

Lecture - 12 : Interval Estimation

Section 1: Confidence Intervals

What we have learned till now is point estimation. In order to estimate the value of a population parameter, we just draw a random sample of a reasonable size and construct a statistic which can best estimate that population parameter, in some sense, either as unbiased or as the best maximum likelihood estimator. But if we are drawing a sample from a probability density function the probability that a value of the statistic will coincide with population parameter is zero. This motivates us to take the value that view that it might be better to be able to construct an interval around the statistic so that the population parameter value falls in that interval with certain probability.

Let us draw a sample from the population with pdf $f_X(x, \theta)$, and let x_1, \dots, x_n be a random sample of size n drawn from it, and let Θ be a statistic, $\Theta = \Theta(x_1, \dots, x_n)$ constructed to estimate, θ . We may for example want to know, if

$$\Theta - c < \theta < \Theta + c$$

with probability say γ ; i.e.

$$P(\Theta - c < \theta < \Theta + c) = \gamma$$

Then $(\Theta - c, \Theta + c)$ forms a confidence interval associated with Θ . This is more formally defined as follows.

Definition: If $\hat{\Theta}_1$ and $\hat{\Theta}_2$ are values of the random variables $\hat{\Theta}_1, \hat{\Theta}_2$, such that which are themselves statistics, such that,

$$P(\hat{\Theta}_1 < \theta < \hat{\Theta}_2) = 1 - \alpha$$

for some specified probability $1 - \alpha$. We shall refer to the interval $(\hat{\Theta}_1, \hat{\Theta}_2)$ as a $(1 - \alpha) 100\%$ confidence interval for the population parameter θ . The value $1 - \alpha$ is called degree of confidence & $\hat{\Theta}_1$ & $\hat{\Theta}_2$ the upper and lower confidence limits.

limits. So if $\alpha = 0.05$, we say the degree of confidence is 0.95 and we get a 95% confidence interval. It is also important to note at the very beginning that the confidence intervals for θ need not be unique as we will just see.

Section 2: Interval Estimation of means

In this section we will study first the interval estimation of the mean μ of a normal population whose variance σ^2 is known. Next we will see how to deal the case when σ^2 is unknown. It will be important to note how $\hat{\theta}_1$ and $\hat{\theta}_2$ will be very different in these two scenarios.

Case 1: σ^2 is known.

Let x_1, \dots, x_n be a random sample of size n drawn from a population $N(\mu, \sigma^2)$, where σ^2 is known but μ is not. In order to estimate the $(1-\alpha)100\%$ confidence interval of μ we proceed as follows.

$$\text{Note that } E(\bar{x}) = \mu \quad \text{and} \quad \text{Var}(\bar{x}) = \frac{\sigma^2}{n}$$

$$\text{where } \bar{x} = \frac{x_1 + \dots + x_n}{n}.$$

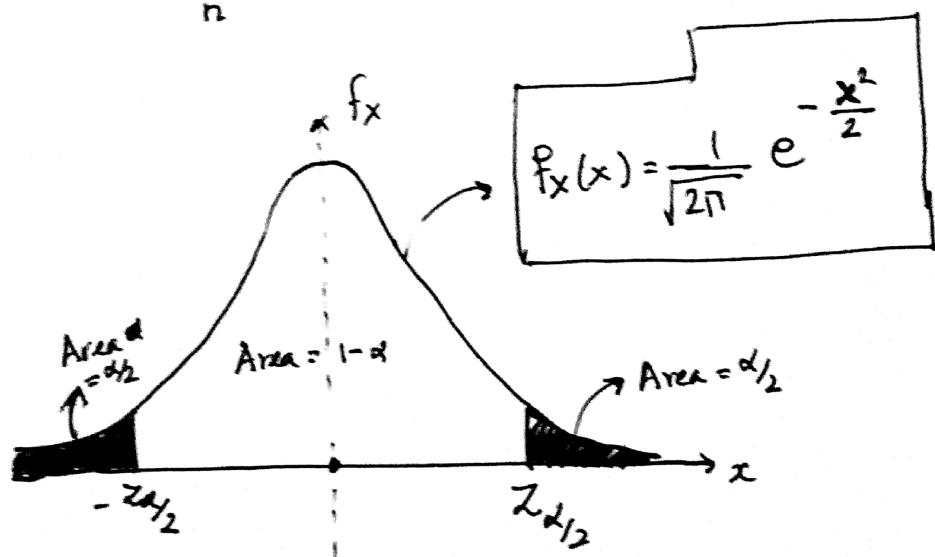


Fig 1: The point $z_{\alpha/2}$

Let us look at diagram above and it is seen how the point $Z_{\alpha/2}$ is defined, i.e.

$$f_{Z_{\alpha/2}} \left[\int_{Z_{\alpha/2}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \alpha/2 \right]$$

From diagram if $Z \sim N(0,1)$ then

$$P(-Z_{\alpha/2} < Z < Z_{\alpha/2}) = 1 - \alpha \rightarrow \textcircled{2}$$

Since we have drawn a normal population

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \quad [\text{Check Chapter Lecture 10}]$$

$$\therefore \text{Set } Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$$

Hence from $\textcircled{2}$ we have

$$P\left(-Z_{\alpha/2} < \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < Z_{\alpha/2}\right) = 1 - \alpha$$

$$\boxed{P\left(\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha}$$

$$\text{Here } \hat{\theta}_1 = \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad \& \quad \hat{\theta}_2 = \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

(3)

Now we will again consider studying a normal population but this time we shall consider that the population variance σ^2 is unknown. In this situation we have to use sample variance s^2 in order to try to develop a statistic using it. Let us set

$$T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

Let us assume that the sample is small, i.e. $n < 30$, then it better to proceed by noting that $T \sim t$ distribution with $n-1$ degrees of freedom (see Lecture 10). To find out the $(1-\alpha)100\%$ confidence interval we must have two points q_1, q_2 , such that

$$P(q_1 < T < q_2) = 1 - \alpha$$

$$\Rightarrow P(q_1 < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < q_2) = 1 - \alpha$$

$$\Rightarrow P\left(\frac{s}{\sqrt{n}}q_1 < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < \frac{s}{\sqrt{n}}q_2\right) = 1 - \alpha$$

$$\Rightarrow P\left(\bar{X} - q_2 \frac{s}{\sqrt{n}} < \mu < \bar{X} - q_1 \frac{s}{\sqrt{n}}\right) = 1 - \alpha$$

Thus in this case the confidence interval for the mean is given by

$$\left(\bar{X} - q_2 \frac{s}{\sqrt{n}}, \bar{X} - q_1 \frac{s}{\sqrt{n}} \right)$$

Hence the length of this interval = $(q_2 - q_1) \frac{s}{\sqrt{n}}$

$$\text{or } L = (q_2 - q_1) \frac{s}{\sqrt{n}}.$$

From a more practical point of view we want $(q_{12} - q_1)$ to be as small as possible. However note that L is in general a random variable. However for a given sample size of size say n , we need to actually solve the optimization problem

$$\min L = (q_{12} - q_1) \frac{S}{\sqrt{n}} \longrightarrow (E)$$

Subject to $\int_{q_1}^{q_{12}} f_T(q) dq = 1 - \alpha \longrightarrow (F).$

where f_T is the pdf of the t-distribution with $(n-1)$ degrees of freedom. Let us rewrite (F) using the dummy variable t , i.e

$$\int_{q_1}^{q_{12}} f_T(t) dt = 1 - \alpha. \longrightarrow (F')$$

Since α is fixed from a theoretical point of view one can always express q_{12} as a function of q_1 . Let us first differentiate (F') with respect to q_1 .

In order to do so we need to use the Leibnitz formula for differentiating under the integral sign. We shall write it in the following box.

Leibnitz rule : For differentiating under the integral sign.

Let $a(x)$ and $b(x)$ are differentiable function. Then consider the integral

$$\int_{a(x)}^{b(x)} \varphi(z, t) dt$$

$$\frac{d}{dx} \int_{a(z)}^{b(z)} \varphi(z, t) dt = \int_{a(z)}^{b(z)} \frac{\partial}{\partial z} \varphi(z, t) dt + \cancel{\int_{a(z)}^{b(z)} \varphi(z, b(z)) \frac{db}{dz}}$$

$$- \varphi(z, a(z)) \frac{da}{dz}$$

Note that if we go back to (P'), f_T depends on t and not on q_1 . Thus we will have

$$\frac{d}{dq_1} \int_{q_1}^{q_2} f_T(t) dt = 0$$

By Leibnitz rule we

$$\int_{q_1}^{q_2} \frac{\partial}{\partial q_1} f_T(t) dt + f_T(q_2) \frac{dq_2}{dq_1} - f_T(q_1) \frac{dq_1}{dq_1} = 0$$

\therefore

\Rightarrow As $\frac{\partial}{\partial q_1} f_T(t) = 0$ we have from the above expression

$$f_T(q_2) \frac{dq_2}{dq_1} - f_T(q_1) = 0 \longrightarrow (a)$$

Now when we replace q_2 as a function of q_1 in L then the problem become unconstrained and to minimize L we first differentiate with respect to q_1 ,

$$\left(\frac{dq_2}{dq_1} - 1 \right) \frac{S}{\sqrt{n}} = 0$$

$$\Rightarrow \left(\frac{f_T(q_1)}{f_T(q_2)} - 1 \right) \frac{S}{\sqrt{n}} = 0$$

$$\text{or } f_T(q_1) = f_T(q_2)$$

Since q_1 is the minimizer then $f_T(q_1) = f_T(q_2)$.

Of course one solution is $q_1 = q_2$, but that will give $1 - 1 = 0$ from (P'), which is a contradiction.

But as the graph of t distribution is symmetric we have $\alpha_1 = \alpha_2$ as the solution. So from the statistical tables we have to find α_1 such that

$$P(T > q_{\alpha}) = \frac{\alpha}{2}, \text{ where } T \sim t \text{ (n-1 degrees of freedom)}$$

In general one writes $q_{\alpha} = t_{\alpha/2, n-1}$

What happens if we have a large sample. Let us draw a sample from a normal distribution. Then using what is known as "central limit theorem", we can deduce that if

$$Z_n = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

where n is the sample size, then as $n \rightarrow \infty$, the distribution function of Z_n coincides with that of the standard normal distribution $N(0,1)$. Thus for very large $Z_n \sim N(0,1)$ "approximately". So I leave it to the reader to deduce that the confidence $(1-\alpha)100\%$ confidence interval of μ is given as

$$\left(\bar{x} - Z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \frac{s}{\sqrt{n}} \right)$$

Of course the next question one may ask is what about the variance. If we are studying a normal population with mean μ and variance σ^2 , what is the confidence interval for σ^2 , assuming of course σ^2 is not known. This is what we will study in the next section.

3. Interval Estimation for the Variance

As before let us draw a sample x_1, \dots, x_n from a normal distribution with mean μ and variance σ^2 . Here the variance is unknown. In order to proceed we will introduce what is known as a pivotal quantity.

Defn 3.1 (Pivotal quantity)

Let x_1, \dots, x_n be a random sample of size n drawn from a density $f(\cdot, \theta)$, where θ as before is the parameter of the distribution. Let $Q = q_f(x_1, \dots, x_n, \theta)$, i.e. Q is a function of the sample observations and θ . If Q is a probability density that does not depend on θ , then Q is called a pivotal quantity.

E.g.: In the example of estimating mean $Z_n = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ is a pivotal quantity.

So the key idea of estimation is to find $q_{1-\alpha}$ such that

$$P[q_{1-\alpha} < Q < q_{\alpha}] = 1 - \alpha.$$

In the particular case of estimating variance, the pivotal quantity is

$$Q = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2}$$

Now $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2$ with $(n-1)$ degrees of freedom.

$$P[q_{1-\alpha} < Q < q_{\alpha}] = 1 - \alpha$$

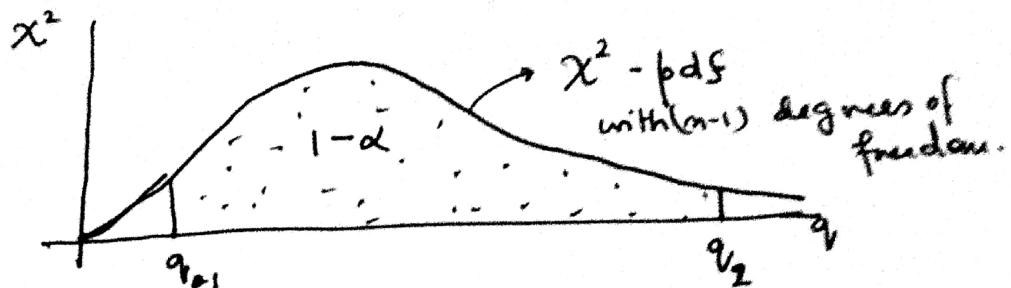
$$\Rightarrow P[q_{1-\alpha} < \frac{(n-1)S^2}{\sigma^2} < q_{\alpha}] = 1 - \alpha \quad (\text{S^2 is the observed value of the r.v. S^2})$$

$$P \left[\frac{(n-1)s^2}{q_2} < \sigma^2 < \frac{(n-1)s^2}{q_1} \right] = 1-\alpha$$

∴ The $(1-\alpha)100\%$ confidence intervals are then given as

$$\left(\frac{(n-1)s^2}{q_2}, \frac{(n-1)s^2}{q_1} \right)$$

Remember that χ^2 with $(n-1)$ degrees of freedom is not a symmetric distribution



The key idea is to find q_1 and q_2 such that

$$P[Q < q_1] = \alpha/2 \quad \text{and} \quad P[Q > q_2] = \alpha/2.$$

This is called equal tails estimation. In practice one should figure this out from the Statistical tables related to χ^2 distribution but as before one may also obtain this by numerically solving the optimization problem.

$$\min \min L = \frac{(n-1)s^2}{q_1} \left[\frac{1}{q_1} - \frac{1}{q_2} \right]$$

Subject to $\int_{q_1}^{q_2} f_Q(q) dq = 1 - \alpha.$

4. Estimation of the Confidence Interval for the difference of the means of two populations

In this case we will be drawing samples from two normal populations.

Sample 1 will be drawn from $N(\mu_1, \sigma_1^2)$

Sample 2 will be drawn from $N(\mu_2, \sigma_2^2)$

We will assume that σ_1^2 and σ_2^2 are known. We want to estimate the difference of the means μ_1 and μ_2 , i.e. $\mu_1 - \mu_2$ to know how different these two normal populations are. Size of sample 1 is n_1 .

The pivotal quantity in this case is

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

where \bar{x}_1 is the mean of the first sample and \bar{x}_2 is the mean of the second sample. One can show that

$$Z \sim N(0, 1)$$

Observe that

$$E(\bar{x}_1 - \bar{x}_2) = E(\bar{x}_1) - E(\bar{x}_2) = \mu_1 - \mu_2$$

$$\text{Var}(\bar{x}_1 - \bar{x}_2) = \text{Var}(\bar{x}_1) + \text{Var}(\bar{x}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

As, \bar{x}_1 & \bar{x}_2 are independent

It can be proved using any approach of finding the sampling distribution of $\bar{x}_1 - \bar{x}_2$, that

$$\bar{x}_1 - \bar{x}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right).$$

Thus $Z \sim N(0, 1)$.

Hence if we are looking for a $(1-\alpha)100\%$ interval confidence interval then

$$(\bar{x}_1 - \bar{x}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

If the samples are large and $\sigma_1^2 + \sigma_2^2$ are not known then replace them by s_1^2 and s_2^2 in the expression of Z and proceed using the central limit theorem.

In that case we have

$$(\bar{x}_1 - \bar{x}_2) - Z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + Z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

However if the sample size is small i.e $(n_1 < 30 \text{ or } n_2 < 30)$ or both, then the confidence interval can be used if both population is assumed to have the same variance $\sigma^2 = \sigma_1^2 = \sigma_2^2$ even though its exact value is not known. Then the key is to consider the pooled sample variance

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

This is nothing but an weighted average of s_1^2 & s_2^2 ; i.e

$$S_p^2 = \frac{\sum_{i=1}^{n_1} (x_i^1 - \bar{x}_1)^2 + \sum_{i=1}^{n_2} (x_i^2 - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

You will show in the assignments that $E(S_p^2) = \sigma^2$

Construct a new random variable

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

Now $\frac{(n_1-1)}{\sigma^2} S_p^2 \sim \chi^2$ with n_1-1 degrees of freedom

and $\frac{(n_2-1)}{\sigma^2} S_p^2 \sim \chi^2$ with n_2-1 degrees of freedom.

Thus $Y \sim \chi^2$ with $n_1 + n_2 - 1$ degrees of freedom. Now construct the pivotal variable

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Of course it can be shown that $\rightarrow T$ follows the t distribution with $n_1 + n_2 - 2$ degrees of freedom. Once this is known we can go ahead and write down the $(1-\alpha)100\%$ confidence intervals. But how did we get T .

Note that let us set

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

of course $Z \sim N(0,1)$. It can also be shown that Z and Y are independent. (Note that we have seen one can show that \bar{X} and S^2 are independent).

$$\therefore T = \frac{Z}{\sqrt{\frac{Y}{n_1 + n_2 - 2}}}$$

Hence $T \sim t$ with $(n_1 + n_2 - 2)$ degrees of freedom. This is known from our study of Sampling distributions. The fact that $\rightarrow T \sim t$ with $(n_1 + n_2 - 2)$ degrees of freedom precisely needs the fact that $\rightarrow Z$ and Y are independent

$$\therefore T = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Hence that $\rightarrow (1-\alpha)100\%$ confidence interval is given as

$$(\bar{x}_1 - \bar{x}_2) - t_{\alpha/2, n-1} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2, n-1} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$