**Marketing Research and Analysis-II**
**Prof. Dr Jogendra Kumar Nayak**
**Department of Management Studies**
**Indian Institute of Technology Roorkee**

**Lecture – 42**
**Correlation in SPSS**

Hello friends. Welcome to the class of marketing research and analysis. In the last lecture we had started with the concept of covariance and correlation where I explained you what is basically the utility the use and the meaning of covariance and correlation how correlation is much powerful technique over the covariance the reason behind it and I said that covariance is slightly becomes awkward due to the units that it uses of for 2 variables.

That means 2 variables could be having very different units and that sometimes looks very difficult very absurd to compare. And that is why we need to have another technique which can standardize them and give us a solution so that it is easier for us to compare. So covariance was basically telling us about the direction and not the degree or intensity where as correlation on the other hand tells us both the degree and the direction. So that is the advantage.

So we will proceed with today's class with the correlation. So how do you know study correlation, so basically from a visual point of view.

**(Refer Slide Time: 01:42)**

You can have a Scatter Diagram Method or a Graphic Method or we have a Karl Pearson's Coefficient of correlation of course it and Method of Least Squares. So these are the ways you know used to measure the correlation or check the correlation.

**(Refer Slide Time: 01:52)**



### Scatter Diagram Method

- The **Scatter Diagram Method** is the simplest method to study the correlation between two variables wherein the values for each pair of a variable is plotted on a graph in the form of dots thereby obtaining as many points as the number of observations. Then by looking at the scatter of several points, the degree of correlation is ascertained.
- Scatter diagrams are unable to give you the exact extent of correlation.
- This chart does not show you the relationship for more than two variables.

Now what is the Scatter Diagram? Scatter Diagram Method is the very simple method to study the correlation between 2 variables wherein the values for each pair is plotted on a graph in the form of dots thereby obtaining as many points as the number of observations. Then by looking at the Scatter of the several points the degree of correlation is ascertained that means it looks something like this.
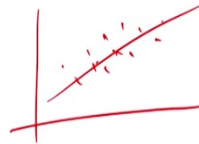
For example, if I plot certain dots for example these are the dots. So dots are each my observations. So, from here let us say this is my X and this is my Y. So can I say what kind of correlationship exists? So if I draw a line in between so I can say well there is some positive correlation it is increasing it seems. But suppose I have another relationship something like this.

So here if I draw what is a correlation ship it is very difficult to ascertain because it is in fact that I can hardly say there is any relationship? So scattered diagrams help you to visually understand that but there are unable to give you the exact extent of correlation how much the which is the most important part of correlation is the how much that; how much is the what you do not get here. This chart does not show you the relationship for more than 2 variables only 2 variables x and y but not if there are 3 variables multiple correlation you cannot do it.

**(Refer Slide Time: 03:18)**

Scatter Diagram Method

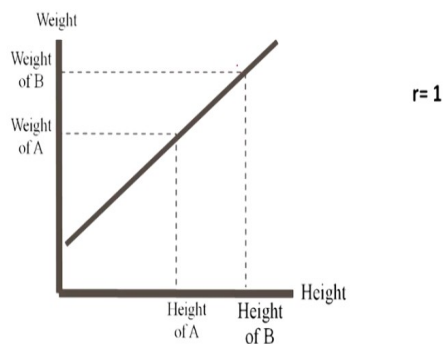The degree to which the variables are related to each other depends on the manner in which the points are scattered over the chart. Yes, that is true. How they are scattered, is there any pattern or not that decides. The more the points plotted are scattered over the chart, the less is the degree of correlation that means is the spread is too large and it is in a systematic pattern then the correlation is weak the more the points plotted are closer to the line that line which I said suppose this is the line right the higher is the degree of the correlation so the closer they are to the line the higher is the degree of correlation between the 2 variables it is denoted by r.
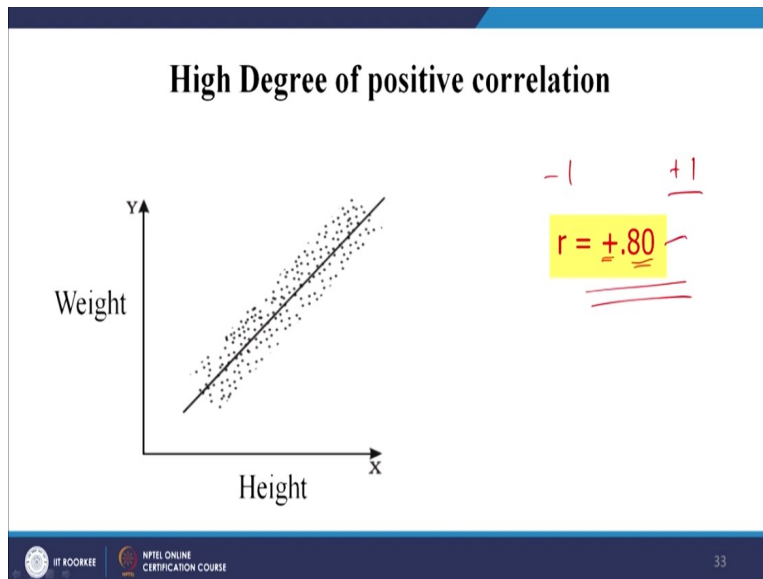
**(Refer Slide Time: 04:00)**



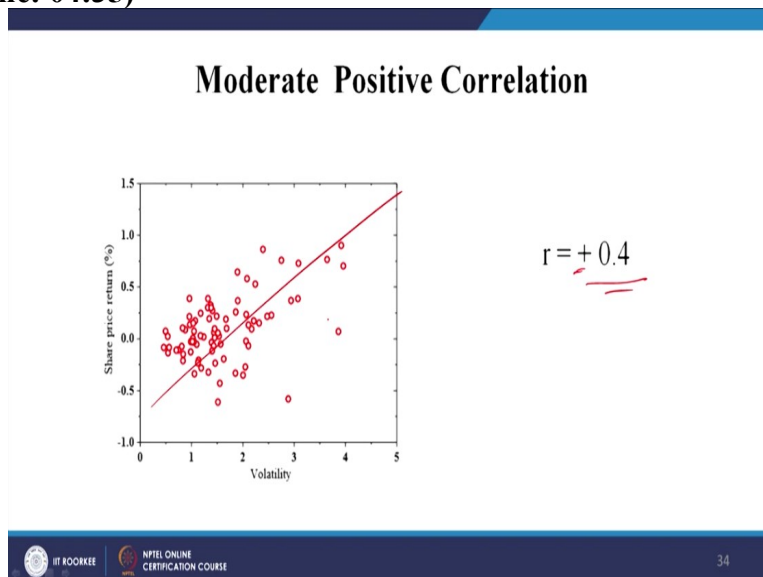So this is then you see the perfect correlation positive correlation where r = 1 so r =1 weight and we are taken height so when we are saying it is a perfect positive correlation.
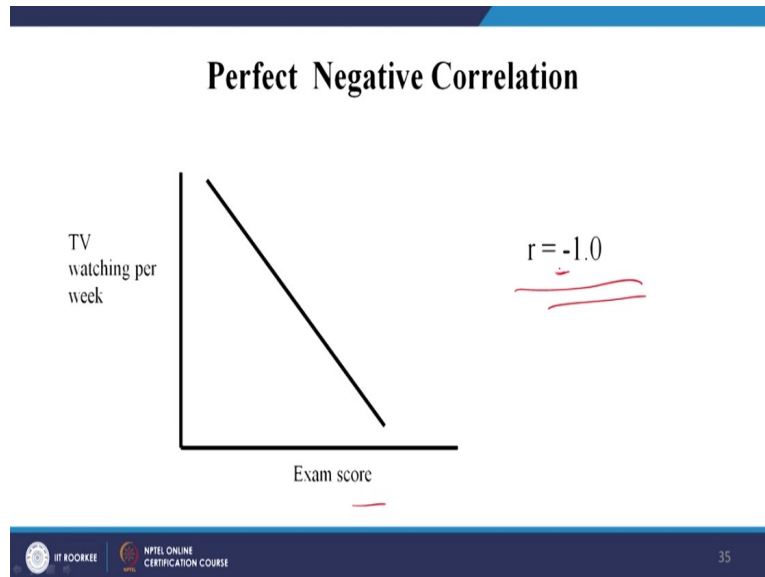
**(Refer Slide Time: 04:13)**

High Degree of positive correlation

When r = .8 when we are measured r we will see how to measure r. But when height and weight we are saying r = .8. So we say it is a substantially strong positive correlation it is a positive because of the sign and because of the value lies between -1 to +1. So it is close to +1 so highly positive.

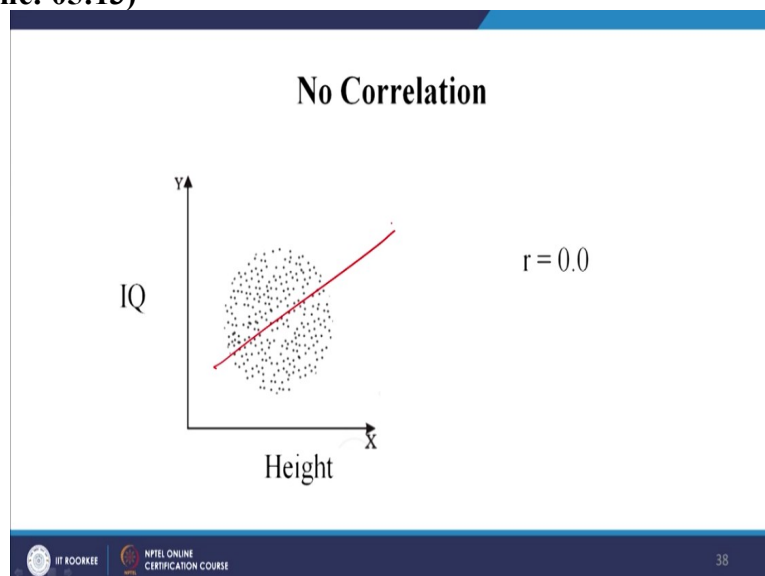**(Refer Slide Time: 04:35)**



Moderate Positive Correlation

Moderately positive now, so this is somewhere sign is positive and 0.4. If I draw a line you see. So this looks little awkward but still it is a positive correlation.

**(Refer Slide Time: 04:47)**

Perfect Negative Correlation

$r = -1.0$

TV watching per week

Exam score

Now this is a Perfect Negative you see TV watching per week and scores in the exam. So, my relationship is negative as TV watching is high. My marks are low. As a TV watching goes down my marks are increasing a Perfect Correlation. But the direction is inverse Negative relationship that means if one increases the other decreases, moderately negative. Moderately negative similarly, so -0.8. No correlation.

**(Refer Slide Time: 05:13)**



No Correlation

$r = 0.0$

IQ

Height

So here you see IQ and height. Is there any relationship between IQ and height .Well if I draw what is the correlation coming nothing.

**(Refer Slide Time: 05:21)**

So the best method of measuring statistically is called the Karl Pearson's method of Coefficient of correlation. So this is also called as product moment correlation. You must have seen somewhere written as simple correlation or product moment correlation. It gives a precise numerical value of the degree of linear relationship between 2 variables X and Y. So what is the degree of the relationship and this relationship may be given by a regression equation which is Y = a + b of X, where Y is my intercept, B is my slope and X is my independent variable.

**(Refer Slide Time: 05:56)**



**(Refer Slide Time: 06:00)**

**Karl Pearson's Coefficient of Correlation**

$$r = r_{xy} = \frac{Cov(x,y)}{S_x \times S_y}$$

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2 \sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

As I said it lies between -1 to +1. So this is the formula, so how do you measure the Karl Pearson Coefficient of Correlation. The formula is very simple. You need to understand. So what we said correlation is a standardized form of the covariance. So if my covariance value, if I divide my covariance value between x and y with the standard deviation of x and standard deviation y. This is what it becomes as the; my coefficient of correlation.

$$= \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\Sigma(x_i - \bar{x})^2} \sqrt{\Sigma(y_i - \bar{y})^2}}$$

The N -1 here and the N-1 here actually are cutting each other. So that is why we are not showing it is not required anymore. So this is the way you write. This is a very simple thing. You must remember this much and nothing you do not have to remember any other thing. This is another way of representing the same thing.

**(Refer Slide Time: 07:16)**



**Interpretation of Correlation Coefficient (r)**

- The value of correlation coefficient 'r' ranges from -1 to +1.

- If r = +1, then the correlation between the two variables is said to be perfect and positive.

- If r = -1, then the correlation between the two variables is said to be perfect and negative.

- If r = 0, then there exists no correlation between the variables.

So let us see, so interpretation of the correlation coefficient. The value this is said, if it is +1 perfect and positive, if it is - 1 perfect and negative. If r = 0, no correlation.

**(Refer Slide Time: 07:29)**



Some assumptions in correlation what are the assumptions? The variables should be measured at the interval or ratio level. First they should be continuous for example time revision time, intelligence, exam performance, weight.
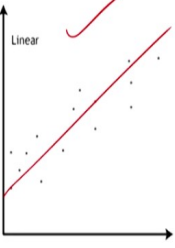
**(Refer Slide Time: 07:46)**



There is a linear relationship between the 2. The assumption this is an assumption in correlation that the relationship is linear. Well there are a number of ways to check a linear relationship creating a scatter plot I guess each other and then visually inspecting is the simplest way. So you see if I draw a line you see I draw a line and here also I draw a line. So, these 2 are nonlinear but this is linear. If the relationship displayed is not linear. You will have to either run in non-parametric equivalent to the Person's correlation, what is the non

parametric equivalent? I have explained it earlier also when I was teaching non parametric techniques.

In which I said there is a technique called Spearman's correlation which is used as a non-parametric test for the correlation. So or you have to transform your data and convert it from non normal into a normally distributed data. If it is non normal you have to convert into normally distributed data. If data is normally distributed the chances are very fair that it will become following a linear pattern.

**(Refer Slide Time: 09:04)**



The third assumption is no significant outlier should be there. If there is an outlier it will distort the entire situation.

**(Refer Slide Time: 09:06)**

Fourth the variables should be approximately normally distributed in order to assess the statistical significance of the correlation. Now what does it mean. If you do not have understanding whether it is statistically significant or not between the relationships between 2 variables, then although you get a degree you cannot say it will happen every time the same way.

Maybe it has happened this time by chance but maybe another time it will not give the same result, so but if it is significant that means it is not by chance it will happen again and again and again in the same way. So but this method is difficult to assess. So, simpler common method is used. This method involves determining the normality of each variable separately. So what it says is you do not check it at the same time but individually you check whether the data or the variables are normally distributed or not.

**(Refer Slide Time: 10:02)**



What are the limitations of Pearson's coefficient? Does it have any limitation? Yes. First of all, assuming in linear relationship in life I said earlier key most of the relationship may not be linear they could be non-linear. Interpreting the value of r is sometimes difficult although we say if it is +1 it is extremely positive -1 extremely negative but it is still interpretation is little challenging.

There effected by the extreme values outliers and there it is time consuming to calculate and understand because in real life there is nothing comes in a plate to you. So suppose the situation is given and then you have to interpret what are the variables and then find out what the relationship is not that simple as it looks when I give you on a pen and paper.

**(Refer Slide Time: 10:51)**

**Example 1 (Pearson correlation)**

Find correlation coefficient between the sales and expenses from the data given below:

| Firm | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|------|---|---|---|---|---|---|---|---|---|----|
| Sales ($ Millions) | 50 | 50 | 55 | 60 | 65 | 65 | 65 | 60 | 60 | 50 |
| Expenses ($ Millions) | 11 | 13 | 14 | 16 | 16 | 15 | 15 | 14 | 13 | 13 |

Now let us solve a problem find the correlation coefficient between the sales and expenses from the data given below the sales is given 50 number of maybe you know different Firms are there. Firm 1 to Firm 10 and the sales is given 50, 50, 55, 60, 65 and goes on. This is all in Millions. The expenses conducted by this or expanded by these companies are 11, 13, 14, 16 to something.

Now I am interested to know what is the direction of the relationship. Is it a positive relationship or negative relationship or there is no relationship? And what is the degree of intensity of the relationship how strong is the relationship. So to do this we if we cannot we will not do covariance, we will do a correlation in this case and luckily these are having same units also. So let us say.

**(Refer Slide Time: 11:43)**

**Calculation of correlation coefficient**

| Firm | Sales (X) | Expenses (Y) | $(X-\bar{X})$ | $(Y-\bar{Y})$ | $(X-\bar{X})^2$ | $(Y-\bar{Y})^2$ | $(X-\bar{X})(Y-\bar{Y})$ |
|------|-----------|--------------|---------------|---------------|-----------------|-----------------|--------------------------|
| 1 | 50 | 11 | -8 | -3 | 64 | 9 | 24 |
| 2 | 50 | 13 | -8 | -1 | 64 | 1 | +8 |
| 3 | 55 | 14 | -3 | 0 | 9 | 0 | 0 |
| 4 | 60 | 16 | +2 | +2 | 4 | 4 | +4 |
| 5 | 65 | 16 | +7 | +2 | 49 | 4 | +14 |
| 6 | 65 | 15 | +7 | +1 | 49 | 1 | +7 |
| 7 | 65 | 15 | +7 | +1 | 49 | 1 | +7 |
| 8 | 60 | 14 | +2 | 0 | 4 | 0 | 0 |
| 9 | 60 | 13 | +2 | -1 | 4 | 1 | -2 |
| 10 | 50 | 13 | -8 | -1 | 64 | 1 | +8 |
| N=10 | $\sum X=580$ | $\sum Y=140$ | $\sum(X-\bar{X})=0$ | $\sum(Y-\bar{Y})=0$ | $\sum(X-\bar{X})^2=360$ | $\sum(Y-\bar{Y})^2=22$ | $\sum(X-\bar{X})(Y-\bar{Y})=70$ |

$$\bar{X}=\sum X/N=580/10=58 \quad \text{and} \quad \bar{Y}=\sum Y/N=140/10=14$$

How do we go about it. So these are my firms. So the 10 Firms sales is given. So my summation of the sales is 580, summation of my expenses 140. Now I need $(X - \bar{X})$ . So what is the formula covariance by standard deviation of X into standard deviation of Y. So first we find out $(X - \bar{X})$, so $\bar{X}$ is how much 580 / 10. So this is my $\bar{X}$ , $\bar{y}$ is 14, so $(X - \bar{X})$ - 8 ,– 8, - 3 . It goes on.

So, this totally 0, similarly, if you look at $(Y - \bar{y})$ is again 0 so, because they adjust each other but for that we will take the X - X square now. So, what is the $(X - \bar{X})^2$ now 360 and $(Y - \bar{y})^2$ is 22. Now the product deviation of $(X - \bar{X})* (Y - \bar{y})$, this is equal to if you multiply 8 x3= 24, 8 * 1= 8, - 3 * 0 = 0 so goes on 70.

**(Refer Slide Time: 13:03)**

We know that

$$r = \frac{\{\sum(X-\bar{X})\ \sum(Y-\bar{Y})\}}{\{\sqrt{\sum(X-\bar{X})^2}\ *\sqrt{\sum(Y-\bar{Y})^2}\ \}}$$

$$r = \frac{70}{\{\sqrt{360}\ *\sqrt{22}\}}$$

$$r = 0.787$$

*Hence, there is a high degree of positive correlation between the two variables i.e. as the value of sales goes up, the expenses also go up.

So what is the formula $(X - \bar{X})$, how much it has come let us check. So $(X - \bar{X})* (Y - \bar{y})$ so this is 70 divided by $(X - \bar{X})^2$, so $(X - \bar{X})^2$ is 360 and this is 22. $\sqrt{360} * \sqrt{22}$. So this is equal to 0.787. So, my correlation between the sales and expense is first of all. Is it positive or negative. It is positive. And that means if sales increases expenses; if expenses increase my sales increases. And it is 787 means it is close to 1.

So it is the high. So there is a high degree of positive correlation between the 2 variables that is as this value of sales goes up the expenses also go up.

**(Refer Slide Time: 13:56)**

**Example**

A researcher wants to know whether a person's height is related to how well they perform in a long jump. The researcher recruited untrained individuals from the general population, measured their height and had them perform a long jump. The researcher then investigated whether there was an association between height and long jump performance by running a Pearson's correlation.

Now I will show you how to do this in the SPSS. Now this is a case where a researcher wants to know whether a person's height. This is an example I have taken but I will use a separate example.

**(Refer Slide Time: 14:06)**



**Test Procedure in SPSS Statistics**

The six steps below show you how to analyse your data using Pearson's correlation in SPSS Statistics when none of the four assumptions in the Assumptions section have been violated.

Step-1 Click **Analyze** > **Correlate** > **Bivariate...** on the main menu, as shown below:

In this case just to show you how to what is the steps. So what you do is you go to analyze go to correlate and go to bivariate.

**(Refer Slide Time: 14:16)**

Test Procedure in SPSS Statistics

Step-2: Transfer the Variables Height and Jump_Dist into the Variables box by dragging-and-dropping them

After taking going to bivariate you take the bivariate why it is called bivariate because it is a 2 variables Bi 2, Bicycle 2 wheels. So, height and suppose the Jumping distance in this case it was; in this example. So and what is the tail and you have to take.

**(Refer Slide Time: 14:35)**



Test Procedure in SPSS Statistics

**Step-3** Make sure that the Pearson checkbox is selected under the –Correlation Coefficients– area (although it is selected by default in SPSS Statistics).

**Step-4** Click the options button and you will be presented with the Bivariate Correlations: Options dialogue box. If you wish to generate some descriptive, you can do it by clicking on the relevant checkbox

**Step-5** Click the continue button. You will be returned to the Bivariate Correlations dialogue box.

**Step-6** Click the OK button. This will generate the results of Pearson's correlation.

So then you have to go for this and you have to seek you whether you require the mean and standard deviation are not all these things.

**(Refer Slide Time: 14:43)**

Output for Pearson's correlation

when running the Pearson's correlation procedure, you will be presented with the Correlations table in the IBM SPSS Statistics Output Viewer. The Pearson's correlation result is highlighted below:

In this example, we can see that the Pearson correlation coefficient, r, is 0.706, and that it is statistically significant (p = 0.005).

Correlations

|  |  | Height | Jump_Dist |
|---|---|---|---|
| Height | Pearson Correlation | 1 | .706** |
|  | Sig. (2-tailed) |  | .005 |
|  | N | 14 | 14 |
| Jump_Dist | Pearson Correlation | .706** | 1 |
|  | Sig. (2-tailed) | .005 |  |
|  | N | 14 | 14 |

** Correlation is significant at the 0.01 level (2-tailed).

**(Refer Slide Time: 14:49)**



And let me show you here first let me show you then will try to analyze. So let me take this case, now what I have done is this case is a case where we have taken values of the correlation we want to find out the correlation between body temperature and the heart rate the heartbeat rate. So these are some other values for another 120 people I think. So, 130 people in fact. So if I want to see what is the correlation between is there any association between body temperature and heart rate. I want to check that. So to do that how do I go into SPSS now first I will go to analyze I will go to correlate and then go to Bivariate.
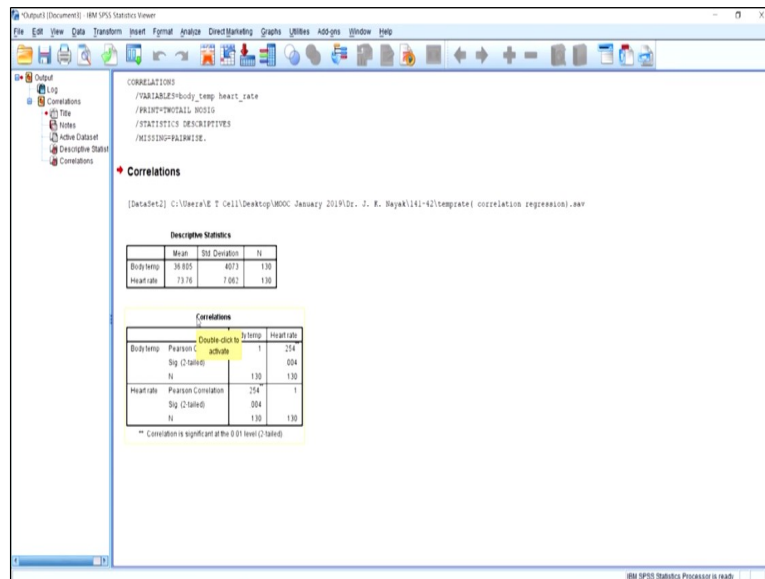
**(Refer Slide Time: 15:29)**

So I want to take my both the variables hot body temperature and heart rate. In options if you go you can go for means and standard deviations if you want and other things remaining the same. Now let us see is there any relationship between the 2, now body temperature with body temperature the correlation will always be 1 is the same thing because and body temperature and heart rate the correlation is .254. But interestingly the important thing is it is significant at what double ** shows it is significant at 0.01 level.

So you can see your .004 which is much less than .01. So the null hypothesis is rejected. What is the null hypothesis in the case of correlation? You should be able to always write the null hypothesis. The null hypothesis is that there is no correlation or there is no relation between the variables. So this null hypothesis is been rejected. So what is an alternate hypothesis that there is a relationship between the variables. The positive or negative that is separate but there is a relationship.
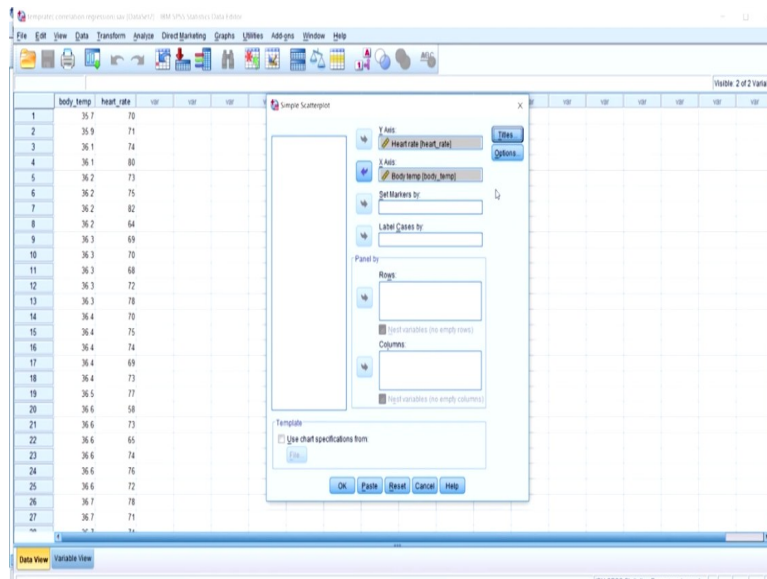
**(Refer Slide Time: 16:40)**

Now look at the body temperature mean and standard division if you want to use this data you can also use this. Now another example we have this is I think you have understood. Now let us take this one. Now somebody wants to see whether there is a relationship between the amount of distance driven and the number of accidents happening by let us say a bus service provider.

Now let us see what is the relationship coming here so I am taking 2 variables. Number of casualties that is happening by let us say particular bus service organization and the number of billions of kilometre they are driving. Now if you look at the output will go to the output and you see the here the correlation between the 2 is .157. But interestingly if you see the significance value it is .253. So what is happening, this .253 is much above the designated value .05 or .01 which we generally take as the significance level value. So the null hypothesis cannot be rejected. So what is our null hypothesis that there is no relationship between accidents and the amount of kilometre driven?

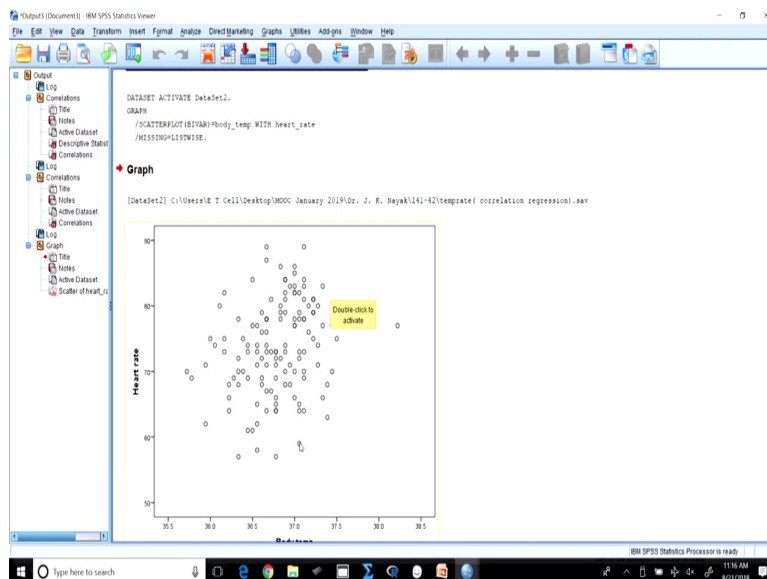And yes that is acceptable. That means you cannot say there is any relationship between the 2. You cannot claim that right. So if I want to show you how this would have looked in the scatter diagram because some people wanted. I have explained scatter diagram so I will explain how you show this in scatter diagram, in this relationship. Now go to this Chart Builder or Legacy and go to scatter plot here.
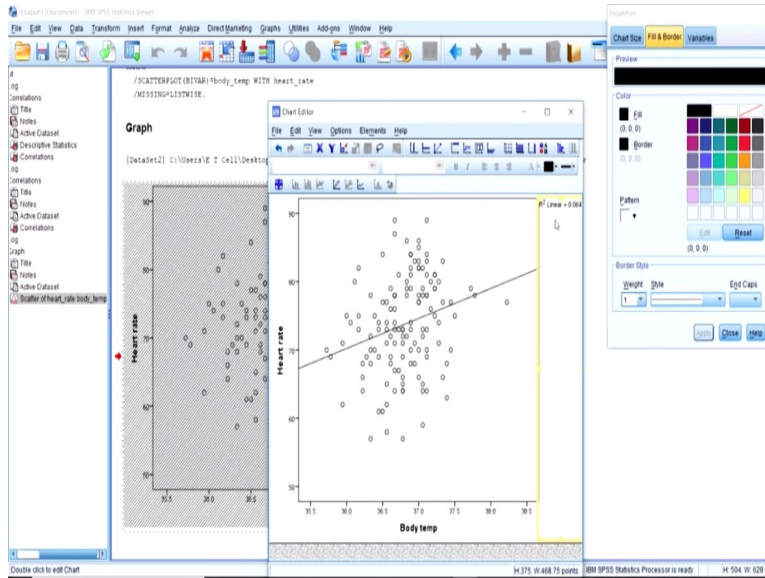
**(Refer Slide Time: 18:29)**

Go to the simple scatter, now take what is your heart rate. Heart rate is my dependent variable so I am taking my heart rate as dependent variable but my body temperature as my X axis my independent variable. If you want to give a title to this you can give it. So suppose some title you want to give you can write here. Now I want to just see it. How does it look.

**(Refer Slide Time: 18:43)**



So let us go to the output file and if you look at it this is how the variables are scattered the variables whatever data as scattered. So my 2 variables are body temperature x axis, heart rate y axis. And this is my how the data is distributed. If you look at this data and you want to understand then you can you may be able to understand but still to make it more clearer, let us do one thing. Make a double click on this.

**(Refer Slide Time: 19:13)**

Now if I make a double click then what I do is I go to this you know fit the total fit. So when I go this I draw a line it draws a line which now helps me to understand the distribution of the data around this line. And if you see here something is written called $r^2$. The r square value is something which I will be explaining you later. It is called the coefficient of determination.

So this $r^2$ is equal to how much 6.4% well. So this tells me well there is some kind of a; if $r^2$ is 6.4. So what is my r? The r you can just say root over of this that means something around. Let us say how much it will be $r^2$ is 6.4 root over of 6.4 is 2.5 or more than that right. So that is you can find out the relationship between r and $r^2$ and thus you say the correlation value is this much. So this is the use of scatter diagram.

**(Refer Slide Time: 20:16)**



So it helps you to give it visual look so output for Pearson's Correlation how it writing running in this case you see so height. Supposing this is just a fictitious example height and

the jumping distance the Correlation was .706 and it is significant. In this example we can see that the Pearson's correlation coefficient is .706 and it is statistically significant at .05 or taken as .005. But usually it is .05. That means 5%.

**(Refer Slide Time: 20:53)**



So that is what it says in our example. What is this coefficient of determination that is what I just said $r^2$. The convenient way of interpreting the value of correlation coefficient is to use the square of coefficient of correlation which is called coefficient of determination. So the coefficient of determination is nothing but this $r^2$. So if r is my 0.9 $r^2$ is 0.81. So similarly if I had a $r^2$ of let us say .064 in my case. So what is my If this is my $r^2$ so what is my r. This is equal to 64 by let us say 1000.

So root over So 8 by let us say how much it 33. So that is somewhere around you are say whatever it comes root over this is this value is equal to 2 point something. So that is what I was saying 2.5 or 2.4 whatever it comes. So this is my r value. So, .24 or .25 whatever it is. Because it cannot be more than – 1 to +1 it has to range in between that.

**(Refer Slide Time: 21:59)**

**Coefficient of Determination**

- The maximum value of $r^2$ is 1 because it is possible to explain all of the variation in y but it is not possible to explain more than all of it.

- Coefficient of Determination = Explained variation / Total variation

Coefficient of determination the maximum value is 1. Because it is possible to explain all of the variation in Y but it is not possible to explain more than all of it. That means what if my r is 1. So r = 1, r square also can be 1 so explained variation coefficient of determination in regression also I will show you. It means it is the ratio of the explained variance by total variance. Now what is this total variance total variance is equal to explained variance plus unexplained variance. So this is my total variance so here I am saying coefficient of determination or

$$r2 = \frac{EV}{EV+UEV}$$

So this is what it seems. Basically theoretically this is the meaning.

**(Refer Slide Time: 22:57)**



**Coefficient of Determination: An example**

- Suppose: r = 0.60

  r = 0.30 It does not mean that the first correlation is twice as strong as the second the 'r' can be understood by computing the value of $r^2$

  When    r = 0.60      $r^2 = 0.36$  -----(1)
  　　　　  r = 0.30      $r^2 = 0.09$  -----(2)

  This implies that in the first case 36% of the total variation is explained whereas in second case 9% of the total variation is explained.

Suppose r = .6, r = .3 it does not mean that the first correlation is twice as strong as the second. The r can be understood by computing the value of $r^2$. Now let us see r is .6. So what is my $r^2$ is .36. Second case if my r is .30 what is my r square .9. So, here it implies that in the first case 36% of the total variation is explained wherein second case 9% of the total variance is explained which is almost 4 times not 2 times.

From a visual here you could have you would have thought it this 2 times but actually it is 4 times the explanation that it is giving you is 4 times. So well what I will do is there is another method. This is for the coefficient of correlation. You have understood.

**(Refer Slide Time: 23:57)**



Thus another technique which is also used for a non-parametric data and this test is called a spearman's correlation which is called Spearman's Rank coefficient of correlation. I have already done it earlier also. But for the benefit of many of the viewers who have not followed that I can just brief you. This Spearman Rank correlation coefficient is nothing but a study in which the data is in a non-parametric nature. So in which the variables are not capable of a quantitative measurement but can be arranged in a serial order.

So they do not follow a normal distribution because they are not in a continuous manner. They are not collected in a continuous manner. So that is one of the assumptions that it should be in ratio an interval scale for a correlation. But here it is not in ratio or an interval scale rather it is in the order that is in a ordinal scale. In such situation the Pearson correlation coefficient cannot be used. And we use this Spearman correlation. So this is the formula

$R = 1-(6\sum D^2)/N(N^2-1)$

## Interpretation of Rank Correlation Coefficient (R)

- The value of rank correlation coefficient, R ranges from -1 to +1.

- If R = +1, then there is complete agreement in the order of the ranks and the ranks are in the same direction.

-

- If R = -1, then there is complete agreement in the order of the ranks and the ranks are in the opposite direction.

-

- If R = 0, then there is no correlation.

So, this other things remaining more or less the same.

## Rank Correlation Coefficient (R)

A) Problems where actual rank are given.

1) Calculate the difference 'D' of two Ranks i.e. (R1 – R2).
2) Square the difference & calculate the sum of the difference i.e. $\sum D^2$
3) Substitute the values obtained in the formula.

So let us see how do you do it. Now calculate the difference D of 2 ranks R1 - R2 and square the difference and calculate that D square.

Substitute the values obtained in the formula that was shown to you. Problems where ranks are not given to you, if the ranks are not given then we need to assign the ranks if ranks are not given just assign ranks. So, by maybe the number the score of somebody else you can assign ranks. The lowest value will be ranked 1 or the highest value can be ranked 1 whatever order you want to take.

We need to follow the same scheme of ranking for the other series also. Suppose there are 2 variables 1 and 2. So if you are using in ascending order for 1 please use ascending order for the other one also. If you are using descending use descending for both.

**(Refer Slide Time: 25:55)**



Then calculate the rank correlation coefficient. So this is the formula.

**(Refer Slide Time: 26:01)**

## Limitation of Spearman's Correlation

- Cannot be used for finding out correlation in a grouped frequency distribution.

- This method should be applied where N exceeds 30.

Let me take it. If I have what the limitation of this, it cannot be used for finding out correlation in a group frequency distribution and this method should be applied where N exceeds 30.

**(Refer Slide Time: 26:11)**



### Example 2 (Spearman correlation)

Two managers are asked to rank a group of employees in order of potential for eventually becoming top managers. The rankings are as follows:

| Employees | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| Ranking by Manager 1 | 10 | 2 | 1 | 4 | 3 | 6 | 5 | 8 | 7 | 9 |
| Ranking by Manager 2 | 9 | 4 | 2 | 3 | 1 | 5 | 6 | 8 | 7 | 10 |

Calculate the coefficient of rank correlation and comment on the value.

Now this is the case ranking by manager one of several employees A to J, there are 10 employees. But Manager1 gives these ranks to these people. Manager2 gives this rank. Suppose it would not have given rank to you and it would have given you some score let us say 35 38 40 41 44 so 36, so what would have done is you would have made them a rank how. Now for example I am starting with the lowest listed 1. Let us say 2. Let us say 3, 4, 5, 6. So this is how you do and Manager II also has given some ranks.

**(Refer Slide Time: 26:50)**

Calculation of Rank Correlation Coefficient

| Employees | Rank by Manager I ($R_1$) | Rank by Manager II ($R_2$) | $D^2 = (R_1 - R_2)^2$ |
|-----------|---------------------------|----------------------------|-----------------------|
| A | 10 | 9 | 1 |
| B | 2 | 4 | 4 |
| C | 1 | 2 | 1 |
| D | 4 | 3 | 1 |
| E | 3 | 1 | 4 |
| F | 6 | 5 | 1 |
| G | 5 | 6 | 1 |
| H | 8 | 8 | 0 |
| I | 7 | 7 | 0 |
| J | 9 | 10 | 1 |
| N=10 | | | $\sum D^2 = 14$ |

Now this is the rank by Manager I and this is rank by Manager II. I am finding the D the difference between the 2 the difference between the first 1 is 1. So, I will take $1^2$, 2-4=- $-2^2$= -1 -2=1 it goes on. So my $D^2$ is 14.

**(Refer Slide Time: 27:12)**



We know that,

$$R = 1 - \frac{6 \sum D^2}{N^3 - N}$$

$$= 1 - \frac{(6*14)}{(990)}$$

$$= 1 - 0.085$$

$$= 0.915$$

*Thus, we find that there is a high degree of positive correlation in the ranks assigned by the two managers.

If I have my $D^2$ we know that R with the formula for the Spearman correlation

$$R = 1-(6\sum D^2)/N^3-N).$$

This is hardly it makes a difference it actually does not make much of a difference at all right. So if I do this is coming .915. So, we say does we find that there is a high degree of positive correlation in the ranks assigned by the 2 managers why it is helpful because it is helpful that it says that 2 managers are thinking alike.

So, if you are thinking alike otherwise we would have said that managers are thinking differently. So, this is about the Spearman's correlation. So earlier we discussed about the

Pearson's Correlation. We saw how to do that Pearson's Correlation. And we also saw how today how to do the Spearman's correlation for a non metric variable or non metric data. So this is all for the day. Thank you so much.