

Instrumental Variables

Econ 2560, Fall 2023

Prof. Josh Abel

(Chapter 12)

Introduction

- We've seen how to estimate a causal effect of a binary treatment using a binary “instrumental variable”

- Ratio of reduced form effect to first stage effect

$$E[Y_i|Z_i] = \beta_0^{RF} + \beta_1^{RF} \cdot Z_i + \beta_2^{RF} \cdot X_{1i} + \dots$$

$$E[T_i|Z_i] = \beta_0^{FS} + \beta_1^{FS} \cdot Z_i + \beta_2^{FS} \cdot X_{1i} + \dots$$

$$\beta = \frac{\beta_1^{RF}}{\beta_1^{FS}}$$

- We will now generalize to think about causal effects of continuous variables
 - The intuition will not change, though!

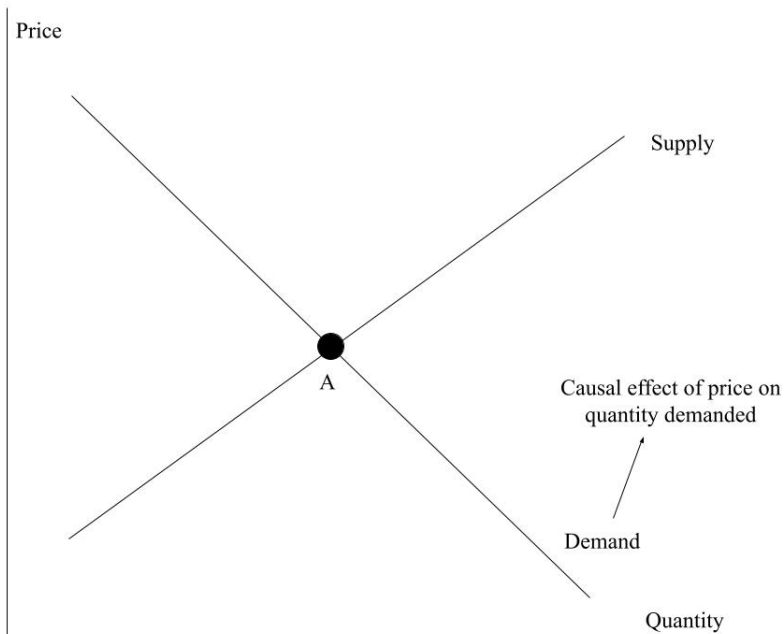
Supply-and-Demand Example

- Suppose we want to know how price affects Demand for fish

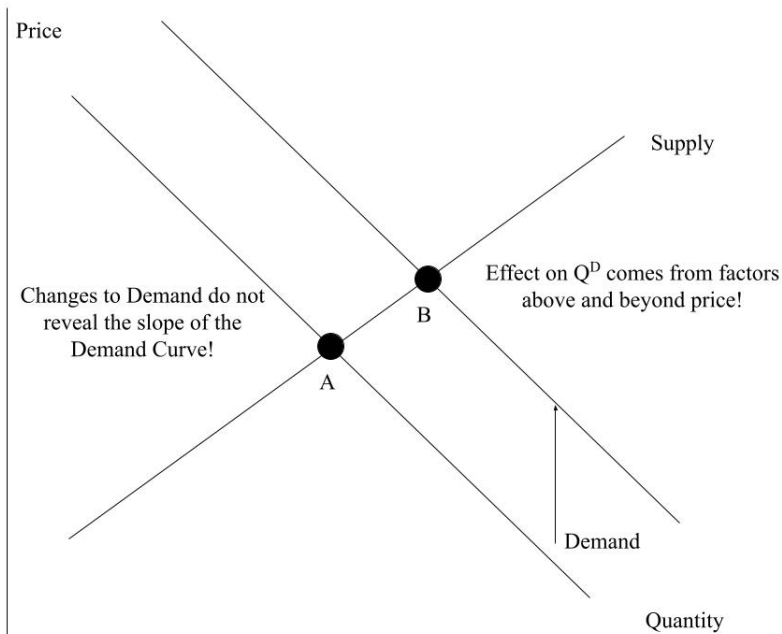
$$Q_i^D = \beta_0 + \beta_1 \cdot P_i + u_i$$

- Can we simply regress Quantity (Q_i^D) on price (P_i)?
- No! u_i likely contains important omitted variables
 - E.g. If there are days when consumers really want to eat fish, that will drive up both price and quantity
 - The regression above could find that Q^D tends to be high when P is high, which can't be the causal effect!
 - Classic Omitted Variable Bias...

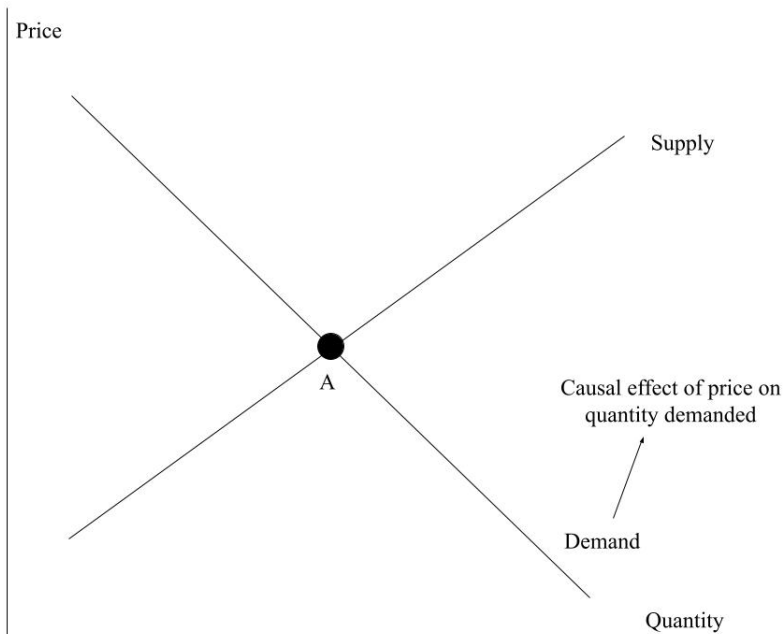
Supply-and-Demand Example, visualized



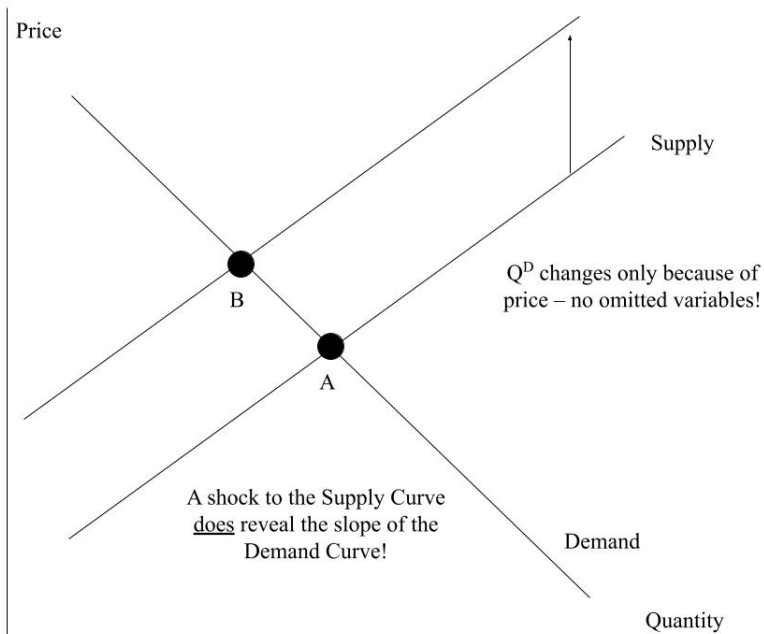
Supply-and-Demand Example, visualized



Supply-and-Demand Example, visualized



Supply-and-Demand Example, visualized



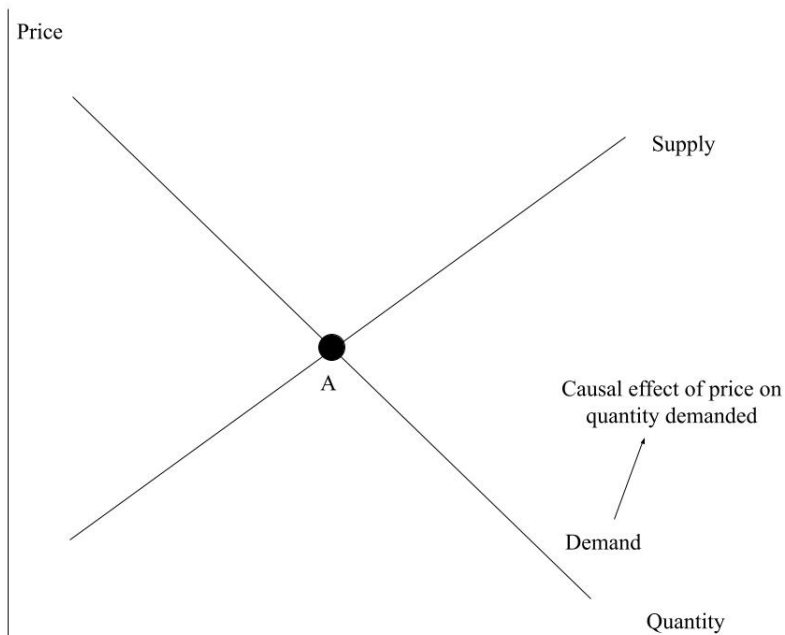
- Graddy argues that windspeed out at sea affects Supply but not Demand
- Consider two regression results

$$\ln(P)_i = -2.32 + 0.74 \cdot \text{WindSpeed}_i;$$

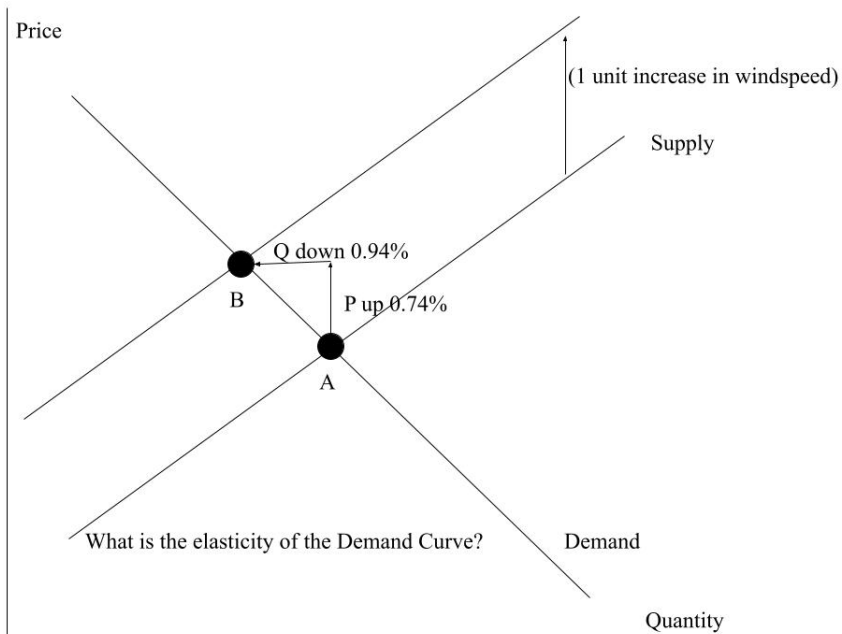
$$\ln(Q)_i = 11.21 - 0.94 \cdot \text{WindSpeed}_i;$$

- A 1-unit increase in windspeed causes a 0.74% increase in price and a 0.94% decrease in quantity...

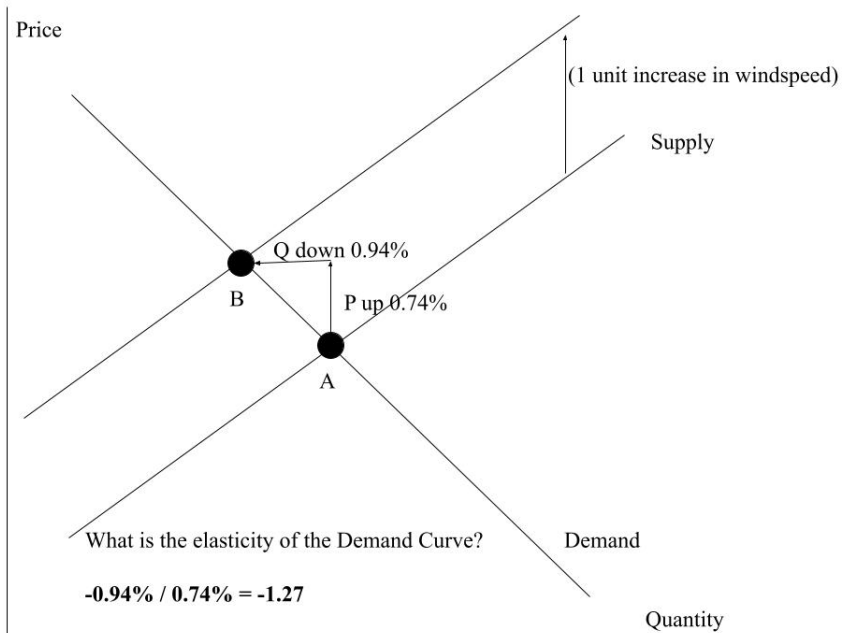
Graddy (2006), visualized



Graddy (2006), visualized



Graddy (2006), visualized



Why did that work?!

- Consider a causal model like the following:

$$\underbrace{Y_i}_{\ln(Q^D)_i} = \beta_0 + \beta_1 \cdot \underbrace{T_i}_{\ln(P)_i} + u_i,$$

where u_i contains omitted variables that affect Y_i and are correlated with T_i

- Suppose you try to estimate β_1 with the following OLS regression:

$$E[Y_i | T_i] = \hat{\beta}_0^{OLS} + \hat{\beta}_1^{OLS} \cdot T_i$$

- There is Omitted Variable Bias, so $\hat{\beta}_1^{OLS}$ will not deliver β_1 :(
- Windspeed is an “Instrumental Variable (IV)” and can solve the problem :)

An instrumental variable

$$Y_i = \beta_0 + \beta_1 \cdot T_i + u_i,$$

- An Instrumental Variable (or, “Instrument”), Z_i , must have two properties:

- 1 Relevance: it is predictive of the treatment (e.g. price of fish)

$$E[T_i|Z_i] = \beta_0^{FS} + \beta_1^{FS} \cdot Z_i, \text{ with } \beta_1^{FS} \neq 0$$

- 2 Exclusion/exogeneity: it does not impact the outcome (e.g. quantity demanded of fish) except through T_i

$$E[u_i|Z_i] = \lambda_0 + \lambda_1 \cdot Z_i, \text{ with } \lambda_1 = 0$$

The IV solution



$$Y_i = \beta_0 + \beta_1 \cdot T_i + u_i,$$

- Consider two regressions:

$$E[T_i|Z_i] = \beta_0^{FS} + \beta_1^{FS} \cdot Z_i$$

$$E[Y_i|Z_i] = \beta_0^{RF} + \beta_1^{RF} \cdot Z_i$$

- Another way of writing the second equation is:

$$E[Y_i|Z_i] = \beta_0 + \beta_1 \cdot E[T_i|Z_i] + E[u_i|Z_i]$$

The IV solution



$$Y_i = \beta_0 + \beta_1 \cdot T_i + u_i,$$

- Consider two regressions:

$$E[T_i|Z_i] = \beta_0^{FS} + \beta_1^{FS} \cdot Z_i$$

$$E[Y_i|Z_i] = \beta_0^{RF} + \beta_1^{RF} \cdot Z_i$$

- Another way of writing the second equation is:

$$E[Y_i|Z_i] = \beta_0 + \beta_1 \cdot E[T_i|Z_i] + E[u_i|Z_i]$$

$$= \beta_0 + \beta_1 \cdot (\beta_0^{FS} + \beta_1^{FS} \cdot Z_i) + \lambda_0$$

$$= \underbrace{(\beta_0 + \beta_1 \cdot \beta_0^{FS} + \lambda_0)}_{\beta_0^{RF}} + \underbrace{(\beta_1 \cdot \beta_1^{FS})}_{\beta_1^{RF}} \cdot Z_i$$

The IV solution



$$Y_i = \beta_0 + \beta_1 \cdot T_i + u_i,$$

- Consider two regressions:

$$E[T_i|Z_i] = \beta_0^{FS} + \beta_1^{FS} \cdot Z_i$$

$$E[Y_i|Z_i] = \beta_0^{RF} + \beta_1^{RF} \cdot Z_i$$

- Another way of writing the second equation is:

$$E[Y_i|Z_i] = \beta_0 + \beta_1 \cdot E[T_i|Z_i] + E[u_i|Z_i]$$

$$= \beta_0 + \beta_1 \cdot (\beta_0^{FS} + \beta_1^{FS} \cdot Z_i) + \lambda_0$$

$$= \underbrace{(\beta_0 + \beta_1 \cdot \beta_0^{FS} + \lambda_0)}_{\beta_0^{RF}} + \underbrace{(\beta_1 \cdot \beta_1^{FS})}_{\beta_1^{RF}} \cdot Z_i$$

- Therefore, $\beta_1 = \beta_1^{RF} / \beta_1^{FS}$.

Understanding the IV solution

- If Z generates “exogenous” variation in T (unrelated to u), then the causal effect of T on Y can be found by:
 - ① Finding the responsiveness of Y to Z
 - ② Dividing by the responsiveness of T to Z
- This is the same logic as “Reduced Form” / “First Stage” from an experiment with imperfect compliance!
 - Reduced form effect (Z on Y) is causal, but it needs to be rescaled because we care about the impact of T , not Z
 - Dividing by the first stage effect rescales the causal effect to properly show the effect of increasing T by a unit

$$\frac{dY}{dT} \Big|_{u \text{ constant}} = \frac{dY/dZ}{dT/dZ}$$

Returning to Graddy (2006)

First Stage (FS): $\ln(P)_i = -2.32 + 0.74 \cdot \text{WindSpeed}_i$

Reduced Form (RF): $\ln(Q)_i = 11.21 - 0.94 \cdot \text{WindSpeed}_i$

- Windspeed \uparrow by 1 $\rightarrow Q^D \downarrow$ by 0.94%
 - $dY/dZ = -0.94\%$
- Windspeed \uparrow by 1 $\rightarrow P \uparrow$ by 0.74%
 - $dT/dZ = 0.74\%$
- $\left. \frac{dY}{dT} \right|_{u \text{ constant}} = \frac{dY/dZ}{dT/dZ} = \frac{-0.94\%}{0.74\%} = -1.27.$

A different two-stage approach

- There is a different two stage process to estimate β_1 :

- 1 Estimate $E[T_i|Z_i] = \hat{\beta}_0^{FS} + \hat{\beta}_1^{FS} \cdot Z_i = \hat{T}_i$

- 2 Estimate $E[Y_i|\hat{T}_i] = \hat{\beta}_0^{SS} + \hat{\beta}_1^{SS} \cdot \hat{T}_i$

- Some math:

$$\begin{aligned} Y_i &= \hat{\beta}_0^{SS} + \hat{\beta}_1^{SS} \cdot \hat{T}_i + \hat{u}_i^{SS} \\ &= \hat{\beta}_0^{SS} + \hat{\beta}_1^{SS} \cdot (\hat{\beta}_0^{FS} + \hat{\beta}_1^{FS} \cdot Z_i) + \hat{u}_i^{SS} \\ &= \underbrace{(\hat{\beta}_0^{SS} + \hat{\beta}_1^{SS} \cdot \hat{\beta}_0^{FS})}_{\hat{\beta}_0^{RF}} + \underbrace{\hat{\beta}_1^{SS} \cdot \hat{\beta}_1^{FS}}_{\hat{\beta}_1^{RF}} \cdot Z_i + \hat{u}_i^{SS} \end{aligned}$$

- Produces the same estimate: as $\beta_1^{SS} = \beta_1^{RF} / \beta_1^{FS} = \beta_1$

Regression table

Dep. var.	OLS ln(Q)	First Stage ln(P)	Reduced Form ln(Q)	Second Stage ln(Q)
ln(P)	-0.54 (0.17)***			
Wind Speed		0.74 (0.17)***	-0.94 (0.29)***	
ln(\hat{P})				-1.27 (0.40)***
Constant	8.42 (0.08)***	-2.32 (0.49)***	11.21 (0.83)***	8.28 (0.11)***

Why did that work?!

- 1 Estimate $E[T_i|Z_i] = \hat{\beta}_0^{FS} + \hat{\beta}_1^{FS} \cdot Z_i = \hat{T}_i$
 - Retain only the variation in the treatment (e.g. price of fish) that can be explained by the instrument (e.g. windspeed) – dump the rest
- 2 Estimate $E[Y_i|\hat{T}_i] = \hat{\beta}_0^{SS} + \hat{\beta}_1^{SS} \cdot \hat{T}_i$
 - Regress outcome (e.g. quantity demanded) on just that “exogenous” variation in the treatment
 - $\hat{\beta}_1^{SS}$ is the effect of “exogenous” movements in treatment
- Like a multivariate regression on steroids: we’ve retained only the “purest” variation in T
 - E.g. only price movements that are generated by a Supply shift – no Demand shifts!

Standard errors for IV

- We've talked about 2 ways to understand how IV generates a point estimate:
 - “Reduced form” effect of Z on Y , rescaled by “first stage” effect of Z on T ($\hat{\beta}_1^{RF} / \hat{\beta}_1^{FS}$)
 - Regression of Y on \hat{T}_i , fitted values from regression of T on Z ($\hat{\beta}_1^{SS}$)
- These methods don't produce the correct SEs
 - Not clear how to compute a SE for $\hat{\beta}_1^{RF} / \hat{\beta}_1^{FS}$
 - $\hat{\beta}_1^{SS}$ ignores that \hat{T}_i is not true data
- It turns out that there is a way to estimate the coefficients in a single step that gets the SEs right...

IV Regression/Two-Stage Least Squares (2SLS)

- OLS works using orthogonality conditions:

$$E[\hat{u}_i^{OLS}] = E[Y_i - \hat{\beta}_0^{OLS} - \hat{\beta}_1^{OLS} \cdot T_i] = 0$$

$$E[T_i \cdot \hat{u}_i^{OLS}] = E[T_i \cdot (Y_i - \hat{\beta}_0^{OLS} - \hat{\beta}_1^{OLS} \cdot T_i)] = 0$$

IV Regression/Two-Stage Least Squares (2SLS)

- OLS works using orthogonality conditions:

$$E[\hat{u}_i^{OLS}] = E[Y_i - \hat{\beta}_0^{OLS} - \hat{\beta}_1^{OLS} \cdot T_i] = 0$$

$$E[T_i \cdot \hat{u}_i^{OLS}] = E[T_i \cdot (Y_i - \hat{\beta}_0^{OLS} - \hat{\beta}_1^{OLS} \cdot T_i)] = 0$$

Problem: True u_i s are not uncorrelated with T_i s, so $\hat{u}_i^{OLS} \neq u_i$, meaning $\hat{\beta}^{OLS} \neq \beta$. This orthogonality condition needs to be replaced – it is forcing us into a bad answer!

IV Regression/Two-Stage Least Squares (2SLS)

- IV Regression/2SLS uses a different equation

$$E[\hat{u}_i^{IV}] = E[Y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} \cdot T_i] = 0$$

$$E[Z_i \cdot \hat{u}_i^{IV}] = E[Z_i \cdot (Y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} \cdot T_i)] = 0$$

IV Regression/Two-Stage Least Squares (2SLS)

- IV Regression/2SLS uses a different equation

$$E[\hat{u}_i^{IV}] = E[Y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} \cdot T_i] = 0$$

$$E[Z_i \cdot \hat{u}_i^{IV}] = E[Z_i \cdot (Y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} \cdot T_i)] = 0$$

- If T_i and u_i are correlated, so not want $E[T_i \cdot \hat{u}_i] = 0$.
- The instrument replaces it with a valid second equation to allow us to solve for β_0 and β_1
 - Recall that the exclusion restriction was that Z and u are uncorrelated
- An “IV Regression” or “Two-Stage Least Squares” (2SLS) solves for $\hat{\beta}_0$ and $\hat{\beta}_1$ that satisfy above equations

Regression Table

Dep. var.	OLS ln(Q)	First Stage ln(P)	Reduced Form ln(Q)	Second Stage ln(Q)	IV ln(Q)
ln(P)	-0.54 (0.17)***				
Wind Speed		0.74 (0.17)***	-0.94 (0.29)***		
$\ln(\hat{P})$				-1.27 (0.40)***	-1.27 (0.49)**
Constant	8.42 (0.08)***	-2.32 (0.49)***	11.21 (0.83)***	8.28 (0.11)***	8.28 (0.13)***

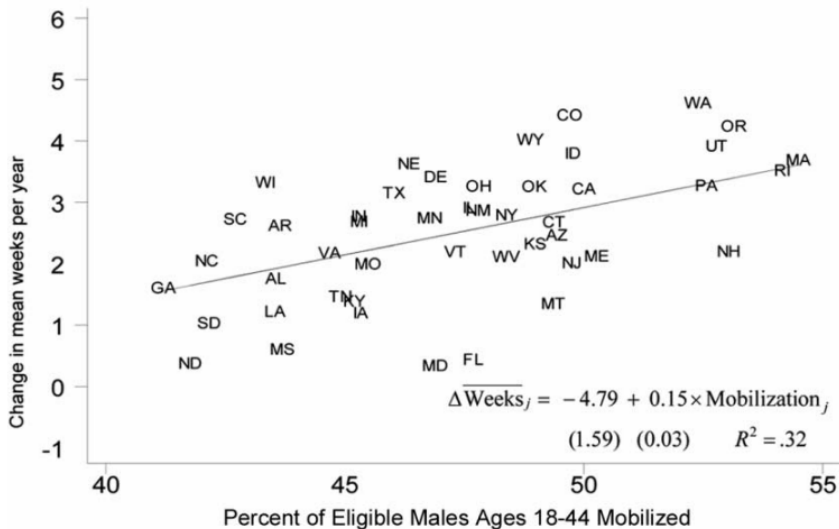


FIG. 4.—State WWII mobilization rates and change in mean female weeks worked per year, 1940–50.

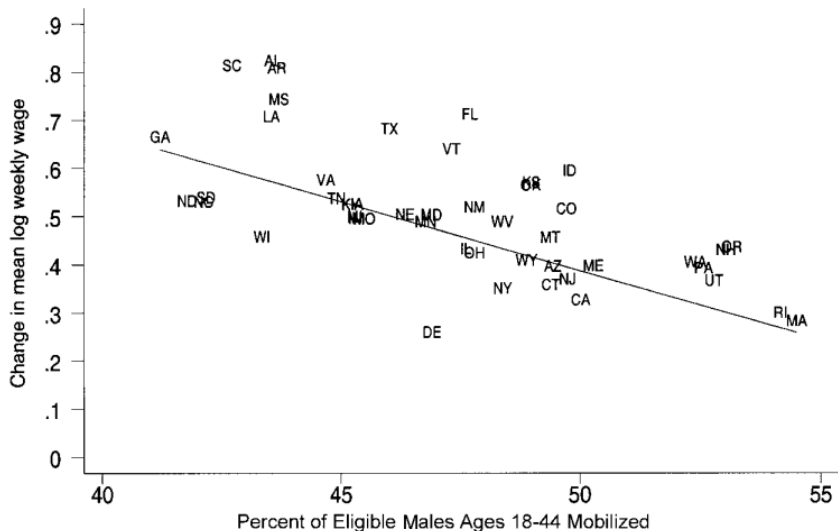


FIG. 6.—State WWII mobilization rates and change in mean log weekly real wages (1959 dollars) of full-time female workers, 1940–50.

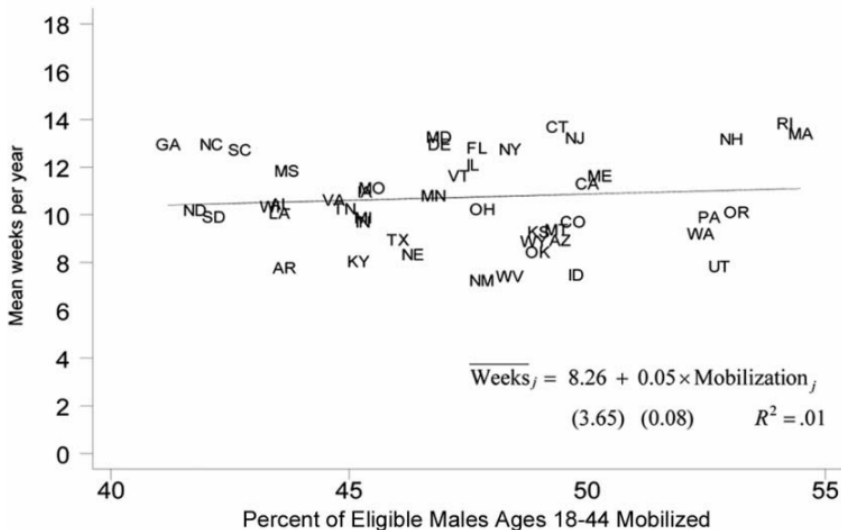
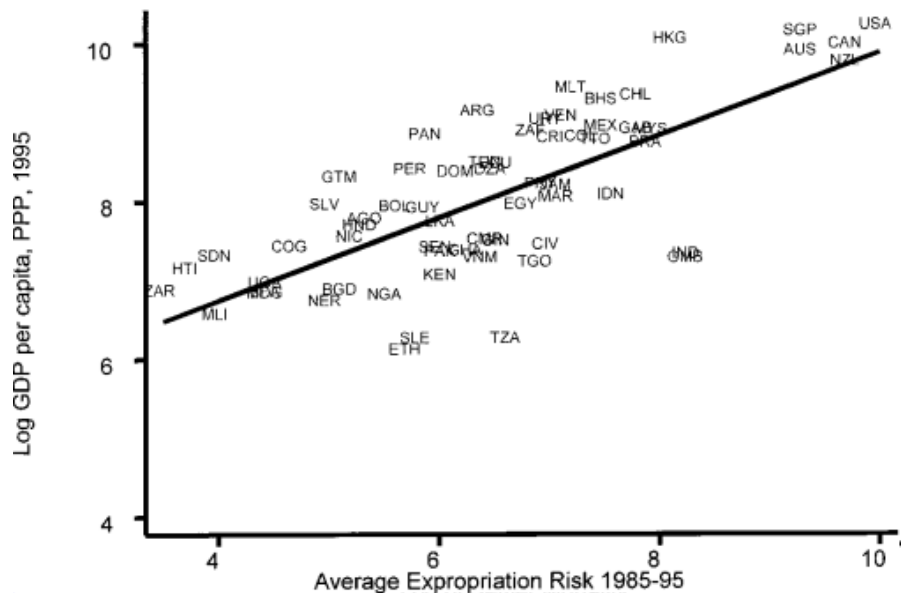
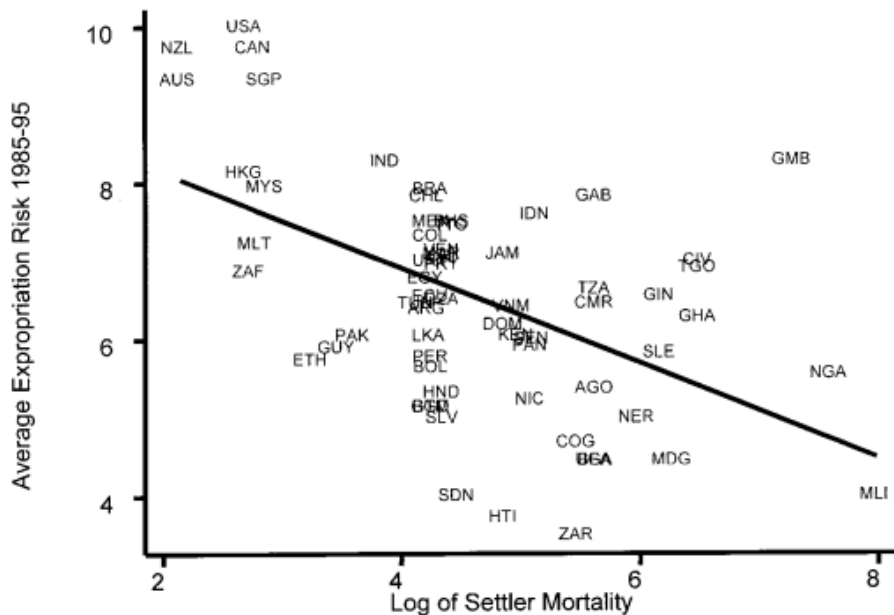


FIG. 3.—State WWII mobilization rates and mean female weeks worked per year, 1940

Acemoglu et al (2001)

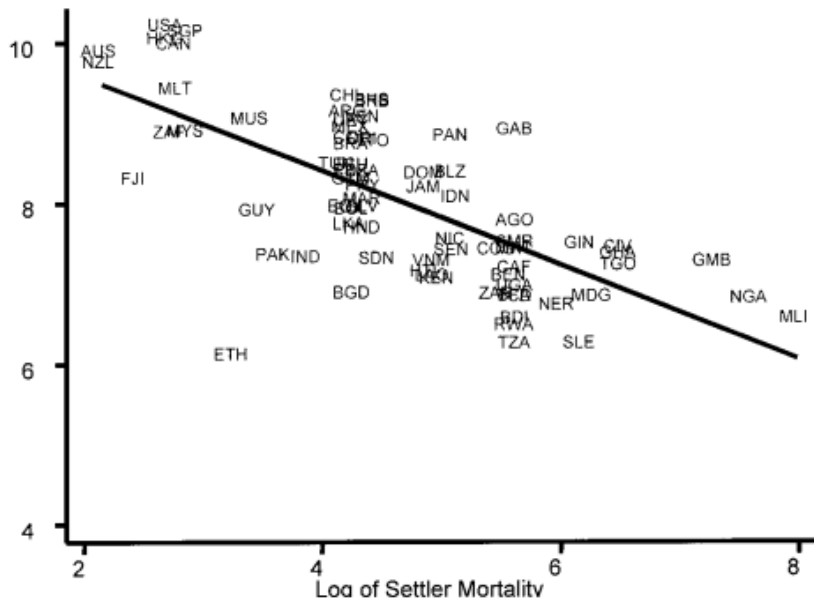


Acemoglu et al (2001)



Acemoglu et al (2001)

Log GDP per capita, PPP, 1995



Acemoglu et al (2001)

TABLE 4—IV REGRESSIONS OF LOG GDP PER CAPITA

	Base sample (1)	Base sample (2)	Base sample without Neo-Europes (3)	Base sample without Neo-Europes (4)	Base sample without Africa (5)	Base sample without Africa (6)	Base sample with continent dummies (7)	Base sample with continent dummies (8)	Base sample, dependent variable is log output per worker (9)
Panel A: Two-Stage Least Squares									
Average protection against expropriation risk 1985–1995	0.94 (0.16)	1.00 (0.22)	1.28 (0.36)	1.21 (0.35)	0.58 (0.10)	0.58 (0.12)	0.98 (0.30)	1.10 (0.46)	0.98 (0.17)
Latitude		-0.65 (1.34)		0.94 (1.46)		0.04 (0.84)		-1.20 (1.8)	
Asia dummy							-0.92 (0.40)	-1.10 (0.52)	
Africa dummy							-0.46 (0.36)	-0.44 (0.42)	
“Other” continent dummy							-0.94 (0.85)	-0.99 (1.0)	
Panel B: First Stage for Average Protection Against Expropriation Risk in 1985–1995									
Log European settler mortality	-0.61 (0.13)	-0.51 (0.14)	-0.39 (0.13)	-0.39 (0.14)	-1.20 (0.22)	-1.10 (0.24)	-0.43 (0.17)	-0.34 (0.18)	-0.63 (0.13)
Latitude		2.00 (1.34)		-0.11 (1.50)		0.99 (1.43)		2.00 (1.40)	
Asia dummy							0.33 (0.49)	0.47 (0.50)	
Africa dummy							-0.27 (0.41)	-0.26 (0.41)	
“Other” continent dummy							1.24 (0.84)	1.1 (0.84)	
R ²	0.27	0.30	0.13	0.13	0.47	0.47	0.30	0.33	0.28
Panel C: Ordinary Least Squares									
Average protection against expropriation risk 1985–1995	0.52 (0.06)	0.47 (0.06)	0.49 (0.08)	0.47 (0.07)	0.48 (0.07)	0.47 (0.07)	0.42 (0.06)	0.40 (0.06)	0.46 (0.06)
Number of observations	64	64	60	60	37	37	64	64	61