

INDIAN INSTITUTE
OF
TECHNOLOGY
KHARAGPUR

NPTEL
National Programme
on
Technology Enhanced Learning

Applied Multivariate Statistical Modeling

Prof. J. Maiti
Department of Industrial Engineering and Management
IIT Kharagpur

Lecture – 14

Topic

Analysis of Variance(ANOVA)

Good afternoon. Today we will discuss Analysis of Variance.

(Refer Slide Time: 00:24)

Analysis of Variance (ANOVA)

No of population (L)	No of variables (P)	Hypothesis	Technique
L=1	P=1	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	t-test
L=1	P>2	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	Hotelling's T^2
L=2	P=1	$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$	t-test
L=2	P>2	$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$	Hotelling's T^2

Handwritten notes:
A circle around 'L=1' with 'L=3' written next to it.
A note 'No of population' with an arrow pointing to the 'L' column.

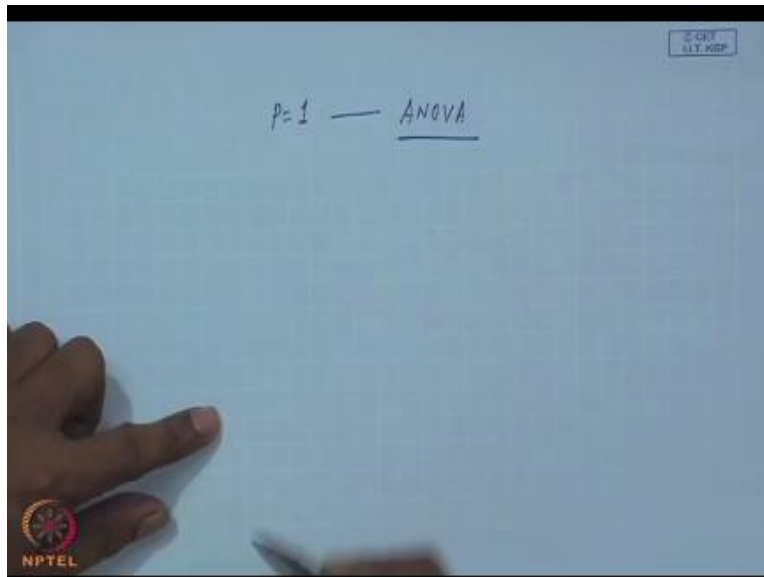
Popularly known as ANOVA. In last class, we have discussed about the difference between two population means and in hypothesis testing we have covered the equality of population means from univariate point of view from multivariate point of view. So if I just recapitulate this what we will find out that, one is number of population, then your number of variables, then hypothesis, then the techniques used, technique what is, what technique is used? If you say number of population l number of variable p and hypothesis as usual.

Then the situation is single population l equal to 1, suppose number of variable is also one $p=1$, what we do in hypothesis testing? We do $H_0, \mu = \mu_0$ and $H_1 \mu \neq \mu_0$ generally we test this and we use for the small sample case t - test, we assume that population is normal. Now l can be 1, and p can be more than 2 all case also that is single population more than or equal to two variables.

And your hypothesis will be again $H_0 \mu = \mu_0$ here μ is a vector quantity and $H_1 \mu \neq \mu_0$ and we have used Hotelling's T^2 for hypothesis testing again for small sample case. Now second issue we have discussed $l=2$ two population case, under this also we have already discussed p equal to one case and there our hypothesis was $\mu_l = \mu_m$ and alternative hypothesis we have framed as $\mu_l \neq \mu_m$ for at least one pair.

And we have used t - test and for $p \geq 2$ case we have hypothesis like this that $\mu_l = \mu_m$ here μ_l and μ_m are vector quantities and $H_1 \mu_l \neq \mu_m$ for at least one pair. We have also used Hotelling's T^2 , okay. So this is what from hypothesis testing point of view we have already completed. What will happen when l is three or more? So our discussion today for $l \geq 3$ case where l stand for population, that is number of population, okay. So, in this case we will first consider the univariate issue.

(Refer Slide Time: 04:39)



That means we are interested in p equal to 1 case and when we test equality of several population means for one response variable the technique used is ANOVA analysis of variance. That means what I said that suppose your case is $p = 1$ that is univariate case there are several groups or population which is definitely three or more, in that case the first t -test is not applicable.

(Refer Slide Time: 05:21)

Analysis of Variance (ANOVA)

No of population (L)	No of variables (P)	Hypothesis	Technique
L=1	P=1	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	t-test
	$P \geq 2$	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	Hotelling's T^2
L=2	P=1	$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$	t-test
	$P \geq 2$	$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$	Hotelling's T^2

Handwritten note: L=3 No of population

SCS IIT KGP

NPTEL

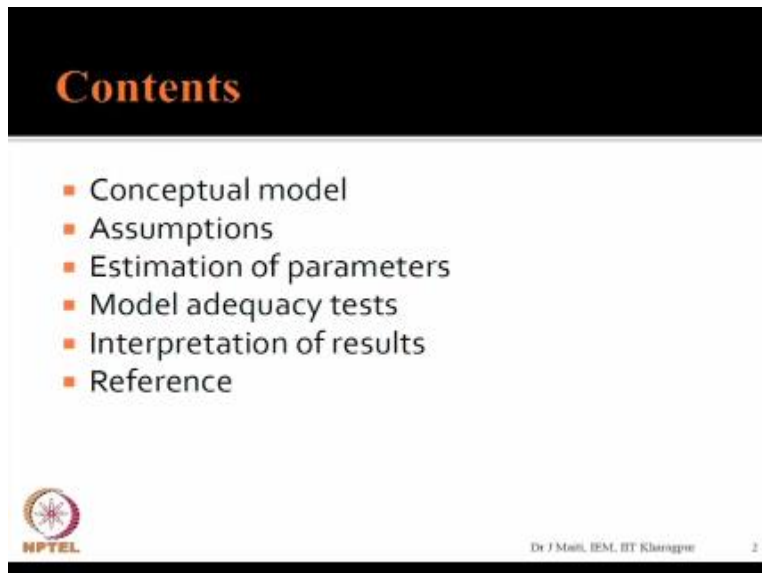
Simultaneously if you want to test the difference.

(Refer Slide Time: 05:33)




You required to use a special technique known as analysis of variance.

(Refer Slide Time: 05:34)



Contents

- Conceptual model
- Assumptions
- Estimation of parameters
- Model adequacy tests
- Interpretation of results
- Reference

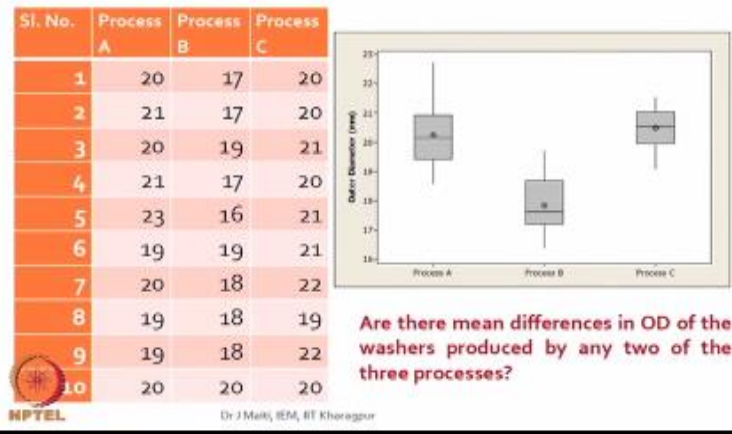


Dr J Mani, BEM, BIT Kharagpur 2

So today we will see in terms of ANOVA for ANOVA these are the contents conceptual model, assumptions, estimation of parameters, model adequacy tests, interpretation of results, followed by references. What is this conceptual model? I will be describing with an example.

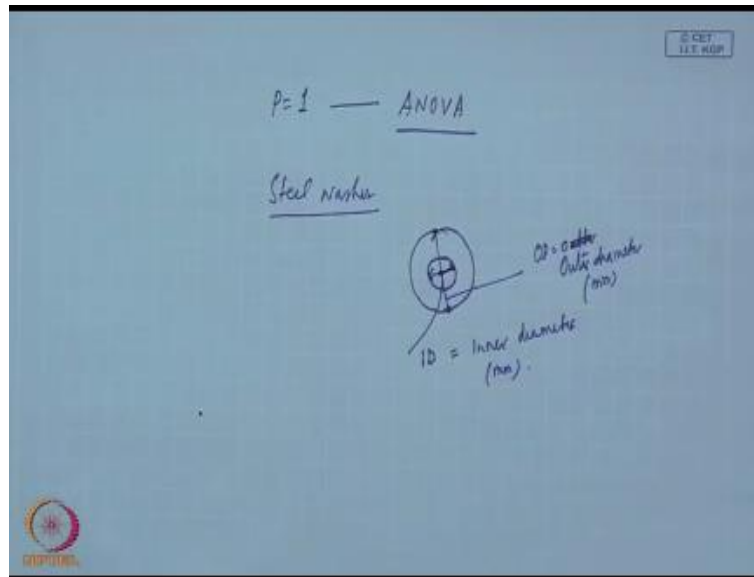
(Refer Slide Time: 05:56)

Conceptual model: An example



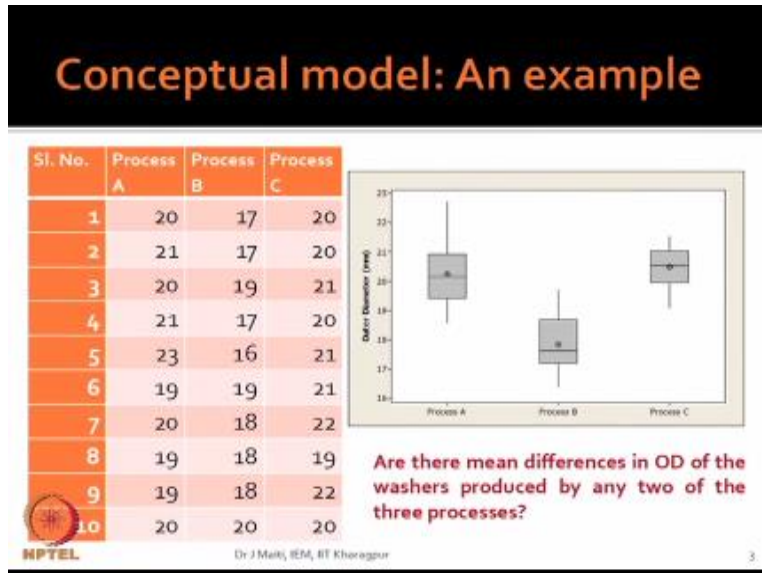
You see this slide here we have three processes A, B and C this three processes are producing let steel washers. So our product is steel washers.

(Refer Slide Time: 06:16)



Steel washer is the product which can be in 2 dimensional figures. We can see like this, suppose this is the case and then this part or otherwise this inner circle part if we consider this is ID inner diameter and if you consider the outer circle this will be OD, which is outer diameter outer. Let both are measured in millimeter is this case.

(Refer Slide Time: 07:04)



Here these values are in millimeter for outer diameter, the problem is that you are producing through three different processes but you are producing same thing that is steel washer and outer diameter is the quality variable and which is specified by the customers, okay. Now as a manufacturer of this steel washers what you require to know, you require to know that whether all the processes are producing at the same level from quality point of view same level or not.

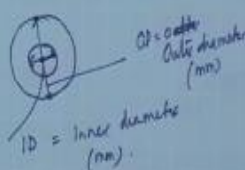
And the manufacturer is interested to test the differences in means of outer diameter for the three processes. So our objective here is are there mean differences in outer diameter of the washers produced by any two of the three processes, okay. Under such situation the most widely used graphical plot is box plot, okay.

(Refer Slide Time: 08:24)

$P=1$ — ANOVA


Steel marker

Box Plot

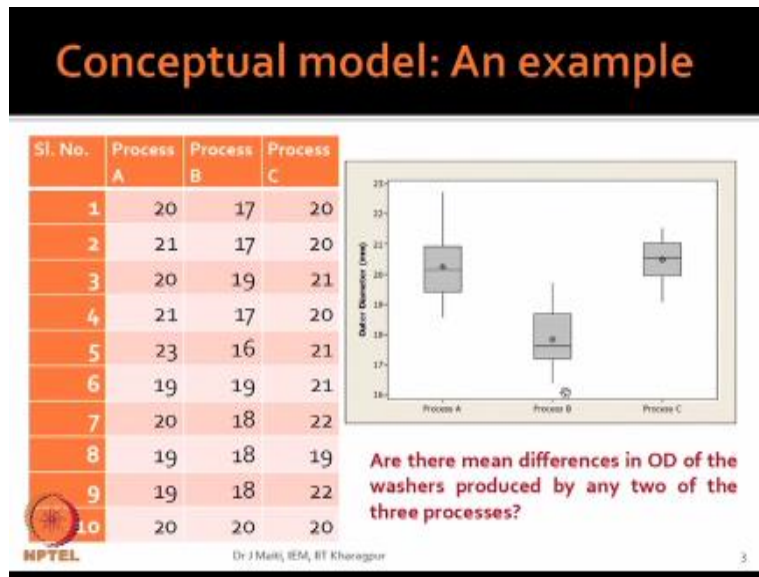


OD = outer Outer diameter (mm)

ID = inner diameter (mm)



(Refer Slide Time: 08:34)



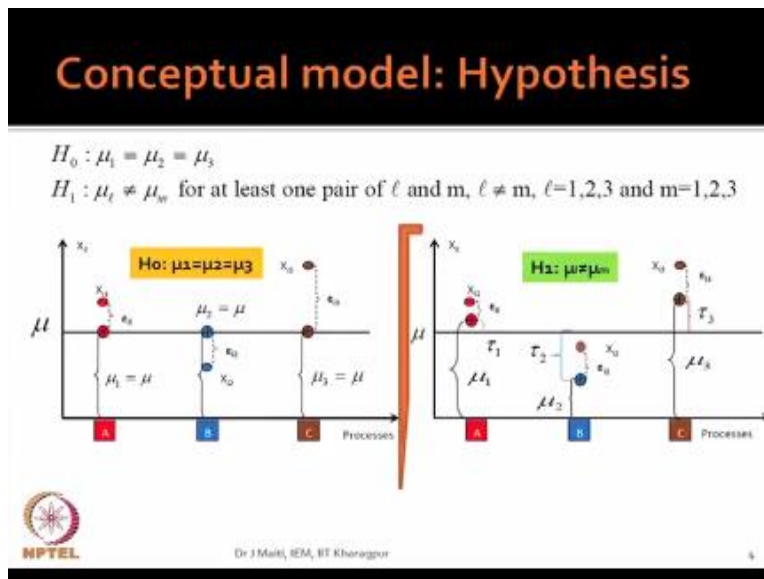
Now see what is box plot this figure shows box plot this box plot is for process A, this is for process B and process C. In the box plot there is one box and two whiskers this one is box the upper portion this straight line is right whisker and this is left whisker, this box is determined by inter quarterly range IQR and the horizontal middle line here this is the median position from the data and the circle plus one that is the mean position, okay.

So if we see the three boxes here for process B and process A as well as process C what we will find out that the box hardly differs in terms of variability except process C but it is to be proven that whether this variability is different from A or B or not apparently it is different, okay. But if you see the mean point the circle plus point this one is quite different apparently it shows that there is mean difference.

You see this is the mean point for process A this is mean point for process B this one is mean point for process C and if I see the value it is 20 point something this one is around 18 and this one is again 20 point something. So even if there is no difference between maybe process A means versus process C means but there is chance that both process A and C means are different from process B.

Which is basically from the graphical plot it resembles like this, but what you want you want a quantitative explanation of this and If we are able to do this then we are able to answer this question are there mean differences in OD of the washers produced by this.

(Refer Slide Time: 10:56)



Then we will use that I told that ANOVA, ANOVA you just two hypothesis, one is null hypothesis and alternative hypothesis what null hypothesis says, null hypothesis says that there is no mean difference as you have taken 3 processes A, B and C. So we are saying that process A mean as μ_1 process B as μ_2 and process C as μ_3 and there is no difference means $\mu_1 = \mu_2 = \mu_3$ and our alternative hypothesis is atleast one of the pairs.

That means either μ_1, μ_2 pair or μ_1, μ_3 pair or μ_2, μ_3 pair is different that means are different. So pictorially the left hand side figure is your null hypothesis and here the right hand side is your alternative hypothesis, if your situation is as the left hand side figure you say what happened here there is one parameter μ please see the left hand side this is μ here this μ is the grand mean.

(Refer Slide Time: 12:22)

$\mu = \text{Grand mean}$

i	Process A	Process B	Process C
1	✓		
2	✓		
⋮			
10			

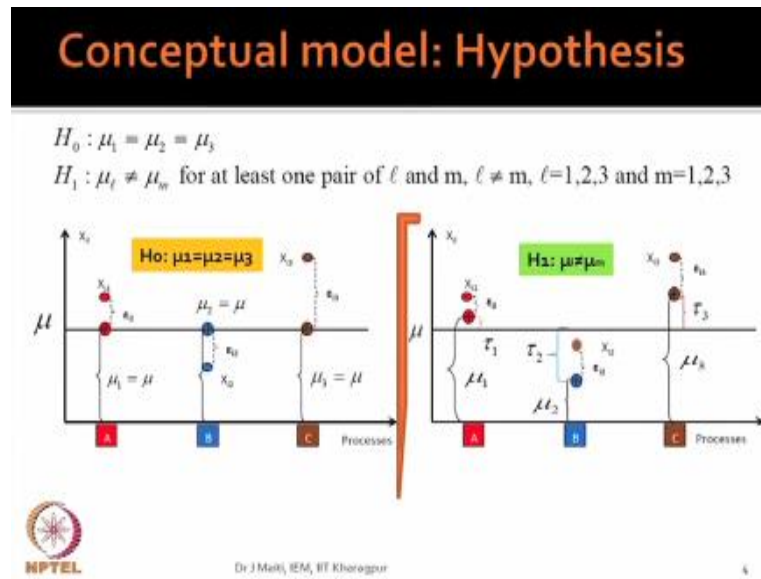
$N = 30$

$\frac{\text{Sum} ()}{30}$

So there is one parameter called μ , which is known as grand mean by grand mean what do I mean, by grand mean we mean that that you have three processes here in this case. So you have produced several items you measure the outer diameter in this case irrespective of the process you take everything as a whole and compute the means, okay, this is what is grand mean the totality.

For example, in this case we have seen that you will collect certain amount of data and we have collected ten data points sample size is 10 for process A process B and process C ten each. So what will happen ultimately total data point if I say N that is 30 and you are considering every data points all sum so sum of all the data points divided by 30 that is your grand mean. If you collect 30 data parts 10 from each of the processes that will become your grand mean, okay.

(Refer Slide Time: 13:45)



Now there is one another mean which is known as process mean.

(Refer Slide Time: 13:51)

$M = \text{Grand mean}$

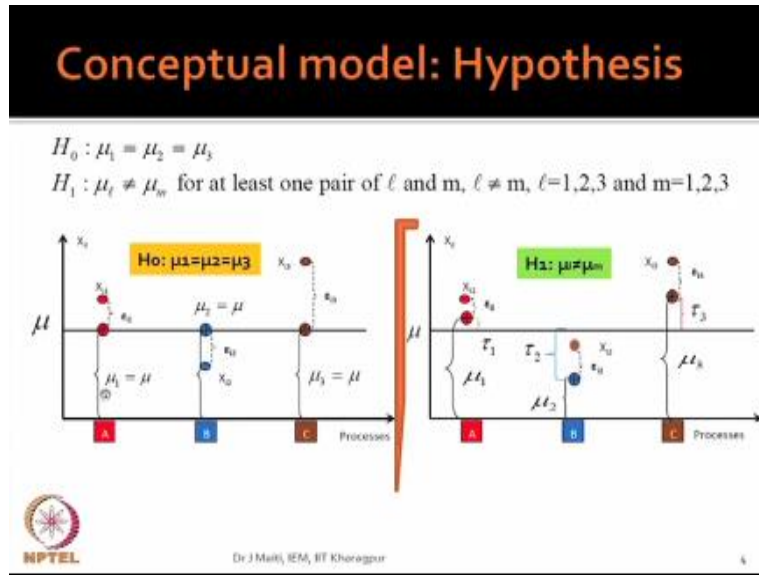
c	Process A	Process B	Process C
1	✓		
2	-		
.			
.			
10			

$N = 30$

$\frac{\text{sum}(x)}{30}$

Or in this case we will say the population mean. Because we have assumed that each of the processes are different as if each of the processes are presenting a population, process A is one population, process B is another process B is the other one.

(Refer Slide Time: 14:13)



μ_1 is the mean of process A, μ_2 is the mean of process B and μ_3 is the mean of process C. When null hypothesis is true I mean $\mu_1 = \mu_2 = \mu_3$ then it will be equal to μ , okay. Because we are collecting some here ten samples from each of the population and when null hypothesis is true, what will happen?

(Refer Slide Time: 14:47)

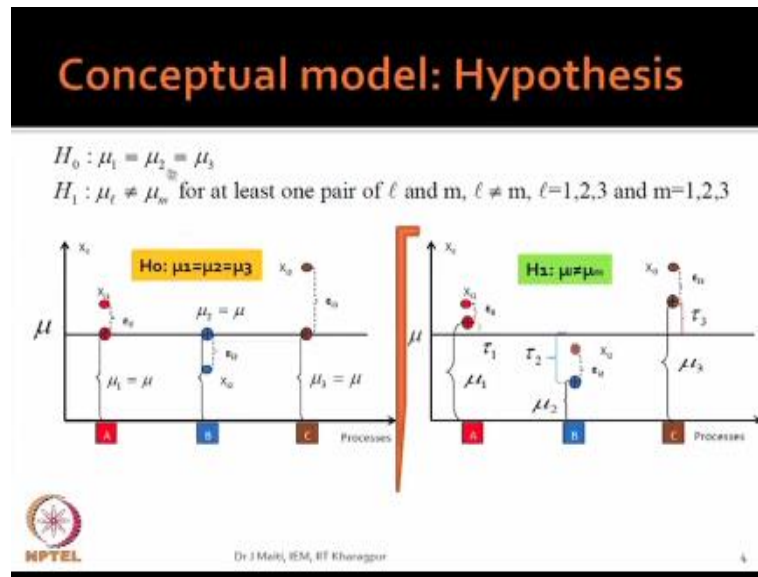
Handwritten table on a blue background. The table is a 2x3 factorial design matrix. The columns are labeled 'Process A', 'Process B', and 'Process C'. The rows are labeled '1' and '2'. A grand mean symbol μ is drawn above the table, with arrows pointing to each column. A small box at the bottom right contains 'N=30' and a sum calculation.

	Process A	Process B	Process C
1	✓		
2	-		

$N=30$ $\frac{\text{Sum}(\quad)}{30}$

You compute the process A mean process B mean process C mean you will also the grand mean we will be finding that they are equal.

(Refer Slide Time: 14:55)



Because it is obvious because what is the group mean calculation do you know that.

(Refer Slide Time: 15:04)

$M = \text{Grand mean}$

	Arabic	Pr/M B	Pr/M C
1	✓		
2	-		
...			
10	✓		

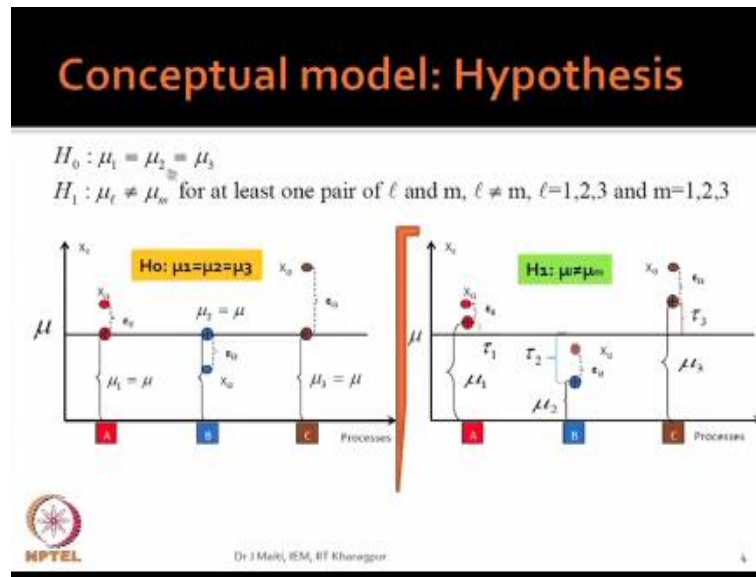
$$\bar{x} = \frac{n_1 f_1 + n_2 f_2 + n_3 f_3}{n_1 + n_2 + n_3}$$

$N = 30$

Sum ()
30

I think all of you know $n_1 f_1 + n_2 f_2 + n_3 f_3 / n_1 + n_2 + n_3$ that is the group mean if you calculate group mean this is the formula, where for the first case frequency is F_1 second case frequency is this is basically $f_1, n_1, f_2, n_2, f_3, n_3$ and n_1 is the number of cases.

(Refer Slide Time: 15:35)



From the that is the I think we have used μ here.

(Refer Slide Time: 15:37)

The image shows handwritten notes on a blue background. At the top left, it says "M = Grand mean". Below this is a table with columns labeled "A", "B", and "C". The rows are labeled "1", "2", and "10". There are checkmarks in the "A" column for rows 1 and 2. To the right of the table is a diagram of a circle with points μ_1, μ_2, μ_3 and a horizontal line representing the grand mean μ . Below the diagram is the formula for the grand mean:
$$\bar{M} = \frac{n_1\mu_1 + n_2\mu_2 + n_3\mu_3}{n_1 + n_2 + n_3} = \mu \left(\frac{n_1 + n_2 + n_3}{n_1 + n_2 + n_3} \right)$$
 Below the formula is a small table with "N = 30" and "Sum ()" above "30".

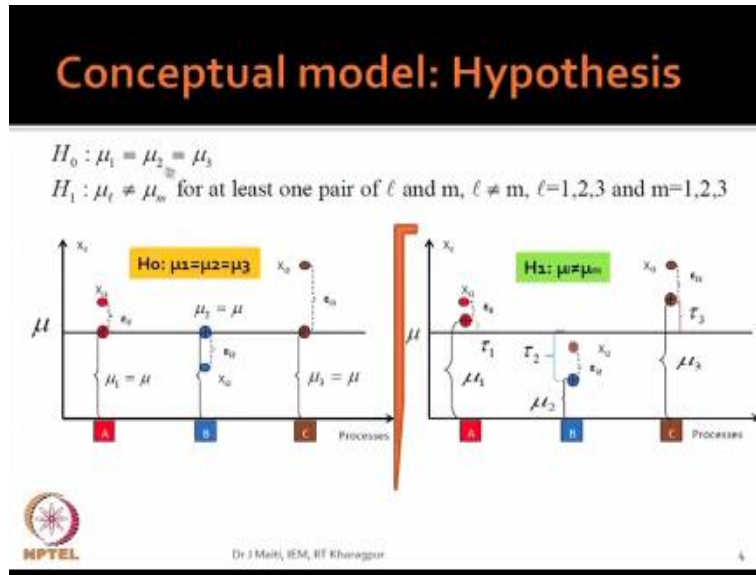
	A	B	C
1	✓		
2	✓		
...			
10			

N = 30

Sum ()
30

$n_1 \times \mu + n_2 \times \mu + n_3 \times \mu$ it is basically μ_1, μ_2 that will be better. So I will write like this $n_1\mu_1 + n_2\mu_2 + n_3\mu_3 / n_1+n_2+n_3$, okay. So this is the general expression we will not consider this one here. So if μ_1, μ_2 and μ_3 are same then what is happening here then it is $\mu n_1 + n_2 + n_3 / n_1+n_2+n_3$ this will be cancelled out this is $\mu = \mu$, so grand mean will be like this.

(Refer Slide Time: 16:26)



Now let us see that what are the other parameters available here, you see the left hand side again there is another parameter which is ϵ_{i1} this one, this is the error quantity error in the sense if you take any observation x_{i1} for the first process from the first process then if we consider that that $\mu_1 = \mu$ then μ_1 will come here the difference between the two that the observed value minus the mean value that is ϵ_{i1} .

Later on we will see this is basically error quantity. So similarly for μ_2 also ϵ_{i2} for μ_3 that is process C case also ϵ_{i3} you are getting. Now come to the right hand side in the right hand side what is the figure resembles, that the μ_1, μ_2 and μ_3 are not same. So they are not coinciding with μ and as a result what will happen you will be getting one more parameter here that is known as τ you see this is τ this one is also τ , this is τ .

So τ_1 is this value what is this, this is the process mean minus grand mean and τ_2 is again the process mean minus grand mean for process B τ_3 is process mean minus grand mean for process C. So one more parameter is coming into picture, so what we will do using all those things?

(Refer Slide Time: 18:04)

Conceptual model: parameters

$$X_{i\ell} = \mu + \tau_{\ell} + \epsilon_{i\ell}$$

$$\tau_{\ell} = \mu_{\ell} - \mu$$


$$\epsilon_{i\ell} = X_{i\ell} - \mu_{\ell}$$

$$\sum_{\ell=1}^L \tau_{\ell} = 0, \text{ for equal sample size}$$

$$\text{and } \sum_{\ell=1}^L n_{\ell} \tau_{\ell} = 0, \text{ for unequal sample size}$$

$$\infty$$

$$H_0 = \tau_{\ell} = 0, \text{ for all } \ell = 1, 2, \dots, L.$$

$$H_1 = \tau_{\ell} \neq 0, \text{ for at least one } \ell.$$


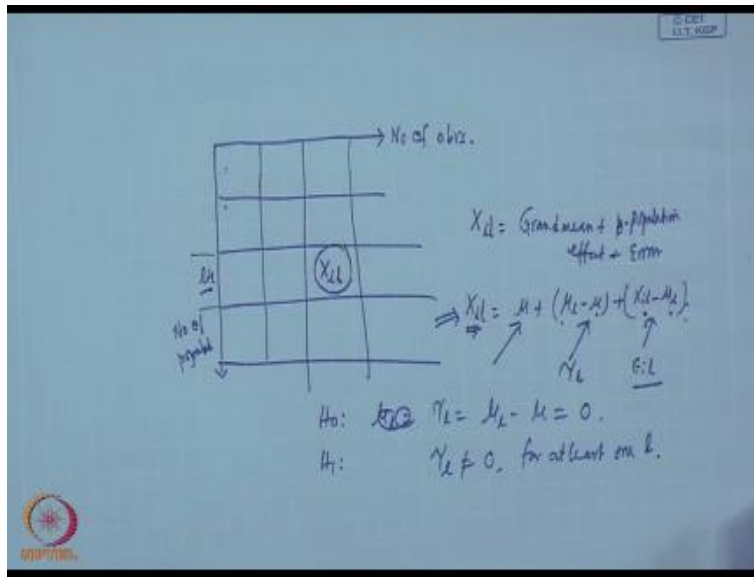
Dr J Malvi, IEM, BT Kharagpur

We will redefine the ANOVA model and that is that any observation can be partitioned into three parts. That is $X_{i\ell}$ is the observation that is the i^{th} observation on the ℓ^{th} population to be collected. μ is the grand mean, τ_{ℓ} is the population effect, and $\epsilon_{i\ell}$ is the error term which is not captured either through grand mean or τ_{ℓ} or combining it is not captured basically, okay. So, what is τ_{ℓ} ? τ_{ℓ} is the population effect.

So, the population effect is τ_{ℓ} which is also known as treatment effect. So, this is $\mu_{\ell} - \mu$ that means the mean of the ℓ^{th} population minus grand mean, and the other one $\epsilon_{i\ell}$ that is the error part is calculated like this. Here, on the right-hand side, there are two conditions: $\sum \tau_{\ell} = 0$ and $\sum n_{\ell} \tau_{\ell} = 0$. For equal sample size, this means when you collect the same sample size from the three populations or L populations, what will happen here?

You will find out that their total effect will be 0 and one is more than the other. Some cases are more than the mean grand mean, some are less than the grand mean, but the total effect will be 0. And if you go for unequal sample size, then your formula will be like this: $\sum_{\ell=1}^L n_{\ell} \tau_{\ell} = 0$, okay. I think I will, let me repeat this.

(Refer Slide Time: 20:02)

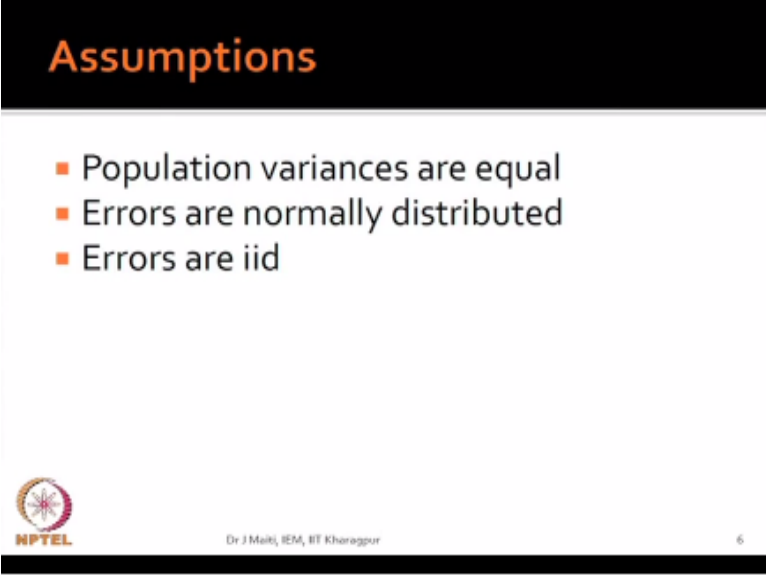


What we are saying in ANOVA case we are saying that you have number of observations in this case, here it is number of population, okay. So you may find out that several populations and also here will be several observations suppose I am saying that one observation is X_{il} that is the general observation that means i^{th} observation on the l^{th} population this is my l population, what we are saying further that X_{il} it can be decomposed as grand mean plus your population effect plus some error quantity.

How do you get it, we will get it like by this manner X_{il} equal to, I can write μ plus let write like this $\mu_l - \mu$. So μ , μ cancelled out μ_l is there but left hand side is X_{il} , so $X_{il} - \mu_l$. See μ and μ will be cancelled out μ_l , μ_l will be cancelled out. So $X_{il} = X_{il}$, so this is what is known as partitioning observation into three components here one is grand mean other one is the τ_l that is the population effect and this one is ϵ_{il} .


So you are partitioning the observation values into three components here and each component is a parameter of ANOVA. That is from conceptual point of view that there are three parameters one is μ another one is τ_l and error component settings are there, okay. If this is the case then your hypothesis can be changed to H_0 that is $\mu_l - \tau_l = \mu_l - \mu = 0$ and H_1 $\tau_l \neq 0$ for at least one l .

(Refer Slide Time: 22:48)



Assumptions

- Population variances are equal
- Errors are normally distributed
- Errors are iid

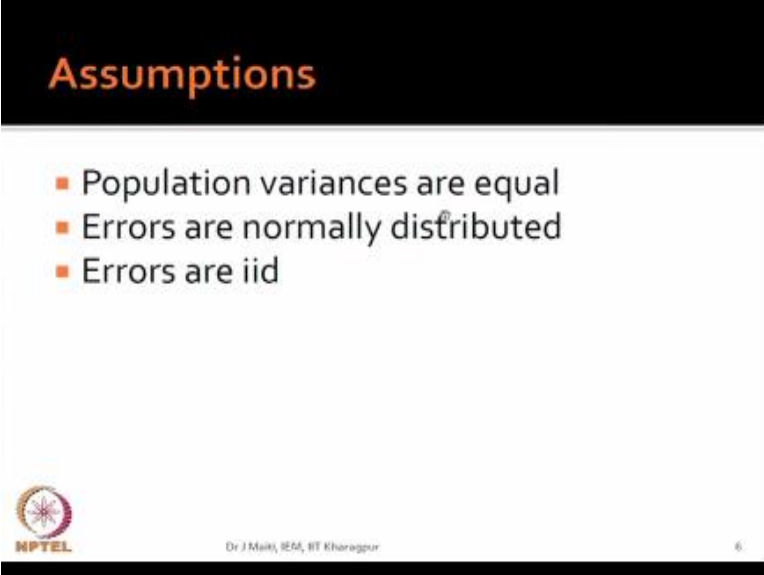
 NPTEL

Dr J Malvi, IEM, BIT Khoragpur

6


Like other multivariate models, ANOVA comes under univariate model like other multivariate model or like any model any statistical model.

(Refer Slide Time: 23:15)



Assumptions

- Population variances are equal
- Errors are normally distributed
- Errors are iid

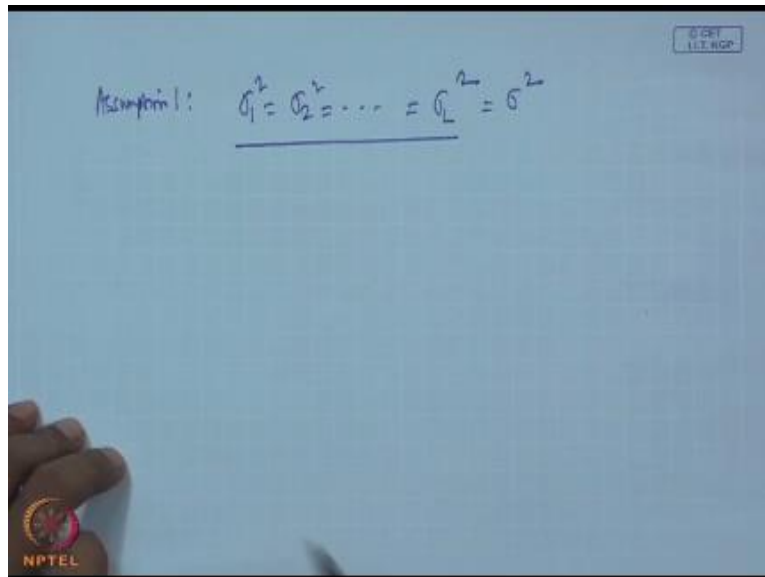
 NPTEL

Dr J Malik, SEM, BIT Kharagpur

6

That ANOVA also requires certain assumptions to be satisfied, these assumptions are populations variances are equal. So you are sampling from 1 populations, the population variances must be equal that means what I mean to say.

(Refer Slide Time: 23:29)



$\sigma_1^2 = \sigma_2^2 = \sigma^2$, okay. This to be satisfied this is our first assumption, assumption one.

(Refer Slide Time: 23:48)

Assumptions

- Population variances are equal
- Errors are normally distributed
- Errors are iid



Dr J Maiti, IISM, IIT Kharagpur

6

Assumption 2 is errors are normally distributed, what is this error you have seen earlier that we have said that.

(Refer Slide Time: 23:55)

Assumption 1: $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2 = \sigma^2$

$\hat{e}_{il} = X_{il} - \hat{\mu}_i$

$\hat{e}_{il} = X_{il} - \hat{\mu}_i$. If you estimate this will be like this, so this error quantity is normally distributed.

(Refer Slide Time: 24:15)

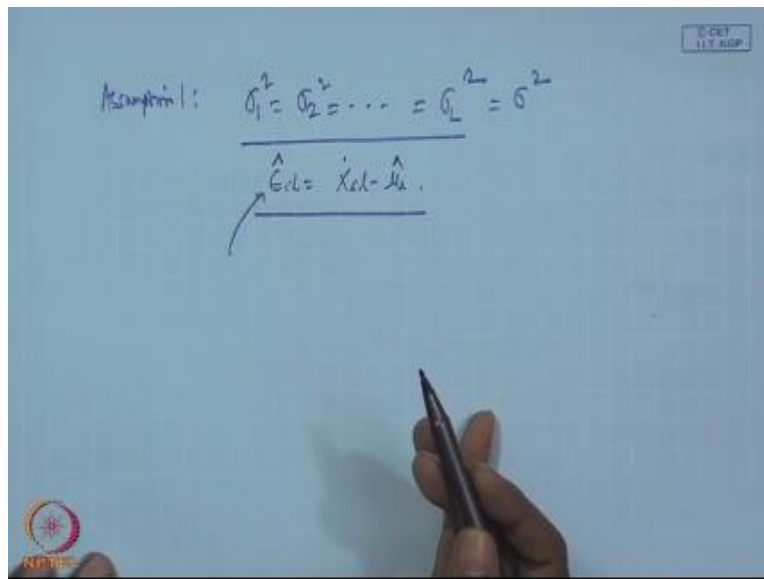
Assumptions

- Population variances are equal
- Errors are normally distributed
- Errors are iid



Errors are IID mean independent and identically distributed.

(Refer Slide Time: 24:19)



Assumption 1: $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$

$\hat{e}_{it} = x_{it} - \hat{\mu}_i$

If I say error one is normally distributed then error N also normally distributed that is identically distributed, independent mean there is no correlation between the errors, okay? So normal distributions as well as IID condition those things you know and how to test the normality we have also discussed earlier in some lectures. Now I will first explain how to test the equality of population variances.

If the one is not satisfied this assumption 1, that is equality of population variances is not satisfied, what will happen? Then you cannot use the traditional ANOVA, okay. It is there some other methods are there but you have to go for different way of doing things.

(Refer Slide Time: 25:18)

Assumptions

- Population variances are equal
- Errors are normally distributed
- Errors are iid




But here this assumption is vital for us.

(Refer Slide Time: 25:22)

Test of equality of population variances

- Bartlett's test

Hypothesis	$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2$ $H_1 : \sigma_i^2 \neq \sigma_m^2, \text{ for at least pair of } (i, m).$
Statistic	$\chi_0^2 = 2.3026 \frac{q}{c}$ $q = (N - L) \log_{10} S_p^2 - \sum_{i=1}^L (n_i - 1) \log_{10} S_i^2$ $c = 1 + \frac{1}{3(L-1)} \left(\sum_{i=1}^L (n_i - 1)^{-1} - (N - L)^{-1} \right), \quad S_p^2 = \frac{\sum_{i=1}^L (n_i - 1) S_i^2}{N - L}$
Decision	Reject H_0 when $\chi_0^2 > \chi_{\alpha, L-1}^2$.



Dr J Malvi, SEM, BIT Kharagpur

So for equality of population variances you will use Bartlett test, okay. Now see the format what we have written here hypothesis statistic and decision and that we have seen earlier also, okay. Now what is our hypothesis? That there is no differences in the population variances and alternative hypothesis is at least one pair of variances are different, okay. Now here Bartlett proposed one statistic, which is $2.3026 \frac{q}{c}$.

This is a statistic where q is $(N - L) \log_{10} S_p^2 - \sum_{i=1}^L (n_i - 1) \log_{10} S_i^2$. I hope you can recollect what is S_p^2 and what is S_i^2 , okay. This is pooled variance and this one is individual population variance. So you require to compute first $S_1^2, S_2^2, S_3^2, S_4^2$ like the all individual population variance values, then also you require to compute the pooled one this pooled one is given in this equation.

$S_p^2 = \frac{\sum_{i=1}^L (n_i - 1) S_i^2}{N - L}$. I think that we have already seen in last class, what we are saying.

(Refer Slide Time: 27:27)

Assumption: $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2 = \sigma^2$

$$\hat{\sigma}_{i,l} = \frac{X_{i,l} - \mu_{i,l}}{\sigma}$$
$$S_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + \dots + (n_L - 1)s_L^2}{n_1 + n_2 + \dots + n_L - L}$$
$$= \frac{\sum_{l=1}^L (n_l - 1)s_l^2}{\sum_{l=1}^L n_l - L}$$


That S_p^2 this is $(N_1 - 1)S_1^2 (N_2 - 1)S_2^2$ like this $N_n - 1S_n^2 / N_1 + N_2 + \dots + N_n$ minus the how many populations are there that is L . So, this quantity is nothing but $N_1 - 1S_1^2 = 1$ to L divided by sum total of $N_1 - 1$ where l equal to L . So this is your S_p^2 you have developed earlier also.

(Refer Slide Time: 28:22)

Test of equality of population variances

- Bartlett's test

Hypothesis	$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2$ $H_1 : \sigma_i^2 \neq \sigma_m^2, \text{ for at least pair of } (i, m).$
Statistic	$\chi_0^2 = 2.3026 \frac{q}{c}$ $q = (N-L) \log_{10} S_p^2 - \sum_{i=1}^L (n_i - 1) \log_{10} S_i^2$ $c = 1 + \frac{1}{3(L-1)} \left(\sum_{i=1}^L (n_i - 1)^{-1} - (N-L)^{-1} \right), \quad S_p^2 = \frac{\sum_{i=1}^L (n_i - 1) S_i^2}{N-L}$
Decision	Reject H_0 when $\chi_0^2 > \chi_{\alpha, L-1}^2$.


Dr J Mall, SEM, BT Khoragpur
7

If you know this then you are in a position to calculate q what is capital N here? Which one, N – 1 this capital N is nothing but this one.

(Refer Slide Time: 28:45)

Assumption 1: $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2 = \sigma^2$

$\hat{\sigma}_k^2 = \frac{1}{n_k} \sum_{i=1}^{n_k} (x_{ki} - \bar{x}_k)^2$

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + \dots + (n_L - 1)s_L^2}{n_1 + n_2 + \dots + n_L - L}$$

$$= \frac{\sum_{k=1}^L (n_k - 1) s_k^2}{\sum_{k=1}^L n_k - L}$$

$$N = \sum_{k=1}^L n_k$$


Capital N is $1 = 1$ to L that means the total sample size from all the populations. Although, it is total sample size is not that meaningful mean from every population. You collected certain sample of size some sizes and that totality we are talking about.

(Refer Slide Time: 29:17)

Test of equality of population variances

- Bartlett's test

Hypothesis	$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2$ $H_1 : \sigma_i^2 \neq \sigma_m^2, \text{ for at least pair of } (i, m).$
Statistic	$\chi_0^2 = 2.3026 \frac{q}{c}$ $q = (N-L) \log_{10} S_p^2 - \sum_{i=1}^L (n_i - 1) \log_{10} S_i^2$ $c = 1 + \frac{1}{3(L-1)} \left(\sum_{i=1}^L (n_i - 1)^{-1} - (N-L)^{-1} \right), \quad S_p^2 = \frac{\sum_{i=1}^L (n_i - 1) S_i^2}{N-L}$
Decision	Reject H_0 when $\chi_0^2 > \chi_{\alpha, L-1}^2$.



Dr J Mall, SEM, BT Khoragpur

7


Then you have to compute C is $1 + 1/3 \sum_{i=1}^L (n_i - 1)^{-1} - (N - L)^{-1}$ this is your C value, okay. So please go through this C there are basically C is another quantity which is basically $\sum_{i=1}^L (n_i - 1) S_i^2$ and N values are defecting this C and this C will be used as the divisor. So then chi square is $2.3026q/c$. So given data set you have to calculate q you have calculate C then you find out the chi square value.

Chi square is $2.3026 q/c$ and what does it mean, when you say chi square because this quantity follows chi square distribution and what will be the degrees of freedom? That you see it is clearly written here in the decision case that reject H_0 when the chi square computed is greater than chi square $1-\alpha$ for a particular alpha value. So what do you mean then that 2.3026 follows chi square distribution with $L-1$ degrees of freedom. I think the similar thing you have seen earlier also.

(Refer Slide Time: 30:54)

Bartlett's test

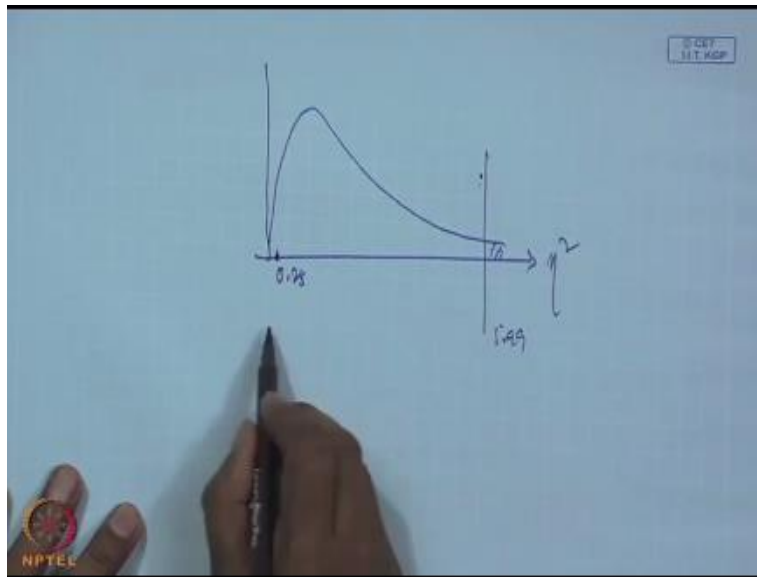
Process A	Process B	Process C		
20	17	20	s1-sq	1.51
21	17	20	s2-sq	1.43
20	19	21	s3-sq	0.93
21	17	20		
23	16	21	Sp-sq	1.29
19	19	21	q	0.26
20	18	22	c	1.05
19	18	19	Chi-sq	0.25
19	18	22	Chi(2, 0.05)	5.99
20	20	20	Failed to reject Ho	

 Dr. J. Mehta, IIM, IT Kharagpur 8

I have been repeating, see the same problem. So we have 30 observations 10 from each of the processes and what we have done here we have calculated the variance for the first process that is 1.5 one variance for the second process that is 1.43 variance for the third process that is 0.93. and then you have computed chi square by that this one is chi square. So you have computed first pool variance then q then C then chi square, which is 0.25.

Now we all know that chi square two with two degrees of freedom and alpha equal to 0.05. This value is 5.99 we have seen earlier we have seen this one earlier also. Now the computed value is less than the tabulated value. So what will be your decision you will accept null hypothesis or reject null hypothesis accept null hypothesis you fail to reject null hypothesis.

(Refer Slide Time: 32:09)




So the same principle this will be my chi square suppose this is your chi square distribution. Your alpha tabulated value here this is 5.99 but your value is how much? 0.25, this is very close to 0. So what does it signify?

(Refer Slide Time: 32:33)

Bartlett's test

Process A	Process B	Process C		
20	17	20	s1-sq	1.51
21	17	20	s2-sq	1.43
20	19	21	s3-sq	0.93
21	17	20		
23	16	21	Sp-sq	1.29
19	19	21	q	0.26
20	18	22	c	1.05
19	18	19	Chi-sq	0.25
19	18	22	Chi(2, 0.05)	5.99
20	20	20	Failed to reject Ho	


 Dr J Maiti, IIM, BT Kharagpur 8

It signifies that there is no population variance differences all the population variances are equal same. So we are fit for ANOVA case, okay.

(Refer Slide Time: 32:50)

Decomposition of total sum of squares

Population	$t = 1, 2, \dots, t$	$\bar{x}_t = \frac{1}{n} \sum_{i=1}^n x_{it}$	Partitioning of observations (x_{it})
1	$x_{11}, x_{21}, \dots, x_{i1}, \dots, x_{n1}$	$\bar{x}_1 = \frac{1}{n} \sum_{i=1}^n x_{i1}$	$x_{it} = \bar{x} + (\bar{x}_t - \bar{x}) + (x_{it} - \bar{x}_t)$
2	$x_{12}, x_{22}, \dots, x_{i2}, \dots, x_{n2}$	$\bar{x}_2 = \frac{1}{n} \sum_{i=1}^n x_{i2}$	
⋮		⋮	
t	$x_{1t}, x_{2t}, \dots, x_{it}, \dots, x_{nt}$	$\bar{x}_t = \frac{1}{n} \sum_{i=1}^n x_{it}$	
L	$x_{1L}, x_{2L}, \dots, x_{iL}, \dots, x_{nL}$	$\bar{x}_L = \frac{1}{n} \sum_{i=1}^n x_{iL}$	
Grand mean		$\bar{x} = \frac{1}{nL} \sum_{t=1}^L \sum_{i=1}^n x_{it}$	



Dr J Mani, BEM, BIT Kharagpur

9

In ANOVA it will not go in the same manner like Hotelling t- square or your other one that t-test the way we have developed in ANOVA it is a different ball game altogether, because here the primary concern is partitioning the observation.

(Refer Slide Time: 33:11)

$$\begin{aligned}
 x_{il} - \bar{x} &= \bar{x}_l - \bar{x} + (x_{il} - \bar{x}_l) \\
 \Rightarrow x_{il} - \bar{x} &= (\bar{x}_l - \bar{x}) + (x_{il} - \bar{x}_l) \\
 \sum_{i=1}^n (x_{il} - \bar{x})^2 &= \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + 2 \sum_{i=1}^n (\bar{x}_l - \bar{x})(x_{il} - \bar{x}_l) + \sum_{i=1}^n (x_{il} - \bar{x}_l)^2 \\
 &= \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{i=1}^n (x_{il} - \bar{x}_l)^2 \\
 \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x})^2 &= \sum_{l=1}^L \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2 \\
 &= \sum_{l=1}^L n (\bar{x}_l - \bar{x})^2 + \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2
 \end{aligned}$$

We have seen earlier, I told you that from the population point of view when I am planning to collect some data let X_{il} this is the observation to be collected from a population where the grand mean is μ and population mean is μ_l then we have partition the observation like this, you have seen this one see the left hand side and right hand side both are same. Now you have collected data.

That means you have a fixed value X_{il} . So you have to go by sample grand mean plus your population mean minus again grand mean plus X_{il} that value minus population sample means not population sample mean. So it is now you do little manipulation $X_{il} - \bar{X} = \bar{X}_l - \bar{X} + X_{il} - \bar{X}_l$ square the term, If you square it $X_{il} - \bar{X}^2 = \bar{X}_l - \bar{X}^2 + 2 \times \bar{X}_l - \bar{X} \times X_{il} - \bar{X}_l + X_{il} - \bar{X}_l^2$ fine, you take summation over I , okay.

So $i = 1$ to N we are assuming that equal size samples are collected from the l populations. So again $i = 1$ to N then this side will be $i = 1$ to N this side will be $i = 1$ to N . Now what will happen to this middle one? The middle one will become zero getting me because of this quantity. So when you write down sum total $i = 1$ to N some X_{il} this will be $N \bar{X}_l - \bar{X}_l$ for N times that will be also $N \bar{X}_l$.

So this value will become zero so your resultant quantity will be $\sum_{i=1}^N \bar{X}_L - \bar{X}^2 + \sum_{i=1}^N X_{i1} - \bar{X}_L^2$, okay. Now see that this one this is in the left hand side some quantity right hand side also, me quantity two squares is given here, okay. This is for a particular population l population yes or no? Because we have taken X_{i1} every where l is there but how many l populations are there?


Capital l populations are there. So that means your observation is spreaded overall l populations. So if take one more summation what will happen? What will you do? Then you will write down l equal to 1 to capital l sum total of i equal to 1 to $N \sum_{i=1}^N X_{il} - \bar{X}_L^2 = \sum_{l=1}^L \sum_{i=1}^N X_{il} - \bar{X}_L^2$ plus sum total of $\sum_{l=1}^L \sum_{i=1}^N X_{il} - \bar{X}_L^2$ this is also square. Now in this case is there any \bar{X}_L and \bar{X} there is no i component.

So that mean the same component for N times, so what you can write like this $\sum_{l=1}^L \sum_{i=1}^N X_{il} - \bar{X}_L^2$ plus sum total $\sum_{l=1}^L \sum_{i=1}^N X_{il} - \bar{X}_L^2$, okay.

(Refer Slide Time: 38:57)

Decomposition of total sum of squares

Population	$i = 1, 2, \dots, n$	$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$	Partitioning of observations (x_{ij})
1	$x_{11}, x_{21}, \dots, x_{n1}$	$\bar{x}_1 = \frac{1}{n} \sum_{j=1}^n x_{1j}$	$x_{ij} = \bar{x} + (\bar{x}_i - \bar{x}) + (x_{ij} - \bar{x}_i)$
2	$x_{12}, x_{22}, \dots, x_{n2}$	$\bar{x}_2 = \frac{1}{n} \sum_{j=1}^n x_{2j}$	
⋮	⋮	⋮	
l	$x_{1l}, x_{2l}, \dots, x_{nl}$	$\bar{x}_l = \frac{1}{n} \sum_{j=1}^n x_{lj}$	
⋮	⋮	⋮	
L	$x_{1L}, x_{2L}, \dots, x_{nL}$	$\bar{x}_L = \frac{1}{n} \sum_{j=1}^n x_{Lj}$	
Grand mean		$\bar{x} = \frac{1}{nL} \sum_{i=1}^L \sum_{j=1}^n x_{ij}$	


Dr J Mohi, SEM, BT Kharagpur

So what will be the left hand side and the right hand side?

(Refer Slide Time: 39:04)

Handwritten derivation on a whiteboard:

$$\sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x})^2 = \sum_{k=1}^L \sum_{i=1}^n (\bar{x}_k - \bar{x})^2 + \sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x}_k)^2$$

$$\frac{\sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x})^2}{SST} = \sum_{k=1}^L \sum_{i=1}^n (\bar{x}_k - \bar{x})^2 + \sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x}_k)^2$$

Additional notes on the whiteboard include "k-th group" and "L" written near the first equation, and an NPTEL logo in the bottom left corner.

If you see when we calculate the variance what we will do? We will subtract the mean value from all the observed values and then square it and then we divide it by the degrees of freedom here we have not divided it into degrees of freedom but actually what is happening here, this $x_{ki} - \bar{x}$ this was the variability this square. So for all the observation case you have subtracted by the grand central average. So this quantity is known as sum square total, okay.

(Refer Slide Time: 39:55)

Decomposition of total sum of squares

$$x_{i\ell} - \bar{x} = \bar{x}_\ell - \bar{x} + x_{i\ell} - \bar{x}_\ell$$

$$\sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x})^2 = n \sum_{\ell=1}^L (\bar{x}_\ell - \bar{x})^2 + \sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x}_\ell)^2$$

$$SST = SSB + SSE$$

$$N-1 = L-1 + N-L$$



Dr. J. Malik, IIS, IIT Kharagpur

1

You see here that I equal to one to I this $i = 1$ to $|X|$ - \bar{X}^2 this is the sum square total and we have already seen that other one N into I equal one to $|X|$ - \bar{X}^2 this is what is the difference the variability part from grand mean to the individual population means that sum we have taken you see this one what you have done here?

(Refer Slide Time: 40:29)

The whiteboard shows the following derivation:

$$\begin{aligned}
 \sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x})^2 &= \sum_{k=1}^L \sum_{i=1}^n (\bar{x}_k - \bar{x})^2 + \sum_{k=1}^L \sum_{i=1}^n (x_{ki} - \bar{x}_k)^2 \\
 \text{SST} &= \text{SSB} + \text{SSE}
 \end{aligned}$$

The derivation starts with the identity $x_{ki} - \bar{x} = (\bar{x}_k - \bar{x}) + (x_{ki} - \bar{x}_k)$. Squaring both sides and summing over all observations i for each population k yields the decomposition of the total sum of squares (SST) into the sum of squares between populations (SSB) and the sum of squares error (SSE).

\bar{x}_k this is the population for k^{th} population sample average and \bar{x} is the grand sample average the difference between the two and then this observations you have squared. So if there is no difference between the population means we assume that the sample average also will become same and it will be equal to the grand sample average, so then this quantity will become 0. So if there is any variability this is because of the variability between the population.

So that is why this quantity is known as SSB and we have seen earlier also $x_{ki} - \bar{x}$ this is nothing but the error part which is not explained by the which is not explained by SSB. So this is SSE sum square error, now sum square total is equal to sum square between populations plus sum square error. So, what is happening here?

(Refer Slide Time: 41:48)

Decomposition of total sum of squares

$$x_{i\ell} - \bar{x} = \bar{x}_\ell - \bar{x} + x_{i\ell} - \bar{x}_\ell$$

$$\sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x})^2 = n \sum_{\ell=1}^L (\bar{x}_\ell - \bar{x})^2 + \sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x}_\ell)^2$$

$$\text{SST} = \text{SSB} + \text{SSE}$$

$$N-1 = L-1 + N-L$$



Dr J Maiti, IEM, BIT Kharagpur

Then we are basically.

(Refer Slide Time: 41:50)

The image shows handwritten notes on a blue background. At the top, the equation $SST = SSB + SSE$ is written, with an arrow pointing to SSE and the label \leftarrow DOF. Below this, $N-1 =$ is written. To the left, a circle contains $N = nL$, with an arrow pointing to the equation $= \sum_{l=1}^L n_l$. To the right, a diagram shows a data matrix with rows labeled 1, 2, ..., L and columns labeled $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_L$. The NPTEL logo is visible in the bottom left corner.

Dividing the total variability in terms of sum square total equal to variability explained by the population plus variability not explained by the population correct. So this is the known as decomposition of decomposing the total variability into between population variability and error variability, okay. So now we have to see the degree of freedom part degrees of freedom. So when we compute SST what we require basically we require \bar{X} .

So there are total N observations and we have sacrificed one to calculate this one \bar{X} . So degree of freedom for SST is N minus capital N minus one where capital N is N into 1 when you sample from 1 population and sample size is equal. If sample size is unequal then this will be your $l = 1$ to N n_l getting me. So this is for equation for the equal sample size case this is the equation for computing capital N for the unequal sample size case.

Now come to the between sum square now in between sum square how many levels you have if you see the computation what we will be finding out? We will be finding out the population 1 to population l and you will be computing here that bar \bar{X} . So here \bar{X}_1, \bar{X}_2 like \bar{X}_l and then you find out the \bar{X} .

(Refer Slide Time: 44:07)

The image shows a whiteboard with handwritten mathematical derivations. The top part shows the decomposition of the total sum of squares (SST) into between-group and within-group components. The middle part shows the decomposition of the between-group sum of squares (SSB) into a term involving the overall mean and a term involving the group means. The bottom part shows the decomposition of the within-group sum of squares (SSW) into a term involving the overall mean and a term involving the group means. The overall mean is denoted by \bar{X} and the group means by \bar{X}_l . The number of groups is L and the number of observations in each group is n .

$$\sum_{i=1}^n (x_{il} - \bar{x})^2 = \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$
$$\sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x})^2 = \sum_{l=1}^L \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$
$$SST = \sum_{l=1}^L n (\bar{x}_l - \bar{x})^2 + \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$

SSB = $\sum_{l=1}^L n (\bar{x}_l - \bar{x})^2$
SSW = $\sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$

And if you see the computation of SSB you are finding out that this is nothing but $\bar{X}_1 - \bar{X}^2$. Then multiplied by N and that summation, essentially you are using \bar{X}_1 .

(Refer Slide Time: 44:25)

Handwritten notes on a whiteboard:

$$SST = SSB + SSE$$

$$\frac{N-1}{1} = \frac{L-1}{1} + \frac{N-L}{1}$$

Annotations:

- Arrows point from the terms in the second equation to the label "DOF".
- A circle around $N-1$ has an arrow pointing to the equation $N = nL = \sum_{l=1}^L n_l$.
- A circle around $N-L$ has an arrow pointing to the equation $N-1$.

1	\bar{x}_1
2	\bar{x}_2
...	...
L	\bar{x}_L
	\bar{x}

So as there are there are 1 populations and there is relationship also with the individual population mean this sum is grand mean. So one degree you are losing here and as a result you have total 1 data points and one degree is lost by computing \bar{X} grand mean. So plus $N - 1$ because this is the rule that the total variability will also be decomposed into the component degrees of freedom.

So total degrees of freedom is $N - 1$ and your SSB degrees of freedom $l-1$ SSL degrees of freedom $N - 1$. So if you make a sum $l-1 + N-1$ it is nothing but $N - 1$ right hand side $N - 1 + 1 - 1$ sum total will be $N-1$ and which is equal to this keep in mind, this one.

(Refer Slide Time: 45:46)

The image shows a handwritten derivation on a whiteboard. At the top, it states:
$$\Rightarrow x_{il} - \bar{x} = (\bar{x}_l - \bar{x}) + (x_{il} - \bar{x}_l)$$
Below this, the total sum of squares is derived:
$$\sum_{i=1}^n (x_{il} - \bar{x})^2 = \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \underbrace{2 \sum_{i=1}^n (\bar{x}_l - \bar{x})(x_{il} - \bar{x}_l)}_0 + \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$
This simplifies to:
$$= \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$
The next step shows the summation over L groups:
$$\sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x})^2 = \sum_{l=1}^L \sum_{i=1}^n (\bar{x}_l - \bar{x})^2 + \sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2$$
Finally, it identifies the components:
$$\frac{N-1}{SST} = \underbrace{\sum_{l=1}^L n (\bar{x}_l - \bar{x})^2}_{SSB} + \underbrace{\sum_{l=1}^L \sum_{i=1}^n (x_{il} - \bar{x}_l)^2}_{SSE}$$

Yes you are finding out N a, what you are you are using you are using 1 populations. So, you have 1 mean values. Second one what is happening N into 1 this is the total observations. So we are talking about $N = n \times 1$ this is the total observation this side is the total observation is this, but you have already computed \bar{X} . So this one this portion see how many \bar{X}_l you have computed?

Bar L populations. So total is $N - 1$, so there is 1 \bar{X}_l to compute this one you require to know all those \bar{X}_l . So already 1 degrees of freedom you have lost you have capital N degrees of observations 1 is lost.

(Refer Slide Time: 47:06)

Handwritten notes on a blue background showing the decomposition of degrees of freedom. The notes include the following equations and a table:

$$SST = SSB + SSE \leftarrow$$
$$N-1 = \frac{L-1}{L} + \frac{N-L}{L} \leftarrow \text{DOF}$$

Annotations include a circle around $N-1$ with an arrow pointing to $N = nL$, and another circle around $N-L$ with an arrow pointing to $N-1$. Below the equations, a table is drawn with L rows and 2 columns. The first column contains the numbers 1, 2, and L . The second column contains the symbols $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_L$. A horizontal line is drawn below the last row of the table, and a vertical line is drawn to the right of the last column. The overall mean \bar{x} is written below the horizontal line.

1	\bar{x}_1
2	\bar{x}_2
\vdots	\vdots
L	\bar{x}_L
\bar{x}	

So $N - 1$ so this is what is the decomposing the degree of freedom, right.

(Refer Slide Time: 47:17)

Decomposition of total sum of squares

$$x_{i\ell} - \bar{x} = \bar{x}_\ell - \bar{x} + x_{i\ell} - \bar{x}_\ell$$

$$\sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x})^2 = n \sum_{\ell=1}^L (\bar{x}_\ell - \bar{x})^2 + \sum_{\ell=1}^L \sum_{i=1}^n (x_{i\ell} - \bar{x}_\ell)^2$$

$$\text{SST} = \text{SSB} + \text{SSE}$$

$$N-1 = L-1 + N-L$$



Dr J Malvi, BSM, IIT Kharagpur

This total concept whatever we have discussed so far.

(Refer Slide Time: 47:24)

Decomposition of total sum of squares: easier computation

Equal sample size

$$A = \sum_{\ell=1}^L \sum_{i=1}^n x_{i\ell}, \text{ and } A_{\ell} = \sum_{i=1}^n x_{i\ell}$$

$$N = nL$$

Un-equal sample size

$$A_{\ell} = \sum_{i=1}^{n_{\ell}} x_{i\ell}, A = \sum_{\ell=1}^L A_{\ell}$$

$$N = \sum_{\ell=1}^L n_{\ell}$$

$$SST = \sum_{\ell=1}^L \sum_{i=1}^n x_{i\ell}^2 - \frac{A^2}{N}$$

$$SSB = \sum_{\ell=1}^L \frac{A_{\ell}^2}{n} - \frac{A^2}{N}$$

$$SSE = SST - SSB$$



Dr J Malvi, IEM, IIT Kharagpur

Can be seen in a table.

(Refer Slide Time: 47:27)

Hypothesis testing

Sources of variation	Sums square (SS)	Degrees of freedom	Mean square (MS)	F	Reject H_0
Population (treatment)	SSB	L-1	$MSB = \frac{SSB}{L-1}$	$MSE = \frac{SSE}{nL-L}$	$F > F_{L-1, nL-L}^{(\alpha)}$
Error (random component)	SSE	nL-L	$F = \frac{MSB}{MSE}$		
Total	SST	nL-1			



Which is known as ANOVA table, whenever you go through any textbook on ANOVA or any software you use for ANOVA what will happen you will be finding out that one table will be formulated and this table is known as ANOVA table which is very popular table, later on even in regression in other cases also we will be using ANOVA table. This ANOVA table for the problem we have considered now.

They there are few items one is sources of variation then the sum square what we have calculated SST and SSB all those things degrees of freedom also. We have seen then you have to compute mean square you have to compute f value then based on f value we will be accepting or rejecting the null hypothesis, okay. So we have considered 1 population no other factors we have considered in this present case.

So sources of variation is the different populations and apart from this you cannot nullify the random effects. So random effect is coming under error, so another source of variability is error and this population and error this two sources are making the total variability. So as a result your table looks like this population error and total then SSB SSE and SST you have already seen that,

what is the degrees of freedom case 1- 1for SSB for SSE $N - 1$ and for SST $N - 1$. This N multiplied by capital 1 this is the capital N , what I have discussed earlier that is capital N .

(Refer Slide Time: 49:39)

Handwritten notes on a whiteboard:

$$SST = SSB + SSE$$

$$N-1 = L-1 + N-L$$

Annotations: $N = nL$ and $\sum_{i=1}^L n_i$ are circled. The term $N-L$ in the second equation is also circled and labeled as $N-1$. The label \leftarrow DOF is placed to the right of the second equation.

	\bar{x}_1
1	\bar{x}_1
2	\bar{x}_2
\vdots	\vdots
\bar{x}_L	\bar{x}_L
\bar{x}	\bar{x}

This capital N this is N into 1, but if you take unequal sample size this is the different one. Now what is MSB?

(Refer Slide Time: 49:54)

Hypothesis testing

Sources of variation	Sums square (SS)	Degrees of freedom	Mean square (MS)	F	Reject H_0
Population (treatment)	SSB	L-1	$MSB = \frac{SSB}{L-1}$	$MSE = \frac{SSE}{nL-L}$	$F > F_{L-1, nL-L}^{(\alpha)}$
Error (random component)	SSE	nL-L	$F = \frac{MSB}{MSE}$		
Total	SST	nL-1			



MSB is mean square between populations this is nothing but sum square between population divided by its degrees of freedom, okay. Then what will be your MSE? MSE is SSE by its degrees of freedom. So this MSE will come here you write like this MSB is this MSE here under error you write under error MSE. So, I am writing again what will happen here?

(Refer Slide Time: 50:35)

Source	SS	DF	MS	F
Population	\underline{SSB}	$L-1$	$\frac{SSB}{L-1} = MSB$	$\frac{MSB}{MSE} = F_{(L-1, N-L)}$
Error	\underline{SSE}	$N-L$	$\frac{SSE}{N-L} = MSE$	
Total	\underline{SST}	$N-1$		

$N = nL$

$F > F_{L-1, N-L}$

Population then error then total these are the sources of variation sources we are computing sum square that is sum square between populations sum square error and sum square total what is our degree of freedom? Degree of freedom it will be $l - 1$ and this will be your $N - 1$ and this will be $N - 1$ where N is nothing but capital $l \times n$ and then we are talking about MS mean square which is $SSB / N - 1 - 1$ and this be MSE, $SSE / N - 1$.

Then what is this, $SSB / l - 1$ is MSB, $SSE / N - 1$ is MSE, okay. Now you are finding out a statistics called f which is MSB / MSE , why it is f distributed? Ratio of 2 chi square variable because SSB is the square this is sum square and sum square this is f distributed, okay and what will be the degrees of what will be the degrees of numerator and denominator degrees of freedom for f here?

Definitely what will be there in numerator SSB by degrees of freedom $l - 1$. So, we will be comparing this with $l - 1$ and what is the denominator degrees of freedom $N - 1$ and you will find out certain value of α and if your computed f that f computed here is greater than $f_{l - 1, N - 1}$ α you will reject the null hypothesis, okay. So once you reject null hypothesis you will accept the alternative hypothesis that is the issue.

(Refer Slide Time: 53:07)

Hypothesis testing

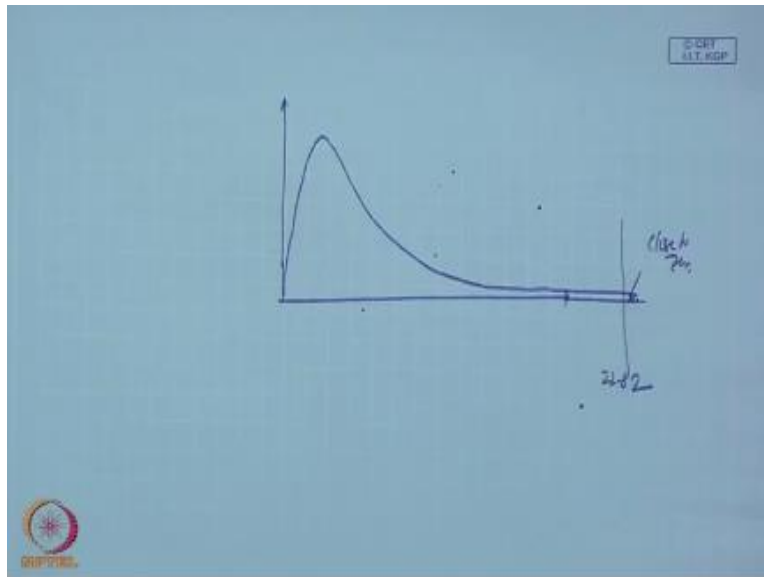
Sources of variation	Sums square (SS)	Degrees of freedom	Mean square (MS)	F	Reject H_0
Population (treatment)	43.415	2	21.708	22.82	$p = 0.000$
Error (random component)	25.686	27	0.951		
Total	69.101	29			



The problem in the same problem we have done the hypothesis testing and population that variability sum square between population is 43.415. Error SSE is 25.686 and total is 69.101 degrees of freedom because there are three processes. So 3-1 that is 2, 3-1 is 2 and there are thirty total observations that minus one is 29 and you the difference between 29 and 2 is 27 that will be definitely the error's degree of freedom.

Then mean square is $43.415 / 2$ that is 21.708 and for error that is again $25.686 / 27$ which is 0.951 you see the mean square value is very less it is, so less because there is effect. So then f is MSE / MSB , MSB / MSE . This is 21.708 is MSB and 0.951 is MSE and that ratio is 22.82 and which is far away from 0.

(Refer Slide Time: 54:46)




If you see the f table you will find the probability is almost close to 0. That mean if I say like this is my f distribution that mean you are somewhere here this one which is your 22.82. This one is close to 0.

(Refer Slide Time: 55:21)

Hypothesis testing

Sources of variation	Sums square (SS)	Degrees of freedom	Mean square (MS)	F	Reject H_0
Population (treatment)	43.415	2	21.708	22.82	p = 0.000
Error (random component)	25.686	27	0.951		
Total	69.101	29			



Dr. J. Malik, IEM, IIT Kharagpur

32

So there is difference this is what is tested through tested through ANOVA first level test is this there is difference in population our this one the means population means that is what you are testing, okay. Now here when you test this f test when you are doing here your what is your null hypothesis? Null hypothesis no population no differences in population means correct. What is alternative hypothesis? At least one pair is different.

So when you complete this one and if you find that you are rejecting H_0 . It simply says that there is population difference in terms of population means, but you do not know which they are, so we have to know which pair are different, okay. So next class I will explain this which pair is different and other things.

NPTEL Video Recording Team

NPTEL Web Editing Team

Technical Superintendents

Computer Technicians

A IIT Kharagpur Production

www.nptel.iitm.ac.in

Copyrights Reserved