

# Economics 468

October 26, 2015

R. Davidson

## Midterm Examination

*Do not be upset if this exam seems too long for you! Do not waste time on questions for which you do not see how to obtain the answer. Rather, answer as much as you can. Everything you do will be taken into account.*

1. If two random variables  $X_1$  and  $X_2$  are statistically independent, show that  $E(X_1 | X_2) = E(X_1)$ .

The **covariance** of two random variables  $X_1$  and  $X_2$ , which is often written as  $\text{cov}(X_1, X_2)$ , is defined as the expectation of the product of  $X_1 - E(X_1)$  and  $X_2 - E(X_2)$ . Consider a random variable  $X_1$  with mean zero. Show that the covariance of  $X_1$  and any other random variable  $X_2$ , whether it has mean zero or not, is just the expectation of the product of  $X_1$  and  $X_2$ .

Show that the variance of the random variable  $X_1 - E(X_1 | X_2)$  cannot be greater than the variance of  $X_1$ , and that the two variances are equal if  $X_1$  and  $X_2$  are independent.

Let a random variable  $X_1$  be distributed as  $N(0, 1)$ . Now suppose that a second random variable,  $X_2$ , is constructed as the product of  $X_1$  and an independent random variable  $Z$ , which equals 1 with probability  $1/2$  and  $-1$  with probability  $1/2$ . What is the (marginal) distribution of  $X_2$ ? What is the covariance between  $X_1$  and  $X_2$ ? What is the distribution of  $X_1$  conditional on  $X_2$ ?

2. Show that the  $t^{\text{th}}$  residual from running the regression

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \alpha \mathbf{e}_t + \mathbf{u},$$

is zero. Here  $\mathbf{e}_t$  is the unit basis vector, all of whose elements are zero, except for element  $t$ , which is equal to one. Use this fact to demonstrate that, as a result of omitting observation  $t$ , the  $t^{\text{th}}$  residual from the regression  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  changes by an amount

$$\hat{u}_t \frac{h_t}{1 - h_t},$$

where  $h_t$  is the  $(t, t)$ -diagonal element of the orthogonal projection matrix  $\mathbf{P}_{\mathbf{X}}$ .

Show that the leverage measure  $h_t$  is the square of the cosine of the angle between the unit basis vector  $\mathbf{e}_t$  and its projection on to the span  $\mathcal{S}(\mathbf{X})$  of the regressors.

3. Consider the following linear regression:

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u},$$

where  $\mathbf{y}$  is  $n \times 1$ ,  $\mathbf{X}_1$  is  $n \times k_1$ , and  $\mathbf{X}_2$  is  $n \times k_2$ . Let  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  be the OLS parameter estimates from running this regression.

Now consider the following regressions, all to be estimated by OLS:

- (a)  $\mathbf{y} = \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (b)  $\mathbf{P}_1\mathbf{y} = \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (c)  $\mathbf{P}_1\mathbf{y} = \mathbf{P}_1\mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (d)  $\mathbf{P}_X\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (e)  $\mathbf{P}_X\mathbf{y} = \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (f)  $\mathbf{M}_1\mathbf{y} = \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (g)  $\mathbf{M}_1\mathbf{y} = \mathbf{M}_1\mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (h)  $\mathbf{M}_1\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{M}_1\mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (i)  $\mathbf{M}_1\mathbf{y} = \mathbf{M}_1\mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{M}_1\mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ ;
- (j)  $\mathbf{P}_X\mathbf{y} = \mathbf{M}_1\mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ .

Here  $\mathbf{P}_1$  projects orthogonally on to the span of  $\mathbf{X}_1$ , and  $\mathbf{M}_1 = \mathbf{I} - \mathbf{P}_1$ . For which of the above regressions are the estimates of  $\boldsymbol{\beta}_2$  the same as for the original regression? Why? For which are the residuals the same? Why?

4. Show that the difference between the unrestricted estimator  $\tilde{\boldsymbol{\beta}}$  of the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2\mathbf{I}),$$

and the restricted estimator  $\hat{\boldsymbol{\beta}}$  of the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2\mathbf{I}),$$

is given by

$$\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{M}_Z\mathbf{M}_X\mathbf{y}.$$

**Hint:** In order to prove this result, it is easiest to premultiply the difference by  $\mathbf{X}^\top\mathbf{M}_Z\mathbf{X}$ . Suppose that the random variable  $z$  follows the  $N(0, 1)$  density. If  $z$  is a test statistic used in a two-tailed test, the corresponding  $P$  value is shown to be  $p(z) \equiv 2(1 - \Phi(|z|))$ . Show that  $F_p(\cdot)$ , the CDF of  $p(z)$ , is the CDF of the uniform distribution on  $[0, 1]$ . In other words, show that

$$F_p(x) = x \quad \text{for all } x \in [0, 1].$$

5. Two econometricians obtained data on the sale prices and various characteristics of 546 houses sold in Windsor, Ontario in 1987. The characteristics are:

<code>ls</code>	= the lot size of a property in square feet
<code>bdms</code>	= the number of bedrooms
<code>baths</code>	= the number of full bathrooms
<code>sty</code>	= the number of stories excluding basement
<code>drv</code>	= 1 if the house has a driveway
<code>rec</code>	= 1 if the house has a recreational room
<code>ffin</code>	= 1 if the house has a full finished basement
<code>ghw</code>	= 1 if the house uses gas for hot water heating
<code>ca</code>	= 1 if there is central air conditioning
<code>gar</code>	= the number of garage places
<code>reg</code>	= 1 if the house is located in the preferred neighbourhood of the city

Note that all of the characteristics except the first three are indicator, or dummy, variables, that can take on only the values of zero or one. The sale prices (the variable `value`) were regressed on a constant and all of the characteristics. The results of the regression are as follows:

Ordinary Least Squares:

Variable	Parameter estimate	Standard error	T statistic
constant	-3890.912944	3316.422971	-1.173226
ls	3.535404	0.348470	10.145510
bdms	1775.929742	1034.454519	1.716779
baths	14418.793223	1483.883283	9.716932
sty	6571.231722	921.765886	7.128960
drv	6575.674659	2029.200053	3.240526
rec	4573.381783	1892.691281	2.416338
ffin	5510.572018	1582.046133	3.483193
ghw	12831.052275	3205.286305	4.003091
ca	12431.136062	1552.170690	8.008872
gar	4240.271919	836.987405	5.066112
reg	9469.694699	1662.941237	5.694546

Number of observations = 546    Number of estimated parameters = 12  
 Mean of dependent variable = 67929.289377  
 Sum of squared residuals = 1.260551e+11  
 Explained sum of squares = 2.785264e+12  
 Estimate of residual variance  
   (with d.f. correction) = 2.360582e+08  
 Mean of squared residuals = 2.308701e+08  
 Standard error of regression = 15364.186869  
 R squared (uncentred) = 0.956702    (centred) = 0.678319

This type of regression is called a **hedonic regression**. It attempts to estimate the values that consumers place on various characteristics of a good, in this case a house. By how much, on average, does the price of a house increase when the number of full bathrooms increases by one?

After running the regression, the econometricians plotted the residuals against the fitted values, and obtained the plot shown on the next page. Why does this plot suggest that the regression is not well specified? Can you make suggestions as to how a better fitting regression might be obtained?

