# IIT Madras

## ONLINE DEGREE

**Statistics for Data Science 1**
**Professor Usha Mohan**
**Department of Management Studies**
**Indian Institute of Technology, Madras**
**Lecture 4.5**
**Association between Two Numerical Variables: Describing Association**

(Refer Slide Time: 00:15)



So how do we describe the association between the 2 variables using a scatter plot? So when we are describing association between 2 variables, there are 4 key questions that I need to answer. The first question is, is there a direction? What do I mean by a direction? Does the pattern trend up or down or does it exhibit some sort of a trend? Is it linear or does it curve? The third is, are the points tightly clustered around the pattern or are they spread? Do I find anything unexpected? We will look at each one of these questions in detail now.
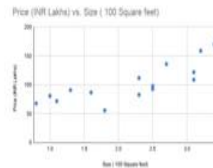
So, the first question we want to know is, does the pattern trend up or down? Let us look at an example again, if you look, go back to this example where I had plotted the size of a whole house on my x axis, and I wanted to know the price versus size, we see that there is a pattern where I can easily say that as the sizes of my homes increase.
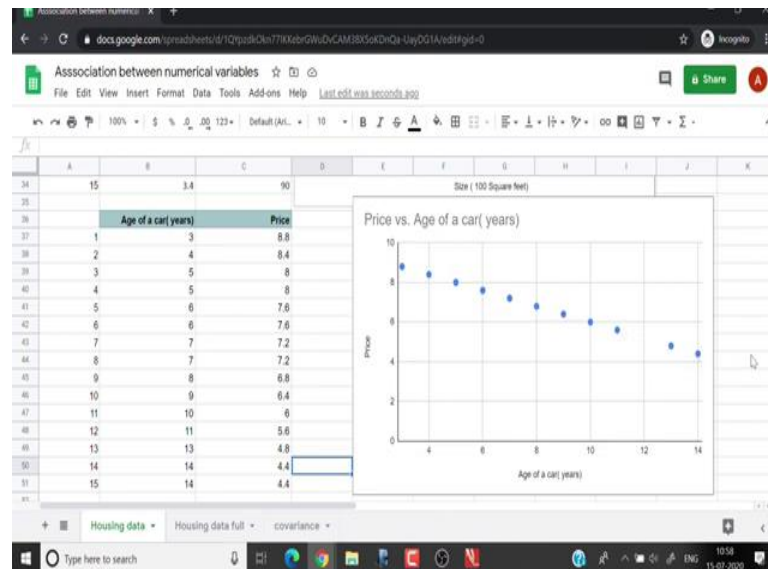
As size increases, price increases. I can see there is an upward trend, so the trend is up, okay. So, let us look at another example. We know that as car ages, so here instead of looking at an age of a person, my explanatory variable is age of a car. And my response variable is the price of the

car. So what I mean by this is as a car becomes older or the older the car is the price I am going to get for that car again reduces.

(Refer Slide Time: 02:32)



So let us look at a data here. So you can look at the data for a 3 year old car if I am getting the 8.8 lakhs again prices in lakhs of rupees for a 4 year old car I might get 8.4, for a 5 I might get 8, for 6 I would get 7.6 lakhs. So you can see that as the car becomes older my prices showing is coming down.
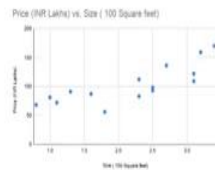
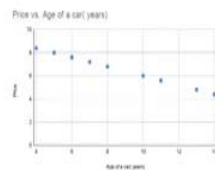 And that is what is shown by the scatter plot here. So what is the question or what is the pattern we see here? As the age increases, I see the price decreases. In this case I saw as size increases price increases. So here I can describe my pattern to have a decreasing or down trend. So the first thing which we need to see when we look at our scatter plot says whether it is showing an upward direction or a downward direction.
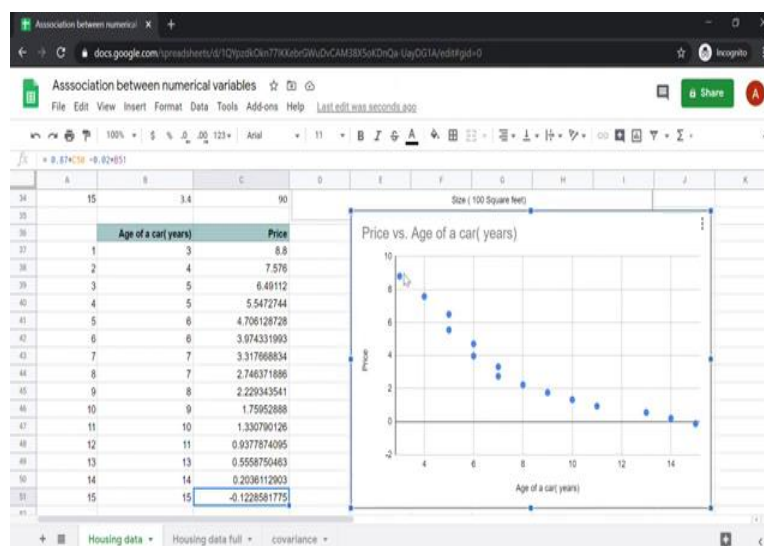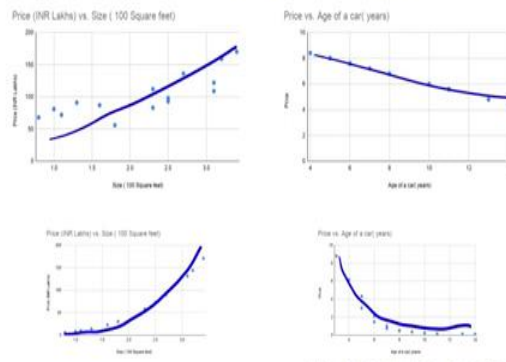
The next thing is, does the pattern appear to be linear or does it curve? Is it a linear pattern or is it a curved pattern. So let us look at the second example here. So now what you can see here is when I look at the price of the car versus age of a car here, now that you see notice a different pattern in this case.

So here you can see that so what you will notice here is even though the price of a car is decreasing, or it is showing a downward pattern as the car gets older, but you can see here there is a steep curvature, and it is not appearing linear as it was appearing earlier.

So in this case, it appeared to be a linear trend whereas here it is appearing to have curvature. Similarly, when I look at price of a house versus size of a house this graph appears to be a linear upward trend that is as my size increases, the price appears to increase linearly in terms of the size. Whereas here you can see that there is a curvature or I can say that as my size increases, the price increases according to a curved pattern.
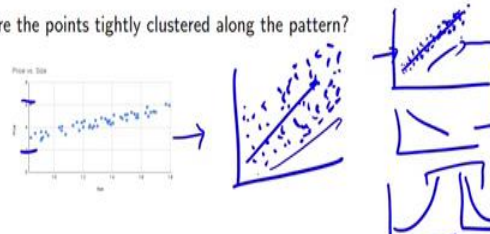
So the next question we ask is, look at the scatter plot and check whether it is linear, or whether it is a curve. For this lecture, we are focusing on linear relationships.

(Refer Slide Time: 05:32)

The next question which we want to answer is whether points are tightly clustered around the pattern? Now, if you look at this example, again, I have actually about 100 homes here. I have done the same thing, we have about 100 homes, okay. There are 100 homes I have taken the data of 100 homes here. So, you can see that the data we are talking about is 100 homes, and they are actually tightly clustered around each other.

So, if you look at this scatter plot which we are describing here, these are about 100 homes, and you can see that the data is tightly clustered in the sense that each of these 100 homes are between this range of my prices which could be between 2, 20 or 2 lakhs to 60 lakhs depending on what is the unit of the prices.

So, this is what we refer to as tightly clustered along the pattern. So, what do we mean by a pattern? First, we saw whether it was an upward pattern or a downward pattern, then we checked whether it was a curved pattern or whether it was a linear pattern, this is these 2 are linear patterns, whereas these 2 are curved patterns. Now, we are interested in knowing that along this pattern is my data tightly clustered.

This is the first question or along the pattern is my data variable very high. So, this first example gave me an example of a tightly clustered data whereas the second example is giving me an example of more variable data, both the cases I have an increasing trend which appears linear whereas this is more tightly clustered along the pattern, whereas this is more variable along the pattern.
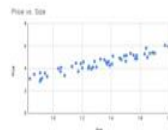
(Refer Slide Time: 07:44)



So, that is the next thing which we check this example was tightly clustered. Whereas you can see that this example is more variable or it is not tightly clustered. So, this is the third thing which is check is, how varied is one of the variable or what is the variation exhibited in the variable? This is the next important thing we checked.

(Refer Slide Time: 08:13)



The last thing which we check is what we refer to as the presence of an outlier. Again refer to the housing data. Now, look at these 3 points. So, there is this in fact, look at these 2 points which I am circling in red and there is this point, this is not an outlier, but now for now, let me just focus

on these 2 points, you can see that all the other points are actually behaving according to a particular pattern whereas these 2 points are away from the regular pattern does exhibited by the other points.

Now, if you look at this point, this tells us about a house which is between 1000 square feet and say 1250 square feet which is actually priced higher than the usual houses in this interval. And what this says is a large sized house which is priced lower than the smallest house also. So, these 2 points are referred to as outliers. So, outlier means or it is not following the pattern which other points exhibit.

(Refer Slide Time: 09:52)



So, when we are looking at association, in summary, the key things which we are trying to look at when we want to talk about association is. First we understand the direction whether it is an upward direction or a downward direction, whether the plot has an upward or a downward direction? Whether it is linear or it is curved? Whether it exhibits variation whether it is tight or whether it is varied?

And finally, whether there are presence of outliers? So, by this time we should know how to plot a scatter plot and look for the association between the variables through visual inspection when we are doing visual inspection these are the 4 key things which we need to take into account.