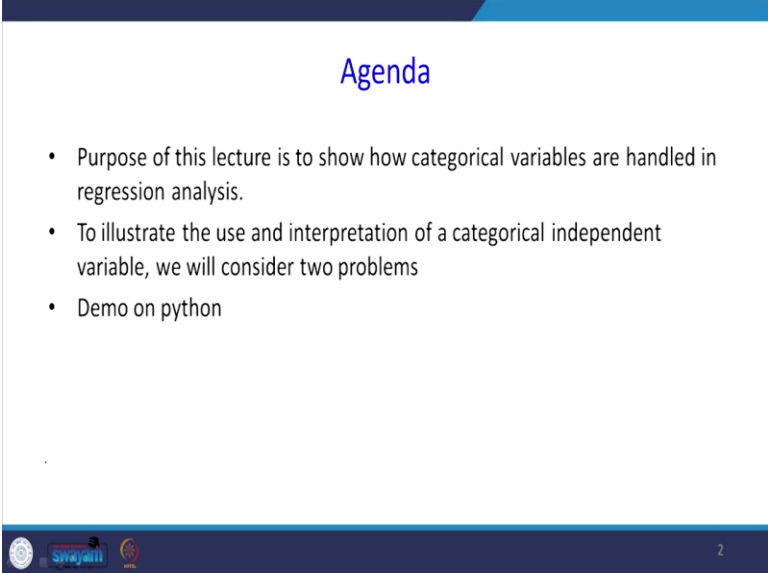


Data Analytics with Python
Prof. Ramesh Anbanandam
Department of Management Studies
Indian Institute of Technology – Roorkee

Lecture – 35
Categorical Variable Regression

Dear students, in this lecture, we will see how to handle Categorical Variable linear regression analysis. Whenever we do a linear regression analysis, the assumption is the nature of independent and dependent variable has to be continuous variable. Sometimes what will happen we have to include the categorical variable into independent variable category? How to handle that kind of regression analysis that we will see in this class?

(Refer Slide Time: 00:57)



The slide is titled "Agenda" in blue text. It contains a bulleted list of three items: "Purpose of this lecture is to show how categorical variables are handled in regression analysis.", "To illustrate the use and interpretation of a categorical independent variable, we will consider two problems", and "Demo on python". The slide has a blue header and footer. The footer contains logos for IIT Roorkee, Synergy, and a small circular logo, along with the number 2.

Agenda

- Purpose of this lecture is to show how categorical variables are handled in regression analysis.
- To illustrate the use and interpretation of a categorical independent variable, we will consider two problems
- Demo on python

The Agenda of this lecturer is, to show how categorical variable are handled in regression analysis. Illustrate and will interpret how to do the categorical independent regression analysis. The same problem we will do in Python will explain how to code and how to do this categorical regression in Python programming.

(Refer Slide Time: 01:17)

What are dummy variables?

- Dummy variables, also called indicator variables allow us to include categorical data (like Gender) in regression models
- A dummy variable can take only 2 values, 0 (absence of a category) and 1 (presence of a category)



3

Another name for categorical variable is called dummy variable dummy variable also called indicator variable. It allows us to include categorical nature in regression analysis. For example, gender is one of a categorical data where there is only two levels are possible male or female. If dummy variable can take only two values, when it is gender category, for example, zero means absence of category and one means the presence of category. Here zero will be taken as their reference. With respect to zero, we will compare what will happen to another level of the categorical variable.

(Refer Slide Time: 01:51)

Example 1: Problem / Background

- Johnson Filtration, Inc., provides maintenance service for water-filtration systems.
- Customers contact Johnson with requests for maintenance service on their water-filtration systems
- To estimate the service time and the service cost, Johnson's managers want to predict the repair time necessary for each maintenance request
- Hence, repair time in hours is the dependent variable
- Repair time is believed to be related to two factors,
 - the number of months since the last maintenance service
 - the type of repair problem (mechanical or electrical).



Source: Statistics for Business & Economics, David R. Anderson, Dennis J. Sweeney, Thomas A. Williams, Jeffrey D. Camm, James J. Cochran, Cengage Learning, 2013



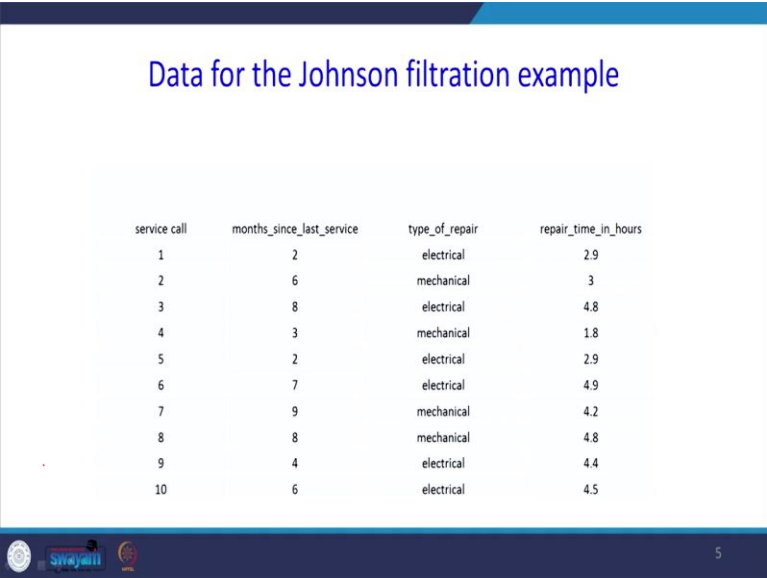
4

We will take a problem with the help of problem I will explain how to use categorical variable into the regression analysis and how to interpret it. This problem is taken from statistics for

Business and Economics from David Anderson, Sweeney and Williams. It is Syncage Publication in 2003 to 2013 edition. Johnson filtration Incorporation provides maintenance service for water filtration systems.

Customers contact Johnson's with a request for maintenance service on their water filtration system. To estimate the service time and the service cost Johnson's managers want to predict the repair time necessary for each maintenance request. Hence, the repair time in hours is the dependent variable. Repair time is believed to be related to two factors. One factor is number of months since the last maintenance service was done; second factor is the type of repair problem. Here the type of repair problem, mechanical or electrical is the categorical variable.

(Refer Slide Time: 02:53)



service call	months_since_last_service	type_of_repair	repair_time_in_hours
1	2	electrical	2.9
2	6	mechanical	3
3	8	electrical	4.8
4	3	mechanical	1.8
5	2	electrical	2.9
6	7	electrical	4.9
7	9	mechanical	4.2
8	8	mechanical	4.8
9	4	electrical	4.4
10	6	electrical	4.5

This is the given data. What is there is a Column 1 is the service call, the column 2 says months since the last service was done, in terms of month. Column 3 says the type of repair whether it is the repairs with respect to electrical system or mechanical system. The last column is repair time in hours. How much time it is taken for doing repairing?

(Refer Slide Time: 03:18)

```
In [23]: import pandas as pd
import matplotlib as mpl
import statsmodels.formula.api as sm
from sklearn.linear_model import LinearRegression
from scipy import stats
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
import statsmodels.api as s

In [24]: tbl = pd.read_excel('dummy.xlsx')
tbl

Out[24]:
```

	servicecall	months_since_last_service	type_of_repair	repair_time_in_hours
0	1	2	electrical	2.9
1	2	6	mechanical	3.0
2	3	8	electrical	4.8
3	4	3	mechanical	1.8
4	5	2	electrical	2.9
5	6	7	electrical	4.9
6	7	9	mechanical	4.2
7	8	8	mechanical	4.8
8	9	4	electrical	4.4
9	10	6	electrical	4.5

I have taken the screenshot of our python code get so I have to import necessary libraries like import Pandas as pd, import matplotlib as mpl, import statsmodels dot formula dot api as sm from sklearn.linear underscore model import LinearRegression from scipy import stats import seaborn sns, import numpy as np, import matplotlib.pyplot as plt, import statsmodels dot api as s. First we will load this regression file it is a data file I have saved in the name of dummy dot xlsx that we are going to save any object called tv1.

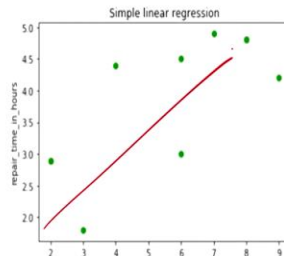
When you execute this one, we can see this is a data file. At the end of the class going to give the demo for this what are the codes which ever done in it. There also we can understand the steps. Here this is the data, display the data.

(Refer Slide Time: 04:10)

Linear Regression

```
In [41]: plt.scatter(tbl['months_since_last_service'], tbl['repair_time_in_hours'], color = "green")
plt.ylabel('repair_time_in_hours')
plt.title('Simple linear regression')
```

```
Out[41]: Text(0.5,1,' Simple linear regression')
```



But first will do the scatter plot between the months since last service and repair time in hours. When we look at this scatter plot, you see that there seems to be positive trend because when the month since last services more the repair time in hours also getting more. This is a simple linear regression considering only one independent variable. Here independent variable is continuous variable.

(Refer Slide Time: 04:41)

OLS Summary

```
In [44]: from statsmodels.formula.api import ols
Reg = ols(formula="repair_time_in_hours ~ months_since_last_service", data = tbl)
Fit1 = Reg.fit()
print(Fit1.summary())
```

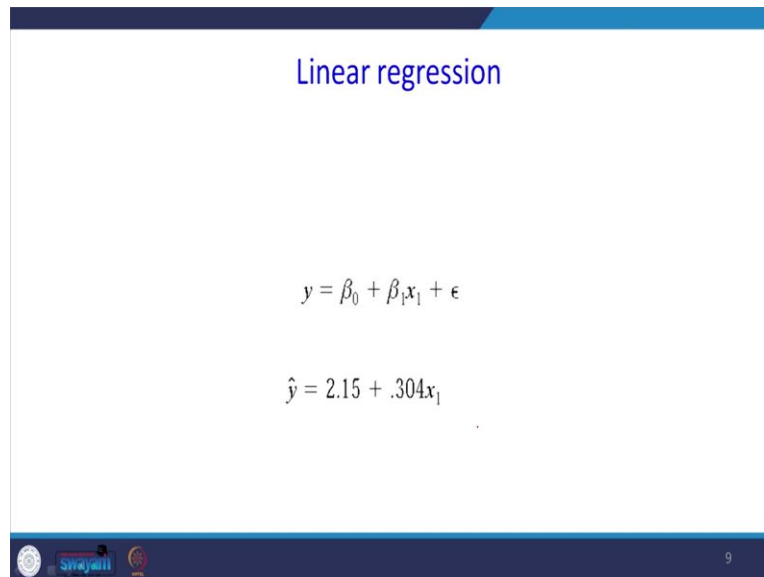
```
OLS Regression Results
=====
Dep. Variable:  repair_time_in_hours  R-squared:      0.534
Model:            OLS  Adj. R-squared:    0.476
Method:            Least Squares  F-statistic:    9.174
Date:            Sat, 07 Sep 2019  Prob (F-statistic):  0.0163
Time:            13:26:03  Log-likelihood: -10.602
No. Observations:  10  AIC:                25.20
Df Residuals:      8  BIC:                25.81
Df Model:          1
Covariance Type:    nonrobust
=====
                    coef    std err          t      Pr>|t|    [0.025    0.975]
-----
Intercept          2.1473      0.405      5.309      0.000      1.252      3.042
months_since_last_service  0.3041      0.100      3.029      0.016      0.097      0.516
=====
Omnibus:          0.907  Durbin-Watson:      2.154
Prob(Omnibus):    0.635  Jarque-Bera (JB):    0.751
Skew:             -0.501  Prob(JB):             0.687
Kurtosis:         2.107  Cond. No.             15.1
=====
```

$$y = 2.1473 + 0.3041x$$

When you do the regression analysis, this is output of python. So, from statsmodels.formula .api import ols, ols is used for doing regression analysis. Here, the dependent variable is repair underscore time in hours tilde sign independent variables months since last service. When you look at this series, y intercept, I can write y equal to 2.1473 + 0.3041 because this is independent

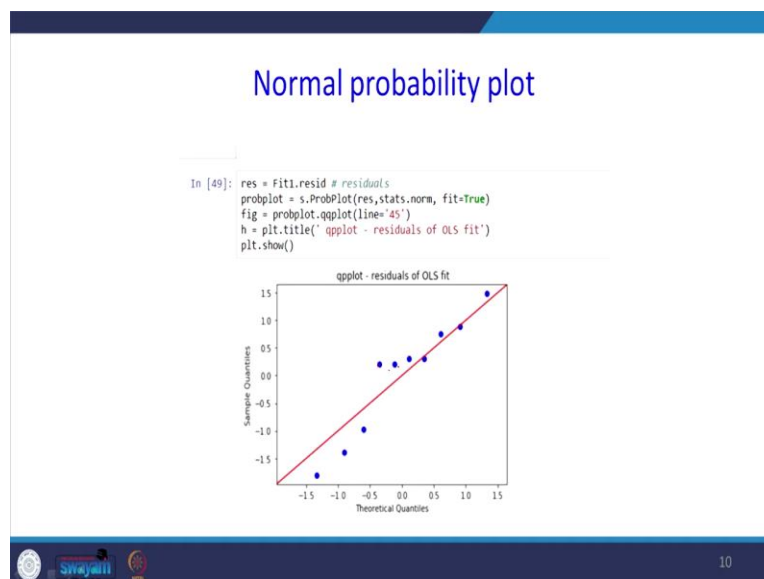
variables months since the last service was done. Look at the R square. R square is 53.4 % look at the P value of this independent variable here. Here, it is significant because it is less than 0.05. Now what we are going to do residual plots for this problem?

(Refer Slide Time: 05:38)



You look at this is 2.15. + 0.304 X1 is our regression model.

(Refer Slide Time: 05:47)



When we use this to do regression model, When you do the normal probability plot, look at this it has to all the probability points has to align with is red point. What is happening is there are so many points it is away from the red line. So, we can say that even if the, the residual plot is not appropriate, so the data, the error is not following normal distribution.

(Refer Slide Time: 06:10)

Creating dummies

```
In [34]: just_dummies = pd.get_dummies(tbl['type_of_repair'])
just_dummies
```

Out[34]:

	electrical	mechanical
0	1	0
1	0	1
2	1	0
3	0	1
4	1	0
5	1	0
6	0	1
7	0	1
8	1	0
9	1	0

$y = a + b_1x_1 + b_2x_2$
 $y = a + b_1x_1 + b_2(1)$
 $y = a + b_1x_1 + b_2(0)$

First, we will create a dummy variable for the categorical data. How to create a dummy variable for this categorical data? so that new dummy variable I going to call it is just underscore dummies equal to `pd.get_dummies` where the filename which column has to be converted into Dummies. So, the type of repair that is the value where we have written, whether the problem is related to mechanical or electrical.

So, when we display the just dummies see that that one variable is know, it is taken into 2 parts. 1 is for Electrical so, the presence of one says electrical; the absence of one says mechanical. There are 2 columns is there which is the dummy variable. So what happened both are same whether we can use this variable interval into our new regression model or this variable for Our new regression model, if you take electrical equal to 1.

So the equation be written as y equal to $a + b_1x_1 + b_2x_2$. Here, x_1 is independent variable.

The b_2 value will be 1 if it is suppose we write if the problem the, this is the common regression equation. In this regression equation, when you substitute x_2 equal to 1 that equation for Electrical problem related to electrical repair $a + b_1x_1 + b_2(1)$ this equation for repair due to electrical problem. Instead of this y equal to $a + b_1x_1 + b_2(0)$ this is what problem related to mechanical. You can reverse also, no problem.

Mechanical can be taken as 1 and electrical can be taken as zero. There will not be problem in the interpretation.

(Refer Slide Time: 08:06)

DATA FOR THE JOHNSON FILTRATION EXAMPLE WITH TYPE OF REPAIR
INDICATED BY A DUMMY VARIABLE ($x_2 = 0$ FOR MECHANICAL; $x_2 = 1$
FOR ELECTRICAL)

Customer	Months Since Last Service (x_1)	Type of Repair (x_2)	Repair Time in Hours (y)
1	2	1	2.9
2	6	0	3.0
3	8	1	4.8
4	3	0	1.8
5	2	1	2.9
6	7	1	4.9
7	9	0	4.2
8	8	0	4.8
9	4	1	4.4
10	6	1	4.5

This was the data which we have converted into dummy variable. Month since last service 1 represents problem related to electrical Zero represents problem related to mechanical. This was our Y is our dependent variable.

(Refer Slide Time: 08:22)

Adding dummies to table

```

In [18]: just_dummies = pd.get_dummies(tbl['type_of_repair'])
step_1 = pd.concat([tbl, just_dummies], axis=1)
step_1
step_1.drop(['type_of_repair', 'mechanical'], inplace=True, axis=1)

# to run the regression we want to get rid of the strings 'mechanical' and 'electrical'
# and we want to get rid of one dummy variable to avoid the dummy variable trap
# arbitrarily chose 'mechanical', coefficients on 'electrical' would show effect of 'electrical'
# relative to 'mechanical'

In [19]: step_1
Out[19]:

```

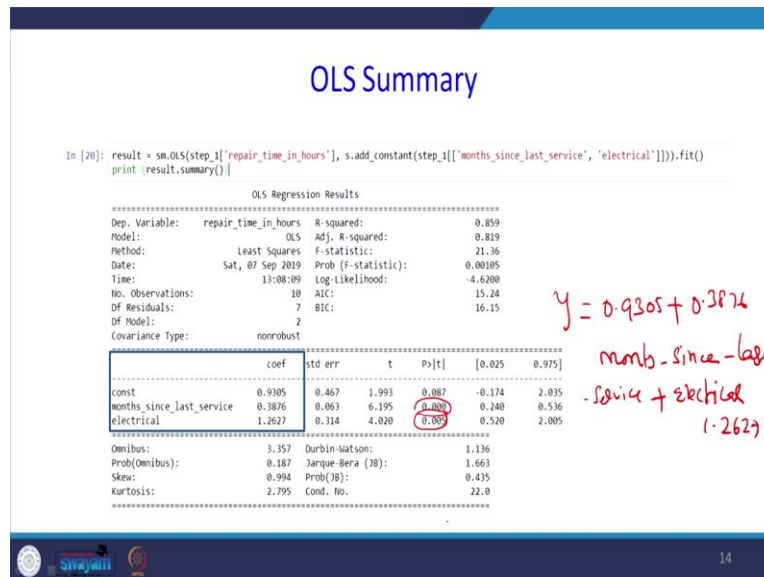
	servicecall	months_since_last_service	repair_time_in_hours	electrical
0	1	2	2.9	1
1	2	6	3.0	0
2	3	8	4.8	1
3	4	3	1.8	0
4	5	2	2.9	1
5	6	7	4.9	1
6	7	9	4.2	0
7	8	8	4.8	0
8	9	4	4.4	1
9	10	6	4.5	1

When you do the regression analysis, see that just_dummy is pd.get underscore dummies p1 is a type of repair. So here what I have done? I have displayed, I have dropped the certain columns what column I have dropped, I have dropped the column that is type of repair then I have added

only dummy variable with respect to the electrical repair. That is why this column has come. So, now this is going to this is the last column that is under electrical heading.

It is going to be taken as independent variable that will do the regression analysis.

(Refer Slide Time: 08:54)



Result equal to `sm.OLS(step_one['repair_time_in_hours']` is taken as a dependent variable. Months underscore since underscore last underscore service taken as independent variable. So this electrical is taken as reference because that column where 1 means electric repair 0 means mechanical repair. When you look at this, you see this equation can be written as Y equal to 0.9305 Plus months since last underscore service ,See coefficient for this one is 0.3876 + electrical 1.2627.

So look at R square it is 0.85 previously, the R square was when there is only one independent variable I m going back previously asked for R is only for 0.534 when we introduce another variable what has happened, the R square is increased to 0.859. So, F statistics corresponding probability the p-value is very low 0.005 so as a whole this regression model is significant. When we look at the individual independent variable, for example, months_since there is independent variable 1, the P value is less than 0.01.

So we can say this variable is significant. Similarly, for the second one the type of repair, where electrical is taken as the reference this also less than 0.05, so, this is also a significant variable.

(Refer Slide Time: 10:42)

Dummy regression

$$x_2 = \begin{cases} 0 & \text{if the type of repair is mechanical} \\ 1 & \text{if the type of repair is electrical} \end{cases}$$
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$
$$\hat{y} = .93 + .388x_1 + 1.26x_2$$

$x_2 = 1$ - electrical
 $x_2 = 0$ - mechanical

15

Now, this is the regression equation \hat{y} equal to $0.93 + 0.388 X_1 + 1.26 x_2$. If x_2 equal to 1 means electrical if I say x_2 equal to one it is related problem related to electrical if x_2 be 0 it is related to mechanical.

(Refer Slide Time: 11:15)

Interpreting the Parameters

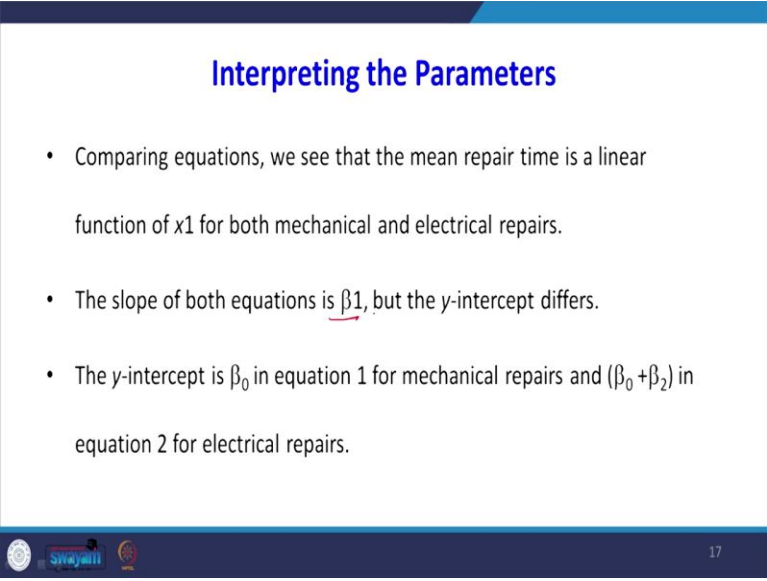
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$
$$E(y \mid \text{mechanical}) = \beta_0 + \beta_1 x_1 + \beta_2(0) = \beta_0 + \beta_1 x_1 \quad \text{Equation 1}$$
$$E(y \mid \text{electrical}) = \beta_0 + \beta_1 x_1 + \beta_2(1) = \beta_0 + \beta_1 x_1 + \beta_2$$
$$= (\beta_0 + \beta_2) + \beta_1 x_1 \quad \text{Equation 2}$$

16

The most important part, that is, interpreting the parameters. We know the expected value of Y equal to $\beta_0 + \beta_1 x_1 + \beta_2 x_2$ when you substitute equal to 1, when you substitute this $x_2 = 0$ that equation for mechanical, problem related to mechanical. So, $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$ so this term will become there $\beta_0 + \beta_1 x_1$ will be there. When substitute this x_2 equal to one that equations for the problem related to electrical.

So, E expected value y electrical equal to $\beta_0 + \beta_1 x_1$ so, $\beta_2(1)$ what is happening so, $\beta_0 + \beta_2$ that can be grouped that this will be $\beta_1 x_1$. See, both equations are same, both equation having the same slope β_1 only it differs by this extra value in our Y intercept how much with β_2 .

(Refer Slide Time: 12:16)



The slide is titled "Interpreting the Parameters" in blue text. It contains three bullet points: "Comparing equations, we see that the mean repair time is a linear function of x_1 for both mechanical and electrical repairs.", "The slope of both equations is β_1 , but the y-intercept differs.", and "The y-intercept is β_0 in equation 1 for mechanical repairs and $(\beta_0 + \beta_2)$ in equation 2 for electrical repairs." The slide has a blue header and footer. The footer contains logos on the left and the number 17 on the right.

- Comparing equations, we see that the mean repair time is a linear function of x_1 for both mechanical and electrical repairs.
- The slope of both equations is β_1 , but the y-intercept differs.
- The y-intercept is β_0 in equation 1 for mechanical repairs and $(\beta_0 + \beta_2)$ in equation 2 for electrical repairs.

Comparing equations 1 and 2 we see that the mean repair time is linear function of X_1 for both mechanical and electrical repair. The slope of both equation is β_1 , but the y-intercept differs. The y intercept is β_0 in equation 1 for mechanical repairs and $\beta_0 + \beta_2$ in equation 2 for Electrical repairs.

(Refer Slide Time: 12:45)

Interpreting the Parameters

- The interpretation of β_2 is that it indicates the difference between the mean repair time for an electrical repair and the mean repair time for a mechanical repair.
- If β_2 is positive, the mean repair time for an electrical repair will be greater than that for a mechanical repair; if β_2 is negative, the mean repair time for an electrical repair will be less than that for a mechanical repair.
- Finally, if $\beta_2 = 0$, there is no difference in the mean repair time between electrical and mechanical repairs and the type of repair is not related to the repair time.

The interpretation of Beta 2 is that it indicates the difference between the mean repair time of electrical repair and the mean repair time of mechanical repair. So the time differs by with this unit of this Beta 2. Beta 2 is positive the mean repair time for electrical repair will be greater than that of the mechanical repair. In our problem, it is beta 2 is positive, if the beta 2 is negative the mean repair time for an electrical repair will be less than that of mechanical repair. If finally you Beta 2 equal to zero there is no difference in the mean repair time between electrical and mechanical repairs.

And the type of repair is not related to repair time. This is most important because after doing a dummy variable regression you have to interpret it. The interpretation is this way. The first thing is you have to look at what is the sign of this Beta 2. Beta 2 is positive or negative. Then in case the beta 2 is 0, we can save the type of the time taken to repair that filter is nothing to do with the type of problem it has occurred. Whether it is problems related to mechanical repair or problem related to electrical repair.

(Refer Slide Time: 14:01)

Interpreting the Parameters

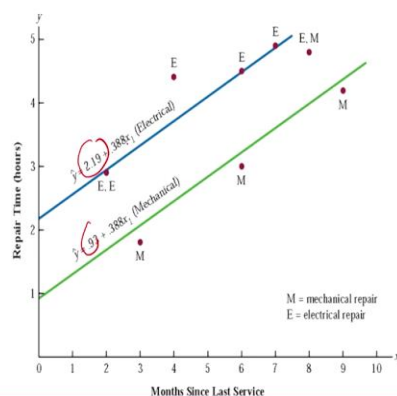
- In effect, the use of a dummy variable for type of repair provides two estimated regression equations that can be used to predict the repair time, one corresponding to mechanical repairs and one corresponding to electrical repairs.
- In addition, with $\beta_2 = 1.26$, we learn that, on average, electrical repairs require 1.26 hours longer than mechanical repairs.

In effect, the use of dummy variable for type of repair provides 2 estimated regression equation that can be used to predict the repair time, one corresponding to mechanical repair and another corresponding to electrical repairs, in addition, beta 2 = 1.26 we are getting this 1.26, going back, this 1.26. This 1.26 we learnt that the average electrical repairs required 1.26 longer than the mechanical repairs because for electrical repairs we have taken $x_1 = 1$, for mechanical repair, we have taken $x_1 = 0$.

So, the electrical repair is taken as the reference. What is the meaning of that is that the 1.26 time units the electrical repair is taking longer time than mechanical repairs. Look at this picture.

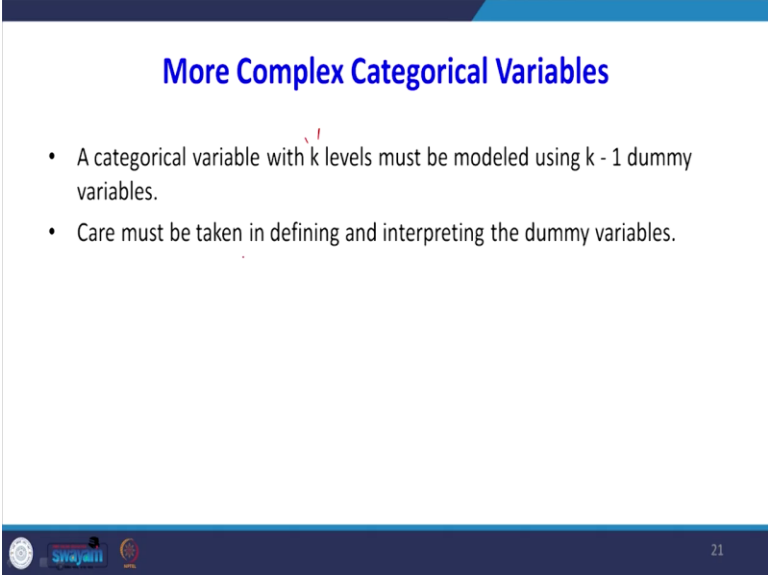
(Refer Slide Time: 14:55)

Interpreting the Parameters



The green one is for mechanical repair when substitute $x_2 = 0$ here, the blue one is for electrical repair, very extreme cold one. Look at this one. 2.19, this is 0.19. Both the slopes are same. This slope is 0.388 for this equation and this equation. Only the intercept is differs.

(Refer Slide Time: 15:13)



The slide has a blue header and footer. The title 'More Complex Categorical Variables' is in blue. The content area is white with two bullet points. The footer contains logos on the left and the number '21' on the right.

More Complex Categorical Variables

- A categorical variable with k levels must be modeled using $k - 1$ dummy variables.
- Care must be taken in defining and interpreting the dummy variables.

21

What is the logic is that here we have we have seen only two levels. Sometimes, there may be more than two levels. So, the number of a categorical variable with k levels must be modeled using $k-1$ dummy variable. What happened previously there was a 2 level, so we have taken only one dummy variable x_2 . So there are three levels you have to take $3 - 1$ that is a 2 dummy variable. Care must be taken in defining and interpreting the dummy variable.

What is the care here is what is the value we have assigned is equal to 1. For example, electrical repair, you take an equal to one that equation is integrated with respect to $x_2 = 1$.

(Refer Slide Time: 15:54)

Example 2: Problem / Background

- The manager of a small sales force wants to know whether average monthly salary is different for males and females in the sales force.
- He obtains data on monthly salary and experience (in months) for each of the 9 employees as shown on the next slide.



22

We will go for another problem. This problem is taken from statistics for management from Lemen N Rubeen. The manager of a small sales force wants to know whether the average monthly salary is different for males and females in a sales force. He obtained a data on monthly salary and experience for each of 9 employees as shown in the next slide.

(Refer Slide Time: 16:20)

Data

Employee	Salary	Gender	Experience
1	7.5	Male	6
2	8.6	Male	10
3	9.1	Male	12
4	10.3	Male	18
5	13	Male	30
6	6.2	Female	5
7	8.7	Female	13
8	9.4	Female	15
9	9.8	Female	21

Look at this. This is there are nine employees their salary, there is gender, there is experience. Now what you are going to do in this example, what is the salary of the females even though they have equal experience with the male, whether females are discriminated or not when we can say that they are getting discriminated, even though they have equal experience with male, they are getting lesser salary that means the females are discriminated.

(Refer Slide Time: 16:49)

```
In [50]: tbl2 = pd.read_excel('dummy2.xlsx')
tbl2
```

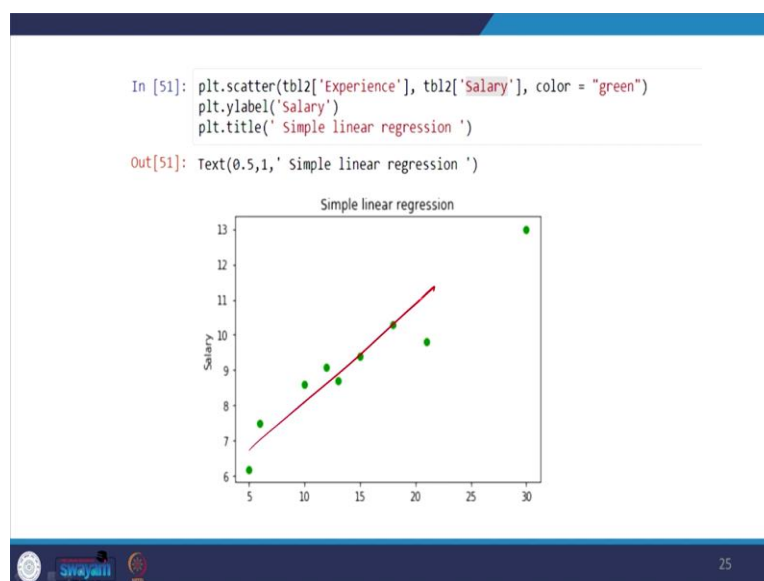
```
Out[50]:
```

	Employee	Salary	Gender	Experience
0	1	7.5	Male	6
1	2	8.6	Male	10
2	3	9.1	Male	12
3	4	10.3	Male	18
4	5	13.0	Male	30
5	6	6.2	Female	5
6	7	8.7	Female	13
7	8	9.4	Female	15
8	9	9.8	Female	21

24

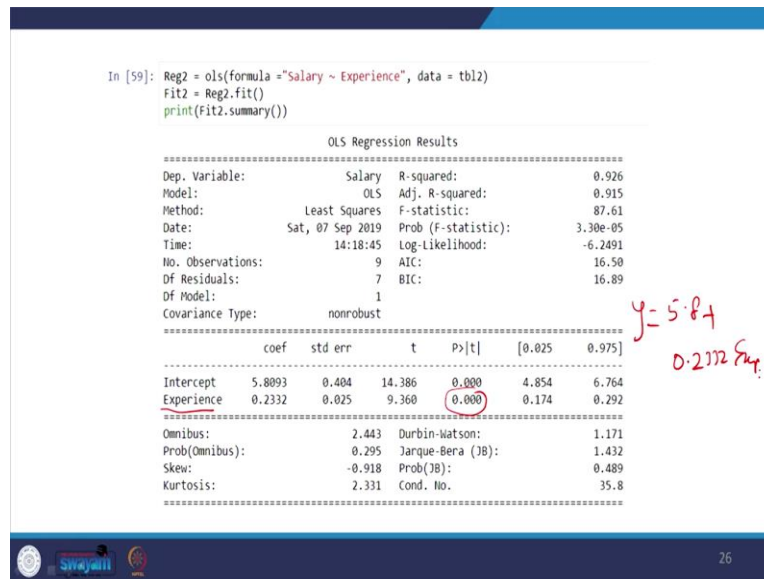
First, we will import the data. Here are imported in the object called `tbl2 = pd.read_excel`. The excel data where I have stored this problem is in the filename called `dummy2`. So, when I show this. Look at this, this is the employee salary, gender and experience. Next, what we are going to do?

(Refer Slide Time: 17:14)



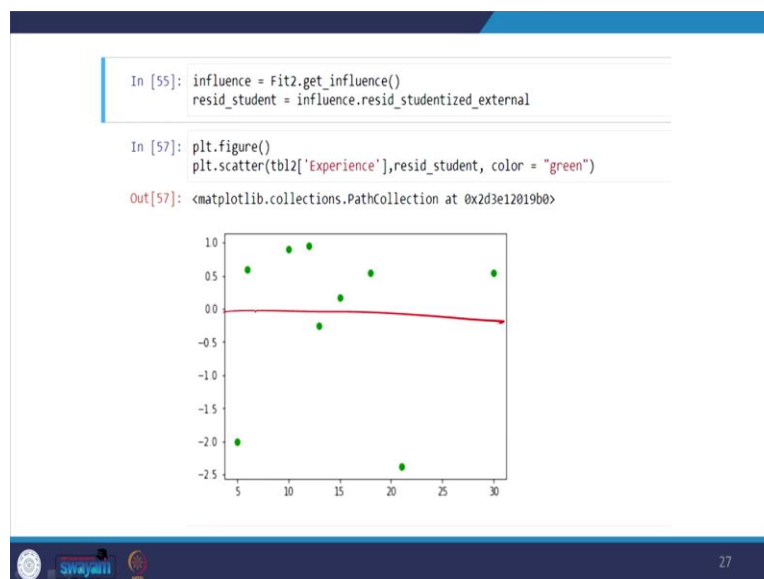
We are going to find out the scatter plot or is there any trend between the experience and the salary? It seems to be there is a positive trend. But look at the residual plot. What is this equation?

(Refer Slide Time: 17:29)



Y equal to see, R square is 0.926. See, the experience is the independent variable. Experience is the because p value is less than 0.05 we can say as the significant value. So we can write Y equal to $5.8 + 0.2332 \text{ experience}$. This is a regression equation. Ok now let us do the residual plot. For this we will do the error analysis.

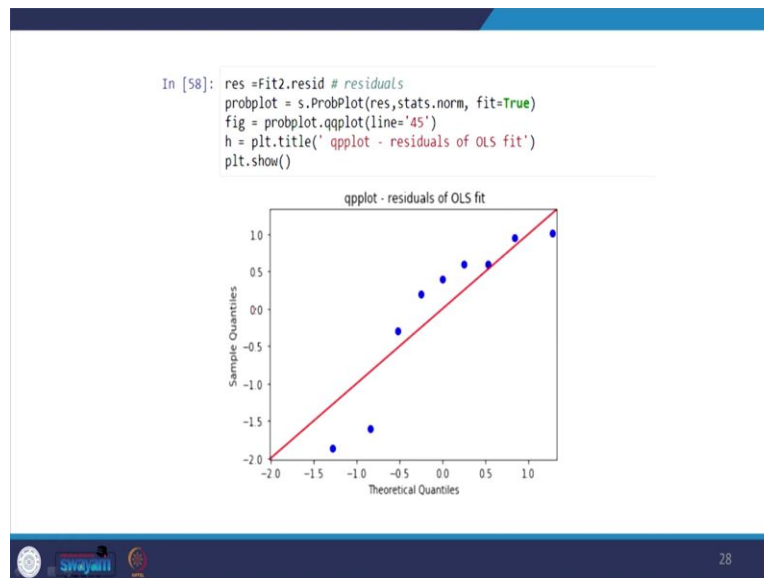
(Refer Slide Time: 17:58)



We will do the Residual analysis. You see that most of the points support, taken as the reference. This is a standardized residuals. Most of the points are you should be randomly it has to be distributed. Most of the points are above this way, there is a zero line. That means there is a

problem in assumption. Otherwise there might be some other variable that may affect the salary apart from experience.

(Refer Slide Time: 18:25)



Look at the see that quantile plot. You see that here also most of the points are above the pointers it has to sit on this red line, but it is not sitting on red Line then there is a problem in the assumption of that equal variance. That means error is not following equal variance.

(Refer Slide Time: 18:42)

Creating a dummy variable for gender

- Categorical data is included in regression analysis by using dummy variables
- For example, we can assign a value of 0 for males and 1 for females in our data so that a MR model can be developed

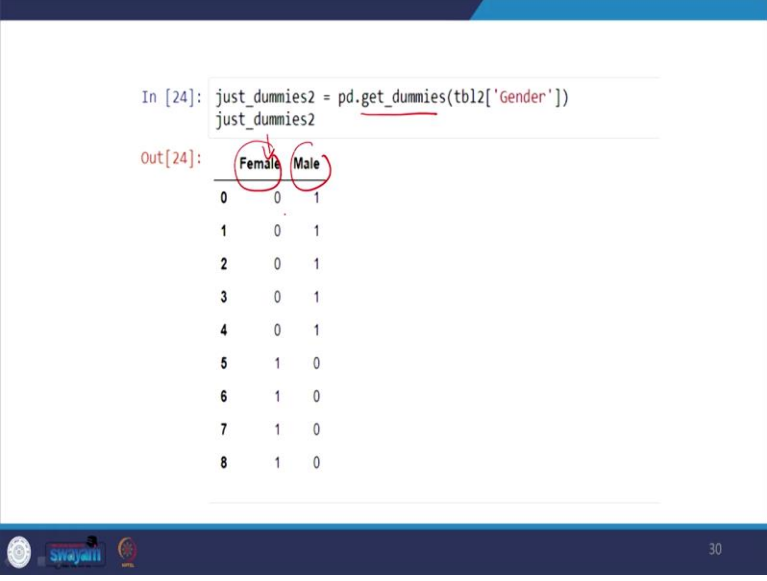
$x_1 = 0$ males
 $x_2 = 1$ females

Employee	Salary	Gender
1	7.5	0
2	8.6	0
3	9.1	0
4	10.3	0
5	13	0
6	6.2	1
7	8.7	1
8	9.4	1
9	9.8	1

Now, what we have done in this data. Categorical data is included in the regression analysis by using dummy variable here what you have done? Zero for males, 1for females. What has taken

zero also as reference or one also reference? So, one for male, female is taken as a reference now data, so that a multiple regression model can be developed. We will do that one.

(Refer Slide Time: 19:08)



The screenshot shows a Jupyter Notebook interface. The input cell contains the code: `just_dummies2 = pd.get_dummies(tbl2['Gender'])`. The output cell displays a table with two columns, 'Female' and 'Male', and eight rows of data. Red circles are drawn around the column headers 'Female' and 'Male' in the output table.

	Female	Male
0	0	1
1	0	1
2	0	1
3	0	1
4	0	1
5	1	0
6	1	0
7	1	0
8	1	0

From the given data, I have converted into dummy variable, one dummy variable for female because there are two level female and male. So, male is taken as one female taken is zero. The coding is that zero is taken as male one is taken female. So in this we are going to take this column for our further analysis. So, how to interpret this 0 means female one means male.

In creating a dummy variable for gender, we are going to follow this notation $x_2 = 0$ means male x_2 equal to 1 is taken as a female. So, after creating dummy variable first how to create a dummy variable in Python just _ dummies to that is a variable which I have given, `pd.get_dummies`. This was the command for making dummy variable. So we are going to take female column for further analysis. Zero means male one means female.

(Refer Slide Time: 20:10)

```
In [62]: step_1 = pd.concat([tbl2, just_dummies2], axis=1)
step_1.drop(['Gender', 'Male'], inplace=True, axis=1)
# to run the regression we want to get rid of the strings 'male' and 'female'
# and we want to get rid of one dummy variable to avoid the dummy variable trap
# arbitrarily chose "male", coefficients on "female" would show effect of "female"
# relative to "male"

result = sm.OLS(step_1['Salary'], s.add_constant(step_1[['female']])).fit()
print(result.summary())
```

OLS Regression Results						
Dep. Variable:	Salary	R-squared:	0.107			
Model:	OLS	Adj. R-squared:	-0.020			
Method:	Least Squares	F-statistic:	0.8426			
Date:	Sat, 07 Sep 2019	Prob (F-statistic):	0.389			
Time:	14:23:57	Log-likelihood:	-17.455			
No. Observations:	9	AIC:	38.91			
Df Residuals:	7	BIC:	39.30			
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	9.7000	0.853	11.367	0.000	7.682	11.718
Female	-1.1750	1.280	-0.918	0.389	-4.202	1.852
Omnibus:	0.387	Durbin-Watson:	1.912			
Prob(Omnibus):	0.824	Jarque-Bera (JB):	0.280			
Skew:	0.330	Prob(JB):	0.869			
Kurtosis:	2.441	cond. No.	2.51			

$y = 9.7 - 1.1750x_1$

This was our Python output for that regression analysis. When you look at this, R square is 0.107 but look at here, first I will write the regression equation. $Y = 9.7 - 1.1750 x_1$. How to interpret this result you see that in the x_1 is not the significant value here it is not significant. At the sample data level, what is the meaning of x_1 R? Look at this. If you write $x_2=1$ here x_1 equal to 1 is not x_2 it is x_1 .

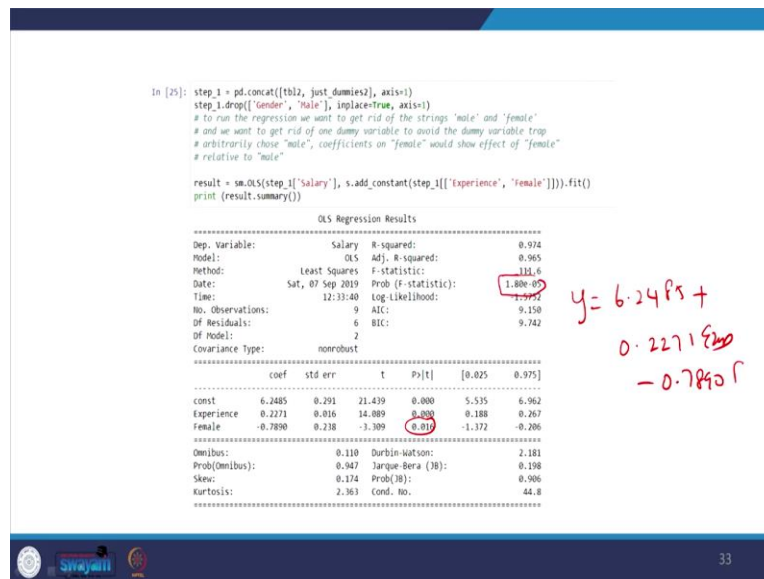
When you substitute $x_1 = 1$ this one, this coefficient says, it is negative. What is the meaning of these negative is this female is getting lesser salary when compared to male by this much unit because it is a negative sign, we go for interpretation.

(Refer Slide Time: 21:23)

More on the intercept and slope

- The value of the intercept, 9.70, is the average salary for males (as we coded gender=1 for females and 0 for males)
- The value of the slope, -1.175, tells us that the average females salary is lower than the average male salary by 1.175

The value of the intercept is 9.7 the average salary for males has been coded a gender 1 for female ok then, 1 for female and 0 for males. So, the value of the slope is - 1.175 tells us that the average salary is lower than the average male salary by 1.175. What is the meaning of this? Females are getting 1.175 units they are getting lesser salary when compared to male. If it is a positive then we can interpret that. When compared to male females are getting more salary because the negative you are saying that when compared to male, females are getting less salary. (Refer Slide Time: 22:09)



Now what we are going to do? We are going to introduce the previously we considered only gender. That is a female is taken as a reference. Now, we are going to introduce the experience also. When you introduce the experience also know the regression equation is y equal to 6.2485 + 0.2271, experience - 0.7890 female. Now look at the p-value these p-values now less than 0.05. Now, here the gender is significant variable.

In your previous slide, when you go back in this slide, you see that the p value is not significant. So we cannot say there is a gender discrimination. We can write a regression equation with the help of sample data, but at the population level, there is no connection between Gender and their salary because the relation between x_1 that is there gender and the salary there is no relationship. That means here the both female and male are getting same salary.

But when we introduce our experience was one of the variable now the general also is significant, so by considering experience and the gender, now gender also one of their significant because you are the P value is less than 0.05. We look at the f value the f value is very low. The probability value also very low as a whole model, this model is significant individually also all the variables are significant.

(Refer Slide Time: 23:47)

What would have happened if we had used 0 for females and 1 for males in our data? Would our results be any different?

```

In [6]: step_1 = pd.concat([tbl2, just_dummies2], axis=1)
        step_1.drop(['gender', 'female'], inplace=True, axis=1)

result = sm.OLS(step_1['Salary'], step_1[['Male']]).fit()
print(result.summary())

```

OLS Regression Results

	coef	std err	t	Pr> t	[0.025	0.975]
Dep. Variable:	Salary			R-squared:	0.107	
Model:	OLS			Adj. R-squared:	-0.009	
Method:	Least Squares			F-statistic:	0.0426	
Date:	Sat, 07 Sep 2019			Prob (F-statistic):	0.309	
Time:	14:27:56			Log-likelihood:	-17.455	
No. Observations:	9			AIC:	38.91	
Df Residuals:	7			BIC:	39.30	
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	Pr> t	[0.025	0.975]
const	8.5250	0.954	8.935	0.000	6.269	10.781
Male	1.1750	1.280	0.918	0.389	-1.852	4.202
Omnibus:	0.387	Durbin-Watson:			1.912	
Prob(Omnibus):	0.824	Jarque-Bera (38):			0.280	
Skew:	0.130	Prob(38):			0.869	
Kurtosis:	2.441	Cond. No.			2.77	

What would happen if we used zero for females and one for males in our data. Would our results be any different right? So for that purpose we have done some modification here. For example, gender female it is just reversed. You see that there is a difference in intercept but the slope is same but the slope sign is different. So what is the meaning here? The male right, because the male is 1, the males are getting 1.175 unit of higher salary when compared to females.

(Refer Slide Time: 24:27)

Male = 1, female = 0

- Not really – With coding as above, the intercept would change to 8.525 (the average female salary), the slope for gender would still be 1.175, but now it would have a positive sign (reflecting that average male salary is higher than average female salary by 1.175).
Predicted salaries from the model for males / females would not change no matter how dummy variable is coded

So what happened is there any difference in the result not really. With the coding as above, the intercept change to 8.525 see that 8.525 the slope of the gender would still 1.175, but it would have a positive sign reflecting that the average male salary is higher than average female salary by 1.175. So predicted salaries from the model for males and females would not change no matter how the dummy variable is coded.

(Refer Slide Time: 25:00)

More on dummy variables

- For gender, we had only 2 categories – female and male – thus we used a single 0/1 variable for this
- When there are more than 2 categories, the number of dummy variables that should be used equals the number of categories minus 1
- No. of Dummy Variables = No. of levels - 1

Sometimes, what will happen, they may be more than one dummy variable, how that in our problem. We have only two levels, sometimes there are three levels. We should have to Dummy variable will see that example. For gender, we had only two categories female and male does we used a single dummy variable 0,1 variable for this. When there are more than two categories the

number of dummy variables that should be used = the number of categories -1. So, the number of dummy variable is number of levels -1.

(Refer Slide Time: 25:33)

Example: Salary vs. Job Grade

- In this example, the categorical variable job grade has 3 levels, 1 (lowest grade), 2, and 3 (highest job grade)

Employee	Job Grade	Salary (\$000)
1	1	7.5
2	3	8.6
3	2	9.1
4	3	10.3
5	3	13
6	1	6.2
7	2	8.7
8	2	9.4
9	3	9.8

37

You see that there is one example where the job grade is there are three levels. 1, 2, 3, in this example, the categorical variable job grade as three level so, 1, 2, 3, 1 means lowest Grade, 2 means medium and 3 means highest grade. We are going to have three levels in our categorical data three levels are level 1, level 2, level 3.

(Refer Slide Time: 25:54)

Representing 3-level Job Grade using dummy variables Job_1 and Job_2

Employee's Job Grade	Job Grade	Dummy Variables	
		Job_1	Job_2
1	1	1	0
2	2	0	1
3	3	0	0

Job Grade 3 is the reference category

38

There are 3 levels and we are going to have only 2 Dummy variables. Job 1 say taken as 1, 0 job2 taken as 0,1 job3 is 0,0. So now, we can say this 0, 0 is taken as a reference, ok. So, the

presence of 1,0 will explain category 1; 0,1 will explain category 2; 0,0 will explain category 3. So here what is happening is there are 3 levels. But we are going to have only two dummy variable dummy variables. Dummy variable 1 and dummy variable 2.

(Refer Slide Time: 26:33)

Data file with dummy variables for job grade

Employee	Job Grade	Salary	Job_1	Job_2
1	1	7.5	1	0
2	3	8.6	0	0
3	2	9.1	0	1
4	3	10.3	0	0
5	3	13	0	0
6	1	6.2	1	0
7	2	8.7	0	1
8	2	9.4	0	1
9	3	9.8	0	0

Now, this is a new data set how this data set can be used for doing dummy variable regression. The interpretation is already I have explained to you now will go for demo of this code which I have shown in our, this presentation.

(Refer Slide Time: 26:49)

```

In [1]: import pandas as pd
import matplotlib as mpl
import statsmodels.formula.api as sm
from sklearn.linear_model import LinearRegression
from scipy import stats
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
import statsmodels.api as sm

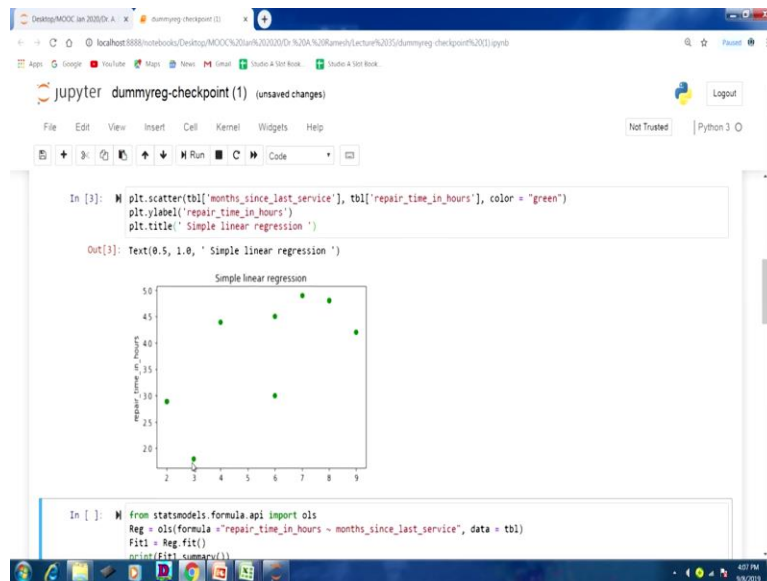
In [2]: tbl = pd.read_excel('dummy.xlsx')
tbl

Out[2]:
  servicecall  months_since_last_service  type_of_repair  repair_time_in_hours
0          1                2          electrical              2.9
1          2                6          mechanical              3.0
2          3                8          electrical              4.8
3          4                3          mechanical              1.8
4          5                2          electrical              2.9
5          6                7          electrical              4.9
  
```

I have prepared already code for that person. First I am going to remove this output by clicking kernel restart and clear output. I have cleared the output now I am going to run this one. So as

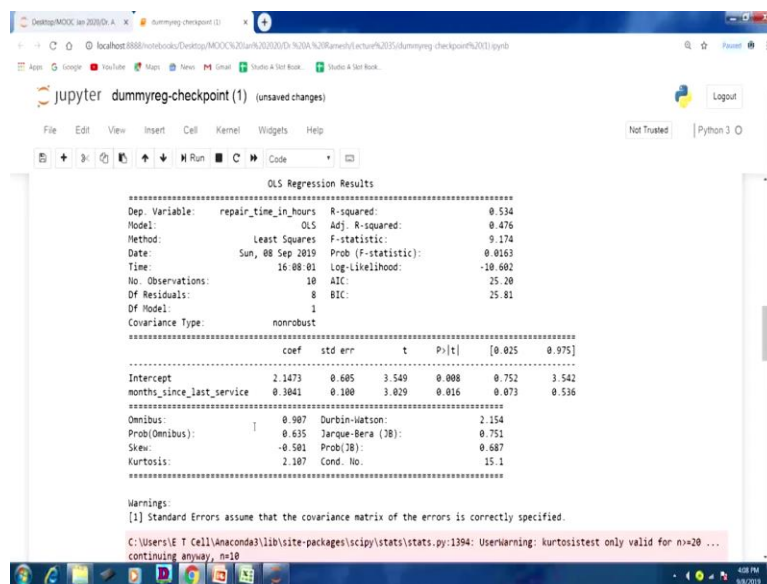
you know, this is Shift Enter so again shift enter this is the data. This data shows service call months since last service type of repair. Next we will go for scatter plot.

(Refer Slide Time: 27:17)



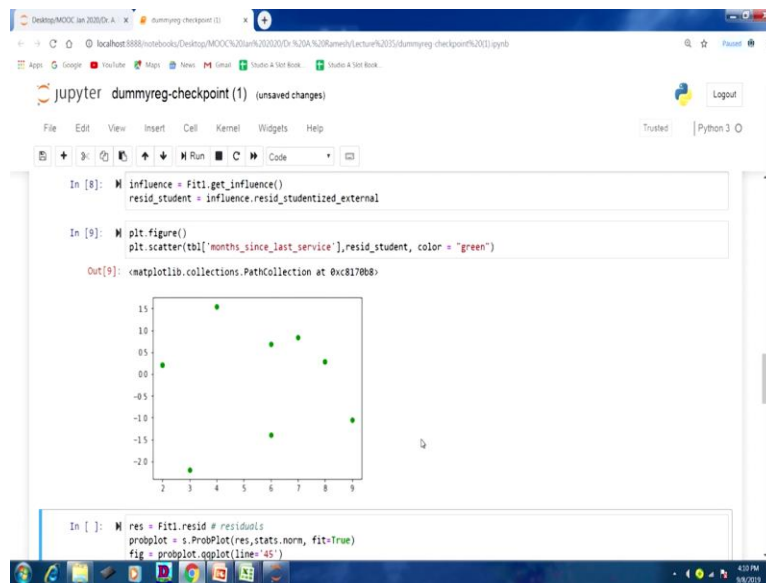
Scatter plot shows that there is a correlation between month since last service and repair time in hours next will go for simple linear regression where were taken only one independent variable.

(Refer Slide Time: 27:35)



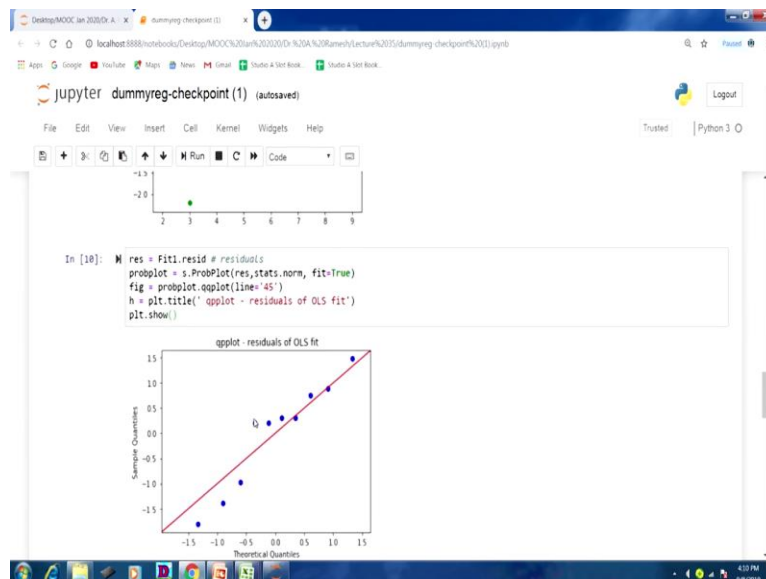
When you look at this here, this is equal to $2.14 + 0.3041 \times x_1$, suppose these variables x_1 . Look at the p-value this p-value is less than 0.05. So this variable is significant value variable, R square also it is good above than 0.5. See, when there is f statistic this is also less than 0.05. So, as a whole model it is valid.

(Refer Slide Time: 28:04)



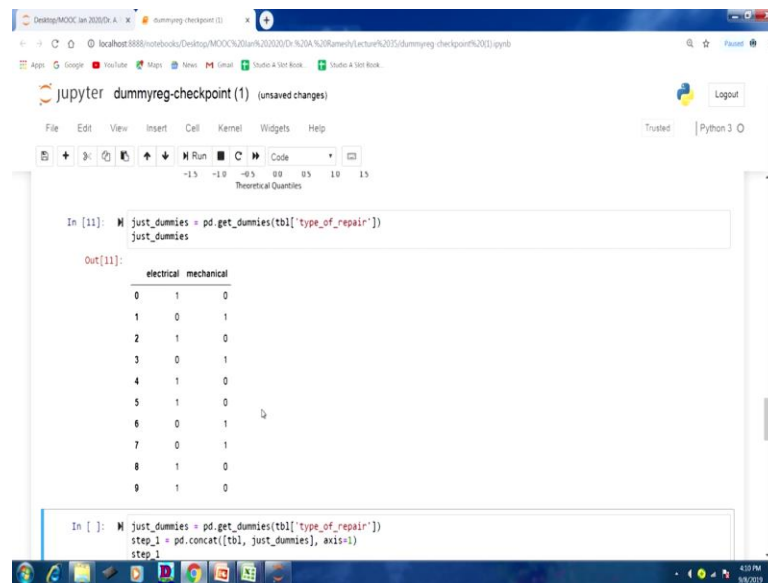
Now, we will plot standardized residual plot. When you look at the standardized residual plot this is the pattern. See there are that some points which are going above -2 how to interpret the standardized residual plot all the points should be between -2 to $+2$. But it seems that there are some variable which goes beyond -2. So it is violating our model assumptions. Now we will go for this q plot.

(Refer Slide Time: 28:38)



These also see that these into some pattern continuously three lines are below this line. There are so many points are above this line. There are also problems in variance of the error variable also.

(Refer Slide Time: 28:51)



The screenshot shows a Jupyter Notebook interface with a code cell and its output. The code cell contains the following Python code:

```
In [11]: just_dummies = pd.get_dummies(tbl["type_of_repair"])
just_dummies
```

The output cell displays a DataFrame with two columns: 'electrical' and 'mechanical'. The rows represent individual observations, indexed from 0 to 9. The 'electrical' column has values 0 or 1, and the 'mechanical' column has values 0 or 1.

	electrical	mechanical
0	1	0
1	0	1
2	1	0
3	0	1
4	1	0
5	1	0
6	0	1
7	0	1
8	1	0
9	1	0

Below the output, there is another code cell that concatenates the original data with the dummy variables:

```
In [ ]: just_dummies = pd.get_dummies(tbl["type_of_repair"])
step_1 = pd.concat([tbl, just_dummies], axis=1)
step_1
```

Now we will convert the data into dummy variable. This is dummy variable electrical is taken as one mechanical is taken as 0, now after converting into dummy variable to drop the column dummy variable belongs to mechanical. After Dropping we can see this output, for duplicate this now. There is no mechanical column only electrical column is there.

I will do for this data set will go for regression analysis two independent variable one is months underscore since last service another one is type of repair that is electrical is taken as reference. When we look at the p-value the p-value are both independent variables less than 0.05, so the significant model in this equation when you substitute x_2 equal to 1 will get a regression equation for problem related to electrical. When you substitute x_2 equal to 0 will get a regression equation for problem related to mechanical repairing system.

Now we will be going for another problem. This is our second problem, where the salary is the dependent variable. Experience is independent variable gender also independent variable. When we plot that between experience and salary there is a positive relationship. Now will take salary and experience, experience is an independent variable. You see that experiences a significant because less than 0.05. R square is 0.26.

There is no problem in this. Now, we look at the standardized residual plot that most of the points there is not equally plotted, most of the point above zero there is no randomness in the distribution. There seems to be some pattern in the residuals. We will go for checking the normality of the variance error term. See that it is this also following some kind of a pattern and then also not sitting on the exactly the diagonal line.

Now will go for create a dummy variable for the gender. There is one is for female another one is male now will drop this one. So, female is taken as 1 male is taken as zero when you do there is regression analysis where Gender is taken as a female now, you see that Y equal to 9.7 plus (-1.175) female. So, females are getting less salary than the male but look at the P value, when you consider only the gender, the P value is more than 0.05. So this gender variable is not significant.

When you, when you bring another variable is an experience when you look at our previous code, it is only gender is taken gender also it is not significant because the p-value 0.389. Now will take Gender and experience together, let us see what is happening. When you take Gender and experience together, you see that the P value for female is less than 0.05. The experience also, listen 0.05. Both the variables are significant, but the female is getting less salary, when compared to male even though they have equal experience.

Now, what will happen when you reverse the code? Suppose, we have taken female equal to 1 male equal to zero now, what will happen? When you reverse that code send male equal to 1 and female equal to zero what will happen if there will not be any change in the result. Only the sign of usually the male is taken that was - 1.17. Now female is taken as reference. So we are getting only the positive value of 1.17.

Only the difference in the Y intercept, otherwise all interpretations are same. In this lecture by using dummy variable regression I have taken 2 problems with the help of python code I have explained how to do a dummy variable regression and I have also interpreted the result. We know what is the dummy variable regression is sometime the gender is one example for dummy variable regression because there are two possibilities, male and female.

Similarly the job category, Category 1, category 2, category 3, these are dummy variable. For this purpose we have learnt how to do a regression analysis, the next class very important topic that is logistic regression, we are going to see that one before seeing Logistic regression. There is a one principle called maximum likelihood principle. I will explain what is the maximum likelihood principle? With the help of some examples, then we will go per Logistic regression in the next class. Thank you very much.