

Practitioners Course in Descriptive, Predictive and Prescriptive Analytics
Prof. Deepu Philip
Dr. Amandeep Singh Oberoi
Mr. Sanjeev Newar
Department of Industrial and Management Engineering
Indian Institute of Technology, Kanpur
National Institute of Technology, Jalandhar

Lecture – 24
Analysis of Variance (ANOVA) – (Part-2)

So, welcome to the part 2 of analysis of variance session.

(Refer Slide Time: 00:21)

Conducting One-Way ANOVA
Decompose the total variation

$$SS_{\text{total}} = SS_{\text{X}} + SS_{\text{error}}$$

$$SS_{\text{total}} = \sum_{i=1}^N (y_i - \bar{y})^2$$

$$SS_{\text{X}} = \sum_{j=1}^C n_j (\bar{y}_j - \bar{y})^2$$

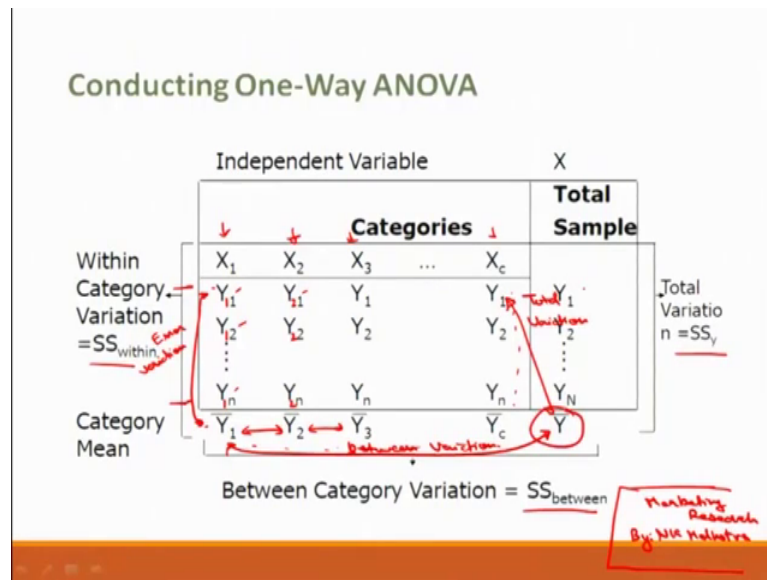
$$SS_{\text{error}} = \sum_{j=1}^C \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_j)^2$$

$N = \text{Sample Size}$
 $n = \text{Group Size}$
 $N = C \times n$
 $C = \text{No. of categories}$

n (5)
(1, 2, 3)
..... C
No. of categories

So, where were we? We were discussing these relationships for total variation, between variation and error variation.

(Refer Slide Time: 00:33)



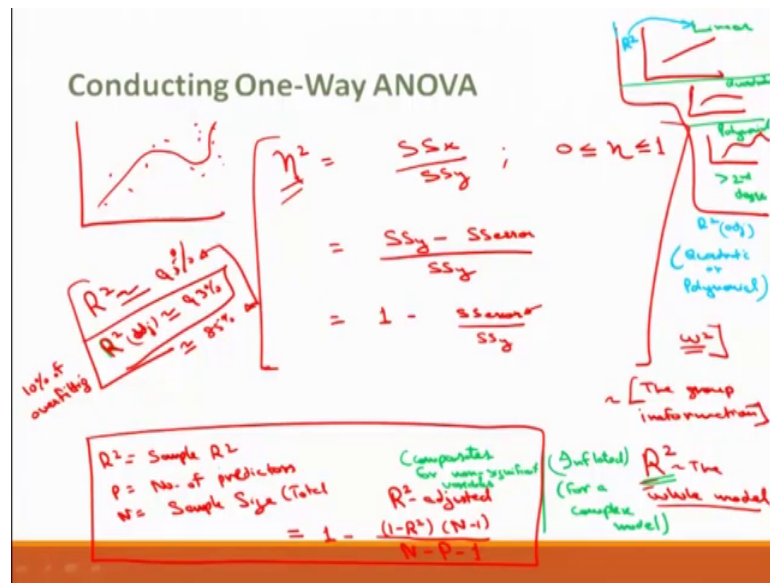
So, let me put it in a tabular form here. So, we have this table. So, this data is taken from the book by NK Malhotra. This illustration is taken from here it is marketing research the example which I will discuss here is also taken from the same book. So, you can read further if you like to have more information. So, what we have? We have these categories. So, mind it we had capital N here; capital N here is sample size not the population size.

So, small n is my group size; where capital N is divided into divided into C groups and therefore, C here is the number of categories. So, we have c categories 1 to C categories here 1 to n is my group size we can see the all group sizes are same here and we will see that the 1 variation here is between my element Y and Y 1 bar you can see this here the error $Y_{ij} - \bar{Y}_j$ you can see this here the error $Y_{ij} - \bar{Y}_j$ is $Y_{11} - \bar{Y}_1$ this minus this is error variation ok.

So, between these between these implies between \bar{Y}_1 and \bar{Y} this is my between variation and for each element for Y_{11} Y_{12} shown up to Y_{1n} and. So, on for category 2 Y_{21} why I can even put this is Y_{21} , Y_{22} , Y_{2n} this can be Y_{21} Y_{22} Y_{2n} for each element here for each Y_i i varies from 1 to n here, we have the total variation I can put it here this is total variation. So, total variation is SS_y between category variation is $SS_{between}$ and within category variation is SS_{within} . So, we have this category means here and this is the total sample mean \bar{Y} .

So, let us see how to conduct this ANOVA analysis using this table.

(Refer Slide Time: 04:16).



So, to conduct ANOVA we need to see the eta square. Eta square is SS x by SS y that is between variation over total variation, which varies from 0 to 1 this can also be put as SS x can be put as SS y minus SS error over SS y, we can say this 1 minus SS error by SS y. So, if you see this we have come SS error here. So, eta square and 1 more statistic known as omega square are more focused on the group information, it tells us about the within group variation and there is statistic known as R square which is nothing, but square of the correlation coefficient that tells us about the whole model.

So, R square and eta square are varying two different things, R square is the contribution of the entire model the whole model in explaining this study variation, eta square and omega square measure the contribution of individual model term. So, they are different eta square has a positive bias; So, it over estimates the individual contributions. So, from a practical stand point the bias is small. So, it depends upon your situation the data behavior how we have obtained that.

So, R square can be significantly inflated by adding non significant terms in the model. So, there is a term known as R square adjusted; R square adjusted is for compensating the non significant terms, which have inflated. So, we have inflation here in R square it is inflated for a complex model and R square adjusted compensates with the compensate note with compensate for, non significant terms.

So, the R square value has a bias it over fits the model, there are chance correlations here that trends in (Refer Time: 07:49) and certain variable forms can also influence this one. So, R square adjusted is calculated as I will put it here, R square adjusted is $1 - \frac{1 - R^2}{n - p - 1}$ where R square is the sample R square, p is the number of predictors and n is the sample size a sample size I mean that total sample size.

So, this is R square adjusted that gives the more close or more accurate information about the fit of the model. So, by fit of the model I mean if R square value is let me say if this is 95 percent, but it is over fitting memory it is its telling that the model is 95 percent close to the original data means for instance I have the data points here, and this is the model here. This is the kind of a cubical model the way I have drawn it.

So, we have various kinds of models, this is given linear model, this would be a quadratic model and this would be a cubical or some polynomial model linear quadratic and polynomial that is more than second degree more than second degree polymer. So, in this case if R square value is telling 95 percent value, R square adjusted value is telling R square adjusted value is somewhat 93 percent than this is more close or sometimes if R square adjusted value is may be 85 percent and this difference this 10 percent difference is over fitting; that means, 10 percent of over fitting. So, this R square is used for the linear model for the quadratic and the polynomial model, we used R square adjusted this is for linear model and this is for quadratic or polynomial.

So, these terms we need to calculate to see the closeness to see the fit of our model ANOVA tower the model we have generated ANOVA model would generate the regression models here and even we can have certain other advanced models for which we du conduct the ANOVA beforehand advance models might be new networks genetics algorithm artificial B column though all those techniques are kind of advanced techniques, but ANOVA can give the beforehand and the overall field that how would our model look like, how would over model behave what is the curve fitting in case of regression fitting we are talking about.

(Refer Slide Time: 12:12)

Conducting One-Way ANOVA

In one-way analysis of variance, the interest lies in testing the null hypothesis that the category means are equal in the population.

$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_c$

Under the null hypothesis, SS_x and SS_{error} come from the same source of variation. In other words, the estimate of the population variance of Y_i

$S_y^2 = SS_x / (c - 1)$ *degrees of freedom*
= Mean square due to X
= MS_x

or

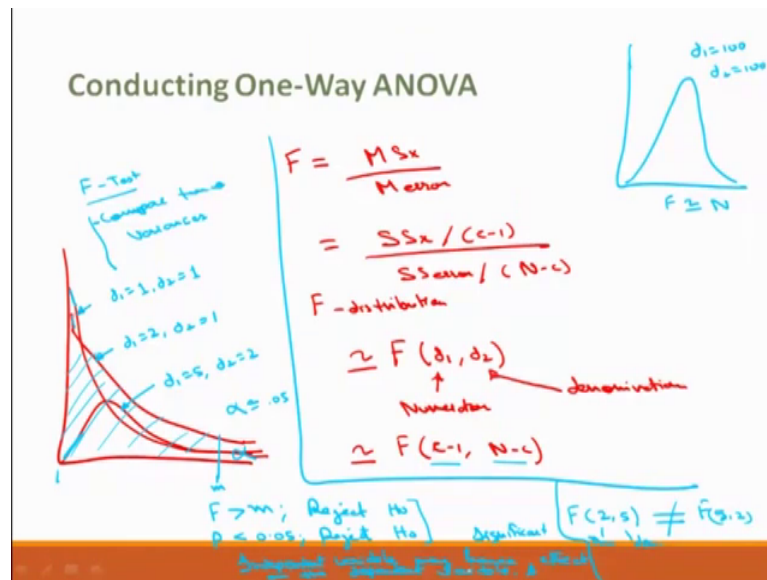
$S_y^2 = SS_{error} / (N - c)$ *degrees of freedom*
= Mean square due to error
= MS_{error}

So, this is what ANOVA would help us. So, in one way analysis of variance, the interest lies in testing hypothesis that category measure all equal in the population and as I mentioned before; that means, all category means are equal; that means, the all come from the same population.

So, under this null hypothesis SS_x and SS_{error} come from the same source of variation, in other words the estimate of the population variance of Y this S_y square root can be calculated as SS_x by the degrees of freedom. So, this SS_x . So, this $c - 1$ is the degrees of freedom this is also my degrees of freedom.

So, if you see this estimate of population variance this is nothing, but mean square due to x , which we denote as ms_x and the second one is mean square due to error in which we have this error degrees of freedom and if we have the category or group degrees of freedom.

(Refer Slide Time: 13:38)



So, this is ms error; then we can have the F statistic value. Here the F which I just said is Ms x by Ms error this is equal to SS x by its degree of freedom and SS error by its degrees of freedom. So, this statistic follow F distribution. F distribution we write it like this F degrees of freedom 1 and degrees of freedom 2. Please make a note this d 1 is always numerator and this is denominator. So, in this case our degrees of freedom would be F C minus 1 and N minus C. So, F distribution is a writes queued distribution it is like this, something like this or may be like this or may be even like this, it depends upon the degrees of freedom for which we are trying to represent this F distribution.

Here we can have the alpha value that is, if F value if F statistic value is greater than this specific value then let me say this value is m. If F is greater than m we reject H naught and if it lies in this region this is excepted or in other words you can say if the probability this is p, alpha is my significance level alpha it may say if it is 0.05 for 95 percent confidence level or 5 percent significance level. If P value is less than 0.05, we reject H naught of the equal variances.

So, F test is the catch of all term for any test that it is F distribution. In most cases when people talk about F test they may are actually talking out the F test to compare two variances. So, F test is used to compare two variances and please note here that F let me say 2 5 is different this is numerator this is denominator this is not equal to F 5 2.

So, taking the freedom from numerator comes first for denominators comes second and why are these lines 3 lines 1 2 3 why are these different? This depend upon the kind of the degrees of freedom for the numerator and denominator like this line can be for d 1 is equal to 1, d 2 is equal to 1. The second line this which is having a less slope this can be for d 1 is equal to 2 and d 2 is equal to 1 the third line that is kind of a writes clued normal distribution, this third line is for t 1 is equal to 5 and t 2 is equal to 2.

So, for higher degrees of freedom we can even see the normal trend here, for let me say it would behave like this. If I say d 1 is equal to 100 and d 2 is equal to 100. So, very high degrees of freedom F is an approximation of normal. So, if this null hypothesis of equal category means is not rejected, then independent variables does not have significant effect on dependent variable.

This is rejected here; that means, the independent variable may have effect on the dependent variable. So, in this case the null hypothesis is rejected; that means, that the independent variable is significant, I can say here significant effect. It has significant effect on the dependent variable the comparison of category mean values will indicate the nature of effect of independent variable. So, let us try to compare category means here.

(Refer Slide Time: 19:55)

Conducting One-Way ANOVA

Illustrative example

| Store Number | Coupon Level | In-Store Promotion | Sales |
|--------------|--------------|--------------------|-------|
| 1 | 1.00 | 1.00 | 4.00 |
| 2 | 1.00 | 1.00 | 9.00 |
| 3 | 1.00 | 1.00 | 10.00 |
| 4 | 1.00 | 1.00 | 8.00 |
| 5 | 1.00 | 1.00 | 9.00 |
| 6 | 1.00 | 2.00 | 8.00 |
| 7 | 1.00 | 2.00 | 8.00 |
| 8 | 1.00 | 2.00 | 7.00 |
| 9 | 1.00 | 2.00 | 9.00 |
| 10 | 1.00 | 2.00 | 6.00 |
| 11 | 1.00 | 3.00 | 5.00 |
| 12 | 1.00 | 3.00 | 7.00 |
| 13 | 1.00 | 3.00 | 6.00 |
| 14 | 1.00 | 3.00 | 4.00 |
| 15 | 1.00 | 3.00 | 5.00 |
| 16 | 2.00 | 1.00 | 8.00 |
| 17 | 2.00 | 1.00 | 9.00 |
| 18 | 2.00 | 1.00 | 7.00 |
| 19 | 2.00 | 1.00 | 7.00 |
| 20 | 2.00 | 1.00 | 6.00 |
| 21 | 2.00 | 2.00 | 4.00 |
| 22 | 2.00 | 2.00 | 5.00 |
| 23 | 2.00 | 2.00 | 5.00 |
| 24 | 2.00 | 2.00 | 6.00 |
| 25 | 2.00 | 2.00 | 4.00 |
| 26 | 2.00 | 3.00 | 2.00 |
| 27 | 2.00 | 3.00 | 3.00 |
| 28 | 2.00 | 3.00 | 2.00 |
| 29 | 2.00 | 3.00 | 1.00 |
| 30 | 2.00 | 3.00 | 2.00 |

Handwritten notes on the slide:

- Dependent variable: Sales
- Independent variable: Coupon level, In-store promotion
- Scale: Ordinal
- 1 - Low
- 2 - Medium
- 3 - High

So, we have this data here, in which we have this store sales the in store promotion and coupon level and this is store number.

So, let me recall this test that we discussed. First we will try to identify the dependent and independent variable then will decompose the total variation total variation in to between groups and within groups, then we measure the effects that is eta square, omega square, R square then will test this significance for the specific level of significance for the specific significance will be set 5 percent significance level, then we interpret the results. In interpretation of the results we see that which is the most influential factor here. So, first let us try to identify, which is my dependent variable and what are my independent variables ok.

So, just please make a focused look on this, coupon level, in store promotion and sales, which is our dependent variable what is a store trying to look here. I think you are able to identify this the dependent variables would be the sales. So, dependent variable is sales independent variable. There are two independent variables, this store number is just a number of store you can call it as serial number or this is just the nominal data ok. Then independent variable is the coupon level and in store promotion and in store promotion.

So, if we see this we can put the types of this scales here, is it metric or non-metric scales dependent variables yeah the sales is a Indian metric scale, I will put the scale here it is interval scale and this one coupon level is nothing, but these are categories and in store promotion is also my categories. You can see 1 and 2 there are 2 categories here coupon level, and in store promotion they are 1 2 3 for coupon 1 and 1, 2, 3 for coupon 2, coupon 1 and this is for coupon 2 this is coupon 1 and 2.

So, if I say for in store promotion this 1 is low, 2 is medium and 3 is high. So, this is my ordinal scale here. So, I will better put it here this is the ordinal scale low medium and high level of promotion.

(Refer Slide Time: 23:49)

Conducting One-Way ANOVA

EFFECT OF IN-STORE PROMOTION ON SALES

| Store No. | Level of In-store Promotion | | |
|-----------------------------|-----------------------------|---------------|---------------|
| | High (1) | Medium (2) | Low (3) |
| 1 | 10 | 8 | 5 |
| 2 | 9 | 8 | 7 |
| 3 | 10 | 7 | 6 |
| 4 | 8 | 9 | 4 |
| 5 | 9 | 6 | 5 |
| 6 | 8 | 4 | 2 |
| 7 | 9 | 5 | 3 |
| 8 | 7 | 5 | 2 |
| 9 | 7 | 6 | 1 |
| 10 | 6 | 4 | 2 |
| Column Totals | 83 | 62 | 37 |
| Category means: \bar{y}_j | $83/10 = 8.3$ | $62/10 = 6.2$ | $37/10 = 3.7$ |
| Grand mean, \bar{y} | $(83 + 62 + 37)/30 = 6.067$ | | |

Handwritten notes:
 - Group size = 10
 - No. of groups = 3
 - $df = 3 - 1 = 2$
 - Sample Size = $10 \times 3 = 30$
 - $df = N - c = 30 - 3 = 27$

So, could I divide this in to these stores we have total 30 stores 1 2 30 stores? So, I have divided this is to 3 groups of in store promotion, that is low this is low and medium and number 3 this 3 three is high so; that means, we have got 3 categories here 3 categories high, medium and low and the category size is 10. So, this is my group size group size is 10 and this size that is number of groups this is equal to 3, this is equal to 10.

So, if I see the degrees of freedom for number of groups that would be the df would be 3 minus 1 is equal to 2 and we have sample size is equal to 10, that is group size N into C that is 10 into 3 is equal to 30 and degrees of freedom for this sample if you recall those were N minus C that would local to 30 minus 3 is equal to 3 minus 3 is equal to 27. So, we will use these 2 terms in the calculations.

So, this is the group mean. Group mean is 8.3, group total is 83 group mean for high, for medium, for low we have group means and we have grand mean of the whole data of the overall data this is the overall mean.

(Refer Slide Time: 26:12)

Conducting One-Way ANOVA

To test the null hypothesis, the various sums of squares are computed as follows:

$$\begin{aligned}
 SS_y &= (10-6.067)^2 + (9-6.067)^2 + (10-6.067)^2 + (8-6.067)^2 + (9-6.067)^2 \\
 &+ (8-6.067)^2 + (9-6.067)^2 + (7-6.067)^2 + (7-6.067)^2 + (6-6.067)^2 \\
 &+ (8-6.067)^2 + (8-6.067)^2 + (7-6.067)^2 + (9-6.067)^2 + (6-6.067)^2 \\
 &+ (4-6.067)^2 + (5-6.067)^2 + (5-6.067)^2 + (6-6.067)^2 + (4-6.067)^2 \\
 &+ (5-6.067)^2 + (7-6.067)^2 + (6-6.067)^2 + (4-6.067)^2 + (5-6.067)^2 \\
 &+ (2-6.067)^2 + (3-6.067)^2 + (2-6.067)^2 + (1-6.067)^2 + (2-6.067)^2 \\
 \\
 &= (3.933)^2 + (2.933)^2 + (3.933)^2 + (1.933)^2 + (2.933)^2 \\
 &+ (1.933)^2 + (2.933)^2 + (0.933)^2 + (0.933)^2 + (-0.067)^2 \\
 &+ (1.933)^2 + (1.933)^2 + (0.933)^2 + (2.933)^2 + (-0.067)^2 \\
 &+ (-2.067)^2 + (-1.067)^2 + (-1.067)^2 + (-0.067)^2 + (-2.067)^2 \\
 &+ (-1.067)^2 + (0.9333)^2 + (-0.067)^2 + (-2.067)^2 + (-1.067)^2 \\
 &+ (-4.067)^2 + (-3.067)^2 + (-4.067)^2 + (-5.067)^2 + (-4.067)^2 \\
 &= 185.867
 \end{aligned}$$

So, let us see the calculations SS total SS total SS y what was that? That was the variation of each value 10 minus the grand mean, this 9 minus grand mean 10 minus grand mean this 8 minus grand mean, 9 minus 10 grand mean so on we have 30 items and we keep on taking this squares of this, this is the sum of squares. So, this is for 10 minus the grand mean, 9 minus grand mean, 10 minus grand mean.

So, this is for group 1, group 2 and group 3 we got these values and we got the final total variation as 185.867. So, this is the calculation this is how the mechanism of ANOVA work. Actually when you just know this mechanism, when you work on the software you work on the packages like r package excel or if you have some spss or minute dial software's, you need not to do all this calculations this all would be done by the computers, but you should know that what is happening in the GUI what is the software doing what is the graphic use interface in that.

So, this thing would help you in that understanding how ANOVA works.

(Refer Slide Time: 27:42)

Conducting One-Way ANOVA

$$\begin{aligned}SSx &= 10(8.3-6.067)^2 + 10(6.2-6.067)^2 + 10(3.7-6.067)^2 \\ &= 10(2.233)^2 + 10(0.133)^2 + 10(-2.367)^2 \\ &= 106.067\end{aligned}$$
$$\begin{aligned}SSerror &= (10-8.3)^2 + (9-8.3)^2 + (10-8.3)^2 + (8-8.3)^2 + (9-8.3)^2 \\ &+ (8-8.3)^2 + (9-8.3)^2 + (7-8.3)^2 + (7-8.3)^2 + (6-8.3)^2 \\ &+ (8-6.2)^2 + (8-6.2)^2 + (7-6.2)^2 + (9-6.2)^2 + (6-6.2)^2 \\ &+ (4-6.2)^2 + (5-6.2)^2 + (5-6.2)^2 + (6-6.2)^2 + (4-6.2)^2 \\ &+ (5-3.7)^2 + (7-3.7)^2 + (6-3.7)^2 + (4-3.7)^2 + (5-3.7)^2 \\ &+ (2-3.7)^2 + (3-3.7)^2 + (2-3.7)^2 + (1-3.7)^2 + (2-3.7)^2 \\ &= (1.7)^2 + (0.7)^2 + (1.7)^2 + (-0.3)^2 + (0.7)^2 \\ &+ (-0.3)^2 + (0.7)^2 + (-1.3)^2 + (-1.3)^2 + (-2.3)^2 \\ &+ (1.8)^2 + (1.8)^2 + (0.8)^2 + (2.8)^2 + (-0.2)^2 \\ &+ (-2.2)^2 + (-1.2)^2 + (-1.2)^2 + (-0.2)^2 + (-2.2)^2 \\ &+ (1.3)^2 + (3.3)^2 + (2.3)^2 + (0.3)^2 + (1.3)^2 \\ &+ (-1.7)^2 + (-0.7)^2 + (-1.7)^2 + (-2.7)^2 + (-1.7)^2 \\ &= 79.80\end{aligned}$$

Similarly, we have SS between. SS between is the thing if you see this was the between means the grand mean and this means 10 minus 6.067, this 10 minus 6.067 or this is nothing, but difference in the categories and the grand mean; that is 8.3 this mean minus this mean. Then 6.2 minus 6.067 then 3.7 minus 6.067 whole square we have only 3 categories and we have total this is my group sizes 10 10 10.

So, this value comes down to 106.067. So, we have taken it for all the values you can see between. So, SS error. SS error is again the category mean here is for category one that is high level 8.3, category mean that difference of this mean from the elements 10 minus 8.3, this 9 minus 8 point 3, 10 minus 8.3 whole square this 8 minus 8.3 whole square so on we then move to second category 8 minus 6.2 this 8 minus 6.2 whole square plus 7 minus 6.2 whole square and so on then we move to category 3, similarly and we get the SS error and we got the total value as this much.

So, you might have noted that 106.067 plus 79.80 is equal to 185.867. So, that our decomposition of total variation that we saw before holds good here

(Refer Slide Time: 30:00)

Conducting One-Way ANOVA

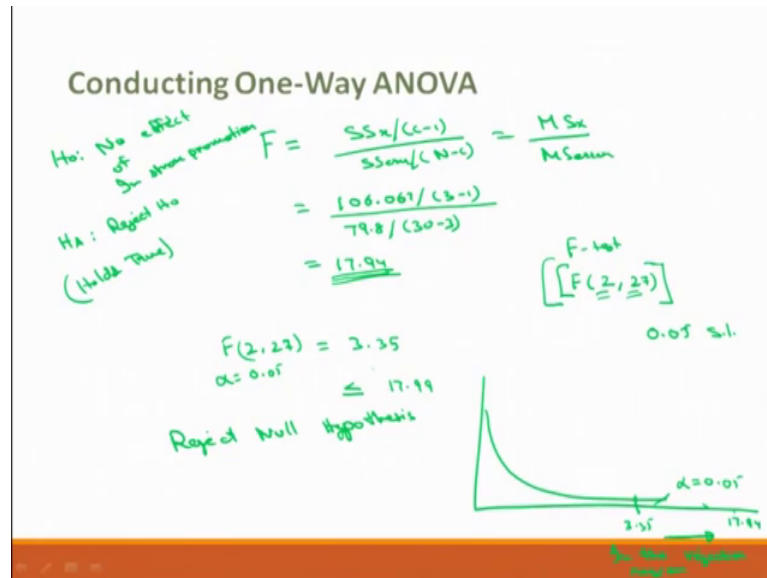
$$SS_y = SS_x + SS_{error}$$
$$185.867 = 106.067 + 79.80$$
$$\eta^2 = \frac{SS_x}{SS_y}$$
$$= \frac{106.067}{185.867} = 0.571$$

57.1% of Sales is (due to) in-store promotion' accounts for

So, SS y is equal to SS x plus SS error so; that means, 185.867 is equal to 106.067 plus 79.80. So, the strength effects that is eta square that we can calculate by sd x by SS y this comes down to SS x value is 106.067 and SS y value is 185.867. So, this value is 0.571. So, in other words we can say that 55, 57.1 percent of variation in sales is counted by in store promotion here. 57.1 percent of sales is due to the variable this categorical variable in store promotion.

So, now, let us test the null hypothesis, we have got this statistic here and this is the value that test the null hypothesis using F statistic here.

(Refer Slide Time: 31:48)



So, F statistic is SS x that is mean square for x that is SS x minus degrees of SS x over this degrees of freedom, then SS error by its degrees of freedom. So, this is I will put it again this is Msx and Ms error. This value comes down to SS x value is 106.067 SS error value is 79.8. 106.067 by 3 minus 1 and ms error is 79.8 by 30 minus 3.

So, this value comes down to 17.94. So, we can test this F statistic we can conduct the F test we have the F test using F distribution with F what is the degrees of freedom from numerator here? 3 minus 1 2 and degrees of freedom for denominator is 30 minus 3 27 for this we can see the value of F in the table and we will see whether this value is 17.94 is greater than the tab tabulated value here.

If it is greater than we will reject the null hypothesis, then we say we can say that yes whatever the hypothesis we are saying here this sale is due to I can put here better word is accounts for. So, let me test the significance. So, for 0.05 significance level we can see the F table here.

(Refer Slide Time: 34:08)

Conducting One-Way ANOVA

F distribution

| F | Degrees of Freedom in the Numerator | | | | | | | | | | | | | | | | | | | |
|------|-------------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 12 | 15 | 20 | 25 | 30 | 40 | 50 | 100 | 200 | 1000 |
| 1 | 39.86 | 49.50 | 53.59 | 55.83 | 57.24 | 58.20 | 58.91 | 59.44 | 59.86 | 60.19 | 60.71 | 61.22 | 61.74 | 62.05 | 62.26 | 62.53 | 62.69 | 63.01 | 63.17 | 63.30 |
| 2 | 8.53 | 9.00 | 9.16 | 9.24 | 9.29 | 9.33 | 9.35 | 9.37 | 9.38 | 9.39 | 9.41 | 9.42 | 9.44 | 9.45 | 9.46 | 9.47 | 9.47 | 9.48 | 9.49 | 9.49 |
| 3 | 5.54 | 5.46 | 5.39 | 5.34 | 5.31 | 5.28 | 5.27 | 5.25 | 5.24 | 5.23 | 5.22 | 5.20 | 5.18 | 5.17 | 5.17 | 5.16 | 5.15 | 5.14 | 5.14 | 5.13 |
| 4 | 4.54 | 4.32 | 4.19 | 4.11 | 4.05 | 4.01 | 3.98 | 3.95 | 3.94 | 3.92 | 3.90 | 3.87 | 3.84 | 3.83 | 3.82 | 3.80 | 3.80 | 3.78 | 3.77 | 3.76 |
| 5 | 4.06 | 3.78 | 3.62 | 3.52 | 3.45 | 3.40 | 3.37 | 3.34 | 3.32 | 3.30 | 3.27 | 3.24 | 3.21 | 3.19 | 3.17 | 3.16 | 3.15 | 3.13 | 3.12 | 3.11 |
| 6 | 3.78 | 3.46 | 3.29 | 3.18 | 3.11 | 3.05 | 3.01 | 2.98 | 2.96 | 2.94 | 2.90 | 2.87 | 2.84 | 2.81 | 2.80 | 2.78 | 2.77 | 2.75 | 2.73 | 2.72 |
| 7 | 3.59 | 3.26 | 3.07 | 2.96 | 2.88 | 2.83 | 2.78 | 2.75 | 2.72 | 2.70 | 2.67 | 2.63 | 2.59 | 2.57 | 2.56 | 2.54 | 2.52 | 2.50 | 2.48 | 2.47 |
| 8 | 3.46 | 3.11 | 2.92 | 2.81 | 2.73 | 2.67 | 2.62 | 2.59 | 2.56 | 2.54 | 2.50 | 2.46 | 2.42 | 2.40 | 2.38 | 2.36 | 2.35 | 2.32 | 2.31 | 2.30 |
| 9 | 3.36 | 3.01 | 2.81 | 2.69 | 2.61 | 2.55 | 2.51 | 2.47 | 2.44 | 2.42 | 2.38 | 2.34 | 2.30 | 2.27 | 2.25 | 2.23 | 2.22 | 2.19 | 2.17 | 2.16 |
| 10 | 3.29 | 2.92 | 2.73 | 2.61 | 2.52 | 2.46 | 2.41 | 2.38 | 2.35 | 2.32 | 2.28 | 2.24 | 2.20 | 2.17 | 2.16 | 2.13 | 2.12 | 2.09 | 2.07 | 2.06 |
| 12 | 3.18 | 2.81 | 2.61 | 2.48 | 2.39 | 2.33 | 2.28 | 2.24 | 2.21 | 2.19 | 2.15 | 2.10 | 2.06 | 2.03 | 2.01 | 1.99 | 1.97 | 1.94 | 1.92 | 1.91 |
| 15 | 3.07 | 2.70 | 2.49 | 2.36 | 2.27 | 2.21 | 2.16 | 2.12 | 2.09 | 2.06 | 2.02 | 1.97 | 1.92 | 1.89 | 1.87 | 1.85 | 1.83 | 1.79 | 1.77 | 1.76 |
| 20 | 2.97 | 2.59 | 2.38 | 2.25 | 2.16 | 2.09 | 2.04 | 2.00 | 1.96 | 1.94 | 1.89 | 1.84 | 1.79 | 1.76 | 1.74 | 1.71 | 1.69 | 1.65 | 1.63 | 1.61 |
| 25 | 2.92 | 2.53 | 2.32 | 2.18 | 2.09 | 2.02 | 1.97 | 1.93 | 1.89 | 1.87 | 1.82 | 1.77 | 1.72 | 1.68 | 1.66 | 1.63 | 1.61 | 1.56 | 1.54 | 1.52 |
| 30 | 2.88 | 2.49 | 2.28 | 2.14 | 2.05 | 1.98 | 1.93 | 1.88 | 1.85 | 1.82 | 1.77 | 1.72 | 1.67 | 1.63 | 1.61 | 1.57 | 1.55 | 1.51 | 1.48 | 1.46 |
| 40 | 2.84 | 2.44 | 2.23 | 2.09 | 2.00 | 1.93 | 1.87 | 1.81 | 1.79 | 1.76 | 1.71 | 1.66 | 1.61 | 1.57 | 1.54 | 1.51 | 1.48 | 1.43 | 1.41 | 1.38 |
| 50 | 2.81 | 2.41 | 2.20 | 2.06 | 1.97 | 1.90 | 1.84 | 1.80 | 1.76 | 1.73 | 1.68 | 1.63 | 1.57 | 1.53 | 1.50 | 1.46 | 1.44 | 1.39 | 1.36 | 1.33 |
| 100 | 2.76 | 2.36 | 2.14 | 2.00 | 1.91 | 1.83 | 1.78 | 1.71 | 1.69 | 1.66 | 1.61 | 1.56 | 1.49 | 1.45 | 1.42 | 1.38 | 1.35 | 1.29 | 1.26 | 1.22 |
| 200 | 2.73 | 2.33 | 2.11 | 1.97 | 1.88 | 1.80 | 1.75 | 1.70 | 1.66 | 1.63 | 1.58 | 1.52 | 1.46 | 1.41 | 1.38 | 1.34 | 1.31 | 1.24 | 1.20 | 1.16 |
| 1000 | 2.71 | 2.31 | 2.09 | 1.95 | 1.85 | 1.78 | 1.72 | 1.68 | 1.64 | 1.61 | 1.55 | 1.49 | 1.43 | 1.38 | 1.35 | 1.30 | 1.27 | 1.20 | 1.15 | 1.08 |

So, this is F distribution this is for specific alpha value, in this case alpha is 10 percent; 0.1 is 10 percent 0.1 significance level. So, in this case we have degrees of freedom in the numerator in the columns and degrees of freedom in the denominator in the rows, but this table is not for our test we are conducting because this is for specific alpha value. The second table yes this second table is for 0.05 that we have considered the alpha value here. So, for this value let me say two and 27 degrees of freedom.

(Refer Slide Time: 34:57)

Conducting One-Way ANOVA

F distribution

| F | Degrees of Freedom in the Numerator | | | | | | | | | | | | | | | | | | | |
|------|-------------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 12 | 15 | 20 | 25 | 30 | 40 | 50 | 100 | 200 | 1000 |
| 1 | 161.4 | 199.5 | 215.7 | 224.8 | 230.2 | 234.0 | 236.8 | 238.9 | 240.5 | 241.9 | 243.9 | 245.9 | 248.0 | 249.3 | 250.1 | 251.1 | 251.8 | 253.0 | 253.7 | 254.2 |
| 2 | 18.51 | 19.00 | 19.16 | 19.25 | 19.30 | 19.33 | 19.35 | 19.37 | 19.38 | 19.40 | 19.41 | 19.41 | 19.45 | 19.46 | 19.46 | 19.47 | 19.48 | 19.49 | 19.49 | 19.49 |
| 3 | 10.13 | 9.55 | 9.28 | 9.12 | 9.01 | 8.94 | 8.89 | 8.85 | 8.81 | 8.79 | 8.74 | 8.70 | 8.66 | 8.63 | 8.62 | 8.59 | 8.58 | 8.55 | 8.54 | 8.53 |
| 4 | 7.71 | 6.94 | 6.59 | 6.39 | 6.26 | 6.16 | 6.09 | 6.04 | 6.00 | 5.96 | 5.91 | 5.86 | 5.80 | 5.77 | 5.75 | 5.72 | 5.70 | 5.66 | 5.65 | 5.63 |
| 5 | 6.61 | 5.79 | 5.41 | 5.19 | 5.05 | 4.95 | 4.88 | 4.82 | 4.77 | 4.74 | 4.68 | 4.62 | 4.56 | 4.52 | 4.50 | 4.46 | 4.44 | 4.41 | 4.39 | 4.37 |
| 6 | 5.99 | 5.14 | 4.76 | 4.53 | 4.39 | 4.28 | 4.21 | 4.15 | 4.10 | 4.06 | 4.00 | 3.94 | 3.87 | 3.83 | 3.81 | 3.77 | 3.75 | 3.71 | 3.69 | 3.67 |
| 7 | 5.59 | 4.74 | 4.35 | 4.12 | 3.97 | 3.87 | 3.79 | 3.73 | 3.68 | 3.64 | 3.57 | 3.51 | 3.44 | 3.40 | 3.38 | 3.34 | 3.32 | 3.27 | 3.25 | 3.23 |
| 8 | 5.32 | 4.46 | 4.07 | 3.84 | 3.69 | 3.58 | 3.50 | 3.44 | 3.39 | 3.35 | 3.28 | 3.22 | 3.15 | 3.11 | 3.08 | 3.04 | 3.02 | 2.97 | 2.95 | 2.93 |
| 9 | 5.12 | 4.26 | 3.86 | 3.63 | 3.48 | 3.37 | 3.29 | 3.23 | 3.18 | 3.14 | 3.07 | 3.01 | 2.94 | 2.89 | 2.86 | 2.83 | 2.80 | 2.76 | 2.73 | 2.71 |
| 10 | 4.96 | 4.10 | 3.71 | 3.48 | 3.33 | 3.22 | 3.14 | 3.07 | 3.02 | 2.98 | 2.91 | 2.85 | 2.77 | 2.73 | 2.70 | 2.66 | 2.64 | 2.59 | 2.56 | 2.54 |
| 12 | 4.75 | 3.89 | 3.49 | 3.26 | 3.11 | 3.00 | 2.91 | 2.85 | 2.80 | 2.75 | 2.69 | 2.62 | 2.54 | 2.50 | 2.47 | 2.43 | 2.40 | 2.35 | 2.32 | 2.30 |
| 15 | 4.54 | 3.68 | 3.29 | 3.06 | 2.90 | 2.79 | 2.71 | 2.64 | 2.59 | 2.54 | 2.48 | 2.40 | 2.33 | 2.28 | 2.25 | 2.20 | 2.18 | 2.12 | 2.10 | 2.07 |
| 20 | 4.35 | 3.49 | 3.10 | 2.87 | 2.71 | 2.60 | 2.51 | 2.45 | 2.39 | 2.35 | 2.28 | 2.20 | 2.12 | 2.07 | 2.04 | 1.99 | 1.97 | 1.91 | 1.88 | 1.85 |
| 25 | 4.24 | 3.39 | 2.99 | 2.76 | 2.60 | 2.49 | 2.40 | 2.34 | 2.28 | 2.24 | 2.16 | 2.09 | 2.01 | 1.96 | 1.92 | 1.87 | 1.84 | 1.78 | 1.75 | 1.72 |
| 30 | 4.17 | 3.32 | 2.92 | 2.69 | 2.53 | 2.42 | 2.33 | 2.27 | 2.21 | 2.16 | 2.09 | 2.01 | 1.93 | 1.88 | 1.84 | 1.79 | 1.76 | 1.70 | 1.66 | 1.63 |
| 40 | 4.08 | 3.23 | 2.84 | 2.61 | 2.45 | 2.34 | 2.25 | 2.18 | 2.12 | 2.08 | 2.00 | 1.91 | 1.84 | 1.78 | 1.74 | 1.69 | 1.66 | 1.60 | 1.57 | 1.54 |
| 50 | 4.03 | 3.18 | 2.79 | 2.56 | 2.40 | 2.29 | 2.20 | 2.13 | 2.07 | 2.03 | 1.95 | 1.87 | 1.78 | 1.73 | 1.69 | 1.63 | 1.60 | 1.52 | 1.48 | 1.45 |
| 100 | 3.94 | 3.09 | 2.70 | 2.46 | 2.31 | 2.19 | 2.10 | 2.03 | 1.97 | 1.93 | 1.85 | 1.77 | 1.68 | 1.62 | 1.57 | 1.52 | 1.48 | 1.39 | 1.34 | 1.30 |
| 200 | 3.89 | 3.04 | 2.65 | 2.42 | 2.26 | 2.14 | 2.06 | 1.98 | 1.93 | 1.88 | 1.80 | 1.72 | 1.62 | 1.56 | 1.52 | 1.46 | 1.41 | 1.32 | 1.26 | 1.21 |
| 1000 | 3.85 | 3.00 | 2.61 | 2.38 | 2.22 | 2.11 | 2.02 | 1.95 | 1.89 | 1.84 | 1.76 | 1.68 | 1.58 | 1.52 | 1.47 | 1.41 | 1.36 | 1.26 | 1.19 | 1.11 |

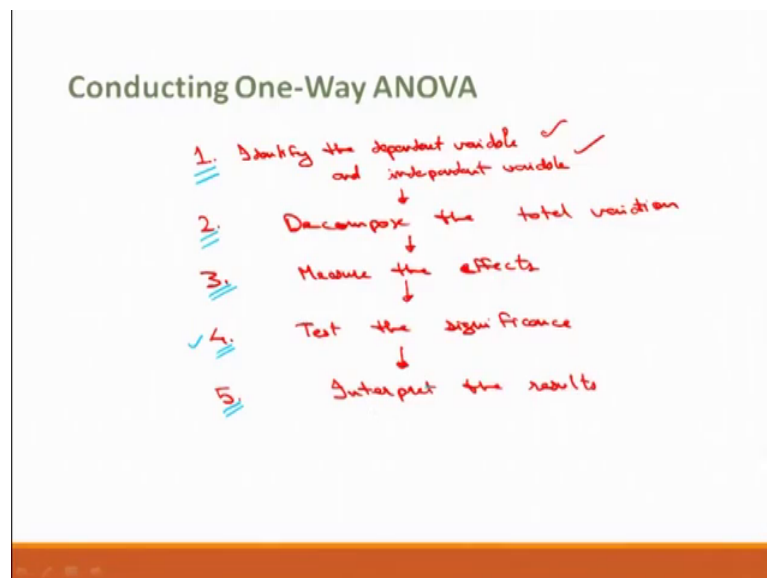
$F \approx 3.35$

So, for two degrees of freedom we will select this column and for 27 degrees of freedom we have 25 and 30. We select these columns for 25 and 30 and we can extrapolate the value for 27 here. So, this value is 39 3.39 and 3.32 the in between value if we say that value comes down to F value, it is tabulated is 3.35. So, we found that the tabulated value that is F 2 27 value for 0.05 alpha is equal to 3.35.

So, then we see that this value is far less than 17.94. So, if I say this is my let me suppose my F distribution here would be this something like this and this is my alpha this is alpha is equal to 0.05 and this alpha value here is 3.35 and my F value the calculated F value is something very far that is 17.94 that is far from my significance level here. So, this is in the rejection region now we can conclude that we reject null hypothesis. Null hypothesis would be that no effect of in store promotion. So, this hypothesis is rejected and alternative hypothesis that is reject H naught.

So, this holds true; that means, we can say that the calculated value because is greater than the critical value the tabulated value therefore, the in store promotion has significant effect on these sales. So, this is the next step we have tested in null hypothesis now we move to the final step that is interpretation of the results.

(Refer Slide Time: 38:08)



So, this is the significance level testing that is we rejected the null hypothesis and the interpretation of results is the same which I just said the in store promotion has significant effect on the sales here. So, we have completed the 5 steps identify, dependent

and independent variables, then we decomposed the total variation then we calculated the statistic value, then we test the null hypothesis. After null hypothesis tested we see that is the null hypothesis rejected or not. It is rejected, then we have interpretation of results that is our independent variable has significant effect on our dependent variable. So, this was our analysis of variance the mechanism here.

(Refer Slide Time: 39:04)

Conducting One-Way ANOVA

One-Way ANOVA: Effect of In-store Promotion on Store Sales

$$df = (C-1) + (N-C) = N-1$$

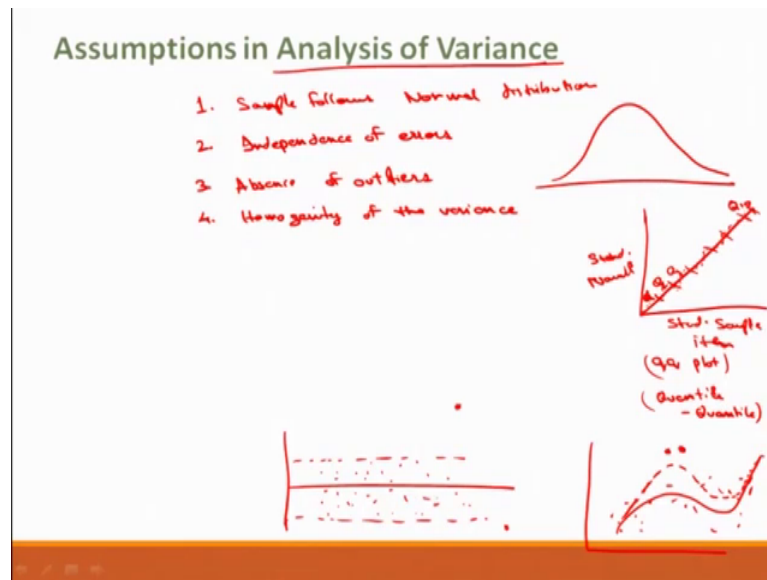
| Source of prob. Variation | Sum of squares | df | Mean square | F ratio | F |
|----------------------------|----------------|------------|-------------|---------|-------|
| Between groups (Promotion) | → 106.067 | → 2 | 53.033 | 17.944 | 0.000 |
| Within groups (Error) | → 79.800 | → 27 | 2.956 | | |
| TOTAL | → 185.867 | → 29 (N-1) | 6.409 | 17.944 | 3.35 |

| Level of Promotion | Count | Mean |
|--------------------|-------|-------|
| High (1) | 10 | 8.300 |
| Medium (2) | 10 | 6.200 |
| Low (3) | 10 | 3.700 |
| TOTAL | 30 | 6.067 |

So, when you will do this in spss window or will do this in r package will show the examples of there as well we call it gran r you will see these kinds of windows. So, what do they say you just read the data from the excel or csv file, then we put the syntax for ANOVA and we get these values between groups within groups total variation degrees of freedom for these.

So, this is actually total degrees of freedom. Total degrees of freedom is nothing, but n minus 1 say it is actually this plus this is 2 plus 27. So, it is c minus 1 plus N minus C this is a category. So, if I plus this one degrees of freedom, total this comes down to n minus 1 then we have mean square values and we have F ratio. So, this was one way ANOVA. Effect of in store promotion on this store sales is significant because it tabulated on a critical value was 3 point 3 five. So, this 17.94 is greater than 3.35. So, this cell means also it is giving here.

(Refer Slide Time: 40:39)



So, there are certain assumptions in analysis of variance these assumptions are bound to be met if they are not met then analysis of variance that it does not hold good. So, what are assumptions? Number 1 sample follows normal distribution that is if I plot all the sample items, it should follow this normal distribution how do we test this? We actually do the quartile test we put the standardized normal distribution here and we have standardized sample items. And we see that what is the correlation between them, do the line is close to this 45 degree line if it lies like this one? Yes we can say that they are very much following the normal distribution.

So, this is known as qq plot, qq is nothing, but the quantile quantile plot. This note I have put quantile quantile here because we have we can divided into number of parts here. If we say the quartiles like in the box plot we have 3 quartiles q_1 , q_2 and q_3 we have divided the whole data into four components here and the q_2 is was the median if you remember.

So, in this case we quantile quantile plot we can say each let me say if they are 10 different parts here this is q_1 to q_{10} and equal number of data points are here in each slot here q_2 q_3 . So, this sample follows normal distribution we will see this graphical this visualization, when will see the assumptions of ANOVA in r package we will see r is the assumptions being met or not. So, next is independence of errors that says that the errors between the cases are independent of one another there is no dependence of errors

the errors are totally independent and there is some random chance that errors are coming.

So, next assumption is that absence of outliers; absence of outliers. For instance if this is my data and this is the total variance plot here this could be the outlier these points might the outlier; So, absence of outlier because this outlier would affect my data. So, for example, if this is my model, this is the data, and this is the model here. So, if I have a few outliers here 1 or 2 points which are far behind the model would come like this. Like these outliers would affect the model it will try to attract the model the line to r itself this outliers. So, that would affect now this point I was mentioning here this is the variance here.

So, it states that the variance has to be homogenous number four is homogeneity of the variance that is population variance in different levels of each independent variable are equal. So, they are equally distributed here.

(Refer Slide Time: 44:57)

N-Way Analysis of Variance

More than one factor

$$SS_{total} = \underbrace{SS_{x_1} + SS_{x_2}}_{\text{Main effects}} + \underbrace{SS_{x_1 x_2} + SS_{error}}_{\text{Interaction effect within category}}$$

Multiple $\eta^2 = \frac{(SS_{x_1} + SS_{x_2} + SS_{x_1 x_2})}{SS_{yy}}$; $\frac{\sum (SS_{x_1} + SS_{x_2} + SS_{x_1 x_2})}{df_n}$

$\omega^2 = \frac{MS_{x_1, x_2, x_1 x_2}}{MS_{error}}$

$$df_n = (c_1 - 1) + (c_2 - 1) + (c_1 - 1)(c_2 - 1)$$

$$= c_1 c_2 - 1$$

$$df_d = N - c_1 c_2$$

So, next comes the n way analysis of variance; n way analysis of variance is nothing, but we have more than ONE factor like we had only ONE factor that is in store promotion in the previous example. in n way analysis of variance it is often concerned with the effect of more than ONE factor simultaneously for example, the advertising levels and price levels both how does it depend upon the brand sales for example, they are the we divide the advertisement the level of advertisement into 3 categories high medium low, it

divides the price of product in to 3 categories ONE high medium and low and there is a 2 categories here, that is a 2 factors we can say and what is the effect of these two factors on this sales.

In this case 2 way analysis of variance would happen another example I can take here I can extend the example which I took before the familiarity with the NPTEL courses and NPTEL portal and number of students who register if I say this was the 1 factor familiarity high medium low was 1 factor the number of courses taken the number of students in the courses is the measure dependent variable here; If I take another factor that familiarity with the NPTEL courses second factor if I take reputation of this specific institution specific IITs. IIT Kanpur, IIT Madras, IIT Kharagpur, IIT Delhi these are 2 factors reputation, high and low depending upon the quality of the courses which are delivered by the specific institutions.

So, we have the quality of the courses depending upon the institution and the familiarity with the NPTEL portal, the NPTEL courses are there. These two factors if are there in again we have 2 way ANOVA. So, I will try to continue my previous example here and would calculate the SS total as SS due to x 1 and SS due to x 2 and SS due to x 1 and x 2 combined this is known as interaction effect, these are known as main effects.

So, we had 2 factors here coupon level and in store promotion. We did first in store promotion and sales. Now will do the in store promotion and coupon level on sales they have coupon level 1 that 2, 2 levels of categories here they have 3 levels of categories here in store promotion, how does this effect our sale we will see this. So, not to 4, but here this is actually the perfect model if we do not have the error here we need to put SS error, that is SS within category.

So, in this case we have the overall effect that is known as multiple eta square then we have omega square as well here that is the individual effects here. So, this eta square for if I said 2 ANOVA x 1 and x 2 here, this eta square is calculated as $SS_{x_1} + SS_{x_2} + SS_{x_1, x_2}$ by SS_{total} or SS_y . And in this case F statistic would be $\frac{SS_{x_1} + SS_{x_2} + SS_{x_1, x_2}}{SS_{error}}$ by the degrees of freedom for n this is F statistic here over this is mean square for x 1 x 2 and x 1 x 2 combined this over the SS error by degrees of freedom for error. So, this is nothing, but mean square for x 1 comma x 2 comma x 1 x 2 by mean square for error.

So, degrees of freedom for numerator here would be, the degrees of freedom for category 1 plus degrees of freedom for category 2 plus combined degrees of freedom for category 1 and category 2 category 1 category 2 combined category 1 category 2 that comes down to $C_1, C_2 - 1$. And degrees of freedom for denominator for this F test would be degrees of freedom for total all the values minus C_1, C_2 ok. So, we can calculate this F statistic to see the overall effect we can even see the F value for the individual effects.

(Refer Slide Time: 50:40)

N-Way Analysis of Variance

$$F_{x_1} = \frac{S_{x_1} / df_n}{SS_{error} / df_d}$$

$$F_{x_1 \times x_2} = \frac{S_{x_1 \times x_2} / df_n}{SS_{error} / df_d}$$

$$= \frac{M S_{x_1 \times x_2}}{M S_{error}}$$

$df_n = C_1 - 1$
 $df_d = N - C_1 C_2$

$df_n = (C_1 - 1)(C_2 - 1)$
 $df_d = N - C_1 C_2$

For example interaction effect $F_{1 \text{ and } x_1 \times x_2}$ that would be $S_{x_1 \times x_2}$ by degrees of freedom for a numerator here by SS_{error} by degrees of freedom for denominator. So, this is $MS_{x_1 \times x_2}$ by MS_{error} . So, for this is this was for the interaction effect for the main effect let me say for x_1 , this would be again S_{x_1} by degrees of freedom for the specific x_1 and SS_{error} by degrees of freedom for the denominator.

In this case the degrees of freedom for $x_1 \times x_2$ is degrees of freedom for numerator here is $C_1 - 1$ into $C_2 - 1$ and degrees of freedom for denominator here is $N - C_1 C_2$. And in this case degrees of freedom for numerator would be for this category 1 $C_1 - 1$ and degrees of freedom for denominator here would be $N - C_1 C_2$.

(Refer Slide Time: 52:14)

Conducting Two-Way ANOVA

| Source of Variation | Sum of squares | df | Mean square | F | Sig. of F |
|---------------------|----------------|----|-------------|--------|-----------|
| Main Effects | | | | | |
| Promotion | 106.067 | 2 | 53.033 | 54.862 | 0.000 |
| Coupon | 53.333 | 1 | 53.333 | 55.172 | 0.000 |
| Combined | 159.400 | 3 | 53.133 | 54.966 | 0.000 |
| Two-way interaction | 3.267 | 2 | 1.633 | 1.690 | 0.226 |
| Model | 162.667 | 5 | 32.533 | 33.655 | 0.000 |
| Residual (error) | 23.200 | 24 | 0.967 | | |
| TOTAL | 185.867 | 29 | 6.409 | | |

Handwritten notes:

- Interpretation:**
 - Promotion and Coupon level significantly affect sales.
 - Coupon level is having larger influence on sales.
 - Interaction effect is not significant.
- Significance levels:**
 - 0.000 < 0.05 (Highly significant)
 - 0.226 > 0.05 (Not significant)
- Model:** 0.000 < 0.05 (Highly significant)

Handwritten formula:

$$SS_{\text{total}} = (df_{\text{total}} \times MS_{\text{error}}) + (df_{\text{model}} \times MS_{\text{model}})$$

Handwritten notes for formula:

- $df_{\text{total}} = 29$
- $df_{\text{model}} = 5$
- $df_{\text{error}} = 24$

So, if I consider the two categories here, this kind of table you can do the all the calculations if you want you can do that in excel or may be manually you can do, but I will show you the table that is an SPSS output here. So, we have the promotion the coupon and the combined this is category 1 this is category 2 and this is actually x 1 x 2 and this is x 1 x 2 combined.

So, we had for promotion this is the sum of squares and for coupon level if you see this is the sum of square we had three categories this is 3 minus 1 is equal to 2, we had two categories 2 minus 1 is equal to 1, total is 2 plus 1 is 3 and combined is nothing, but sum of this, this is sum of these 2 interaction we have c this in to this 2 in to 1 is equal to this 2, 2 into 1 is equal to 2.

So, what we can see we have the F values for 2 factors and combined F value and we also have the 2 way interaction and if total for the model here. So, we can see here that the F value as in the significance of F the probability of F is less than 0 less than 0.05 here this one, this one and this one. These two that is the promotion and coupon level both are significant the interpretation would be number 1 promotion it is the in store promotion and coupon level significantly affect the sales this is number 1.

Number 2 we can see that which has more significant, which has more influence on this in that case again F is the criteria that would decide. In this case coupon level is having larger value, but these 2 values are close 54 and 55 are close, we can say that coupon

level is having higher influence on sales. If I suppose that this for promotion, they are value let me say it would have been some value like 14, couple level it would identify.

So, we can check this significance also is the difference between these 2 values significant or not. Then that third interpretation here is 1 2 and 3; third interpretation here is that the interaction this is the interaction, this value is greater than 0.05 therefore, no interaction effect. So, the third interpretation I will put here, interaction effect is not significant.

So, these are the interpretations or the results of our ANOVA table also we have omega square here. Now what is omega square here? Omega square indicates what portion of variation in the dependent variable is related to a particular independent variable, this has the proportion ok. So, this is calculated only for this significant variables here, because this is less than 0.05. So, this is calculated omega is calculated omega square is equal to the SS_x minus degrees of freedom for x in to MS_{error} by SS_{total} plus MS_{error} .

So, if in this case if I calculate the value of omega square for promotion and for coupon the values comes from to 0.557 and 0.280 respectively. So, as a guide to interpret this term, we can say that in large experiment effect produces index of 0.15 or greater. And a medium effect produces an index around 0.06 and a small effect produces an index of 0.01.

So, this is large, this is medium and this is small, we are talking about the omega square here ok. If it is less than 0.01 it is small, if it is between 0.06 is less than omega square is less than 0.15. So, in the present example 0.557 and 0.280 are both greater than 0.15. Therefore, it has large effect, large better to put large experimental effect. So, this was analysis of variance for 2 factors or we call it 2 way analysis of variance. Similarly we can conduct n way analysis of variance and multiple factors could be identified we can have the list of factors, but please mind it.

The more we add, the more factors we have more is the value of R square that is inflated. So, value of R square would be higher four number of factors, in that case R square adjusted value has to be seen to see the effect of the model and also we have R square predicted, and also we have R square we have R square predicted, R square adjusted and the normal r square.

(Refer Slide Time: 59:22)

Conducting Two-Way ANOVA

R^2 - predicted \leftarrow
 R^2 - adjusted \leftarrow
 R^2 -

So, this R square looks good but they can be serious problem with an over fit of the model. For one thing the regression coefficients represent the noise that rather than the general relationship of the population, that is the un accounted factors could effect this one.

So, this adjusted R square can compensate that. More over the predicted R square value is designed to detect the over fitting as well. So, next thing is analysis of covariance.

(Refer Slide Time: 60:06)

Conducting One-Way ANOVA

Illustrative example

| Store Number | Coupon Label | In-Store Promotion | Sales | Clientel Rating |
|--------------|--------------|--------------------|-------|-----------------|
| 1 | 1.00 | 1.00 | 10.00 | 9.00 |
| 2 | 1.00 | 1.00 | 9.00 | 10.00 |
| 3 | 1.00 | 1.00 | 10.00 | 8.00 |
| 4 | 1.00 | 1.00 | 8.00 | 4.00 |
| 5 | 1.00 | 1.00 | 9.00 | 6.00 |
| 6 | 1.00 | 2.00 | 8.00 | 8.00 |
| 7 | 1.00 | 2.00 | 8.00 | 4.00 |
| 8 | 1.00 | 2.00 | 7.00 | 10.00 |
| 9 | 1.00 | 2.00 | 9.00 | 6.00 |
| 10 | 1.00 | 2.00 | 6.00 | 9.00 |
| 11 | 1.00 | 3.00 | 5.00 | 8.00 |
| 12 | 1.00 | 3.00 | 7.00 | 9.00 |
| 13 | 1.00 | 3.00 | 6.00 | 6.00 |
| 14 | 1.00 | 3.00 | 4.00 | 10.00 |
| 15 | 1.00 | 3.00 | 5.00 | 4.00 |
| 16 | 2.00 | 1.00 | 8.00 | 10.00 |
| 17 | 2.00 | 1.00 | 9.00 | 6.00 |
| 18 | 2.00 | 1.00 | 7.00 | 8.00 |
| 19 | 2.00 | 1.00 | 7.00 | 4.00 |
| 20 | 2.00 | 1.00 | 6.00 | 9.00 |
| 21 | 2.00 | 2.00 | 4.00 | 6.00 |
| 22 | 2.00 | 2.00 | 5.00 | 8.00 |
| 23 | 2.00 | 2.00 | 5.00 | 10.00 |
| 24 | 2.00 | 2.00 | 6.00 | 4.00 |
| 25 | 2.00 | 2.00 | 4.00 | 9.00 |
| 26 | 2.00 | 3.00 | 2.00 | 4.00 |
| 27 | 2.00 | 3.00 | 3.00 | 6.00 |
| 28 | 2.00 | 3.00 | 2.00 | 10.00 |
| 29 | 2.00 | 3.00 | 1.00 | 9.00 |
| 30 | 2.00 | 3.00 | 2.00 | 8.00 |

So, let us see this example unit we have this two factors coupon level and in store promotion these were the factors categories, and there is another variable that is cliental rating what is the rating of the client and this is in interval scale. And then as I mentioned earlier if we have the interval scale as well as categorical scale, the analysis of variance is then called as analysis of covariance.

So, when examining the mean values of the dependent variable related to the effect of the controlled independent variable, it is often necessary to take into an account the influence of the uncontrolled independence variable. For example, in determining how different groups exposed to different commercials evaluative brands it may be necessary to control the prior knowledge of the customer to of the groups here.

In determining how the price level will effect, a house hold consumption it may be sensual to take the how house hold size into account. For instance in this case the cliental rating if the client is having high rating, the purchasing by a this client might be higher this might effect. So, in this case we conduct analysis of co variance.

(Refer Slide Time: 61:27)

Analysis of Covariance

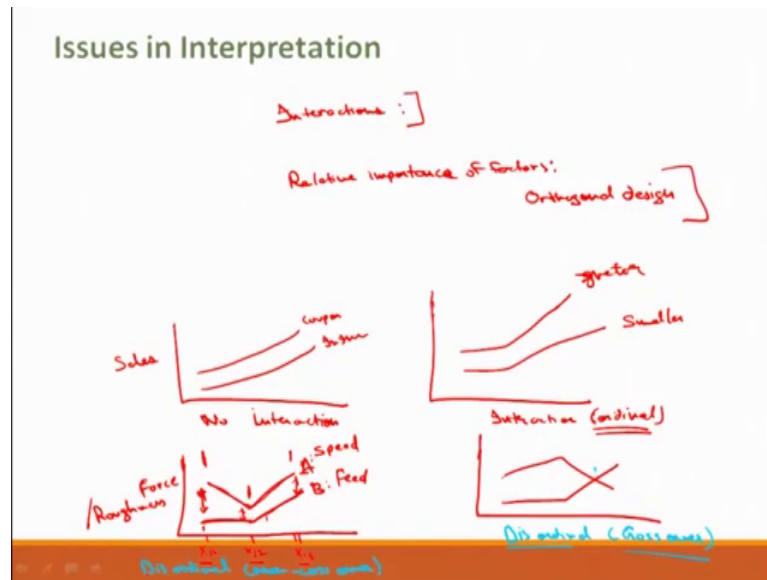
| Source of Variation | Sum of Squares | df | Mean Square | F | Sig. of F |
|--------------------------|---------------------------|-----------|--------------|--------|-----------|
| Covariance | | | | | |
| Clientele | 0.838 | 1 | 0.838 | 0.862 | 0.363 |
| Main effects | | | | | |
| Promotion | 106.067 | 2 | 53.033 | 54.546 | 0.000 |
| Coupon | 53.333 | 1 | 53.333 | 54.855 | 0.000 |
| Combined | 159.400 | 3 | 53.133 | 54.649 | 0.000 |
| 2-Way Interaction | | | | | |
| Promotion* Coupon | 3.267 | 2 | 1.633 | 1.680 | 0.208 |
| Model | 163.505 | 6 | 27.251 | 28.028 | 0.000 |
| Residual (Error) | 22.362 | 23 | 0.972 | | |
| TOTAL | 185.867 | 29 | 6.409 | | |
| Covariate | | | | | |
| Clientele | Raw Coefficient -0.078 | | | | |

So, this is the table here, these are the main effects here, this is the covariance, this is actually the interval scale. So, these are in my categorical scale the factor factors here.

So, in this case we have again the F statistic value. So, we can do the similar interpretations from this table as well. In this case the cliental rating is not significant,

this value is greater than 0.05; that means, that cliental rating does not have or does not influence my sales here. So, next there are certain issues interpretation, there are interactions and there are relative importance of the factors.

(Refer Slide Time: 62:18)



First I will take the relative importance of factors experimental design are usually balanced. In that each cell contains the same number of respondents, as I had 10 stores in each cell here, but sometimes this is not possible. So, this kind of design when each cell is having same number of respondents, that is known as orthogonal design. So, it is possible to determine the relative importance of each factor in explaining the variation independent variable in this case.

So, interactions are when the 2 variables are interacting each other. So, for instant if I plot this here, let me say sales and the 2 variables here my coupon and in store promotion, this is not interacting with each other; like if with increase in this also increasing there is no interaction.

We can have some interaction in certain cases for instance, if it is like this one. In this case we have a interaction that is ordinal, it has this one greater this one is smaller. That is the lines are not exactly parallel here. But we can see that is the interactions significant or not. All though interaction is there, but significance also needs to be tested, then we have dis ordinal interaction also. This is the interaction for instance if its like this one. See if this is increasing if I say they are 3 values x 1 1, x 1 2 and x 1 3 this 1 and I could

put A and B, this is increasing at this point and again increasing the B. B is increasing and again the slope is higher and it is increasing A is first decreasing and then increasing. So, there is a difference.

So, this is large, this distance is large, this distance is small again this distance is large. So, there is an interaction. So, these things for example, if I do machining, if you know the mechanical people do machining they cut the material. They change the tool cutting speed, they change the tool speed cutting speed is the for instance the rotation of the work piece, which is being cut here it is if the work piece is rotating it is cut by the tool.

The feed is the rate of cutting here of this tool this can have a transient behavior for instance a certain at lower speeds at lower speeds, the cutting; that means, am trying to see the cutting force here or may be roughness of the work piece and you response I can choose here, speed feed also I have put A as speed and B as feed.

So, here I can see the feed is increasing from level x_1^1 to level x_1^2 , but increase from level x_1^2 and x_1^3 is higher this slope is higher; however, at lower speeds this is low speed this is high speed and lower speed the force is higher, at intermediate speed force is lower and at high speed the force is again higher. So, this has an interaction the previous. So, the previous one was ordinal.

So, this is dis ordinal dis ordinal, but non cross over. Sometimes we can have the case when if I say this same case for different machines or for different material for different tool combination tool work combination or for different environment, it can behave something like this as well. So, this is again dis ordinal, but we have cross over. So, this is dis ordinal and also we have cross over here, it is crossing. So, this transient behaviors can be like this as well.

(Refer Slide Time: 67:58)

The slide features the title "Nonmetric Analysis of Variance" in green text at the top left. To its right, the handwritten text "Y: is ordinal" is written in red. Below the title, the handwritten text "[K-sample median test]" is written in red. The slide has a white background and a blue footer bar.

The for the non-metric analysis of variance there are certain tests like K median test, this is actually known as K sample median test. So, this is non metric analysis of variance this examiner difference in the central tendencies of more than 2 groups, when the dependent variable is measured on an ordinal scale in this case Y or dependent variable is ordinal. So, it is non metric.

So, in this cases we can have these kinds of test, which are I can say that out of the scope of this course that we can provide notes to you for this.

(Refer Slide Time: 68:47)

The slide features the title "Multivariate Analysis Of Variance" in green text at the top left. Below the title, the handwritten text "MANOVA" is written in red. Underneath that, the handwritten text "More than one dependent variables (correlated)" is written in red. At the bottom, the handwritten text "H0: Vectors of means are equal across groups" is written in red. The slide has a white background and a blue footer bar.

Also we have multivariate analysis of variance that is MANOVA; so, what if MANOVA multivariate analysis of variance is similar to ANOVA, except that we have more than one dependent variable. So, in MANOVA the null hypothesis this is that the vectors of means of people depend with peoples are equal across groups; vectors of means are equal across groups. So, multivariate analysis of variance is appropriate when there are two or more dependent variables that are correlated.

So, next we will meet in the r session for analysis of variance we will take the certain examples, here we will do ANOVA test for 2 factors for 1 factor for 2 factors and for multiple factors and we will see how to interpret results, we will see how the assumptions are being met and what to do if assumptions are not being met here.

Thank you.