

# Dummy variables and Multiple Regression

Alessandro Cipriani

Federica Romei

October 20, 2010

1. For this exercise we are going to use the dataset `Openness.dta` that contains inflation and import share data from 1973. The dataset contains all non-centrally planned countries with flexible and fixed exchange rate listed by Summers and Heston for whom data on inflation and openness are available. Inflation is measured as average annual change in the log GDP or GNP deflator since 1973. Openness is measured as the average share of imports in GDP or GNP over the years beginning in 1973. This dataset contains:

- *open*: imports as % GDP, '73-'80
- *inf*: avg. annual inflation, '73-
- *oil*: dummy variable :
  - =1 if major oil producer,
  - =0 otherwise;

This dataset have been used by Romer in order to test if countries more open have lower level of inflation. The data are collecting from 1973 for a specific reason : before the '73, there was the Bretton Woods <sup>1</sup> system that limited the possibility to pursue independent monetary policy.

You will not use other variables hence you don't need the description.

- (a) Construct a statistical procedure to test whether oil countries and non-oil countries have the same level of inflation, against the Alternative Hypothesis that they have different level.

---

<sup>1</sup>The chief features of the Bretton Woods system were an obligation for each country to adopt a monetary policy that maintained the exchange rate of its currency within a fixed value—plus or minus one percent—in terms of gold and the ability of the International Monetary Funds to bridge temporary imbalances of payments. Then, on August 15, 1971 the United States unilaterally terminated convertibility of the dollar to gold. This action created the situation whereby the United States dollar became the sole backing of currencies and a reserve currency for the member states. In the face of increasing financial strain, the system collapsed in 1971.

**Solution:**

First of all, let's state the Null Hypothesis ( $H_0$ ) and the alternative ( $H_1$ ):

$$H_0 : E(inflation|oil = 1) = E(inflation|oil = 0),$$

vs

$$H_1 : E(inflation|oil = 1) \neq E(inflation|oil = 0).$$

We are considering a dual side alternative, but it could also be possible to entertain as alternative the situation in which test whether oil countries tend to have lower or higher inflation than non oil ones. Our variable of interest is

$$\Delta i = E(inflation|oil = 1) - E(inflation|oil = 0).$$

Unfortunately we don't have the real expected values, but under regularity condition we can replace  $\Delta i$  with his sample counterpart, that is an estimator. To be more precise if *inflation* is i.i.d (independent means that level of inflation of the observation  $i$  doesn't depend on the level of inflation of another observation in the sample while identically distributed means that we are assume that all the observations have the same distribution) and  $E(inflation) < \infty$ , we can apply the Law of Large Number, i.e.

$$\widehat{\Delta i} = \overline{inflation_{oil}} - \overline{inflation_{non\ oil}} \xrightarrow{p} \Delta i$$

where  $\overline{inflation_{oil}}$  and  $\overline{inflation_{non\ oil}}$  are the sample averages of the inflation level of oil and non oil countries, respectively. Moreover, if we assume that  $Var(inflation) < \infty$ , we can apply the Central Limit Theorem, which implies that under the null hypothesis,

$$t = \frac{\widehat{\Delta i}}{SE(\widehat{\Delta i})} \xrightarrow{D} N(0, 1)$$

where

$$SE(\widehat{\Delta i}) = \sqrt{SE(inflation_{oil})^2 + SE(Inflation_{non\ oil})^2}.$$

We will reject  $H_0$  if  $|t| > 1.96$  and not reject otherwise. From the data we get:

```
mean inf if oil==1
```

```
Mean estimation      Number of obs   =      7
```

	Mean	Std. Err.	[95% Conf. Interval]	
inf	10.64286	1.325907	7.398485	13.88724

```
. mean inf if oil==0
```

```
Mean estimation      Number of obs   =    107
```

	Mean	Std. Err.	[95% Conf. Interval]	
inf	17.69721	2.387897	12.96297	22.43144

We have all the possible information to compute  $t$ , which is

$$t = \frac{10.64 - 17.69}{\sqrt{1.32^2 + 2.38^2}} = -2.59.$$

Thus, we can reject the null hypothesis at 5%.

- (b) Regress  $inf$  on  $oil$ . How can you interpret  $\beta_1$ ? And  $\beta_0$ ? Do you notice some similarities with the point (a)?

**Solution:** Under the assumptions:

- $inf$  and  $u$  are i.i.d.;
- $E(u|oil) = 0$ ;
- fourth moment of  $inf$  and  $oil$  are well defined;

we can build a model like:

$$inf_i = \beta_0 + \beta_1 oil_i + u_i.$$

We can now interpret the coefficients.  $\beta_0$  and  $\beta_1$  have a specific meaning in this type of regression, indeed :

$$\beta_0 = E(inf|oil = 0)$$

and

$$\beta_1 = E(inf|oil = 1) - E(inf|oil = 0).$$

Testing that  $\beta_1 = 0$  is equivalent to test that  $E(inf|oil = 1) - E(inf|oil = 0) = 0$ , which means that  $\hat{\beta}_1$  is  $\widehat{\Delta i}$  of the previous exercise. While  $\hat{\beta}_0$  corresponds to  $\overline{inflation_{non\ oil}}$ .

You can check easily writing on Stata

reg inf oil, r:

Linear regression

Number of obs = 114

F( 1, 112) = 6.83

Prob > F = 0.0102

R-squared = 0.0050

Root MSE = 24.044

		Robust				
inf		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
oil		-7.054343	2.698782	-2.61	0.010	-12.40163 -1.707051
_cons		17.69721	2.39784	7.38	0.000	12.94619 22.44822

$\hat{\beta}_1$  is  $-7.05$  and means that oil producer countries have on average  $-7.05$  points of inflation than non oil producer countries.

$\hat{\beta}_0$  is the average level of inflation of the non producer oil countries.

$t$  - statistic of  $\beta_1$  is  $-2.61$  and corresponds to the  $t$  - statistic of the previous exercise.

**Pay Attention:**  $t$  is not exactly the same because the variance in the previous exercise have been divided by  $n_{oil}$  and  $n_{non\ oil}$  respectively. Now, indeed, they have been divided by  $n - k$  where  $n = n_{oil} + n_{non\ oil}$  and  $k$  is the number of regressors(included the constant).

We can reject the Null Hypothesis that  $\beta_1$  is equal to zero at 5% significant level. This means that oil producer countries tend to have lower level of inflation than non oil producer countries.

$p$  - value is very low both for  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . However at 1% significant level we won't reject the Null Hypothesis that  $\beta_1 = 0$  (remember that in order to reject  $p$  - value  $< \alpha$  ).

$R^2$  is very low which means that we need further variables to explain the inflation.

2. For this exercise we use Crime.dta file that contains data on arrests during the year 1986 and other information on 2725 men born in either 1960 or 1961 in California. Each man in the sample was arrested at least once prior to 1986. The variables of interest are:

- *narr86*: number of times the man was arrested during 1986, it is zero for most of the sample (72.29%) and it varies from 0 to 12;
- *pcnv*: is the proportion of arrests prior to 1986 that lead to conviction;
- *avgsen*: it is the average sentence length served for prior convictions;
- *ptime86* is the months spent in prison in 1986;
- *qemp86*: is the number of quarters during which the man was employed in 1986 ( from zero to four).

*pcnv* should be interpreted as a proxy for the likelihood to be convicted of a crime while *avgsen* is a measure of expected severity of punishment, if convicted. *ptime86* captures the incarcerative effects of crime: if an individual is in prison, he cannot be arrested for a crime outside the prison. Labor market opportunities are crudely captured by *qemp86*.

- (a) Regress *narr86* on the other variables except *avgsen*. Interpret the coefficients and their significance.

**Solution:** Under the assumptions:

- *pcnv*, *ptime86*, *qemp86* and *u* are iid;
- $E(u|pcnv, ptime86, qemp86) = 0$ ;
- Fourth moments are well defined;

a linear model explaining the arrest can be:

$$narr86 = \beta_0 + \beta_1 pcnv + \beta_2 ptime86 + \beta_3 qemp86 + u.$$

In Stata we can write

```
reg narr86 pcnv ptime86 qemp86, r
```

and visualize this output:

Linear regression

Number of obs = 2725  
 F( 3, 2721) = 36.45  
 Prob > F = 0.0000  
 R-squared = 0.0413  
 Root MSE = .8416

-----						
		Robust				
narr86		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
-----						
pcnv		-.1499274	.0339896	-4.41	0.000	-.2165754 -.0832794
ptime86		-.0344199	.0057706	-5.96	0.000	-.0457352 -.0231047
qemp86		-.104113	.0115364	-9.02	0.000	-.126734 -.0814921
_cons		.7117715	.0409092	17.40	0.000	.6315552 .7919878
-----						

As in the univariate regression case, the first column shows the variables you are using as regressors. The constant,  $\beta_0$  is always in the last row.

The second column shows the coefficients. In the multivariate case the interpretation of the coefficients should be more accurate.

- $\beta_1$ , is -0.14. Suppose to have 2 men that have been not employed and not in prison in 1986. The first man, Mr Brown, has  $pcnv = 0$ , which means that he has never been convicted once arrested and the other, Mr White, has been always convicted, once arrested, which means,  $pcnv = 1$ . Then Mr White should have been arrested -.14 less in 1986 than Mr. Brown. This may be unusual because the arrest cannot decrease of a fraction. Think to have 100 Mr Brown and 100 Mr White. Then the "Mr White" group should have been arrested 14 times less than "Mr Brown" group. This means that to be convicted after the arrest lowers the level of next crime.
- $\beta_2$  is -0.03. Assume to have 2 men who as  $pcnv = qemp86 = 0$ . Assume that the first, Mr Pink, spent 1 month in prison in 1986, while the second, Mr Red, spent 0 month in prison in 1986. Then Mr Pink should have been arrested -0.03 times than Mr Red in 1986. As before, arrests can't decrease of a fraction. Hence, as before, assume to have 100 Mr Pink and 100 Mr Red. Than in the group of Mr Pink there should be 3 arrest less than in Mr Red group.
- $\beta_3$  is -0.1. Now assume to have 2 "groups" made by 100 members. For both group

$pcnv = ptime86 = 0$ . The first group, Mr Green group, works for 1 quarter in 1986 while the second, Mr Black, didn't work. Then in the second group we will observe 10 arrest more in 1986.

The third column shows the Standard Error of the coefficients.

The fourth column shows the  $t$  - *statistic*, as in the univariate case. This  $t$  is built under the Null Hypothesis that each  $\beta_i = 0$ , for  $i = 1, 2, 3$  ( **We are performing 4 different tests, one for each  $\beta_i$** ). If  $|t| > 1.96$ , you will reject the Null Hypothesis at 5% and we can say that the coefficient is significant at 5% level. As you can see all coefficients are significant at 5%.

The fifth column shows the  $p$  - *value*. This  $p$  - *value* is built under the Null Hypothesis that that each  $\beta_i = 0$ , for  $i = 1, 2, 3$ ( **We are performing 4 different tests, one for each  $\beta_i$** ). If  $p - value < \alpha$  you can reject the Null Hypothesis at  $\alpha\%$  significant level.

The sixth column show the Confidence Interval. If this interval contains zero then you cannot say that  $\beta_i$  is significant different from zero at 5% significant level.

In the top-right column you can find  $R^2$ . In this case it is very low, which means that there should be some omitted variables that can help us in explaining the number of times that a person has been arrested in 1986.

- (b) Now do the same regression of point (a) adding *avgsen*. Does this new variable help you in explaining *narr86*? Why?

**Solution:** Adding the new variable we should have this output:

```
reg narr86 pcnv ptime86 qemp86 avgseu ,r
```

Linear regression

Number of obs = 2725  
 F( 4, 2720) = 28.39  
 Prob > F = 0.0000  
 R-squared = 0.0422  
 Root MSE = .84138

		Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
narr86							
pcnv		-.1508318	.0338918	-4.45	0.000	-.2172881	-.0843756
ptime86		-.0373908	.0061239	-6.11	0.000	-.0493987	-.0253829
qemp86		-.103341	.011567	-8.93	0.000	-.1260221	-.0806599
avgseu		.0074431	.0052076	1.43	0.153	-.0027681	.0176543
_cons		.7067564	.0411722	17.17	0.000	.6260246	.7874883

We have to interpret only the  $\hat{\beta}_4$ . Suppose, as before, we have 2 groups of individual. Both group has  $pcnv = ptime86 = qemp86 = 0$  and are composed by 1000 individuals. The first group, Mr Magenta, has a sentence length for prior convictions equal to 1, while the second group, Mr Yellow, has a sentence length equal to zero. The regression predicts that in the first group there should be 7 arrests more. This is not the sign we expected.

$|t - statistic| < 1.96$  , then you cannot reject the Null Hypothesis at 5% significant level.

$p - value > 0.05$  and you cannot reject at 5% significant level that  $\beta_4$  is different from 0. You can also test at 10% and not reject  $H_0$ .

The Confidence Interval contains zero.

$R^2$  doesn't increase too much respect to the other regression, then the new variable doesn't help in explaining the number of arrest in 1986.

(c) Do you think that there is a causal effect between  $narr86$  and  $pcnv$ <sup>2</sup>?

**Solution:** We can find some variables that have correlation with  $narr86$  and  $pcnv$ . For example the income of your family can have a negative effect on both. Also the

<sup>2</sup>Causal effect means that  $E(pcnv, u) = 0$ . There is no variables that has an effect on the proportion of arrests prior to 1986 and contemporary on  $narr86$  except the variables we put in the regression



neighborhood or the district in which you live. This means that there is no a causal effect between *narr86* and *pcvn*.