

# Problem Set 3: Instrumental Variables Key

Claire Duquennois

*Name:*

## 1 Empirical Analysis using Data from Ananat (2011, AEJ:AE)

This exercise uses data from Elizabeth Ananat’s paper, “The Wrong Side(s) of the Tracks: The Causal Effects of Racial Segregation on Urban Poverty and Inequality,” published in the *American Economic Journal: Applied Economics* in 2011. This paper studies how segregation has affected population characteristics and income disparity in US cities using the layout of railroad tracks as an instrumental variable.

## 2 Finding the data

The data can be found by following the link on the AEJ: Applied Economics’ website which will take you to the ICPSR’s data repository. You will need to sign in to get access to the data files. Once logged in, you will find the set of files that are typically included in a replication file. These include several datasets, several .do files (which is a STATA command file). For this assignment we will be using the `aej_maindata.dta` file.

### 3 Set up and opening the data

3.1 Question: Load any packages you will need and the data contained in the `aej_maindata.dta` file. How many observations are contained in the data. What is the level of an observation?

Code and Answer:

```
#install.packages('haven',repos = "http://cran.us.r-project.org")  
#install.packages("dplyr",repos = "http://cran.us.r-project.org")
```

```
library(haven)  
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 4.1.2
```

```
library(stargazer)  
library(lfe)  
library(ggplot2)
```

```
mydata<-read_dta("../../PSdata/PS3_IVdata//data/aej_maindata.dta")  
mydata<-as.data.frame(mydata)  
nrow(mydata)
```

```
## [1] 121
```

The data contains 121 observations where each observations is for one city/town.

## 4 Data Description

**4.1 Question:**The dataset contains many variables, some of which are not used in this exercise. Keep the following variables in the final dataset (Hint: use the `select` function in `dplyr`).

Name	Description
dism1990	1990 dissimilarity index
herf	RDI (Railroad division index)
lenper	Track length per square km
povrate_w	White poverty rate 1990
povrate_b	Black poverty rate 1990
area1910	Physical area in 1910 (1000 sq. miles)
count1910	Population in 1910 (1000s)
ethseg10	Ethnic Dissimilarity index in 1910
ethiso10	Ethnic isolation index in 1910
black1910	Percent Black in 1910
passpc	Street cars per capita 1915
black1920	Percent Black 1920
lfp1920	Labor Force Participation 1920
incseg	Income segregation 1990
pctbk1990	Percent Black 1990
manshr	Share employed in manufacturing 1990
pop1990	Population in 1990

You can find the detailed description of each variable in the original paper.

**Code:**

```
mydata2<-mydata[, c("dism1990", "herf", "lenper", "povrate_w", "povrate_b", "area1910",  
                    "count1910", "ethseg10", "ethiso10", "black1910", "passpc", "black1920",  
                    "lfp1920", "incseg", "pctbk1990", "manshr", "pop1990" )]
```

## 5 Summary Statistics:

5.1 Question: Report summary statistics of the following variables in the dataset: “dism1990”, “herf”, “lenper”, “povrate\_w”, “povrate\_b”. Present these summary statistics in a formatted table, you can use `stargazer` or other packages.

Code:

```
stargazer(mydata2[,c("dism1990", "herf", "lenper", "povrate_w", "povrate_b")])
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
% Date and time: Wed, Oct 18, 2023 - 12:52:47 PM

Table 2:

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
dism1990	121	0.569	0.135	0.329	0.457	0.673	0.873
herf	121	0.723	0.141	0.238	0.638	0.830	0.987
lenper	121	0.001	0.001	0.0002	0.0004	0.001	0.013
povrate_w	121	0.095	0.035	0.035	0.069	0.114	0.216
povrate_b	121	0.264	0.080	0.093	0.209	0.313	0.504

## 6 Reduced Form:

**6.1 Question:** We are interested in understanding how segregation affects population characteristics and income disparity in US cities. We will focus on two outcome variables: the poverty rate for blacks and whites. Regress these two outcome variables on segregation in 1990, our explanatory variable, and interpret your results. Report robust standard errors. Make sure you specify the units of measurement in your interpretation.

Code:

```
reg1<-felm(povrate_w~dism1990, mydata2)
reg2<-felm(povrate_b~dism1990, mydata2)

stargazer(reg1,reg2,se = list(reg1$rse, reg2$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Wed, Oct 18, 2023 - 12:52:47 PM

Table 3:

	<i>Dependent variable:</i>	
	povrate_w	povrate_b
	(1)	(2)
dism1990	-0.073*** (0.019)	0.182*** (0.045)
Constant	0.136*** (0.012)	0.161*** (0.029)
Observations	121	121
R <sup>2</sup>	0.081	0.095
Adjusted R <sup>2</sup>	0.074	0.088
Residual Std. Error (df = 119)	0.033	0.076
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

**Answer:** For the explanatory variable, one standard deviation is about 0.14 as reported in the summary statistics. It is helpful for interpretation to interpret these effects in terms of standard deviations of the explanatory variable. Thus, we can see that a one standard deviation increase in the segregation index is associated with a  $(0.14 * (-0.073) = -0.0102)$ , one percentage point decrease in white poverty and a  $(0.14 * (0.182) = 0.025)$ , 2.5 percentage point increase in black poverty.

**6.2 Question: Explain the problem with giving a causal interpretation to the estimates you just produced. Give examples of specific factors that might make a causal interpretation of your result problematic.**

**Answer:** There are many problems with giving these estimates a causal interpretation. Omitted Variable Bias is a particular concern. There are many variables that are jointly correlated with both segregation and poverty, such as political corruption or the presence of industrial sectors that are associated with occupational segregation to name just a few. Omission of any of these variables from our regression could lead to biased causal estimates. In addition, there is also the concern of selection as segregation may induce selective migration. Finally, reverse causality could also be a concern as cities with greater black poverty could elect to implement more segregating policies.

## 7 Validity of the instrument:

### 7.1 Question: Estimate the following regression and interpret its coefficients,

$$\text{dism1990}_i = \beta_0 + \beta_1 \text{RDI}_i + \beta_2 \text{tracklength}_i + \epsilon.$$

Code:

```
regfirststage<-felm(dism1990~herf+lenper, mydata2)
stargazer(regfirststage,se = list(regfirststage$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
% Date and time: Wed, Oct 18, 2023 - 12:52:47 PM

Table 4:	
	<i>Dependent variable:</i>
	dism1990
herf	0.357*** (0.088)
lenper	18.514* (10.731)
Constant	0.294*** (0.064)
Observations	121
R <sup>2</sup>	0.203
Adjusted R <sup>2</sup>	0.189
Residual Std. Error	0.122 (df = 118)
Note:	*p<0.1; **p<0.05; ***p<0.01

**Answer:** For the explanatory variable, one standard deviation is, here too, about 0.14 as reported in the summary statistics. It is helpful for interpretation to interpret these effects in terms of standard deviations of the explanatory variable. Thus, we can see that a one standard deviation increase in the RDI is associated with a  $(0.14 * (0.357) = 0.049)$ , a 5 point increase in the segregation index, which is about 0.37 standard deviations.

**7.2 Question:** Re-estimate the specification above using the `scale()` command around the variables you wish to standardize in the regression. What do you notice?

Code:

```
regfirststagesd<-felm(scale(dism1990)~scale(herf)+lenper, mydata2)
stargazer(regfirststagesd,se = list(regfirststagesd$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Wed, Oct 18, 2023 - 12:52:47 PM

Table 5:	
	<i>Dependent variable:</i>
	scale(dism1990)
scale(herf)	0.374*** (0.092)
lenper	136.911* (79.355)
Constant	-0.123 (0.106)
Observations	121
R <sup>2</sup>	0.203
Adjusted R <sup>2</sup>	0.189
Residual Std. Error	0.901 (df = 118)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

**Answer:**

This makes interpretation in terms of standard deviations much more straight forward.



**7.3 Question: In the context of instrumental variables, what is this regression referred to as and why is it important?**

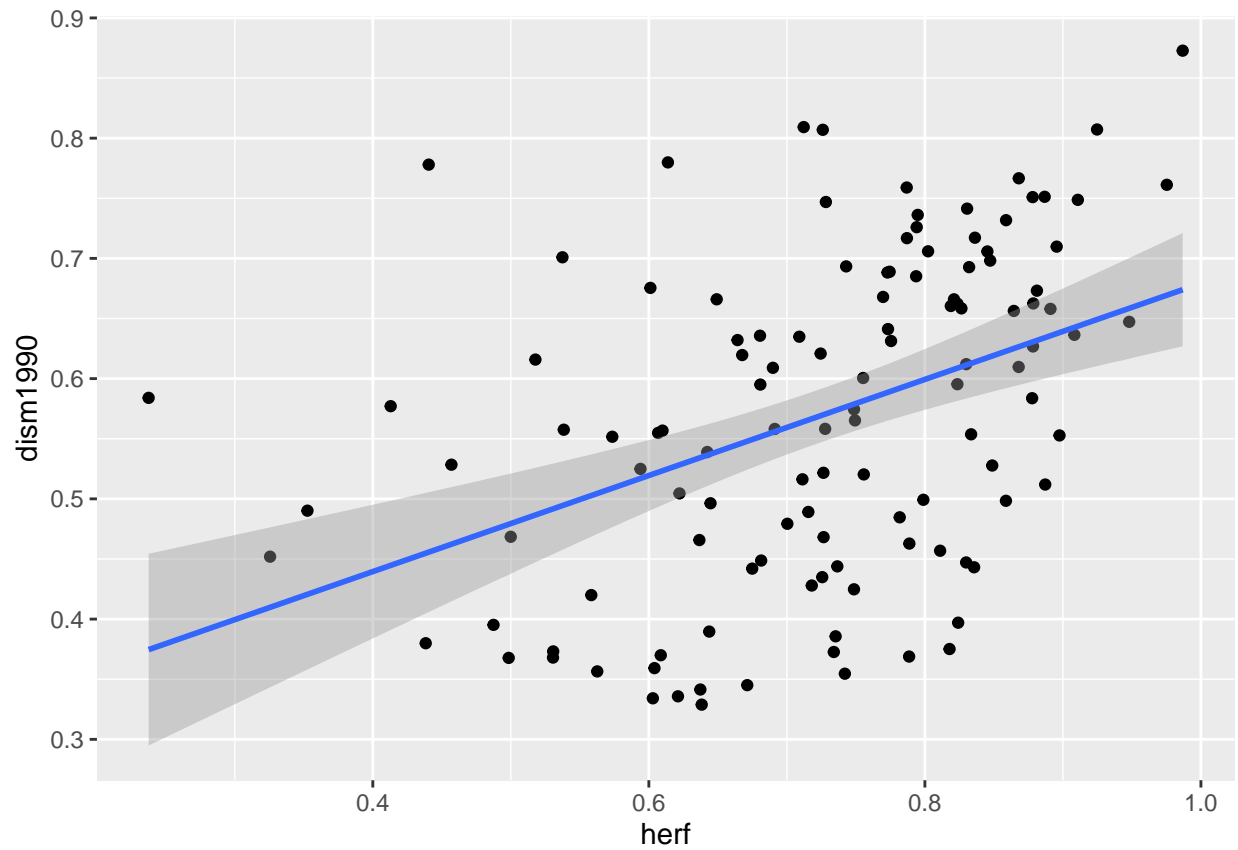
**Answer:** This regression is referred to as the first stage regression. It is a regression of the instrument on the endogenous explanatory variable. It is important since for the RDI to work as an instrument, it needs to have some explanatory power over the endogenous variable, segregation.

#### 7.4 Question: Illustrate the relationship between the RDI and segregation graphically.

Code:

```
myplot1<-ggplot(mydata, aes(x=herf, y=dism1990)) +  
  geom_point()+  
  geom_smooth(method='lm', aes(y = dism1990))  
myplot1
```

## 'geom\_smooth()' using formula 'y ~ x'



**7.5 Question: Is there a concern that this might be a weak instrument? Why would this be a problem?**

**Answer:**

```
summary(regfirststagesd)$fstat
```

[1] 14.98272

The F-statistic for the first stage model is greater than 10, the benchmark used to detect weak instruments so the weak instrument problem likely does not apply here. Weak instruments are a problem because if the instrumental variable does not have good predictive power over the endogenous variable, any small bias that results from a violation of the exclusion restriction gets magnified to generate large bias.

## 7.6 Question: Select a number of relevant city characteristics in the data to regress on the RDI and track length. Present your results and interpret your findings. Why do these results matter for answering our question of interest?

Code and Answer:

```

regex1<-felm(area1910~herf+lenper, mydata2)
regex2<-felm(count1910~herf+lenper, mydata2)
regex3<-felm(black1910~herf+lenper, mydata2)
regex4<-felm(incseg~herf+lenper, mydata2)
regex5<-felm(lfp1920~herf+lenper, mydata2)

stargazer(regex1,regex2,regex3,regex4,regex5,se = list(regex1$rse,regex2$rse,regex3$rse,regex4$rse,regex5$rse))

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Wed, Oct 18, 2023 - 12:52:48 PM

```

Table 6:

	<i>Dependent variable:</i>				
	area1910 (1)	count1910 (2)	black1910 (3)	incseg (4)	lfp1920 (5)
herf	-3,992.637 (11,986.490)	665.751 (1,362.964)	-0.001 (0.010)	0.032 (0.032)	0.028 (0.024)
lenper	-574,401.000 (553,669.000)	75,553.190 (134,814.900)	9.236*** (0.650)	-2.504 (1.626)	-3.427** (1.500)
Constant	18,409.570** (8,612.320)	976.876 (927.189)	0.007 (0.007)	0.196*** (0.025)	0.401*** (0.018)
Observations	58	121	121	69	121
R <sup>2</sup>	0.007	0.006	0.290	0.028	0.015
Adjusted R <sup>2</sup>	-0.029	-0.011	0.278	-0.001	-0.002
Residual Std. Error	15,050.340 (df = 55)	1,903.415 (df = 118)	0.018 (df = 118)	0.032 (df = 66)	0.042 (df = 118)

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

The RDI does not have a statistically significant effect on any of these city characteristics, holding the length of rail lines in the city constant. These regressions serve to support our assumption that the exclusion restriction holds. For our instrument to be valid, the only way that the RDI affects current poverty must be through its effect on segregation. This is an assumption that cannot be tested directly. To determine whether or not we believe that it holds, we must evaluate the arguments presented by the author and evidence such as the results of these regressions. These regressions suggest that the RDI is not associated with city characteristics prior to the Great Migration and were not associated with different population characteristics at the beginning of the great migration, before segregation began to shape population characteristics.

**7.7 Question:** What are the two conditions necessary for a valid instrument? What evidence do you have that the RDI meet these conditions? Be specific in supporting this claim.

**Answer:** The two conditions are that we need a first stage and for the exclusion restriction to be satisfied:

1)  $cov(z, x_i) \neq 0$

2)  $cov(z, \nu) = 0$

We have seen above that there is a first stage and that the first stage is quite strong. Regarding the exclusion restriction, Ananat provides some evidence that it should hold through her discussion of the history of railroads and the regressions presented in table 1 though some concerns may still be valid.

## 7.8 Question: Do you believe the instrument is valid? Why/why not?

**Answer:** Being able to study how segregation impacts communities is an important and interesting question but finding a source of exogenous variation for segregation is not easy. The instrument proposed by Ananat may not be perfect but does provide plausibly exogenous variation with which to tackle this question. We have to be convinced that the way the railroads were laid out was not influenced by socio-economic characteristics that could explain patterns seen today and we also have to believe that the rail configuration is not correlated with who migrated to the city during the great migration in a way that she cannot detect. If, for instance, certain types of manufacturing plants are associated with black-white disparity were more likely to be located in cities with certain types of rail configurations (due to topography for example) this could violate the exclusion restriction. While there are certainly stories that could be told that would challenge the validity of this instrument, it is a plausible source of exogenous variation in segregation, which is no small feat.

## 7.9 Question: Generate a table that estimates the effect of segregation on the poverty rate for blacks and whites by OLS and then using the RDI instrument. Make sure you report robust standard errors. How does the use of the RDI instrument change the estimated coefficients?

Code and Answer:

```
regmain1<-felm(povrate_w~dism1990, mydata2)
regmain2<-felm(povrate_b~dism1990, mydata2)
regmain3<-felm(povrate_w~lenper|0|(dism1990~herf+lenper), mydata2)
```

```
## Warning in chol.default(mat, pivot = TRUE, tol = tol): the matrix is either
## rank-deficient or indefinite
```

```
regmain4<-felm(povrate_b~lenper|0|(dism1990~herf+lenper), mydata2)
```

```
## Warning in chol.default(mat, pivot = TRUE, tol = tol): the matrix is either
## rank-deficient or indefinite
```

```
stargazer(regmain1,regmain2,regmain3,regmain4, se = list(regmain1$rse,regmain2$rse,regmain3$rse,regmain4$rse))
```

```
% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Wed, Oct 18, 2023 - 12:52:48 PM
```

Table 7:

	<i>Dependent variable:</i>			
	povrate_w	povrate_b	povrate_w	povrate_b
	(1)	(2)	(3)	(4)
dism1990	-0.073*** (0.019)	0.182*** (0.045)		
lenper			0.602 (1.970)	-4.780 (3.067)
‘dism1990(fit)’			-0.196*** (0.065)	0.258** (0.108)
Constant	0.136*** (0.012)	0.161*** (0.029)	0.205*** (0.037)	0.121** (0.061)
Observations	121	121	121	121
R <sup>2</sup>	0.081	0.095	-0.150	0.084
Adjusted R <sup>2</sup>	0.074	0.088	-0.170	0.068
Residual Std. Error	0.033 (df = 119)	0.076 (df = 119)	0.037 (df = 118)	0.077 (df = 118)

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

Using the RDI instrument leads to much larger estimated effects than when using the simple OLS model. The effect of 1 unit on the segregation index on white poverty rate goes from -0.073 to -0.196 whereas the effect on black poverty rate goes from 0.182 to 0.258, a substantial change in both cases.

### 7.10 Question: What is the reduced form equation?

**Answer:** The reduced form equation is the regression of the outcome variable directly on the instrument and any other exogenous variables. In this case it is

$$Y_i = \pi_0 + \pi_1 RDI_i + \eta$$



## 7.11 Question: (2 pages) For the two poverty rates, estimate the reduced form on all the cities and illustrate the reduced form relationships graphically.

Code:

```
regrf1<-felm(povrate_w~herf+lenper, mydata2)
regrf2<-felm(povrate_b~herf+lenper, mydata2)

stargazer(regrf1,regrf2,se = list(regrf1$rse,regrf2$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Wed, Oct 18, 2023 - 12:52:48 PM

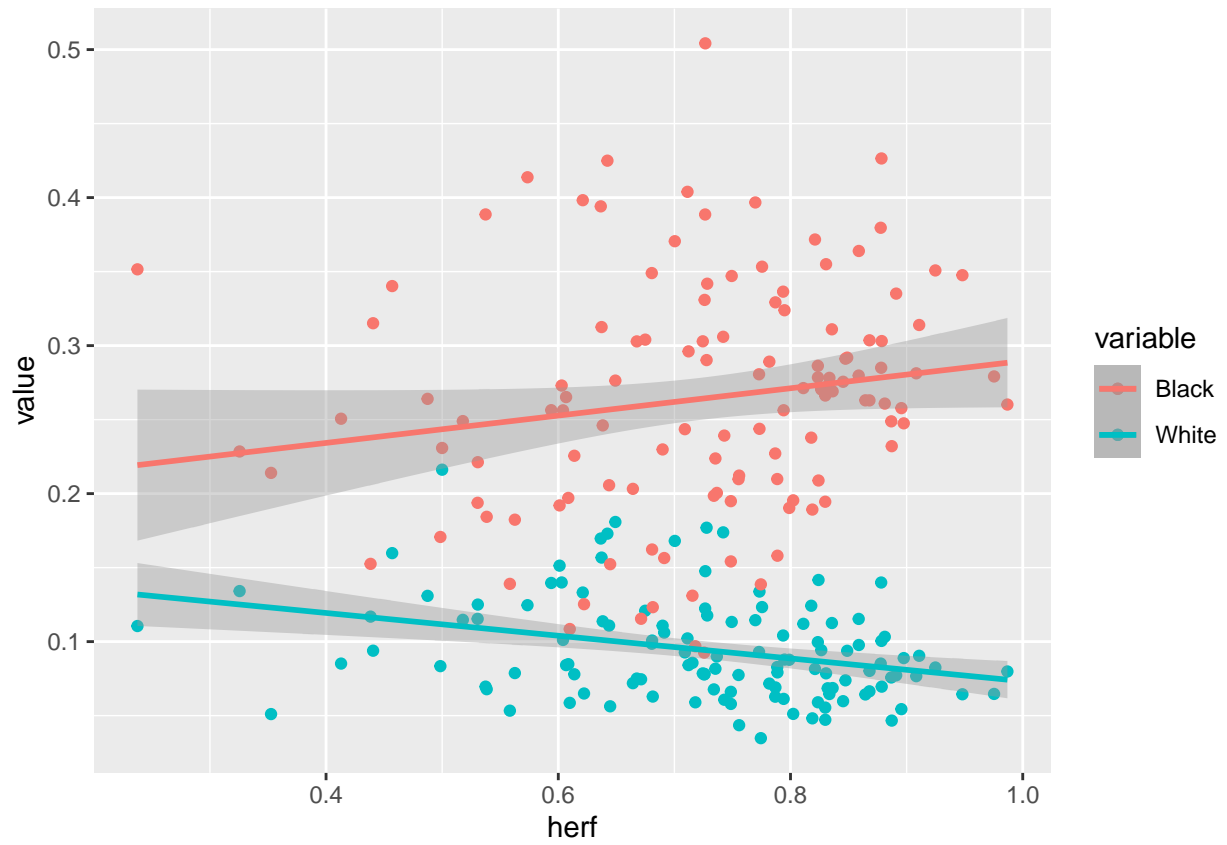
Table 8:

	<i>Dependent variable:</i>	
	povrate_w	povrate_b
	(1)	(2)
herf	−0.070*** (0.021)	0.092* (0.048)
lenper	−3.022*** (1.011)	0.004 (4.398)
Constant	0.148*** (0.017)	0.197*** (0.036)
Observations	121	121
R <sup>2</sup>	0.111	0.027
Adjusted R <sup>2</sup>	0.096	0.010
Residual Std. Error (df = 118)	0.033	0.079
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

```
plotted<-ggplot(mydata2, aes(herf, y = value, color = variable)) +
  geom_point(aes(y =povrate_w , col = "White")) +
  geom_point(aes(y = povrate_b, col = "Black"))+
  geom_smooth(method='lm', aes(y = povrate_w, col = "White"))+
  geom_smooth(method='lm', aes(y = povrate_b, col = "Black"))

plotted
```

```
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
```



## 7.12 Question: Generate a table with at least six estimations that checks whether the main results are robust to adding additional controls for city characteristics. What do you conclude?

Code:

```
regcont1<-felm(povrate_w~lenper+pctbk1990|0|(dism1990~herf+lenper+pctbk1990), mydata2)
regcont2<-felm(povrate_b~lenper+pctbk1990|0|(dism1990~herf+lenper+pctbk1990), mydata2)
regcont3<-felm(povrate_w~lenper+pctbk1990+manshr|0|(dism1990~herf+lenper+pctbk1990+manshr), mydata2)
regcont4<-felm(povrate_b~lenper+pctbk1990+manshr|0|(dism1990~herf+lenper+pctbk1990+manshr), mydata2)
regcont5<-felm(povrate_w~lenper+pctbk1990+manshr+pop1990|0|(dism1990~herf+lenper+pctbk1990+manshr+pop1990), mydata2)
regcont6<-felm(povrate_b~lenper+pctbk1990+manshr+pop1990|0|(dism1990~herf+lenper+pctbk1990+manshr+pop1990), mydata2)

stargazer(regcont1,regcont2,regcont3,regcont4,regcont5,regcont6,se = list(regcont1$rse,regcont2$rse,regcont3$rse,regcont4$rse,regcont5$rse,regcont6$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Wed, Oct 18, 2023 - 12:52:49 PM

Table 9:

	<i>Dependent variable:</i>				
	povrate_w	povrate_b	povrate_w	povrate_b	povrate_w
	(1)	(2)	(3)	(4)	(5)
lenper	-0.479 (1.801)	-2.331 (2.402)	1.214 (3.586)	-3.274 (4.590)	1.649 (3.661)
pctbk1990	0.211 (0.153)	-0.478* (0.246)	0.254 (0.215)	-0.480* (0.281)	0.233 (0.219)
manshr			0.135 (0.158)	-0.063 (0.239)	0.153 (0.166)
pop1990					0.000 (0.000)
‘dism1990(fit)’	-0.241** (0.097)	0.360** (0.141)	-0.338* (0.188)	0.344 (0.238)	-0.354* (0.195)
Constant	0.219*** (0.048)	0.091 (0.068)	0.248*** (0.068)	0.116 (0.082)	0.252*** (0.070)
Observations	121	121	111	111	111
R <sup>2</sup>	-0.254	0.108	-0.535	0.086	-0.592
Adjusted R <sup>2</sup>	-0.286	0.085	-0.593	0.052	-0.667
Residual Std. Error	0.039 (df = 117)	0.076 (df = 117)	0.043 (df = 106)	0.071 (df = 106)	0.044 (df = 105)

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

**Answer:** Controlling for additional city characteristics does not significantly alter the magnitude of the point estimates. If anything it makes the point estimates larger. This is evidence that confirms that the effect of RDI is operating through segregation and not via some other city characteristic.

## 8 Why Two Stage least squares?

Because the estimates in this paper only feature one endogenous regressor and one instrument, it is an excellent example with which to illustrate build intuition and see what the instrumental variables regressor is actually doing because in this scenario the IV estimator is exactly equal to the two stage least squares estimator ( $\hat{\beta}_{IV} = \hat{\beta}_{2SLS}$ ).

**8.1 Question: Estimate the first stage regression and use your estimates to generate the predicted values for the explanatory variable for all the observations.**

**Code:**

```
regfs<-lm(dism1990~herf+lenper, mydata2)
mydata2$pred_dism1990<-predict(regfs)
```

**8.2 Question:** If our instrument is valid, the step above “removed” the “bad” endogenous variation from the predicted explanatory variable, keeping only the exogenous variation that is generated by the instrument. Now run the second stage by regressing our outcome variable on the predicted values generated above and the relevant controls. Compare your estimates from this regression to those generated earlier. How do they compare?

**Code:**

```
reg2s1<-lm(povrate_w~pred_dism1990+lenper, mydata2)
reg2s2<-lm(povrate_b~pred_dism1990+lenper, mydata2)

stargazer(reg2s1,reg2s2, se = list(reg2s1$rse,reg2s2$rse))
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Wed, Oct 18, 2023 - 12:52:49 PM

Table 10:

	<i>Dependent variable:</i>	
	povrate_w	povrate_b
	(1)	(2)
pred_dism1990	−0.196*** (0.061)	0.258* (0.148)
lenper	0.602 (2.961)	−4.780 (7.153)
Constant	0.205*** (0.033)	0.121 (0.081)
Observations	121	121
R <sup>2</sup>	0.111	0.027
Adjusted R <sup>2</sup>	0.096	0.010
Residual Std. Error (df = 118)	0.033	0.079
F Statistic (df = 2; 118)	7.336***	1.629
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01		

**Answer:** This approach using two separate regressions returns the exact same point estimates (though the standard errors are different).

## 9 Yet another IV trick: Taking the “Good” variation and scaling it

9.1 Question: Take the coefficient from you reduced form estimate and divide it by your first stage estimate. How does this value compare your earlier estimate for the main result?

**Answer:**

The coefficient from my reduced form estimate is -0.070 for white poverty and 0.092 for black poverty. The coefficient for the first stage is 0.357. We thus get -0.196 and 0.258, the same as our earlier point estimates.

## 10 Submission instructions:

- Make sure the final version of your assignment is knit in pdf format and uploaded to gradescope. Make sure you have one question response per page (unless otherwise indicated) so that question positions align with the template in gradescope. The final PDF should be 22 pages long.