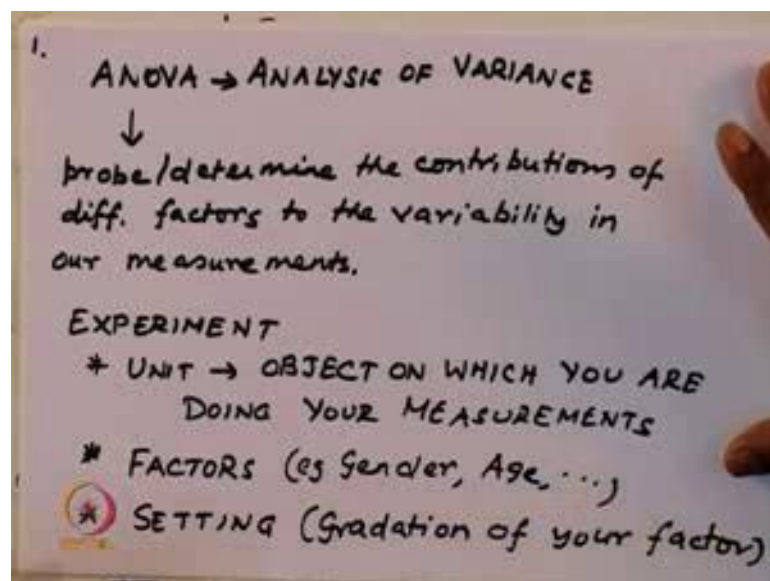


Introduction to Biostatistics
Prof. Shamik Sen
Department of Bioscience and Bioengineering
Indian Institute of Technology, Bombay

Lecture - 38
ANOVA

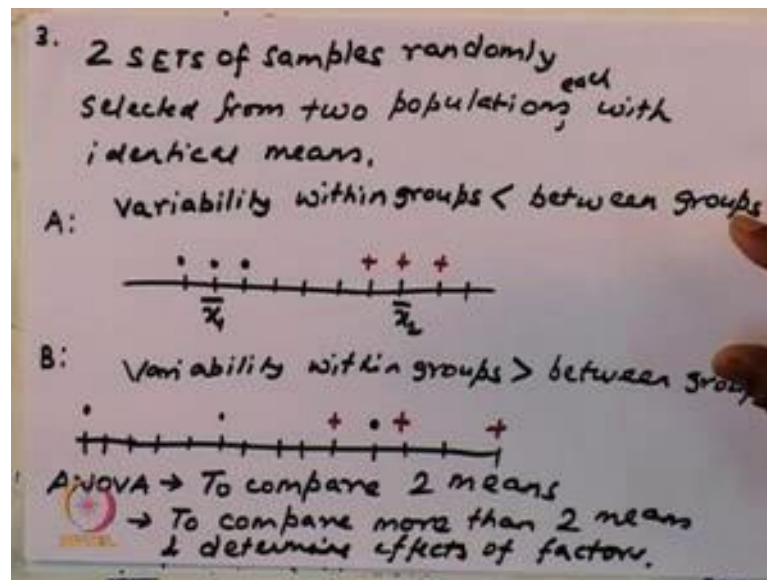
Hello and welcome to our lecture today. So, we would briefly recap what we had discussed in last lecture which was analysis of variance right or ANOVA.

(Refer Slide Time: 00:27)



So, ANOVA allows us to probe or determine the contributions of different factors to the variability in our measurements. So, when you do an experiment you have the experimental unit which is the object on which you are doing your measurements. You have your factors or example gender is a factor, or let us say height age can be a factor so and so forth. You have settings, which is the gradation of your factor. So, also something called treatment.

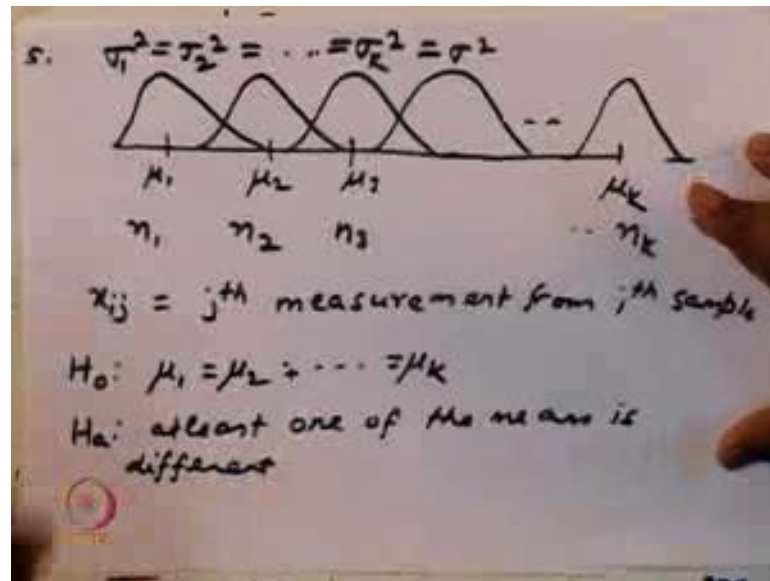
(Refer Slide Time: 05:26)



So, imagine you have 2 sets of samples randomly selected from 2 populations, with identical means each with identical means. So, let me draw one particular case. So, imagine one population; one population one of the samples gives you these value. The other of the samples other of the samples gives you these value. So, this is your scenario A.

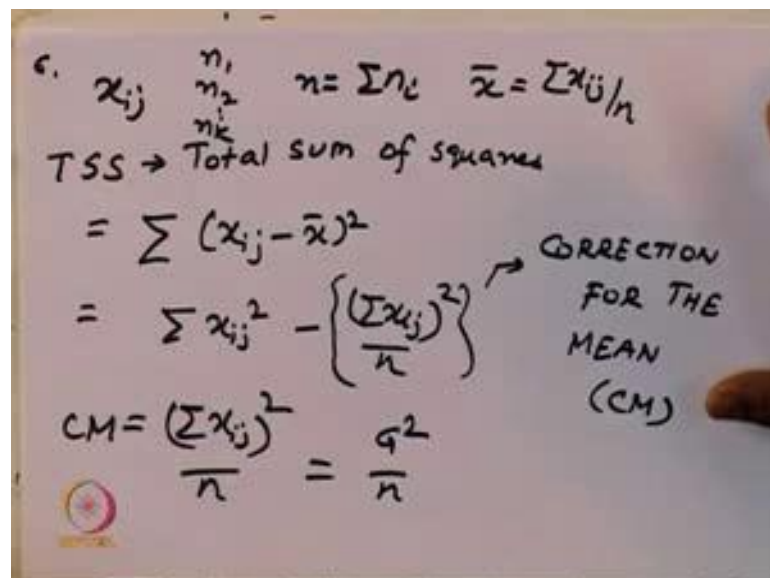
Imagine the other scenario. So, this is your x 1 bar and this is your x 2 bar. The other case let us say you have one population here. So, these are your 2 situations. So, what you observe is in this situation A there is the variability within the groups is much less than between groups, versus in this case the variability within groups is much greater than that between groups. So, this is where the ANOVA approach is important because ANOVA can be used to compare 2 means. It can also be used to compare more than 2 means and determine effects of various factors. So, in ANOVA you can do one of the 2. So, as we said let us take an example. So, there are various ways in which you can draw the sample. So, one of the experimental ways.

(Refer Slide Time: 11:21)



So, here you have $\mu_1, \mu_2, \mu_3, \mu_k$ each of these populations your variances are same. So, any want to test the heights. So, let us say your sample sizes are n_1, n_2, n_3, n_k . And you can have x_{ij} is the j th measurement from i th sample. So, you want to test the hypothesis, that $H: \mu_1 = \mu_2 = \mu_3$ so and so forth. So, your null hypothesis is $\mu_1 = \mu_2 = \mu_k$, and alternate hypothesis is at least one of the means is different.

(Refer Slide Time: 13:09)



So, what you do? So, as I said x_{ij} is your j th measurement from i th sample. So, you calculate the following quantities TSS is total sum of squares. And you have samples drawn from n_1, n_2, \dots, n_k you define n as summation of all n_i . So, TSS is defined as total sum of squares, and let us say \bar{x} is summation x_{ij} by n . So, TSS is defined as summation of x_{ij} minus \bar{x} whole square. So, you can show if you expand this you can show this is nothing, but summation of x_{ij} square minus. So, this term is called the correction for the mean or CM. So, CM is summation x_{ij} whole square by n and this you can write as G^2 by n . So, G represents the basically sum of all the terms.

(Refer Slide Time: 14:44)

7. $TSS = SST + SSE \rightarrow$ Sum of squares of errors

↓

Sum of squares for treatments

$$SST = \sum n_i (\bar{x}_i - \bar{x})^2$$

$$SST = \sum \frac{T_i^2}{n_i} - CM$$

$$SSE = (n_1 - 1)S_1^2 + (n_2 - 1)S_2^2 + \dots + (n_k - 1)S_k^2$$

$$TSS = SST + SSE$$

So, in ANOVA you distributed the TSS into 2 fragments, you distribute that TSS into 2 fragments one you call as SST. So, TSS equal to SST plus SSE sum of squares of treatments. So, this stands for sum of squares for treatments and this term MST or sum of squares of errors. So, you define SST as summation of x_i bar is the sample average for the i th sample n_i is a sample size of the i th sample. So, this you can again expand and show you can show that SST is same as t_i square by n_i , t_i is the sum of all in i th sample total sum of i th sample. So, an SSE is given by, so, you can show that TSS is equal to SST plus SSE. So, what can I write over the degree of freedom or degrees of freedom for each of these terms.

(Refer Slide Time: 16:36)

B. Degrees of freedom (df)

$$\begin{aligned}TSS &\rightarrow n-1 \\SST &\rightarrow k-1 \\SSE &\rightarrow (n_1-1) + \overbrace{(n_2-1) + \dots + (n_k-1)}^{k \text{ times}} \\&= \sum n_i - \sum 1 \text{ (k times)} \\&= n-k\end{aligned}$$
$$TSS = SST + SSE \quad MSE = \frac{SSE}{n-k}$$
$$n-1 = k-1 + n-k$$

~~MSE~~ / MSE = Mean $MST = \frac{SST}{k-1}$

For TSS you have n terms which you square and add up. So, TSS has to be n minus 1 for SST you have k terms. So, this for SST you do this over all k terms. So, SST is the degree of freedom for SST is k minus 1. And for the errors SSE so, you haven't 1 minus 1 square plus so and so forth. So, your SSE has contribution from 1 minus 1 plus n_2 minus 1 plus dot dot dot n_k minus 1. So, this is k times. So, this you can simplify equal to $n-1$. So, summation n_i you have n terms all together and minus 1 which is k times summation one k times this is nothing, but n minus k .

So, given that TSS is equal to SST plus SSE. So, degrees of freedom of n minus 1 is equal to k minus 1 plus n minus k . So, the degrees of freedom also add up. And for corresponding to these terms you can define MSE which is the mean square either you can define MST mean square for treatments or mean square for error it is. So, MST is defined as SST by this degree of freedom that is SST by k minus 1. And MSE is the mean square of error is given by SSE by n minus k .

(Refer Slide Time: 18:50)

9. ANOVA TABLE

SOURCE	df	SS	MS	F-value
TREATMENTS	$k-1$	SST	MST	f
ERROR	$n-k$	SSE	MSE	

$TSS = \sum x_{ij}^2 - CM$ $CM = \frac{(\sum x_{ij})^2}{n}$

$SST = \sum \frac{T_i^2}{n_i} - CM$

$T_i \rightarrow$ Total of sample 'i'

So, when you plot all these together you generate an ANOVA table which look something like this. Your typical ANOVA table will look something like this - source, degrees of freedom, sum of squares, mean squares and f value. Your source correspondence to each of the treatments, if you have 2 treatments you will have 2 if you have 4, you will have 4 of these conditions and you have accumulation from error is degrees of freedom is k minus 1 and n minus k. This is your SST or SSE your MST or MSE and you will get some value some f value.

So, again, you have TSS is summation x_{ij} square minus CM SST. So, CM is CM is given by summation x_{ij} whole square by n SST is summation of t_i square by n_i minus CM, where t_i is a total of sample i. So, let us do a case.

(Refer Slide Time: 20:34)

10. EXPERIMENT
→ nutrition on attention spans in class

ATTENTION TIMES			$n_1 = 5$
NB	LB	HB	$n_2 = 5$
8	14	10	$n_3 = 5$
7	16	12	$n = \sum n_i = 15$
9	12	16	'k' → <u>$k = 3$</u>
13	17	15	
10	11	12	
<hr/>			
$T_1 = 47$	$T_2 = 70$	$T_3 = 65$	

So, imagine, you are doing an experiment where you wish to study the effect of nutrition on attention spans in class. In other words, you want to see that is there a difference between students who have their breakfast and then come to class whether they pay more attention compared to students who have light breakfast or no breakfast. So, I have a plot of attention times and I have 3 categories. Let us say 3 treatments students who did not have any breakfast. So, no breakfast the values are who had light breakfast and who had heavy breakfast.

So, for each of these cases I can calculate t_i . So, t_i is what is total of sample i . So, you haven't 1 it is a sample size for condition 1 which is equal to 5, n_2 is also 5 n_3 equal to 5. So, n is summation of n_i is equal to 15. So, what is the value total sum of t_i , you can sum these up. 15 this is 47, this is 70 and this, this is t_1 this is t_2 and t_3 . I can find out as 65.

So, what is the k , k value is the number of populations. So, in our case k is equal to 3 you have 3 different populations. So, summation x_{ij} is you add up all these values. So, we should come out to be 47 plus 70 plus 65 equal to 182. So, your CM.

(Refer Slide Time: 23:01)

The image shows a whiteboard with handwritten mathematical formulas for ANOVA. The calculations are as follows:

$$\sum x_{ij} = 47 + 70 + 65 = 182$$
$$CM = \frac{(\sum x_{ij})^2}{n} = \frac{182^2}{15}$$
$$SST = \sum \frac{T_i^2}{n_i} - CM$$
$$= \frac{47^2}{5} + \frac{70^2}{5} + \frac{65^2}{5} - CM$$
$$\approx 130.58.5$$
$$SST = TSS - CM$$
$$= 129.7$$
$$SSE = TSS - SST = 71.2$$

The correction for the mean is summation x_{ij} whole square by n should come out to be 182 square by 15. Your TSS is summation t_i square by n_i minus CM. So, which you can calculate as 47 square by 5 plus 70 square, 70 square by 5 plus 65 square by 5 minus, CM and you will get a value of this equal to roughly 130.

SST is given by summation t_i square by n_i minus CM, your TSS, So, this is actually gives you a value of SST will give you a value of roughly 55 58.5. TSS you have to add up all the squares.

So, in other words you have to do 8 square summation x_{ij} whole square. So, TSS is equal to 8 square plus 7 square plus 9 square you square up all the terms, in for each of these conditions and minus you do CM. So, this comes out to be value of 127.7. TSS returns your value of SST. So, you can calculate SSE is equal to TSS minus SST and this gives you a value of roughly 58.53, no SSE gives you a value of 71.2.

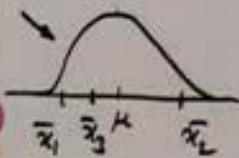
(Refer Slide Time: 25:13)

12

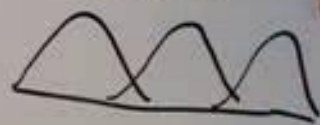
ANOVA TABLE

Source	df	SS	MS
B/f	2	58.5	29.27
Error	12	71.2	5.9
Total	14	129.7	

$H_0: \mu_1 = \mu_2 \dots = \mu_k$



$H_0 = \text{false}$



So, based on these values we can calculate the ANOVA table. So, your source is you know is your meal or breakfast your df you have SS you have MS and you have error. So, for b/f your degrees of freedom is 2 because k is equal to 3, for error you have equal to 12 n minus k will give you a value of n equal to 15. So, n minus k will give you a value of 12. This we calculate as 58.5 this you calculate as 71.2. You can accordingly calculate MS and ME and you can do the total this is 14 this is 129.7. So, this is your ANOVA table.

Now, what do you need to do to calculate your test your hypothesis? So, your H_0 is $\mu_1 = \mu_2 = \dots = \mu_k$. So, if these means for all same then, you would have had a distribution let us say hypothetically, you can have \bar{x}_1 here \bar{x}_2 here \bar{x}_3 here. So, in this case you might have agreed to say H_0 is true; however, if your case was something like this here your H_0 is false.

(Refer Slide Time: 27:09)

13. $\sigma^2 = \text{Common Variance}$

$MSE = \frac{SSE}{n-k} \rightarrow \text{estimate of } \sigma^2$

If $H_0 = \text{true}$,

$MST = \frac{SST}{k-1} \rightarrow \text{unbiased estimate of } \sigma^2$

TEST STATISTIC:

$F = \frac{MST}{MSE}$

The graph shows an F-distribution curve with a vertical line at F_α and the area to the right of this line labeled α .

So, sigma square your assumption is sigma square is common variance for all k. So, your MSE is given by SSE by n minus k. It is an estimate of sigma square and if your H_0 was true. So, your MST, MST which is SST by k minus 1, it should give you an unbiased estimate of sigma square.

So, you can use the test statistic as f equal to MST by MSE. And as before you can use your f test calculate f of alpha. And then see whether you can test whether if your f value is greater than f of alpha then you know your hypothesis is not true.

With that I conclude my talk today. So, you get an idea of how ANOVA can be used instead of repeatedly using students to test for calculating whether means are same. You can have come to this conclusion using a single approach which is your ANOVA.

So, you create your ANOVA table by calculating by distributing the total sum of squares or TSS into SST or sum of squares of treatments and SSE which accounts for random error. So, sum of squares of errors from there you calculate the mean square error or the mean square sum of t and then use the statistic MST by MSE to find out whether your means are same or they are distinct.

Thank you for your attention.