

## **Econometric Analysis of Panel Data**

Spring 2007 – Tuesday, Thursday: 1:00 – 2:20

**Professor William Greene** 

## **Midterm Examination**

This examination has four parts. Weights applied to the four parts will be 15, 15, 30 and 40. This is an open book exam. You may use any source of information that you have with you. You may not phone or text message or email or Bluetooth (is that a verb?) to "a friend," however.

#### Part I. Fixed and Random Effects

Define the two basic approaches to modeling unobserved effects in panel data. What are the different assumptions that are made in the two settings? What is the benefit of the fixed effects assumption? What is the cost? Same for the random effects specification. Now, extend your definitions to a model in which all parameters, not just the constant term, are heterogeneous. For the random parameters case, describe the estimators that one would use under the two assumptions.

Two approaches are fixed effects and random effects. In the "effects model,"  $y_{it} = x_{it}'\beta + c_i + \varepsilon_{it}$ ,  $x_{it}$  is exogenous with respect to  $\varepsilon_{it}$ . FE:  $c_i$  may be correlated with  $x_{it}$ . Benefits: General approach, Robust – estimator of  $\beta$  is consistent even if RE is the right model. Cost: Many parameters, inefficient if RE is correct. Precludes time invariant variables. RE:  $c_i$  is uncorrelated with  $x_{it}$ Benefits: Tight parameterization - only one new parameter Efficient estimation - use GLS Allows time invariant parameters Cost Unreasonable orthogonality assumption Inconsistent if RE is the right model. Random parameters case. Replace the model statement with  $y_{it} = x_{it}'\beta_i + \varepsilon_{it}$ ,  $\beta_i = \beta + w_i$ . Case 1: wi may be correlated with xit. This is the counterpart to FE. In this case, it is necessary to fit the equations one at a time. Requires that there be enough observations to do so, so T > K. The efficient estimator is equation by equation OLS. Same benefits (robustness) and costs (inefficiency) as FE Case 2;  $w_i$  is uncorrelated with  $x_{it}$ . This RP model can be fit An efficient estimator will be the matrix weighted FGLS estimator. (Swamy et al.) This would be a two step estimator, just like FGLS for the RE model. This model can also be fit by simulation - we mentioned this briefly in class, will return to it later this semester.

#### Part II. Minimum Distance Estimation

I have data on 10 firms for 25 years of production. Variables are  $y_{it} = \log of$  value added, and  $\mathbf{x}_{it} = (\log K_{it}, \log L_{it}, \log L_{it})$  where K, L and E are capital, labor and energy. I also have a variable  $d_{it}$  which equals 1 if the firm is in a service industry and 0 if the firm is a manufacturing firm. Note that  $d_{it}$  is time invariant. The model I propose is  $y_{it} = \alpha_i + \mathbf{x}_{it}' \mathbf{\beta} + \delta d_{it} + \varepsilon_{it}$ 

where  $E[\varepsilon_{it}|\mathbf{x}_{js}, d_{js}] = 0$  for all *i*,*t* and *j*,*s*.

 $\mathbf{E}[\varepsilon_{it}\varepsilon_{js}] = \sigma_{ij} \text{ if } t = s \text{ and } 0 \text{ if } t \neq s.$ 

(I.e., firms are correlated but there is no correlation across time.)

I propose to fit this model by the following strategy:

1. Estimate the equation separately for each firm

2. Use a minimum distance estimator to reconcile the 10 competing estimators of  $\beta$ 

1. Does this procedure produce 10 sets of consistent estimators of the parameters of the model? 2. Assuming that  $\sigma_{ij}$  equals zero when  $i \neq j$ . (i.e., no correlation across firms, but different variances), show how to compute the minimum distance estimator.

3. How does the strategy in 2 change if I do not make the assumption that  $\sigma_{ii} = 0$  when  $i \neq j$ .

1. (This is not a trick question – it is the difference between FE and RE.) As the model is laid out above, each single equation treatment can only estimate a firm specific constant ( $\alpha_i + \delta$ ). Since  $d_{it}$  is time invariant, for each firm,  $d_{it}$  is a (second) constant. So, you can regress  $y_{it}$  on  $(1, x_{it})$  to estimate ( $\alpha_i + \delta$ ) and  $\beta$ . This OLS estimator will be consistent for  $\gamma_i = (\alpha_i + \delta)$  and  $\beta$ . It will not be efficient for  $\beta$ , since we have 10 estimators of the same  $\beta$ . That is the point of the question. In fact, it will be an efficient estimator of  $\gamma_i$  under the assumption of part 2, but not of part 3.

2. The constants can now be ignored. We have 10 estimators of  $\beta$ ,  $\mathbf{b}_i$  which has covariance matrix Var[ $\mathbf{b}_i$ ] =  $\sigma_{ii}(\mathbf{X}_i \mathbf{M}^0 \mathbf{X}_i)^{-1} = \Omega_{ii}$ . We can estimate  $\Omega_{ii}$  with  $V_{ii}$  by using  $s_{ii}$  to estimate  $\sigma_{ii}$ . Now, how to reconcile the 10 competing estimators? I proposed an MDE. The MDE would minimize with respect to  $\beta$ 

$$q = \begin{bmatrix} (\mathbf{b}_1 - \boldsymbol{\beta}) & (\mathbf{b}_2 - \boldsymbol{\beta}) & \dots & (\mathbf{b}_{10} - \boldsymbol{\beta}) \end{bmatrix}' \begin{bmatrix} \mathbf{W} \end{bmatrix} \begin{bmatrix} (\mathbf{b}_1 - \boldsymbol{\beta}) \\ (\mathbf{b}_2 - \boldsymbol{\beta}) \\ \dots \\ (\mathbf{b}_2 - \boldsymbol{\beta}) \end{bmatrix}'$$

for some appropriate choice of weighting matrix, **W**. The obvious choice (though inefficient) would be **W** = **I**. As shown in class, if we use this weighting matrix, then the solution to the minimization problem would be  $\hat{\boldsymbol{\beta}} = \overline{\boldsymbol{b}} = \frac{1}{10} \sum_{i=1}^{10} \boldsymbol{b}_i$ . However, this is not an efficient weighting matrix. What we are looking for is the counterpart to Passmore's model where he averaged 4 competing estimators with different variances. The efficient weighting matrix, assuming that there is no correlation across firms will be the inverse of a block diagonal matrix in which the 10 diagonal blocks are  $\Omega_{ii}$ . Since we don't know  $\Omega_{ii}$ , we use  $V_{ii}$ . So,

$$\mathbf{W} = \begin{bmatrix} \mathbf{V}_{11} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_{11} & \dots & \mathbf{0} \\ \dots & \mathbf{0} & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{V}_{11} \end{bmatrix}^{-1}$$

The analytic solution in this case is  $\hat{\boldsymbol{\beta}} = \sum_{i=1}^{10} \mathbf{A}_i \mathbf{b}_i$  where  $\mathbf{A}_i = \left[\sum_{i=1}^{10} \mathbf{V}_{ii}^{-1}\right] \mathbf{V}_{ii}^{-1}$ . Note this is a weighted average in which the weights are matrices that sum of **I**.

3. If there is correlation across firms, then the weighting matrix in 2 is not the most efficient choice, but the strategy is right. We would now have to fill in the off diagonal blocks of **W**. But, again, we would use the same estimator. The off diagonal blocks are  $\mathbf{\Omega}_{ij} = \sigma_{ij} (\mathbf{X}_i' \mathbf{M}^0 \mathbf{X}_j)^{-1} \mathbf{X}_i' \mathbf{M}^0 \mathbf{X}_j (\mathbf{X}_j' \mathbf{M}^0 \mathbf{X}_j)^{-1}$ . (In grading

this, I looked for mention of the off diagonal blocks. The actual  $\Omega_{ij}$  need not appear in the answer for it to be correct.)

#### Part III. Dynamic Model

Consider the dynamic, linear, cross country, random effects regression model

 $y_{it} = \alpha + \beta x_{it} + \delta_i z_{it} + \gamma y_{i,t-1} + u_i + \varepsilon_{it}, t = 1,...,4$  (and  $y_{i,0}$  is observed data).

in which *i* is a country and *t* is a year;  $y_{it}$  is national income per capita,  $z_{it}$  is domestic investment and  $x_{it}$  is a measure of national labor input. You have 30 countries and 4 years of data. Note that the coefficient on  $z_{it}$  is allowed to differ across countries.

(1) Assuming for the moment that  $\delta_i$  is constant across countries, show that the pooled ordinary least squares estimator is inconsistent.

(2) Continuing to assume that  $\delta_i$  is the same for all countries, show two approaches, (1) Anderson and Hsiao and (2) Hausman and Taylor, could be used to obtain consistent estimators of  $\beta$ ,  $\delta$  and  $\gamma$ .

(3) I propose to use a different strategy. Let  $w_{it} = (y_{it} - \alpha + \beta x_{it} + \delta z_{it} + \gamma y_{i,t-1})$ . Consider the set of instruments  $\mathbf{f}_{it} = (1, x_{it}, z_{it}, x_{i,t-1}, z_{i,t-1})$ .

(a) Does the simple strategy of pooling the panel and simply using two stage least squares with **F** as the set of instruments produce a consistent estimator of the parameters? Explain.

(b) I propose to use a GMM estimator based on the moment conditions corresponding to  $E[\mathbf{f}_{it}w_{it}]=0, t=2,3,4.$  Describe in detail how the GMM estimator will proceed. How will this differ from the estimator in part (3a)?

(c) Suppose I extend the strategy in (b) by further assuming "strict exogeneity," that is,  $E[\mathbf{f}_{it}w_{is}] = \mathbf{0}$ , t=2,3,4 and s=2,3,4. How does this change the computations in (b)? (Note and hint: the constant term in  $\mathbf{f}_{it}$  creates some redundant moment conditions. E.g.,  $(1/n)\Sigma_i \mathbf{f}_{i4}w_{i4} = 0$  and  $(1/n)\Sigma_i \mathbf{f}_{i3}w_{i4} = 0$ , both include a term that is  $(1/30)\Sigma_i w_{i4} = 0$ . For purposes of your answer to this question, ignore this fact – in practice, it would be necessary to reduce the set of moment conditions appropriately.)

(4) Now allowing  $\delta_i$  to differ across countries, comment on the consistency of the estimator you used in part (3a). Is it consistent? Can you propose a consistent estimator of this model when  $\delta_i$  varies across countries?

(1) Assuming  $\delta_i$  is constant across countries, the regression is a linear model in which one of the independent variables,  $y_{i,t-1}$  is correlated with the disturbance,  $w_{it} = (u_i + \varepsilon_{it})$ .  $u_i$  is part of the disturbance in the equation for  $y_{i,t-1}$  as well. So, this is a familiar case of an endogenous variable – OLS is inconsistent.

(2) In the Anderson and Hsiao approach, we can use an instrumental variable estimator, as usual. There are many available instruments using lagged values of  $x_{it}$  and  $z_{it}$ , say  $(x_{i,t-1}, z_{i,t-1})$ , or additional lags. A&H suggested taking first differences.  $\Delta y_{it} = \beta(\Delta x_{it}) + \delta(\Delta z_{it}) + \gamma(\Delta y_{i,t-1}) + \Delta \epsilon_{it}$ . This eliminates  $u_i$  from the equation, so in addition to the lags of  $x_{it}$  or lags of  $\Delta x_{it}$  we can use sufficiently lagged values of  $y_{it}$  or  $\Delta y_{i,t-1}$ . For example, if we go back to  $y_{i,t-2}$ , (or  $\Delta y_{i,t-2}$ ) that is far enough that the instrument is not correlated with anything in the differenced equation. The model as stated is also a candidate for the Hausman and Taylor approach. The variable that is correlated with the effect is  $y_{i,t-1}$ . The rest of the model fits precisely into the H and T framework.

(3) (a) It does provide a consistent set of estimators. This is what was suggested at the begining of part (2) above.

(3) (b) The estimator in (3)(a) is equivalent to using GMM while assuming homoscedasticity of the disturbances. The empirical moment condition is

 $E[\mathbf{f}_{it}w_{it}] = \mathbf{0}$  – note this is 5 equations in 4 unknowns

$$\begin{split} & E[1(y_{it} - \alpha + \beta x_{it} + \delta_{i}z_{it} + \gamma y_{i,t-1})] = 0, \\ & E[z_{it}(y_{it} - \alpha + \beta x_{it} + \delta_{i}z_{it} + \gamma y_{i,t-1})] = 0, \\ & E[x_{it}(y_{it} - \alpha + \beta x_{it} + \delta_{i}z_{it} + \gamma y_{i,t-1})] = 0, \\ & E[z_{i,t-1}(y_{it} - \alpha + \beta x_{it} + \delta_{i}z_{it} + \gamma y_{i,t-1})] = 0, \\ & E[x_{i,t-1}(y_{it} - \alpha + \beta x_{it} + \delta_{i}z_{it} + \gamma y_{i,t-1})] = 0, \end{split}$$

and the empirical moment proposed is simply  $\mathbf{m}(\boldsymbol{\beta}) = \sum_i \sum_t \mathbf{f}_{it} w_{it} = \mathbf{0}$ . When we pool the data in this fashion and minimize  $\mathbf{m}(\boldsymbol{\beta})'\mathbf{m}(\boldsymbol{\beta})$ , the resulting estimator is simply 2SLS. The proposed estimator suggests that we use the moment conditions separately for three periods. You can think of this as if we were using periods

2, 3 and 4 separately to estimate the parameters, which we could do using 2sls in each, then averaging the estimators. The suggestion is that we use the moments for the three periods separately. This would imply 15 moment equations,

$$\begin{split} & E[1(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0, \\ & E[x_{i2}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0, \\ & E[z_{i2}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0, \\ & E[x_{i1}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0, \\ & E[z_{i1}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0, \\ & E[1(y_{i3} - \alpha + \beta x_{i3} + \delta_{i}z_{i3} + \gamma y_{i,2})] = 0, \\ & E[x_{i3}(y_{i3} - \alpha + \beta x_{i3} + \delta_{i}z_{i3} + \gamma y_{i,2})] = 0, \\ & E[z_{i3}(y_{i3} - \alpha + \beta x_{i3} + \delta_{i}z_{i3} + \gamma y_{i,2})] = 0, \\ & E[z_{i2}(y_{i3} - \alpha + \beta x_{i3} + \delta_{i}z_{i3} + \gamma y_{i,2})] = 0, \\ & E[z_{i2}(y_{i3} - \alpha + \beta x_{i3} + \delta_{i}z_{i3} + \gamma y_{i,2})] = 0, \\ & E[z_{i2}(y_{i3} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0, \\ & E[1(y_{i4} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0, \\ & E[z_{i4}(y_{i4} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0, \\ & E[x_{i3}(y_{i4} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0, \\ & E[z_{i3}(y_{i4} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0, \\ & E[z_{i3}(y_{i4} - \alpha + \beta x_{i4} + \delta_{i}z_{i4} + \gamma y_{i,3})] = 0. \end{split}$$

The proposed GMM estimator would proceed as follows: We need a preliminary estimator of the parameters, which we computed before using 2SLS. We now need to compute the weighting matrix. We can simply compute  $\mathbf{W} = (1/30)\Sigma_i \mathbf{m}_i \mathbf{m}_i'$  where  $\mathbf{m}_i$  is the 15×1 vector shown explicitly above. Then, the two step GMM estimator is

$$\hat{\boldsymbol{\theta}} = [\mathbf{X}'\mathbf{Z}(\mathbf{W}^{-1})\mathbf{Z}'\mathbf{X}]^{-1}[\mathbf{X}'\mathbf{Z}(\mathbf{W}^{-1})\mathbf{Z}'\mathbf{y}]$$

(In grading this question, I am looking for generalities that suggest the approach. The preceding is more detailed than most of you would be providing even if you had enough time.)

(3) (c) The extended approach would add many additional moment equations. In addition to the preceding, consider just the equations added by  $E[\mathbf{f}_{i3}w_{i2}] = 0$ . These would be

$$E[1(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i3} + \gamma y_{i,1})] = 0,$$
  

$$E[x_{i3}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0,$$
  

$$E[z_{i3}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0,$$
  

$$E[x_{i2}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0,$$
  

$$E[z_{i2}(y_{i2} - \alpha + \beta x_{i2} + \delta_{i}z_{i2} + \gamma y_{i,1})] = 0,$$

Notice that the first of these is already in the set of 15 - it is the first one. But, this adds 4 new moment equations. If we do this for each pair (t,s), we have 4 new moment equations for each of (t,s) = (2,3),(2,4),(3,2),(3,4),(4,2),(4,3), or 6 new sets of 4 moments, for a total of 39. In principle, this would now proceed exactly as we did before, using a  $39 \times 39$  weighting matrix. There is a problem, however. We have only 30 observations. There are not enough observations to proceed in this fashion.

(4) If  $\delta$  differs across countries, then none of the GMM estimators suggested will be consistent, since they estimate only a single  $\delta$ . The only hope is to estimate an equation for each country,

$$y_{it} = \alpha + \beta x_{it} + \delta_i z_{it} + \gamma y_{i,t-1} + u_i + \varepsilon_{it}, t = 1,...,4$$
 (and  $y_{i,0}$  is observed data).

With only 4 observations, this does not look promising. Suppose you could assume that  $\delta_i = \delta + w_i$  where  $w_i$  is orthogonal to the other variables in the model. Then,

$$y_{it} = \alpha + \beta x_{it} + \delta z_{it} + \gamma y_{i,t-1} + u_i + w_i z_i + \varepsilon_{it}, t = 1,...,4$$
 (and  $y_{i,0}$  is observed data).

This is the same model as above, except there is now heteroscedasticity in the random effect. All the same problems as before exist, but the GMM estimators suggested do work – they may be inefficient – in the presence of heteroscedasticity. If, however, it cannot be assumed that  $w_i$  is uncorrelated with everything else, then the cause is lost. There is no consistent estimator.

#### Part IV. Analysis of Panel Data

The following analysis is based on a panel of data on the Swiss railroad network. The data are a panel of observations on 50 railway companies, with numbers of observations per company ranging from 1 to 13. (Frequencies are: 37:13 obs; 8:12 obs; 1:10 obs; 2:7 obs; 1:3 obs; 1:1 obs.) The variables in the data set that are used in the regressions below are as follows:

ID: Company ID from 1 to 51 (50 companies, 605 obs)

YEAR: Year (1985 to 1997)

TOTCOST: Total cost (in 1000 CHF)

NI: Number of years for each company

CT: Total costs adjusted for inflation (1000 CHF)

Q1: Total output in train-kilometers

Q2: Total passenger-output in passenger-kilometers

Q3: Total goods-output in ton-kilometers

PL: Labor price adjusted for inflation (

PK: Capital price using the total number of seats as a proxy for capital stock (CHF per seat)

PE: Price of electricity (CHF per kWh)

STOPS: Number of stations on the network

NARROW\_TRACK: Dummy for the networks with narrow track (1 m wide) The usual width is 1.435m.

RACK: Dummy for the networks with RACK RAIL (cremaillere) in at least some part (used to maintain a slow movement of the train on high slopes)

TUNNEL: Dummy for networks that have tunnels with an average length of more than 300 meters. VIRAGE: Dummy for the networks whose minimum radius of curvature is 100 meters or less.

In the regressions below,

 $\begin{array}{ll} lnCT &= log(totcost/pE) \\ lnpk &= log(pK/pE) \\ lnpl &= log(pL/pE) \\ lnq2 &= log(q2) \\ lnq3 &= log(q3) \\ t &= time trend for year, coded Year - 1984 = 1,2,... \\ \end{array}$ 

The essential model is

$$\begin{split} logC_{it} = & \beta_1 logQ_{2,it} + \beta_2 logQ_{3,it} + \beta_3 logP_{K,it} + \beta_4 logP_{L,it} + \beta_5 logP_{E,it} + \beta_6 t + \beta_7 Stops_{it} \\ & \gamma_1 Virage_i + \gamma_2 Tunnel_i + \gamma_3 Narrow\_T_i + \gamma_4 Rack_i + c_i + \epsilon_{it} \end{split}$$

The constraint that  $\beta_3 + \beta_4 + \beta_5 = 1$  has been built into the estimated model by dividing  $C_{it}$ ,  $P_{K,it}$  and  $P_{L,it}$  all by  $P_{E,it}$  then using logs of the normalized variables in the regression. This is the constraint that imposes **linear homogeneity in the input prices** on the cost function.

Two sets of results are given below. The first set is based on a restricted model in which  $\gamma_1,...,\gamma_4$  all equal zero. That is, the time invariant variables are not included in the model. The second set of results includes the time invariant variables.

1. How would you test the restriction of linear homogeneity in the input prices in the context of the pooled linear regression model? Do the results given below provide the statistics you need to carry out the test? If yes, show how to do it. If not, explain why not -i.e., what you need that is not provided.

The results given do not provide any way to test this hypothesis. We would need either (1) The unrestricted regression, which would give us the sum of squares or  $R^2$  so we could carry out an F test (2) The covariance matrix for the unrestricted regression so we could carry out a Wald test or (3) enough results to carry out an LM test using the restricted regression, which we do not have either.

2. Using the pooled least squares results, test the hypothesis that  $\gamma_1 = \gamma_2 = \gamma_3 = \gamma_4 = 0$ . Can you carry out this test using the fixed effects results? Explain? How would you carry out this test using the random effects results?

With the restriction imposed, the  $R^2$  is .9151612. In the unrestricted regression,  $R^2$  =.9546219. So, F = [(.9546219 - .9151612)/4]/[(1 - .9546219)/(605 - 11)] = 129.135. The critical F statistic with (4,594) degrees of freedom is about 2.39, so the hypothesis is rejected.

3. Based on the results given, which model do you think the analyst should report as their best estimates, the pooled least squares results, the fixed effects results or the random effects results? Justify your answer with the statistical evidence.

In the model without the time invariant regressors, the LM statistic reported is 2941.2 = chi squared with 1 degree of freedom. This is very large and would rule out OLS – the classical model. Then, the Hausman statistic is 63.52 = chi squared with 6 degrees of freedom. This is large. The critical value is about 12.59, so this supports the fixed effects approach.

4. Notice that in the first set of results, the sum of squared residuals for the fixed effects estimator is 3.097795. In the second set of results, where the time invariant variables are added to the regression, the sum of squared residuals given for the fixed effects regression is 3.097795 again!!. Shouldn't the sum of squared residuals decline when variables are added to the regression? Can you explain this strange outcome?

The time invariant variables cannot contribute to the fit of a fixed effects model, since they are all linear combinations of variables that are already in the model – the fixed effects. So, their coefficients will be zero, and they will not change the sum of squares.

5. Using the first set of regression results, test the hypothesis that all the constant terms in the fixed effects model are equal to each other.

The  $R^2$  in the pooled regression with one constant term is 0.9151612. The  $R^2$  in the fixed effects regression is 0.9957743. So, the F statistic is

```
F(49,605-50-6) = [(.9957743 - .9151612)/49]/[(1 - .9957743)/(605 - 50-6)] = 213.740
(This is given in the regression results)
```

The critical F is 2.403, so the hypothesis is rejected.

6. The hypothesis of constant returns to scale is  $\beta_2 + \beta_3 = 1$ . Using the first regression, carry out a test of this hypothesis using the model that you chose in part 3.

We would use the fixed effects model. In the model shown, these are the coefficients on Inq2 and Inq3. The test statistic, using the Wald, or chi squared, would be

 $(.21431433 + .02548159-1)^{2}/(.000878247 + .0000346396 - 2 \times .0000112294) = 649.03$ This is much larger than the critical value of 3.84. Translating to a t statistic, the value would be - 25.48. 7. In a cost function such as this, the assumption that the output variables are exogenous is sometimes justified by an appeal to the regulatory environment in which some regulatory body sets the prices for the firm and they must accept all demand that is forthcoming. The argument works for electricity or gas providers. It probably doesn't work for profit maximizing railroads. In general terms, how would you want to change your estimation strategy to deal with the possibility that these two variables are endogenous in the model.

We would need to find two instrumental variables. It's not clear what these might be. We could only speculate. Wherever they come from, call them  $z_1$  and  $z_2$ , the next step would be 2 stage least squares. Nothing in the statement of the problem suggests that GMM provides any additional benefit.

8. The random effects model in the first results embodies an undesirable assumption of uncorrelatedness of  $c_i$  and the independent variables. The fixed effects model has many coefficients and is inefficient (possibly). The Mundlak approach represents a compromise of these two. Describe how to use Mundlak's estimator in this model.

The Mundlak is based on the proposition that we can project the effects on the means of the exogenous variables, that is,

#### $c_i = means'\gamma + w_i$ .

If we insert this in the fixed effects model, we come up with a random effects model in which the variables are the original time varying variables plus the group means of these variables. In this particular setting, it might make sense to think about  $c_i$  also depending on the time invariant variables listed, which would put them back in the (now random effects) model.

9. After computing the fixed effects model in the first set of results below, I computed the 50 railroad specific intercept terms,  $a_i = \overline{y}_i - \overline{x}'_i \mathbf{b}_{LSDV}$ , i = 1,...,50. This gives me a sample of 50 observations. I then regressed this  $a_i$  on a constant and the railroad specific values of the four time invariant variables listed above. The results were as shown below. How (if at all) does this two step procedure relate to computing the fixed effects estimator and the random effects estimator in the <u>second</u> set of results below? Or, does this regression make no sense at all? What is your interpretation of this model? Is this two step procedure a valid estimator in the context of a particular model? Explain.

It has nothing to do with the fixed effects estimator, since the fixed effects embody all the time invariant information about each railroad. The regression suggests a sort of Mundlak approach to the random effects model, however, based on 8 above. The model that seems to be suggested by the procedure is

$$y_{it} = \alpha_i + x_{it}'\beta + \varepsilon_{it}$$
  
$$\alpha_i = \alpha + z_i'\gamma + u_{j}.$$

We will have to assume that  $u_i$  and  $(z_i, x_{it})$  are uncorrelated. This does define the random effects model. However, note that estimation of the model in two steps is not the same as fitting the model by GLS. Inserting the second equation in the first produces the RE model, which we know consistently estimates  $\alpha, \beta, \gamma$ . Doing this in two steps obtains a consistent estimator of  $\beta$  and an unbiased estimator of  $\alpha_i$ . The implication is that

$$a_i = \alpha_i + w_i$$

where  $E[w_i|z_i] = 0$ . It follows that  $\alpha$  and  $\gamma$  are estimable by OLS. So, this is an alternative, less efficient way to estimate the parameters of the random effects model.

+			+	
Ordinary least squares	regression			
LHS=AI Mean	=	4.227825		
Standard dev	riation =	.5606726	İ	
WTS=none Number of ob	servs. =	50		
Model size Parameters	=	6	İ	
Degrees of f	reedom =	44	İ	
Residuals Sum of squar	res = 2	10.04725	İ	
Standard err	or of e =	.4778563	İ	
Fit R-squared	=	.3477224	i	
Adjusted R-s	guared =	.2735999		
Model test $F[5, 44]$	l(prob) = 4	4.69 (.0016	5)	
+	·		+	
++++		-+	++	+
Variable  Coefficient   S	tandard Error	t-ratio	P[ T >t]	Mean of X
++	1 4 0 0 4 0 6 0	-+	++	+
Constant -4.44420462	.14204060	-31.288	.0000	
STOPS .01119768	.00398614	2.809	.0074	21.1800000
VIRAGE23585127	.36712913	642	.5239	.70000000
TUNNEL .37900057	.19465879	1.947	.0579	.18000000
NARROW_T02249261	.36654719	061	.9513	.66000000
1122722				

# FIRST SET OF RESULTS: TIME INVARIANT VARIABLES OMITTED FROM THE MODEL

-								-		
	OLS With	nout	Group Dummy	/ Variables						
	Ordinary   LHS=LNCI	7 [	least squar Mean	res regress	10n = 1	1.30622				
			Standard o	leviation	= 1	.101691				
	WTS=none	3	Number of	observs.	=	605				
	Model Si	lze	Degrees of	freedom	=	598				
	Residual	s	Sum of squ	lares	= б	2.19436				
			Standard e	error of e	= .	3224964				
	Fit		R-squared	aguarod	= .	9151612				
	   Model te	est	F 6, 5	(dorq) [893	 =1075	.11 (.000	0)			
	Diagnost	ic	Log likeli	Lhood	= -1	70.2812	- /			
			Restricted	l(b=0)	= -9	16.5494				
	   Info ari	tor	Chi-sq [	6] (prob) Prd Crt	=1492	.54 (.000	0)			
		LUEL.	Akaike Inf	to. Criter.	= -2	.251823				
4	' +						+	+		
4	+1						+	F		
	Panei Da 	ita A IIn	nalysis of conditional	LNCT LNCT (No	LONE	wayj ssors)				
	Source	011	Variation	Deg. Free		Mean Squa	re			
	Between		720.242	49	•	14.6988				
	Residual	_	12.8465	555	•	.231468E-	01			
4	10Lai +		/33.089	604	• 	1.213/2 		-		
н	+4		+	+		+	+	+		-+
	Variable	Coe	fficient	Standard	Error	b/St.Er.	P[  +	Z >z]	Mean of	X
	LNQ2		.58153570	.014	63150	39.745	· .	.0000	16.31758	81
	LNQ3		.05791869	.006	40043	9.049		.0000	12.49438	68
	LNPK		.254/59/7	.031	30808	8.137		.0000	13 21935	36
	T I		.00435867	.003	72354	1.171		.2418	5.915702	48
	STOPS		.00892884	.000	96104	9.291		.0000	20.47603	31
	Constant	-	6.99824742	1.166	91897	-5.997		.0000		
1	Least Sc	quare	s with Grou	up Dummy Va	riable	 s				
	Ordinary	7	least squar	res regress	ion					
	LHS=LNCT		Mean	landatian	= 1	1.30622				
	   WTS=none	2	Number of	observs	= 1	.IUI691 605				
	Model si	ze	Parameters	3	=	56				
			Degrees of	freedom	=	549				

Residual     Fit     Model te	Ls Sum of Standar R-squar Adjuste Est F[ 55,	squares d error of e ed d R-squared 549] (prob	= 3.09779 = .751173 = .995774 = .995351 ) =2352.20 (.	5   3E-01   3   0   0000)			
+   Panel:Gr     +	roups Empty Small Avera	0, V est 1, I ge group size	Valid data Largest	50   13   12.10			
++  Variable	Coefficient	+   Standard	Error  b/St.	++ Er. P[ Z >z]	Mean of X		
LNQ2 LNQ3 LNPK LNPL T STOPS	.214314 .025481 .315512 .616695 .003757 .016476	33 .029 59 .009 54 .01 50 .039 73 .000 99 .000	263523         7.           588554         4.           781963         17.           576588         17.           112156         3.           248814         6.	232         .0000           330         .0000           706         .0000           243         .0000           350         .0008           622         .0000	16.3175881 12.4943868 10.1795011 13.2193536 5.91570248 20.4760331		
🎟 Matrix	k - Cov.Mat.						×
[6, 6]	Cell:						
	LNQ2	LNQ3	LNPK	LNPL	Т	STOPS	~
LNQ2	0.000878247	-1.12294e-005	-7.17965e-005	-0.00018965	-1.0367e-005	-1.77992e-005	
LNQ3	-1.12294e-005	3.46396e-005	-8.14176e-006	-1.33476e-006	3.12326e-006	-7.3875e-008	
LNPK	-7.17965e-005	-8.14176e-006	0.000317539	-0.000192832	8.49508e-007	-1.81445e-006	
LNPL	-0.00018965	-1.33476e-006	-0.000192832	0.0012792	-1.16043e-005	-4.58712e-006	
T	-1.0367e-005	3.12326e-006	8.49508e-007	-1.16043e-005	1.2579e-006	-3.51269e-008	
STOPS	-1.77992e-005	-7.3875e-008	-1.81445e-006	-4.58712e-006	-3.51269e-008	6.19086e-006	~
Mc  (1) Cons  (2) Grou  (3) X -  (4) X ar	odel stant term on up effects on variables on nd group effe	Log-Likel: ly -916.9 ly 306.8 ly -170.2 cts 737.0	Lhood Sum of 54938 .73308 32066 .12846 28114 .62194 08990 .30977	Squares F 86930D+03 45922D+02 35608D+02 94979D+01	e-squared   .0000000   .9824763   .9151612   .9957743		
(2) vs (1 (3) vs (1 (4) vs (1 (4) vs (2 (4) vs (2 (4) vs (3	Likelihood R. Chi-squared L) 2446.740 L) 1492.536 L) 3307.279 2) 860.538 3) 1814.742	Hypothesis atio Test d.f. Prob 49 .0000 6 .0000 55 .0000 6 .0000 49 .0000	s Tests F Tes 6 635.027 00 1075.110 00 2352.201 00 287.948 00 213.740	ts num. denom. 49 555 6 598 55 549 6 549 49 549	P value .00000 .00000 .00000 .00000 .00000		
Random E Estimate Lagrange ( 1 df, (High va Baltagi- Fixed vs ( 6 df, (High (1	Effects Model es: Var[e] Var[u] Corr[v(i Multiplier prob value = alues of LM f. Li form of L s. Random Eff prob value = low) values o Sum of S R-square	<pre>: v(i,t) = e ; ,t),v(i,s)] = Test vs. Mode .000000) avor FEM/REM M Statistic = ects (Hausman .00000) f H favor FEN quares d</pre>	<pre>(i,t) + u(i) = .564261D- = .983613D- = .945746 el (3) = 2941 over CR mode = 1802 h) = 63 4 (REM).) .108770D+ .851628D+</pre>	02 01 .42 1.) .20 .52 03 00			
++  Variable	Coefficient	+   Standard	Error  b/St.	Er. P[ Z >z]	Mean of X		
LNQ2 LNQ3 LNPK	.342609 .032116 .302540	63 .024 34 .005 35 .01	450705 13. 562987 5. 753417 17.	980 .0000 705 .0000 254 .0000	16.3175881 12.4943868 10.1795011		

LNPL	.58213153	.03528976	16.496	.0000	13.2193536
т	.00278970	.00109781	2.541	.0110	5.91570248
STOPS	.01802224	.00187232	9.626	.0000	20.4760331
Constant	-5.84523658	.52101033	-11.219	.0000	

Matrix - Cov.Mat.

[7, 7]	Cell:							
	LNQ2	LNQ3	LNPK	LNPL	Т	STOPS	ONE	^
LNQ2	0.000600596	-1.8565e-005	-4.58468e-005	-0.000125011	-7.65049e-006	-1.90247e-005	-0.0070151	8
LNQ3	-1.8565e-005	3.16955e-005	-7.1754e-006	3.2751e-006	2.96026e-006	-7.86842e-007	-6.42196e-005	2
LNPK	-4.58468e-005	-7.1754e-006	0.000307447	-0.000194747	5.30471e-007	-1.6234e-007	0.000282487	8_
LNPL	-0.000125011	3.2751e-006	-0.000194747	0.00124537	-1.19631e-005	-9.09148e-007	-0.0123894	0=
Т	-7.65049e-006	2.96026e-006	5.30471e-007	-1.19631e-005	1.20518e-006	2.09577e-008	0.000233131	8
STOPS	-1.90247e-005	-7.86842e-007	-1.6234e-007	-9.09148e-007	2.09577e-008	3.5056e-006	0.000260404	N
ONE	-0.0070151	-6.42196e-005	0.000282487	-0.0123894	0.000233131	0.000260404	0.271452	
			<u>(((((((</u>			<u>(((((((</u>	(((((((	~

### SECOND SET OF RESULTS: TIME INVARIANT VARIABLES INCLUDED IN THE MODEL

OLS With Ordinary LHS=LNCT WTS=none Model si Residual Fit Model te	hout Group Dummy Variables y least squares regress y Mean Standard deviation e Number of observs. Lze Parameters Degrees of freedom s Sum of squares Standard error of e R-squared Adjusted R-squared est F[ 10, 594] (prob)	sion = 11.30622 = 1.101691 = 605 = 11 = 594 = 33.26614 = .2366508 = .9546219 = .9538580 =1249.60 (.0000)	- - - - - - - - - - - - - - - - - - -
+Danel Da Source Between Residual Total	ata Analysis of LNCT Unconditional ANOVA (No Variation Deg. Free 720.242 49 12.8465 555 733.089 604	[ONE way] o regressors) e. Mean Square 0. 14.6988 5231468E-01 4. 1.21372	+
+  Variable	Coefficient   Standard	Error  b/St.Er. P[	Z >z] Mean of X
LNQ2 LNQ3 LNPK T STOPS VIRAGE TUNNEL NARROW_T RACK Constant	.60397404 .012 .05675679 .006 .43028007 .024 .48044792 .066 .00125984 .002 .01164985 .000 05855252 .053 17749327 .032 18639735 .056 .58275984 .025 -10.1709783 .872	291133       46.779         662610       8.566         171100       17.412         529234       7.247         277632       .454         177905       14.954         349910       -1.094         117998       -5.516         562731       -3.292         598474       22.427         292761       -11.652	.0000 16.3175881 .0000 12.4943868 .0000 10.1795011 .0000 13.2193536 .6500 5.91570248 .0000 20.4760331 .2738 .71570248 .0000 .18842975 .0010 .67603306 .0000 .23471074 .0000
Least So Ordinary LHS=LNCT WTS=none Model si Residual Fit Model te Diagnost	<pre>guares with Group Dummy Va y least squares regress y Mean Standard deviation y Number of observs. y Parameters Degrees of freedom s Sum of squares Standard error of e R-squared Adjusted R-squared est F[ 59, 545] (prob) ic Log likelihood Restricted(b=0)</pre>	riables sion = 11.30622 = 1.101691 = 605 = 60 = 545 = 3.097795 = .7539249E-01 = .9957743 = .9953169 = 2176.75 (.0000) = 737.0899 = -916.5494	

Info c	riter	LogAmemiya Prd Crt	= 3	-5 075537
11110 0.	11001.	Akaike Info. Criter.	=	-5.076191
Estd. 2	Autoco	rrelation of e(i,t)		.663021

\_\_\_\_\_ Panel:Groups Empty 0, Valid data 50 Smallest 1, Largest 13 Average group size 12.10 There are 4 vars. with no within group variation. VIRAGE TUNNEL NARROW\_T RACK Look for huge standard errors and fixed parameters. F.E. results are based on a generalized inverse. They will be highly erratic. (Problematic model.) Unable to compute std.errors for dummy var. coeffs. ·----+ \*----\*---\* |Variable| Coefficient | Standard Error |b/St.Er.|P[|Z|>z]| Mean of X| .0000 16.3175881 .0000 12.4943868 .21431433 .02974378 7.205 .02548159 .00590710 4.314 LNO2 LNO3 
 .00590710
 4.314
 .0000

 .01788490
 17.641
 .0000

 .03589689
 17.180
 .0000

 .00112567
 3.338
 .0008

 .00249726
 6.598
 .0000
 LNPK .31551254 .0000 10.1795011 .61669550 LNPL 13.2193536 .00375773 5.91570248 Т .01647699 20.4760331 STOPS .000000 .....(Fixed Parameter)...... .000000 .....(Fixed Parameter)...... VIRAGE TUNNET. .....(Fixed Parameter)...... .000000 NARROW\_T RACK .000000 .....(Fixed Parameter)..... +----\_\_\_\_\_ Test Statistics for the Classical Model -----Model Log-Likelihood Sum of Squares R-squared .0000000 (1) Constant term only -916.54938 .7330886930D+03 
 (2)
 Group effects only
 306.82066
 .1284645922D+02

 (3)
 X - variables only
 19.00043
 .3326613925D+02

 (4)
 X and group effects
 737.08990
 .3097794979D+01
 .9824763 .9546219 .9957743 Hypothesis Tests Hypotnesis TestsLikelihood Ratio TestF TestsChi-squared d.f. Prob.F num. denom.(2) vs (1) 2446.74049.00000635.02749555(3) vs (1) 1871.10010.000001249.60310594(4) vs (1) 3307.27959.000002176.75459545(4) vs (2)860.53810.00000171.51010545(4) vs (3)1436.17949.00000108.31849545 P value .00000 .00000 .00000 .00000 \_\_\_\_\_ \_\_\_\_\_ \_\_\_\_\_ Error 425: REGR; PANEL. Could not invert VC matrix for Hausman test. -----+ Random Effects Model: v(i,t) = e(i,t) + u(i)Estimates: Var[e] = .568403D-02 Var[u] = .503196D-01 Corr[v(i,t),v(i,s)] = .898506 Lagrange Multiplier Test vs. Model (3) = 2330.30 ( 1 df, prob value = .000000) (High values of LM favor FEM/REM over CR model.) Baltagi-Li form of LM Statistic = 1427.77 Fixed vs. Random Effects (Hausman) = .00 (10 df, prob value = 1.000000) (High (low) values of H favor FEM (REM).) Sum of Squares .664209D+02 R-squared .909412D+00 \_\_\_\_\_ |Variable| Coefficient | Standard Error |b/St.Er.|P[|Z|>z]| Mean of X| .02314079 16.826 .00564514 6.030 .01743458 17.478 .03520706 16.155 .00110142 2.167 .00164611 10.811 LNO2 .38937211 .0000 16.3175881 12.4943868 10.155 .03403754 .0000 LNO3 .30471515 .56876577 .00238707 LNPK .0000 .0000 LNPL 13.2193536 .0302 Т 5.91570248 .0000 STOPS .01779681 20.4760331 .3046 .17329440 -1.027 .71570248 VTRAGE -.17791631 .20298377 TUNNEL .09501443 2.136 .0327 .18842975 -.04830923 .17275122 NARROW\_T -.280 .7797 .67603306 .08205919 .0000 .43134224 5.256 .23471074 RACK Constant -6.44689245 50866766 -12.674 .0000