# ECONOMETRICS I - MIE

## Take-Home Exam 1

Prof. Julia Cajal-Grossi - TA Viktoria Vidahazy

**Deadline**: Final submissions must be sent over the email to the TA and the professor no later than Thursday 27 October at 5 PM.

**Teams**: This take-home exam can be solved in groups of no more than three students. You are welcome to choose your team partners freely. You are also welcome to work on your own if you wish to do so. The list of team members needs to the clearly stated on the front cover of the submission. Note that collaboration across work teams is not permitted in this take home exam.

**Formatting**: Please submit a single .pdf file with the answers to all four questions. If you are not familiar with LaTex, feel free to type in Word and convert to .pdf. The document (excluding a cover page) cannot exceed 10 pages with standard margins, single spacing and a font size of 12.

**Grading**: This take-home exam accounts for 20% of the overall grade of the course. The grading criteria do not depend on the size of the team – 1, 2 or 3 people.

**Plagiarism**: Please, observe the rules against plagiarism at the Institute.

*Contain all variables. Without the one with collinearity problem. ⇒ co variance*

*Test significance.*

**Question 1** (40%)

*Collinarity ⇒ what's the problem of it*

You are an applied economist working for the research division of an international organization. Your team leader asks you to propose an econometric approach for answering the following research questions: *are factory workers more productive when their work environment meets good social compliance standards? if so, how much more? is this different for female and male workers?* In order to do so you are given access to a dataset that was obtained through a survey of 10,000 line operators in garment plants in Dhaka, Bangladesh (i.e. N=10,000). Workers may be employed at the same factory but no two workers work on the same division or production line. All survey data were collected during September 2022, after an audit of the performance of all factories in Dhaka along several dimensions of social compliance. The dataset you are given contains the following variables:

*• All the variables*

*• correlated*

*Discussion: in this question we have this data*
*⇒ cause problem*
*⇒ This model first: results*
*This another model*
*⇒ Do this test*

- worker_id: a numeric identifier for the worker; this variable uniquely identifies an observation in the data
- factory_id: a numeric identifier for the factory of the worker (there are several thousand different factories in the data)
- productivity: the number of pieces of garment the worker produced per hour of work the day before the survey (which was a normal day and should be taken as representative of the worker's productivity)
- product: a categorical variable that contains the specific product the worker was producing the previous day (t-shirts, jeans, sweaters, etc.)
- female: a dummy that takes value one if the worker is a female (zero if male)
- age: the age of the worker in years (between 18 and 65)
- marital status: categorical variable indicating the marital status of the worker (single, married, widowed, etc.)
- education: categorical variable containing the maximum level of education attained by the individual (none, incomplete primary, complete primary, etc.)
- experience: years of experience in their current role in the factory
- household_size: number of people leaving in the same household as the worker
- children: number of children (zero if none)
- factory_size: the number of employees in the factory
- division_size: the number of employees in the division of the worker
- happiness_score: a discrete variable ranging from 1 (extremely unhappy) to 5 (extremely happy) collecting answers from workers to the question *how happy are you with the working conditions in your plant?*
- harassment_pass: dummy that takes value one if the audit showed that the factory complies with good standards around the harassment of workers (zero if the factory is below a good standard)
- voice_pass: dummy that takes value one if the audit showed that the factory complies with good standards around workers' voice and dialogue in the factory (zero if the factory is below a good standard)
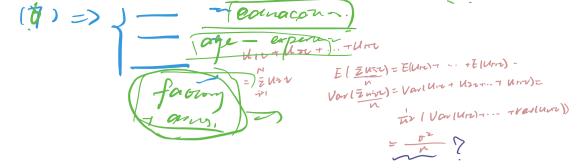- maternity_pass: dummy that takes value one if the audit shows that the factory complies with good standards around maternity and childcare provisions (zero if the factory is below a good standard)
- investment_sc: dollars invested by the factory in 2021 on social compliance activities (note that only 10% of the workers are in factories with positive investments; the rest of the sample shows investments equal to zero)

*• Standardize variables.*

*• Correlations*
*• Isolate effects*

Write a model (or set of models) and an estimation approach for answering the research question you were tasked with. Note that you can only use the dataset at hand and the estimators and tests studied in this course. Please be specific in indicating the variables that you include / exclude in your model(s), how you would transform the variables (if at all), the expected signs and units of measurement of the relationships you are targeting, the construction of any tests you may want to perform, the hypothesis, rejection rules, and degrees of

*y = xβ + c*
*y = xβ + ...*
*xα β + ε*

2

*[handwritten annotations in margins throughout]*

$(\bar{q}) \Rightarrow$ ... $\overline{\text{eqnacomm.}}$
$\overline{age - experair}$
$u_{it} + (u_{zt} +) \dots + u_{nt}$
$= \frac{N}{n}\sum u_{it}$

$\overline{\text{factor}}$
$\overline{\text{array}}$

$E(\bar{z}u_{it}) = E(u_{it}) + \dots + E(u_{nt}) =$
$Var(\bar{z}u_{it}) = \frac{Var(u_{it}) + u_{2t} + \dots + u_{nt})}{n} =$
$\frac{1}{n^2}(Var(u_{it}) + \dots + Var(u_{nt}))$
$= \frac{\sigma^2}{n}$ ?

freedom, etc.. Justify both the economics and the econometrics of your proposed approach. Describe any shortcomings of your exercise and be sure to explain your team leader what caveats they need to have in mind when interpreting the results that your approach would yield.

## Question 2 (20%)

*[handwritten: What happens for a model ⇒ To a model aggregation over i to t specific]*

Suppose that for the population of firms in Switzerland, the relationship over time between dividends and some observable regressors (such as the firm's size) follows the assumptions of the Classical Linear Regression Model, conditional on the realized values of the regressors. Let the model be:

$$y_{it} = x'_{it}\beta + u_{it}$$

*[handwritten: $y_{\Rightarrow t} = \begin{bmatrix} y_1 \\ \vdots \\ y_T \end{bmatrix}$   wher $y_1 = \sum_{i} y_{i1}$]*

with $u_{it}$ i.i.d. (over $i$ and $t$), $\sim \mathcal{N}(0, \sigma^2)$, and independent of $x_{it}$

$$i = 1 \dots N \text{ and } t = 1, \dots T$$

*[handwritten: $\bar{i}, t \rightarrow$ time ; individual]*

Rather than a random sample of firms over time, assume you have a sample of T averages (across all firms) of dividends and regressors – i.e., you have T observations. You estimate an alternative model adapted to the data you have at hand.

*[handwritten: Across firm in each time period ⇒ Only indexed by t]*

*[handwritten: Aagregation over i for T specific]*

1. What would this model look like? Please spell out the model in as much detail as above.

2. Would the OLS estimator of the parameters in the adapted model be unbiased for $\beta$ in the model presented in the introduction? Discuss your reasoning in detail (you do not need to formalize a proof if not needed for your argument).

3. Would it be efficient? Discuss your reasoning in detail (you do not need to formalize a proof if not needed for your argument).

## Question 3 (20%)

*[handwritten: For some observations, observe both x & y]*
*[handwritten: For some observations, only observe x but noty]*

Let a dataset containing $N$ observations be formed by $N_c$ *complete* observations (this is, we have information both for $X$ and $y$) and $N_m$ observations for which we have data on $X$ but $y$ is *missing*, $N = N_c + N_m$. We are interested in the parameters of the model $y = X\beta + u$. A colleague suggests that you proceed as follows: (i) regress $y_c$ on $X_c$ to obtain $\hat{\beta}_c$; (ii) predict $\hat{y}_m$ using $X_m\hat{\beta}_c$; (iii) regress vector $(y_c, \hat{y}_m)$ onto $(X_c, X_m)$ to obtain $\hat{\beta}_{c,m}$.

1. Is $\hat{\beta}_{c,m}$ unbiased for $\beta$? If biased, in what direction? Discuss your reasoning in detail (you do not need to formalize a proof if not needed for your argument).

*[handwritten: $\begin{pmatrix} y_c \\ \hat{y}_m \end{pmatrix} = \begin{pmatrix} X_c \\ X_m \end{pmatrix}\beta + \begin{pmatrix} u_c \\ u_m \end{pmatrix}$]*

*[handwritten: (i) $y_c = X_c\beta + u_c$]*

3

*[handwritten: $\hat{\beta}_c = (X_c'X_c)^{-1}X_c'y_c$   $\hat{y}_m = X_m\hat{\beta}_c$]*

$$\hat{\beta}_{cm} = \left[ (X_c \ X_m)' \begin{pmatrix} X_c \\ X_m \end{pmatrix} \right]^{-1} \left[ \begin{matrix} X_c \\ X_m \end{matrix} \right]' \left[ \begin{matrix} y_c \\ y_m \end{matrix} \right]$$

$$= \ \cdot \ - \ - \ \cdot \ \cdot$$

$$= \hat{\beta}_c$$

✗ biased

2. Is $\hat{\sigma}^2_{c,m}$, the estimator of the variance of the error in the 'imputed' regression unbiased? And, if biased, in what direction? Discuss your reasoning in detail (you do not need to formalize a proof if not needed for your argument).

3. Are there any implications of the approach proposed by your colleague, if you are considering testing the significance of an individual coefficient in your regression?

**Question 4** (20%)

Consider a partitioned regression model, whose normal equations are:

$$X_1' X_1 \hat{\beta}_1 + X_1' X_2 \hat{\beta}_2 = X_1' y$$

$$X_2' X_1 \hat{\beta}_1 + X_2' X_2 \hat{\beta}_2 = X_2' y$$

Let $\hat{\epsilon}$ denote the residuals of the full regression of $y$ on $X$. Let $\tilde{X}_2 = M_1 X_2$ and $\tilde{y} = M_1 y$.

1. Is it true that the residuals in the regression of $\tilde{y}$ on $\tilde{X}_2$ numerically equal $\hat{\epsilon}$? Provide a formal proof.

2. Provide an intuitive explanation for why the statement in (1) is true / is false. *Need to be true or not .because partition regression means . . . ,*

3. Is it true that $\tilde{y}'\tilde{y} - \hat{\epsilon}'\hat{\epsilon} = \tilde{y}' X_2 (X_2' M_1 X_2)^{-1} X_2' \tilde{y}$ ? Provide a formal proof. *after residualize*

4. Provide an intuitive explanation for why the statement in (3) is true / is false. *this over that, there is no variation . . .*

$$\bar{P} = \tilde{X}_2' (\tilde{X}_2' \tilde{X}_2)^{-1} \tilde{X}_2^* \ \tilde{X}_2' \ (\tilde{X} \sim \tilde{X}_2)^{-1} \tilde{X}_2'$$

4