

UNIVERSITY OF OSLO
DEPARTMENT OF ECONOMICS

Term paper in: **ECON4136 – Applied statistical analysis for the social sciences**

Handed out: Monday, October 01, 2012

To be submitted by: Monday, October 22, 2012 at 09:00 a.m.

To be submitted electronically to: submissions@econ.uio.no

Further instructions:

- This term paper is **compulsory**.
- You must use a printed front page, which will be found at the course semester page.
- **Note:** Students are allowed (but not required) to prepare the term paper in groups of max 2 candidates. Candidates who submit term paper together, should submit only one front page – with both names, but separate declarations.
- It is of importance that the term paper is delivered by the deadline (see above). Term papers delivered after the deadline, **will not be corrected**.*)
- You must hand in a declaration form with your term paper. You will find this on the course semester page. **Term papers without declaration forms will not be corrected.**
- Information about citing and referring to sources:
<http://www.uio.no/english/studies/about/regulations/sources/>
- **Information about consequences of cheating:**
<http://www.uio.no/english/studies/admin/examinations/cheating/index.html>
- All term papers must be submitted to the address given above. You must not submit your term paper to the course teacher.
- If the term paper is not accepted, you will be given a new attempt. If you still not succeed, you will not be permitted to take the exam in this course. You will then be withdrawn from the exam, so that this will not be an attempt.

*) If a student believes that she or he has a good cause not to meet the deadline (e.g. illness) she or he should discuss the matter with the course teacher and seek a formal extension. Normally extension will only be granted when there is a good reason backed by supporting evidence (e.g. medical certificate).

Term paper - ECON 4136, fall 2012

For this exercise, use the data set -DahlLochner2012AER.dta- available on the course homepage. You should solve the exercise using Stata. Include your Stata do-file after the main text, tables and figures. Please be brief, but precise, in your answers. Note that you do not have to report more in the text than is asked for. You are allowed to prepare the term paper alone or in groups of 2.

In a recent study published in the *American Economic Review* 2012, 102(5): 1927–1956, Dahl and Lochner (hereafter, DL) study how children’s school performance depends on family income.¹ They posit the following model of the relationship

$$y_{ia} = \mathbf{x}'_i \boldsymbol{\alpha}_a + \mathbf{w}'_{ia} \boldsymbol{\beta} + \delta I_{ia} + u_{ia} \quad (1)$$

where y_{ia} and $I_{i,a}$ are the performance and family income, respectively, of child i at age a . \mathbf{x}_i and \mathbf{w}_{ia} are permanent and time-varying characteristics listed below, while u_i reflects unobserved determinants of school performance.

1. There are three performance measures in the data set -math-, -readingcomp- and -readingrecog-. Create a new variable -score- as the average of these variables, and standardize it to mean equal zero and standard deviation equal one.
2. How much of the variation in -score- and -faminc- is coming from comparisons across individuals and how much is coming from comparisons within individuals over time?
3. Generate new variables for the upper and the lower end of a 95 % confidence interval for -score- and -faminc-, and graph the mean and confidence interval of these variables over time.
4. Estimate model (1) using OLS with -score- as the dependent variable, controlling for variables 9–26 below (i.e. -black- through -sib3-). Use robust standard errors. Interpret the coefficient on -faminc-.
5. Do you think that the OLS-estimates may be biased? Explain your answer. In which direction do you think δ is biased? The Ramsey RESET may test for model misspecification, and can be implemented using -estat ovtest-. Perform the test and explain briefly why it may reveal misspecification.

We have panel data with information on school performance of each child in several years. Assume that the error term above has an individual-specific component μ_i that is fixed over time, such that

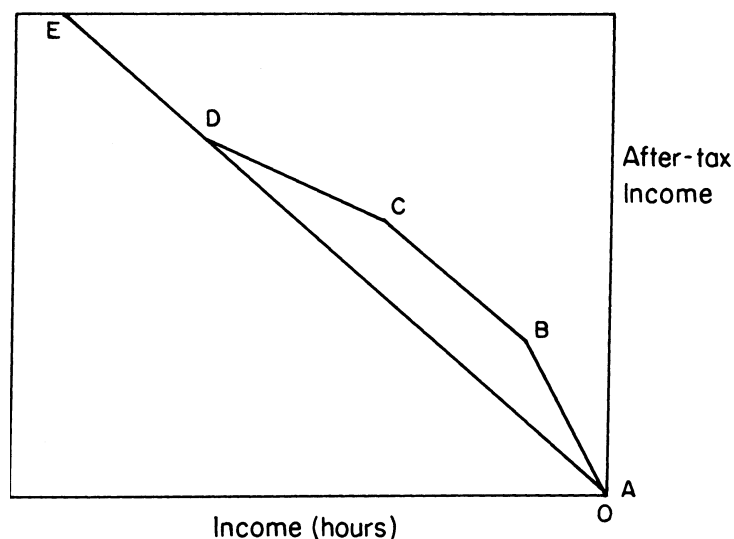
$$u_{ia} = \mu_i + \varepsilon_{ia}$$

¹Notice that the results from these estimations will not match the estimates in the paper, both because part of the data is classified, and because we have simplified their model somewhat.

where ε_{ia} is random residual.

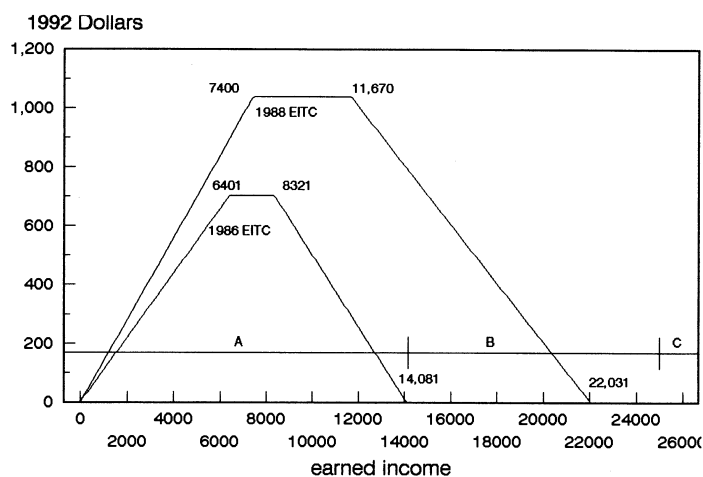
6. Explain how you can use the panel structure of the data to get a more reliable estimate of δ . Estimate this model using first differences for `-score-` and `-faminc-`. Include as control variables `-black-`, `-hispanic-`, `-male-`, `-age-`, `-sib1-`, and `-sib3-` (not differenced).
7. Estimate the model with fixed effects using `-xtreg, fe-`, including the same controls. Why does Stata exclude the variables `-black-`, `-hispanic-`, and `-male-`? How would you interpret the coefficient on these variables in the model in first differences?
8. Why may we be worried about omitted variables bias also in the panel data models? (Hint: What is driving changes in family income?)

The Earned Income Tax Credit (EITC) is a major US transfer program, that provides direct transfers to working families depending on their income, and the number of children. The following figure shows how the EITC changes the budget constraint:



While the EITC and other tax schedules do not generally vary with the child's age in any given year, they do sometimes change over time (that is: with the age of the child, a).

The following figure illustrates this for the 1986 and 1988 EITC in the US:



Total net family income is therefore given by

$$I_{ia} = P_{ia} + \chi_{ia}(P_{ia}) - \tau_{ia}(P_{ia})$$

where P_{ia} is family income prior to taxes and transfers, and χ_{ia} and τ_{ia} are the EITC and tax schedules, respectively.

9. Explain why $\Delta\chi_{ia}(P_{i,a-1}) = \chi_{ia}(P_{i,a-1}) - \chi_{i,a-2}(P_{i,a-2})$ may be an instrument for I_{ia} . Do you think $\Delta\chi_{ia} = \chi_{ia}(P_{i,a}) - \chi_{i,a-2}(P_{i,a-2})$ would be a better or worse instrument for I_{ia} ?
10. In the data, $\chi_{ia}(P_{i,a}) = \text{eitc}$ and $\chi_{ia}(P_{i,a-1}) = \text{eitc_sim}$. Estimate the model in first differences using $\Delta\chi_{ia}(P_{i,a-1})$ as an instrument.
11. Should we be worried about $\Delta\chi_{ia}(P_{i,a-1})$ being a weak instrument?

We may be worried that also $P_{i,a-1}$ is endogenous, since it may be associated with $P_{i,a}$ by e.g. serially correlated shocks. By including in our IV-model flexible controls for $P_{i,a-1}$, we may more plausibly incorporate in our instrument only the changes in I_{ia} deriving from changes in EITC, and avoid incorporating general changes in family income.

12. Reestimate the IV-model in 10 above, including as control variables the dummy -laborpart- and a fifth-order polynomial in -faminc_L1-. Compare the estimates to those you got above.
13. Using this final model, create a loop that estimates the model repeatedly, setting as the dependent variable one of the test score-variables: -mathread-, -math-, -readingcomp-, and -readingrecog-.

Data description:

The file `-DahlLochner2012AER.dta-` contains *biannual* data on school performance and family income in the years 1987–1999, in addition to a number of observable characteristics of the children and their families. Each child is observed at least twice and at most four times. The number of observations equals 7,280, covering 3,692 children. Because data are biannual, it will prove very useful to apply the `-S2-` operator, see `-help tsvarlist-`.

The file includes the following variables:

	Variable	Label
1	<code>id</code>	Id
2	<code>year</code>	Period
3	<code>faminc</code>	Family income (in \$1000, 2000-dollars)
4	<code>eitc</code>	EITC
5	<code>eitc_sim</code>	EITC, simulated
6	<code>math</code>	Mathematics
7	<code>readingrecog</code>	Reading recognition
8	<code>readingcomp</code>	Reading comprehension
9	<code>black</code>	Black
10	<code>hispanic</code>	Hispanic
11	<code>male</code>	Male
12	<code>age</code>	Age
13	<code>agemom</code>	Mother's age
14	<code>ed1age23</code>	Mother high school dropout
15	<code>ed2age23</code>	Mother high school graduate
16	<code>ed3age23</code>	Mother attended college
17	<code>ed4age23</code>	Mother graduated college
18	<code>afqt</code>	Mother's AFQT-score (normalized)
19	<code>afqt_miss</code>	Mother's AFQT-score missing
20	<code>married</code>	Married
21	<code>spouseage</code>	Father's age
22	<code>spouseage_miss</code>	Father's age missing
23	<code>famsize</code>	No. of siblings
24	<code>famsize_miss</code>	No. of siblings missing
25	<code>sib1</code>	One sibling
26	<code>sib3</code>	Two or more siblings
27	<code>laborpart</code>	Labor participation
28	<code>faminc_L1</code>	Family income, 1 year previous
29	<code>faminc_L2</code>	Family income, 2 years previous