

Statistical Analysis of Dummy Variable Models and Testing for Seasonal Fluctuations
Part-5
Professor Sabuj Kumar Mandal
Department of Humanities and Social Sciences
Lecture 42
Indian Institute of Technology, Madras

Welcome once again to our discussion on dummy variable models basically we were discussing about the applications of dummy variable models. So, in our last class if you recall then what we were doing? We were trying to understand basically whether a dummy variable model can be applied in a context where our variable, so firstly we discussed about on the saving income relationship.

And then later on we talked about how to use that saving income relationship in a dummy variable context. That means to test the structural break how the dummy variable model can be an alternative of the Chow test that we learned earlier. And we discussed that dummy variable model can actually overcome the several limitations of Chow test model.

(Refer Slide Time: 1:56)

Dummy variable model for testing seasonal fluctuation

$$\text{Sales} = \alpha + \beta (\text{dummy variable}) + u_t$$

ANOVA: $\text{Sales} = \alpha + \beta_2 D_{2i} + \beta_3 D_{3i} + \beta_4 D_{4i} + u_i$

Obj: Whether sales in a particular quarter is significantly different from the base quarter, which is the 1st quarter here.

Initial obj: Whether sales are sig in all the quarters or not?

$D_{2i} = 1$ if 2nd quarter
 $= 0$ otherwise
 $D_{3i} = 1$ if 3rd quarter
 $= 0$ otherwise
 $D_{4i} = 1$ if 4th quarter
 $= 0$ otherwise
 D_{1i} : 1st quarter is the base category

$T = S + C + T + U$
 ↓ ↓ ↓
 seasonal trend random

Today again we will discuss another application of dummy variable model. And this application is basically to test whether there is any seasonal fluctuation in a time series data. So dummy variable model for testing seasonal fluctuation.

Now many times what we see that the time series variable give some kind of fluctuation based on which particular season we are observing that particular variable. For example, if you think

about the sales data let us say we have sales data on refrigerator, sales data on refrigerator and other house hold appliances, let us say washing machine, dish washer then refrigerator so on and so forth.

And the many a times what we observe that the sales of these house hold appliances do show some kind of seasonal fluctuation. For example, sales of refrigerator will be quite high during the summer season and quite low in winter. Similarly, the home appliances like washing machine, T.V so on and so forth, you might be seeing that the sales of these appliances are quite large during the Diwali season.

So, if you have let us say quarterly data that means in a year basically you have four quarters. And your objective here is to see which particular quarter in a year gives you the significant sale. Whether all the quarter sales in all the quarters are significant or it is happening in some of the quarters.

And based on that as a salesman, as a sales person you may devise some kind of policies, that is our objective. So, we have data on sales for these household appliances and you want to see whether sales in in a particular quarter is significant. Now why do we do so? This is very important for time series data.

For example, here if you look at the sales basically sales is actually a function of let us say $\alpha + \beta_1$ is the expenditure on durables. Let us say this is my durable expenditure plus this is let us say t ut. This is sales which is sales on household appliances which is a function of the expenditure on durables.

And if you try to run this type of regression then what will happen? Firstly, what you need to do, you need to remove some kind of seasonal fluctuations from this data. Why this is so? Because a time series data t is actually a summation of three types of fluctuations.

What are those? Firstly, a time series data will show you some kind of seasonal fluctuation. So, let us say that is called seasonal, then second is cyclical, then third is t and the fourth is the random one. So, this is seasonal, this is cyclical, then this is trend any time series data will show you some kind of trend over a period of time and this is purely random.

Now this random component what you are already adding in the model you see there is u_t but s , c and t the seasonal, cyclical, and trend this variable should be removed from the time series variables. For example, here sales or durable expenditure otherwise you may get a spurious regression in time series data. And there are specific ways by which you can remove the cyclical and the trend component of the time series.

In this example we are basically interested in how to remove the seasonal component from the time series data. So, if you do so sales is the function of durable expenditure that means to remove the seasonal component from the data, we need to consider that these four quarters are actually different. And that is why while modeling, we need to accommodate the fact that these quarters are actually different. If you do not consider that fact then your model will suffer from the seasonal fluctuation.

So, what we need to do here? Here let us say for the time being I am removing the durable expenditure part, this variable so we basically will start our analysis with a simple ANOVA model, simple ANOVA model. So, as you know that ANOVA model will not consider any quantitative covariate as explanatory variable. So, in that sales is a function of let us say you have four quarters, then you have α plus let us say β_1 , β_2 , β_3 , β_4 plus u_i .

And how we defined D_2 , D_3 and D_4 ? So D_{2i} equals to 1 if second quarter and 0 otherwise. D_{3i} equals to 1 if third quarter and 0 otherwise; and D_{4i} equals to 1 if fourth quarter and 0 otherwise. This is how we have defined the dummy. And we have not defined what is our base category here. Since we have not assigned any dummy for the first quarter that means D_{1i} is a base category.

So that means all this β_2 , β_3 , and β_4 they are differential intercept. So that means they will basically tell you that whether the sales in second quarter is significantly different from first quarter. D_3 will tell whether the sales in the third quarter is significantly different from the first

quarter. And D_4 will be β_4 will tell whether the sales in fourth quarter is significantly different from the first quarter.

But see our objective initially when we pose the problem our objective was not to examine whether sales in a particular quarter is different from any of the other quarters any with the base quarter which is first quarter here. Rather we were interested in out of these four quarters which particular quarter gives you the significantly higher sale. It might be second quarter, third quarter, fourth quarter, first quarter or all the quarters together.

Are you getting my problem? Our initial research question was whether the sales in refrigerator or likewise any other household appliances, sales in household appliances are significant in all the quarters or in one or two quarters. So basically, we were interested in only seasonal presence of seasonal fluctuation on the sales of the durables. That is our objective.

If that is the case if you said this type of model, dummy variable model, then your objective is not fulfilled because in this model they will only tell you whether the sales in a particular quarter is different from the base quarter. So, what is this model tell you? This model's objective is whether sales in a particular quarter is significantly different from the base quarter which is the first quarter here.

But our initial objective, what was our initial objective? Our initial objective was to see whether sales are significant in all the quarters or not, all the quarters or not. So, we were not interested in checking the significance of sales in a particular quarter with the base. Now if this is your objective then you cannot set this type of model.

(Refer Slide Time: 15:58)

Windows Journal window showing handwritten notes:

$$\text{Sales} = \alpha + \beta_1 D_{1i} + \beta_2 D_{2i} + \beta_3 D_{3i} + \beta_4 D_{4i} + u_i$$

Problem? We have assigned four dummies to represent four quarters.

Rule: (n-1) dummy
- Dummy var. trap

$D_{1i} = 1$ if 1st quarter
 $= 0$ otherwise

So, what type of model then you have to set? You have to basically set this type of model. Sales equals to alpha plus beta1 D1i plus beta2 D2i plus beta3 D3i plus beta4 D4i plus ui. And now beta D2, D3, D4 we have already defined. And D1i is now equals to 1 if first quarter and 0 otherwise. So, this model should be specified. But if you specify this model what would be the problem in estimation? What would be the problem in estimation? Can think of? If I set this type of model what would be our problem in estimation?

The problem is, see here we have assigned four dummies to represent four quarters. But what is the rule says? The rule says that there should be n-1 dummy, that means if you have n categories, here it is a four categories, four quarters so you should have introduced four minus one, three dummies.

Now our problem is if we assign three dummies then that model will tell you what would be the interpretation of the dummy variable coefficient. They will only tell you whether sales in a particular quarter is significantly different from base quarter. So that cannot serve our initial purpose and if we fit this model to serve our initial purpose that means to see whether sales in all the quarters are significant. Then we will end up with dummy variable trap.

(Refer Slide Time: 19:21)

	dish	frig	wash	dur	d1	d2	d3	d4
1	841	1317	1271	252.6	1	0	0	0
2	957	1615	1295	272.4	0	1	0	0
3	999	1662	1313	276.9	0	0	1	0
4	968	1295	1158	275.9	0	0	0	1
5	894	1271	1289	268.9	1	0	0	0
6	851	1555	1245	262.9	0	1	0	0
7	863	1639	1278	276.9	0	0	1	0
8	878	1238	1383	263.4	0	0	0	1
9	792	1277	1273	268.6	1	0	0	0
10	589	1258	1831	231.9	0	1	0	0
11	657	1417	1343	242.7	0	0	1	0
12	699	1185	1381	248.6	0	0	0	1
13	675	1196	1381	258.7	1	0	0	0
14	652	1418	1316	248.4	0	1	0	0
15	628	1417	1318	255.5	0	0	1	0
16	539	919	1325	248.4	0	0	0	1
17	488	943	1316	247.7	1	0	0	0
18	538	1175	1111	241.1	0	1	0	0
19	557	1269	1111	241.1	0	0	1	0
20	682	973	1111	241.1	0	0	0	1
21	658	1182	1111	241.1	1	0	0	0
22	749	1344	1111	241.1	0	1	0	0
23	827	1643	1111	241.1	0	0	1	0
24	858	1643	1111	241.1	0	0	0	1

Now what is the exact problem of dummy variable trap? That first we will try to understand using the data. So, look at the data set. This is our data. So, we have quarterly data. So, in each year we have four data points. And look at, see here dish means dishwasher sale 841 all in thousands. This data is also taken from the United States and this is taken from your textbook Gujarati that basic econometric textbook and this is table number 9.3. This data is from table 9.3.

So, this is the sales data of refrigerator which is represented as frig, and then this is washing machine and this is actually quantitative covariate which shows the total expenditure on the durables. And then what we did? We have introduced four dummies. So, this data is from first quarter that is why D1 equals to 1 and see all other quarters are 0. D2, D3, D4 are 0.

Similarly, the second data point is from the second quarter that is why first, third and fourth quarter dummies are 0. This data is from the third quarter that is why first, second, and fourth quarter they take the value 0. And similarly, the dummy for the fourth quarter is 1 and that is why first, second, and third quarter take the value 0. So here we have assigned four dummies for four quarters. And let us see if we do so that means if we violate the rule of assigning dummy what type of estimation problem, we may end up due to dummy variable trap.

(Refer Slide Time: 21:24)

StataSE 16.0

File Edit Data Graphics Statistics User Window Help

History

Statistics/Data Analysis 16.0 Copyright 1985-2019 StataCorp LLC
StataCorp
4905 Lakeway Drive
College Station, Texas 77845 USA
800-STATA-PC http://www.stata.com
979-696-4600 stata@stata.com
979-696-4601 (fax)

Special Edition

Single-user Stata perpetual license:
Serial number: 40160623548
Licensed to: NPTEL
NPTEL

Notes:

1. Unicode is supported; see help unicode_advice.
2. Maximum number of variables is set to 5000; see help set_maxvar.

.*(8 variables, 32 observations pasted into data editor)

Variables

Name	Label
dish	DISH
frig	FRIG
wash	WASH
dur	DUR
d1	D1
d2	D2
d3	D3
d4	D4

Properties

Variables

Data

StataSE 16.0

File Edit Data Graphics Statistics User Window Help

History

NPTEL

Notes:

1. Unicode is supported; see help unicode_advice.
2. Maximum number of variables is set to 5000; see help set_maxvar.

1. reg frig d1 d2 d3 d4

.*(8 variables, 32 observations pasted into data editor)

. reg frig d1 d2 d3 d4
note: d4 omitted because of collinearity

Source	SS	df	MS	Number of obs =
Model	915635.844	3	305211.948	32
Residual	806142.375	28	28790.7991	10.60
Total	1721778.22	31	55541.2329	F(3, 28) = 0.0001
				Prob > F = 0.0001
				R-squared = 0.5318
				Adj R-squared = 0.4816
				Root MSE = 169.68

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
d1	62.125	84.83926	0.73	0.470	-111.6603 235.9103
d2	307.5	84.83926	3.62	0.001	133.7147 481.2853
d3	409.75	84.83926	4.83	0.000	235.9647 583.5353
d4	0 (omitted)				
_cons	1160	59.99041	19.34	0.000	1037.115 1282.885

Variables

Name	Label
dish	DISH
frig	FRIG
wash	WASH
dur	DUR
d1	D1
d2	D2
d3	D3
d4	D4

Properties

Variables

Data

Let us see we will try to estimate. So, what we will do? We will regress reg let us say frig means sales in refrigerator which is a function of D1, D2, D3, and D4, D1, D2, D3 and D4 four dummies we have introduced. And then enter. Now look at this, look at this Stata has reported that D4 omitted because of collinearity. And here also in the, in the, in the output set we see the coefficient corresponding to D4 is 0 and it is written in the bracket omitted, omitted.

Now what is the problem that means we have assigned four dummies but while reporting the result Stata is basically omitting one dummy here, we have omitted the fourth dummies. Now why Stata is dropping one dummy that means is there any problem in estimation. If we introduce

four dummies for four quarters to understand that we have to once again go back to the way, we have represented our multiple linear regression model using that matrix Algebra. That will make things clear. So, you will know what exactly is the problem of dummy variable trap.

(Refer Slide Time: 23:07)

$$\text{Sales} = \alpha + \beta_1 D_{1i} + \beta_2 D_{2i} + \beta_3 D_{3i} + \beta_4 D_{4i} + U_i$$

$$D_{ii} = 1 \text{ if 1st quarter}$$

$$= 0 \text{ otherwise}$$

Problem? We have assigned four dummies to represent four quarters.

Rule: (n-1) dummy - Dummy var. trap

$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + U_i \quad \text{for } i=1, 2, \dots, n; \text{ } i \text{ is indexed for individual}$$

$$y_1 = \alpha + \beta_1 x_{11} + \beta_2 x_{21} + \beta_3 x_{31} + \dots + \beta_k x_{k1} + U_1$$

$$y_2 = \alpha + \beta_1 x_{12} + \beta_2 x_{22} + \beta_3 x_{32} + \dots + \beta_k x_{k2} + U_2$$

$$\vdots$$

$$y_n = \alpha + \beta_1 x_{1n} + \beta_2 x_{2n} + \beta_3 x_{3n} + \dots + \beta_k x_{kn} + U_n$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{21} & x_{31} & \dots & x_{k1} \\ 1 & x_{12} & x_{22} & x_{32} & \dots & x_{k2} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & x_{3n} & \dots & x_{kn} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix} + \begin{pmatrix} U_1 \\ U_2 \\ \vdots \\ U_n \end{pmatrix}$$

So, for that we will go back to our model. So, let us say that this was y_i equals to alpha plus beta 1 x_{1i} plus beta 2 x_{2i} plus dot dot dot beta k x_{ki} plus u_i and where i is basically equals to 1 to n . So that means this equation is for i^{th} individual, i is indexed for individual.

So, what will happen? You will see y_1 if I substitute the values of i for 1, 2, 3, n what will happen? This would become alpha plus beta 1 x_{11} plus beta 2 then you will get x_{21} plus beta 3 x_{31} plus dot dot beta k x_{k1} plus u_1 . Similarly, y_2 would be alpha plus beta 1 then x_{12} plus beta 2 x_{22} plus beta 3 x_{32} plus beta k x_{k2} plus u_2 this will happen.

So that means ultimately you will have if you represent this, what will happen? y_n would become alpha plus beta 2 x_{1n} plus beta sorry this is beta 1, this is beta 1, so this is beta 1, beta 1 x_{1n} plus beta 2 x_{2n} plus beta 3 x_{3n} plus beta k x_{kn} plus u_n for the n^{th} . So, in matrix notation how we represented this?

We represented as $y_1 \ y_2 \ \text{dot dot dot } y_n$ equals to what we did? So, we will take this alpha as common. So, this will become 1 1 1 1 and then we will have $x \ x_{11} \ x_{21} \ x_{31}$ and then x_{k1} then you have x_{12} dot dot dot $x_{1n} \ x_{21}$ sorry this is x_{22} , x_{22} then x_{2n} then x_{32} then you have x_{3n} and

x_2 then x_n . And here you will have α and then $\beta_1, \beta_2 \dots \beta_k, \beta_{k+1}$ plus you have $u_1, u_2, u_3 \dots u_n$. So, this we have discussed.