamma_1}{\gamma_2} F)}$$

$$\therefore p(F) = \int_0^{\infty} p(F, y) dy$$

$$\therefore p(F) = \frac{\gamma_1^{\frac{1}{2}\gamma_1} \gamma_2^{\frac{1}{2}\gamma_2}}{B(\frac{\gamma_1}{2}, \frac{\gamma_2}{2})} \frac{F^{\frac{1}{2}\gamma_1 - 1}}{(\gamma_2 + \gamma_1 F)^{\frac{1}{2}(\gamma_1 + \gamma_2)}}$$

$$0 < F < \infty$$

The moments

$$\mu'_F = C \int_0^{\infty} F^{\frac{1}{2}\gamma_1 + r - 1} (\gamma_2 + \gamma_1 F)^{-\frac{1}{2}(\gamma_1 + \gamma_2)} dF$$

where

$$C = \frac{\gamma_1^{\frac{1}{2}\gamma_1} \gamma_2^{\frac{1}{2}\gamma_2}}{B(\frac{\gamma_1}{2}, \frac{\gamma_2}{2})}$$

write $z = \frac{\gamma_1 F}{\gamma_2 + \gamma_1 F}$

$$\therefore 1 - z = \frac{\gamma_2}{\gamma_2 + \gamma_1 F}$$

$$\therefore F = \frac{\gamma_2}{\gamma_1} \cdot \frac{Z}{1-Z}$$

$$\therefore \mu'_r = \left(\frac{\gamma_2}{\gamma_1}\right)^r \frac{B\left(\frac{\gamma_1}{2} + r, \frac{\gamma_2}{2} - r\right)}{B\left(\frac{\gamma_1}{2}, \frac{\gamma_2}{2}\right)}$$

Exercises

(1) Prove that

$$\Pr\left\{t_{\gamma} > t_0\right\} = \frac{1}{2} \left\{1 - I_k\left(\frac{1}{2}, \frac{1}{2}\gamma\right)\right\}$$

where

$$k = \frac{t_0^2}{1 + t_0^2/\gamma}$$

(2) Find β_1 & β_2 for F-distribution

(3) Show that

$$p(t_{\gamma}^2) = p(F_{1,\gamma})$$

(VIII) Testing statistical hypothesis
and Confidence interval

Testing Statistical Hypothesis
and Confidence Limits.

1) Testing statistical hypothesis around the mean

(i) Large samples

Theorem:- The sampling distribution of the means of independent random samples drawn from a normal population $N(\mu, \sigma^2)$ is also normal with expected value μ and variance σ^2/n .

Cor.(1) Usually the population is unknown and hence the above theorem can be stated again as:

The sampling distribution of the means of independent large samples drawn at random from an infinite population (μ, σ^2) approaches a normal distribution with expected value μ & variance σ^2/n as n increases.

Cor.(2) If σ is unknown, then the sampling distribution of the means of independent large samples drawn at random from infinite population (μ, σ^2) approaches a normal distribution with expected value μ & variance $s^2/(n-1)$ as n increases.

(Remember $\hat{\sigma}^2 = \frac{n}{n-1} s^2$).

Hence to test the null hypothesis $\mu = \mu_0$ against $\mu \neq \mu_0$, we use the test criterion

$$u = \frac{\bar{x} - \mu_0}{\sigma_{\bar{x}}}$$

where $\sigma_{\bar{x}}$ is the standard error of the mean and

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad \text{if } \sigma \text{ is known}$$

$$\text{or} \quad = \frac{s}{\sqrt{n-1}} \quad \text{if } \sigma \text{ " unknown}$$

The above statistic is distributed as unit normal $N(0,1)$ and hence we compare u with u_{α} looked in the tables at $\alpha\%$ level of significance

(ii) Small samples

To test the null hypothesis $\mu = \mu_0$ against $\mu \neq \mu_0$ we use the test criterion

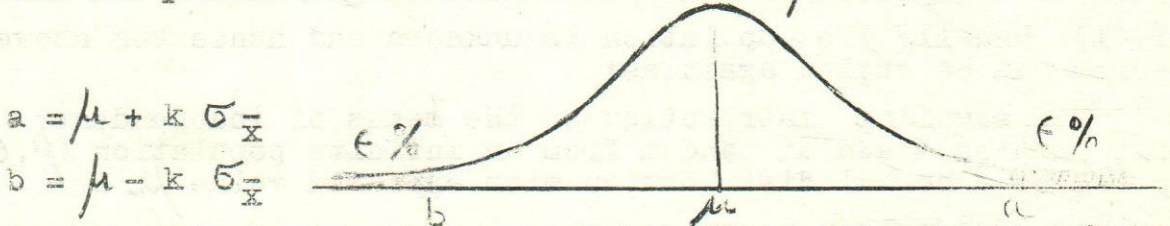
$$t = \frac{\bar{x} - \mu_0}{\sigma_{\bar{x}}}$$

where $\tilde{\sigma}_{\bar{x}} = s/\sqrt{n-1}$. This statistic is distributed as t-distribution with $n-1 = \gamma$ degrees of freedom. Hence we compare t with $t_{\gamma, \alpha}$ at $\alpha\%$ level of significance

2) Confidence limits for the population mean

It is required to find two values between which we expect that the population mean will lie, with confidence coefficient $\gamma = (100 - 2\epsilon)\%$

Consider the sampling distribution of the mean and that a, b are two points on the sides of the mean μ where



$$a = \mu + k \sigma_{\bar{x}}$$

$$b = \mu - k \sigma_{\bar{x}}$$

k is determined such that area below the curve between a, b is $(100 - 2\epsilon)\%$. Then the probability that \bar{x} lies between a, b is given as

$$(1) \Pr \{ \mu - k \sigma_{\bar{x}} \leq \bar{x} \leq \mu + k \sigma_{\bar{x}} \} = (100 - 2\epsilon)\%$$

$$\text{Now } \therefore \mu - k \sigma_{\bar{x}} \leq \bar{x}, \therefore \mu \leq \bar{x} + k \sigma_{\bar{x}}$$

$$\therefore \mu + k \sigma_{\bar{x}} \geq \bar{x}, \therefore \mu \geq \bar{x} - k \sigma_{\bar{x}}$$

Hence $\mu - k \sigma_{\bar{x}} \leq \bar{x} \leq \mu + k \sigma_{\bar{x}}$ is equivalent to $\bar{x} - k \sigma_{\bar{x}} \leq \mu \leq \bar{x} + k \sigma_{\bar{x}}$

Consequently

$$(2) \Pr \{ \bar{x} - k \sigma_{\bar{x}} \leq \mu \leq \bar{x} + k \sigma_{\bar{x}} \} = (100 - 2\epsilon)\%$$

\therefore The upper Confidence limit for μ is

$$\bar{x} + k \sigma_{\bar{x}}$$

and the lower Confidence limit is

$$\bar{x} - k \sigma_{\bar{x}}$$

In large samples, $k = u_{\alpha}$ where

$u_{\alpha} = 1.96$ at 95 % confidence coefficient

& $u_{\alpha} = 2.58$ at 99 % confidence coefficient

$\sigma_{\bar{x}} = \sigma/\sqrt{n}$ or $s/\sqrt{n-1}$ according to as σ is known or unknown.

In small samples, $k = t_{\alpha, \gamma}$ and we need to look in tables of t at $\gamma = n-1$ degrees of freedom and at $\alpha = 100 - \gamma$ level of significance. But $\tilde{\sigma}_{\bar{x}} = s/\sqrt{n-1}$ where $s^2 = \sum (x_i - \bar{x})^2/n$.

Example

A standardized intelligence examination has been given for several years with an average score of 80 and a standard deviation 7. A group of 36 students are taught with special emphasis on reading skill. If the 36 students obtained a mean grade of 83 on this examination, is there reason to believe that the special emphasis changes the results on the test?

- 1- The null hypothesis is $\mu_0 = 80$
 The alternative " " $\mu_0 \neq 80$

2- Test criterion

$$u = z = \frac{-\mu_0}{\tilde{\sigma}_{\bar{x}}} = \frac{83 - 80}{7/\sqrt{35}} = 2.53$$

3- Comparing this $|u| = 2.53$ with $|u_{1\%}| = 2.58$ we accept the null hypothesis. In other words, there is no evidence that special emphasis on reading skill changes the results on the test.

Example:

A fertilizer mixing machine is set to give 10 pounds of nitrate for every 100 pounds of fertilizer. Ten bags (each of 100 pounds) are examined. The percentages of nitrate are as follows:

9, 12, 11, 10, 11, 9, 11, 12, 9, 10

Is there reason to believe that the mean is not equal to 10%.

$$\text{Mean } (\bar{x}) = 10.4$$

$$\text{Variance } (s^2) = 1.24$$

$$\tilde{\sigma}_{\bar{x}} = \frac{s}{\sqrt{n-1}} = 0.37$$

Null hypothesis $\mu_0 = 10$

alternative hypothesis $\mu_0 \neq 10$

The test criterion is

$$t = \frac{\bar{x} - \mu_0}{\tilde{\sigma}_{\bar{x}}} = \frac{10.4 - 10}{0.37} = 1.11$$

From the tables $t_{9, 1\%} = 3.25$

∴ We accept the null hypothesis

There is an evidence that the population mean is 10.

3) Testing statistical hypothesis around the difference between two means

3-1) Independent samples

(i) Large samples: Given two independent random samples. The 1st sample is of size n_1 , has a mean \bar{x}_1 , a standard deviation s_1 and is drawn from a normal population $N(\mu_1, \sigma_1^2)$. The 2nd sample is of size n_2 , has a mean \bar{x}_2 , a standard deviation s_2 and is drawn from a normal population $N(\mu_2, \sigma_2^2)$. Then the sampling distribution of $\bar{x}_1 - \bar{x}_2$ is exactly normal with expected value $\mu_1 - \mu_2$ and standard deviation

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

If the populations are infinite and the samples are large, then the sampling distribution of $\bar{x}_1 - \bar{x}_2$ approaches normal with mean $\mu_1 - \mu_2$ and standard deviation

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Usually σ_1^2 & σ_2^2 are unknown and hence the sampling distribution of $\bar{x}_1 - \bar{x}_2$ approaches normal with mean $\mu_1 - \mu_2$ and standard deviation

$$\sqrt{\frac{s_1^2}{n_1-1} + \frac{s_2^2}{n_2-1}}$$

Now to test the null hypothesis $\mu_1 - \mu_2 = 0$ against the alternative hypothesis $\mu_1 - \mu_2 \neq 0$ we use the test criterion

$$u = \frac{\bar{x}_1 - \bar{x}_2}{\sigma_{\bar{x}_1 - \bar{x}_2}} \sim N(0,1)$$

where

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

$$\text{or } \tilde{\sigma}_{\bar{x}_1 - \bar{x}_2}^2 = \frac{s_1^2}{n_1 - 1} + \frac{s_2^2}{n_2 - 1}$$

according to σ_1^2, σ_2^2 are known or unknown. Hence we compare u with u_α at $\alpha\%$ level of significance

(ii) Small samples:- In large samples we know that the standard error of $\bar{x}_1 - \bar{x}_2$, i.e.

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

Assuming that the two populations have the same variance i.e.

$$\sigma_1^2 = \sigma_2^2 = \sigma^2$$

$$\therefore \sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

But σ^2 is usually unknown and can be replaced by its estimate $\tilde{\sigma}^2$ then

$$\tilde{\sigma}_{\bar{x}_1 - \bar{x}_2} = \sqrt{\tilde{\sigma}^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

where

$$\tilde{\sigma}^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2} = \frac{\sum (x_1 - \bar{x}_1)^2 + \sum (x_2 - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

Hence, the test criterion

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\tilde{\sigma}_{\bar{x}_1 - \bar{x}_2}}$$

is distributed as t-distribution with $\nu = n_1 + n_2 - 2$ degrees of freedom. The calculated t will be compared with $t_{\nu, \alpha}$ at $\alpha\%$ level of significance.

(13-2) Related samples

In some problems we have to take one random sample and call its values-before treatment - the control group and its values-after treatment - the experimental group. For examples, to study the effect of a certain diet on the weight of Children, the effect of a certain training course on the achievement of students and the effect of a certain medicine on the recovery from a given disease.

Let us denote the observations for the control group by x_c and those for the experimental group by x_e . Then the differences will be $d = x_e - x_c$.

Now the null hypothesis is $D = 0$
the alternative " is $D \neq 0$

The test criterion

$$\text{for large samples is } u = \frac{\bar{d}}{s_d/\sqrt{n-1}}$$

$$\& \text{ for small samples is } t = \frac{\bar{d}}{s_d/\sqrt{n-1}}$$

with $V = n-1$ degrees of freedom

Example Two independent random samples are drawn from large populations A, B respectively. The following information is known

	n	\bar{x}	s
<u>Sample A</u>	122	105	33
<u>Sample B</u>	50	115	28

Test that the two populations have identical means.

$$\sigma_{\bar{x}_1 - \bar{x}_2}^2 = \frac{s_1^2}{n_1 - 1} + \frac{s_2^2}{n_2 - 1} = \frac{33^2}{121} + \frac{28^2}{49} = 25$$

$$\therefore \sigma_{\bar{x}_1 - \bar{x}_2} = 5$$

The null hypothesis is

$$\mu_1 - \mu_2 = 0$$

The alternative " is

$$\mu_1 - \mu_2 \neq 0$$

The test criterion is

$$u = \frac{\bar{x}_1 - \bar{x}_2}{\sigma_{\bar{x}_1 - \bar{x}_2}} = \frac{115 - 105}{5} = 2$$

Comparing with $|u_{5\%}| = 1.96$, we reject the null hypothesis.

In other words, there is no evidence that the two samples are drawn from populations having the same mean.

Example:- A certain stimulus is to be tested for its effect on blood pressure. Twelve men have their blood pressure measured before and after the stimulus. The results are shown below:

Men:-	1	2	3	4	5	6	7	8	9	10	11	12
Before:	120	124	130	118	140	128	140	135	126	130	126	127
After:	128	131	131	127	132	125	141	137	118	132	129	135

Is there reason to believe that the stimulus world on the average raise the blood pressure.

The differences d_i are :-

8, 7, 1, 9, -8, -3, 1, 2, -8, 2, 3, 8

$$\bar{d} = \frac{22}{12} = 1.83$$

$$s_d^2 = \frac{434}{12} - \left(\frac{11}{6}\right)^2 = 32.8056$$

$$s_d = 5.18$$

$$t = \frac{\bar{d}}{s_d / \sqrt{n-1}} = \frac{1.83 \times 11}{5.18} = 1.17$$

From the table $t_{11,5\%} = 2.20$

\therefore We accept the null hypothesis that $\mu_1 - \mu_2 = 0$. In other words there is no evidence that the stimulus would on the average raise the blood pressure.

Exercise The following table shows the distribution of households according to the annual consumption expenditure of two random samples drawn from the rural and urban parts of the country.

Intervals	Frequencies	
	Urban	Rural
15-	6	43
25-	66	280
50-	167	474
75-	261	453
100-	621	760
150-	520	446
200-	367	238
250-	262	127
300-	373	126
400-	295	68
600-	125	15
800-	31	3
1000-1400	51	4
Total	3145	3037

Is the difference between the averages of annual consumption expenditure in rural and urban parts significant?

Exercise :- We have the weights (\bar{x}) of 10 children before eating a certain food and their weights (\bar{y}) after eating this food.

Children:	1	2	3	4	5	6	7	8	9	10
x :	142	140	143	158	149	140	134	124	116	157
y :	146	139	148	161	151	138	135	127	115	162

Is there reason to believe that the food would on the average increase the weights of the Children?

-4) Testing the hypothesis around the proportion and the difference between two proportions

(6-4.1) If we draw a large sample of size n from an infinite dichotomus population and if the proportion of successes is p , then, to test the null hypothesis that the sample is drawn from a population where the proportion is $\emptyset = \emptyset_0$ against $\emptyset \neq \emptyset_0$ we use the test criterion

$$u = \frac{p - \emptyset_0}{\tilde{\sigma}_p}$$

where

$$\tilde{\sigma}_p = \sqrt{\frac{\emptyset_0(1-\emptyset_0)}{n}}$$

Example: A superintendent of schools has stated that at least 60% of high school seniors expect to attend college. In a random sample of 200 cases only 96 say they are planning for college. Does this refute the superintendent's statement?

$$\tilde{\sigma}_p = \frac{0.6 \times 0.4}{200} = 0.1096$$

The null hypothesis is $\emptyset_0 = 0.6$

The alternative " " $\emptyset_0 \neq 0.6$

$$p = \frac{96}{200} = 0.48$$

The test criterion is

$$u = \frac{0.48 - 0.60}{0.1096} = 1.09$$

Comparing this with $|u_{5\%}| = 1.96$, we accept the null hypothesis. In other words there is an evidence that the superintendent's statement is accepted

(-4.2) A large random sample of size n_1 with proportion of successes p_1 , is drawn from an infinite dichotomus population where the proportion \emptyset_1 is unknown. Another large random sample of size n_2 with proportion of successes p_2 is drawn from another infinite dichotomus population where the proportion \emptyset_2 is unknown. It is required to test the null hypothesis $\emptyset_1 = \emptyset_2 = \emptyset$ against the alternative $\emptyset_1 \neq \emptyset_2$.

Since \emptyset is unknown (common proportion), an estimate based on the sample values is given as $\hat{\emptyset}$ where

$$\hat{\emptyset} = (n_1 p_1 + n_2 p_2) / (n_1 + n_2)$$

Now, if the null hypothesis is true, the test criterion is

$$u = \frac{p_1 - p_2}{\tilde{\sigma}_{p_1 - p_2}}$$

where

$$\tilde{\sigma}_{p_1 - p_2} = \sqrt{\hat{\emptyset}(1 - \hat{\emptyset}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

Example In an experiment, it is found that 56 out of 6815 inoculated persons were attacked by a certain disease while 272 out of 11668 not inoculated persons were attacked by the same disease. Test the significance of the difference between the proportions of the attacked persons in the two populations.

$$p_1 = 56/6815 = 0.0082 = 0.008$$

$$p_2 = 272/11668 = 0.0233 = 0.023$$

$$\emptyset = \frac{56 + 272}{18483} = 0.0177 = 0.018$$

$$\tilde{\sigma}_{p_1 - p_2}^2 = 0.18(0.82) \left(\frac{1}{6815} + \frac{1}{11668} \right) = 0.000034$$

$$\tilde{\sigma}_{p_1 - p_2} = 0.0059$$

$$u = \frac{0.023 - 0.008}{0.0059} = \frac{.015}{.0059} = 2.54$$

Comparing with $|u_{5\%}| = 1.96$, we reject the null hypothesis. In other words there is an evidence that the difference between the proportions of dichotomus in the two populations are significant

Exercise:- A random sample of size 450 is drawn from an infinite dichotomus population and we find that 50 of them were attacked by a certain disease. Test the null hypothesis that the proportion of attacked persons in the population is 0.5.?

Exercise:- A random sample of 40 individuals is drawn from an infinite population and we find that in the sample 10 individuals were unemployed. Another random sample of 30 individuals is drawn from another infinite population and we find that in this sample 10 were unemployed. Test the hypothesis that the proportions of unemployment in the two populations are the same?

-5) Testing the hypothesis around the population variance

(6-5.1) Given a random variable x_i having a normal distribution with mean μ and standard deviation σ , then the standardized variate

$$u_i = \frac{x_i - \mu}{\sigma} \quad (1)$$

will have a normal distribution with mean 0 and unit standard deviation. Then by definition $\sum_{i=1}^n u_i^2$ will be distributed as χ^2 with n degrees of freedom i.e.

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} = \chi_n^2 \quad (2)$$

If we replace μ by \bar{x} , then

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \chi_{n-1}^2 \quad (3)$$

$$\text{i.e.} \quad \frac{ns^2}{\sigma^2} = \chi_{n-1}^2 \quad (4)$$

(4) is used as a test for the population variance, where s^2 is the sample variance.

5.2) we can also determine the confidence limits between which we expect the variance of the normal population will lie with a given confidence coefficient $\gamma = 100(1-\alpha)\%$ where α is the level of significance.

The lower limit $\underline{\sigma^2}$ is determined from

$$\chi_{\text{upper}}^2 = \frac{ns^2}{\underline{\sigma^2}}$$

The upper limit $\bar{\sigma^2}$ is determined from

$$\chi_{\text{lower}}^2 = \frac{ns^2}{\bar{\sigma^2}}$$

Example (1) A random sample of size 5 and variance 40 is drawn from a normal population. Test the hypothesis that the population variance is 25.

The null hypothesis is $\sigma^2 = 25$

The alternative " " $\sigma^2 \neq 25$.

$$\therefore \chi^2 = \frac{ns^2}{\sigma^2} = \frac{5 \times 40}{25} = 8$$

From the tables $\chi^2_{0.975, 4} = 0.484$

$$\chi^2_{0.025, 4} = 11.143$$

\therefore The obtained value of χ^2 lies between the two tabulated values. Hence we accept the null hypothesis. In other words there is an evidence that the population variance is 25.

Example (2) A random sample of size 10 and variance 30, is drawn from a normal population, determine the 99% confidence limits of the population variance.

$$\chi^2_{0.995, 9} = 1.735$$

$$\chi^2_{0.005, 9} = 23.589.$$

The lower limit $\underline{\sigma^2} = 10 \times 30 / 23.589 = 12.72$

The upper limit $\overline{\sigma^2} = 10 \times 30 / 1.735 = 172.91$

Remark When the sample size is large, we know that the sampling distribution of the statistic

$$x = \sqrt{2\chi^2} - \sqrt{2\gamma - 1}$$

- where γ is the no. of d.f. - is very near to a normal curve $N(0,1)$.

Example (3) A random sample of size 146 & $s^2 = 197.26$ is drawn from a normal population. Test the hypothesis that population variance is 256.

$$\chi^2 = ns^2 / \sigma^2 = \frac{146 \times 197.26}{256} = 112.5$$

$$\therefore x = -2$$

But at 5% level of significance $|u| = 1.96$.

\therefore We reject the null hypothesis

Exercise (1) A random sample of 65 individuals is drawn from an infinite population. If the mean of the sample is 95 and its standard deviation is 16. Can we say that this sample is drawn from a population with mean equals to 100?

Exercise (2) A certain type of rat shows a mean gain in weight of 65 grams during the first 3 months of life. Twelve rats were fed a particular diet from birth until age 3 months and the following weight gains were observed:

55, 62, 54, 58, 65, 64, 60, 62, 59, 67, 62, 61.

Is there reason to believe at the 5% level of significance that the diet causes a change in the amount of weight gained?

Exercise (3) Two astronomers recorded observations on a certain star. The 120 observations obtained by the first astronomer have a mean reading of 1.20. The 80 observations obtained by the second astronomer have a mean 1.15. Past experience has indicated that these astronomers obtain readings with a variance of about 0.40. Does the difference between the two results seem reasonable?

Exercise (4) Two new types of rations are fed to pigs. It is desired to test whether one or the other of the types is better. A sample of 12 pigs is fed type (A) ration, and another sample of 12 pigs is fed type (B) ration. The gains in weight are recorded below

Type A	31	34	29	26	32	35	38	34	30	29	32	31
Type B	26	24	28	29	30	29	32	26	31	29	32	28

Exercise (5) A random sample of 10 observations

11, 11, 10, 12, 8, 10, 15, 8, 10, 8, is drawn from a normal population. Test that the population variance is 2.5.

Chi-Square & its application

In the field of research, the experimenter may face one of the following situations:-

1- The experimenter may believe that the students should be distributed equally regarding the preference of one of the three different subjects A,B,C. He draws a sample of 90 students and asking them which subject they mostly prefer, he get the following result

29	students	mostly	prefer	subject	A
28	"	"	"	"	B
33	"	"	"	"	C

According to what he believes, the expected number of students who will prefer any of the subjects should be 30.

The data can be represented in the following table

Subjects	observed frequencies	Expected frequencies
A	29	30
B	28	30
C	33	30

Are the observed frequencies Consistent with the expected frequencies? The null hypothesis to be tested is "there is consistency".

2- In fitting a certain theoretical model to a given observed frequency distribution, the experimenter may wish to test whether the model is a good fit or not. In other words, it is required to test the goodness of fit. For example, given the distribution of students according to their scores, the experimenter would like to test that this distribution will follow a normal curve. The null hypothesis to be tested is "there is a good fit".

3- Sometimes the experimenter may be faced with the data arranged in a bivariate tables where the two variates are graded into categories e.g.

Marital Educational Status Status	Single	married	divorced	widow
illiterate				
read & write				
medium culificate				
higher "				

Such tables are called "Association tables" or "contingency tables"

The hypothesis to be tested is "there is no association" or "there is independence"

Now, in the above situations, we have:

Observed frequencies denoted by f_o
 Expected frequencies denoted by f_e

The test criterion is

$$T = \sum \frac{(f_o - f_e)^2}{f_e} \quad (1)$$

This test statistic is a discrete variable and has a sampling distribution which was proved to be very near to that of the continuous variate χ^2 defined as

$$\chi^2 = \sum u^2$$

where u is a standardized normal variate. Hence we refer to T as χ^2 i.e.

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} \quad (2)$$

and we use the corresponding tables for χ^2 defined by (2)

Note f_e in any cell should not be less than 5

Applications

Example (1) In a questionnaire distributed to 200 students it is asked, "which subject A, B, C, or D do you prefer?" It is found that 45 prefer A, 54 prefer B, 46 prefer C and 55 prefer D. Test the hypothesis that there is consistency with what is believed that the four subjects are equally preferable.

Now the data can be represented in the following table

Subjects	f_o	f_e
A	45	50
B	54	50
C	46	50
D	55	50

$$\begin{aligned}
 \therefore \chi^2 &= \sum \frac{(f_o - f_e)^2}{f_e} \\
 &= \frac{(45-50)^2}{50} + \frac{(54-50)^2}{50} + \frac{(46-50)^2}{50} + \frac{(55-50)^2}{50} \\
 &= \frac{82}{50} = 1.64
 \end{aligned}$$

From the tables

$$\chi^2_{3,0.05} = 7.82$$

\therefore We accept the null hypothesis. In other words, there is an evidence of existence of consistency.

Example (2) A random sample of 2970 students is drawn from the 3rd grade secondary schools. They are given an achievement test. The observed frequency distribution of their scores is as follows:-

<u>Intervals:</u>	40	40-	50-	60-	70-	80-	90-
f :	1	4	7	48	130	338	592

<u>Intervals</u> :	100-	110-	120-	130-	140-	150-	160-	200
f :	719	610	330	120	55	13		

Do we consider this sample as drawn from a normal population having the same mean and standard deviation?

The mean (\bar{x}) = 104.56

The standard deviation (s) = 16.99

The following table summerises the results of fitting a normal

Upper limits	f_o	Standardized Values	Area below	Δ	f_e
39.5	1	- 3.83	0.0001	0.0001	0.30
49.5	4	- 3.24	0.0004	0.0005	1.49
59.5	7	- 2.65	0.0040	0.0036	10.10
69.5	48	- 2.06	0.0197	0.0157	46.63
79.5	140	- 1.47	0.0708	0.0511	151.77
89.5	338	- 0.89	0.1867	0.1159	344.22
99.5	522	- 0.30	0.3821	0.1954	580.34
109.5	719	0.29	0.6141	0.2320	689.04
119.5	610	0.88	0.8106	0.1965	583.61
129.5	330	1.47	0.9292	0.1186	352.24
139.5	120	2.06	0.9803	0.0511	151.77
149.5	55	2.65	0.9960	0.0157	46.63
159.5	13	3.23	0.9994	0.0034	10.10
199.5	3	5.59	1.0000	0.0006	1.78
	2970				2970.02

In order that none of the expected frequencies should be less than 5, we add the 1st three expected frequencies (i.e.: $0.30+1.49+10.10 = 11.89$) and the corresponding observed frequencies (i.e. $1+4+7 = 12$) and consider each as one class. Similarly we add the last two expected frequencies (i.e. $10.10+1.78 = 11.88$) and the corresponding observed frequencies (i.e. $13+3 = 16$). Calculating χ^2 we find that

$$\chi^2 = \frac{\sum (f_o - f_e)^2}{f_e} = 16.97$$

The number of degrees of freedom $\nu = 11-3 = 8$ where 11 is the number of classes and 3 is the number of constraints (the total, the mean and the standard deviation).

At 1% level of significance we find from the table

$$\chi_{8,1\%}^2 = 20.09$$

\therefore We accept the null hypothesis. In other words there is an evidence that the distribution of the scores in the population is normal with mean 104.56 and standard deviation 16.99

Example (3) Five coins were tossed 100 times. The frequency distribution for the number of coins on which the head occurs is given below

x :	0	1	2	3	4	5	Total
f :	2	14	20	34	22	8	100

Fit a binomial distribution and test for the goodness of fit.

$$\bar{x} = 2.84 = np$$

$$5p = 2.84 \quad \therefore p = 0.57$$

The probability binomial distribution is given by

$$p(x) = \binom{n}{x} p^x q^{n-x}$$

In our example

$$p(x) = \binom{5}{x} (0.57)^x (0.43)^{5-x}$$

for $x = 0, 1, 2, 3, 4, 5$.

The table below gives the observed frequencies and the expected frequencies calculated from the above formula.

x	f_o	Probabilities	Expected freq. f_e
0	2	0.0147	1.47
1	14	0.0974	9.744
2	20	0.25832	25.832
3	34	0.34242	34.242
4	22	0.22695	22.695
5	8	0.06017	6.017
	100		100.000

$$\chi^2 = 4.035$$

Degrees of freedom = $5 - 2 = 3$

From the tables $\chi^2_{3,5\%} = 7.815$

\therefore We accept the null hypothesis. In other words, there is an evidence that the distribution of the number of the coins on which the head occurs (in the population) is a binomial distribution with mean 2.84.

Contingency tables: In general the table can be represented as follows where we have "r" rows and "s" columns

A \ B	1	...	j	...	s	
1	f_{11}	...	f_{1j}	...	f_{1s}	$f_{1.}$
⋮	⋮		⋮		⋮	⋮
i	f_{i1}	...	f_{ij}	...	f_{is}	$f_{i.}$
⋮	⋮		⋮		⋮	⋮
r	f_{r1}	...	f_{rj}	...	f_{rs}	⋮
	$f_{.1}$...	$f_{.j}$...	$f_{.s}$	$f_{..}$

f_{ij} denotes the frequency in the (i,j) cell

$f_{i.} = \sum_{j=1}^s f_{ij}$ denotes the frequency in the i^{th} row

$f_{.j} = \sum_{i=1}^r f_{ij}$ denotes the frequency in the j^{th} column

$f_{..} = \sum_i \sum_j f_{ij} = \sum_i f_{i.} = \sum_j f_{.j}$ denotes the total

frequency.

It should be noticed that no association or independence occurs in two cases:--

- (i) if the frequencies in the cells are all equal
or (ii) if the proportions of the corresponding frequencies in any two rows (or columns) do not differ from column to column (or row to row).

The expected frequency f_e for the observed $f_o = f_{ij}$ is calculated as

$$f_e = f_{i.} \cdot f_{.j} / f_{..}$$

$$\text{Now } \sum_i f_{e.} = \sum_i f_{i.} f_{.j} / f_{..} = f_{.j}$$

$$\& \sum_j f_{e.} = \sum_j f_{i.} f_{.j} / f_{..} = f_{i.}$$

This means that the marginal frequencies for the observed and expected frequencies are the same. Consequently the number of degrees of freedom γ is given as

$$\gamma = rs - (s-1) - (r-1) = (r-1)(s-1)$$

Example:- In a questionnaire filled by a random sample of 2026 students, we ask, "Do you like to work in the government or on your own?" The students are also classified into three groups A, B, C according to their father's job. The results are given in the following table:-

		Father's Job			Total
		A	B	C	
Attitude toward working in the government	In the government	120 (58.9)	208 (242.0)	13 (40.1)	341
	On their own	230 (291.1)	1230 (1196.0)	225 (197.9)	1685
Total		350	1438	238	2026

The expected frequencies are given between brackets

$$\text{Now } \chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} = 103.904$$

The number of degrees of freedom is $\gamma = (2-1)(3-1) = 2$. At 1% level of significance we find that $\chi^2 = 9.21$, with 2 d.f. Hence we reject the null hypothesis. In other words there is an evidence that there is an association between attitude towards working in the government and father's Job.

2X2 tables

Let us consider the case where each of the two variates A&B is classified into two categories. The tables are known as 2X2 tables and take the following form:-

		Variate B		Total
		B ₁	B ₂	
Variate A	A ₁	a	b	a+b
	A ₂	c	d	c+d
	Total	a+c	b+d	N

where $N = a + b + c + d$

The expected frequency for "a" is $(a+b)(a+c)/N$

∴ The difference between "a" and its expected values will be

$$a - \frac{(a+b)(a+c)}{N} = \frac{ad - bc}{N}$$

for

Similarly the differences, b, c, d and their expected values are

$$-\frac{ad - bc}{N}, -\frac{ad - bc}{N}, \frac{ad - bc}{N}$$

respectively

$$\begin{aligned} \therefore \chi^2 &= \frac{(f_o - f_e)^2}{f_e} \\ &= \frac{(ad - bc)^2}{N^2} \left[\frac{1}{(a+b)(a+c)} + \frac{1}{(a+b)(b+d)} + \frac{1}{(c+d)(a+c)} \right. \\ &\quad \left. + \frac{1}{(c+d)(b+d)} \right] \\ &= \frac{N(ad - bc)^2}{(a+b)(c+d)(a+c)(b+d)} \end{aligned}$$

Example

The distribution of a sample of 1464 individuals according to sex and attitude towards smoking is given in the following table

	Smoke	Do not Smoke	Total
Men	295	462	757
Women	183	524	707
Total	478	966	1464

The hypothesis to be tested is that there is no association between smoking and sex. Using the above formula the value of χ^2 is 28.5. For $\gamma = 1$, and at 1% level of significance we have from the tables $\chi^2 = 6.6$. Hence we reject the hypothesis. In other words there is an evidence that in the population, attitude towards smoking depends upon sex.

Note (1) when the sample size is < 500 , we use Yate's correction which is only appropriate for χ^2 with 1 degree of freedom. The procedure is to change the frequency in each cell by 0.5, keeping the marginal total unchanged and reducing the size of χ^2 . Hence

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

Note (2) It is known that for a large sample if χ^2 has 1 degree of freedom, $\sqrt{\chi^2}$ has a distribution which is the right hand half of a normal distribution. Hence, we can use χ^2 with 1 degree of freedom as a test for the hypothesis that the proportions are the same in the two populations. In the above example, the hypothesis to be tested may be stated as follows: the proportion of persons who smoke is the same among men as among women. The outcome of the test gives an evidence that equal proportions in the two populations is untenable.

In testing the difference between two proportions, the test criterion is

$$u = \frac{p_1 - p_2}{\tilde{\sigma}_{p_1 - p_2}}$$

$$\text{where } \tilde{\sigma}_{p_1 - p_2} = \sqrt{\hat{\theta}(1-\hat{\theta})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

In terms of the 2 X 2 table,

$$p_1 = \frac{a}{a+b}, \quad p_2 = \frac{c}{c+d}, \quad \hat{\theta} = \frac{a+c}{N}$$

$$\therefore (p_1 - p_2)^2 = \left(\frac{a}{a+b} - \frac{c}{c+d} \right)^2 = \frac{(ad - bc)^2}{(a+b)^2 (c+d)^2}$$

$$\begin{aligned} \hat{\phi} (1 - \hat{\phi}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right) &= \frac{a+c}{N} \frac{b+d}{N} \left(\frac{1}{a+b} + \frac{1}{c+d} \right) \\ &= \frac{(a+c)(b+d)}{N(a+b)(c+d)} \end{aligned}$$

$$\therefore u^2 = \frac{N (ad - bc)^2}{(a+b)(c+d)(a+b)(c+d)} = \chi_1^2$$

Note (3) In $2 \times k$ tables, χ_{k-1}^2 is used to test the null hypothesis that the k proportions are identical.

(112)

(X) Simple Analysis of Variance

Simple Analysis of Variance

1. Introduction: Analysis of variance is a technique applicable to testing the hypothesis that several independent samples have been drawn at random from a common normal population. The development of the analysis of variance as a powerful test in experimental and research work is largely the responsibility of R.A. Fisher and his Co-workers.

The analysis of variance has proved to be not only convenient method but also a powerful method of analysis for the research worker is testified to by the extent to which it is being used in the planning design and analysis of research in a variety of fields.

Basically the analysis of variance is a simple arithmetical method of sorting out the components of variation in a given set of results. Whenever there is heterogeneity of variation, more than one component is present. Suppose that we have a population in which the people are all of the same race and the variable studied is height in inches. Four groups are drawn from the population. One group of men, another group of women, a third group of boys from 13-15 years of age and a fourth group of girls of a similar age. When the four groups are combined, the frequency distribution might appear reasonably normal but, we know that two components of variation are actually present, one representing variation within groups and another between the groups. The arithmetical procedure of the analysis of variance enables us to sort out and evaluate the components of variation for such mixed population.

In the present chapter, we are going to deal with analysis of variance in its simplest form, namely, the on-way classification or sometimes called simple randomized design. The importance of this design in experimental work cannot be overemphasized. Not only is the design widely employed by itself, but it constitutes a basic unit in nearly all of the more complex designs employed in experimental research. In this design each treatment is independently administered to a different group of subjects, all groups having been originally drawn at random from the same parent population. After the treatments have been administered, these groups may be regarded as random samples from a single population, only if the treatments all have identical effects on the distribution of criterion measures for the population. Otherwise, the group receiving treatment T_1 say, may be regarded as a random sample from a hypothetical population which is like the parent population except that all its members have received treatment T_1 . The sample that received treatment T_2 may, likewise, be regarded as a random sample from a population like the original except that all members of this population have received T_2 etc.

The hypothesis to be tested is that the criterion means of these treatment populations are identical. The design can be represented as follows:-

T_1	T_2	T_j	T_k
x_{11}	x_{12}	x_{1j}	x_{1k}
x_{21}	x_{22}	x_{2j}	x_{2k}
.
.
.
x_{i1}	x_{i2}	x_{ij}	x_{ik}
.
.
.
x_{n1}	x_{n2}	x_{nj}	x_{nk}
$\bar{x}_{.1}$	$\bar{x}_{.2}$	$\bar{x}_{.j}$	$\bar{x}_{.k}$

$i = 1, 2, \dots, n$ & $j = 1, 2, \dots, k$

x_{ij} denotes the value of the i^{th} observation in the group of the j^{th} treatment

$\bar{x}_{.j}$ denotes the mean of the values of the j^{th} treatment

$\bar{x}_{..}$ denotes the over all mean

n denotes the number of observations in each group.

It may differ from one treatment to another i.e. n_j where $j = 1, 2, \dots, k$.

If all the treatments have identical effects on the distribution of criterion measures for the parent population, that is, if the distribution of criterion measures is the same for all treatment populations, these populations may be regarded as just one population. In this case the various treatment groups may be all regarded as simple random samples from the same population whose variance we denote by σ^2 .

2. Estimates for σ^2 and measure of discrepancy

From the experimental data we can derive two independent estimates for σ^2 . One estimate is based on the differences among the observed treatment (group) means, the other upon the variance of the measures within the individual treatment groups. We can form the ratio of the first of these estimates to the second. If the treatments have identical effects, the 1st of these estimates will exceed the other only by chance. If the treatments really differ in effectiveness, then the differences among the observed treatment means will be larger than they would be and the ratio will then tend to be larger than 1.0. Accordingly we can use the ratio as the measure of the discrepancy between hypothesis and observations. If we find the sampling distribution of this ratio, we can use it in a statistical test of the hypothesis.

The 1st estimate (based on the means of the groups) If the k treatments have identical effects on the distribution of the criterion measures for the parent population the various treatment groups are all random samples from the same population and the means of these treatment groups are a simple random sample of k values from a population of such means. Hence an unbiased estimate is given as

$$\tilde{\sigma}_{\bar{x}}^2 = \sum_{j=1}^k (\bar{x}_j - \bar{\bar{x}})^2 / (k-1)$$

But

$$\tilde{\sigma}_{\bar{x}}^2 = \tilde{\sigma}^2 / n$$

$$\tilde{\sigma}_{(1)}^2 = n \tilde{\sigma}_{\bar{x}}^2 = \frac{n \sum (\bar{x}_j - \bar{\bar{x}})^2}{k-1} \quad (1)$$

The 2nd estimate (based on the individuals within groups) If the treatments have identical effects on the criterion distribution, each treatment group is a simple random sample from the same population. For the j^{th} treatment group, an unbiased estimate for σ^2 is given as

$$\tilde{\sigma}_j^2 = \sum_j (x_{ij} - \bar{x}_{.j})^2 / (n-1)$$

A better estimate can be secured by averaging

$$\tilde{\sigma}_1^2, \tilde{\sigma}_2^2, \dots, \tilde{\sigma}_k^2, \text{ i.e., } \tilde{\sigma}_{(2)}^2 = \sum_{j=1}^k \tilde{\sigma}_j^2 = \frac{\sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2}{k(n-1)} \quad (2)$$

3. Interpretation of $\frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(2)}^2}$

$$\frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(2)}^2} = \frac{n \sum_j (\bar{x}_{.j} - \bar{x})^2 / (k-1)}{\sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2 / k(n-1)} \quad (3)$$

We know that if the population has a normal distribution, then the sampling distribution of the means of random samples drawn from this population is also normal having the same mean as the population mean and variance σ^2/n .

$$u = \frac{(\bar{x}_{.j} - \bar{x})}{\sigma / \sqrt{n}} \sim N(0,1)$$

$$\therefore \sum_{j=1}^k \frac{n(\bar{x}_{.j} - \bar{x})^2}{\sigma^2} = \sum_{j=1}^k u^2 = \chi_{k-1}^2 \quad (4)$$

Also, for the j^{th} group

$$u = \frac{x_{ij} - \bar{x}_{.j}}{\sigma} \sim N(0,1)$$

$$\sum_{i=1}^n \frac{(x_{ij} - \bar{x}_{.j})^2}{\sigma^2} = \sum_{i=1}^n u^2 = \chi_{n-1}^2$$

Hence

$$\sum_{j=1}^k \sum_{i=1}^n \frac{(x_{ij} - \bar{x}_{.j})^2}{\sigma^2} = \sum_{j=1}^k \chi^2_{n-1} = \chi^2_{k(n-1)} \quad (5)$$

(applying the additive property of χ^2)Dividing both numerator and denominator by σ^2 and using (4), (5) we get

$$\frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(2)}^2} = \frac{\chi^2_{k-1} / (k-1)}{\chi^2_{k(n-1)} / k(n-1)} \sim F_{k-1, k(n-1)} \quad (6)$$

Note: The above derivation is valid even if the number of observations (n_j) differs from one group to another

$$1. \quad \frac{x_{ij} - \bar{x}_{.j}}{\sigma} = u$$

$$\sum_{i=1}^{n_j} \frac{(x_{ij} - \bar{x}_{.j})^2}{\sigma^2} = \sum_{i=1}^{n_j} u^2 = \chi^2_{n_j-1}$$

$$\sum_{j=1}^k \sum_{i=1}^{n_j} \frac{(x_{ij} - \bar{x}_{.j})^2}{\sigma^2} = \sum_j \chi^2_{n_j-1} = \chi^2_{N-k}$$

where

$$N = \sum_{j=1}^k n_j$$

$$2. \quad \text{Also } \frac{(\bar{x}_{.j} - \bar{x})}{\sigma / \sqrt{n_j}} = u_j$$

$$\therefore \sum_{j=1}^k \frac{n_j (\bar{x}_{.j} - \bar{x})^2}{\sigma^2} = \sum_{j=1}^k u_j^2 = \chi^2_{k-1}$$

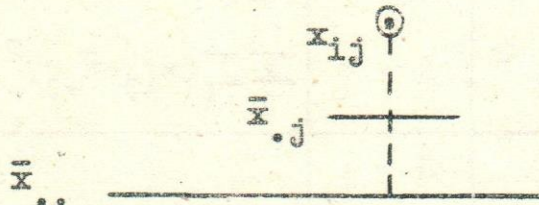
$$\begin{aligned} \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(2)}^2} &= \frac{\sum_{j=1}^k n_j \frac{(\bar{x}_{.j} - \bar{x})^2}{\sigma^2} / (k-1)}{\sum_{j=1}^k \sum_{i=1}^{n_j} \frac{(x_{ij} - \bar{x}_{.j})^2}{\sigma^2} / (N-k)} = \frac{\chi^2_{k-1} / (k-1)}{\chi^2_{N-k} / (N-k)} \\ &= F_{k-1, N-k} \end{aligned}$$

4. Analysis of Variance tablePartition of total sum of squares

Again let x_{ij} denote the i^{th} observation in the j^{th} group or (sample)

$\bar{x}_{.j}$ denote the mean of the j^{th} group

& $\bar{x}_{..}$ denote the over all mean



Hence

$x_{ij} - \bar{x}_{..}$ represents the deviation of the i^{th} observation in the j^{th} group from the over all mean.

$x_{ij} - \bar{x}_{.j}$ represents the deviation of the i^{th} observation in the j^{th} group from the mean of the j^{th} group.

$\bar{x}_{.j} - \bar{x}_{..}$ represents the deviation of the mean of the j^{th} group from the over all mean

$$\text{Now } x_{ij} - \bar{x}_{..} = (x_{ij} - \bar{x}_{.j}) + (\bar{x}_{.j} - \bar{x}_{..})$$

Squaring both sides

$$\sum_j \sum_i (x_{ij} - \bar{x}_{..})^2 = \sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2 + \sum_j n_j (\bar{x}_{.j} - \bar{x}_{..})^2$$

i.e.

Total sum of squares (SST) = Sum of squares within groups (SSW)
+ Sum of Squares between groups (SSB)

Where SST has $N-1$ degrees of freedom

SSW has $N-k$ " " "

SSB has $k-1$ " " "

For numerical calculations of the sums of Squares we have

$$\text{SST} = \sum_j \sum_i (x_{ij} - \bar{x}_{..})^2 = \sum_j \sum_i x_{ij}^2 - \frac{(\sum_j \sum_i x_{ij})^2}{N}$$

$$\text{SSB} = \sum_j n_j (\bar{x}_{.j} - \bar{x}_{..})^2 = \sum_j \frac{(\sum_i x_{ij})^2}{n_j} - \frac{(\sum_j \sum_i x_{ij})^2}{N}$$

and by subtraction we get SSW.

$(\sum_j \sum_i x_{ij})^2 / N$ is called the correction term.

A summary can be arranged in the following analysis of variance table

Source of variation	Sum of Squares	Degrees of freedom	Mean Sum of Squares
Between groups	$\sum_j n_j (\bar{x}_j - \bar{x}_{..})^2$	$k - 1$	$\sum_j n_j (\bar{x}_j - \bar{x}_{..})^2 / (k-1)$
Within groups	$\sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2$	$N - k$	$\sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2 / (N-k)$
Total	$\sum_j \sum_i (x_{ij} - \bar{x}_{..})^2$	$N - 1$	

5- Expected values of sums of squares

Let x_{ij} denote the value of the character x for the i^{th} observation in the j^{th} sample (or group) we may write the model as

$$x_{ij} = \mu_j + z_{ij}$$

where z_{ij} is a random variable with expected value 0 & standard deviation σ_j . ($j=1,2,\dots,k$).

We suppose that the sampling is random within each group and x_{ij} from different group are independent. Further we assume that the x_{ij} 's are independent in the same group. We also assume that

$\sigma_1 = \sigma_2 = \dots = \sigma_k$ i.e. the groups from which the samples are drawn have common variance.

In the above model the statistical hypothesis to be tested is

$$\mu_1 = \mu_2 = \dots = \mu_j = \dots = \mu_k.$$

This model can be modified replacing μ_j by $\mu + \lambda_j$ where μ is a parameter representing the over all average level of the character measured. λ_j represents the deviations of the average of each group from the whole average. Without any loss of generality we can assume that $\sum n_j \lambda_j = 0$

The model becomes

$$x_{ij} = \mu + \lambda_j + z_{ij} \quad (1)$$

This can be presented in another way. Assume that the x_{ij} are independently $N(\mu_j, \sigma^2)$. Consider the identity

$$x_{ij} = \mu + (\mu_j - \mu) + (x_{ij} - \mu_j)$$

Let $\lambda_j = \mu_j - \mu$ and $\mu = \sum_{j=1}^k n_j \mu_j / N$

where $N = \sum_{j=1}^k n_j$

so that
$$\sum_{j=1}^k n_j \lambda_j = \sum_{j=1}^k n_j (\mu_j - \mu) = \sum_{j=1}^k n_j \mu_j - N \mu = N - N = 0$$

Also define $z_{ij} = x_{ij} - \mu_j$ which has mean 0 since the mean of x_{ij} is μ_j . But both z_{ij} & x_{ij} have the same variance σ^2 . Finally z_{ij} are normally distributed since the x_{ij} came from a normal distribution.

$$\therefore x_{ij} = \mu + \lambda_j + z_{ij} \quad \begin{cases} i = 1, 2, \dots, n_j \\ j = 1, 2, \dots, k \end{cases}$$

z_{ij} are independently $N(0, \sigma^2)$

$$\sum_{j=1}^k n_j \lambda_j = 0$$

Now
$$\bar{x}_{.j} = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{ij} = \mu + \lambda_j + \bar{z}_{.j} \quad (2)$$

$$\bar{x}_{..} = \frac{1}{N} \sum_j \sum_i x_{ij} = \frac{1}{N} \sum_{j=1}^k n_j \bar{x}_{.j} = \mu + \bar{z}_{..} \quad (3)$$

$$\therefore \bar{x}_{.j} - \bar{x}_{..} = \lambda_j + (\bar{z}_{.j} - \bar{z}_{..}) \quad \& \quad x_{ij} - \bar{x}_{.j} = z_{ij} - \bar{z}_{.j}$$

Hence

1- Sum of Squares between groups (SSB) =
$$\sum_{j=1}^k n_j (\bar{x}_{.j} - \bar{x}_{..})^2$$

$$= \sum_{j=1}^k n_j [\lambda_j + (\bar{z}_{.j} - \bar{z}_{..})]^2$$

$$\begin{aligned}
\therefore E(SSB) &= \sum_{j=1}^k n_j \lambda_j^2 + E\left\{ \sum_{j=1}^k n_j (\bar{z}_{.j} - \bar{z}_{..})^2 \right\} \\
&= \sum_{j=1}^k n_j \lambda_j^2 + (k-1) n_j \frac{\sigma^2}{n_j} \\
&= (k-1) \sigma^2 + \sum_{j=1}^k n_j \lambda_j^2 \quad (4)
\end{aligned}$$

$$\begin{aligned}
2- \text{Sum of squares within groups (SSW)} &= \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{.j})^2 \\
&= \sum_{j=1}^k \sum_{i=1}^{n_j} (z_{ij} - \bar{z}_{.j})^2 \\
\therefore (SSW) &= \sum_{j=1}^k (n_j - 1) \sigma^2 = (N - k) \sigma^2 \quad (5)
\end{aligned}$$

These can be represented in the following summary table

Source of variation	Sum of Squares	Degrees of freedom	Expected value of mean sum of Squares
Between groups	$\sum_{j=1}^k n_j (\bar{x}_{.j} - \bar{x}_{..})^2$	$k-1$	$\sigma^2 + \sum_{j=1}^k n_j \sigma_j^2 / (k-1)$
Within groups	$\sum_j \sum_i (x_{ij} - \bar{x}_{.j})^2$	$N-k$	σ^2
Total	$\sum_j \sum_i (x_{ij} - \bar{x}_{..})^2$	$N-1$	

If the hypothesis tested is valid i.e. $\lambda_1 = \lambda_2 = \dots = \lambda_k = 0$, (4) becomes $(k-1) \sigma^2$, while (5) is still $(N-k) \sigma^2$. A large value of the mean square ratio obtained from (4), (5) by $(k-1)$, $(N-k)$ the corresponding degrees of freedom respectively may be regarded as evidence of difference between groups independently of any assumption about the distribution of the z_{ij} 's. That is we have no system of significance limits against which the observed values may be judged. To obtain such significance limits, we assume that each of the z_{ij} is normally distributed i.e. $z_{ij} \sim N(0, \sigma)$. On the null hypothesis i.e. $\lambda_j = 0$, the mean square ratio is

$$\frac{\sum_j n_j (\bar{z}_{.j} - \bar{z}_{..})^2 / (k-1)}{\sum_j \sum_i (z_{ij} - \bar{z}_{.j})^2 / (N-k)}$$

But $\frac{1}{k-1} \sum_{j=1}^k n_j (\bar{z}_{.j} - \bar{z}_{..})^2 \sim \chi_{k-1}^2 \sigma^2 / (k-1)$

$$\frac{1}{N-k} \sum_j \sum_i (z_{ij} - \bar{z}_{.j})^2 \sim \chi_{N-k}^2 \sigma^2 / (N-k)$$

∴ Sum of squares within groups independent of sum of squares between groups, then

$$\text{Mean square ratio} = \frac{\text{MSSB}}{\text{MSSW}} = \frac{\chi_{k-1}^2 / (k-1)}{\chi_{N-k}^2 / (N-k)} \sim F_{k-1, N-k}$$

Example Random samples of dials are drawn from 4 different machines. The values of their diameters are given as follows

machines

A	B	C	D
25	24	23	24
20	26	23	23
28	30	32	33
15	25	24	26
	35	28	29
	34		

Taking 25 as an arbitrary origin we have

A	B	C	D
0	-1	-2	-1
-5	1	-2	-2
3	5	7	8
-10	0	-1	1
	10	+3	4
	9		

n_j	4	6	5	5	20 = N
$\sum_i x_{ij}$	-12	24	5	10	27 = $\sum_j \sum_i x_{ij}$
$\sum_i x_{ij}^2$	134	208	67	86	495 = $\sum_j \sum_i x_{ij}^2$
$(\sum_i x_{ij})^2 / n_j$	36	96	5	20	157 = $\sum_j \frac{(\sum_i x_{ij})^2}{n_j}$

$$\begin{aligned}\text{Correction term} &= \frac{(27)^2}{20} \\ \text{SST} &= 495 - 36.45 = 458.55 \\ \text{SSB} &= 157 - 36.45 = 120.55\end{aligned}$$

Analysis of variance table

Source of Variation	Sum of Squares	Degrees of freedom	Mean sum of Squares.
Between groups	120.55	3	40.183
Within groups	338.00	16	21.125
Total	458.55	19	

$$\text{Mean Ratio} = \frac{40.183}{21.125} = 1.90$$

From the tables we have $F_{3,16} = 3.24$

∴ We accept the null hypothesis, i.e., there is an evidence that the population means of the diameters of the discs produced by the 4 machines are identical

Note If F proves to be significant this means that $\lambda_j \neq 0$. In other words the population means are not identical. To test the significance of the difference between the means of any two populations from which samples h, g are drawn we apply

$$t = \frac{\bar{x}_h - \bar{x}_g}{\sqrt{\sigma^2 \left(\frac{1}{n_h} + \frac{1}{n_g} \right)}}$$

n_h, \bar{x}_h are the size and the mean of the sample "h", n_g, \bar{x}_g are the size and the mean of the sample "g" & σ^2 is the mean sum of squares of within groups.

Bartlett's test for homogeneity of variances

In analysis of variance it is necessary to assume that a number of population variances are equal. If we have reason to doubt that this is the case, we may want to test the null hypothesis

$$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$$

gainst the alternative

H_1 : at least two variances differ Under the assumptions:-

- (i) k random samples are drawn from k populations
- (ii) the k populations are normal, the statistic

$$\frac{2.3026}{C} \left[(N-k) \log_{10} s_p^2 - \sum_{j=1}^k (n_j-1) \log_{10} s_j^2 \right]$$

is approximately distributed as X^2 with $k-1$ degrees of freedom if H_0 is true. Here $s_1^2, s_2^2, \dots, s_k^2$ are the k samples variances. Also

$$s_p^2 = \frac{1}{N-k} \sum_{j=1}^k (n_j-1) s_j^2,$$

$$s_j^2 = \frac{1}{n_j-1} \sum (x_{ij} - \bar{x}_j)^2$$

and
$$C = 1 + \frac{1}{3(k-1)} \left[\sum_{j=1}^k \frac{1}{n_j-1} - \frac{1}{N-k} \right]$$

Example To test the homogeneity of the populations from which the following samples are drawn

- a. 8, 4, 7, 5, 6
- b. 6, 12, 14, 8
- c. 4, 2, 2, 1, 2, 1
- d. 3, 2, 4
- e. 10, 7, 15, 17, 13, 8, 14

a	b	c	d	e
8	6	4	3	10
4	12	2	2	7
7	14	2	4	15
5	8	1		17
6		2		13
		1		8
				14

n_j	5	4	6	3	7	25 = N
$\sum_i x_{ij}$	30	40	12	9	84	$175 = \sum_j \sum_i x_{ij}$
$\sum_i x_{ij}^2$	190	440	30	29	1092	$1781 = \sum_j \sum_i x_{ij}^2$
$(\sum_i x_{ij})^2 / n_j$	180	400	24	27	1008	$1639 = \sum_j \frac{(\sum_i x_{ij})^2}{n_j}$
$\sum_i (x_{ij} - \bar{x}_{.j})^2$	10	40	6	2	84	

Samples	$\sum_i (x_{ij} - \bar{x}_{.j})^2$	$n_j - 1$	s_j^2	$\log_{10} s_j^2$	$(n_j - 1) \log_{10} s_j^2$
a	10	4	2.5	0.3979	1.5916
b	40	3	13.3	1.1239	3.3717
c	6	5	1.2	0.0792	0.3960
d	2	2	1.0	0.0000	0.0000
e	84	6	14.0	1.1461	6.8766
	142	20			12.2359 = $\sum_j (n_j - 1) \log_{10} s_j^2$

$$s_p^2 = 142/20 = 7.1$$

$$\log_{10} s_p^2 = \log_{10} 7.1 = 0.8513$$

$$\begin{aligned} \therefore (N-k) \log_{10} s_p^2 &= 20 \times 0.8513 = 17.026 \\ &= 1.117 \end{aligned}$$

$$\therefore \chi_4^2 = \frac{2.3026}{1.117} [17.0260 - 12.2359] = 9.88$$

But from the tables $\chi_{4,5\%}^2 = 9.49$

\therefore We reject the null hypothesis $H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$

In other words there is no evidence that the populations from which these samples are drawn - have common variance

Exercises

1- A professor is trying to select a good text book for his elementary course in statistics from 4 different ones which are available. He has 37 students whom he distributed at random to 4 different groups. The assignment of textbooks to groups is also done at random. After the course is over all students who are still enrolled take the same examination. The results are

Textbook			
A	B	C	D
68	41	54	44
68	47	44	51
69	54	51	69
60	65	56	59
73	32	47	59
64	73	61	55
71	44	59	66
67	48	49	
75	64	41	
	54	61	
		73	

Assuming that differences attributed to using a number of textbooks are the only variables that need be considered, what conclusions can be drawn?

2- Worms are classified into three groups by a structural characteristic. Three random samples of 11 are taken from each group and the length of each worm is measured. These data are recorded as follows

Group (1) 8.9, 9.7, 11.5, 8.2, 10.5, 10.8, 11.0, 8.0, 9.9, 11.0, 11.0
Group (2) 12.2, 12.0, 11.5, 8.7, 10.5, 9.0, 10.5, 13.0, 13.0, 11.0, 11.1
Group (3) 9.5, 8.0, 8.3, 10.0, 9.5, 10.0, 11.3, 10.5, 8.0, 8.0, 9.2

Test the hypothesis that the mean length of each group is the same

3- Jenkins & Snedecor compared the yields of a number of varieties of corn. Each variety being represented by several inbred lines. Six varieties with yields (bushels per acre) of their inbred lines are follows:

(V_1) : 7.4, 7.3, 4.5, 7.4, 5.0, 5.9, 6.4, 6.3, 5.0, 6.1, 7.9, 5.7
 V_2 : 7.7, 5.4, 5.2, 4.0
 V_3 : 6.9, 6.8, 7.6, 8.1, 9.4, 12.0, 15.9, 7.4, 9.0, 5.2, 9.2, 8.6
 V_4 : 9.6, 7.8, 9.6, 7.7, 8.2, 7.3, 11.3, 9.5, 8.8, 8.4, 6.8
 V_5 : 4.8, 9.2, 8.5, 8.8, 7.9, 5.9, 9.2
 V_6 : 4.3, 8.4, 6.6, 4.9, 5.8, 7.6, 5.7

Given n pairs of values for the two variates X & Y , i.e.

$X : x_1, x_2, \dots, x_n$

$Y : y_1, y_2, \dots, y_n$,

it is required to find the best fitted line to the above data. The equation for such a line is not just a mathematical form where we can calculate the value of Y given the value of X and vice versa. In fact we are going to have two equations. From one equation we can predict the value of Y given the value of X with the least possible error in Y . This is called the regression line of Y on X and usually written in the form

$$Y = a + b(X - \bar{X}) .$$

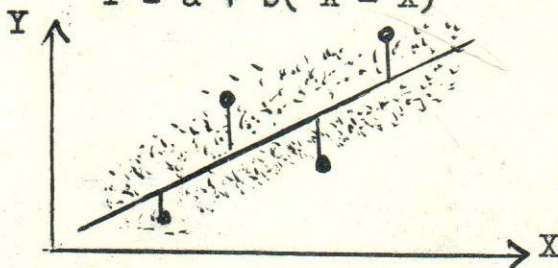
From the the other equation, we can predict the value of X given the value of Y with the least possible error in X . This is called the regression line of X on Y and usually written in the form

$$X = a_1 + b_1(Y - \bar{Y}) .$$

1- The regression line of Y on X :

Let the required equation be in the form

$$Y = a + b(X - \bar{X}) \quad (1)$$



If y_i denotes the observed value and \hat{y}_i denotes the predicted value then the error is represented by

$$e_i = y_i - \hat{y}_i \quad (i=1,2,\dots,n) \quad (2)$$

The over all error around the fitted line can be measured by

$$E = \sum e_i^2 / n = \sum (y_i - \hat{y}_i)^2 / n$$

But

$$\hat{y}_i = a + b(x_i - \bar{x})$$

then

$$E = \sum (y_i - a - b(x_i - \bar{x}))^2 / n$$

Hence E is minimum if the quantity

$$Q = \sum (y_i - a - b(x_i - \bar{x}))^2 \quad (3)$$

is minimum. (n is a given value)

It is obvious that this quantity depends on "a" & "b".
Now the question is : what are the estimates of "a" & "b" which make Q minimum ? Using the Least Square method, the required conditions are

$$\partial Q / \partial a = 0, \quad \partial Q / \partial b = 0 \quad (4)$$

From (3) we have

$$\left. \begin{aligned} \partial Q / \partial a &= -2 \sum (y_i - a - b(x_i - \bar{x})) \\ \partial Q / \partial b &= -2 \sum (x_i - \bar{x})(y_i - a - b(x_i - \bar{x})) \end{aligned} \right\} \quad (5)$$

From (4) & (5) the required conditions become

$$\left. \begin{aligned} \sum (y_i - a - b(x_i - \bar{x})) &= 0 \\ \sum (x_i - \bar{x})(y_i - a - b(x_i - \bar{x})) &= 0 \end{aligned} \right\} \quad (6)$$

Solving (6), we get

$$a = \sum y_i / n = \bar{y} \quad (7)$$

&

$$b = \sum (x_i - \bar{x})(y_i - \bar{y}) / \sum (x_i - \bar{x})^2 \quad (8)$$

Substituting in (1) we get the regression line of Y on X , and "b" is called the coefficient of regression.

Similarly it can be shown that the regression line of X on Y is

$$X = a_1 + b_1(Y - \bar{Y}),$$

where

$$a_1 = \bar{x}$$

&

$$b_1 = \sum (x_i - \bar{x})(y_i - \bar{y}) / \sum (y_i - \bar{y})^2$$

Remarks :

1- The two regression lines of Y on X and of X on Y intersect at the point (\bar{x}, \bar{y})

2- The regression coefficients of Y on X and of X on Y can be rewritten as

$$b = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

and

$$b_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum y_i^2 - (\sum y_i)^2}$$

Both formulae are preferable for numerical calculations.

3-Again the coefficients of regression can be shown to have the following forms:

$$b = r \frac{s_y}{s_x} \quad \text{and} \quad b_1 = r \frac{s_x}{s_y}$$

4- Introducing the above values for b & b_1 in the regression lines the two equations become

$$\frac{y - \bar{y}}{s_y} = r \frac{x - \bar{x}}{s_x}$$

and

$$\frac{x - \bar{x}}{s_x} = r \frac{y - \bar{y}}{s_y}$$

respectively.

5- From remark 3, we get

$$b b_1 = r^2$$

Hence b , b_1 and r have the same sign, i.e. they may be all positive or may be all negative.

6- If θ is the angle between the two regression lines then

$$\tan \theta = \frac{s_x s_y (r^2 - 1)}{r (s_x^2 + s_y^2)}$$

Consequently, if the two regression line coincide then $r = 1$ and if they are perpendicular then $r = 0$

7- If

$$d_1 = (x - c_1)/l_1 \quad \& \quad d_2 = (y - c_2)/l_2$$

then

$$b_{y.x} = b_{d_2.d_1} \frac{l_2}{l_1}$$

and

$$b_{x.y} = b_{d_1.d_2} \frac{l_1}{l_2}$$

2- Testing the significance of "a" & "b"

Let the fitted line be of the form

$$y_i = a + b (x_i - \bar{x}) \quad (9)$$

where "a" & "b" are given by (7) & (8)

Let also the theoretical model be of the form

$$y_i = \alpha + \beta (x_i - \bar{x}) + z_i \quad (10)$$

where z_i are independent random variates with

$$E(z_i) = 0 \quad \& \quad \text{Var.}(z_i) = \sigma^2 \quad (11)$$

From (9) & (10)

$$\bar{y} = a \quad (12)$$

$$\bar{y} = \alpha + \bar{z} \quad (13)$$

$$\therefore a = \alpha + \bar{z}$$

$$\text{Hence } (a) = \quad (14)$$

i.e. "a" is an unbiased estimator for α

$$\begin{aligned} \text{Also } \text{Var}(a) &= E(a - E(a))^2 = E(a - \alpha)^2 \\ &= E(\bar{z}^2) = \sigma^2/n \end{aligned} \quad (15)$$

Now

$$y_i = \alpha + \beta (x_i - \bar{x}) + z_i$$

$$\bar{y} = \alpha + \bar{z}$$

\therefore

$$y_i - \bar{y} = \beta (x_i - \bar{x}) + z_i - \bar{z}$$

Hence

$$\begin{aligned} b &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x})(\beta (x_i - \bar{x}) + z_i - \bar{z})}{\sum (x_i - \bar{x})^2} \\ &= \beta + Z \end{aligned} \quad (16)$$

where

$$\begin{aligned} Z &= \sum (x_i - \bar{x}) z_i / \sum (x_i - \bar{x})^2 \\ &= \sum c_i z_i \end{aligned} \quad (17)$$

The c_i are constants and given as

$$c_i = (x_i - \bar{x}) / \sum (x_i - \bar{x})^2 \quad (18)$$

$$\begin{aligned}\therefore E(Z) &= E\left(\sum c_i z_i\right) = 0 \\ &\& \text{Var}(Z) = \sum c_i^2 \text{Var}(z_i) = \sigma^2 \sum c_i^2 \\ &= \sigma^2 / \sum (x_i - \bar{x})^2\end{aligned}\quad (19)$$

$$\therefore E(b) = E(\beta + Z) = \beta \quad (20)$$

$$\begin{aligned}\text{Var}(b) &= (b - E(b))^2 \\ &= E(b - \beta)^2 = E(Z^2) \\ &= \text{Var}(Z) = \sigma^2 / \sum (x_i - \bar{x})^2\end{aligned}\quad (21)$$

Now, to test the significance of "a", we test $\alpha = \alpha_0$ against $\alpha \neq \alpha_0$. The test criterion is

$$T = \frac{a - \alpha_0}{\sigma / \sqrt{n}}$$

To test the significance of "b", we test $\beta = 0$ against $\beta \neq 0$. The test criterion is

$$T = \frac{b}{\sigma / \sqrt{\sum (x_i - \bar{x})^2}}$$

In the above two test criteria, σ is usually unknown. An estimate for σ^2 is discussed in the next section.

3- Regression analysis table

Partition of total sum of squares i.e. $(y_i - \bar{y})^2$

$$\begin{aligned}\sum (y_i - \bar{y})^2 &= \sum (y_i - \hat{y} + \hat{y} - \bar{y})^2 \\ &= \sum (y_i - \hat{y})^2 + \sum (\hat{y} - \bar{y})^2 + 2 \sum (y_i - \hat{y})(\hat{y} - \bar{y})\end{aligned}$$

$$\begin{aligned}\text{But } \sum (y_i - \hat{y})(\hat{y} - \bar{y}) &= \sum (y_i - \bar{y} - b(x_i - \bar{x}))(b(x_i - \bar{x})) \\ &= b \left[\sum (y_i - \bar{y})(x_i - \bar{x}) - b \sum (x_i - \bar{x})^2 \right] \\ &= 0\end{aligned}$$

$$\therefore \sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y})^2 + \sum (\hat{y}_i - \bar{y})^2$$

In other words the total sum of squares ($\sum (y_i - \bar{y})^2$) is partitioned into
 sum of squares due to regression ($\sum (\hat{y}_i - \bar{y})^2 = b^2 \sum (x_i - \bar{x})^2$)
 and
 sum of squares of residuals ($\sum (y_i - \hat{y})^2$)

Expected values for sum of squares due to regression and sum of squares of residuals

$$\begin{aligned} 1- \text{Sum of squares due to regression} &= b^2 \sum (x_i - \bar{x})^2 \\ &= (\beta + Z)^2 \sum (x_i - \bar{x})^2 \end{aligned}$$

$$\begin{aligned} \text{Expected value of SS(Reg.)} &= E[(\beta^2 + 2\beta Z + Z^2) \sum (x_i - \bar{x})^2] \\ &= \beta^2 \sum (x_i - \bar{x})^2 + \sum (x_i - \bar{x})^2 E(Z^2) \\ &= \sigma^2 + \beta^2 \sum (x_i - \bar{x})^2 \end{aligned}$$

$$2- \text{Sum of squares of residuals} = \sum (y_i - \hat{y})^2$$

$$\begin{aligned} &= \sum (y_i - \bar{y} - b(x_i - \bar{x}))^2 \\ &= \sum [(z_i - \bar{z}) - (b - \beta)(x_i - \bar{x})]^2 \\ &= \sum (z_i - \bar{z})^2 + \sum Z^2 (x_i - \bar{x})^2 \\ &\quad - 2 \sum (z_i - \bar{z}) Z (x_i - \bar{x}) \\ &= \sum (z_i - \bar{z})^2 - \sum Z^2 (x_i - \bar{x})^2 \end{aligned}$$

$$\begin{aligned} \text{Expected value of SS(Resd.)} &= E[\sum (z_i - \bar{z})^2 - \sum Z^2 (x_i - \bar{x})^2] \\ &= (n-1) \sigma^2 - \sigma^2 = (n-2) \sigma^2 \end{aligned}$$

In other words

$$E\left\{\sum (y_i - \hat{y})^2 / (n-2)\right\} = \sigma^2$$

i.e. $\sum (y_i - \hat{y})^2 / (n-2)$ is an unbiased estimator for σ^2

$$\therefore \hat{\sigma}^2 = \sum (y_i - \hat{y})^2 / (n-2)$$

= Mean sum of squares of residual.

Source of Variation	Sum of Squares	Degrees of Freedom	Mean sum of Squares	Expected value of Mean-Sum of Squares
Due to regression	$b^2 \sum (x_i - \bar{x})^2$	1	$b^2 \sum (x_i - \bar{x})^2$	$\sigma^2 + \beta^2 \sum (x_i - \bar{x})^2$
Residual	$\sum (y_i - \hat{y})^2$	n- 2	$\sum (y_i - \hat{y})^2 / (n-2)$	σ^2
Total	$\sum (y_i - \bar{y})^2$	n- 1		

Regression Analysis Table

To test the significance wheather $\rho = 0$, we have the criterion

$$\frac{b^2 \sum (x_i - \bar{x})^2}{\sum (y_i - \hat{y})^2 / (n-2)} = F_{1, n-2}$$

$$= t_{n-2}^2$$

Note : From the above criterion, we can deduce the distribution of the coefficient correlation r given as

$$r = \sum (x_i - \bar{x})(y_i - \bar{y}) / ns_x s_y$$

$$\begin{aligned} \text{Sum of squares due to linear regression} &= b^2 \sum (x_i - \bar{x})^2 \\ &= r^2 \sum (y_i - \bar{y})^2 \end{aligned}$$

$$\begin{aligned} \text{Sum of squares of residual} &= \sum (y_i - \hat{y})^2 \\ &= (1-r^2) \sum (y_i - \bar{y})^2 \end{aligned}$$

$$\therefore \frac{b^2 \sum (x_i - \bar{x})^2}{\sum (y_i - \hat{y})^2 / (n-2)} = \frac{r^2 (n-2)}{(1-r^2)}$$

$$\text{i.e. } \frac{r^2 (n-2)}{(1-r^2)} = t_{n-2}^2 \quad \therefore \frac{r \sqrt{n-2}}{\sqrt{(1-r^2)}} = t_{n-2}$$

$$\therefore \text{ If } u = r \sqrt{n-2} / \sqrt{1-r^2} \text{ \& } \rho = 0, \text{ then } u = t_{n-2}$$

$$\text{But } p(t_{n-2}) = c \left(1 + \frac{t^2}{n-2} \right)^{-(n-1)/2}$$

$$\begin{aligned} \text{Hence } p(r) &= p(t) \left| \frac{dt}{dr} \right| \\ &= c (1 - r^2)^{(n-4)/2} \end{aligned}$$

where

$$c = 2^{3-n} / B \left[\frac{1}{2}(n-2), \frac{1}{2}(n-2) \right]$$

EXAMPLES

1-

<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>	<u>x</u>	<u>y</u>
3	2	3	1	3	1	3	1	5	3	7	4		
5	5	5	4	5	4	5	4	7	5	9	3		
7	5	7	5	7	5	7	5	9	4	5	3		
9	5	9	4	9	4	9	4	5	3	7	4		
11	3	11	2	11	2	11	2	7	4	7	4		
7	4	9	3	5	2	3	7	9	3	5	2		
7	4	7	4	7	4	7	3	9	2	7	3		
7	3												

$$n = 50$$

$$\sum x = 350$$

$$\sum y = 168$$

$$\sum xy = 1202$$

$$\sum x^2 = 2690$$

$$\sum y^2 = 632$$

$$\bar{x} = 7$$

$$\bar{y} = 3.36$$

$$\therefore a = \bar{y} = 3.36 \quad \& \quad b = 0.108$$

$$\therefore y = 3.36 + 0.108(x - 7)$$

$$= 2.60 + 0.11x$$

$$\text{Total sum of squares} = \sum (y - \bar{y})^2 = \sum y^2 - (\sum y)^2/n$$

$$= 67.52$$

Due to regression sum of squares

$$= b^2 \sum (x - \bar{x})^2 = b^2 \left[\sum x^2 - (\sum x)^2/n \right]$$

$$= 2.82$$

Regression analysis table

Source of Variation	Sum of Squares	Degrees of freedom	Mean SSQ
Due to regression	2.82	1	2.82
Residual	64.70	48	1.35
Total	67.52	49	

$$R = 2.82/1.35 = 2.09$$

$$\text{From the tables } F_{1,48, 5\%} = 4.00$$

\therefore We accept the null hypothesis that $= 0$. In other words the regression is not significant

2- Distribution of 100 labours according to their age and number of working hours

No. of working hours (Y) \ Age(X)	25-	35-	45-	Total
27-		10	10	20
32-	10	40		50
37-	20	10		30
Total	30	60	10	100

X-table

Intervals	f	d_1	fd_1	fd_1^2
25-	30	-1	-30	30
35-	60	0	0	0
45-55	10	1	10	10
Total	100		-20	40

Y-table

Intervals	f	d_2	fd_2	fd_2^2
27-	20	-1	-20	20
32-	50	0	0	0
37-42	30	1	30	30
Total	100		10	50

$$N = 100 \quad \sum fd_1 d_2 = -30$$

$$\sum f(x-\bar{x})^2 = 100(40-4) = 3600$$

$$\sum f(y-\bar{y})^2 = 25(50-1) = 1225$$

$$\sum f(x-\bar{x})(y-\bar{y}) = 50(-30 + 2) = -1400$$

$$\therefore b = -1400/3600 = -0.389$$

$$\bar{x} = 10(-20/100) + 40 = 38$$

$$\bar{y} = 5(10/100) + 34.5 = 35$$

\therefore The regression line is

$$y = 35 - 0.389(x-38)$$

$$\text{Now Total sum of squares } \sum f(y-\bar{y})^2 = 1225$$

$$\text{Due to regression SSQ } b^2 \sum f(x-\bar{x})^2 = 544.44$$

$$\text{Residual sum of squares} = 680.56$$

$$\therefore R = 544.44 / (680.56/98) = 78.4$$

Comparing the value of R with that of $F_{1,98,5\%} = 4.00$, we reject the null hypothesis ($=0$).
In other words the regression is significant.

Exercise

The following is the bivariate distribution of 311 skulls according to their length & breadth. Carry out the required regression analysis

y \ x	Breadth in mms.										
	133-	136-	139-	142-	145-	148-	151-	154-	157-	160-	163-
160-				1	2						
165-	2		1	1		2	1				
170-			2	2	2	1		2			
175-	2	2	2	7	9	4	5	3			
180-		1	1	10	15	11	9	2	1		
185-		1	4	13	12	18	7	4	1		
190-				4	13	17	14	8	3	2	1
195-			3	1	13	14	5	6	2	1	
200-				2	6	4	4	3			1
205-				1	3	2	1	2	1	1	
210-				1	1	3	1	2			
215-							1	1			
Total	4	4	13	43	76	76	48	33	8	4	2

Exercise

The followings are the measures of the wing length and tongue length both in millimeters for 25 bees:-

Wing :	9.68	9.81	9.59	9.68	9.84	9.59	9.61	
Tongue:	6.53	6.71	6.70	6.69	6.70	6.62	6.59	
Wing :	9.55	9.25	9.08	9.70	9.60	9.50	9.74	9.72
Tongue:	6.55	6.35	6.25	6.61	6.51	6.55	6.74	6.75
Wing :	9.64	9.73	9.77	9.72	9.54	9.65	9.74	9.59
Tongue:	6.45	6.75	6.70	6.65	6.68	6.77	6.44	6.54
Wing :	9.71	9.56						
Tongue:	6.64	6.55						

Carry out the regression analysis required.



(XII) Multiple Regression

(1). Introduction: This is an extension for simple linear regression, where we have more than two characters. We are going to deal with the Case where we can express one of our variables (dependent) in terms of the other variables (independent). Suppose X_0 is the dependent variable & X_1, X_2, \dots, X_k are the independent variables. Let X_0 take the values x_{0i} , X_1, X_2, \dots, X_k take the values $x_{1i}, x_{2i}, \dots, x_{ki}$ and $i = 1, 2, \dots, n$.

The data can be arranged in the following form

X_0	X_1	X_2	\dots	X_k
x_{01}	x_{11}	x_{21}	\dots	x_{k1}
x_{02}	x_{12}	x_{22}	\dots	x_{k2}
\vdots	\vdots	\vdots	\dots	\vdots
\vdots	\vdots	\vdots	\dots	\vdots
x_{0n}	x_{1n}	x_{2n}	\dots	x_{kn}

The expected value of x_{0i} given the values of $x_{1i}, x_{2i}, \dots, x_{ki}$ is

$$x_{0i} = a + b_1(x_{1i} - \bar{x}_1) + \dots + b_r(x_{ri} - \bar{x}_r) + \dots + b_k(x_{ki} - \bar{x}_k) \quad (1)$$

The mathematical model is

$$x_{0i} = \alpha + \beta_1(x_{1i} - \bar{x}_1) + \dots + \beta_r(x_{ri} - \bar{x}_r) + \dots + \beta_k(x_{ki} - \bar{x}_k) + \varepsilon_i \quad (2)$$

We are going to treat these independent variables as constants. Moreover, we assume that the ε_i are independent random variables, $E(\varepsilon_i) = 0$ & $\text{Var}(\varepsilon_i) = \sigma^2$. The β 's are noticed to be the partial regression Coefficients i.e. $\beta_{0.12\dots(r)\dots k}$. The question is to derive estimators for β_r and test the null hypothesis $\beta_r = 0$ against the alternative $\beta_r \neq 0$.

(2-1) Let us try first as an example the simplest case where we have two independent variates. The mathematical model is

$$x_{oi} = \alpha + \beta_1 (x_{1i} - \bar{x}_1) + \beta_2 (x_{2i} - \bar{x}_2) + \epsilon_i \quad (3)$$

Let the best fitted line be of the form

$$x_{oi} = a + b_1 (x_{1i} - \bar{x}_1) + b_2 (x_{2i} - \bar{x}_2) \quad (4)$$

Using least square method, we minimize

$$S = \sum_{i=1}^n \left[x_{oi} - a - b_1 (x_{1i} - \bar{x}_1) - b_2 (x_{2i} - \bar{x}_2) \right]^2 \quad (5)$$

To get a, b_1, b_2 we should have

$$\left. \begin{aligned} \frac{\partial S}{\partial a} &= 0 \\ \frac{\partial S}{\partial b_1} &= 0 \\ \frac{\partial S}{\partial b_2} &= 0 \end{aligned} \right\} \quad (6)$$

In other words

$$a = \bar{x}_0 \quad (7)$$

$$b_1 \sum_i (x_{1i} - \bar{x}_1)^2 + b_2 \sum_i (x_{2i} - \bar{x}_2)(x_{1i} - \bar{x}_1) = \sum_i (x_{oi} - \bar{x}_0)(x_{1i} - \bar{x}_1) \quad (8)$$

$$\& \quad b_1 \sum_i (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) + b_2 \sum_i (x_{2i} - \bar{x}_2)^2 = \sum_i (x_{oi} - \bar{x}_0)(x_{2i} - \bar{x}_2) \quad (9)$$

If we write $\lambda_{rs} = \sum_{i=1}^n (x_{ri} - \bar{x}_r)(x_{si} - \bar{x}_s)$, then (8), (9) become

$$b_1 \lambda_{11} + b_2 \lambda_{21} = \lambda_{01} \quad (10)$$

$$b_1 \lambda_{12} + b_2 \lambda_{22} = \lambda_{02} \quad (11)$$

Solving (10), (11) we get

$$b_1 = \frac{\lambda_{01} \lambda_{22} - \lambda_{02} \lambda_{12}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} \quad (12)$$

$$b_2 = \frac{\lambda_{02} \lambda_{11} - \lambda_{01} \lambda_{21}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} \quad (13)$$

(2-2) We want to express b_1, b_2 in terms of our theoretical model. In other words, what are $\lambda_{01}, \lambda_{02}$ in terms of the β 's?

$$x_{0i} = \alpha + \beta_1 (x_{1i} - \bar{x}_1) + \beta_2 (x_{2i} - \bar{x}_2) + \beta_3 \quad (14)$$

$$\bar{x}_0 = \alpha + \bar{\beta} \quad (15)$$

$$\therefore x_{0i} - \bar{x}_0 = \beta_1 (x_{1i} - \bar{x}_1) + \beta_2 (x_{2i} - \bar{x}_2) + \beta_3 - \bar{\beta} \quad (16)$$

Now

$$\begin{aligned} \lambda_{01} &= \sum_i (x_{0i} - \bar{x}_0) (x_{1i} - \bar{x}_1) \\ &= \sum_{i=1}^n [\beta_1 (x_{1i} - \bar{x}_1) + \beta_2 (x_{2i} - \bar{x}_2) + (\beta_3 - \bar{\beta})] (x_{1i} - \bar{x}_1) \\ &= \beta_1 \lambda_{11} + \beta_2 \lambda_{21} + \lambda_{\beta_3 1} \end{aligned} \quad (17)$$

Similarly

$$\lambda_{02} = \sum_i (x_{0i} - \bar{x}_0) (x_{2i} - \bar{x}_2) = \beta_1 \lambda_{12} + \beta_2 \lambda_{22} + \lambda_{\beta_3 2} \quad (18)$$

$$\text{But } \lambda_{\beta_3 1} = \sum_i (\beta_3 - \bar{\beta}) (x_{1i} - \bar{x}_1) = \sum_i (x_{1i} - \bar{x}_1) \beta_3$$

$$\therefore E(\lambda_{\beta_3 1}) = 0 \quad (19)$$

$$\text{Also } \lambda_{\beta_3 2} = \sum_i (\beta_3 - \bar{\beta}) (x_{2i} - \bar{x}_2) = \sum_i (x_{2i} - \bar{x}_2) \beta_3$$

$$\therefore E(\lambda_{\beta_3 2}) = 0 \quad (20)$$

Again $E(\lambda_{z_1}^2) = E\left[\sum_i z_i^2 (x_{1i} - \bar{x}_1)^2\right]$
 $= \sum_i (x_{1i} - \bar{x}_1)^2 = \lambda_{11} \sigma^2$ (21)

Similarly $E(\lambda_{z_2}^2) = \lambda_{22} \sigma^2$ (22)

& $E(\lambda_{z_1} \lambda_{z_2}) = \lambda_{12} \sigma^2$ (23)

But $b_1 = \frac{\lambda_{01} \lambda_{22} - \lambda_{02} \lambda_{12}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2}$
 $= \beta_1 + \frac{\lambda_{22} \lambda_{z_1} - \lambda_{12} \lambda_{z_2}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} = \beta_1 + z_1$

Similarly $b_2 = \beta_2 + \frac{\lambda_{11} \lambda_{z_2} - \lambda_{12} \lambda_{z_1}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} = \beta_2 + z_2$

$E(b_1) = \beta_1$

& $E(b_2) = \beta_2$

Also $\text{Var}(b_1) = E(b_1 - \beta_1)^2 = E(z_1^2) = \frac{\lambda_{22}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} \sigma^2$ (24)

& $\text{Var}(b_2) = E(b_2 - \beta_2)^2 = E(z_2^2) = \frac{\lambda_{11}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} \sigma^2$ (25)

and $\rho(b_1, b_2) = \frac{E(b_1 - \beta_1)(b_2 - \beta_2)}{\sigma_{b_1} \sigma_{b_2}} = -\frac{\lambda_{12}}{\lambda_{11} \lambda_{22} - \lambda_{12}^2} \sigma^2$
 $= -\rho(x_1, x_2)$ (26)

(2-3) Partition of Total Sum of Squares

Total sum of squares $= \sum_i (x_{0i} - \bar{x}_0)^2 = \lambda_{00}$

$= \sum_{i=1}^n \left[\{x_{0i} - \bar{x}_0 - b_1(x_{1i} - \bar{x}_1) - b_2(x_{2i} - \bar{x}_2)\} + \{b_1(x_{1i} - \bar{x}_1) + b_2(x_{2i} - \bar{x}_2)\} \right]^2$

$= \sum_{i=1}^n \{x_{0i} - \bar{x}_0 - b_1(x_{1i} - \bar{x}_1) - b_2(x_{2i} - \bar{x}_2)\}^2$

$+ \sum_{i=1}^n \{b_1(x_{1i} - \bar{x}_1) + b_2(x_{2i} - \bar{x}_2)\}^2$

$+ 2 \sum_{i=1}^n \{x_{0i} - \bar{x}_0 - b_1(x_{1i} - \bar{x}_1) - b_2(x_{2i} - \bar{x}_2)\} \{b_1(x_{1i} - \bar{x}_1) + b_2(x_{2i} - \bar{x}_2)\}$

The 1st term is about multiple regression sum of squares (residual sum of squares) i.e.

$$\begin{aligned}
 & \sum_{i=1}^n \left\{ x_{i0} - \bar{x}_0 - b_1(x_{i1} - \bar{x}_1) - b_2(x_{i2} - \bar{x}_2) \right\}^2 \\
 &= \lambda_{00} + b_1^2 \lambda_{11} + b_2^2 \lambda_{22} - 2b_1 \lambda_{01} - 2b_2 \lambda_{02} + 2b_1 b_2 \lambda_{12} \\
 &= \lambda_{00} + b_1(b_1 \lambda_{11} + b_2 \lambda_{12} - \lambda_{01}) + b_2(b_1 \lambda_{12} + b_2 \lambda_{22} - \lambda_{02}) \\
 &\quad - \lambda_{01} b_1 - \lambda_{02} b_2 \\
 &= \lambda_{00} - b_1 \lambda_{01} - b_2 \lambda_{02}
 \end{aligned}$$

The 2nd term is due to regression sum of squares i.e.

$$\begin{aligned}
 & \sum_{i=1}^n \left\{ b_1(x_{i1} - \bar{x}_1) + b_2(x_{i2} - \bar{x}_2) \right\}^2 \\
 &= b_1^2 \lambda_{11} + b_2^2 \lambda_{22} + 2b_1 b_2 \lambda_{12} \\
 &= b_1(b_1 \lambda_{11} + b_2 \lambda_{12}) + b_2(b_2 \lambda_{22} + b_1 \lambda_{12}) \\
 &= b_1 \lambda_{01} - b_2 \lambda_{02}
 \end{aligned}$$

The last term i.e.

$$\begin{aligned}
 & 2 \sum_{i=1}^n \left\{ x_{i0} - \bar{x}_0 - b_1(x_{i1} - \bar{x}_1) - b_2(x_{i2} - \bar{x}_2) \right\} \left\{ b_1(x_{i1} - \bar{x}_1) + b_2(x_{i2} - \bar{x}_2) \right\} \\
 &= 2(b_1 \lambda_{01} - b_1^2 \lambda_{11} - b_1 b_2 \lambda_{12} + b_2 \lambda_{02} - b_1 b_2 \lambda_{12} - b_2^2 \lambda_{22}) \\
 &= 2 b_1(\lambda_{01} - b_1 \lambda_{11} - b_2 \lambda_{12}) + b_2(\lambda_{02} - b_1 \lambda_{12} - b_2 \lambda_{22}) \\
 &= 0
 \end{aligned}$$

In other words the last term vanishes and we are left with

Total sum of squares = Residual sum of squares + Due to regression sum of squares.

i.e.

$$\lambda_{00} = (\lambda_{00} - b_1 \lambda_{01} - b_2 \lambda_{02}) + (b_1 \lambda_{01} + b_2 \lambda_{02})$$

(2-4) The Expected Values

$$\therefore x_{0i} - \bar{x}_0 = \beta_1(x_{1i} - \bar{x}_1) + \beta_2(x_{2i} - \bar{x}_2) + z_i - \bar{z}$$

$$\therefore \lambda_{00} = \sum_{i=1}^n (x_{0i} - \bar{x}_0)^2$$

$$= \beta_1^2 \lambda_{11} + \beta_2^2 \lambda_{22} + \lambda_{zz} + 2\beta_1 \lambda_{z1} + 2\beta_2 \lambda_{z2} + 2\beta_1 \beta_2 \lambda_{12}$$

$$\therefore E(\lambda_{00}) = \beta_1^2 \lambda_{11} + \beta_2^2 \lambda_{22} + 2\beta_1 \beta_2 \lambda_{12} + (n-1)\sigma^2$$

$$\begin{aligned} \text{(i) Due to regression sum of squares} &= b_1 \lambda_{01} + b_2 \lambda_{02} \\ &= (\lambda_{22} \lambda_{01}^2 - 2 \lambda_{12} \lambda_{01} \lambda_{02} + \lambda_{11} \lambda_{02}^2) / (\lambda_{11} \lambda_{22} - \lambda_{12}^2) \end{aligned}$$

$$\begin{aligned} \text{But } E(\lambda_{01}^2) &= E(\beta_1 \lambda_{11} + \beta_2 \lambda_{21} + \lambda_{z1})^2 \\ &= (\beta_1 \lambda_{11} + \beta_2 \lambda_{21})^2 + \lambda_{11} \sigma^2 \end{aligned}$$

$$\text{Also } E(\lambda_{02}^2) = (\beta_1 \lambda_{21} + \beta_2 \lambda_{22})^2 + \lambda_{22} \sigma^2$$

$$\& E(\lambda_{01} \lambda_{02}) = (\beta_1 \lambda_{11} + \beta_2 \lambda_{21})(\beta_1 \lambda_{21} + \beta_2 \lambda_{22}) + \lambda_{12} \sigma^2$$

\therefore Expected value of due to regression sum of squares

$$= E(b_1 \lambda_{01} + b_2 \lambda_{02}) = \beta_1^2 \lambda_{11} + \beta_2^2 \lambda_{22} + 2\beta_1 \beta_2 \lambda_{12} + 2\sigma^2$$

$$\text{(ii) Residual sum of squares} = \lambda_{00} - (b_1 \lambda_{01} + b_2 \lambda_{02})$$

$$\therefore \text{Expected value of Residual sum of squares} = (n-3)\sigma^2$$

Regression analysis table

Source of Variation	Sum of Squares	Degrees of freedom	Expected value of mean sum of squares
Due to multiple regression	$b_1 \lambda_{01} + b_2 \lambda_{02}$	2	$\sigma^2 + \frac{1}{2}(\beta_1^2 \lambda_{11} + \beta_2^2 \lambda_{22} + 2\beta_1 \beta_2 \lambda_{12})$
Residual	$\lambda_{00} - b_1 \lambda_{01} - b_2 \lambda_{02}$	$n-3$	σ^2
Total	λ_{00}	$n-1$	

From the above table, the mean sum of squares of residual is an unbiased estimator for σ^2 i.e.

$$\hat{\sigma}^2 = ms_{\text{resid.}}$$

(2-5) Tests for the significance of b_1, b_2

Again the normal equations (10), (11) can be written in a matrix form

$$\begin{bmatrix} \lambda_{11} & \lambda_{21} \\ \lambda_{12} & \lambda_{22} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} \lambda_{01} \\ \lambda_{02} \end{bmatrix} \quad (27)$$

$$\therefore \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} \lambda_{11} & \lambda_{21} \\ \lambda_{12} & \lambda_{22} \end{bmatrix}^{-1} \begin{bmatrix} \lambda_{01} \\ \lambda_{02} \end{bmatrix}$$

$$= \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix} \begin{bmatrix} \lambda_{01} \\ \lambda_{02} \end{bmatrix}$$

where $c_{11} = \frac{\lambda_{22}}{\Delta}$

$$c_{22} = \lambda_{11} / \Delta$$

$$c_{12} = c_{21} = -\lambda_{12} / \Delta$$

$$\& \Delta = \lambda_{11} \lambda_{22} - \lambda_{12}^2$$

Hence equations (24), (25), (26) can be rewritten as

$$\begin{aligned}\text{Var}(b_1) &= C_{11} \sigma^2 \\ \text{Var}(b_2) &= C_{22} \sigma^2\end{aligned}\tag{28}$$

Now to test the null hypothesis $\beta_1 = 0$ against the alternative $\beta_1 \neq 0$, we apply the test criterion

$$t_{b_1} = b_1 / \tilde{\sigma}_{b_1}$$

where $\tilde{\sigma}_{b_1} = C_{11} \text{ms}_{\text{resd.}}$

Similarly to test the null hypothesis $\beta_2 = 0$ against the alternative $\beta_2 \neq 0$ we apply

$$t_{b_2} = b_2 / \tilde{\sigma}_{b_2}$$

where $\tilde{\sigma}_{b_2} = C_{22} \text{ms}_{\text{resd.}}$

To test for the significance of the difference between b_1, b_2 we have

$$\begin{aligned}\tilde{\sigma}_{b_1 - b_2}^2 &= (C_{11} + C_{22} - 2C_{12}) \sigma^2 \\ &= (C_{11} + C_{22} - 2C_{12}) \text{ms}_{\text{resd.}}\end{aligned}$$

$$\therefore t = (b_1 - b_2) / \tilde{\sigma}_{b_1 - b_2}$$

The no. of degrees of freedom for the statistic "t" is $n-3$, the same as that for the residual sum of squares.

(3) The above procedure can be generalized to more than two independent variables.

Let

$$x_{oi} = a + b_1(x_{1i} - \bar{x}_1) + \dots + b_r(x_{ri} - \bar{x}_r) + \dots + b_k(x_{ki} - \bar{x}_k) \quad (1)$$

where $i = 1, 2, \dots, n$ denotes the number of observations
and $r = 1, 2, \dots, k$ denotes the number of independent variables.

The mathematical model is

$$x_{oi} = \alpha + \beta_1(x_{1i} - \bar{x}_1) + \dots + \beta_r(x_{ri} - \bar{x}_r) + \dots + \beta_k(x_{ki} - \bar{x}_k) + \beta_i \quad (2)$$

where $E(\beta_i) = 0$ & $\text{Var}(\beta_i) = \sigma^2$

Minimizing

$$S = \sum_{i=1}^n \left[x_{oi} - a - b_1(x_{1i} - \bar{x}_1) - b_2(x_{2i} - \bar{x}_2) - \dots - b_k(x_{ki} - \bar{x}_k) \right]^2$$

we get

$$a = \bar{x}_0$$

(3)

$$b_1 \lambda_{11} + b_2 \lambda_{21} + \dots + b_r \lambda_{r1} + \dots + b_k \lambda_{k1} = \lambda_{01}$$

$$b_1 \lambda_{12} + b_2 \lambda_{22} + \dots + b_r \lambda_{r2} + \dots + b_k \lambda_{k2} = \lambda_{02}$$

$$\begin{array}{ccccccc} \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ b_1 \lambda_{1k} + b_2 \lambda_{2k} + \dots & b_r \lambda_{rk} & \dots & b_k \lambda_{kk} & = & \lambda_{0k} \end{array}$$

Where λ_{rs} is as defined before

Now the above set of equations (4) can be rewritten in a matrix form as

$$\begin{bmatrix} \lambda_{11} & \lambda_{21} & \dots & \lambda_{r1} & \dots & \lambda_{k1} \\ \lambda_{12} & \lambda_{22} & \dots & \lambda_{r2} & \dots & \lambda_{k2} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \lambda_{1k} & \lambda_{2k} & \dots & \lambda_{rk} & \dots & \lambda_{kk} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix} = \begin{bmatrix} \lambda_{o1} \\ \lambda_{o2} \\ \vdots \\ \lambda_{ok} \end{bmatrix} \quad (5)$$

This is simply written

$$[\lambda]^{(b)} = (\lambda_o) \quad (6)$$

where $[\lambda]$ is a square matrix, $(b), (\lambda_o)$ are column matrices.

$$\therefore (b) = [\lambda]^{-1} (\lambda_o) \quad (7)$$

From (2)

$$\bar{x}_o = \alpha + \bar{z}$$

$$\therefore x_{oi} - \bar{x}_o = \beta_1(x_{1i} - \bar{x}_1) + \beta_2(x_{2i} - \bar{x}_2) + \dots + \beta_k(x_{ki} - \bar{x}_k) + (\bar{z}_i - \bar{z})$$

Multiplying both sides by $x_{1i} - \bar{x}_1, x_{2i} - \bar{x}_2, \dots, (x_{ki} - \bar{x}_k)$

successively and summing over i in each case, we get

$$\begin{aligned} \lambda_{o1} &= \beta_1 \lambda_{11} + \beta_2 \lambda_{21} + \dots + \beta_k \lambda_{k1} + \lambda_{z1} \\ \lambda_{o2} &= \beta_1 \lambda_{12} + \beta_2 \lambda_{22} + \dots + \beta_k \lambda_{k2} + \lambda_{z2} \\ &\vdots \\ \lambda_{ok} &= \beta_1 \lambda_{1k} + \beta_2 \lambda_{2k} + \dots + \beta_k \lambda_{kk} + \lambda_{zk} \end{aligned} \quad (8)$$

These set of equations can be written in a matrix form

$$\begin{bmatrix} \lambda_{o1} \\ \lambda_{o2} \\ \vdots \\ \vdots \\ \lambda_{ok} \end{bmatrix} = \begin{bmatrix} \lambda_{11} & \lambda_{21} & \dots & \lambda_{k1} \\ \lambda_{12} & \lambda_{22} & \dots & \lambda_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{1k} & \lambda_{2k} & \dots & \lambda_{kk} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \lambda_{z1} \\ \lambda_{z2} \\ \vdots \\ \vdots \\ \lambda_{zk} \end{bmatrix} \quad (9)$$

This is simply written as

$$(\lambda_o) = [\lambda](\beta) + (\lambda_z) \quad (10)$$

$$\text{where } \bar{\delta}(\lambda_z) = 0 \quad (11)$$

Substituting from (10) in (7)

$$\begin{aligned} (b) &= [\lambda]^{-1} \{ [\lambda](\beta) + \lambda_z \} \\ &= (\beta) + [\lambda]^{-1}(\lambda_z) \end{aligned} \quad (12)$$

$$\text{Hence } \bar{\delta}(b_r) = \beta_r \quad (13)$$

From (12)

$$(b - \beta) = [\lambda]^{-1}(\lambda_z) \quad (14)$$

$$\therefore (b - \beta)' = (\lambda_z)' [\lambda]^{-1}$$

$$(b - \beta)(b - \beta)' = [\lambda]^{-1} \lambda_z \lambda_z' [\lambda]^{-1} \quad (15)$$

But

$$\lambda_z \lambda_z' = \begin{pmatrix} \lambda_{z1} \\ \lambda_{z2} \\ \vdots \\ \lambda_{zk} \end{pmatrix} (\lambda_{z1} \quad \lambda_{z2} \quad \dots \quad \lambda_{zk})$$

(152)

$$= \begin{bmatrix} \lambda_{z1}^2 & \lambda_{z1}\lambda_{z2} & \dots & \lambda_{z1}\lambda_{zk} \\ \lambda_{z2}\lambda_{z1} & \lambda_{z2}^2 & \dots & \lambda_{z2}\lambda_{zk} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{zk}\lambda_{z1} & \lambda_{zk}\lambda_{z2} & \dots & \lambda_{zk}^2 \end{bmatrix} \quad (16)$$

Since $\lambda_{zr} = \sum_{i=1}^n (\bar{z}_i - \bar{z}) (x_{ri} - \bar{x}_r)$

$$= \sum_{i=1}^n (x_{ri} - \bar{x}_r) \bar{z}_i$$

$$\therefore E(\lambda_{zr}^2) = \lambda_{rr} \sigma^2 \quad (17)$$

Also $E(\lambda_{zr}\lambda_{zs}) = \lambda_{rs} \sigma^2 \quad (18)$

Consequently

$$E(\lambda_z \lambda_z') = \begin{bmatrix} \lambda_{11}\sigma^2 & \lambda_{12}\sigma^2 & \dots & \lambda_{1k}\sigma^2 \\ \lambda_{21}\sigma^2 & \lambda_{22}\sigma^2 & \dots & \lambda_{2k}\sigma^2 \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{k1}\sigma^2 & \lambda_{k2}\sigma^2 & \dots & \lambda_{kk}\sigma^2 \end{bmatrix} = [\lambda] \sigma^2 \quad (19)$$

$$\begin{aligned} E(b - \beta)(b - \beta)' &= E[\lambda^{-1} \lambda_z \lambda_z' \lambda^{-1}] \\ &= \lambda^{-1} \lambda \sigma^2 \lambda^{-1} = [\lambda]^{-1} \sigma^2 \end{aligned} \quad (20)$$

(153)

But

$$(b-\beta)(b-\beta)' = \begin{pmatrix} b_1 - \beta_1 \\ b_2 - \beta_2 \\ \vdots \\ b_k - \beta_k \end{pmatrix} (b_1 - \beta_1 \quad b_2 - \beta_2 \quad \dots \quad b_k - \beta_k)$$

$$= \begin{bmatrix} (b_1 - \beta_1)^2 & (b_1 - \beta_1)(b_2 - \beta_2) & \dots & (b_1 - \beta_1)(b_k - \beta_k) \\ (b_2 - \beta_2)(b_1 - \beta_1) & (b_2 - \beta_2)^2 & \dots & (b_2 - \beta_2)(b_k - \beta_k) \\ \vdots & \vdots & \ddots & \vdots \\ (b_k - \beta_k)(b_1 - \beta_1) & (b_k - \beta_k)(b_2 - \beta_2) & \dots & (b_k - \beta_k)^2 \end{bmatrix} \quad (21)$$

If $[\lambda]^{-1}$ is written as

$$[\lambda]^{-1} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1k} \\ c_{21} & c_{22} & \dots & c_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ c_{k1} & c_{k2} & \dots & c_{kk} \end{bmatrix} \quad (22)$$

then, substituting in (20) and comparing with the expected values of (21), we get

$$\sum (b_r - \beta_r)^2 = c_{rr} \sigma^2$$

$$\text{i.e.} \quad \text{Var}(b_r) = c_{rr} \sigma^2 \quad (23)$$

Also
$$\mathcal{E} (b_r - \beta_r)(b_s - \beta_s) = c_{rs} \sigma^2 \quad (24)$$

Similarly

$$\begin{aligned} \text{Var} (b_r - b_s) &= \mathcal{E} [b_r - b_s - (\beta_r - \beta_s)]^2 \\ &= \mathcal{E} [(b_r - \beta_r) - (b_s - \beta_s)]^2 \\ &= \mathcal{E} (b_r - \beta_r)^2 + \mathcal{E} (b_s - \beta_s)^2 - 2 \mathcal{E} (b_r - \beta_r)(b_s - \beta_s) \\ &= (c_{rr} + c_{ss} - 2 c_{rs}) \sigma^2 \end{aligned} \quad (25)$$

The regression analysis table

Source of variation	Sum of squares	Degrees of freedom
Due to multiple Regression	$b_1 \lambda_{01} + b_2 \lambda_{02} + \dots + b_k \lambda_{0k}$	k
Residual	$\lambda_{00} - (b_1 \lambda_{01} + b_2 \lambda_{02} + \dots + b_k \lambda_{0k})$	n - k - 1
Total	00	n-1

Following the same procedure as in the case of two independent variates, it can be shown that

$$\mathcal{E} (\text{Mean sum of Squares of residuals}) = \sigma^2$$

$$\text{i.e. } \hat{\sigma}^2 = \text{MS}_{\text{resid.}} = \frac{\text{Sum of Squares of Residuals}}{n - k - 1} \quad (26)$$

To test the null hypothesis $\beta_r = 0$ against the alternative $\beta_r \neq 0$, we apply the test criterion

$$t = \frac{b_r - \beta_r}{\hat{\sigma}_{b_r}} = b_r / \hat{\sigma}_{b_r} \quad (27)$$

where $\tilde{\sigma}_{b_r} = \sqrt{C_{rr} MS_{resid}}$ (28)

and with $n - k - 1$ degrees of freedom

Also to test the null hypothesis $\beta_r - \beta_s = 0$ against the alternative $\beta_r - \beta_s \neq 0$ we apply the test criterion.

$$t = \frac{b_r - b_s}{\tilde{\sigma}_{b_r - b_s}} \quad (29)$$

where $\tilde{\sigma}_{b_r - b_s} = \sqrt{(C_{rr} + C_{ss} - 2C_{rs}) MS_{resid}}$ (30)

Example (1) Consider the following data.

x_0	:	12	1	2	4	12	13	3	1
x_1	:	15	1	6	23	26	27	20	2
x_2	:	6	11	9	6	4	2	7	11

$$\sum x_0 = 48, \quad \sum x_1 = 120, \quad \sum x_2 = 56$$

Sum of Squares and Sum of products
in deviate form

	x_0	x_1	x_2
x_0	200	290	-105
x_1		800	-227
x_2			72

$$\Delta = \begin{vmatrix} 800 & -227 \\ -227 & 72 \end{vmatrix} = 6071$$

$$C_{11} = 72/6071 = 0.01186$$

$$C_{22} = 800/6071 = 0.13177$$

$$C_{12} = 227/6071 = 0.03739$$

$$b_1 = C_{11} \ o_1 + C_{12} \ o_2 = -0.4866$$

$$b_2 = C_{21} \ o_1 + C_{22} \ o_2 = -2.9927$$

Sum of squares due to multiple regression

$$= b_1 \lambda_{o1} + b_2 \lambda_{o2} = 173.1195$$

Total Sum of squares = 200.

Regression analysis table

Source of variation	Sum of Squares	Degree of freedom	Mean Sum of Squares
Due to multiple regression	173.1195	2	86.5597
Residual	26.8805	5	5.2761
Total	200.0000	7	

$$R = 86.5597 / 5.2761 = 16.4$$

$$F_{2,5,10/0} = 13.27$$

The regression is significant

$$\text{Var}(b_1) = C_{11} \text{MS}_{\text{resid}} = 0.0625$$

$$\text{Var}(b_2) = C_{22} \text{MS}_{\text{resid}} = 0.6952$$

$$\tilde{\sigma}_{b_1} = 0.25 \quad \& \quad \tilde{\sigma}_{b_2} = 0.83$$

$$\text{Hence } t_{b_1} = 1.96$$

$$t_{b_2} = 3.61$$

$$\text{But } t_{5, 5\%} = 2.571$$

There is an evidence that $\beta_1 = 0$ while $\beta_2 \neq 0$.

Example (2) The following data represent: the score on the final exam (x_0), number of daily study hours (x_1), intelligent quotient (x_2) and the average score through the year (x_3) for a sample of 10 students.

x_0	:	7	8	5	11	17	18	13	11	9	10
x_1	:	5	6	4	7	8	9	7	6	5	4
x_2	:	10	10	11	14	15	13	10	12	14	11
x_3	:	7	10	6	15	16	18	15	10	7	12

$$x_0 = 109, \quad x_1 = 61, \quad x_2 = 120, \quad x_3 = 116$$

Sum of squares and sum of products in
derivate form

	x_0	x_1	x_2	x_3
x_0	154.9	55.1	38.0	143.6
x_1		24.9	13.0	54.4
x_2			32.0	28.0
x_3				162.4

To get the inverse of the matrix formed from the squares and products of x_1, x_2, x_3 we use the Square root method as follows

x_1	x_2	x_3	I		
24.9	13.0	54.4	1	0	0
	32.0	28.0	0	1	0
		162.4	0	0	1
4.98999	2.60522	10.90183	0.20040	0	0
	5.02124	-0.07999	-0.10398	0.19915	0
		6.59877	-0.33234	0.00241	0.15154
$[\lambda]^{-1} =$			0.16142	-0.02151	-0.05036
				0.03966	0.00036
					0.02296

where $[\lambda] = \begin{bmatrix} 24.9 & 13.0 & 54.4 \\ & 32.0 & 28.0 \\ & & 162.4 \end{bmatrix}$

$$[\lambda]^{-1} = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix} = \begin{bmatrix} 0.16142 & -0.02151 & -0.05036 \\ & 0.03966 & 0.00036 \\ & & 0.02296 \end{bmatrix}$$

$$b_1 = c_{11}\lambda_{01} + c_{12}\lambda_{02} + c_{13}\lambda_{03} = 0.8451$$

$$b_2 = c_{21}\lambda_{01} + c_{22}\lambda_{02} + c_{23}\lambda_{03} = 0.3736$$

$$b_3 = c_{31}\lambda_{01} + c_{32}\lambda_{02} + c_{33}\lambda_{03} = 0.5359$$

Sum of Squares due to regression

$$= b_1 \lambda_{01} + b_2 \lambda_{02} + b_3 \lambda_{03} = 137.7197$$

$$\& \text{ Total sum of Squares } \lambda_{00} = 154.9$$

Regression analysis table

Source of variation	Sum of Squares	Degrees of freedom	Mean sum of squares
Due to regression	137.7197	3	45.9066
Residual	17.1803	6	2.8634
	154.9000	9	

$$R = 45.9066 / 2.8634 = 16.03$$

$$F_{3,6,1} \% = 9.78$$

The regression is Significant

$$\text{Var}(b_1) = C_{11} \text{MS}_{\text{resid}} = 0.4622$$

$$\text{Var}(b_2) = C_{22} \text{MS}_{\text{resid}} = 0.1136$$

$$\text{Var}(b_3) = C_{33} \text{MS}_{\text{resid}} = 0.0657$$

$$\tilde{G}_{b_1} = 0.680, \tilde{G}_{b_2} = 0.337, \tilde{G}_{b_3} = 0.256$$

$$t_1 = 1.24, t_2 = 1.11, t_3 = 2.09$$

$$\text{But } t_{6,5\%} = 2.447$$

There is an eridence that each of $\beta_1, \beta_2, \beta_3$ is equal to zero.

Exercise In a problem to study the factors which may be used in estimating cotton and cotton seeds production index, the data collected are given in the following table

Y	D	A _r	A _s	L	F _N	F _p
1945	636	504	525	1016	234	65
46	734	585	655	1014	189	1
47	769	627	683	1009	483	177
48	1074	796	805	1002	677	214
49	1051	852	931	999	765	258
50	1928	1154	1116	999	924	703
51	976	1108	1133	1001	999	1003
52	1200	1115	1117	998	997	1264
53	857	759	745	999	1011	1037
54	938	864	889	1003	993	755
55	900	962	992	1010	943	1047
56	875	891	926	1017	840	1254
57	1091	990	1042	1027	934	1427
58	1200	1142	1093	1032	1184	1522
59	1229	1017	996	1040	993	1475
60	1285	1065	1035	1048	981	1941

where Y the years

P Cotton & Cotton seeds production index

A_r Area cultivated with cotton weighted with the relative product of 1 feddan in upper, middle and lower Egypt.

A_s Area cultivated with cotton weighted with the relative product of 1 feddan cultivated with long, long medium and medium Staple cotton

L Labour index

F_N, F_P nitrate and phosphate index

Make the required regression analysis.

Exercise The following represents data on pigs. Thus animals were kept on adequate vations during the period of approximately uniform growth rate. We wish to inquire the amount of information about rate of gain which is furnished in advance by the two independent variates

Pig No.	Initial age x_1 (days)	Weight x_2 (pounds)	Rate of gain Y (pounds/day)
1	78	61	1.40
2	90	59	1.79
3	94	76	1.72
4	71	50	1.47
5	99	61	1.26
6	80	54	1.28
7	83	57	1.34
8	75	45	1.57
9	62	41	1.57
10	67	40	1.26
11	78	74	1.61
12	99	75	1.31
13	80	64	1.12
14	75	48	1.35
15	94	62	1.29
16	91	42	1.24
17	75	52	1.29
18	63	43	1.43
19	62	50	1.29
20	67	40	1.26
21	78	80	1.67
22	83	61	1.41
23	79	62	1.73
24	70	47	1.23
25	85	59	1.49
26	83	42	1.22
27	71	47	1.39

Pig No.	Initial age x_1 (days)	weight x_2 (pounds)	Rate of gain Y (pounds/day)
28	66	42	1.39
29	67	40	1.56
30	67	40	1.36
31	77	62	1.40
32	71	55	1.47
33	78	62	1.37
34	70	43	1.15
35	95	57	1.22
36	96	51	1.48
37	71	41	1.31
38	63	40	1.27
39	62	45	1.22
40	67	39	1.36

Exercises

(I) The mean coefficient of volume of expansion of water was determined by 21 different observers at 3 different temperatures. It is known that the determination of this coefficient will become more difficult as the temperature increases owing to the formation of bubbles in the water. Fine out from the data given below whether the determinations of the coefficient are significantly more variable as the temp. range of the water increases. Estimate a mean coefficient for each temp. and give confidence limits for the true coefficient

Temp. observer	Coefficient of Volume Expansion of Water (10^6)						
	0-40	40-60	60-80	Temp Observer	0-40	40-60	60-80
1	239	522	517	11	249	386	612
2	274	474	584	12	301	404	437
3	295	512	606	13	175	504	561
4	256	417	447	14	260	361	580
5	354	327	542	15	301	491	285
6	213	361	627	16	445	449	453
7	264	504	614	17	335	436	520
8	236	572	632	18	166	346	372
9	302	390	505	19	270	449	566
10	269	526	667	20	254	560	607
				21	343	468	560

2- Five independent samples of 25 individuals each were drawn from a normal population of unknown mean and variance. The 98% confidence intervals for the population mean were obtained as follows from each sample

Sample No.Confidence limits

1
2
3
4
5

105.7 - 116.4
104.3 - 113.8
107.2 - 116.0
100.6 - 115.4
101.7 - 118.1

Obtain a 98% confidence interval for the population standard deviation

3- Using the symbols given in the table below, derive expressions for the probabilities that of four men aged exactly 75, 80, 85 & 90 respectively

- (i) all will attain age 95
- (ii) all will die before attaining age 95
- (iii) at least one will survive 10 years
- (iv) none will die between age 90 & 95.

4- Rutherford and Geiger counted the number of alpha - particles emitted from a disc in 2608 periods of 7.5 seconds duration. The frequencies are given below

Number per period	frequency	Number per period	frequency
0	57	8	45
1	203	9	27
2	383	10	10
3	525	11	4
4	532	12	2
5	408	13	0
6	273	14	0
7	139		

Compare the relative frequencies with the corresponding probabilities of the fitted "Poisson distribution"

5- Show that the sum of two Poisson variates is itself a Poisson variate with mean equal to the sum of the separate means

6- Number of individual incomes in different ranges of net income assessed in 1945-46:

Range of Income after tax (x)	Number of Incomes
150 - 500	13,175,000
500 - 1000	652,000
1000 - 2000	137,500
2000 and over	
Total	140,000,000

@ Exact age	75	80	85	90	95
Probability of surviving 5 years	P_0	P_1	P_2	P_3	P_4

Assume that this distribution of incomes $f(x)$ is linked with the normal distribution

$$N(t) = \frac{1}{2} e^{-\frac{1}{2} t^2}$$

by the relationship

$$\int_{-\infty}^t N(t) dt = \int_{150}^x f(x) dx$$

where $t = a \log (x - 150) + b$

Obtain estimates for "a" & "b" from the data, and find the number of incomes between 250 and 500.

7- The following results were obtained in 4 independent samplings

(1)	6	14	12	6	2	5
(2)	10	17	6	19	19	16
(3)	11	11	19	23	8	17
(4)	19	2	29	16	14	20

Carryout an analysis of variance on these data

8- Twelve dice were thrown 26306 times and a "5" or a "b" was counted as a success. The number of successes in each throw was noted, with the following results

Number of successes	frequency	Number of successes	frequency
0	185	6	3067
1	1149	7	1331
2	3265	8	403
3	5475	9	105
4	6114	10	18
5	5194		
		Total	26306

Is there evidence that the dice are biased?

9- The following table gives the distribution of the length, measured in cm. of 294 eggs

Length	frequency	length	
3.5	1	4.2	54
3.6	1	4.3	34
3.7	6	4.4	12
3.8	20	4.5	6
3.9	35	4.6	1
4.0	53	4.7	2
4.1	69		
		Total	294

Test whether these results are consistent with the hypothesis that egg length is normally distributed

10- A certain type of surgical operation can be performed either with a local anaesthetic or with a general anaesthetic. Results are given below

	alive	Dead
Local	511	24
General	173	21

Test for any difference in the mortality rates associated with the different types of anaesthetic.

11- A Ministry of Labour Memorandum on carbon Monoxide Poisoning gives the following data on accidents due to gassing by carbon monoxide

	1941	1942	1943	Total
At blast furnaces	24	20	19	63
At gas producers	28	34	41	103
At gas ovens and works	26	26	10	62
In distribution and use gas	80	108	123	311
Miscellaneous sources	68	51	32	151
Total	226	239	225	690

Is there significant association between the site of the accident and the year.

(12) From the following data find the regression equation of x_1 on x_2 & x_3

x_1	x_2	x_3
5	2	21
3	4	21
2	2	15
4	2	17
3	3	20
1	2	13
8	4	32

(13) If s_1^2 , s_2^2 are the variances in two independent samples of the same size taken from a common normal population, determine the distribution of $s_1^2 + s_2^2$

(14) A random variate x is known to have the distribution

$$p(x) = C(1 + \frac{x}{a})^{m-1} e^{-\frac{mx}{a}} \quad -a \leq x \leq \infty$$

Find the constant C and the 1st four moments of x . Prove that

$$8\beta_2 - 9\beta_1 = 4$$

(15) If $u = ax + by$ and $v = bx - ay$, where x, y represent deviations from respective means and if ρ is the correlation coefficient between x, y , but u, v are uncorrelated, show that

$$\sigma_u \sigma_v = (a^2 + b^2) \sigma_x \sigma_y \sqrt{1 - \rho^2}$$

(16) The following table gives the distribution according to age of becoming a widower

Age	freq.	Age	freq.
18 - 22	5	53 - 57	88
23 - 27	43	58 - 62	76
28 - 32	108	63 - 67	82
33 - 37	107	68 - 72	59
38 - 42	130	73 - 77	41
43 - 47	109	78 - 82	16
48 - 52	115	83 - 87	8

Calculate β_1 & β_2

(17) Examine the following data showing how age is associated with size of farm

size of farm	Percentage of occupiers whose ages are						
	under 25	25-34	35-44	45-54	55-64	65-74	over 74
	%	%	%	%	%	%	%
20-	2	12	25	28	23	8	2
50-	1	7	21	36	26	7	2
100-	1	10	29	28	24	6	2
150-	1	9	21	36	21	10	2
over 300	0	5	19	39	22	11	4

(18) Four machines A, B, C, D, are producing large numbers of small articles. It is known that the average proportions of defective articles produced by the machines are

$$A = 1\%, \quad B = 1\frac{1}{2}\%, \quad C = 1\frac{1}{2}\%, \quad D = 2\frac{1}{2}\%$$

In a group of 20 articles known to have been produced by the same machine, one defective article is found. What is the probability that this group was produced by machine D.

State carefully any assumption which you make.

(19) x_1, x_2, \dots, x_n are random variables each with the same expected value μ , and the same variance σ^2 . The correlation between any two of the x 's is ρ . Show that

$$(i) \quad \left\{ \sigma(\bar{x}) \right\}^2 = \frac{\sigma^2}{n} + (1 - \frac{1}{n})\rho\sigma^2$$

$$(ii) \quad E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] = (n-1)(1-\rho)\sigma^2$$

$$(iii) \quad \rho > -\frac{1}{n-1}$$

(20) Suppose that the probability of death within a year for a person aged exactly x is q_x . If there are n_x persons aged exactly x , what is the probability that r of these persons will die before reaching exact age $(x+1)$?

(21) Assuming that in 1000 items, 20 are defective, what is the probability that in a random sample of 100 items x items will be defective.

(22) Adie has $n+1$ faces, numbered $0, 1/n, 2/n, \dots, (n-1)/n, 1$ respectively. Assuming it is equally likely to fall with any one face uppermost, find the expected value and standard deviation of the random variable corresponding to the number on the uppermost face.

(23) For a high voltage network, uniform cables of great tensile strength are required. Each cable is composed of wires which are manufactured in one length. In order to examine the tensile strength a sample is taken from each wire and tested. The table shows the tensile strength of each wire of 5 cables each with 10 wires

Cable 1:	345	327	335	338	330	334	335	340	337	342
Cable 2:	329	327	332	348	337	328	328	330	335	334
Cable 3:	340	330	325	328	338	332	335	340	336	339
Cable 4:	328	344	342	350	335	332	328	340	335	337
Cable 5:	347	341	345	340	350	346	345	342	340	339

Carry out the analysis of variance (subtract an arbitrary origin)

(24) The following table gives the number of yeast cells in 400 squares of a haemocytometer:

No. of cells	0	1	2	3	4	5	Total
Frequency	213	128	37	18	3	1	400

Fit a negative binomial

(Hint: mean=0.68, variance=0.81. $np=0.68$, $npq=0.81$, then $q=1.19$ & $p=-0.19$ and $n=-3.59$)

(25) Fit a normal distribution to the distribution given in the following table and test for the goodness of fit.

Height	≤59	60-	61-	62-	63-	64-	65-	66-	67-	68-
Freq.	23	169	439	1030	2116	3947	5965	8012	9089	8763
Height	69-	70-	71-	72-	73-	74-	75-	76-	77&over	
Freq.	7132	5314	3320	1884	876	383	153	63	25	

REFERENCES

- 1- David, F.N. Probability theory for statistical methods
- 2- Edwards, Allen L. Experimental design in Psychological research
- 3- Fisher, R.A. Statistical methods for research workers
- 4- Gouldin, Cyril H. Methods of statistical analysis
- 5- Hald, A. Statistical theory with engineering applications
- 6- Johnson, N.L. "Statistics" An intermediate text book
& Tetley, H. Vol. 1 & II
- 7- Kendall, M.G. The advanced theory of statistics
& Stuart, A. Vol. 1 & II
- 8- Lindquist, E.F. Design and analysis of experiments in Psychology & Education.
- 9- Snedecor, George W. "Statistical methods"
- 10- Rosander, A.C. Elementary principles of statistics
- 11- Walker, Helen M. "Statistical inference"
& Lev, Joef
- 12- Wilks, S.S. Mathematical statistics

Remark

The material in this note is still in a tentative form. Corrections, criticisms and expressions of other points of view on the teaching and presentation of the course in this note will be gratefully received.

Moharram W. Mahmoud