

What is Econometrics

Example Questions:

1. How does an additional year of school effect earnings
2. Does sunscreen really prevent skin cancer

Econometrics uses data to learn about Economic relationships (models)

Econometric Process

1. Question 4. Estimation
2. Model 5. Inference
3. Data

Causal Questions, if I change one thing how does other change

The Model

$Y \rightarrow$ Dependent Variable

$X \rightarrow$ Independent Variable

↳ for each question we must determine which is which

Intro to Linear Bivariate Model

say that,

$$\text{wage}(\text{schooling}) = A + r * \text{schooling}$$

$$\ln(\text{wage}) = \ln[A] + r * \text{schooling}$$

$$\ln(\text{wage}) = E[\ln A] + r * \text{schooling}$$

$$y_i(x)$$

↓

$$y_i(x) = \beta_0 + \beta_1 x + \epsilon_i$$

↳ Linear Bivariate Model

Important: If you can transform relation into linear form you can use OLS

$y_i \rightarrow$ Dependent/outcome Variable

$\beta_0 \rightarrow$ Intercept

$\beta_1 \rightarrow$ Slope or Marginal effect

$\epsilon_i \rightarrow$ Error term or Disturbance

$x_i \rightarrow$ regressor, variable of interest

14.32 Mid term

Empirical Models

we make the assumption

$$E[\epsilon_i | X_i = x] = 0$$

Assumption justified if

1. X_i is randomly assigned
2. Economic Theory implies conditional Moment Restriction
3. control variables related to X unrelated to y

Why Matters

* On average, People w/ different ϵ don't have different y

* Proof of Assumption 2 *

Data

1. Cross-Section → 1 time, many individuals

2. Time series → one thing, many times

3. Panel Data → Many time periods, many individuals

↳ we assume each observation is an i.i.d. Draw from some joint distribution $F(x, y)$

Simple (Bivariate) Regression

If we knew (x_i, y_i) of all individuals we'd know exactly β_1

Sample S

we use a random sample as a proxy for population

Formula for Sample Mean

* Ordinary Least Squares *

Estimate line that minimizes

Average Squared Error

$$\min_{\beta_0, \beta_1} \hat{E}\left[(y_i - (\beta_0 + \beta_1 x_i))^2\right]$$

$$\hat{\beta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Properties of OLS Coeffs

$(\hat{\beta}_0, \hat{\beta}_1)$ are Random Variables

↳ it'll take on different values when we consider different samples

Proofs that $\hat{\beta}'s = \beta'$

(4)

An estimate is an RV who's centered about Pop. value

Predicted or fitted value

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Residual

$$\hat{u}_i = y_i - \hat{y}_i$$

Regression breaks up y into two parts

$$y_i = \hat{y}_i + \hat{u}_i$$

(\hat{y}_i = explained part of y)

\hat{u}_i = unexplained part of y

Multivariate Regression Analysis

Omitted Variable Bias

Derivation of OVB (5)

Factors that influence $\hat{\beta}$ correlate w/ Y cause Bias

control strategies

1. Randomized controlled trials (unethical)

2. Hold other factors fixed (Matching) → decimates sample

Multiple Regression

Multile Regression Analysis

- Putting Omitted Variables into Regression solves problem of bias

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \epsilon_i$$

$$Y = X\beta + \epsilon$$

Derivation of $\hat{\beta}$ matrix

$$\hat{\beta} = (X'X)^{-1}X'Y$$

Predicted values:

$$\hat{y} = X\hat{\beta}$$

Residuals:

$$\hat{u} = y - \hat{y}$$

Properties of Residuals

- Average is zero
- Residual uncorrelated w/ every X

R^2

$$R^2 = 1 - (\text{SSR}/\text{SST})$$

↳ Note: will always increase if we add another variable

Alternate Form of OLS Coeff

$$\hat{\beta}_i = \frac{\text{Cov}(y_i, \hat{r}_{1:i})}{\text{Var}(\hat{r}_{1:i})}$$

↳ $\hat{r}_{1:i}$ is residual from regressing X_i on other included variables

Projections

- Regression is a projection of the outcome variable on the column space of X -variables

$$\begin{aligned}\hat{Y} &= X\hat{\beta} \\ &= X(X'X)^{-1}X'Y \\ &= P_X Y\end{aligned}$$

$$\begin{aligned}\hat{U} &= Y - P_X Y \\ &= (I - P_X) Y \\ &= M_X Y\end{aligned}$$

Multile Regression Beta Form

$$\begin{aligned}\hat{\beta}_i &= \frac{\sum_{j=1}^n y_j \hat{r}_{ij}}{\sum_{j=1}^n \hat{r}_{ij}^2}\end{aligned}$$

↳ \hat{r}_{ij} is residual from regressing X_i on other included variables

$$\begin{aligned}\sum_{i=1}^n \hat{r}_i &= 0 \\ \sum_{i=1}^n \hat{r}_i X_i &= 0\end{aligned}$$

Properties of OLS Estimators

(Gauss-Markov Assumptions)

A0 No perfect collinearity

↳ w/ this assumption $E[\hat{\beta}|X] \rightarrow \beta$

7 $(X'X)^{-1}$ exists

A1: Conditional expectation 0

$$E[\epsilon|X_1, \dots, X_k] = 0$$

A2: Homoskedasticity

A3: No serial correlation **⑧**

$$V(\epsilon|X) = \sigma^2 I$$

$$V(\hat{\beta}|X) = \sigma^2 (X'X)^{-1}$$

$$V(\hat{\beta}_i) = \frac{\sigma^2}{SST_{X_i}(1-R_i^2)}$$

Gauss - Markov theorem

if **A0** \rightarrow A3 hold true

OLS is BLUE

↳ "Best Linear Unbiased Operator"

Proof of Gauss Markov theorem **W**

|| Refer to sheet ||

Inference

- We have $\hat{\beta}_i$, we want to know more about β_i

Hypothesis tests

1. A null hypothesis **H₀**

2. Alternative Hypothesis **H₁**

3. a test statistic, **↑**

4. A distribution of test statistic if **H₀** true

Specific Mechanics of t-test

- we run regression & obtain

$$\hat{\beta} = (X'X)^{-1}X'Y$$

↳ $(k+1) \times 1$ column vector of estimates

- say we want to test if $\beta_j = 0$ or 0.25 or -1

Step 1: Null Hypothesis

$$H_0: \beta_j = \beta_{0j}$$

$$H_1: \beta_j \neq \beta_{0j}$$

Step 2: test statistic

Given,

$$\hat{\sigma}^2 = \frac{SSR}{n-(k+1)}$$

$$\text{s.e}(\hat{\beta}_j) = \sqrt{\frac{\hat{\sigma}^2}{SST_{X_j}(1-R_j^2)}}$$

Step 3: t-stat

$$t_{\hat{b}_j} = \frac{\hat{b}_j - b_0}{\text{s.e}(\hat{b}_j)}$$

Step 4: Distribution

$$\hat{t}_{\hat{b}_j} \sim t(n-(k+1))$$

↳ Distribution

$$P = F_{t(n-(k+1))}(-|t_{\hat{b}_j}|) + 1 - F_{t(n-(k+1))}(|t_{\hat{b}_j}|)$$

↳ Finally you test

if $P \leq \alpha$ "significance level"

↳ Don't reject if true

↳ reject if False

Note: as $n-(k+1) \geq 0, .05$

critical value is around 2

confidence intervals

$$\frac{\hat{b}_j - b_0}{\text{s.e}(\hat{b}_j)} = \pm c_\alpha \begin{pmatrix} \text{critical} \\ \text{value of} \\ \text{significance} \end{pmatrix}$$

$$b_0 = \hat{b}_j \pm c_\alpha (\text{s.e}(\hat{b}_j))$$

- Means out of 100 draws, correct value will be contained 95 times out of 100

Test of Hypotheses on linear combo

- What if we want to test the relationship between more than just one type of parameter

Procedure:

Step 1

$$H_0: \beta_j - \beta_n = 0$$

$$H_1: \beta_j - \beta_n \neq 0$$

Step 2: t-stat

$$t = \frac{\hat{\beta}_j - \hat{\beta}_n}{\text{s.e}(\hat{\beta}_j - \hat{\beta}_n)}$$

Step 3: Distribution

$$\hat{t} \sim t(n-(k+1))$$

standard error (different)

$$\text{s.e}(\hat{\beta}_j - \hat{\beta}_n) = \sqrt{\hat{V}(\hat{\beta}_j - \hat{\beta}_n)}$$

$$= \sqrt{\hat{V}(\hat{\beta}_j) + \hat{V}(\hat{\beta}_n) - 2 \text{Cov}(\hat{\beta}_j, \hat{\beta}_n)}$$

All pieces can be found in Var-Cov matrix

$$\hat{V}(\hat{\beta}) = \frac{\text{SSR}}{n-(k+1)} (X'X)^{-1}$$

Hypotheses on linear combo: General

$$H_0: a_0 B_0 + a_1 B_1 + \dots + a_k B_k = 0$$

$$\hat{a}' \hat{B} = 0$$

$$\hat{T} = \hat{a}' \hat{b} / \sqrt{\hat{a}' \hat{V}(\hat{b}) \hat{a}}$$

Testing Multiple Linear Restriction

t-tests are good for F-test
a test of one relationship

- To test Multiple Linear relationships we perform an F-test

Procedure

1. Step 1

$$H_0: \beta_1 = 0, \beta_m = 0, \beta_n = 0$$

H_1 : one or more of parameters is not zero

Step 2 (F-test)

$$\hat{F} = \frac{(SSR_{\text{restrict}} - SSR_{\text{unrestrict}})/q}{SSR_{\text{unrestrict}}/(n-(k+1))}$$

$$\hat{F} \sim F(q, n-(k+1))$$

Violation of A4: Errors Not Normal

Assumption: Errors Distributed Normally

↳ If normality doesn't hold we still get this condition in large samples

Large Sample Asymptotic OLS Properties

- \hat{b}_j is consistent as sample size gets larger

$$\lim_{n \rightarrow \infty} \Pr(|\hat{b}_j - \beta_j| > \epsilon) = 0$$

- \hat{b} is asymptotically normally distributed

Violation of A2: Heteroskedasticity

Assumption: $\text{Var}(\epsilon_i | X) = \sigma^2$

when it fails: when Error fans out depending on X

Test: Breusch-Pagan Test

steel: run a regression to achieve SSR

step 2: Regress SSR on X

step 3: F-test that all

$$\beta_i^2 | X_i = 0$$

call heteroskedastic if F test fails

Test: White t-test

↳ also add X_1^2 & $X_1 X_2$

$$\hat{V}_{\text{white}}(\hat{b}) = (X'X)^{-1} \left(\sum_{i=1}^n x_i x_i \epsilon_i^2 \right) (X'X)$$

Fix: to make OLS Blue use GLS

Violation of A1: Functional Form

Proposed Model

$$\ln(w_i) = \beta_0 + \beta_1 \text{educ}_i + \beta_2 x_{pi} + \eta_i$$

True Model:

$$\ln(w_i) = \beta_0 + \beta_1 \text{educ}_i + \beta_2 x_{pi} + \underline{\beta_3 x_{pi}^2} + \varepsilon_i$$

- OLS estimate no longer consistent
(return is $\underline{\beta_2 + 2\beta_3 x_p}$)
for experience
- To test: Throw in Polynomial terms & conduct F-test

Violation of A2: Serial Correlation

- When outcome of one observation influenced by outcome of another

when can go wrong:

- Time series data

- Panel Data

- Clustering

Consequences:

- OLS still unbiased & consistent

- Standard errors are wrong

- Not efficient

Suppose: $\hat{H}_t = \rho H_{t-1} + \varepsilon_t$

Test: Durbin-Watson Test

- ① Run original Regression of Y on X

- ② calculate \hat{u}_i & regress on \hat{u}_{i-1}

- ③ T-test $P = 0$

Cov(ε, x) = 0

causes

1. omitted variables

$$y_i = \alpha + \beta x_i + \gamma w_i + \eta_i$$

$$\Leftrightarrow \text{if } \gamma \neq 0, C(x_i, w_i) \neq 0$$

Simultaneous Equations Model

Problem:

$$\text{Cov}(X_i, \varepsilon_j) = 0$$

↳ possible cause: simultaneity

Demand Function

$$Q_t^d = \alpha_t + \beta P_t + Z_t^d + \lambda^d + H_{1t}$$

Q_t^d = quantity of good demanded

P_t = price of good

→ w/ LOG specification β = elasticity

Z_t^d = Demand shifters (something that effects demand but not supply)

Supply Function

$$P_t = \gamma + \delta Q_t + Z_t^s + \lambda^s + H_{2t}$$

→ $1/\delta$ is supply elasticity

→ Z_t^s = Supply shifters

• Equations are called structural Equations

• Coefficients called structural Parameters

Solution: use exogenous shifters

* See equation sheet: Simultaneity Problem

Consequence: OLS is biased (upward prob.)

Step 1:

→ Solve for endogenous variables in terms of exogenous variables

Reduced Form Equations

$$\begin{aligned} \Pi_{11} &= \alpha_1 + \beta_1 P_1 + Z_1^d + \lambda_1^d + H_{11} \\ \Pi_{12} &= \alpha_2 + \beta_2 P_1 + Z_2^d + \lambda_2^d + H_{12} \\ \vdots & \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \Pi_{1n} &= \alpha_n + \beta_n P_1 + Z_n^d + \lambda_n^d + H_{1n} \\ \Pi_{21} &= \alpha_1 + \beta_1 P_2 + Z_1^s + \lambda_1^s + H_{21} \\ \vdots & \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \Pi_{2n} &= \alpha_n + \beta_n P_2 + Z_n^s + \lambda_n^s + H_{2n} \end{aligned}$$

14.32 Final

• Coefficients → Reduced Form Parameters

$$\hat{\Pi}_{12} \rightarrow \lambda^s \frac{\beta}{1-\beta\delta}$$

$$\hat{\Pi}_{22} \rightarrow \lambda^s \frac{1}{1-\beta\delta}$$

$$\hat{\beta} = \hat{\Pi}_{12} / \hat{\Pi}_{22} \rightarrow \boxed{\beta}$$

→ Known as method of indirect least squares

Summary

• Supply shifter Z_t^s serves as an instrument to endogenous variable P_t

conditions on Z_t^s

1. Exclusion from demand Equation

2. Has to be relevant to supply equation

$$\hat{\beta} = \frac{\hat{c}(Z_t^s, Q_t)}{\hat{c}(Z_t^s, P_t)}$$

Instrumental Variables

Problem: Endogenous explanatory variables

Causes: • measurement error
• Omitted Variables
• Simultaneity

* Endogeneity of X_i makes the OLS estimators biased & inconsistent

Simple Bi-variate IV setup

Causal Model: $y_i = \beta_0 + \beta_1 x_i + u_i$

Problem: $\text{cov}(x_i, u_i) \neq 0$

what makes a good instrument?

1. Exclusion $\text{cov}(z_i, u_i) = 0$

2. Relevance $\text{cov}(z_i, x_i) \neq 0$

• can test "Relevance" on t-test of regression

$$X_i = \Pi_{20} + \Pi_{22} Z_i + \varepsilon_i$$

$$H_0: \Pi_{22} = 0$$

• Exclusion just determined by common sense

Identification and Estimation

→ Start w/ $\text{cov}(z_i, y_i)$ and solve

$$\text{for } \beta_1 = \text{cov}(z_i, y_i) / \text{cov}(z_i, x_i)$$

$$\hat{b}_1^{\text{IV}} = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})}$$

$$\text{S.e.}(\hat{b}_1^{\text{IV}}) = \frac{\hat{\sigma}^2}{SST_X \cdot R_{X,Z}^2}$$

→ construct t-tests and confidence intervals as usual

2 Comments

• If X_i not endogenous, and used as instrument for self

$$\hat{b}_1^{\text{IV}} \rightarrow \hat{b}_1^{\text{OLS}}$$

• If X_i not endogenous, but use z any ways still consistent but larger S.e.

IV w/ Multiple Regressors

Regression $\left\{ \begin{array}{l} \times \rightarrow \text{endogenous} \\ Z_{ij} \rightarrow \text{exogenous regressor} \end{array} \right.$

$$y_i = \beta_0 + \beta_1 x_i + Z_i' \lambda_i + \varepsilon_i$$

• Propose Z_2 as indicator

$$-\text{cov}[Z_{2i}, Z_{1i}] = 0 \text{ (exclusion)}$$

$$-\text{cov}[Z_{2i}, X] \neq 0 \text{ (relevance)}$$

Matrices

$$X = \begin{bmatrix} 1 & X_1 & Z_{11} \\ \vdots & \vdots & \vdots \\ 1 & X_n & Z_{1n} \end{bmatrix}, Z = \begin{bmatrix} 1 & Z_{21} & Z_{11} \\ \vdots & \vdots & \vdots \\ 1 & Z_{2n} & Z_{1n} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

Notation of Model

$$\text{Model: } Y = X\beta + \varepsilon$$

$X \rightarrow$ matrix including all RHS variables in model

$Z \rightarrow$ matrix w/ all exogenous variables (included in model or not)

$$\hat{\beta}_{IV} = (Z'X)^{-1} Z'Y$$

↑
see IV estimator
Derivation

Identification

- Need $\# \text{ of RHS Variables} = \# \text{ of Exogenous Variables}$
- $\# \text{ of excluded instruments} = \# \text{ of endogenous variables}$ (order condition)

Two Stage Least Squares (2SLS)

Reasons to use more instruments than needed

- 1. More Precise estimation
- 2. Allows for statistical testing of Assumptions

of instruments $\frac{n}{r}$ \Rightarrow just identified
of variables

of instruments \sqrt{n} \Rightarrow over identified
of variables

SET UP

$$X = \begin{bmatrix} x_{11} & \dots & x_{1q} & x_{p+1,1} & \dots & x_{1n} \\ \vdots & & \vdots & & \vdots & \vdots \\ x_{in} & \dots & x_{ip} & x_{p+1,n} & \dots & x_{nn} \end{bmatrix}$$

$$Z = \begin{bmatrix} z_{1,1} & \dots & z_{1n} & x_{p+1,1} & \dots & x_{kn} \\ \vdots & & \vdots & & \vdots & \vdots \\ z_{pn} & \dots & z_{pn} & x_{p+1,n} & \dots & x_{kn} \end{bmatrix}$$

- P endogenous regressors
- R excluded exogenous instruments
- $R-P$ exogenous regressors

Order condition: $r = p$

Computing 2SLS Estimator

- computing simple IV estimator want work $(Z'X)$ not square

* We must take our $K+r-p$ instrumental variables and combine them to form one instrument for each of the regressors

Exogenous regressors

↳ can act as its own instrument

Endogenous Regressors

- What combination of exogenous instruments would work as best instrument for any variable?

$$W_1 = \Pi_0 + \Pi_1 Z_1 + \dots + \Pi_r Z_r + \Pi_{p+1} X_p + \Pi_K X_K$$

↳ w_i = instrument for x_i

Instrument coefficients are OLS

conditions to make w_i good

1. $\text{cov}(w_i, \varepsilon) = 0$ All z 's uncorrelated

2. $\text{cov}(w_i, x) = 0$

What set of coefficients results in highest R^2 value between x and z 's?

↑
OLS coefficients!

$$W_1 = Z(Z'Z)^{-1} Z'X.$$

$$W_2 = Z(Z'Z)^{-1} Z'X_2$$

Putting it together

$$W = Z(Z'Z)^{-1} Z'X$$

$$= \hat{X}$$

- Now that we have an instrument for each X

$$\hat{\beta}_{2SLS} = (\hat{X}'\hat{X})^{-1}\hat{X}'Y$$

$$\text{or } = (X'P_Z X)^{-1} X' P_Z Y$$

$$P_Z = \begin{bmatrix} & & \\ & I & \\ & & Z(Z'Z)^{-1} Z' \end{bmatrix}$$

Reasons why its called 2SLS

1. Regress endogenous variables on ALL exogenous instruments

2. Regress y on predicted values from first stage

Proof: simple IV = Just Identified
Testing Endogeneity

- Only want to use IV if we have to, will be less accurate

Test of Endogeneity

step1: Regress Potentially Endogenous variables on all exogenous variables

calculate residuals \hat{V}_2

step2: include residual \hat{V}_2 + OLS regression of structural Equation

step3: test if its coefficient is zero
reject → endogenous
Don't Reject → exogenous

Testing Validity of Instruments	Time Series	Trends and Seasonality
<p><u>1. Hausman T test</u></p> <p>↳ requires over-identification</p> <p>Idea: (a) 2SLS (overidentified) (b) Just-identified (subset of all possible instruments)</p> <p>These two should be consistent</p> $\hat{T} = \frac{\hat{b}_{\text{Just identified}} - \hat{b}_{\text{2SLS}}}{\text{S.e.}(\hat{b}_{\text{JI}} - \hat{b}_{\text{2SLS}})}$ <p>$H_0: \hat{T} = 0$ instruments valid $H_1: \hat{T} \neq 0$ one or more not valid</p>	<p>We now look at how outcome variable evolves over time (index by time t)</p> <p>Model: $y_t = \beta_0 + \beta_1 x_{1t} + \dots + \beta_k x_{kt} + u_t$ $y_t' = \beta + u_t, t=1, \dots, n$</p> <p>Important difference between time series and cross-section</p> <p>We allow one observation to depend on past observations</p> <p>$y_t = \text{houseprice}, x_{1t} = \text{int}_t, x_{2t} = \text{int}_{t-1}$</p> <p>Dynamic Models Dependent Data create problems</p>	<ul style="list-style-type: none"> Time series can grow over time → model this by including Time in Regression $y_t = d_0 + d_1 t + \dots + H_t$ <ul style="list-style-type: none"> Trends also exhibit seasonality → Model this by including dummy variables for months of the year $y_t = d_0 + d_1 \text{jan}_t + d_2 \text{feb}_t + d_3 \text{mar}_t + \dots + d_{12} \text{dec}_t + H_t$
<p><u>2. Omnibus test</u></p> <p>(a) estimate via 2SLS to form residuals</p> <p>(b) Regress \hat{u}_t on all instruments and get R^2</p> <p>(c) R^2 will follow χ^2 with $r-a-p$ degrees of freedom</p>	<p>Types of time series Models</p> <p><u>1. Static Model</u> $y_t = \beta_0 + \beta_1 z_t + H_t$</p> <p><u>2. Finite Distributed Lag (FDL)</u> $y_t = d_0 + \delta_0 z_t + \delta_1 z_{t-1} + H_t$</p>	<p>Observations from one period often depend on observations from another</p> <p>Moving Average Process of order one</p> $x_t = e_t + d_1 e_{t-1}$ <p>Autoregressive Process AR(1)</p> $y_t = \rho y_{t-1} + e_t$ <p>→ to be stable, $\rho < 1$</p> <p>weak/strong dependence processes</p>
<p>Examples of IV Estimation</p> <p><u>1. Demand Elasticity</u> Endogenous: simultaneity Instrument: harvest-yield</p> <p><u>2. military service on latent earnings</u> Endogenous: high self-selection Instrument: draft-eligibility</p> <p><u>3. effect of many districts on student performance</u> Endogenous: weather areas have more exogenous: # of major streams in area</p>	<p>Assumptions for Time Series</p> <p><u>A 0</u> same as before, No perfect collinearity $(X'X)^{-1}$ exists</p> <p><u>A 1</u> $E[u_t X] = 0, H_t$ must be uncorrelated across all time periods</p> <p>↳ strict exogeneity</p> <p><u>A 2</u> same as before: Homoskedasticity</p> <p><u>A 3</u> same as before: No serial correlation ↳ much more likely to fail consequences</p> <p>A0-A2 \Rightarrow unbiased</p>	<p>weakly dependent processes</p> <p>↳ correlation between two observations dies down as time between them increases</p> <p>Strongly dependent (AR(1) where $\rho \neq 1$)</p> <p>$y_t = y_{t-1} + e_t$ ↳ effectively a random walk</p> <p>$\text{Var}(y_t) = \sigma_e^2 t \rightarrow \infty$</p> <p>$\text{Corr}(y_t, y_{t+h}) = \sqrt{t/(t+h)} \rightarrow 1$</p>

Solution for high persistence processes

- Random Walk w/ Drift

$$y_t = \alpha_0 + y_{t-1} + \epsilon_t$$

- Solution: First-Difference is weakly dependent

\Downarrow (transformed)

$$\Delta y_t = y_t - y_{t-1} = \epsilon_t$$

↳ nicely behaved

Panel Data

- Many Observations over Many time Periods
- We label individuals using index i , and time periods using index t
- # of individuals = N
of time Periods = T

What Panel is used for?

- Learn about dynamic relationships
- control for time-invariant Variables

Panel Data: Workhorse Model

$$y_{it} = \beta_0 + W_i' \beta_1 + S_i' \beta_2 + Z_i' \beta_3 + d_i + v_{it}$$

W_i → characteristics of individuals that don't change over time

S_i → variables that change over time across individuals

Z_i → variables that change over time but are the same for all individuals

d_i ** individual effect

Estimation: Panel Models

- What would happen if we ran regression with all variables

↳ if d_i : uncorrelated (consistent
unbiased)

↳ if d_i : correlated (inconsistent
biased)

First Differencing & Fixed effects

- Imagine forming a new variable (year-to-year change in y)

$$\begin{aligned} \Delta y_{it} &= y_{it} - y_{it-1} \\ &= \Delta S_i' \beta_2 + \Delta Z_i' \beta_3 + \Delta d_i + \Delta v_{it} \end{aligned}$$

↳ using this we can at least get β_2 and β_3

- Option 2: Run OLS with indicator Variable for each individual

$$y_{it} = \beta_0 + W_i' \beta_1 + S_i' \beta_2 + Z_i' \beta_3 + \sum_{j \neq i} d_j + v_{it}$$

identical to making this transformation v_{it}

$$\tilde{y}_{it} = y_{it} - \bar{y}_i;$$

(Running fixed effects) $= \tilde{S}_i' \beta_2 + \tilde{Z}_i' \beta_3 + \tilde{d}_i + \tilde{v}_{it}$

Random Effects

- Assume d_i are uncorrelated with the regressors

Random Effects

$$\tilde{y}_{it} = y_{it} - \hat{\lambda} \bar{y}_i;$$

$$\tilde{x}_{it} = x_{it} - \hat{\lambda} \bar{x}_i;$$

$$\hat{\lambda} = 1 - \sqrt{\frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + T \hat{\sigma}_\alpha^2}}$$

- Run Pooled OLS on transformed variables

Important consequences

1. Will be consistent
2. does allow estimation of coefficients on time-invarying variables

3. Proof: Prove $\hat{\beta}_{RE}$ is Blue by proving A3 holds

Difference in Differences

example: effect of NYC soft drink ban on obesity rates

• People of NYC: Treated = 1

• People of Westchester = 0

• Cross-section time series alone will fail

Model:

$$\text{obesity}_{it} = \beta_0 + \beta_1 \text{Treated}_i + \beta_2 \text{Post}_t + \beta_3 \text{Treated}_i * \text{Post}_t + \beta_4 \text{D}_it$$

$$\Delta$$

→ Difference in Post/Pre-Difference between NYC and Westchester

Non-linear Models

- Suppose we explore age of marriage vs. Probability of Divorce

• $y_i = \text{Divorce}(0 or 1)$

$x_i = \text{Age of Marriage}$

Problem: Probability has to fall between 0 and 1 but linear function of x doesn't have to

- We model effect of X on a Binary Dependent variable as

$$\Pr(y_i=1|x_i) = G(x_i; \beta)$$

$\hookrightarrow G(\cdot)$ link function that lies between zero and one

Maximum Likelihood Estimation

Data

obs	y_i	x_i
1	1	x_1
2	0	x_2
3	1	x_3
4	1	x_4
n	1	x_n

- Probability that $y_i = 1$

$$G(x_i; \beta)$$

- Probability that $y_i = 0$

$$1 - G(x_i; \beta)$$

- In general $\Pr(y_i = y)$

$$f(y_i|x_i; \beta) = G(x_i; \beta)^y (1 - G(x_i; \beta))^{1-y}$$

Likelihood Function

$$L(Y|X; \beta) = \prod_{i=1}^n f(y_i|x_i; \beta)$$

$$\text{Log}(L) = \sum_{i=1}^n y_i \ln(G(x_i; \beta)) + (1-y_i) \ln(1 - G(x_i; \beta))$$

$$\hat{\beta}_{\text{MLE}} = \arg \max_{\beta} \left(\ln [L(\beta)] \right)$$

Link Functions:

Probit $\rightarrow G(\cdot) \rightarrow \Phi(\cdot)$

Logit $\rightarrow G(\cdot) \rightarrow \frac{\exp(\cdot)}{1 + \exp(\cdot)}$

Quantile Regression

- The Models we've looked at have explored how explanatory variables effect the mean of the outcome variable

Setup

$$\text{Quantile Function: } Q_Y(\tau) = \alpha_\tau : \Pr(Y \leq \alpha_\tau) = \tau$$

also,

$$Q_Y(\tau) = F_Y^{-1}(\tau)$$

Model

$$Y_i = X_i' \beta(\tau) + \varepsilon_i(\tau)$$

$$Q_{\varepsilon|X}(\tau | X_i) = 0$$

Quantile Regression Estimation

$$\hat{\beta}_{QR} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n p_+(\varepsilon_i - X_i' \beta)$$

$\hookrightarrow p_+$ is a check function

$$p(\varepsilon) = (\tau - \mathbb{I}(\varepsilon < 0)) \varepsilon$$