

**Data Analysis and Decision Making – I**  
**Prof. Raghu Nandan Sengupta**  
**Department of Industrial & Management Engineering**  
**Indian Institute of Technology, Kanpur**

**Lecture – 17**  
**Statistical Inference**

Very good morning, good afternoon, good evening my dear friends. Welcome, this to this Data Analysis and Decision Making - I course on the NPTEL, MOOC series. And this is a 12 week course of 30 hours, each week as you know we have 5 lectures and each lecture being of half an hour. So, we are in the 17th lecture which is we have just started the 4th week and we are in the second class for the 4th week. And I am Raghu Nandan Sengupta from the IME Department, IIT, Kanpur.

So, if you remember we were just in the last class we were discussing about interval estimation, and I said that to wrap it up I said that for the mean value for the interval estimation given the parameter estimate is known from the sample we will use the z distribution or the t distribution. And if for the standard deviation to be found out we will use the chi-square distribution and the f distribution. And the other distributions considering can also be used using the concept of central limit theorem and all these things.

So, we will try to do a problem and then discuss. So, the problem is basically what we left in the discussion in the last class which is the 16th the last 3- 2 slides. So, you will first cover the concept then do a problem. And in case if it is needed we will just discuss the concept and then go further on continue, but with the promise that we will utilize those concepts where for where we have not utilized the problems to basically cleared them in more details, will clear them in a much more bigger scope where things would be much holistic in their concepts or that I am sure the students would understand it much better. So, this is the example for the mean with standard deviation variance known.

(Refer Slide Time: 02:13)

**Statistical Inference: Interval Estimation**  
**(Example for mean with SD/variance **known**)**

The height,  $X$ , of boy studying in class II of any school in the city of Bhopal, is normally distributed with mean  $\mu=125$  cms and variance  $\sigma^2=100$  cms<sup>2</sup>. We also know the heights (in cms) of 5 such boys who have been selected as the sample are, 120, 100, 110, 140 and 130, then what is the probability that the height of any boy selected at random ((i) from the whole population and (ii) from this sample) will be between 120 cms and 130 cms.

Data Analysis & Decision Making      R.N.Sengupta, IIM Dept., IIT Kanpur      196

The height  $X$  of boy of boys studying in class II of any school in city of Bhopal is normally distributed with mean 125 centimeters and variance sigma square is 100 centimeter square. So, remember the variance for the whole population for all the students who are studying in the city of Bhopal in class II, the whole population variance is given as sigma square which is there 100. I am only talking about the values. So, I am not going into the units.

We also know the height in centimeters of 5 such boys who have been selected as the sample are given as 120, 100, 110, 140 and 130. Then we want to find out that what is the probability, there are two questions what is the probability that the height of any boy selected at random number 1, from the whole population and from this been selected will be between 120 and 130. So, basically we have between 120 and 130 we have to find out that the boy being selected from the whole population would be in that group and boy selected from the sample would be in that group. So, what is that group? That group is basically the interval 120 to 130. So, you have to basically solve it.

Before I give you the solution first let us see the problem what are the informations given. So, the informations given are basically the population mean and the population variance. So, population mean and population variance are given and you want to find out something to do with the mean so obviously, we know that as the population variance is given we have to use the z distribution. So obviously, it will be t if the population

variance was not known, not given. So, for population mean we they either we use the z or the t.

So, the continuation is the first bullet point says then X is basically normally distributed with 125 as the mean and the variance is 100. Now, we need to find out that what is the distribution of 5 observations taken in a group and what is the distribution of their sample mean.

(Refer Slide Time: 04:13)

**Statistical Inference: Interval Estimation**  
(Example for mean with SD/variance **known**  
contd...)

- $X \sim N(\mu=125, \sigma^2=100)$
- $X_{\text{mean},5} \sim N(\mu=125, \sigma^2/5=20)$ , where  $n=5$ ,  
 $X_{\text{mean},5}=125$

Thus from the given information we know that

- $\Pr\{120 \leq X \leq 130\} = \Pr\left\{\frac{(120-125)}{10} \leq X \leq \frac{(130-125)}{10}\right\} = \Pr\{-0.5 \leq X \leq 0.5\} = 0.3830$
- $\Pr\{120 \leq X_{\text{mean},5} \leq 130\} = \Pr\left\{\frac{(120-125)}{\sqrt{20}} \leq X_{\text{mean},5} \leq \frac{(130-125)}{\sqrt{20}}\right\} = \Pr\{-1.12 \leq X_{\text{mean},5} \leq 1.12\} = 0.7372$

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 197

So, here the suffix X mean means and comma 5 means that we are taking a bunch of 5 finding or they mean and basically trying to find out what is the distribution of that mean.

So, X mean comma 5 means bar 5. So, this is also obviously, we know that will be normally distributed, but 125 mean value, but now in this case the total variance would be if you remember the sigma square by n. So, it will be 100 by 5 it would be 20. So, where n is equal to 5, and the sample mean is given by 125, the mean yeah as we found out.

And now from the information what we need to know is that, what is the probability that the value of X would be lying between 120 and 130, and in the second case that what is the probability that the sample mean would be lying between 120 and 130. So, let us do the problems accordingly in the first case X is the random variable. So, let us ask that

what is the distribution of  $X$ ; we know the distribution of  $X$  is normal. What is the mean? I my answer is 125 and sigma square is 100. So, we you need to convert that into a standard normal and proceed with our calculations.

So, the probability that it  $X$  is between 120 and 130 is given by this let me highlight. So, I need to convert  $X$  into the standard normal, as I do that this is basically the  $X_1$  and  $X_2$  which I have the corresponding LCL and UCL for the  $X$ , they would also be converted into LCL and UCL for the  $z$  distribution. So, they would also be transformed accordingly.

So, here how do we convert the  $X$  into a standard normal deviate? It will be  $X$  minus the expected value of  $X$  divided by the standard deviation exactly what we do is basically this  $X$  is now converted into  $Z$ . So, we have not written it, so my apologies for that. So, this would basically be converted into a  $Z$ . So, technically what we are doing is  $X$  minus  $\mu$  by  $\sigma$ . So,  $\mu$  is basically given as 125 which is here the  $\sigma$  is basically 10 we convert that into a  $Z$  value the capital  $Z$  standard normal and these are basically  $Z_1$  and this is also  $Z_2$ . So, these  $Z_1$  and  $Z_2$  are corresponding to 120 and 130.

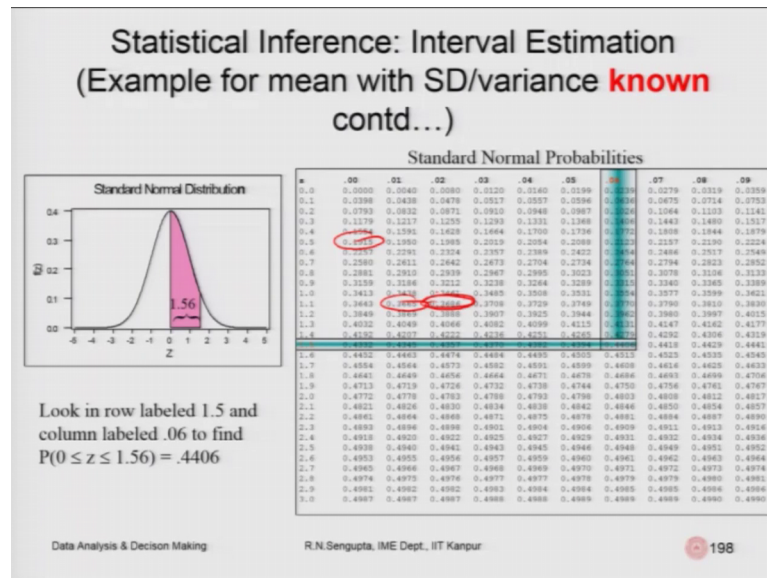
So, hence for the  $Z$  it is basically bounded between minus 5 to plus 5 the value basically comes out with the probability. So, here you have the distribution the mean value is 0, this is 0 points minus 0.5, plus 0.5 the whole area bounded between them is about 38 percent at is given and hence the probability is 0.38.

Now, go let us go back to finding out that what is the probability of the mean value being between that. So obviously, again we find out the same stand standard formula probability that  $X$  mean of that corresponding to the sample size of 5 you will be between 120 and 130. We know the sample mean is distributed with normal distributed with a mean value as given as 125 and the standard deviation or let us talk about the variance sigma square by the sample numbers which is 5 which comes out to be 20. So, hence it will be the; it will be bounded between 120 and 130.

Again convert this  $X$  mean 5 into a standard normal deviate which is basically the  $Z$  value and the corresponding lower control and the upper control will be calculated accordingly which will be, in the first case the lower control will be 130 minus 120 minus 125 that is and divided by square root of sigma by root of  $n$  which is technically

square root of sigma square by n that which comes out to be square root of 20 put that values. And on the right hand side again you will have that value of 130 minus 125 which is plus 5 divided by square root of sigma square by n. So, put those values and that value comes out to be 0.732 in this case. That means this is deviating more on to the left more on to the right hence the overall coverage would be much more than 0.38. So, it is coming out to be almost double.

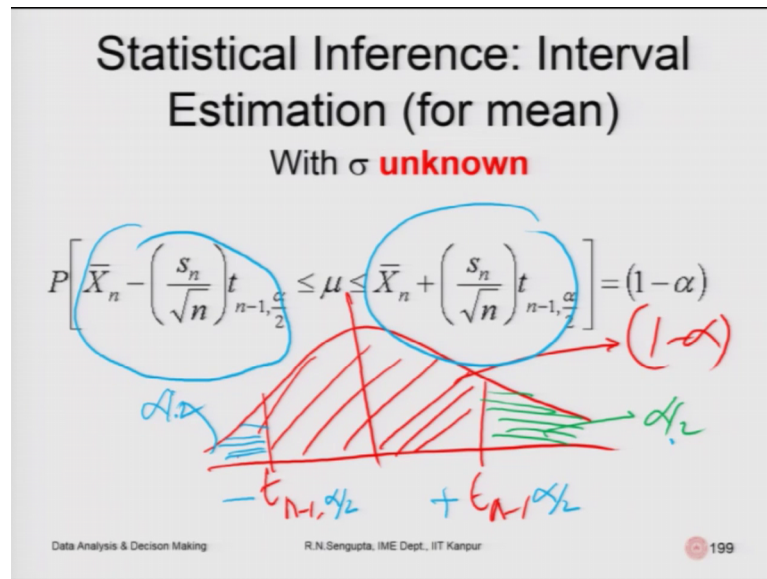
(Refer Slide Time: 09:12)



So, these values which I have seen I am just giving the standard normal table also in order to understand so minus 5 and plus 5. So, minus 5 value is basically 0.195, so double of that because you are going equally disposed onto the left hand side on the right hand side it comes out to be 0.38 as the value.

For the case when you are using 0.7, so let me go back till the value it is coming out to be 0.732. So, it is plus minus 1.12. So, 1.12 would come out to be this. So, this would be make about this, so 0.3686. So, both 0.3686 on the left hand side and the right hand side add them up or multiply them by 2 comes out to be 0.72 which is 72 percent.

(Refer Slide Time: 10:15)



Now, consider the case of the statistical inference interval estimation for the mean with the standard deviation being unknown. So, none the first question comes is that if you are interested to find out something with this with the mean value for the population. So, what distribution you should be used. To answer as we know is basically the z or the t.

Now, why we are using t because the case is that the standard deviation of the or the variance of the population is not known, if it is not known mean we know we have to use one of the s s. And when we use this one of the s s, I am not going to come into that detail what I am basically highlighting the main points is that we will basically lose one degrees of freedom hence you will basically be utilizing the t distribution. Hence t distribution with me n minus 1 degrees of freedom, but remember one thing here if I draw the t distribution let me draw it. So, this is not a normal, this is t.

So, this is the mean value and we want to find out this limits. So, this is basically so called t n minus 1 and this would be t n minus 1 but the and this whole area as we know is equal to 1 minus alpha by 2; alpha. This area let me change the color, let me change it to blue or green here. This will be basically alpha by 2; this area is also alpha by 2.

So, hence some of alpha by 2, alpha by 2 and 1 minus alpha would basically be me the total area under that curve is 1. Now, remember t distribution have I mentioned is symmetric. So, if it is symmetric the values on to the left hand side and right hand side

would basically given by alpha by 2 here, alpha by 2 here, but only remember the mean value being 0. So, any value on to the right hand side would be positive anywhere on to the left hand side would should be negative.

So, we use the t distribution, but in this case remember also we would not be using my s dash, because s dash is only used for the case when the sample the population mean is known, but in our case the population mean is not known. So, you will be utilizing s without the dash. So, hence it will be divided by 1 by n minus 1. Why 1 n minus 1? Please bear with me I will be repeating these things time and again because we lose one degree of freedom considering that  $X_1$  to  $X_n$  have been obtained you will use them for the first time to find out what is the best estimate of the population mean given the sample mean is the best estimate which we know.

So, again we will we will basically formulate the problem in this way where  $\bar{X} - n$  minus this term and  $\bar{X} + n$  plus this term would basically with the corresponding lower control and the upper control correspondingly.

(Refer Slide Time: 13:10)

**Statistical Inference: Interval Estimation**  
(Example for mean with SD/variance **unknown**)

A food inspector examines 10 jars of certain brand of butter and obtained the following percentages of impurities, the values of which are, 2.3, 1.9, 2.1, 2.8, 2.3, 3.5, 1.8, 1.4, 2.0 and 2.1. Form a 90% level of confidence (loc) for the estimate of the mean of the impurity level, where you can assume the population distribution of the level of impurity as normal, i.e.,  $X \sim N(\mu, \sigma^2)$

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 200

So, we will do a again a problem, small simple problems. A food inspector examines and this is the case for the population variance being unknown. A food inspector examines 10 jars of certain bread of butter, brand of butter and obtain the following percentage of the impurities the values of which are given as 2.3, 1.9, 2.1 so on and so forth. You want to

basically form a 90 percent level of confidence loc, for the estimate of the mean of the impurity levels such that where you can assume the population distribution of the level of impurities basically given in a normal and we want to basal comment intelligently according to that.

So, your values or impurities are given, the level of confidence 90 is given. Level of confidence means the overall area between the lower control limit and the upper control me. So, if I basically have and the distribution; again I am not trying to specify with this a normal distribution, t distribution chi-square distribution or f distribution. So, these are the limits the overall area we in between is basically given by level of confidence which is 1 minus alpha.

(Refer Slide Time: 14:17)

**Statistical Inference: Interval Estimation**  
(Example for mean with SD/variance **unknown**  
contd...)

- $n=10, \bar{x}_{\text{mean},10}=2.22, s_{10}=0.5789$

Thus from the given information we know that

- $\Pr\{\bar{X}_{\text{mean},n}-(s_n/\sqrt{n})t_{n-1,\alpha/2}\leq\mu\leq\bar{X}_{\text{mean},n}+(s_n/\sqrt{n})t_{n-1,\alpha/2}\}=(1-\alpha)$
- $\Pr\{2.22-(0.5789/\sqrt{10})1.833\leq\mu\leq 2.22+(0.5789/\sqrt{10})1.833\}=0.90$
- LCL=1.8445 and UCL=2.5556

Data Analysis & Decision Making
R.N.Sengupta, IME Dept., IIT Kanpur
201

So, from this n value is given as 10, the mean value X X mean of with the sample mean of 10 is given by 2.22, and this s without the dash because we are utilizing the sample mean as the best estimate of the population mean based on which when we calculate we find out the standard error comes out to be 0.5789. Thus, from this given information we know that the lower control limit and the upper control limit are in such a way that X mean of this ends observations.

And on the left hand side if we remember to be s n divided by the square root of n because here again you will basically be dividing by square root of n depending on the



set of observation which is there for the sample. Into  $t_{n-1}$  into  $\alpha$  and  $t_n$  on the right hand side also it will be  $t_{n-1}$  into  $\alpha$ ;  $\alpha$  by 2. But remember technically they can be written later on we will see for the chi and the z, chi and the f distribution on the right hand side values would be given by  $1 - \alpha$  by 2 depending on whatever area you have to cover on the right hand side I am going to come to that later.

That is for the given information we know that probabilities being such that that the lower control limit which is  $\bar{X}$  mean for this  $n$  observations and  $s_n$  being the square root  $s_n$  divided by square root of  $n$  into  $t$  the  $t$  distribution the value. This  $t_{n-1}$   $\alpha$  by 2,  $\alpha$  by 2 is basically depending on the level of confidence which you have, so that and on the right hand side the upper control limit would be  $\bar{X}$  mean  $n$  and again the same thing, but with the same a plus sign here. So, once you put this values the lower control limit and upper control limit you can check that comes out to be 1.8445 and upper control has 2.5556.

(Refer Slide Time: 16:18)

**Statistical Inference: Interval Estimation**  
(Example for mean with SD/variance **unknown** contd...)

df/p	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	0.314267	1.000000	3.077684	6.313752	12.70620	31.82082	63.65734	636.6192
2	0.287595	0.816497	1.885618	2.919986	4.302450	6.964560	9.246444	31.52611
3	0.274777	0.764882	1.637744	2.353363	3.182446	4.540750	5.840811	12.92400
4	0.270722	0.740008	1.533206	2.131847	2.776445	3.746950	4.604090	8.61033
5	0.267191	0.720877	1.475884	2.015048	2.575829	3.364910	4.033143	6.96456
6	0.264335	0.711558	1.439756	1.943180	2.449911	3.142672	3.707432	5.95885
7	0.262167	0.711142	1.414924	1.894579	2.364622	2.997950	3.489488	5.40779
8	0.260527	0.708827	1.396715	1.855837	2.309987	2.909945	3.352636	5.04131
9	0.259355	0.707222	1.383029	1.831113	2.262165	2.821444	3.248844	4.78039
10	0.258585	0.699812	1.372184	1.812461	2.228134	2.762277	3.169227	4.58689
11	0.258056	0.697445	1.363430	1.798850	2.200891	2.718088	3.105811	4.43370
12	0.257683	0.695483	1.356217	1.789288	2.178911	2.681160	3.054544	4.31778
13	0.257401	0.693829	1.350171	1.780933	2.160372	2.650231	3.012228	4.22088
14	0.257173	0.692417	1.345030	1.773210	2.144279	2.624469	2.976884	4.14005
15	0.257000	0.691197	1.340606	1.766050	2.129140	2.602348	2.949711	4.07228
16	0.256879	0.690132	1.336757	1.759384	2.115931	2.583449	2.926078	4.01580
17	0.256797	0.689195	1.333379	1.753067	2.103982	2.566953	2.905223	3.96511
18	0.256742	0.688364	1.330291	1.747064	2.100022	2.552228	2.878444	3.92116
19	0.256703	0.687621	1.327428	1.741313	2.096052	2.539048	2.860063	3.88284
20	0.256673	0.686954	1.325341	1.735718	2.092066	2.527398	2.845434	3.84905
21	0.256650	0.686352	1.323888	1.730273	2.087961	2.517195	2.831368	3.81933
22	0.256632	0.685805	1.322727	1.724944	2.083727	2.508262	2.818776	3.79271
23	0.256619	0.685306	1.321800	1.719672	2.079366	2.499877	2.807324	3.76775
24	0.256610	0.684845	1.321036	1.714422	2.074872	2.491916	2.796944	3.74454
25	0.256605	0.684420	1.320395	1.709141	2.070254	2.484351	2.787444	3.72251
26	0.256603	0.684023	1.319822	1.703878	2.065513	2.477083	2.779711	3.70166
27	0.256602	0.683650	1.319303	1.703288	2.061632	2.470080	2.772880	3.68190
28	0.256602	0.683293	1.318827	1.701131	2.058441	2.463214	2.766226	3.67239
29	0.256602	0.682944	1.317838	1.699127	2.055222	2.456462	2.760229	3.66308
30	0.256602	0.682600	1.317041	1.697261	2.052222	2.449726	2.754880	3.65400
$\chi^2$	0.253347	0.674480	1.281552	1.644854	1.959966	2.202635	2.575832	3.29095
Ct			80%	90%	95%	98%	99%	99.9%

90  
70%  
5% 8%  
005

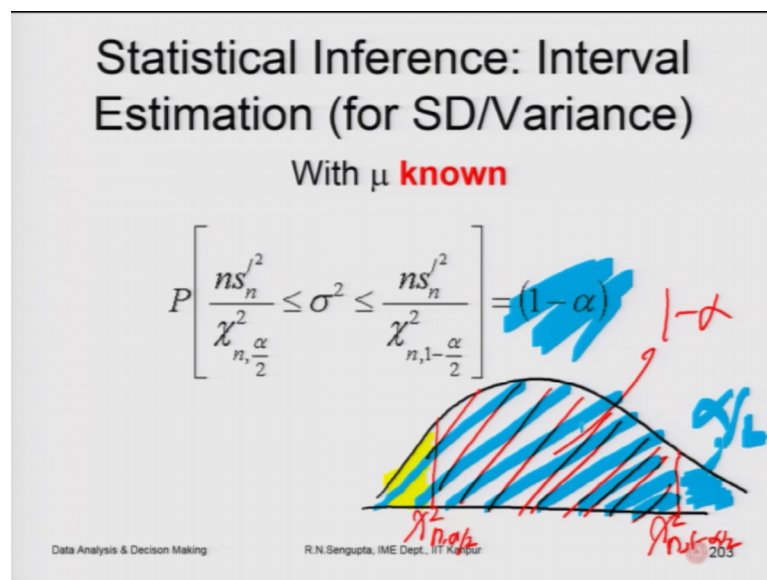
Data Analysis & Decision Making R.N. Sengupta, IME Dept., IIT Kanpur 202

So, again I have added, the this is the first time you is seeing the  $t$  distribution table on the  $t$  distribution table on the leftmost column you have the degrees of freedom, on the top most value on the row you have the  $p$  values and the inside values are basically the  $t$  values.

So, remember one very very interesting thing when we using the z distribution it gives you the actually the overall coverage of the probabilities. Hence, if you add up all the values starting from minus infinity to plus infinity the probability values for the z distribution there would be 1. But the values which are inside here are the corresponding t values based on which you are going to do with the calculation it does not give you the cdf value. So, be careful about that that would also be true for the f distribution that should also be true for the chi-square distribution.

So, for a value of 29 and if you see the value being 29 because we have lost 1 degree no this would be basically sorry my apologies it would be 9. So, basically we go to 9. So, the first we basically mark this, then let us see what is the degrees of freedom, degrees of freedom is basically given by 9 which we know 1 minus alpha is known. So, if we have, so basically it was 90 percent so the overall area left and right is 10 percent. So, this is would be 10 percent is the total coverage. So, it will be divided into 2 equal halves 5 percent, 5 percent. So, if you see 5 percent, 5 percent technically 0.05. So, for the 0.05 value it comes out to be 1.83313 utilize that value to do the calculation.

(Refer Slide Time: 18:19)



Consider the case when you want to find or do some statistical inference, but with respect to finding or something further for the variance of the standard deviation given the mean value is known. So, this is the first instance where we want to find out the interval estimation or the population variance.

Now, if the mean value is known immediately it will occur to your mind that we should be utilizing  $\bar{s}$  and not  $s$ , because  $\bar{s}$  basically means that the population mean is known, so we do not use the sample mean and the degrees of freedom remains as  $n$ , point one noted. Now, when we utilize that remember that the chi-square distribution which is being formed would also would not lose 1 degrees of freedom. So, it will basically be chi-square  $n$  comma  $\alpha$  by 2 and the right hand side it will be  $1 - \alpha$  by 2 and this is the first time I am trying to basically explain why it is like this, and see that this is the chi-square.

So, you will have chi-square on the left hand side this is  $n$  this is  $\alpha$  by 2, this is the second one chi-square  $n$   $1 - \alpha$  by 2 this whole area inside is  $1 - \alpha$ . Now, see how the nomenclature is this, if you see the suffix it gives me gives me  $\alpha$  by 2 that means,  $\alpha$  by 2,  $\alpha$  by 2; yes,  $\alpha$  by 2 is the area on to the left of the chi-square first value and  $1 - \alpha$  by 2 is the overall area which is on to the left.

So, we are basically following the same nomenclature as considering what is the overall area covered on to the left. See if I see this, so  $1 - \alpha$  by 2 technically means the whole area we have covered from the minimum value chi-square to that value where it is chi-square this whole area is  $1 - \alpha$  by 2. And the right hand side which  $1 - \alpha$  by 2 would be for the case wait,  $\alpha$  by 2 should be  $1 - \alpha$  by 2. This would be  $\alpha$  by 2, yes, this area should be also  $\alpha$  by 2 because the total sum of the areas on the left hand side the middle portion and the right hand side should be  $\alpha$  by 2. So, I will basically highlight this in more details because let me do the problem I will come to this conceptual framework in more details later on.

So, in this case we use the chi-square with without losing one degrees of freedom to recap important things. We will use  $\bar{s}$  with the dash that means, we are not losing any degrees of freedom, so which is basically  $\mu$  is  $\mu$  not  $\bar{X}$  because the sample mean is not to be utilized. And once we have that the overall rule is probability of  $n \bar{s}^2$  by  $n$  divided by,  $n$  means the sample size divided by chi-square with degrees of freedom  $n$  and  $\alpha$  by 2 b where the position and is depending on the level of confidence which we have. And the right and the right limited the UCL would be again same thing, but the chi-square value would be  $n$  into  $1 - \alpha$  by 2 depending on the nomenclature is and that overall area would be in between would be level of confidence which was  $1 - \alpha$  which is this.

(Refer Slide Time: 22:25)

Statistical Inference: Interval Estimation  
(Example for SD/Variance with mean **known**)

Suppose a sample of 30 school students are given an IQ test. If the sample has a standard error of 12.25 points, find a 90% confidence interval for the population standard deviation. Assume population mean is **known**

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 204

So, let us do a problem and I will highlight it and come to the basic concepts more in more details. Suppose a sample of 30 students, school students are given an IQ test. If the sample has a standard error of 12.25, so we are basically considered the stand the sample standard error is 12.25. Find a 90 percent confidence interval for the population standard deviation assume that the population mean is known. So, if the population mean is known then this 12.25 is basically  $s/\sqrt{n}$  and being 30. That means, given the standard error for the sample and provided the population mean is known we know that is the standard error such that we have not lost any degrees of freedom.

(Refer Slide Time: 23:24)

**Statistical Inference: Interval Estimation**  
(Example for SD/Variance with mean **known**)

$$\chi^2_{n,1-\alpha/2} = \chi^2_{30,0.95} = 18.493 \quad \text{and} \quad \chi^2_{n,\alpha/2} = \chi^2_{30,0.05} = 43.733$$

Then the confidence interval is

$$(ns^2_n/\chi^2_{n,\alpha/2}) \leq \sigma^2 \leq (ns^2_n/\chi^2_{n,1-\alpha/2})$$

$$(30 * 12.25^2 / 43.733) \leq \sigma^2 \leq (30 * 12.25^2 / 18.493)$$

$$102.939 \leq \sigma^2 \leq 243.437$$

$$10.145 \leq \sigma \leq 15.605$$

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 205

So, in this case continuing from the table we will come to the table, chi-square n and 1 minus alpha by 2 this is being on the right hand side would basically be chi-square, what is n? N is 30 because you are not losing any degree freedom and 1 minus alpha by 2 basically be 0.95 that comes out to 18.493. And chi-square n alpha by 2 on the left hand side basically comes out to be 43.75. Then the confidence interval is found out in the same way as we have just noted down, and the confidence interval comes out to be for sigma comes out to be 10.145 and 15.605.

(Refer Slide Time: 24:02)

**Statistical Inference: Interval Estimation**  
(Example for SD/Variance with mean **unknown**  
contd...)

P	P(X ≤ x)							
	0.010	0.025	0.050	0.100	0.500	0.950	0.975	0.990
1	0.000	0.001	0.004	0.016	2.706	3.841	5.024	6.635
2	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210
3	0.115	0.215	0.352	0.584	6.251	7.815	9.348	11.345
4	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277
5	0.554	0.831	1.145	1.610	9.236	11.070	12.833	15.086
6	0.872	1.237	1.635	2.204	10.645	12.592	14.449	16.812
7	1.239	1.690	2.167	2.833	12.017	14.067	16.013	18.475
8	1.646	2.180	2.733	3.490	13.362	15.507	17.535	20.090
9	2.088	2.700	3.325	4.168	14.684	16.919	19.023	21.666
10	2.558	3.247	3.940	4.865	15.987	18.307	20.483	23.209
11	3.053	3.816	4.575	5.578	17.275	19.675	21.920	24.725
12	3.571	4.404	5.228	6.304	18.549	21.026	23.337	26.217
13	4.107	5.009	5.892	7.042	19.812	22.362	24.736	27.688
14	4.660	5.629	6.571	7.790	21.064	23.685	26.119	29.141
15	5.229	6.262	7.261	8.547	22.307	24.996	27.488	30.578
16	5.812	6.908	7.962	9.312	23.542	26.296	28.845	32.000
17	6.408	7.564	8.672	10.085	24.769	27.587	30.191	33.409
18	7.015	8.231	9.399	10.865	25.989	28.869	31.526	34.805
19	7.633	8.907	10.117	11.651	27.204	30.144	32.852	36.191
20	8.260	9.591	10.881	12.443	28.412	31.410	34.170	37.566
21	8.897	10.283	11.661	13.240	29.615	32.671	35.479	38.932
22	9.542	10.982	12.466	14.041	30.813	33.924	36.781	40.289
23	10.196	11.689	13.291	14.848	32.007	35.172	38.076	41.638
24	10.856	12.401	14.134	15.659	33.196	36.415	39.364	42.980
25	11.524	13.120	14.911	16.473	34.382	37.652	40.646	44.314
26	12.198	13.844	15.670	17.292	35.563	38.885	41.923	45.642
27	12.879	14.573	16.411	18.114	36.741	40.113	43.195	46.963
28	13.565	15.308	17.128	18.939	37.916	41.337	44.461	48.279
29	14.256	16.048	17.926	19.767	39.087	42.557	45.722	49.581
30	14.950	16.791	18.693	20.599	40.256	43.773	46.979	50.890

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 206

So, this is the value of 30 and this is the value of 29 we come to that later on. So, here if it is basically 30 let me go back what was we wanted we wanted a value at n is equal to 30 and 0.05; and 0.95; 0.95 is this 43.75 and 0.05 is basically 18.493. Check these values and this is basically 18.493, and 43.733 you utilize these values to solve your problems accordingly.

(Refer Slide Time: 24:50)

**Statistical Inference: Interval Estimation (for SD/Variance)**  
With  $\mu$  **unknown**

$$P \left[ \frac{(n-1)s_n^2}{\chi_{n-1, \frac{\alpha}{2}}^2} \leq \sigma^2 \leq \frac{(n-1)s_n^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} \right] = (1-\alpha)$$

Data Analysis & Decision Making R.N. Sengupta, IME Dept., IIT Kanpur 207

Now, let us consider the problem. The same time, but with the population mean being not given. So, consider that statically inference problem is their interval estimation is their for the standard deviation of the variance, but with the population mean being unknown. So, here again there we loose 1 degrees of freedom because the population mean is not known we use the sample mean to find out the population mean as we do that we have a chi-square which has 1 degrees of freedom less hence it is n minus 1.

In the numerator in place of s dash it now would be s without the dash because you are utilizing the different standard error where mean values of the population are not known and repeating things time and again please listen to me carefully. And obviously, in place of n which was pre multiplied before s it will be n minus 1 depending on the formula. So, you one can check any good book for the formulas. So, the formulas are basically to be utilized not to be derived.

(Refer Slide Time: 25:53)

**Statistical Inference: Interval Estimation**  
(Example for SD/Variance with mean **unknown**)

Suppose a sample of 30 school students are given an IQ test. If the sample has a standard error of 12.23 points, find a 90% confidence interval for the population standard deviation.

**Note:** As no information is given we implicitly assume the population mean is **unknown**

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 208

Here again the same problem suppose a sample of 30 students exactly if the sample has a standard deviation of 12.23, now it is not 12.25. Find a 90 percent confidence interval for the population standard deviation again and now nothing is mentioned about the population mean. So, as no information is given we implicitly assume that the population mean is unknown.

(Refer Slide Time: 26:14)

**Statistical Inference: Interval Estimation**  
(Example for SD/Variance with mean **unknown**  
contd...)

$\chi^2_{1-\alpha/2} = \chi^2_{0.95,29} \approx 17.708$  and  $\chi^2_{\alpha/2} = \chi^2_{0.05,29} \approx 42.557$

Then the confidence interval is:

$$\frac{(n-1)s^2}{\chi^2_{\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2}}$$
$$\frac{(30-1)12.23^2}{42.557} < \sigma^2 < \frac{(30-1)12.23^2}{17.708}$$
$$101.9249 < \sigma^2 < 244.9472$$
$$10.10 < \sigma < 15.65$$

Data Analysis & Decision Making R.N.Sengupta, IME Dept., IIT Kanpur 209

So, here again we find out chi-square 1 minus alpha for n, and now being n minus 1 because it will be your last 1 degrees of freedom. So, it will be n minus 1 will be 29. So,

that value comes out to be for a chi-square 29 and 0.95 comes out to be 17.0708 and chi-square with degrees of freedom of 29, but alpha by 2 now being 0.05. If you look into the table I am going to come to that again it comes out to be 42.557 you utilize that in the formula and solve your problems accordingly. So, let me check the table.

(Refer Slide Time: 26:54)

**Statistical Inference: Interval Estimation**  
(Example for SD/Variance with mean **unknown**  
contd...)

r	P(X ≤ r)							
	0.010	0.025	0.050	0.100	0.500	0.950	0.975	0.990
1	0.000	0.001	0.004	0.016	2.706	3.841	5.024	6.635
2	0.000	0.001	0.103	0.211	4.605	5.991	7.378	9.210
3	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345
4	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277
5	0.554	0.831	1.145	1.610	9.236	11.070	12.833	15.086
6	0.872	1.237	1.635	2.204	10.645	12.592	14.449	16.812
7	1.239	1.690	2.167	2.833	12.017	14.067	16.013	18.475
8	1.646	2.180	2.733	3.490	13.362	15.507	17.535	20.090
9	2.088	2.700	3.325	4.168	14.684	16.919	19.023	21.666
10	2.558	3.247	3.940	4.865	15.987	18.307	20.483	23.209
11	3.053	3.816	4.575	5.578	17.275	19.675	21.920	24.725
12	3.571	4.404	5.226	6.304	18.549	21.026	23.337	26.217
13	4.107	5.009	5.892	7.042	19.812	22.362	24.736	27.688
14	4.660	5.629	6.571	7.790	21.064	23.685	26.119	29.141
15	5.229	6.262	7.261	8.547	22.307	24.996	27.488	30.578
16	5.812	6.908	7.962	9.312	23.542	26.296	28.845	32.000
17	6.408	7.564	8.672	10.085	24.769	27.587	30.191	33.409
18	7.012	8.231	9.390	10.865	25.989	28.869	31.526	34.805
19	7.633	8.907	10.117	11.651	27.204	30.144	32.852	36.191
20	8.269	9.591	10.851	12.443	28.412	31.410	34.170	37.566
21	8.927	10.283	11.591	13.240	29.615	32.671	35.479	38.932
22	9.542	10.982	12.338	14.041	30.813	33.924	36.781	40.289
23	10.196	11.689	13.091	14.848	32.007	35.172	38.076	41.638
24	10.890	12.401	13.848	15.659	33.196	36.415	39.364	42.980
25	11.624	13.120	14.611	16.473	34.382	37.652	40.646	44.314
26	12.398	13.844	15.379	17.292	35.563	38.885	41.923	45.642
27	13.219	14.573	16.151	18.114	36.741	40.113	43.195	46.963
28	13.985	15.308	16.928	18.939	37.916	41.282	44.461	48.278
29	14.799	16.047	17.709	19.768	39.087	42.437	45.722	49.588
30	15.653	16.791	18.493	20.599	40.256	43.278	46.979	50.892

Data Analysis & Decision Making R.N.Sengupta, IIM Dept., IIT Kanpur 210

So, here 29 this is 42.557 for the value of 0.9, 0.95 and the value is basically for 0.05 it is coming out to be 17.708. So, let me check, 42.557, for 0.9 for 0.95 depending on the total level of confidence being 90 because you have divided it accordingly. So, you can use these values to solve it.



(Refer Slide Time: 27:35)

**Statistical Inference: Interval Estimation**

Confidence interval (CI) can be formulated for

- Difference of  $\mu_1$  and  $\mu_2$ , provided  $\sigma_1$  and  $\sigma_2$  are **known**
- Difference of  $\mu_1$  and  $\mu_2$ , provided  $\sigma_1$  and  $\sigma_2$  are **unknown**, but **equal**
- Difference of  $\mu_1$  and  $\mu_2$ , provided  $\sigma_1$  and  $\sigma_2$  are **unknown**, but **unequal**

Data Analysis & Decision Making R.N.Sengupta, IIME Dept., IIT Kanpur 211

Statistical inference estimation can be there, confidence interval it can be there when you are trying to find out the difference of the means provided sigma 1 and sigma 2 are known. So, if it is difference on the mean and standard deviation obviously, we will use the z distribution. I am going to be just repeat it bullet points come back to the problem solving later on.

When you are trying to find out the difference of the means provide a sigma 1 and sigma 2 are unknown both being equal case an unequal case we will use the t distribution with degrees of freedom being lost, one from the first and one from the second. So, be careful it would not be now no more n minus 1; it will be n minus 2 depending on a number of odd degrees you have basically lost from the first population and the second population.

(Refer Slide Time: 28:20)

**Statistical Inference: Interval Estimation**

Confidence interval (CI) can be formulated for

- Ratio of  $(\sigma^2_1/\sigma^2_2)$  provided  $\mu_1$  and  $\mu_2$ , are **known**
- Ratio of  $(\sigma^2_1/\sigma^2_2)$  provided  $\mu_1$  and  $\mu_2$ , are **unknown**

Data Analysis & Decision Making R.N.Sengupta, IIME Dept., IIT Kanpur 212

Similarly, you can use the to find of the ratios of the standard deviations or square which is the variances provided mu 1 and mu 2 is known, we will use the f distribution without losing any degrees of freedom it will be m and n. And in the case when you have basically the finding of the ratios of the variances provided mu 1 and mu 2 are unknown you will use 1 degrees of freedom for the first population and 1 degrees of freedom from the second population. So, it will be f distribution with m minus 1 and n minus 1. I will come to that in more details later on.

So, with this I will end this 17th lecture and continue more discussion about the hypothesis testing and go in to the multiple linear regression. And then definitely go into the multivariate statistics which would be to one of the main part on the discussion how they can utilize for different type of problem solving.

Have a nice day and thank you very much.