

Estimation

①

Estimation: An important problem of statistical inference is the estimation of population parameters from corresponding sample statistics.

Estimation is required to use the statistic obtained from the sample as estimate of the unknown parameter of the population from which sample is drawn.

Estimate: An estimate is a statement made to find an unknown population parameter.

Estimator: The procedure or rule to determine an unknown population parameter.

Ex Sample mean is an estimator of population mean.

Because, the sample mean is a method of determining population mean and it is a statistic of sample.

A parameter can have one or two or many estimators.

Types of estimation:

Basically, there are two types of estimate to determine the statistic of the population, namely.

(a) point estimation (b) Interval estimation.

(a) point estimation:

A point estimation of a parameter is a single numerical value, which is computed from given sample and serves as an approximation of the unknown exact values of the parameter.

Ex The sample mean \bar{x} is a point estimator of population mean.

→ point estimator: A point estimator is a statistic for estimating the population parameter θ and will be denoted by $\hat{\theta}$.

ex 'sample mean \bar{x} ' is an point estimator of μ .

Properties of ^{Good} Estimator:

i) Consistency
An estimator is not expected to estimate the population parameter without errors. But it should give a value which is very close to the true value of unknown parameter. This is called consistency of estimator.

ii) Consistency

Thus, for smaller values of ϵ & η \exists $m \gg n$ such

$$\text{that } P[|\hat{\theta}_n - \theta| < \epsilon] > 1 - \eta$$

$$\text{ie } \hat{\theta}_n \rightarrow \theta \text{ as } n \rightarrow \infty.$$

iii) Unbiased estimator:

A statistic or a point estimator $\hat{\theta}$ is said to be unbiased estimator of the parameter θ if

$$E(\hat{\theta}) = \theta.$$

Otherwise it is biased.

Ex. Sample mean \bar{x} is an unbiased estimator of the population mean. (2)

Q1 Show that $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ is an unbiased estimator of the parameter σ^2 .

Sol Let x_1, x_2, \dots, x_n be the elements of sample of size n and μ and σ^2 are the population mean and variance. So, $E(x_i) = \mu$ and $\text{Var}(x_i) = E[(x_i - \mu)^2]$

Let us write. $= \sigma^2; i=1, 2, \dots, n.$

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n [(x_i - \mu) - (\bar{x} - \mu)]^2 \\ &= \sum_{i=1}^n (x_i - \mu)^2 - 2(\bar{x} - \mu) \sum_{i=1}^n (x_i - \mu) + n(\bar{x} - \mu)^2 \\ &= \sum_{i=1}^n (x_i - \mu)^2 - n(\bar{x} - \mu)^2. \quad \text{--- (1)} \end{aligned}$$

Now consider $E(s^2) = E\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right]$

$$\Rightarrow E(s^2) = \frac{1}{n-1} \left[\sum_{i=1}^n E(x_i - \mu)^2 - n E(\bar{x} - \mu)^2 \right]$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^n \sigma_{x_i}^2 - n \sigma_{\bar{x}}^2 \right] \quad \text{Using (1)}$$

However $\sigma_{x_i}^2 = \sigma^2, i=1, 2, \dots, n$ and $\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$.

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

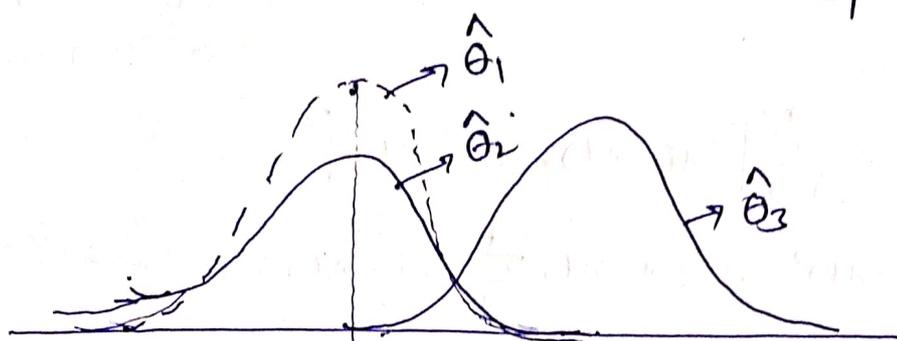
$\therefore E(s^2) = \sigma^2$ hence $\frac{1}{n-1} \sum (x_i - \bar{x})^2$ is an unbiased estimator for σ^2 .

Note $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ is a biased estimator.

ii) Most efficient estimator:

Consider all possible unbiased estimators of parameter θ . The one with the smallest variance is called the most efficient estimator.

Let $\hat{\theta}_1$ and $\hat{\theta}_2$ are two unbiased estimators of the same population parameter θ . Then $\hat{\theta}_1$ is more efficient of θ than $\hat{\theta}_2$ if $\sigma_{\hat{\theta}_1}^2 < \sigma_{\hat{\theta}_2}^2$.



Sampling distribution of estimators.

In the above figure, we illustrate the sampling distribution of three different estimators $\hat{\theta}_1$, $\hat{\theta}_2$ and $\hat{\theta}_3$ all estimating θ . It is clear that only $\hat{\theta}_1$ and $\hat{\theta}_2$ are unbiased. Since their distributions are centered θ i.e. $E(\hat{\theta}_1) = E(\hat{\theta}_2) = \theta$. But, variance of estimator $\hat{\theta}_1$ is smaller than variance of $\hat{\theta}_2$ i.e. $\sigma_{\hat{\theta}_1}^2 < \sigma_{\hat{\theta}_2}^2$. Therefore $\hat{\theta}_1$ is more efficient among the three considered estimators.

iv) Sufficient estimator

An estimator is said to be sufficient for a parameter, if it contains all the information in the sample regarding the parameter. Thus a good estimator is one which is close to true value of the parameter as possible and should have i) consistency ii) unbiased iii) Efficient and iv) sufficient in all directions.

Interval estimation

Even the most efficient unbiased estimator cannot estimate the population parameter exactly. It is true that accuracy increases with large samples, but there is still no reason point estimate from a given sample to be exactly equal to the population parameter, it is supposed to estimate.

Therefore, in many situations it is preferable to determine an interval within which we would expect exact value of parameter.

An estimator of a population parameter given by two magnitudes within which a parameter can lie is called interval estimator of the parameter.

Let θ denote a population parameter. An interval estimator of θ is an interval of the form $\hat{\theta}_L < \theta < \hat{\theta}_U$, where $\hat{\theta}_L$ & $\hat{\theta}_U$ depend on the value of θ for a particular sample and also on the sampling distribution of $\hat{\theta}$.

Let $\hat{\theta}$ be the random variable of estimator. Then samples will generally yield different values of $\hat{\theta}$.

From the sampling distribution of $\hat{\theta}$, we shall able to find $\hat{\theta}_L$ & $\hat{\theta}_U$ such that

$$P(\hat{\theta}_L < \theta < \hat{\theta}_U) = \text{percentage of confidence.}$$

If $1-\alpha$ is a confidence about interval of $\hat{\theta}$ then we can find $\hat{\theta}_L$ and $\hat{\theta}_U$ such that

$$P(\hat{\theta}_L < \theta < \hat{\theta}_U) = 1-\alpha.$$

where $0 < \alpha < 1$ (i.e. critical value).

Here the interval $\hat{\theta}_L < \theta < \hat{\theta}_U$ is called a $(1-\alpha)100\%$ confidence interval.

The fraction $1-\alpha$ is confidence coefficient or degree of confidence. Also the end points, θ_L and θ_U are called lower and upper confidence limits.

Thus, when $\alpha = 0.05$, we have 95% confidence interval and $\alpha = 0.01$ we obtain a confidence interval with 99% confidence.

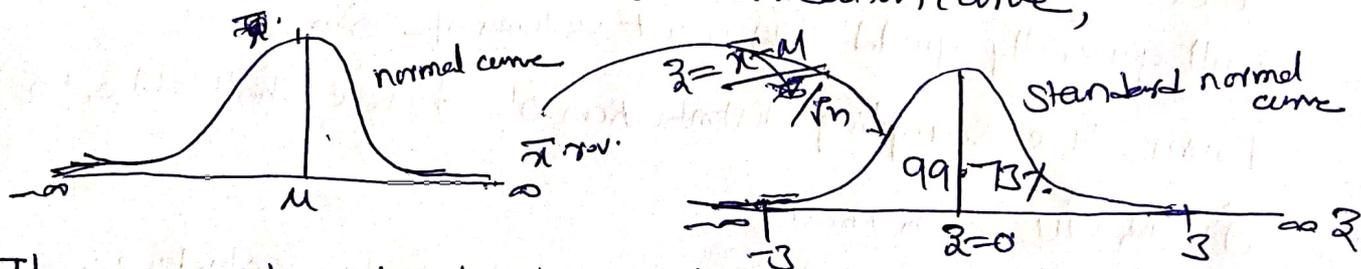
Confidence limits for population mean (μ)

Let \bar{x} be the sample mean of sampling distribution of means of random sample of size n ($n \geq 30$) drawn from a population having mean μ and S.D. σ is approximated by normal distribution with mean $\mu_{\bar{x}} = \mu$ and S.D. $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$.

ie By central limit theorem $\bar{x} \sim N(\mu_{\bar{x}}, \sigma_{\bar{x}})$ for large samples.

Hence, $Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$ has standard normal distribution with mean 0 and S.D. 1. and it is used for μ .

From the standard normal distribution curve,



The area under standard normal curve between $Z = -3$ and $Z = 3$ is 99.73%.

ie $P(-3 \leq Z \leq 3) = 0.9973$ where $Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$

Then the inequality $-3 < \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < 3$ — (1) shows that 99.73% of cases are favourable to μ and 0.27% cases are not favourable.

From ①, we have

$$-3 < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < 3$$

$$\Rightarrow -3 \frac{\sigma}{\sqrt{n}} < \bar{x} - \mu < 3 \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow \bar{x} - 3 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 3 \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow \bar{x} - 3 \text{ S.E.} < \mu < \bar{x} + 3 \text{ S.E.} \quad \square$$

Since S.E. of \bar{x} is $\frac{\sigma}{\sqrt{n}}$.

$\therefore [\bar{x} - 3 \text{ S.E.}, \bar{x} + 3 \text{ S.E.}]$ is called 99.73%.

confidence interval for the population mean μ .

Similarly $[\bar{x} - 2 \text{ S.E.}, \bar{x} + 2 \text{ S.E.}]$ is 95.45% confidence interval for μ .

$[\bar{x} - \text{S.E.}, \bar{x} + \text{S.E.}]$ is 68.26% confidence interval.

In general, confidence limits for population mean μ are

$\bar{x} \pm z_{\alpha}^* (\text{S.E. of } \bar{x})$, where z_{α}^* value of z .

ii) Confidence limits for population mean with:

i) 99% confidence are $\bar{x} \pm 2.58 \frac{\sigma}{\sqrt{n}}$.

ii) 95% confidence are $\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}}$.

iii) 90% confidence are $\bar{x} \pm 1.64 \frac{\sigma}{\sqrt{n}}$.

Thus $(1-\alpha)$ 100% confidence interval of population mean is

$$\bar{x} \pm z_{\frac{\alpha}{2}}^* \frac{\sigma}{\sqrt{n}} \quad \text{where } P(-z_{\frac{\alpha}{2}}^* \leq z \leq z_{\frac{\alpha}{2}}^*) = 1-\alpha.$$

where $z_{\frac{\alpha}{2}}^*$ is a value of z leaving an area of $\frac{\alpha}{2}$ to right ~~and left~~ tail of standard normal distribution curve.

Confidence limits for population proportion:

By central limit theorem, the sampling distribution of the sample proportions p will be approximated by normal distribution with a statistic $Z = \frac{p - P}{\sqrt{\frac{PQ}{n}}}$ for a large sample.

Here Z is a S.D variable with mean 0 and S.D 1.

Then the confidence limits for population proportion with

i) 99% of confidence are $p \pm 2.58$ (S.E of p)

ii) 95% of confidence are $p \pm 1.96$ (S.E of p)

iii) 90% of confidence are $p \pm 1.64$ (S.E of p)

Here S.E of p is $\sqrt{\frac{PQ}{n}}$.

In general, the $(1-\alpha)100\%$ confidence interval for population proportion is $p \pm Z_{\alpha/2}$ (S.E of p).

Here $Z_{\alpha/2}$ is a value of Z leaving an area of $\alpha/2$ to right tail of standard normal curve.

Question

① A random sample of 200 measurements from a large population gave mean value of 50 and S.D of 9. Determine the 95% confidence interval for the mean μ of the population.

sol Given that a sample of 200 is drawn from normal population having mean 50 and S.D 9.

i.e. $n = 200$, $\bar{x} = 50$, $\sigma = 9$.

The 95% confidence interval for ~~pop~~ population mean is

$$\bar{x} \pm 1.96 (\text{S.E of } \bar{x}) \text{ i.e. } \bar{x} \left(50 \pm 1.96 \frac{9}{\sqrt{200}} \right)$$

②. * Maximum Error of estimate σ with $(1-\alpha)100\%$ confidence is

$$E_{\max} = |\sigma - E(\sigma)|$$

①. * Maximum Error in sampling distribution of mean is

$$E_{\max} = Z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right) \text{ for a normal population.}$$

i.e. one can assert a confidence interval for population mean with $(1-\alpha)100\%$ confidence that the error will not exceed $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$.

$$\text{Then } n = \left(\frac{Z_{\alpha/2} \sigma}{E_{\max}} \right)^2$$

②. Maximum error of sample proportion is

$$E_{\max} = Z_{\alpha/2} \sqrt{\frac{pq}{n}}, \quad Q = 1 - P$$

Questions

- ① To estimate the average amount of time visitors take to move from one building to another in an office complex, the mean of a random sample of size n is used. Given $\sigma = 1.40$ minutes, determine how large should be the sample size if it is ascertained with 99% confidence that the maximum error is at most 0.25.
- ② A random sample of size 100 has a standard deviation of 5. What can you say about the maximum error with 95% confidence.
- ③ Construct a 99% confidence interval for the true mean weight loss. If 40 persons on diet control after one month had mean weight loss 5 Kgs with S.D of 1.2 Kgs.