**Biostatistics and Design of Experiments**
**Prof. Mukesh Doble**
**Department of Biotechnology**
**Indian institute of Technology, Madras**

**Lecture - 07**
**T distribution/confidence interval**

Welcome to the next class on Biostatistics and Design of Experiments. We will talk more about t distribution, Z distribution and also confidence interval. That means, I had mentioned before if I am collecting a sample and getting a mean and I want to find out what is the confidence interval on the population mean, I can use some equations to get that sort of information actually. So, that is what we are going to talk about in this class.
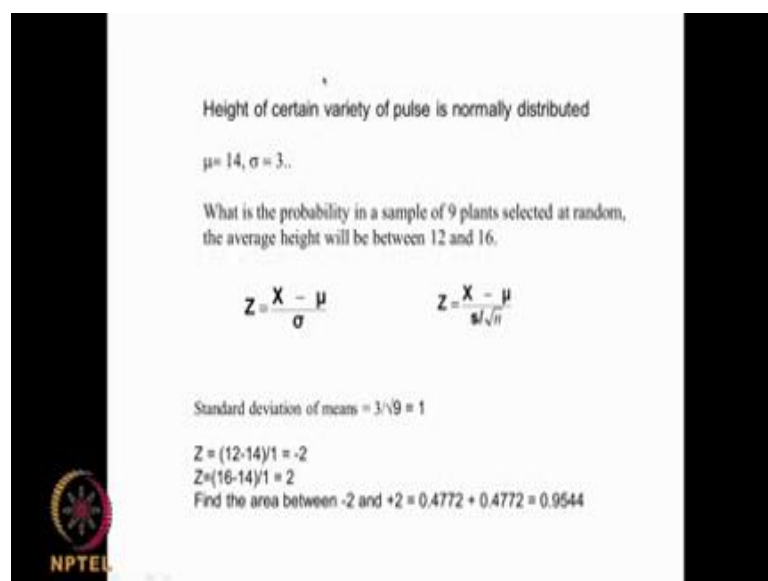
(Refer Slide Time: 00:43)



Let me again recall, suppose I have a population and I have sample, so that means I collect a few pieces from the population and I get a mean. Now when I take another setup of pieces, I may get a different mean. If I keep on repeating I will get large number of means. The means are representation of the population mean. That means, the means of the sample are representation of the population mean. But the standard deviation that is going to happen because the means will be distributed in a normal way, the means of the samples are going to be distributed in a normal way and there is going to be a standard deviation associated with that. So, from that standard deviation we should be

able to get the overall estimate of the standard deviation.

So generally, the sample mean is a usual estimator of the population mean and if I take different samples I will get a large number of means which will have its own mean. That means, mean of all the means and it also will have a distribution, a normal distribution with its own variance. So, the standard error of the mean, we call it standard error, is the standard deviation of these sample means estimate of the population. So how do you calculation standard error? We get the standard deviation of the sample, divided by square root of n. This is called the standard error or standard error of the mean.

So, this s is the sample standard deviation, n is the sample size. So, if I am taking a set of samples, I will get a mean. Now I put it back then again take another set of samples, I will get another mean, if I keep on repeating I will get large number of means. Now these means will be normally distributed, it will have its own mean that is the mean of all these means and it will have its own standard deviation. Because the whole data will be a normally distributed. So, I can calculate the standard error of this mean or it is also called standard error based on the sample standard deviation. So, s is the sample standard deviation, n is the number of samples you take and you do a square root and divide it actually. This is very important standard error of the mean, is very important because we use this and calculating confidence internal on whatever data we collect. We will see that time as we go along.

(Refer Slide Time: 03:24)



Height of certain variety of pulse is normally distributed

$\mu = 14, \sigma = 3..$

What is the probability in a sample of 9 plants selected at random, the average height will be between 12 and 16.

$$Z = \frac{X - \mu}{\sigma} \qquad Z = \frac{X - \mu}{s/\sqrt{n}}$$

Standard deviation of means $= 3/\sqrt{9} = 1$

$Z = (12\text{-}14)/1 = -2$
$Z = (16\text{-}14)/1 = 2$
Find the area between -2 and +2 $= 0.4772 + 0.4772 = 0.9544$

Now let us look at a problem. Height of certain variety of pulse is normally distributed. This is given a μ 14 and a σ is 3. What is the probability in a sample of nine plants? Suppose I take nine plants in a random, that the heights will lie between 12 and 16. What will be the statistics? So, you remember this,

$$Z = \dfrac{X - μ}{σ}$$

, now this σ we do not know but we know only the standard deviation of the samples, right? So instead of σ we need to substitute by

$s/\sqrt{n}$ If you remember in the previous equation right, $s/\sqrt{n}$ , because we have the estimate of the standard deviation from the sample, we do not know the population. So, we use this here and σ is given as 3, √9. We have got 9 plants, so 3 by 9 square root which will be one, so Z = (12-14)/1

because we want find out from X of 12 to x of 16, so Z = 16, we will get - 2 and + 2. So, I want to find out the area between - 2 and + 2, if you remember long time back we did that from the Z table.

(Refer Slide Time: 04:54)



This is the Z table. 2 we will get it as this 0.0228. But the catch here is it estimates this region, so if I want know this region I will subtract from 0.5 minus 0.0228. That is the

**s/**$\sqrt{9}$

(Refer Slide Time: 06:04)



Now again going back to the t distribution, the confidence interval of the mean I take a set of samples and I get the mean that is $\overline{X}$ and I get a standard deviation of that sample that is **s/**$\sqrt{n}$, if you understand this, this is the standard error multiplied by t, t is the t value and it will vary depending upon whether it is 95 % confidence interval or 99 % confidence interval. So generally, if n is large 95 % confidence interval t value will be 1.96, if n is large 99 % confidence will become 2.58.

If you remember long time back I did mention how this 1.96 and 2.58 can logically come. So, for very large N if we take a normally distribution, + or - 2 $\sigma$ occupies approximately 95 % of the area and + or - 3 $\sigma$ occupies in 99 % area. Instead of 2 and 3 here, we are using 1.96 and 2.58 because we are using t distribution. As I said it is impossible to get population so we can only get a sample, and from the sample we make

an estimate of the population values that is why these numbers instead of 2 and 3 it became 1.96 and 2.58. If I have a small sample and I get a mean and I get a standard deviation, then I can get a confidence interval for the population mean using these data by using this formula. This s by $s/\sqrt{n}$ is called the standard error, t is the t value from the t table. I will show you the t table and for large N the t value for 95 % confidence will become 1.96 and for 99 % confidence it will become 2.58. And I also mentioned, what is the logic of this 1.96 and 2.58. As I said, in a normal distribution for a 95 % of the area is covered inside + or - 2 $\sigma$ and 99 $\sigma$ of the area is covered by + or - 3 $\sigma$, so instead of using 2 and 3 because it is a t distribution and sample is always much smaller than entire population, 2 has become 1.96 and 3 has become 2.58. And that is how you would use that formula.

So, you have

$$\mu = \overline{X} \pm \frac{t_{df}\,S}{\sqrt{n}}$$

If n is small there is a table called t table for corresponding degrees of freedom. Where degrees of freedom is equal to n -1, we can get the t value and you multiplied with the standard error and from the $\overline{X}$ that is the mean of the sample, you get the confidence interval for the population mean.

(Refer Slide Time: 09:18)



EXCEL

CONFIDENCE(alpha,standard_dev,size)

**Alpha**   is the significance level used to compute the confidence level. The confidence level equals 100*(1 - alpha)%, or in other words, an alpha of 0.05 indicates a 95 percent confidence level.

**Standard_dev**   is the population standard deviation for the data range and is assumed to be known.

**Size**   is the sample size.

NPTEL

In excel there is function called Confidence. It is alpha, alpha is the significance level. So, it can be 0.05 or 0.01, 0.05 indicates 95 % confidence and 0.01 indicates 99 %. So, you have this standard deviation, is the population standard deviation for the data and size is the n actually. Basically, it gives you this side of it if you know the $\overline{X}$ you add and then you subtract to get the confidence interval. So, we can use the excel function called Confidence also here.

(Refer Slide Time: 09:58)



There is another terminology that is called the coefficient of variation CV, generally we call it, it is a relative standard deviation. It is a standardized measure of dispersion of a probability distribution. So, CV is given by

$$c_v = \frac{\sigma}{\mu} \times 100$$

that is the standard deviation divided by mean. It tells you what is a relatively the spread in the data, because you have the standard deviation in the numerator and you have the mean in the denominator.

If CV is very large, that means my spread is very large, if CV is very small that means my spread is very small. So, CV is a quick measure of telling how bad or how good the spread of the data is? and If I have 2 sets of the data and I know the CV's I can tell which

data set is most spread out or which data set is more tighter and less spread out. So, that way CV is also a very useful property which we can calculate and make use of in comparing different data sets.

(Refer Slide Time: 10:58)



Heart beat of 113 students. Calculate the 99% confidence interval for the mean heart beat? Calculate cv

| heart beat | number |
| --- | --- |
| 57 | 2 |
| 62 | 4 |
| 67 | 11 |
| 72 | 23 |
| 77 | 20 |
| 82 | 21 |
| 87 | 13 |
| 92 | 7 |
| 97 | 5 |
| 102 | 1 |
| 107 | 1 |
| 112 | 1 |
| 117 | 4 |

Let us look at a problem.

Heart beat of 113 students is given there. That is 2 students have 57 beats per minute, 4 students seems to have 62, 11 students seem to have 67 and so on, if you add up all these comes to 113. So, it goes right up to 117, that means 4 students have 100 heart beat of 117. Now we can get a mean for the entire data set, but that is not the population mean. We can get a confidence region for the population mean based on the sample mean. And you are supposed to calculate, what is that confidence interval for the population mean for 99 % data? And also calculate CV. It is quite simple.

What we do is we can calculate $\overline{X}$ from this data then we can calculate standard deviation that is s from this data and I know n and then for 113 students the degree of freedom is 112. So, I go to the table generally the data is large, shown for 99 %, t value will be two 2.58 so I multiply 2.58, multiplied by the standard deviation of this data and divided by square root of n that will be the variation in my $\overline{X}$, $\overline{X}$ is my average. So, we use this equation, understood.

Heart beat of 113 students. Calculate the 99% confidence interval for the mean heart beat? Calculate cv

| heart beat | number |
|------------|--------|
| 57 | 2 |
| 62 | 4 |
| 67 | 11 |
| 72 | 23 |
| 77 | 20 |
| 82 | 21 |
| 87 | 13 |
| 92 | 7 |
| 97 | 5 |
| 102 | 1 |
| 107 | 1 |
| 112 | 1 |
| 117 | 4 |

$$\mu = \bar{X} \pm \frac{t_{df} \, S}{\sqrt{n}}$$

So, we calculate the $\bar{X}$ that is mean of the entire data, and then we calculate the standard deviation of this data, entire data and n is 113. Now t is my t value of which I can determine for 112 degrees of freedom because I have 113 data, so generally the t value here will be 2.58 because the data set is large. Let us do the problem using excel, it is not very difficult.

| heart beat | number | HB*n | | n*(HB-Hb_avg)^2 | |
|------------|--------|------|--|-----------------|--|
| 57 | 2 | 114 | | 1079.274023 | |
| 62 | 4 | 248 | | 1329.344506 | |
| 67 | 11 | 737 | | 1925.387658 | |
| 72 | 23 | 1656 | | 1557.890203 | |
| 77 | 20 | 1540 | | 208.6694338 | |
| 82 | 21 | 1722 | | 65.7843214 | |
| 87 | 13 | 1131 | | 595.8121231 | |
| 92 | 7 | 644 | | 969.7157178 | |
| 97 | 5 | 485 | | 1406.149659 | |
| 102 | 1 | 102 | | 473.9290469 | |
| 107 | 1 | 107 | | 716.628162 | |
| 112 | 1 | 112 | | 1009.327277 | |
| 117 | 4 | 468 | Avg | 5408.105568 | |
| | 113 | 9066 | 80.23009 | 16746.0177 | 149.51 |
| | | | | Sample var (=149.5/113) | 1.32 |
| | | | | sample std error (=√1.32) | 1.150 |

t= 2.58 (for 99% CI)

| | |
|---|---|
| LL | 77.262 |
| UL | 83.197 |
| CV=100*σ/μ | 15.240% |

First I need to, then total number is 113. Now I need to find out the average mean of the entire data set, so what do I do 57 x 2, 62 x 4 and so on, then I add up all of them divided

by 113 so I get an average of 80.23, this is the average heart beat of all the students. But if I want to get a confidence interval on this data because this is a sample, obviously I need to calculate S. If you remember the previous equation I need to calculate S, because S is important and t will be 2.58 in this case square root of 113 that is what is going to happen here.

So how do I calculate this? This is the mean. So, every time I will subtract from the mean and then I will arise square, right? You know how to calculate standard deviation that is not very difficult. So, you get like this, the sample variance will be divided by 113 that is 1.32 and then when I take square root of that, that will be 1.15 because here I need to take square root, right. So, I am taking the variance itself and then taking the overall square root that is what I am doing, so I get 1.15. Now t as I said is 2.58,

so 2.58 x 1.15 that will be + or -  on 80, 2.58 x 1.1 is around 3. That is why the lower limit will be 80 - approximately 3 and the upper limit will be 80 + approximately 3, that is 83. So, the confidence interval gives you the estimate of where the population mean will lie, based on the sample mean. The sample mean is 80.2 and the sample standard error is 1.15 that is s by square root of n is called as I mentioned here the standard error of the mean  $s/\sqrt{n}$  here. I multiplied by 2.58 because t is given by 2.58.
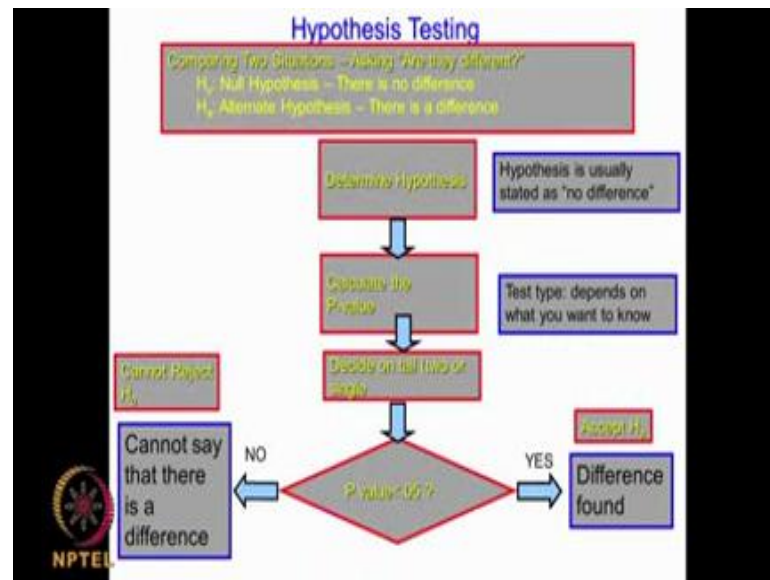
(Refer Slide Time: 15:31)



And that is for 112 degrees of freedom, for 99 % so I use that. Now CV, CV is my coefficient of variation. Formula is  $CV = 100 * \sigma/\mu$

, so the μ is a 80.2 and then we have σ given here, so we can multiply by σ divided by 80.2 and that will give you my what, the coefficient of variation. Quite a simple problem, we can calculate. So, this is the upper and lower limits for the heart beat with 99 % confidence based on a sample of 113 students, understand?

(Refer Slide Time: 16:58)



So, there are many statistical test that are available which we need to use to compare 2 sets of data and say whether one set of data is different from another set of data or they are similar and so on actually, so there are many many test here. Before you perform any test we need to create the hypothesis. As I mentioned you have the null hypothesis and the alternate hypothesis, $H_o$ and $H_a$ or some people use $H_o$ and H 1 and so on.

So, $H_o$ is no difference status and so on actually, $H_a$ is there is a difference or drug a is better than drug b or drug a is less than drug b and so on that is the alternative. So that is the hypothesis. So, you create the hypothesis first and then after that you have the data sets, you calculate the p value, then you decide on whether it is a one-tail or two-tail, I did talk about these tails also. If I am comparing 2 drugs and I am saying there is no difference between drugs or alternative there is the different between drug then we use a two-tail test.

So if I am comparing the heights of student in a class 11 a, with students in height in 11 b and I am saying there is no statistical difference in their heights, my null hypothesis will be $H_o$ null hypothesis there is no difference in the heights, $H_a$ will be there is a difference

in the heights. So in such situations we are not bothered whether the height of class a is higher or more than the higher of students in the class b or height of students in a class a is less than that of the height of the students in the class b. We are not bothered about greater or less, but we are just saying they are different. In such situations we use the two-tailed test, whereas if we are talking about greater or less things like that then we use a single tail test.

So you decide on the tail, then you decide on the p value am I looking a 95 % confidence interval or I am looking at 99 % confidence interval. So 95 % p will be less than 0.05, 99 %  p will be less than 0.01. I do that calculation and then s p is less than 0.05 or p is not different and so on actually. Either I will reject the null hypothesis or I will not reject the null hypothesis. So if there is no difference p is less than 0.05, no then I cannot reject the null hypothesis, s p is less than 0.05 then yes I need to reject the null hypothesis then only I will accept the alternative hypothesis.

Initially you start with the hypothesis $H_o$ and $H_a$ then you decide on your p value should be 95 or 99, then you decide on a tail should it be single tail or double tail, then you calculate the p value and then you compare p value less than 0.01 or 0.05, then if it is less then obviously will reject the null hypothesis so you will accept in a alternative hypothesis, p value is not less than 0.05 or 0.01 so you do not reject null hypothesis. So there is no reason for you to reject the null hypothesis. This is how any statistical analysis is carried out and this is called the hypothesis testing. So how do you calculate the p value, that depends upon what you are comparing, am I comparing means between data, am I comparing mean of a data with some population mean, am I comparing variation of one data set the variation of other data set. So you have different types of test t test, f test, chi square test so many different test and that is what you do here to calculate the p value.
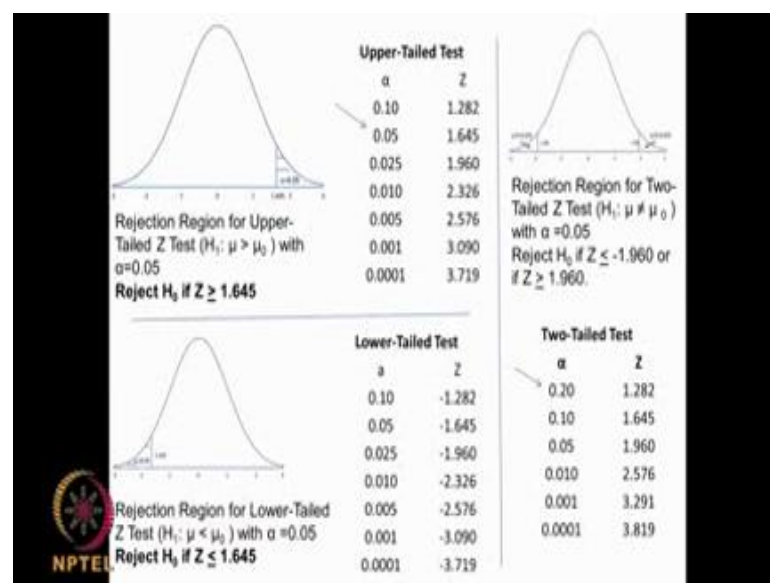
(Refer Slide Time: 21:30)



There are many, many types of test. For a two tailed and a single tailed test for example, in a null hypothesis suppose you are looking at sleeping pattern because of the drug we can say there is no difference, whereas alternative could be there is a change at a $\alpha$ of 0.5 that means 95 % confidence interval. Then it becomes a two tailed tests. There is a change, we are not a saying greater change or less change, but there is a change. We use two tail both sides, but there is a change that means a $\mu$ which you calculate is different from $\mu_0$ so there is the difference. So, it is a two tail test.

Whereas if you say, the $\mu$ sleeping pattern is greater than the old sleeping pattern, $\mu_0$ could be 7 hours, so there is an increase in this sleeping patterns then it is in upper tail test. That is a single tail test or $\mu$ is less than $\mu_0$ that means, sleeping pattern is less than 7 hours there is a decrease then it is a lower tail again it is a single tailed test. So in these 2 type of situation your $H_o$ will be $\mu$ equal to $\mu_0$, $H_1$ will be mu equal to greater than $\mu_0$ or $\mu$ equal to less than $\mu_0$, depending upon whether I am looking at the drug enhancing sleeping pattern or drug reducing sleeping pattern, then it is a single tail test. Whereas if you have $H_o$ equal to $\mu$ is equal to $\mu_0$ $H_1$ is equal to that is H alternative hypothesis is equal to $\mu_0$ is naught equal to $\mu_0$ , then we use a two tail test. Do you understand, this very very important, how to 0 in on the tail, should I 0 in on single tail or should I 0 in on two tailed test. That is a very important when you do a statistical analysis as is showed here, we need to decide on the tail because most of the tables whether it is the t table or any other $Z$ table, area and the curve which a single tail is only 1 side, whereas if it is a

double tail we need to consider both the sides of that area, if you remember that very clearly in our previous lectures. So you formulate the hypothesis and then you calculate your p value based on the type of equations we used, then you decide on whether it is a single tail or a two tailed test, then you decide on a am I going to look at it at 95 % or 99 %, then you say the p you are calculated is it less than 0.05 or it is greater than 0.05, if it is less than 0.05. And if it is a greater than, obviously, we cannot say null hypothesis can be rejected. Whereas if it is very small value we can say there is no reason for rejecting the null hypothesis.

(Refer Slide Time: 24:27)



So these upper tailed and lower tailed let me again spend some time. You know the normal distribution and you know the area, as the outside area is more important than the inside area. If you are talking about 95 % confidence the outside area will be 5 % or 0.05. So one side will be 0.025 other side will be 0.025. For a t value 1.96 for a two tailed this area will be 0.025 and that area will be 0.025 that is called two tailed test. So for a two tailed test for different alpha as we can see that is this is 95 %, this is 90 percent, this is the 80 %nt and this is 99 p% and so on actually.

So for a 95 % two tailed test z is equal to 1.96, just like t is equal to 1.96. This area and this area are equal amount actually, so it is 2.5 % and 2.5 %. Whereas if you are taking single tailed test whether its upper tailed or lower tailed, you are considering only one side of the area. For a upper tail you are considering this side of the area, for lower tail

you are considering only this side of the area, remember that. So for a 95 % confidence interval upper tail test you have 1.6425, sorry 1.645. As you can see when upper tailed 0.025 you get 1.986, if you again go back here 0.05, 1.96 because one side 0.025 other side 0.025 together is 0.05, which is 1.96. You see this, you see this.

For as upper tailed test your alternative hypothesis could be $\mu > \mu_0$ and the null hypothesis $\mu = \mu_0$. So for a lower tailed test again it is the same things here, a Z is equal to - 1.96 here will give 0.025. A Z is equal to - 1.645 will give you 0.05. It is the same, but only thing is here this Z even the signs have changed, area is the same but signs are changed. So for a lower tailed test your null hypothesis could be $\mu = \mu_0$ and alternative is $\mu \leq \mu_0$. For a $\alpha = 0.05$, if we get $H_0$, if Z is equal to less than 1.645 so you should get $Z \geq t$ equal to. So Z and t seem to be analogous to each other as you can see in this particular statistical area.

You have 3 types of situations one is called the two tailed region that means, if it is 95 % outside area is 5 %. So you divide equally on both the sides so you get 2.5 % and 2.5 %. For a two tailed test if we look at that table for a 0.05 you get this value as 1.96 or t as 1.96. Whereas in a single tailed test you are considering only one side of the area, if it is an upper tail you have area on the upper side, if you have the lower tail you have the area on the lower tail. So for a 95 % upper tailed test the Z value or t value will be 1.645. For a 95 % lower tail test again the Z value be - 1.645 and as I said if you look at 0.025 $\alpha$ you get one 1.96 which is similar to this because, when I take 0.05 we are considering 0.025 plus 0.025 to get 0.05 and that is why this term here and this term here matches. Do you understand the entire logic here. This is how you go about doing a test of significance for any data set.

So again, let me go back. So we create the null hypothesis null, hypothesis is there is no difference or status co the alternative hypothesis could be there is a statistical difference between the new data set and the old data set, then we use a two-tailed test because when we are saying there is a difference we do not say whether the difference is more or less. Whereas if the alternative hypothesis, at the new data set is greater than the old data set could be height, could be iq, could be drug, then we use single tailed test this is called a upper tailed test because greater. When we say the new data set is lower than the old data set than we use a lower tailed test or again that is a single tailed test.

So you decide on the type of test and then you decide on what should be my p value at which I am going to study? I may going to study at 95 % confidence or 99 %. So for a 95 percent, p = 0.05, so from that 0.05 I can get a t value or a Z value. Then I calculate my t or a Z value based on the type of test I am going and then I tell whether this calculated p value is much less then what is given in the table, if it is less yes, I can reject the null hypothesis and accept the alternative hypothesis.

If the p value which I calculate is not then obviously, I cannot reject the null hypothesis. The t value which I calculate is much large than the table value, then I need to reject null hypothesis. If the t value I calculate is much less than the table value, then I cannot reject the null hypothesis and that is how you go over doing statistical comparison of data sets either comparison of means or comparison of variance or comparison of ratios and so on actually.

Thank you very much.

Key words - null hypothesis, Alternate Hypothesis, p value, Z, data sets, population, sigma, $H_0$, $\mu_0$, Upper-Tailed Test, lower-Tailed Test, **Two-Tailed Test**, **standard error of the mean** (SEM), Single-Tail Z Table