

UNIT-1

31

Sample Variance

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

$\bar{x} \rightarrow$ mean = $\frac{\text{sum of all obs}}{\text{No of obs}}$
 $x_i \rightarrow$ obs.
 $n \rightarrow$ Total no. of obs

Sample standard Deviation

$$S = \sqrt{\text{variance}}$$

Correlation

↑ true ↑ ne ↑ Perfect
 ↑↑ (equal prop)

Karl Pearson coeff of correlation

$$r = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y} = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$= \frac{n \left[\sum (x_i - \bar{x})(y_i - \bar{y}) \right]}{n \left[\sum (x_i - \bar{x})^2 * \sum (y_i - \bar{y})^2 \right]^{1/2}}$$

$$X = x - \bar{x}$$

$$Y = y - \bar{y}$$

Covariance

$$\text{Cov}_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\text{No. of data values} - 1}$$

$r \in (-1, 1]$ if $r=0$ not correlated
 if $r=1$ +vely perfect correlated
 if $r=-1$ -vely perfect correlated

Direct method for corr coeff

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} * \sqrt{n \sum y^2 - (\sum y)^2}}$$

(small x & y)

Regression

stepping back towards the normal

Regression line Y on X Regression line X on Y

Regression line Y on X

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

Where $r \frac{\sigma_x}{\sigma_y} = b_{yx}$ is the regression coefficient of y on x

σ_x = standard deviation of x
 σ_y = standard deviation of y
 \bar{x} = mean of X
 \bar{y} = mean of Y

2. Regression line X on Y

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

Where $r \frac{\sigma_x}{\sigma_y} = b_{xy}$ is the regression coefficient of X on Y

\bar{x} = mean of x
 \bar{y} = mean of y
 σ_x = SD of x
 σ_y = SD of y

Standard Error

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

Coeff of variability

$$Cv = \frac{\sigma_x^2}{\bar{x}^2} = \frac{(b_{xy})^2}{\bar{x}^2} = \frac{\sigma_x^2}{\bar{x}^2} = \frac{\sigma_y^2}{\bar{y}^2} = \frac{(b_{yx})^2}{\bar{y}^2}$$

we have $b_{xy} = r \sigma_x / \sigma_y$

Rank correlation is measured by spearman rank correlation coeff

$$\rho = 1 - \left[\frac{6 * \left(\sum_{i=1}^n d_i^2 \right)}{n(n^2-1)} \right] \quad [-1, 1]$$

$n \rightarrow$ no. of obs.

$x_i \rightarrow$ Rank of x

$y_i \rightarrow$ Rank of y

$d_i = x_i - y_i$

Repetition of Rank Correlation

$$\rho = 1 - \left[\frac{6 * \left(\sum d_i^2 + C.F. \right)}{n(n^2-1)} \right]$$

CF = Correlation factor

CF = CF in X series + CF in Y series

$$CF = \frac{\sum m(m^2-1)}{12} \quad m \rightarrow \text{no. of times the rank is repeated}$$

Limit of correlation

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\Rightarrow r^2_{xy} = \left(\frac{1}{n} \right)^2 \left[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right]^2 = \frac{\frac{1}{n^2} \left(\sum_{i=1}^n a_i b_i \right)^2}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 * \frac{1}{n} \left(\sum_{i=1}^n (y_i - \bar{y})^2 \right)} = \frac{\frac{1}{n^2} \left(\sum_{i=1}^n a_i b_i \right)^2}{\frac{1}{n^2} \left(\sum_{i=1}^n a_i^2 \sum_{i=1}^n b_i^2 \right)}$$

$$\Rightarrow r^2_{xy} = \frac{\sum_{i=1}^n (a_i b_i)^2}{\sum a_i^2 \sum b_i^2} \quad \text{--- (1)}$$

CAUCHY SWARTZ INEQUALITY

$$\left(\sum a_i b_i \right)^2 \leq \left(\sum a_i^2 \right) \left(\sum b_i^2 \right)$$

sign of equality holding when $\frac{a_1}{b_1} = \frac{a_2}{b_2} = \dots = \frac{a_n}{b_n}$

Using cauchy swartz inequality in eqn (1), we have $r^2(x, y) \leq 1 \Rightarrow r \in [-1, 1]$

Properties of correlation coeff:-

1) $r \in [-1, 1]$ 2) Independent of change of origin & scale

Properties of Regression coeff:-

1) Correlation coeff is Geometric Mean b/w the regression coeff b_{yx} & b_{xy} $b_{yx} b_{xy} = r \frac{\sigma_y}{\sigma_x} * r \frac{\sigma_x}{\sigma_y} = r^2$

$$r = \pm \sqrt{b_{yx} b_{xy}} \quad (\text{sign of } r \text{ is same as regressed coeff})$$

$$2) r^2 \leq 1 \quad \downarrow b_{yx} \leq \frac{1}{b_{xy}}$$

4) regression coefficient independent of change of Origin but not scale
Origin ✓ Scale ✗

$$3) \left| \frac{b_{xy} + b_{yz}}{2} \right| > (r) \quad \text{(not less than)}$$

Angle b/w

Regression lines

$$\theta = \tan^{-1} \left\{ \left(\frac{1-r^2}{r^2} \right) \left(\frac{\sigma_y}{\sigma_x} \right) \right\}$$

Partial Correlation

$$r_{12.3} = \frac{\text{Cov}(X_{1.3}, X_{2.3})}{\sqrt{\text{Var}(X_{1.3}) \text{Var}(X_{2.3})}}$$

$$r_{12.3} = \frac{r_{12} - r_{13} * r_{23}}{\sqrt{(1-r_{13}^2)(1-r_{23}^2)}}$$

$$r_{13.2} = \frac{r_{13} - r_{12} * r_{23}}{\sqrt{(1-r_{12}^2)(1-r_{23}^2)}}$$

$$r_{23.1} = \frac{r_{23} - r_{21} * r_{31}}{\sqrt{(1-r_{21}^2)(1-r_{31}^2)}}$$

Correlation ratio (η)

measure of curvilinear reln b/w 2 variables (concentration of pt about the curve of best fit)

Note: if regression is linear $\eta = r$, otherwise $|\eta| > |r|$

$$\eta^2 = \frac{\sum_{i=1}^m \left[\frac{T_i^2}{n_i} \right] - \frac{T^2}{N}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij} (y_{ij}^2) - \frac{T^2}{N}}$$

$$\eta_{xy} = 1 - \left(\frac{T_j^2}{n_j} \right) - 1$$

(Same as above)

$$1) r=0, \tan \theta = \infty, \theta = \pi/2$$

$$2) r = \pm 1, \tan \theta = 0 \Rightarrow \theta = 0 \text{ or } \pi$$

Used to decide whether to include or not an additional independent variable in regression analysis

Multiple Correlation

$$R_{1.23}^2 = \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{(1-r_{23}^2)}$$

$$R_{3.12}^2 = \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{(1-r_{12}^2)}$$

$$R_{1.23}^2 = \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{1-r_{13}^2}$$

Properties

1) measures b/w var & grp
2) When $\text{MCC of expected value is calculated by MCC method} > \text{MCC linear combination of variables}$

$$3) 0 \leq R_{1.23} \leq 1$$

4) if $R_{1.23} = 0$, Total & Partial Correlation involving x_1 is zero

5) if $R_{1.23} = 1$ Perfect correlation

$$6) R_{1.23} > r_{12}, r_{13}, r_{23}$$

Prop. of Correlation Ratio

1) Independent of change in origin & scale

i.e. $\eta_{yz}^2 = \eta_{xy}^2 + \eta_{uv}^2$

2) $\eta \in [0, 1)$

$N = \sum_{i=1}^m \sum_{j=1}^n f_{ij}$ {Referend for}

$T_j = \sum_{i=1}^m f_{ij} y_{ij}$ η

$T = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij} = \sum_{i=1}^m T_i$
(Last page (corr. Ratio))

Analysis of Variance (ANOVA)

One Way Classification

RSS = Raw Sum of squares

$$RSS = \sum_{i=1}^k \sum_{j=1}^n x_{ij}^2$$

N = Total no. of observations

$G = \text{Grand Total} = \sum_{i=1}^k \sum_{j=1}^n x_{ij}$

TSS = Total sum of squares

$CF = \text{Correction Factor} = \frac{G^2}{N}$

$TSS = RSS - CF$

SST = Sum of squares Treatment

$= \sum_{i=1}^k \left(\frac{T_i^2}{n_i} \right) - CF$

SSE = Error sum of squares (Some of squares due to error)

$SSE = TSS - SST$

Source of Variation	df degrees of freedom	Sum of Squares	mean sum of squares	Variance Ratio
Treatment	$k-1$	SST	$MST = \frac{SST}{k-1}$	$F = \frac{MST}{MSE}$
Error	$n-k$	SSE	$MSE = \frac{SSE}{n-k}$	$\sim F(k-1, n-k)$
TSS	$n-1$	Estimation	Stats	Estimation of parameter

Anova II way

Rows $\rightarrow k$ Column $\rightarrow h$

Formula: $TSS = RSS - CF$
 $CF = \frac{G^2}{N}$

$SST = \frac{1}{h} \sum_{i=1}^k T_i^2 - CF$ (h columns)

$SSV = \frac{1}{k} \sum_{j=1}^h T_j^2 - CF$

$G = \sum_{i=1}^k \sum_{j=1}^h y_{ij}$ $T_i = \sum_{j=1}^h y_{ij}$

$T_j = \sum_{i=1}^k y_{ij}$

$SSE = TSS - SST - SSV$

Estimation - procedure of estimating a population parameter by using sample information or observation.

Parameter - a statistical constant derived from the population values

Source of Variance	d.f	Sum of Squares	Mean sum of squares	Variance Ratio
Treatment	$k-1$	SST	$MST = \frac{SST}{k-1}$	$F = \frac{MST}{MSF}$
Varieties	$h-1$	SSV	$MSV = \frac{SSV}{h-1}$	$\sim F_{(k-1)/(h-1)(k-1)}$
Error	$(k-1)(h-1)$	SSE	$MSE = \frac{SSE}{(k-1)(h-1)}$	$F = \frac{MSV}{MSE}$
Total	$hk-1$			$F_{(h-1)/(h-1)(k-1)}$

Anova II way Table

Properties of Estimators

It's good estimator if it is:

- (I) Unbiased (III) Efficient
- (II) Consistent (IV) Sufficient

1) Consistency - An estimator $T_n = T(x_1, x_2, \dots, x_n)$ based on a random sample of size n is const. if T_n converges to $r(\theta)$ i.e. if $T_n \xrightarrow{p} r(\theta)$ as $n \rightarrow \infty$

OR

In other words, it is consistent estimator of $r(\theta)$ if for every $\epsilon > 0$, $\eta > 0$, there exists a the integer $n \geq m(\epsilon, \eta)$ such that $P\{|T_n - r(\theta)| < \epsilon\} \rightarrow 1$ as $n \rightarrow \infty \Rightarrow P\{|T_n - r(\theta)| < \epsilon\} > 1 - \eta \forall n \geq n$ where n is some very large value
 $r(\theta) \rightarrow$ parameter
 $T_n \rightarrow$ Estimator

1) Unbiased - A statistic $\hat{\theta}$ is said to be unbiased iff the mean of the sample distribution of the estimator $= \theta$ i.e.

$$\text{Mean of } \hat{\theta} = E(\hat{\theta}) = \theta$$

⊗ If the estimator is not unbiased, then difference $E(\hat{\theta}) - \theta$ is called bias of the estimator

Positive Bias	Negative Bias
$E(\hat{\theta}) > \theta$	$E(\hat{\theta}) < \theta$

1) MVUE (minimum variance unbiased estimator)
 OR (most efficient unbiased estimator)

If $\hat{\theta}_1$ & $\hat{\theta}_2$ are 2 unbiased estimator of θ , if σ_1^2 & σ_2^2 are variances of their sampling distribution & $\sigma_1^2 < \sigma_2^2$, $\hat{\theta}_1$ is said to be more unbiased of θ
 Note - if $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ are unbiased estimators of parameter θ & $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$ are variances of their sampling distribution then $\hat{\theta}_1$ is the most efficient estimator

$e = \frac{\text{variance } \theta}{\text{variance } \hat{\theta}}$ where $i=2,3, \dots$ Estimator $\hat{\theta}_i$ of θ if $\sigma_1^2 < \sigma_2^2 < \dots < \sigma_k^2$

$$e = \frac{\text{Var } \theta}{\text{Var } \hat{\theta}_1} \quad e = \frac{\text{Var } \theta}{\text{Var } \hat{\theta}_i}$$

Sufficiency: An estimator is said to be sufficient for a parameter, if it contains all the information in the sample regarding the parameter.

If $T = t(x_1, x_2, \dots, x_n)$ is an estimator of a parameter θ based on sample x_1, x_2, \dots, x_n of size n from the population with density $f(x, \theta)$ such that the conditional distribution of x_1, x_2, \dots, x_n given T , is independent of θ , then t is sufficient estimator of θ .

Method of Estimation

$$L = f(x_1, \theta) f(x_2, \theta) \dots f(x_n, \theta) \\ = \prod_{i=1}^n f(x_i, \theta)$$

- ① Maximum likelihood method
- ② Method of moments
- ③ Method of minimum chi-square
- ④ Method of min. variance
- ⑤ Method of least squares
- ⑥ Method of inverse probability