

Measuring the Social Return of Higher Education

鄭泊聲、倪楷恩、陳致宇、白崇佑

國立台灣大學

May 27, 2022

- 1 **The Idea**
- 2 **Cross Sectional MLR**
 - College Worker Share
 - City Higher Education Level
- 3 **Instrumental Variable**
- 4 **2018-2020 Panel Data**
 - First Differenced
 - Random Effect
- 5 **1982-2020 Time Series Data**
 - Model 1
 - Model 2
 - Model 3
 - Detrending

External versus Internal

- Government subsidies for higher education is high. Is this justified?
- Basic demand and supply model tells us: government subsidy for external benefits maximize welfare.
- How can we measure the external benefits?
 - Education increases personal wage, but that's **internal**.
 - How about average wage in regions with different amount of higher education? This might include **external** effect.

Much of this work is based on Moretti 2004.

Data and Variables

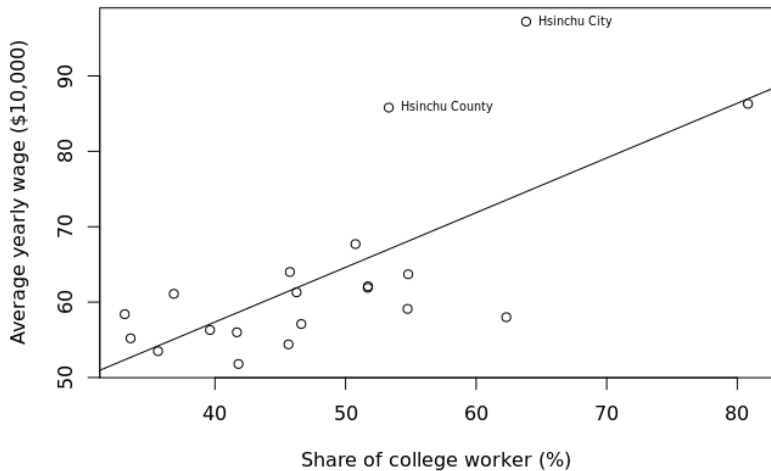
- Wage data: a data based on workplace (instead of household) location city average wage data calculated by DGBAS. Only 2018-2020.
- Education data: How to measure "the amount of higher education"?
 - No. of college graduates in city¹
 - Share of college level worker in city workforce²
 - City population education level - above college share³
- Other city characteristic data as controls

¹Depart. of Statistics, Ministry of Education

²縣市重要統計指標查詢系統, DGBAS

³人口統計資料, Dept. of Household Registration

2020 City Data



Model Specification

Dependent variable is city wage at 2020, $wage_{2020}$. The full MLR model is

$$wage_{2020} = \beta_0 + workforceCollege_{2020} + \beta \mathbf{X} + u \quad (1)$$

with $workforceCollege_{2020}$, the share of college educated worker in city workforce, as main explanatory. \mathbf{X} is a vector containing various city characteristic, including

- *direct*: a dummy for the 6 Special Municipality
- *directEdu2020*: a interaction term between *direct* and *workforceCollege_{2020}* to allow different slope.

All variables

Table: MLR on all variables

<i>Dependent variable:</i>			
wage2020			
workforceCollege_2020	-0.818 (1.814)	eduLevel2020	1.866 (1.746)
direct	81.297** (23.442)	married2020	1.504 (1.141)
hired2020	1.486** (0.608)	expensePerCapita2020	0.003*** (0.001)
manufacture2020	-2.045** (0.850)	unemployment2020	-31.468 (21.214)
service2020	-1.049 (0.915)	directEdu2020	-1.531*** (0.416)
gender2020	2.499** (0.994)	Constant	-268.635 (171.289)
eduExpense2020	1.162* (0.572)	Observations	20
		R ²	0.959
		Adjusted R ²	0.888
		Residual Std. Error	4.055 (df = 7)
		F Statistic	13.604*** (df = 12; 7)
<i>Note:</i>		* p<0.1; ** p<0.05; *** p<0.01	

Joint significance of education

Table:

	<i>Dependent variable:</i>
	wage2020
workforceCollege_2020	0.946 (1.231)
eduExpense2020	0.650 (0.530)
eduLevel2020	-0.333 (1.259)
Constant	9.424 (18.926)
Observations	20
R ²	0.525
Adjusted R ²	0.436
Residual Std. Error	9.118 (df = 16)
F Statistic	5.889*** (df = 3; 16)

Note: *p<0.1; **p<0.05; ***p<0.01

College Worker Share

- 1 The Idea**
- 2 Cross Sectional MLR**
College Worker Share
City Higher Education Level
- 3 Instrumental Variable**
- 4 2018-2020 Panel Data**
First Differenced
Random Effect
- 5 1982-2020 Time Series Data**
Model 1
Model 2
Model 3
Detrending

The Model

We take only reasonable and strong variables to form a compelling model as

$$wage_{2020} = \beta_0 + \beta_1 workforce_{College_2020} + \delta_0 direct + \beta_2 wage_{2018} + \beta_3 manufacture_{2020} + \beta_4 hired_{2020} \quad (2)$$

- *wage*₂₀₁₈: a lagged dependent as proxy to most of the city characteristics
- *manufacture*₂₀₁₈ is the share of manufacturing industry in gross production, *hired*₂₀₂₀ is the share of workforce classified as hired (instead of being employer or self-employed)

MLR

Table:

		Dependent variable:		
		wage2020		
workforceCollege_2020	0.075** (0.029)		Observations	20
			R ²	0.998
direct	-0.194 (0.413)		Adjusted R ²	0.997
			Residual Std. Error	0.689 (df = 14)
wage2018	1.017*** (0.021)		F Statistic	1,175.873*** (df = 5; 14)
manufacture2020	0.025 (0.022)		<i>Note:</i>	* p<0.1; ** p<0.05; *** p<0.01
hired2020	-0.032 (0.038)			
Constant	-1.259 (2.005)			

Heteroskedasticity Robust

Breusch-Pagan test: BP = 10.854, df = 5, p-value = 0.05435

Table:

<i>Dependent variable:</i>	
wage2020	
workforceCollege_2020	0.080*** (0.013)
direct	0.020 (0.240)
wage2018	0.997*** (0.011)
manufacture2020	0.009 (0.022)
hired2020	-0.034 (0.042)
Constant	0.290 (1.613)

Observations	20
R ²	0.998
Adjusted R ²	0.998
Residual Std. Error	0.474 (df = 14)
Note:	*p<0.1; **p<0.05; ***p<0.01

City Higher Education Level

- 1 The Idea**
- 2 Cross Sectional MLR**
 - College Worker Share
 - City Higher Education Level
- 3 Instrumental Variable**
- 4 2018-2020 Panel Data**
 - First Differenced
 - Random Effect
- 5 1982-2020 Time Series Data**
 - Model 1
 - Model 2
 - Model 3
 - Detrending

City Higher Education Level

MLR

Table:

<i>Dependent variable:</i>			
wage2020			
direct	-0.322 (0.421)	Observations	20
manufacture2020	0.022 (0.020)	R ²	0.998
eduLevel2020	0.071** (0.026)	Adjusted R ²	0.997
hired2020	-0.012 (0.033)	Residual Std. Error	0.677 (df = 14)
wage2018	1.015*** (0.021)	F Statistic	1,220.078*** (df = 5; 14)
Constant	-1.967 (1.876)	<i>Note:</i> * p<0.1; ** p<0.05; *** p<0.01	

Heteroskedasticity Robust

Breusch-Pagan test: $BP = 11.97$, $df = 5$, $p\text{-value} = 0.0352$

Table:

	<i>Dependent variable:</i>		
	wage2020		
direct	-0.107 (0.213)	Observations	20
manufacture2020	0.008 (0.021)	R ²	0.999
eduLevel2020	0.080*** (0.015)	Adjusted R ²	0.998
hired2020	-0.017 (0.041)	Residual Std. Error	0.372 (df = 14)
wage2018	0.992*** (0.009)	Note:	*p<0.1; **p<0.05; ***p<0.01
Constant	-0.248 (1.673)		

1 The Idea

2 Cross Sectional MLR

College Worker Share

City Higher Education Level

3 Instrumental Variable

4 2018-2020 Panel Data

First Differenced

Random Effect

5 1982-2020 Time Series Data

Model 1

Model 2

Model 3

Detrending

IV: Lagged Age Structure

Important criteria for an IV, z

- $cov(x, z) \neq 0$: As proportion of college graduates in population grows in time, younger workforce may have more college graduates than older one.
- $cov(u, z) = 0$: Wage is unlikely to be correlated with age.

We use the lagged share of worker aged 15-24, *workforceYoung_2010*, as the IV.

First Stage

With a t-Statistic of -2.199, this choice of IV may not be strong enough.

Table:

	<i>Dependent variable:</i> workforceCollege_2020
workforceYoung_2010	-5.845** (2.658)
Constant	91.140*** (19.525)
Observations	20
R ²	0.212
Adjusted R ²	0.168
Residual Std. Error	10.578 (df = 18)
F Statistic	4.835** (df = 1; 18)
<i>Note:</i>	* p<0.1; ** p<0.05; *** p<0.01



2SLS - College Share

Table:

	<i>Dependent variable:</i>
	wage2020
workforceCollege_2020	0.098 (0.103)
direct	-0.308 (0.650)
wage2018	1.007*** (0.049)
manufacture2020	0.035 (0.048)
hired2020	-0.050 (0.087)
Constant	-0.643 (3.368)

Observations	20
R ²	0.998
Adjusted R ²	0.997
Residual Std. Error	0.704 (df = 14)
Wald test	1125 on 5 and 14 DF, p-value: < 2.2e-16

2SLS - Edu Level

Table:

	<i>Dependent variable:</i>
	wage2020
eduLevel2020	0.097 (0.101)
direct	-0.505 (0.820)
wage2018	1.002*** (0.054)
manufacture2020	0.033 (0.045)
hired2020	-0.026 (0.063)
Constant	-1.483 (2.668)

Observations	20
R ²	0.998
Adjusted R ²	0.997
Residual Std. Error	0.700 (df = 14)
Wald test	1138 on 5 and 14 DF, p-value: < 2.2e-16

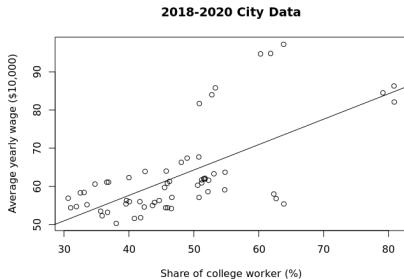
Main Takeaway

- Choosing either share of college worker or share of college city population as main independent yield similar results.
- Either way they have significant positive effect.
- A lagged dependent variable is very explanatory.
- IV analysis suggest that actual effect may be even larger.

- 1 The Idea
- 2 **Cross Sectional MLR**
College Worker Share
City Higher Education Level
- 3 Instrumental Variable
- 4 **2018-2020 Panel Data**
First Differenced
Random Effect
- 5 **1982-2020 Time Series Data**
Model 1
Model 2
Model 3
Detrending

Panel Data

- Same data from 2018-2020.
- Model are specified the same except the the lagged dependent is removed.
- We used an unobserved effect panel data model.



1 The Idea**2 Cross Sectional MLR**

College Worker Share

City Higher Education Level

3 Instrumental Variable**4 2018-2020 Panel Data**

First Differenced

Random Effect

5 1982-2020 Time Series Data

Model 1

Model 2

Model 3

Detrending

First Differenced HAC

Table:

	<i>Dependent variable:</i>
	wageDiff
workforceCollegeDiff	-0.109 (0.086)
direct2	0.341 (0.283)
serviceDiff	0.358*** (0.103)
manufactDiff	-0.028 (0.133)
hiredDiff	-0.092 (0.120)
Constant	0.671*** (0.141)

Observations	40
R ²	0.279
Adjusted R ²	0.173
Residual Std. Error	0.601 (df = 34)
Note:	* p<0.1; ** p<0.05; *** p<0.01

First Differenced with IV

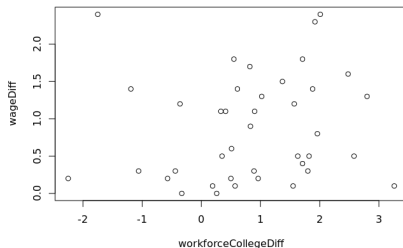
Table:

	<i>Dependent variable:</i>
	wageDiff
workforceCollegeDiff	0.223 (1.203)
direct2	0.294 (0.433)
serviceDiff	0.073 (0.816)
manufactDiff	-0.073 (0.140)
hiredDiff	-0.067 (0.143)
Constant	0.558 (0.552)

Observations	40
R ²	0.024
Adjusted R ²	-0.119
Residual Std. Error	0.748 (df = 34)
Note:	* p<0.1; ** p<0.05; *** p<0.01

First Differenced

Very large standard error in first differencing estimation may be caused by very small variation in explanatory variable.



So we turned to random effect estimator.

Random Effect

- 1 The Idea**
- 2 Cross Sectional MLR**
 - College Worker Share
 - City Higher Education Level
- 3 Instrumental Variable**
- 4 2018-2020 Panel Data**
 - First Differenced
 - Random Effect
- 5 1982-2020 Time Series Data**
 - Model 1
 - Model 2
 - Model 3
 - Detrending

Random Effect

Random Effect

Table:

	<i>Dependent variable:</i>	
	wage2018	
workforceCollege_2018	0.403*** (0.114)	Observations 60
manufacture2018	0.122 (0.127)	R ² 0.539
hired2018	-0.010 (0.169)	Adjusted R ² 0.487
direct	-8.666 (10.479)	F Statistic 62.031***
directEdu2018	0.098 (0.187)	<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01
expensePerCapita2018	0.001*** (0.0002)	
Constant	24.754** (10.832)	

RE with IV - First Stage

The IV is even less significant here.

Table:

	<i>Dependent variable:</i> workforceCollege_2018
workforceYoung_2013	0.082 (0.584)
Constant	47.100*** (4.891)
Observations	60
R ²	0.0003
Adjusted R ²	-0.017
F Statistic	0.020
Note:	* p<0.1; ** p<0.05; *** p<0.01

RE with IV - 2SLS

Table:

	<i>Dependent variable:</i> wage2018		
workforceCollege_2018	1.776 (7.636)	Observations	60
manufacture2018	0.585 (2.097)	R ²	0.352
hired2018	-0.850 (4.705)	Adjusted R ²	0.279
direct	35.758 (264.171)	F Statistic	17.148***
directEdu2018	-0.853 (5.534)	<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	
expensePerCapita2018	0.0003 (0.003)		
Constant	21.315 (22.243)		

Takeaways

- Random effect estimation gives a larger significant coefficient than cross sectional OLS.
- IV using RE also suggest a even larger causal effect.

Other Panel Data Methods

- We find Fixed Effect results vary similar to RE. This can also be seen by the close to one $\hat{\theta}$.
- Pooling independent method finds only the effect of industry structure, *manufacture2018*, has a significant coefficient of 0.4280.

	var	std.dev	share
idiosyncratic	0.5627	0.7501	0.015
individual	37.9300	6.1587	0.985

$\hat{\theta} = 0.9299$

Table: Random effect estimation (without IV)

Remarks

- Our panel data analysis requires strict exogeneity, for which we kind of take it for granted.
- The R^2 in our RE estimation is only 0.487, which definitely leaves room for improvement.

1 The Idea

2 Cross Sectional MLR

College Worker Share

City Higher Education Level

3 Instrumental Variable

4 2018-2020 Panel Data

First Differenced

Random Effect

5 1982-2020 Time Series Data

Model 1

Model 2

Model 3

Detrending

Data

- Between 1982-2020
- National data
- Due to data limitation, after some test we decided to choose:
 - **Average household income** ($income_t$) as the dependent variable.
 - **Government education budget** ($edufund_t$) as the main explanatory variable.

Model 1

- 1 The Idea
- 2 Cross Section MLR
 - College Worker Share
 - City Higher Education Level
- 3 Instrumental Variable
- 4 2018-2020 Panel Data
 - First Differenced
 - Random Effect
- 5 1982-2020 Time Series Data
 - Model 1
 - Model 2
 - Model 3
 - Detrending

Model 1 Specification

$$\begin{aligned} \text{cincome} = & \text{cincome_lag} + \text{cedufund} + \text{cunem} \\ & + \text{cindpd} + \text{cavgGDP} \quad (4) \end{aligned}$$

- *cincome_lag* in a 2 periods lag term of the dependent.
- *cedufund* is our main explanatory.
- *cunem* refers to unemployment rate.
- *cindpd* is the gross production of manufacturing industry, in million NTD.
- *cavgGDP* is the average GDP per capita.

Model 1

Model 1 Results

Table: First Differenced MLR

<i>Dependent variable:</i>			
cincome[3:38]			
cincome_lag	0.555*** (0.121)	Observations	36
cedufund[3:38]	0.0001** (0.00003)	R ²	0.744
cunem[3:38]	-8,050.278*** (1,597.252)	Adjusted R ²	0.702
cindpd[3:38]	-0.001 (0.001)	Residual Std. Error	3,445.707 (df = 30)
cavgGDP[3:38]	0.149** (0.061)	F Statistic	17.451*** (df = 5; 30)
Constant	178.333 (1,437.168)	Note:	* p<0.1; ** p<0.05; *** p<0.01

Serial Correlation

Define u_t as the residuals of the previous model.

Table:

<i>Dependent variable:</i>	
ut	
ut_1	0.023 (0.194)
Constant	-52.015 (553.908)
Observations	35
R ²	0.0004
Adjusted R ²	-0.030
Residual Std. Error	3,265.292 (df = 33)
F Statistic	0.014 (df = 1; 33)
Note:	* p<0.1; ** p<0.05; *** p<0.01

Our model doesn't seem to be affected by SC.

Heteroskedasticity

Breusch-Pagan test:

BP = 5.4892, df = 5, p-value = 0.3591

Our model doesn't seem to be affected by heteroskedasticity.

Model 2

1 The Idea

2 Cross Sectional MLR

College Worker Share
City Higher Education Level

3 Instrumental Variable

4 2018-2020 Panel Data

First Differenced
Random Effect

5 1982-2020 Time Series Data

Model 1
Model 2
Model 3
Detrending

Model 2 Specification

$$\begin{aligned} \text{cincome} = & \text{cincome_lag} + \text{cedufund} + \text{cunem} + \text{cindpd} \\ & + \text{cindpd_lag} + \text{cavgGDP} \quad (5) \end{aligned}$$

- *cincome_lag* is a 2 periods lag term of the dependent.
- *cedufund* is our main explanatory.
- *cunem* refers to unemployment rate.
- *cindpd* is the gross production of manufacturing industry, in million NTD.
- *cindpd_lag* is a 1 period lag term of *cindpd*.
- *cavgGDP* is the average GDP per capita.

Model 2

Model 2 Results

Table: First Differenced MLR model 2

<i>Dependent variable:</i> cincome[3:38]			
cincome_lag	0.553*** (0.112)	cavgGDP[3:38]	0.077 (0.063)
cedufund[3:38]	0.0001*** (0.00003)	Constant	2,159.014 (1,548.312)
cunem[3:38]	-9,850.055*** (1,643.410)	Observations	36
cindpd[3:38]	-0.001* (0.001)	R ²	0.789
cindpd_lag	-0.001** (0.001)	Adjusted R ²	0.745
		Residual Std. Error	3,182.269 (df = 29)
		F Statistic	18.079*** (df = 6; 29)
		Note:	* p<0.1; ** p<0.05; *** p<0.01

Serial Correlation

Define u_t as the residuals of the previous model.

Table:

<i>Dependent variable:</i>	
ut	
ut_1	0.019 (0.198)
Constant	-43.651 (503.820)
Observations	35
R ²	0.0003
Adjusted R ²	-0.030
Residual Std. Error	2,967.877 (df = 33)
F Statistic	0.009 (df = 1; 33)
Note:	* p<0.1; ** p<0.05; *** p<0.01

Our model doesn't seem to be affected by SC.

Heteroskedasticity

Breusch-Pagan test:

BP = 5.3746, df = 6, p-value = 0.4967

Our model doesn't seem to be affected by heteroskedasticity.

Model 3

- 1 The Idea**
- 2 Cross Sectional MLR**
College Worker Share
City Higher Education Level
- 3 Instrumental Variable**
- 4 2018-2020 Panel Data**
First Differenced
Random Effect
- 5 1982-2020 Time Series Data**
Model 1
Model 2
Model 3
Detrending

Model 3 Specification

$$\begin{aligned} \text{cincome} = & \text{cincome_lag} + \text{cedufund} + \text{cunem} \\ & + \text{cservpd} + \text{cavgGDP} \quad (6) \end{aligned}$$

- *cincome_lag* in a 2 periods lag term of the dependent.
- *cedufund* is our main explanatory.
- *cunem* refers to unemployment rate.
- *cservpd* is the gross production of service industry, in million NTD.
- *cavgGDP* is the average GDP per capita.

Model 3

Model 3 Results

Table: First Differenced MLR model 3

<i>Dependent variable:</i>			
cincome[3:38]			
cincome_lag	0.729*** (0.126)	Observations	36
cedufund[3:38]	0.0001** (0.00003)	R ²	0.774
cunem[3:38]	-9,225.299*** (1,609.773)	Adjusted R ²	0.737
cservpd[3:38]	-0.009** (0.004)	Residual Std. Error	3,235.215 (df = 30)
cavgGDP[3:38]	0.201*** (0.063)	F Statistic	20.602*** (df = 5; 30)
Constant	1,517.859 (1,488.989)	Note:	* p<0.1; ** p<0.05; *** p<0.01

Model 3

Serial Correlation

Define u_t as the residuals of the previous model.

Table:

<i>Dependent variable:</i>	
ut	
ut_1	0.119 (0.181)
Constant	-19.639 (517.243)
Observations	35
R ²	0.013
Adjusted R ²	-0.017
Residual Std. Error	3,055.879 (df = 33)
F Statistic	0.435 (df = 1; 33)
<i>Note:</i> * p<0.1; ** p<0.05; *** p<0.01	

Our model doesn't seem to be affected by SC.

Heteroskedasticity

Breusch-Pagan test:

BP = 10.115, df = 5, p-value = 0.07205

Our model doesn't seem to be affected by heteroskedasticity.

Takeaways

- All 3 models yields similar results.
- Government education budget has a positive, significant, but small effect.
- The models are free from heteroskedasticity and serial correlation.
- R^2 s lies in the 0.7 range, still room for improvement.
- The main explanatory accounts for all education level, not just higher.
- The dependent also doesn't directly indicate wage.

Comparisons

In comparison to our cross sectional analysis,

- TS obtains smaller coefficient,
- TS models have lower R^2 .

In comparison to our panel analysis,

- TS obtains smaller coefficient,
- TS models have higher R^2 .

From several models using cross section, panel, time series methods, we find (higher)education does have a positive effect on average wage. How much the effect is distributed to the external is beyond the scope of this work.

Detrending

- 1 The Idea**
- 2 Cross Sectional MLR**
 - College Worker Share
 - City Higher Education Level
- 3 Instrumental Variable**
- 4 2018-2020 Panel Data**
 - First Differenced
 - Random Effect
- 5 1982-2020 Time Series Data**
 - Model 1
 - Model 2
 - Model 3
 - Detrending

Detrending MLR

Table:

<i>Dependent variable:</i>	
income	
edufund	0.0003*** (0.00002)
unem	-4,634.824*** (1,311.579)
indpd	-0.004*** (0.001)
avgGDP	0.275*** (0.066)
t	-2,095.528 (1,398.582)
Constant	32,611.320*** (5,618.118)

Observations	39
R ²	0.996
Adjusted R ²	0.995
Residual Std. Error	6,269.286 (df = 33)
F Statistic	1,575.655*** (df = 5; 33)
Note:	* p<0.1; ** p<0.05; *** p<0.01

Detrended AR

First, we obtain the detrended variables \ddot{y}_t by regressing

$$y_t = \alpha_0 + \alpha_1 t + e_t \quad (7)$$

and specify $\ddot{y}_t = e_t$.

Subsequently, we found detrended \ddot{income}_t highly correlated to \ddot{income}_{t-1} with correlation 0.9570.

So we performed first differencing and define

$$cincome_dt_t = \ddot{income}_t - \ddot{income}_{t-1} \quad (8)$$

Detrending

FD Detrending MLR

Then the regression result is **exactly the same** as the non-detrending FD MLR (in table 16).

Table: First differenced detrending MLR

<i>Dependent variable:</i>			
cincome_dt[4:38]			
cincome_dt[2:36]	0.558*** (0.122)	Observations	35
cedufund_dt[4:38]	0.0001** (0.00003)	R ²	0.741
cunem_dt[4:38]	-8,154.249*** (1,623.290)	Adjusted R ²	0.696
cindpd_dt[4:38]	-0.001 (0.001)	Residual Std. Error	3,482.437 (df = 29)
cavgGDP_dt[4:38]	0.154** (0.062)	F Statistic	16.602*** (df = 5; 29)
Constant	152.261 (602.216)	Note:	* p<0.1; ** p<0.05; *** p<0.01

SC, Heteroskedasticity and Our Question

- Serial correlation and heteroskedasticity result are also the same.
- So, is detrending redundant after first-differencing?

Reference I

- [1] Enrico Moretti. "Estimating the social return to higher education: evidence from longitudinal and repeated cross-sectional data". In: *Journal of Econometrics* 121.1 (2004). Higher education (Annals issue), pp. 175–212. ISSN: 0304-4076. DOI: <https://doi.org/10.1016/j.jeconom.2003.10.015>. URL: <https://www.sciencedirect.com/science/article/pii/S0304407603002653>.