

CENTRO DE ESTUDIOS ECONÓMICOS

Maestría en Economía 2024–2026

Microeconometrics for Evaluation

4 Regression with Controls I

Disclaimer: I AM NOT the original intellectual author of the material presented in these notes. The content is STRONGLY based on a combination of lecture notes (from Aurora Ramirez), textbook references, and personal annotations for learning purposes. Any errors or omissions are entirely my own responsibility.

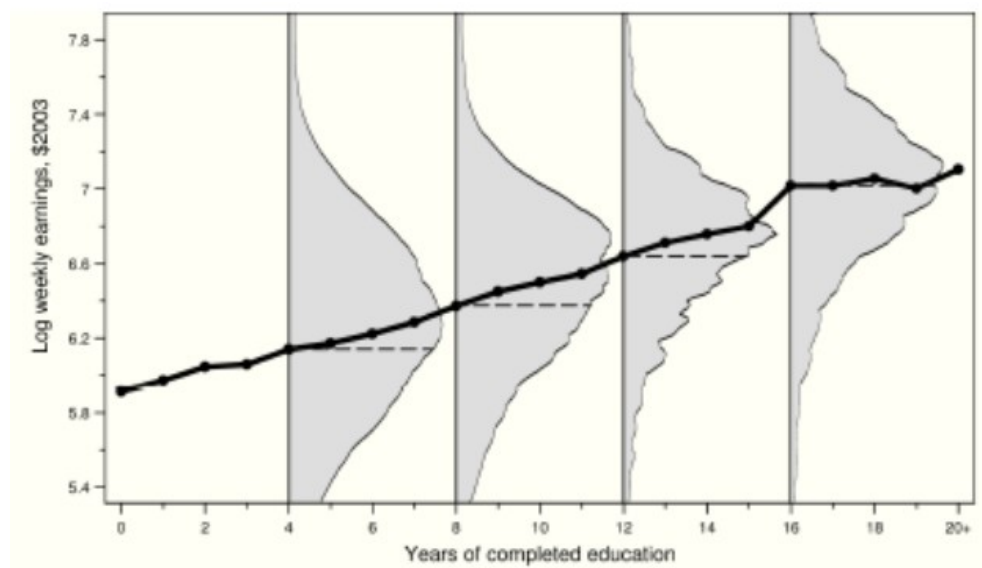
Índice

The Conditional Expectation Function (CEF)	2
CEF of Wages and Education	2
Population Regression	3
Why Regression? Three Reasons	4
The Regression-CEF Theorem (MHE 3.1.6)	4
Proof	4
The CEF is All You Need	6
Saturated Models and Main Effects	8
Saturated Models II: Two Binary Regressors	9
Regression Anatomy: Coefficients as Partial Slopes	9
Omitted Variables Bias (OVB)	11
Regression and Causality	13
Regression and Causality II: Selection Bias	14
Using the CIA	15
Regression and the CIA	16
Regression and the CIA II	16
The College Matching Matrix	18
Private School Effects: Barron's Matches	19
Private School Effects: Average SAT Controls	20
School Selectivity Effects: Average SAT Controls	22
Private School Effects: Omitted Variable Bias	23
Conclusions	24

The Conditional Expectation Function (CEF)

- The **Conditional Expectation Function** for a dependent outcome Y_i , given a $K \times 1$ vector of covariates X_i (with entries x_{ik}), is written $E[Y_i|X_i]$ and is a function of X_i .
- Since X_i is random, the CEF is also random. For a dummy D_i , the CEF has two values: $E[Y_i|D_i = 1]$ and $E[Y_i|D_i = 0]$.
- For a specific X_i , say $X_i = 42$, we denote it as $E[Y_i|X_i = 42]$.
- For continuous Y_i with density $f_Y(\cdot|X_i = x)$, the CEF is $E[Y_i|X_i = x] = \int t f_Y(t|X_i = x) dt$.
- If Y_i is discrete, $E[Y_i|X_i = x]$ equals $\sum_t t f_Y(t|X_i = x)$.
- The CEF residual is uncorrelated with any function of X_i . Let $\epsilon_i = Y_i - E[Y_i|X_i]$. For any function $h(x)$, we have $E[\epsilon_i h(X_i)] = E[(Y_i - E[Y_i|X_i])h(X_i)] = 0$.
- **Intuition:** The CEF is the “best guess” of Y given X . It averages possible Y values conditional on X , whether continuous or discrete, and ensures that prediction errors are unrelated to X itself.

CEF of Wages and Education



Raw data and the **CEF** of mean log weekly earnings by schooling. The sample consists of white males aged 40–49 in the 1990 IPUMS 5 % file.

- The plot shows the link between **years of education** and average *log weekly wages*.
- Each dot traces the **conditional expectation** of earnings for a given education level.

- Shaded areas represent the *distribution* of wages within each schooling group.
- The figure highlights that wages rise steadily with schooling, with **plateaus** around high school and college completion.
- **Intuition:** More education generally leads to higher wages, but gains are not uniform—big jumps often occur when completing key degrees (e.g., finishing high school or college).

Population Regression

Define population regression (from here onwards regression) as the solution to the population least squares problem. Specifically, the $K \times 1$ regression coefficient vector β is defined by solving

$$\beta = \arg \min_b \mathbb{E}[(Y_i - X_i' b)^2] \quad \text{Definición del problema de mínimos cuadrados} \quad (1)$$

$$L(b) = \mathbb{E}[(Y_i - X_i' b)^2] \quad \text{Función de pérdida} \quad (2)$$

$$= \mathbb{E}[Y_i^2 - 2Y_i X_i' b + (X_i' b)^2] \quad \text{Expansión del cuadrado} \quad (3)$$

$$\frac{\partial L(b)}{\partial b} = \frac{\partial}{\partial b} \left(\mathbb{E}[Y_i^2] - 2\mathbb{E}[Y_i X_i' b] + \mathbb{E}[(X_i' b)^2] \right) \quad \text{Derivamos con respecto a } b \quad (4)$$

$$= -2\mathbb{E}[X_i Y_i] + 2\mathbb{E}[X_i X_i'] b \quad \text{Aplicamos reglas de derivación matricial} \quad (5)$$

$$\text{FOC:} \quad -2\mathbb{E}[X_i Y_i] + 2\mathbb{E}[X_i X_i'] b = 0 \quad \text{Condición de primer orden} \quad (6)$$

$$\mathbb{E}[X_i X_i'] b = \mathbb{E}[X_i Y_i] \quad \text{Reordenamos} \quad (7)$$

$$b = \left(\mathbb{E}[X_i X_i'] \right)^{-1} \mathbb{E}[X_i Y_i] \quad \text{Solución explícita para } b \quad (8)$$

$$\beta = b \quad \text{Identificación del estimador poblacional} \quad (9)$$

$$\epsilon_i = Y_i - X_i' \beta \quad \text{Definimos residual poblacional} \quad (10)$$

$$\mathbb{E}[X_i \epsilon_i] = \mathbb{E}[X_i (Y_i - X_i' \beta)] = 0 \quad \text{Propiedad de ortogonalidad} \quad (11)$$

- **Intuition:** Population regression chooses β that minimizes the expected squared error between actual Y and its linear prediction $X\beta$.
- The solution comes from the first-order condition: errors must be *uncorrelated* with regressors.
- In simple words: regression picks coefficients so that on average, X carries no leftover information about residuals.

Why Regression? Three Reasons

- Regression solves the population least squares problem and is therefore the **Best Linear Predictor (BLP)** of Y_i given X_i .
- If the Conditional Expectation Function (CEF) is linear, then regression coincides with it.
- Regression provides the **best linear approximation** to the CEF in general.

The first statement is by definition. The second follows from the orthogonality of the CEF. The third point is formalized in the following theorem:

The Regression-CEF Theorem (MHE 3.1.6)

Regression-CEF Theorem (MHE 3.1.6) The population regression function $X_i'\beta$ delivers the *Minimum Mean Squared Error (MMSE)* linear approximation to $E[Y_i|X_i]$, that is:

$$\beta = \arg \min_b \mathbb{E}[(E[Y_i|X_i] - X_i'b)^2].$$

- **Intuition:** Regression works because it either **exactly equals the CEF** when the relation is linear, or otherwise it finds the **closest linear fit** to the true conditional expectation in terms of mean squared error. In plain words: regression is the best linear “shortcut” to describe how Y depends on X .

Proof

$$\beta = \arg \min_b \mathbb{E}[(Y_i - X_i'b)^2] \quad \text{Starting point} \quad (12)$$

$$Y_i - X_i'b = (Y_i - \mathbb{E}[Y_i|X_i]) + (\mathbb{E}[Y_i|X_i] - X_i'b) \quad \text{Add and subtract } E[Y_i|X_i] \quad (13)$$

$$(Y_i - X_i'b)^2 = \left((Y_i - \mathbb{E}[Y_i|X_i]) + (\mathbb{E}[Y_i|X_i] - X_i'b) \right)^2 \quad \text{Square both sides} \quad (14)$$

$$= (Y_i - \mathbb{E}[Y_i|X_i])^2 + 2(Y_i - \mathbb{E}[Y_i|X_i])(\mathbb{E}[Y_i|X_i] - X_i'b) \quad (15)$$

$$+ (\mathbb{E}[Y_i|X_i] - X_i'b)^2 \quad \text{Expansion of square} \quad (16)$$

$$\mathbb{E}[(Y_i - X_i'b)^2] = \mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])^2] + 2\mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])(\mathbb{E}[Y_i|X_i] - X_i'b)] \quad (17)$$

$$+ \mathbb{E}[(\mathbb{E}[Y_i|X_i] - X_i'b)^2] \quad \text{Take expectation} \quad (18)$$

$$\mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])^2] = \text{Var}(Y_i|X_i) \quad \text{Does not depend on } b \quad (19)$$

$$\mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])(\mathbb{E}[Y_i|X_i] - X_i'b)] = \mathbb{E}\left[\mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])(\mathbb{E}[Y_i|X_i] - X_i'b) \mid X_i]\right] \quad (20)$$

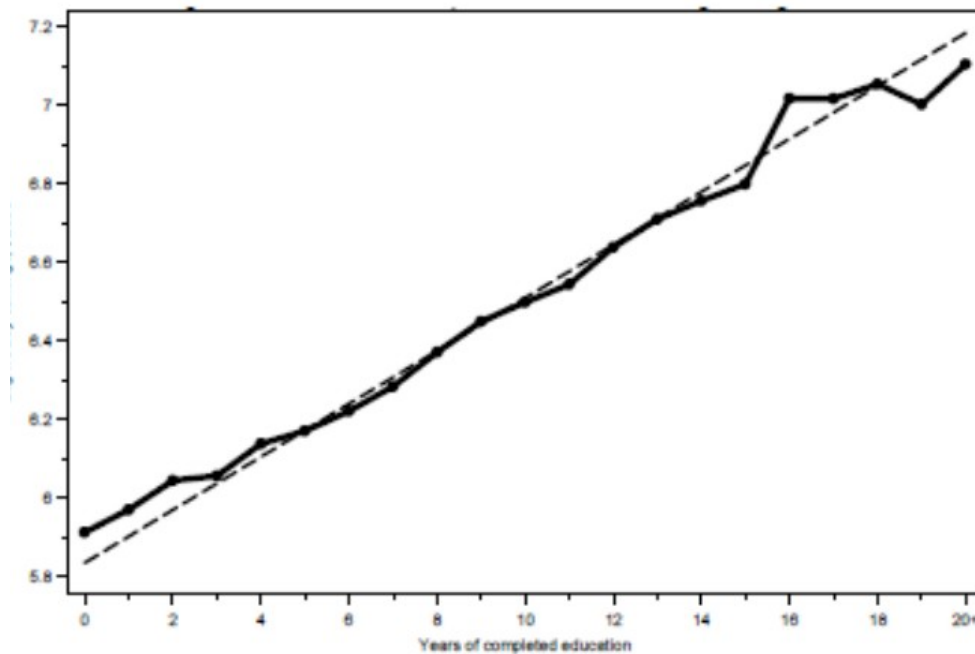
$$= \mathbb{E}\left[(\mathbb{E}[Y_i|X_i] - X_i'b) \mathbb{E}[Y_i - \mathbb{E}[Y_i|X_i] \mid X_i]\right] \quad \text{Law of iterated expectations} \quad (21)$$

$$= \mathbb{E}\left[(\mathbb{E}[Y_i|X_i] - X_i'b) \cdot 0\right] = 0 \quad \text{Property: residual has mean zero} \quad (22)$$

$$\mathbb{E}[(Y_i - X_i'b)^2] = \underbrace{\mathbb{E}[(Y_i - \mathbb{E}[Y_i|X_i])^2]}_{\text{irrelevant for } b} + \mathbb{E}[(\mathbb{E}[Y_i|X_i] - X_i'b)^2] \quad (23)$$

$$\beta = \arg \min_b \mathbb{E}[(\mathbb{E}[Y_i|X_i] - X_i'b)^2] \quad \text{Final result} \quad (24)$$

- **Intuition:** By adding and subtracting $E[Y|X]$, the error splits into two parts: (i) the variance around the conditional mean (independent of b), and (ii) the squared distance between the conditional mean and the linear predictor Xb . Minimizing the full expression is thus the same as making Xb the **closest linear approximation** to $E[Y|X]$ in mean squared error.



Regression line (dashed) compared with the **CEF** of mean log weekly earnings by years of schooling. Sample restricted to white men aged 40–49, from the 1980 Census IPUMS 5% file.

- The solid curve shows the **empirical CEF** of wages by education level.
- The dashed line is the **linear regression fit** to that CEF.
- The regression line smooths out local fluctuations and provides a simple linear summary of the relationship.
- Although the true CEF is not perfectly linear, the regression still captures the main upward trend in wages with education.
- **Intuition:** The regression acts like a “straight thread” through the curved conditional expectation. It cannot capture every bend, but it summarizes the overall slope: more schooling is strongly associated with higher wages on average.

The CEF is All You Need

- The Regression–CEF Theorem shows that we can use $\mathbb{E}[Y_i|X_i]$ directly as the dependent variable instead of Y_i (though the weighting is important).
- Another way to express this result is:

$$\beta = \left(\mathbb{E}[X_i X_i'] \right)^{-1} \mathbb{E}[X_i Y_i] \quad \text{Population regression formula} \quad (25)$$

$$= \left(\mathbb{E}[X_i X_i'] \right)^{-1} \mathbb{E}[X_i \mathbb{E}[Y_i | X_i]] \quad \text{Replace } Y_i \text{ with its CEF} \quad (26)$$

- This equivalence implies that grouped-data regression (using averages of Y_i within X_i categories) is valid when micro-data cannot be analyzed.
- For example, the schooling coefficient in a wage equation can be estimated using conditional means of log earnings across 21 education levels.
- In practice, if we weight groups by the number of individuals at each level, the regression based on grouped data yields the same coefficients as regression on micro-data.
- **Intuition:** The CEF contains all the relevant predictive information about Y given X . Whether you regress Y directly or its conditional expectation, you recover the same β — as long as weights are respected. In plain words: you can “compress” the data into averages without losing information about the slope.

A - Individual-level data						
<code>. regress earnings school, robust</code>						
Source	SS	df	MS	Number of obs = 409435		
Model	22631.4793	1	22631.4793	F(1,409433) =49118.25		
Residual	188648.31	409433	.460755019	Prob > F = 0.0000		
Total	211279.789	409434	.51602893	R-squared = 0.1071		
				Adj R-squared = 0.1071		
				Root MSE = .67879		
earnings	Coef.	Robust Std. Err.	t	Old Fashioned Std. Err.	t	
school	.0674387	.0003447	195.63	.0003043	221.63	
const.	5.835761	.0045507	1282.39	.0040043	1457.38	
B - Means by years of schooling						
<code>. regress average_earnings school [aweight=count], robust</code>						
(sum of wgt is 4.0944e+05)						
Source	SS	df	MS	Number of obs = 21		
Model	1.16077332	1	1.16077332	F(1, 19) = 540.31		
Residual	.040818796	19	.002148358	Prob > F = 0.0000		
Total	1.20159212	20	.060079606	R-squared = 0.9660		
				Adj R-squared = 0.9642		
				Root MSE = .04635		
average_earnings	Coef.	Robust Std. Err.	t	Old Fashioned Std. Err.	t	
school	.0674387	.0040352	16.71	.0029013	23.24	
const.	5.835761	.0399452	146.09	.0381792	152.85	

Comparison of regression results:

- (A) Individual-level data with robust errors,
- (B) Grouped means by schooling with weights.

- Panel A shows the regression of **individual earnings** on years of schooling using the full micro-sample ($N \approx 409,000$).
- The estimated coefficient on schooling is about 0,067, with extremely high statistical significance.
- Panel B repeats the regression using **grouped averages** of earnings by schooling level ($N = 21$ groups), weighting by group size.
- The estimated schooling coefficient is essentially the same (0,067), despite working with far fewer observations.
- This confirms the earlier theorem: regression using grouped CEF data produces coefficients identical to those from the underlying micro-data, provided correct weighting is applied.
- **Intuition:** Whether we use millions of individuals or just a handful of averages, the slope is the same because the CEF already condenses all the information we need. Grouped data is like a “compressed file” — smaller in size but containing the same essential signal about the schooling–wage relationship.

Saturated Models and Main Effects

- A **saturated regression model** is one where the explanatory variable is discrete, and the regression includes a separate parameter for each possible value of that variable.
- In other words, the model perfectly matches group means by allowing one coefficient per category.

Example. Suppose schooling $s_i \in \{0, 1, 2, \dots, \tau\}$. A saturated regression model is:

$$Y_i = \alpha + \beta_1 d_{1i} + \beta_2 d_{2i} + \dots + \beta_\tau d_{\tau i} + \epsilon_i \quad (27)$$

where $d_{ji} = 1[s_i = j]$ is a dummy for being in schooling level j . Here, β_j represents the effect of schooling level j , and α is the reference category mean.

$$\beta_j = \mathbb{E}[Y_i | s_i = j] - \mathbb{E}[Y_i | s_i = 0] \quad \text{Difference in group means} \quad (28)$$

$$\alpha = \mathbb{E}[Y_i | s_i = 0] \quad \text{Reference group mean} \quad (29)$$

- **Intuition:** A saturated regression is basically the same as comparing group averages. Each coefficient β_j tells us how much higher or lower group j is compared to the baseline group. Nothing is approximated — the model exactly fits the sample means for each category.

Saturated Models II: Two Binary Regressors

- Let x_{1i} indicate whether an individual is a **college graduate**.
- Let x_{2i} indicate whether an individual is **female**.
- The CEF conditional on (x_{1i}, x_{2i}) has four possible values:

$$\mathbb{E}[Y_i \mid x_{1i} = 0, x_{2i} = 0] = \alpha \quad \begin{array}{l} \text{baseline:} \\ \text{no college, male} \end{array} \quad (30)$$

$$\mathbb{E}[Y_i \mid x_{1i} = 1, x_{2i} = 0] = \alpha + \beta_1 \quad \begin{array}{l} \text{effect of college,} \\ \text{men only} \end{array} \quad (31)$$

$$\mathbb{E}[Y_i \mid x_{1i} = 0, x_{2i} = 1] = \alpha + \gamma \quad \begin{array}{l} \text{effect of being female,} \\ \text{no college} \end{array} \quad (32)$$

$$\mathbb{E}[Y_i \mid x_{1i} = 1, x_{2i} = 1] = \alpha + \beta_1 + \gamma + \delta_1 \quad \begin{array}{l} \text{college + female} \\ \text{plus interaction} \end{array} \quad (33)$$

These four cases can be summarized with a saturated regression:

$$\mathbb{E}[Y_i \mid x_{1i}, x_{2i}] = \alpha + \beta_1 x_{1i} + \gamma x_{2i} + \delta_1 (x_{1i} \cdot x_{2i}) \quad (34)$$

- The model contains two **main effects** (β_1 for college, γ for female) and one **interaction** (δ_1 for “female college graduate”).
- It is generally odd to include the interaction term without also including the main effects, since the meaning of the interaction depends on them.
- **Intuition:** This regression just assigns an average outcome to each of the four groups (male/non-college, male/college, female/non-college, female/college). The interaction term δ_1 captures whether the combined effect of being both female and college-educated is more (or less) than the sum of the individual effects.

Regression Anatomy: Coefficients as Partial Slopes

Step 1: Bivariate regression recap

$$\beta_1 = \frac{\text{Cov}(Y_i, X_i)}{\text{Var}(X_i)} \quad \begin{array}{l} \text{Slope coefficient} \\ \text{in simple regression} \end{array} \quad (35)$$

$$\alpha = \mathbb{E}[Y_i] - \beta_1 \mathbb{E}[X_i] \quad \begin{array}{l} \text{Intercept from} \\ \text{mean restriction} \end{array} \quad (36)$$

Step 2: Multivariate regression model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_K X_{Ki} + \epsilon_i \quad (37)$$

We want the general formula for the coefficient β_k .

Step 3: Residualizing X_{ki}

$$\tilde{X}_{ki} = X_{ki} - \hat{X}_{ki} \quad \begin{array}{l} \text{Residual from regression} \\ \text{of } X_{ki} \text{ on all other } X_{-k} \end{array} \quad (38)$$

By construction, \tilde{X}_{ki} is orthogonal to all other regressors X_{-k} .

Step 4: Regression anatomy claim

$$\beta_k = \frac{\text{Cov}(Y_i, \tilde{X}_{ki})}{\text{Var}(\tilde{X}_{ki})} \quad \begin{array}{l} \text{Formula for} \\ \text{partial slope} \end{array} \quad (39)$$

Step 5: Verification algebra

$$Y_i = \beta_0 + \beta_1 X_{1i} + \cdots + \beta_K X_{Ki} + \epsilon_i \quad \begin{array}{l} \text{Population regression} \end{array} \quad (40)$$

$$\text{Cov}(Y_i, \tilde{X}_{ki}) = \text{Cov}\left(\beta_0 + \sum_{j=1}^K \beta_j X_{ji} + \epsilon_i, \tilde{X}_{ki}\right) \quad \begin{array}{l} \text{Take covariance} \\ \text{with } \tilde{X}_{ki} \end{array} \quad (41)$$

$$= \beta_0 \cdot \text{Cov}(1, \tilde{X}_{ki}) + \sum_{j=1}^K \beta_j \text{Cov}(X_{ji}, \tilde{X}_{ki}) + \text{Cov}(\epsilon_i, \tilde{X}_{ki}) \quad (42)$$

$$= \beta_k \text{Cov}(X_{ki}, \tilde{X}_{ki}) \quad \begin{array}{l} \text{Only } X_{ki} \text{ survives:} \\ \text{others orthogonal by construction} \\ \text{residual uncorrelated with regressors} \end{array} \quad (43)$$

$$\text{Cov}(X_{ki}, \tilde{X}_{ki}) = \text{Var}(\tilde{X}_{ki}) \quad \begin{array}{l} \text{Key property} \\ \text{of residualization} \end{array} \quad (44)$$

$$\text{Cov}(Y_i, \tilde{X}_{ki}) = \beta_k \text{Var}(\tilde{X}_{ki}) \quad (45)$$

$$\beta_k = \frac{\text{Cov}(Y_i, \tilde{X}_{ki})}{\text{Var}(\tilde{X}_{ki})} \quad \blacksquare \quad (46)$$

- **Intuition:** Each regression coefficient β_k is like a simple slope — but only after “partialling out” the influence of other regressors. By residualizing X_{ki} against all other covariates, we isolate its unique variation. Then the coefficient is just the simple regression of Y on this purified residual. In other words: multivariate regression is just bivariate regression, done on the leftover part of each variable.

Omitted Variables Bias (OVB)

The OVB formula shows how regression coefficients differ when some controls are omitted.

Step 1: Long regression (with ability A_i included)

$$Y_i = \alpha + \rho s_i + A_i' \gamma + \epsilon_i \quad \begin{array}{l} \text{Wages on schooling } s_i \\ \text{controlling for ability } A_i \end{array} \quad (47)$$

Step 2: Short regression (omit A_i)

$$Y_i = \alpha^* + \pi s_i + u_i \quad \text{Coefficient on schooling changes to } \pi \quad (48)$$

Step 3: Expression for short regression slope

$$\pi = \frac{\text{Cov}(Y_i, s_i)}{\text{Var}(s_i)} \quad \begin{array}{l} \text{Definition of slope} \\ \text{in simple regression} \end{array} \quad (49)$$

Step 4: Substitute Y_i from long regression

$$\text{Cov}(Y_i, s_i) = \text{Cov}(\alpha + \rho s_i + A_i' \gamma + \epsilon_i, s_i) \quad (50)$$

$$= \rho \text{Cov}(s_i, s_i) + \text{Cov}(A_i' \gamma, s_i) + \text{Cov}(\epsilon_i, s_i) \quad (51)$$

$$= \rho \text{Var}(s_i) + \gamma' \text{Cov}(A_i, s_i) \quad \begin{array}{l} \text{Assume } \epsilon_i \text{ uncorrelated with } s_i \end{array} \quad (52)$$

Step 5: Plug back into slope formula

$$\pi = \frac{\rho \text{Var}(s_i) + \gamma' \text{Cov}(A_i, s_i)}{\text{Var}(s_i)} \quad (53)$$

$$= \rho + \gamma' \frac{\text{Cov}(A_i, s_i)}{\text{Var}(s_i)} \quad (54)$$

Step 6: Define δ_{As} as regression of A_i on s_i

$$\delta_{As} = \frac{\text{Cov}(A_i, s_i)}{\text{Var}(s_i)} \quad \begin{array}{l} \text{Slope vector from regressing} \\ A_i \text{ on } s_i \end{array} \quad (55)$$

Final Result (OVB formula)

$$\pi = \rho + \gamma' \delta_{As} \quad \begin{array}{l} \text{Short regression} = \text{true effect} \\ + \text{bias from omitted variables} \end{array} \quad (56)$$

- **Intuition:** The coefficient on schooling changes when ability is omitted because part of the schooling effect captures the correlation between schooling and ability.
- If s_i and A_i are uncorrelated, then $\delta_{As} = 0$ and short regression equals the long regression.
- Otherwise, the bias is the product of: (i) the effect of the omitted variable on Y (γ), and (ii) the correlation between the omitted and included regressor (δ_{As}).
- In plain words: omitted variables bias = **effect of the omitted** \times **correlation with the included**.

Table 3.2.1: Estimates of the returns to education for men in the NLSY					
	(1)	(2)	(3)	(4)	(5)
Controls:	None	Age dummies	Col. (2) and additional controls*	Col. (3) and AFQT score	Col. (4), with occupation dummies
	0.132 (0.007)	0.131 (0.007)	0.114 (0.007)	0.087 (0.009)	0.066 (0.010)

Table 3.2.1: Estimates of the returns to education for men in the NLSY (1979 cohort, 2002 survey). The table shows coefficients on years of schooling from regressions of log wages on schooling and various sets of controls. Standard errors are in parentheses. Data are weighted using NLSY sampling weights. Sample size: 2,434. *Additional controls include mother's and father's schooling, plus dummies for race and Census region.

- Column (1): No controls, estimated return $\approx 13,2\%$.
- Column (2): Adding age dummies barely changes the estimate.
- Column (3): Adding parental background and demographics reduces the estimate to $11,4\%$.
- Column (4): Adding AFQT (ability) lowers it further to $8,7\%$.
- Column (5): Adding occupation dummies yields about $6,6\%$.

- This progression shows how including controls changes the coefficient, illustrating omitted variables bias in action.
- **Intuition:** The “raw” estimate overstates the return to schooling because schooling is correlated with ability, family background, and occupation. As controls are added, the coefficient drops, showing how much of the apparent return to schooling was actually due to omitted factors.

Regression and Causality

- Casual regressions are often exploratory or descriptive.
- Causal regressions are more ambitious: they describe counterfactual states of the world and are key for policy analysis.
- Example: Do private colleges increase earnings? Let C_i indicate private attendance ($C_i = 1$ if private, $C_i = 0$ if not).

Potential outcomes framework:

$$Y_{1i} = \text{Earnings if } C_i = 1 \text{ (private college)} \quad (57)$$

$$Y_{0i} = \text{Earnings if } C_i = 0 \text{ (non-private college)} \quad (58)$$

Observed outcome:

$$Y_i = \begin{cases} Y_{1i}, & \text{if } C_i = 1 \\ Y_{0i}, & \text{if } C_i = 0 \end{cases} \quad (59)$$

This can be rewritten algebraically:

$$Y_i = Y_{0i} + (Y_{1i} - Y_{0i})C_i \quad \text{Add and subtract } Y_{0i} \quad (60)$$

Define the **individual treatment effect**:

$$\tau_i = Y_{1i} - Y_{0i} \quad \begin{array}{l} \text{Causal effect for} \\ \text{individual } i \end{array} \quad (61)$$

Then the observed outcome becomes:

$$Y_i = Y_{0i} + \tau_i C_i \quad \begin{array}{l} \text{Observed outcome} = \\ \text{baseline} + \text{effect if treated} \end{array} \quad (62)$$

- We only observe one of $\{Y_{1i}, Y_{0i}\}$ for each individual.
- The parameter of interest is usually an **average treatment effect**.
- Example: the average effect of private college among those who attended:

$$\mathbb{E}[Y_{1i} - Y_{0i} \mid C_i = 1] \quad (\text{Treatment on the Treated, TOT}).$$

- **Intuition:** Causal regressions are about comparing “what happened” with “what would have happened otherwise.” We only observe one reality per person, so we rely on assumptions and design to recover average causal effects.

Regression and Causality II: Selection Bias

Step 1: Observed mean difference

$$\Delta = \mathbb{E}[Y_i | C_i = 1] - \mathbb{E}[Y_i | C_i = 0] \quad \begin{array}{l} \text{Naive comparison of} \\ \text{treated vs untreated} \end{array} \quad (63)$$

Step 2: Substitute potential outcomes

$$\mathbb{E}[Y_i | C_i = 1] = \mathbb{E}[Y_{1i} | C_i = 1] \quad \text{If treated, we see } Y_{1i} \quad (64)$$

$$\mathbb{E}[Y_i | C_i = 0] = \mathbb{E}[Y_{0i} | C_i = 0] \quad \text{If untreated, we see } Y_{0i} \quad (65)$$

$$\Delta = \mathbb{E}[Y_{1i} | C_i = 1] - \mathbb{E}[Y_{0i} | C_i = 0] \quad (66)$$

Step 3: Add and subtract $\mathbb{E}[Y_{0i} | C_i = 1]$

$$\Delta = \left(\mathbb{E}[Y_{1i} | C_i = 1] - \mathbb{E}[Y_{0i} | C_i = 1] \right) + \left(\mathbb{E}[Y_{0i} | C_i = 1] - \mathbb{E}[Y_{0i} | C_i = 0] \right) \quad (67)$$

Step 4: Identify components

$$\Delta = \underbrace{\mathbb{E}[Y_{1i} - Y_{0i} | C_i = 1]}_{\text{Treatment on the Treated (TOT)}} + \underbrace{\left(\mathbb{E}[Y_{0i} | C_i = 1] - \mathbb{E}[Y_{0i} | C_i = 0] \right)}_{\text{Selection Bias}} \quad (68)$$

Final decomposition:

$$\mathbb{E}[Y_i | C_i = 1] - \mathbb{E}[Y_i | C_i = 0] = \text{TOT} + \text{Selection Bias} \quad \blacksquare \quad (69)$$

- The naive comparison mixes causal effect with differences in baseline potential outcomes.
- If those who choose private college would have earned more anyway, the selection bias term is positive.
- In causal models, **selection bias is the analogue of OVB**.
- The Conditional Independence Assumption (CIA) states that after conditioning on X_i , treatment C_i is independent of potential outcomes:

$$(Y_{0i}, Y_{1i}) \perp\!\!\!\perp C_i | X_i.$$

Under CIA, the selection bias term disappears.

- **Intuition:** Comparing outcomes between treated and untreated is misleading because the groups differ in ways beyond treatment. By adjusting for observables (CIA), we “balance” groups so that the only remaining difference is due to treatment itself.

Using the CIA

Step 1: CIA assumption

$$(Y_{0i}, Y_{1i}) \perp\!\!\!\perp C_i \mid X_i \quad \begin{array}{l} \text{Treatment assignment} \\ \text{is as good as random} \\ \text{within } X_i\text{-cells} \end{array} \quad (70)$$

Step 2: Conditional mean comparison

$$\mathbb{E}[Y_i \mid X_i, C_i = 1] - \mathbb{E}[Y_i \mid X_i, C_i = 0] = \mathbb{E}[Y_{1i} \mid X_i, C_i = 1] - \mathbb{E}[Y_{0i} \mid X_i, C_i = 0] \quad (71)$$

$$= \mathbb{E}[Y_{1i} \mid X_i] - \mathbb{E}[Y_{0i} \mid X_i] \quad \begin{array}{l} \text{By CIA, potential outcomes} \\ \text{do not depend on } C_i \text{ given } X_i \end{array} \quad (72)$$

$$= \mathbb{E}[Y_{1i} - Y_{0i} \mid X_i] \quad (73)$$

Step 3: Interpretation

$$\mathbb{E}[Y_{1i} - Y_{0i} \mid X_i, C_i = 1] = \mathbb{E}[Y_{1i} - Y_{0i} \mid X_i] \quad \begin{array}{l} \text{Effect for treated group} \\ \text{equals effect for all with } X_i \end{array} \quad (74)$$

Step 4: Marginalize over X_i

$$\mathbb{E}[\mathbb{E}[Y_i \mid X_i, C_i = 1] - \mathbb{E}[Y_i \mid X_i, C_i = 0]] = \mathbb{E}[\mathbb{E}[Y_{1i} - Y_{0i} \mid X_i]] \quad (75)$$

$$= \mathbb{E}[Y_{1i} - Y_{0i}] \quad \text{Law of iterated expectations} \quad (76)$$

- Conditional-on- X_i , the comparison of treated vs untreated is unbiased.
- The CIA guarantees that within each X_i , treated and untreated differ only by treatment status.
- Averaging over X_i gives us the overall Average Treatment Effect (ATE).
- **Intuition:** The CIA tells us that people with the same X_i are comparable, regardless of whether they attend private college or not. This motivates methods like *matching*, where we compare treated and untreated individuals with identical or very similar characteristics.

Regression and the CIA

- The CIA allows regression to be interpreted causally.
- To focus on selection, assume a **constant treatment effect**.

Step 1: Define potential outcomes

$$Y_{0i} = \alpha + \eta_i \quad \begin{array}{l} \text{Outcome if untreated} \\ \text{with individual noise } \eta_i \end{array} \quad (77)$$

$$Y_{1i} = Y_{0i} + \rho \quad \text{Outcome if treated = baseline + constant effect } \rho \quad (78)$$

Step 2: Observed outcome rule

$$Y_i = Y_{0i} + (Y_{1i} - Y_{0i})C_i \quad \begin{array}{l} \text{From potential outcomes} \\ \text{see previous slides} \end{array} \quad (79)$$

$$= Y_{0i} + \rho C_i \quad \text{Substitute } Y_{1i} - Y_{0i} = \rho \quad (80)$$

$$= (\alpha + \eta_i) + \rho C_i \quad \text{Substitute } Y_{0i} = \alpha + \eta_i \quad (81)$$

$$= \alpha + \rho C_i + \eta_i \quad \text{Collect terms: observed equation} \quad (82)$$

Step 3: Interpretation

- Equation above looks like a bivariate regression of Y_i on C_i :

$$Y_i = \alpha + \rho C_i + \eta_i$$

- But it is not an actual regression equation — here C_i may be correlated with the residual η_i (selection into treatment).
- The CIA assumption is what allows us to treat this as a causal regression: under CIA, $C_i \perp \eta_i$, so ρ is consistently estimated.
- **Intuition:** With constant effects, the observed outcome looks just like a regression of Y on C . The slope is the causal effect ρ . But unless assignment is “as good as random” (CIA), C may be correlated with unobserved factors in η_i , biasing the estimate. The CIA guarantees that the regression coefficient recovers the true causal effect.

Regression and the CIA II

Step 1: CIA in constant-effects setup

$$\mathbb{E}[\eta_i \mid C_i, X_i] = \mathbb{E}[\eta_i \mid X_i] \quad \begin{array}{l} \text{CIA implies residual mean} \\ \text{independent of treatment } C_i \end{array} \quad (83)$$

Step 2: Assume linear form for conditional mean

$$\mathbb{E}[\eta_i | X_i] = X_i' \gamma$$

Model the expectation of η_i
as linear in covariates (84)

Step 3: Compute conditional expectation of Y_i

$$\mathbb{E}[Y_i | X_i, C_i] = \alpha + \rho C_i + \mathbb{E}[\eta_i | X_i] \quad \text{From earlier: } Y_i = \alpha + \rho C_i + \eta_i \quad (85)$$

$$= \alpha + \rho C_i + X_i' \gamma \quad (86)$$

Step 4: Define regression error

$$Y_i = \alpha + \rho C_i + X_i' \gamma + \nu_i \quad (87)$$

$$\nu_i \equiv \eta_i - X_i' \gamma = \eta_i - \mathbb{E}[\eta_i | C_i, X_i] \quad (88)$$

Step 5: Orthogonality condition

$$\mathbb{E}[\nu_i | C_i, X_i] = 0$$

Error term is uncorrelated
with included regressors (89)

- Under CIA, the same ρ that appears in the potential outcomes setup also appears as the regression slope.
- The error ν_i is properly defined so it is orthogonal to both C_i and X_i .
- This establishes the equivalence between regression coefficients and causal effects under CIA.
- **Intuition:** Once we control for X_i , treatment C_i is “as good as random,” so the regression picks up the true causal effect ρ . The regression error is just the part of η_i not explained by X_i , and it no longer contaminates the estimate.

The College Matching Matrix

Applicant group	Student	Private			Public			1996 earnings
		Ivy	Leafy	Smart	All State	Tall State	Altered State	
A	1		Reject	Admit		Admit		110,000
	2		Reject	Admit		Admit		100,000
	3		Reject	Admit		Admit		110,000
B	4	Admit			Admit		Admit	60,000
	5	Admit			Admit		Admit	30,000
C	6		Admit					115,000
	7		Admit					75,000
D	8	Reject			Admit	Admit		90,000
	9	Reject			Admit	Admit		60,000

The college matching matrix shows applicant groups (A–D), their admission decisions across different types of private and public colleges, and their observed earnings in 1996. Note: Enrollment decisions are highlighted in gray.

- Groups A–D represent different applicant pools, with students admitted or rejected across college types (private vs public).
- Highlighted cells show the actual enrollment choices made by each student.
- Average earnings for private college enrollees exceed those for public enrollees by about \$19,500.
- Within matched groups, differences vary: for Group A, private students earn \$5,000 less, while for Group B they earn \$30,000 more, averaging to about \$12,500 overall.
- **Intuition:** Simple comparisons of private vs public students suggest a big payoff, but when we compare students within similar applicant groups (matching on admission profiles), the private advantage shrinks. This illustrates that much of the raw difference comes from selection rather than causal effects.

Private School Effects: Barron's Matches

	No selection controls			Selection controls		
	(1)	(2)	(3)	(4)	(5)	(6)
Private school	.135 (.055)	.095 (.052)	.086 (.034)	.007 (.038)	.003 (.039)	.013 (.025)
Own SAT score ÷ 100		.048 (.009)	.016 (.007)		.033 (.007)	.001 (.007)
Log parental income			.219 (.022)			.190 (.023)
Female			-.403 (.018)			-.395 (.021)
Black			.005 (.041)			-.040 (.042)
Hispanic			.062 (.072)			.032 (.070)
Asian			.170 (.074)			.145 (.068)
Other/missing race			-.074 (.157)			-.079 (.156)
High school top 10%			.095 (.027)			.082 (.028)
High school rank missing			.019 (.033)			.015 (.037)
Athlete			.123 (.025)			.115 (.027)
Selectivity-group dummies	No	No	No	Yes	Yes	Yes

This table reports estimates of the effect of attending a private college or university on earnings. Each column reports coefficients from a regression of log earnings on a private-school dummy and controls. Columns (1)–(3) exclude applicant selectivity controls; columns (4)–(6) include them. Standard errors in parentheses. Sample size: 5,583.

- Without selection controls, the private school effect appears positive and large (around 0,135 log points).
- Once controls are added (SAT scores, family background, demographics, high school rank, etc.), the private coefficient drops toward zero (around 0,007–0,013).
- Including applicant selectivity-group dummies (columns 4–6) confirms the result: little or no private school earnings advantage once comparable students are matched.
- The findings align with a **self-revelation model**: applicants have a good idea of their ability, and admissions outcomes largely reflect that, not private schooling per se.

- **Intuition:** At first glance, private colleges seem to pay off. But once we compare students of similar ability and background, the effect mostly vanishes. The apparent private premium was due to selection — stronger students go to private schools, not because private schools themselves cause higher earnings.

Private School Effects: Average SAT Controls

	No selection controls			Selection controls		
	(1)	(2)	(3)	(4)	(5)	(6)
Private school	.212 (.060)	.152 (.057)	.139 (.043)	.034 (.062)	.031 (.062)	.037 (.039)
Own SAT score ÷ 100		.051 (.008)	.024 (.006)		.036 (.006)	.009 (.006)
Log parental income			.181 (.026)			.159 (.025)
Female			-.398 (.012)			-.396 (.014)
Black			-.003 (.031)			-.037 (.035)
Hispanic			.027 (.052)			.001 (.054)
Asian			.189 (.035)			.155 (.037)
Other/missing race			-.166 (.118)			-.189 (.117)
High school top 10%			.067 (.020)			.064 (.020)
High school rank missing			.003 (.025)			-.008 (.023)
Athlete			.107 (.027)			.092 (.024)
Average SAT score of schools applied to ÷ 100				.110 (.024)	.082 (.022)	.077 (.012)
Sent two applications				.071 (.013)	.062 (.011)	.058 (.010)
Sent three applications				.093 (.021)	.079 (.019)	.066 (.017)
Sent four or more applications				.139 (.024)	.127 (.023)	.098 (.020)

This table reports estimates of the effect of attending a private college or university on earnings. Each column shows coefficients from a regression of log earnings on a private-school dummy and controls. Columns (1)–(3) exclude selection controls, while columns (4)–(6) add them. Additional controls include average SAT score of schools applied to, number of applications, and background covariates. Sample size: 14,238. Standard errors in parentheses.

- Without selection controls (cols. 1–3), the private school effect appears positive (0.21–0.15 log points).

- With selection controls (cols. 4–6), the private school coefficient shrinks sharply (0.034–0.037, sometimes insignificant).
- The **average SAT score of schools applied to** is highly predictive of earnings (0.07–0.11), showing that application choices reveal ability and ambition.
- Other covariates (family income, gender, race, etc.) have smaller effects compared to these selection variables.
- **Intuition:** The raw private-school premium is mostly selection bias. Students who apply to more competitive (high-SAT) schools already have higher potential earnings. Once we control for this revealed-preference measure, the apparent causal effect of attending a private school essentially disappears.

School Selectivity Effects: Average SAT Controls

	No selection controls			Selection controls		
	(1)	(2)	(3)	(4)	(5)	(6)
School average SAT score ÷ 100	.109 (.026)	.071 (.025)	.076 (.016)	-.021 (.026)	-.031 (.026)	.000 (.018)
Own SAT score ÷ 100		.049 (.007)	.018 (.006)		.037 (.006)	.009 (.006)
Log parental income			.187 (.024)			.161 (.025)
Female			-.403 (.015)			-.396 (.014)
Black			-.023 (.035)			-.034 (.035)
Hispanic			.015 (.052)			.006 (.053)
Asian			.173 (.036)			.155 (.037)
Other/missing race			-.188 (.119)			-.193 (.116)
High school top 10%			.061 (.018)			.063 (.019)
High school rank missing			.001 (.024)			-.009 (.022)
Athlete			.102 (.025)			.094 (.024)
Average SAT score of schools applied to ÷ 100				.138 (.017)	.116 (.015)	.089 (.013)
Sent two applications				.082 (.015)	.075 (.014)	.063 (.011)
Sent three applications				.107 (.026)	.096 (.024)	.074 (.022)
Sent four or more applications				.153 (.031)	.143 (.030)	.106 (.025)

This table reports estimates of the effect of attending a more selective institution (measured by the average SAT score of enrolled students) on earnings. Each column shows coefficients from a regression of log earnings on school average SAT and controls. Columns (1)–(3) exclude selection controls, while columns (4)–(6) include them. Sample size: 14,238. Standard errors in parentheses.

- Without selection controls (cols. 1–3), the coefficient on school average SAT is positive and significant (0.109–0.076), suggesting that attending a more selective school is associated with higher earnings.
- With selection controls (cols. 4–6), the coefficient collapses toward zero (–0.021 to –0.008), and is no longer statistically significant.

- This shows that the apparent payoff from school selectivity is explained by student characteristics and preferences (selection), not by causal benefits of attending a higher-SAT peer group.
- **Intuition:** At first glance, selective schools seem to “boost” earnings because they admit stronger students. But once we control for what types of students *apply to and attend* these institutions, the effect of school selectivity disappears. In other words, *it’s not the classmates, it’s the selection*.

Private School Effects: Omitted Variable Bias

	Dependent variable					
	Own SAT score ÷ 100			Log parental income		
	(1)	(2)	(3)	(4)	(5)	(6)
Private school	1.165 (.196)	1.130 (.188)	.066 (.112)	.128 (.035)	.138 (.037)	.028 (.037)
Female		-.367 (.076)			.016 (.013)	
Black		-1.947 (.079)			-.359 (.019)	
Hispanic		-1.185 (.168)			-.259 (.050)	
Asian		-.014 (.116)			-.060 (.031)	
Other/missing race		-.521 (.293)			-.082 (.061)	
High school top 10%		.948 (.107)			-.066 (.011)	
High school rank missing		.556 (.102)			-.030 (.023)	
Athlete		-.318 (.147)			.037 (.016)	
Average SAT score of schools applied to ÷ 100			.777 (.058)			.063 (.014)
Sent two applications			.252 (.077)			.020 (.010)
Sent three applications			.375 (.106)			.042 (.013)
Sent four or more applications			.330 (.093)			.079 (.014)

This table reports regressions of (1)–(3) student ability (own SAT score) and (4)–(6) family background (log parental income) on private school attendance and controls. Large coefficients on “Private school” indicate that private attendance is strongly correlated with student ability and family income, suggesting omitted variable bias in simple regressions of earnings on private schooling.

- Columns (1)–(3) use **Own SAT score** as the dependent variable. Private school attendance is associated with more than a full standard deviation higher SAT score (≈ 1.1), even after adding controls.

- Columns (4)–(6) use **Log parental income** as the dependent variable. Private school attendance is also associated with much higher family income (0,128–0,138), though the effect shrinks with more controls (0,028 in col. 6).
- These results reveal that private schooling is not randomly assigned: it is strongly correlated with both ability and family background.
- **Intuition:** If we simply regress earnings on private school attendance, we will overstate the causal effect. Why? Because part of the “private school premium” comes from omitted variables — smarter students (higher SATs) and richer families are more likely to attend private schools.
- This is a textbook case of *omitted variable bias (OVB)*: the regression coefficient on private schooling loads onto ability and income differences that are not controlled for.

Conclusions

- **Regression** provides the best-in-class approximation to the *Conditional Expectation Function (CEF)*.
- **Regressions** are usually the *first step* in empirical analysis.
- **Regression** is our *primary tool* for confronting the identification problem.
- If the regression you have is not delivering the desired relationship, then the underlying model itself is *unsatisfactory*.
- In such cases: **Instrumental Variables (IV)** is the natural next step.
- **Intuition:** Regression is like the default starting point — it gives us the clearest linear view of how Y relates to X . If this simple picture is misleading or biased, we need stronger tools. That’s where IV comes in: it helps us deal with hidden confounders and extract the true causal effect.