

Revision Notes by Sally Yang

# PRINCIPLES OF ECONOMETRICS

LENT TERM

DR MARCIA SCHAFGANS  
LONDON SCHOOL OF ECONOMICS 2020/21

# LINEAR ALGEBRA

$$(A^{-1})' = (A')^{-1}$$

$$(AB)' = B'A'$$

$(AB)^{-1} = B^{-1}A^{-1}$  (if A and B are both invertible)

$$|A| = |A'|$$

$$|AA^{-1}| = |A||A^{-1}| = 1 \quad (\text{if } A \text{ is invertible})$$

$$\Rightarrow |A^{-1}| = \frac{1}{|A|}$$

$$\Rightarrow |A| \neq 0$$

Let  $X$  be  $n \times k$  matrix.  $X = \begin{pmatrix} x_1 & \dots & x_k \\ x_m & \dots & x_n \end{pmatrix} = \begin{pmatrix} \overbrace{x_1}^{\vdots} & \dots & \overbrace{x_k}^{\vdots} \\ \vdots & \dots & \vdots \\ x_m & \dots & x_n \end{pmatrix}, \quad x_j' = (x_{j1} \dots x_{jk})$

Then  $XX' = \sum_{j=1}^n x_j x_j' = \sum_{j=1}^n \begin{pmatrix} x_{j1} \\ \vdots \\ x_{jk} \end{pmatrix} (x_{j1} \dots x_{jk})$  XX' is symmetric and positive semidefinite.

Positive definite if columns of X lin. indep.

★  $Z'X = \sum_{j=1}^n z_j x_j'$ . Letter in same order but transpose second

Usefulness:  $\frac{XX'}{n} = \frac{\sum x_j x_j'}{n}$  is sample average

★  $(X'X)^{-1}X'y \neq (X^{-1})(X')^{-1}X'y = X^{-1}y$  because  $X$  may not be invertible

★  $X'y = \sum_{j=1}^n x_j y_j$  scalar

**LINERLY INDEPENDENT**  $x_1, \dots, x_k$  are linearly independent if

The linear combination  $\alpha_1 x_1 + \dots + \alpha_k x_k = 0 \Rightarrow \alpha_1 = \dots = \alpha_k = 0$

**ORTHOGONAL**

$x \perp y$  if  $x \cdot y = x'y = 0$  (for vector)

Matrix  $X$  is orthogonal if  $X' = X^{-1}$ . So  $XX' = XX^{-1} = I \Rightarrow |X|^2 = 1 \Rightarrow |X| = \pm 1$

**RANK**  $P(A) =$  Number of linearly independent rows

**EIGENVALUES**  $x$  (non-zero) is an eigenvector of matrix  $A$  if  $Ax = \lambda x$  (Eigenvalue  $\lambda \in \mathbb{R}$ )

Find eigenvalues:  $|A - \lambda I| = 0$

**DIAGONALISING** After finding eigenvectors, eigenvalues. We can write  $A = PDP^{-1}$  hence  $P^{-1}AP = D$

$$D = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}, \quad P = \begin{pmatrix} x_1 & \dots & x_n \\ \vdots & \ddots & \vdots \\ x_m & \dots & x_n \end{pmatrix}$$

$$|A| = |D| = \prod_{i=1}^n \lambda_i$$

**ORTHOGONAL DIAGONALISING**  $P = P^{-1}$ , so  $A = PDP^{-1} = PDP'$  or  $D = P'AP$

Works for any symmetric matrix  $A$ . If  $A$  is also positive definite then  $\exists Q$  st.  $A = QQ'$

**DEFINITENESS** A symmetric matrix  $A$  is positive definite if  $\forall x \in \mathbb{R}^n, x'Ax > 0$

	Eigenvalues	$\forall x, x'Ax > 0$
+ve def	$> 0$	$> 0$
+ve semidef	$\geq 0$	$\geq 0$
-ve def	$< 0$	$< 0$
-ve semidef	$\leq 0$	$\leq 0$
indef	$\exists > 0, \exists < 0$	$\exists > 0, \exists < 0$

**TRACE**  $\text{tr}(A)$  is sum of diagonal elements  $\sum_{i=1}^n a_{ii}$

- $\text{tr}(A) = \text{tr}(A')$
- $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$
- $\text{tr}(AB) = \text{tr}(BA)$   
so  $\text{tr}(\underbrace{ABCD}) = \text{tr}(\underbrace{DABC}) = \text{tr}(\underbrace{DAB}) = \text{tr}(\underbrace{BCDA})$

**IDEMPOTENT** A is idempotent if  $AA = A \Rightarrow A^k = AAA \dots A = A$

- A is always diagonalisable and its eigenvalues are either 0 or 1
- If A is idempotent,  $\text{tr}(A) = p(A) = \dots$

Invertible	Square, $ A  \neq 0$ , full row and column rank (no rows/columns of 0s)
Diagonal	linearly indep. row/columns, $A'$ invertible. $\lambda_i \neq 0 \forall i$
Symmetric	Square, Off-diagonal 0s
Orthogonal	Square, $A = A'$
Idempotent	Square, $AA' = A'A = I$ . $A = A^{-1}$
	Square, $AA = A$

I satisfies all

## VECTOR DIFFERENTIATION

$$f: \mathbb{R}^n \xrightarrow{\substack{\text{nxi vector} \\ \text{output is scalar}}} \frac{\partial f}{\partial \underline{x}} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

### PROPERTIES

linear form (also  $\underline{a} \cdot \underline{x}$ )

- If  $\underline{a}$  and  $\underline{x}$  are nxi vectors,  $\frac{\partial}{\partial \underline{x}} \underline{a}' \underline{x} = \underline{a}$
- $\frac{\partial}{\partial \underline{x}} \underline{x}' A \underline{x} = A \underline{x} + A' \underline{x}$   
 $= 2A \underline{x}$  when A is symmetric
- $\frac{\partial}{\partial A} \underline{x}' A \underline{x} = \underline{x} \underline{x}'$
- $\frac{\partial}{\partial \underline{x}} \underline{x}' A \underline{y} = A \underline{y}$

### SQUARE VECTORS

NOT the inner product  $\underline{m}' \underline{m}$  (gives scalar) but the outer product  $\underline{m} \underline{m}'$  (nxn matrix!)

# MULTIVARIATE STATISTICS

A random vector/matrix  $\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix} / X = \begin{pmatrix} x_1 & \dots & x_n \\ \vdots & \ddots & \vdots \\ x_k & \dots & x_{nk} \end{pmatrix}$  is a vector/matrix where elements are random variables

## EXPECTATIONS

$$E(\underline{x}) = \begin{pmatrix} E(x_1) \\ \vdots \\ E(x_k) \end{pmatrix} \text{ so } E(\underline{\varepsilon}) = 0 \text{ implies } \begin{pmatrix} E(\varepsilon_1) \\ \vdots \\ E(\varepsilon_n) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

Rules of expectation work normally :  $E(a + B\underline{x}) = a + BE(\underline{x})$

## VARIANCE

$$\text{Var}(\underline{x}) = E[(\underline{x} - E\underline{x})(\underline{x} - E\underline{x})'] = \begin{pmatrix} \text{Var}(x_1) & & \\ \text{Cov}(x_1, x_2) & \dots & \\ \vdots & & \\ \text{Cov}(x_m, x_1) & \dots & \text{Var}(x_n) \end{pmatrix} \text{ always positive semidefinite!}$$

Rules of variance :  $\text{Var}(a + B\underline{x}) = \text{Var}(B\underline{x}) = B\text{Var}(\underline{x})B' = B^2\text{Var}(\underline{x})$  if  $B$  is a scalar

## CONDITIONAL EXPECTATIONS

For every  $\underline{x}$ , we have this object  $y|X=\underline{x}=Z$  which is a random variable

$P(Z \leq w) = F_Z(w) = P(y|X=\underline{x} \leq w) = P(y \leq w|X=\underline{x}) = F_{y|X=\underline{x}}(w)$  are all equivalent

Density of  $f_{Y|X=\underline{x}}$  is the density of  $Y$  conditional on  $X=\underline{x}$

$E(y|X=\underline{x}) = \int y f(y|X=\underline{x}) dy$  is a function of  $\underline{x}$ . Same for  $\text{Var}(Y|X=\underline{x})$

## PROPERTIES

$$\begin{aligned} E(f(X) + Y|X) &= E(f(X)|X) + E(Y|X) \\ &= f(X) + E(Y|X) \text{ since conditioning on } X \text{ allows us to think of } X \text{ as non-random anymore.} \end{aligned}$$

$$\text{Var}(f(X) + Y|X) = \text{Var}(Y|X)$$

$$\text{Var}(f(X)Y|X) = f(X)E(Y|X)f(X)'$$

Difference between  $E(y|X=\underline{x})$  and  $E(y|X)$

$E(y|X=\underline{x})$  is a function of  $\underline{x}$ , and is not random.

If we let  $E(y|X=\underline{x}) = h(\underline{x})$ , then  $E(y|X) = h(X)$ .

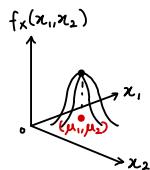
So  $E(y|X)$  is a random variable that takes the values  $E(y|X=\underline{x})$  with density function  $f_X(\underline{x})$

Same for  $\text{Var}(Y|X=\underline{x}) = g(\underline{x}) \Rightarrow \text{Var}(Y|X) = g(X)$

## MULTIVARIATE NORMAL DISTRIBUTION

$$x \sim N(\mu, \text{Var}(x))$$

pdf:  $f_x(x) = \frac{1}{\sqrt{\det(2\pi \text{Var}(x))}} e^{-\frac{1}{2}(x-\mu)' \text{Var}(x)^{-1} (x-\mu)}$



If  $a$  and  $B$  are fixed.

$$a + Bx \sim N(a + B\mu, B\text{Var}(x)B')$$

$\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \sim N(\mu, \text{Var}(z))$  implies its elements are univariate normal  $z_j \sim N_1(\mu_j, (\text{Var}(z))_{jj})$

Converse true only if  $z_1, \dots, z_n$  are independent

just let  $B = \begin{pmatrix} 0 & \dots & 0 & \overset{j^{\text{th entry}}}{\downarrow} & 0 & \dots & 0 \end{pmatrix}$

## UNCORRELATED $\Rightarrow$ INDEPENDENCE

If  $\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \sim N(\mu, \text{Var}(z))$ , and  $z_1, \dots, z_n$  are identically distributed and all uncorrelated  
then it implies independence  $\text{cov} = 0$

# REVIEW

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

## ASSUMPTION

- Mean independence  $E(\varepsilon|X) = 0$

$$\Rightarrow E(Y|X) = E(\alpha + \beta X + \varepsilon|X) = \alpha + \beta X + E(\varepsilon|X) = \alpha + \beta X$$

$$\Rightarrow \beta = \frac{\partial E(Y|X)}{\partial X}$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

"partial effect": effect of  $X_1$  on  $Y$  holding all else ( $X_2$ ) constant

Adding more regressors  $\uparrow R^2$

## ASSUMPTION

- Mean independence  $E(\varepsilon|X_1, X_2) = E(\varepsilon|X) = 0$

$$\Rightarrow E(Y|X) = E(\alpha + \beta_1 X_1 + \beta_2 X_2 + \varepsilon|X) = \alpha + \beta_1 X_1 + \beta_2 X_2 + E(\varepsilon|X) = \alpha + \beta_1 X_1 + \beta_2 X_2$$

$$\Rightarrow \beta_1 = \frac{\partial E(Y|X)}{\partial X_1}$$

## SAMPLING

Obtain a sample  $\{(Y_i, X_i)\}_{i=1}^n$  from population

- Effectively drawing from  $(X, \varepsilon)$  joint dist.
- In **cross-section**, maybe can assume  $(\varepsilon_i, X_i)$  are iid drawings from pop<sup>n</sup>
- In **time-series**,  $(\varepsilon_i, X_i)$  and  $(\varepsilon_j, X_j)$  usually related
- $X_i, \varepsilon_i$  commonly correlated (OVB)

## OLS

$$(\hat{\beta}_0, \hat{\beta}_1) = \underset{\substack{\uparrow \\ b_0, b_1 \\ \text{estimators}}}{\arg \min} \sum_{i=1}^n (Y_i - b_0 - b_1 X_i)^2$$

(realisation of them for particular samples are estimates)

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad \hat{\beta}_1 = \frac{\text{Cov}(Y_i, X_i)}{\text{Var}(X_i)} \quad \begin{matrix} \leftarrow \text{sample cov} \\ \leftarrow \text{sample var} \end{matrix}$$

**RESIDUALS**  $\hat{\varepsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \quad \star \text{NOT same as error } \varepsilon_i$

**FITTED VALUES**  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

## PROPERTIES OF SAMPLING DIST OF $\hat{\beta}_1$

- Unbiasedness  $E(\hat{\beta}_1) = \beta_1$  ("correct on average")
- Standard error  $SE(\hat{\beta}_1) = \sqrt{\text{Var}(\hat{\beta}_1)}$   
 $\uparrow \text{estimated var of } \hat{\beta}_1$

# TERMS

Sample  $\{(x_i, y_i)\}_{i=1}^n$ ,  $x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ik} \end{pmatrix}$  Individual characteristics of person i

Joint density of n observations of the same variable,  $f(x_1, \dots, x_n) \stackrel{\text{vector}}{=} f(\underline{x})$

Marginal density  $f(x_i) = \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_2 \dots dx_n$

Conditional density  $f(\varepsilon|x) = \frac{f(\varepsilon, x)}{f(x)}$   $\leftarrow$  joint  $\leftarrow$  marginal

Under independence,

$f(\varepsilon, x) = f(\varepsilon|x)f(x) = f(\varepsilon)f(x)$  product of marginals

$$E(\varepsilon|x) = \int_{-\infty}^{\infty} \varepsilon f(\varepsilon|x) d\varepsilon = \int_{-\infty}^{\infty} \varepsilon f(\varepsilon) d\varepsilon = E(\varepsilon)$$

Covariance matrix of  $\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$  is  $\text{Var}(\varepsilon) = E[(\varepsilon - E(\varepsilon))(\varepsilon - E(\varepsilon))']$ ,  $E(\varepsilon) = \begin{pmatrix} \varepsilon_1 - E(\varepsilon_1) \\ \vdots \\ \varepsilon_n - E(\varepsilon_n) \end{pmatrix}$

$$\text{matrix with Var on diagonal and cov on off-diagonal}$$

$$= \begin{pmatrix} E((\varepsilon_1 - E(\varepsilon_1))^2) & E((\varepsilon_1 - E(\varepsilon_1))(\varepsilon_2 - E(\varepsilon_2))) & \dots & E((\varepsilon_1 - E(\varepsilon_1))(\varepsilon_n - E(\varepsilon_n))) \\ \vdots & E((\varepsilon_2 - E(\varepsilon_2))^2) & \dots & E((\varepsilon_2 - E(\varepsilon_2))(\varepsilon_n - E(\varepsilon_n))) \\ E((\varepsilon_n - E(\varepsilon_n))^2) & E((\varepsilon_n - E(\varepsilon_n))(\varepsilon_1 - E(\varepsilon_1))) & \dots & E((\varepsilon_n - E(\varepsilon_n))(\varepsilon_n - E(\varepsilon_n))) \end{pmatrix}$$

$E(c + a'x) = c + a'E(x)$   $\leftarrow$  vector of constants  $\downarrow$  vector of variables

$$E(c + Ax) = c + AE(x)$$

$\uparrow$  vector of constants  $\uparrow$  matrix of constants

$$\text{Var}(c + \sum_{i=1}^n a_i x_i) = \sum_{i=1}^n a_i^2 \text{Var}(x_i) + \sum_{i=1}^n \sum_{j=1, j \neq i}^n a_i a_j \text{Cov}(x_i, x_j)$$

$$\text{Var}(c + a'x) = a' \text{Var}(x) a$$

$$\text{Var}(c + Ax) = A \text{Var}(x) A'$$

$\uparrow$  Covariance matrix!

# CLASSICAL (GAUSS-MARKOV) LINEAR REGRESSION

$$y = X\beta + \varepsilon, y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, X = \begin{pmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nk} \end{pmatrix} = (x_{ij}), \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

n observations  
k variables  
k+1 parameters (+constant)

- for each observation i,  $y_i = x_i \beta + \varepsilon_i, \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ik} \end{pmatrix}$

## ASSUMPTIONS

- True model is linear in parameters.  $y = X\beta + \varepsilon, E(\varepsilon) = 0$ 
  - If  $E(\varepsilon) \neq 0$ , build it into systematic part of eqn, adding new vars, make = 0
- No perfect multicollinearity (e.g. dummy variable trap).
  - $X$  has full column rank  $k \leq n$ , lin. indep. rows;  $X'X$  invertible in deriving OLS estimator
- $E(\varepsilon|X) = 0$ . Elab below
- Homoscedasticity, no auto-correlation  $\Rightarrow \text{Var}(\varepsilon|X) = E(\varepsilon\varepsilon'|X) = \sigma^2 I$ 

$$\text{Var}(\varepsilon_i|X) = \sigma^2 \forall i \quad \text{Cov}(\varepsilon_i, \varepsilon_j|X) = 0 \forall i \neq j \Rightarrow \text{Cov}(y_i, y_j|X) = 0$$
  - Reduces to  $\text{Var}(\varepsilon) = \sigma^2 I$  if  $X$  is fixed (non-stochastic)
  - Typically works for cross-section but not time series
  - Heteroscedastic if  $\text{Var}(\varepsilon_i|X)$  is some function of  $X$
- Normality  $\varepsilon|X \sim N(0, \sigma^2 I)$  A1-4 is Gauss Markov. A1-5 (GM+Normality) is classical lin reg model
  - A5 implies A3 and A4; also implies  $\varepsilon \sim N(0, \sigma^2 I)$  since  $E(\varepsilon|X), \text{Var}(\varepsilon|X)$  don't depend on  $X$ 
    - $\hat{\beta}|X \sim N_k(\beta, \sigma^2(X'X)^{-1}) \Rightarrow \hat{\beta}_j|X \sim N_j(\beta_j, \sigma^2(X'X)^{-1}_{jj}) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2(X'X)^{-1}_{jj}}} |X \sim N(0, 1) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2(X'X)^{-1}_{jj}}} \sim N(0, 1)$  H-testing!
  - Even without A5, if  $n$  is large, we can use CLT to say A5 holds "approximately"

## ZERO CONDITIONAL MEAN $E(\varepsilon|X) = 0$

- Regressors strictly exogenous.
- Ensures  $X$  and  $\varepsilon$  uncorrelated,  $\text{Cov}(x_i, \varepsilon_j) = 0$
- A1+A3  $\Rightarrow E(y|X) = X\beta$  and  $\frac{\partial E(y|X)}{\partial x_j} = \beta_j$ .
- Need for assumption arises from regressors being stochastic (random). So the  $x$  should be conditioned on.
  - If  $X$  is fixed (non-stochastic), then it just becomes  $E(\varepsilon) = 0$  and no need for A3
  - If samples are drawn independently ( $\varepsilon_i$  not correlated to  $x_j \forall i \neq j$ ), A3 reduces to  $E(\varepsilon_i|x_i) = 0$
  - Independence  $\Rightarrow$  Mean independence  $\Rightarrow$  Uncorrelatedness 2020 IRDAP 3
 
$$\begin{aligned} E(\varepsilon|X) &= E(\varepsilon)E(X) = 0 & \text{Typically too strong for reality} \\ E(\varepsilon|X) &= E(\varepsilon_i) & \text{Cov}(x_i, \varepsilon_j) = 0 \\ &= E(Y|X) - E(Y)E(X) \\ &= E(Y|X) = E(X(E(Y|X))) \\ &= E(X)E(Y) \end{aligned}$$
- Uncorrelatedness  $\Rightarrow$  Mean independence if we have multivariate normal  $\zeta$  with uncorrelated elements
  - Multivariate normal  $\varepsilon|X$  (A5) + uncorrelated elements  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \forall i \neq j$  (A4)
- $Y = X\beta + \varepsilon, E(\varepsilon) = 0$  (A1) and  $E(\varepsilon|X) = 0$  (A3)  $\Leftrightarrow E(Y|X) = X\beta$

## LAW OF ITERATED EXPECTATIONS $E(h(x,y)) = E(E(h(x,y)|x))$ assuming these expectations exist $x, y$ must both be random

OLS ESTIMATOR

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} = \underset{b}{\operatorname{argmin}} \sum_{i=1}^n (y_i - b_1 x_{i1} - b_2 x_{i2} - \dots - b_k x_{ik})^2 = \underset{b}{\operatorname{argmin}} (y - Xb)(y - Xb)'$$

FOCs: 
$$\begin{cases} \sum_{i=1}^n x_{i1} \hat{\varepsilon}_i = 0 \\ \vdots \\ \sum_{i=1}^n x_{ik} \hat{\varepsilon}_i = 0 \end{cases}, \hat{\varepsilon}_i = y_i - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \dots - \hat{\beta}_k x_{ik}$$

$$\Leftrightarrow \underset{k \times n}{X' \hat{\varepsilon} = 0} \Leftrightarrow \underset{n \times 1}{X'(y - X\hat{\beta}) = 0} \Leftrightarrow \hat{\beta} = \underset{\substack{\text{exists by A2} \\ (\text{otherwise } \beta \text{ not identified})}}{(X'X)^{-1} X'y}$$

# PROPERTIES OF $\hat{\beta}$

## MINIMISES RESIDUAL SUM OF SQUARES

RSS is  $S(b) = \sum_{i=1}^n (y_i - b_1 x_{i1} - \dots - b_k x_{ik})^2 = (y - Xb)'(y - Xb) = y'y - 2\hat{y}'X'y + b'X'Xb$

When  $S(b)$  is at minimum,

$$\frac{\partial S(b)}{\partial b} = \left( \begin{array}{c} \frac{\partial S(b)}{\partial b_1} \\ \vdots \\ \frac{\partial S(b)}{\partial b_k} \end{array} \right) = -2X'y + 2X'Xb = 0 \Rightarrow \hat{\beta} = (X'X)^{-1}X'y$$

## METHOD OF MOMENTS ESTIMATOR

By A1,3,  $E(x_{is}\varepsilon_i) = 0$  and  $E(x_{is}\varepsilon_i) = E(x_{is}(y_i - x_{is}\beta))$ ,  $s = 1, \dots, n$

In MME, you replace the theoretical expected value with sample average

$$\frac{1}{n} \sum_{i=1}^n x_{is}\hat{\varepsilon}_i = 0 \quad \text{with } \hat{\varepsilon}_i = y_i - x_i'\hat{\beta}_{MME}$$

but this coincides with OLS requirements so  $\hat{\beta}_{MME} = \hat{\beta}_{OLS}$

\*  $E(x_{is}\varepsilon_i) = 0$

\* "Existence of moments" If  $n^{th}$  moment exists,  $1 - n-1^{th}$  moment exists

MME :

Step 1: replace  $E$  with  $\frac{1}{n} \sum_{i=1}^n$

Step 2: put hats on all parameters  
replace  $\varepsilon_i$  with  $\hat{\varepsilon}_i$

finite sample property  
i.e. true for any  $n$

## UNBIASED $E(\hat{\beta}) = \beta$

$$\text{Proof: } E(\hat{\beta}) = E[(X'X)^{-1}X'y] = E[(X'X)^{-1}X'(X\beta + \varepsilon)] = E(\beta) + E[(X'X)^{-1}X'\varepsilon] = \beta + E[(X'X)^{-1}X'\varepsilon] = \beta!$$

If  $X$  is fixed,  $E[(X'X)^{-1}X'\varepsilon] = (X'X)^{-1}X'E(\varepsilon) = 0$

If  $X$  is stochastic,

under full <sup>too strong</sup> independence,  $E[(X'X)^{-1}X'\varepsilon] = E[(X'X)^{-1}X']E(\varepsilon) = 0$

under mean independence,  $E[(X'X)^{-1}X'\varepsilon] \stackrel{\text{def}}{=} E(E[(X'X)^{-1}X'\varepsilon|X]) = E((X'X)^{-1}X'E(\varepsilon|X)) = 0$

$E(\varepsilon_i|X) = E(\varepsilon_i|x_i)$  if the sample  $\{(y_i, x_i) : i=1, \dots, n\}$  is random

finite sample property

## VARIABILITY (THEOREM) $\text{Var}(\hat{\beta}|X) = \sigma^2(X'X)^{-1}$

$$\text{Proof: } \text{Var}(\hat{\beta}|X) = \text{Var}(\beta + (X'X)^{-1}X'\varepsilon|X) = (X'X)^{-1}X' \text{Var}(\varepsilon|X) [(X'X)^{-1}X]' = (X'X)^{-1}X'\sigma^2 I X (X'X)^{-1} = \sigma^2 (X'X)^{-1}$$

$$E[\text{Var}(\hat{\beta}|X)] = \text{Var}(\hat{\beta}) \quad \text{PS2}$$

\* Implies that in bivariate regression,  $\text{Var}(\hat{\beta}_2|X) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \equiv \frac{\sigma^2}{n s_{x_2}^2} \leftarrow$  sampling variability  
↑ Sample size (should be  $n-1$  for unbiasedness but wtv)

\* In multivariate regression,  $\text{Var}(\hat{\beta}_j|X) = \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - r_{x_j, x_k}^2)} \equiv \frac{\sigma^2}{n s_{x_j}^2 (1 - r_j^2)}$  sample correlation (deg of collinearity of regressors)  
obtained from regressing  $x_{ij}$  on all other regressors  
 $s_{x_j}^2 = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$  regressors cannot be highly collinear

We can ↑ precision by

↑ sample size,  
↑ variability of regressors,  
↓ variance of errors ( $\sigma^2$ ),  
↓ sample correlation  $r$

**ESTIMATOR FOR  $\text{VAR}(\hat{\beta})$**  The unbiased estimator for  $\sigma^2$ ,  $S^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{\hat{\varepsilon}^2}{n-k}$ ,  $E(S^2) = \sigma^2$

$$\widehat{\text{Var}}(\hat{\beta}) = S^2 (X'X)^{-1}, \quad \text{SE}(\hat{\beta}_j) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_j)} \leftarrow j^{\text{th}} \text{ diagonal element of } \widehat{\text{Var}}(\hat{\beta})$$

proof:  $\hat{\varepsilon} = Me$

# $\hat{\beta}$ IS BLUE

finite sample property

## BEST LINEAR UNBIASED ESTIMATOR (BLUE)

i.e.  $\hat{\beta}$  has lower variance than any other estimator  $\tilde{\beta}$  that is unbiased ( $E(\tilde{\beta}) = \beta$ )

Proof

Let an unbiased estimator be  $\tilde{\beta} = Cy = CX\beta + C\varepsilon$ .

Then  $E(\tilde{\beta}) = \beta \Leftrightarrow E(CX\beta + C\varepsilon) = \beta \Leftrightarrow CX\beta + CE(\varepsilon) = \beta \stackrel{E(\varepsilon)=0}{\Leftrightarrow} CX\beta = \beta \Leftrightarrow CX = I$

\*  $C^{-1}$  and  $X^{-1}$  may not exist!!!

Thus,  $\tilde{\beta} = CX\beta + C\varepsilon = I\beta + C\varepsilon = \beta + C\varepsilon$

Then  $\text{Var}(\tilde{\beta}) = \text{Var}(\beta + C\varepsilon) = C\text{Var}(\varepsilon)C' = C\sigma^2 I C' = \sigma^2 CC'$

$$\begin{aligned} \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}) &= \sigma^2(CC' - (XX')^{-1}) = \sigma^2(CC' - I(XX')^{-1}I) = \sigma^2(CC' - CX(XX)^{-1}X'C) = \sigma^2 C(I - X(XX)^{-1}X'C)C' \\ &= \sigma^2 CMM'C' = \sigma^2 \underbrace{D'D}_{\text{positive semidefinite by PSI}} \geq 0 \end{aligned}$$

(PSI)  $= M$  Symmetric idempotent

$$\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta}) \Leftrightarrow \hat{\beta} \text{ is efficient}$$

## COROLLARIES

- ① Any linear transformation of  $\hat{\beta}$  is also the BLUE of that linear transformation of  $\beta$  (BLUE of  $c'\beta$  is  $c'\hat{\beta}$ )
- ② Any particular element of  $\hat{\beta}$ ,  $\hat{\beta}_j$ , estimates  $\beta_j$  at least as efficiently as any other linear unbiased estimator

# $\hat{\beta}$ IS MVUE

minimum variance unbiased estimator

Assume  $\text{normality } \varepsilon | X \sim N(0, \sigma^2 I)$  to obtain sampling distribution for  $\hat{\beta}$  and enable hypothesis testing

Then by A1-5,  $\hat{\beta} | X \sim N(\beta, \sigma^2 (XX)^{-1})$

and  $\hat{\beta}_j | X \sim N(\beta_j, \sigma^2 c_{jj})$ ,  $c_{jj} = [(XX)^{-1}]_{jj}$

$\hat{\beta}$  IS MLE Makes OLS best amongst linear and non-linear unbiased estimators ( $\Rightarrow$  MVUE)

Given normality assumption A5,  $y | X \sim N(X\beta, \sigma^2 I)$ , then by independence, joint distribution of  $y_1, \dots, y_n$  is

$$L(\alpha, \beta, \sigma^2) = f(y_1, \dots, y_n; \beta, \sigma^2) = \prod_{i=1}^n \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2}(y_i - x_i\beta)^2}$$

$$\log L(\alpha, \beta, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i\beta)^2 = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} S(\hat{\beta})$$

↑ maximise  $-S(\hat{\beta}) \Leftrightarrow$  minimise  $S(\hat{\beta})$  (OLS)

To maximise log-likelihood, FOCs:

$$-\frac{\frac{1}{2\hat{\sigma}_{MLE}^2} \left( \frac{\partial S(\hat{\beta})}{\partial \hat{\beta}} \right)}{n} = 0 \leftarrow \text{Same as OLS, so } \hat{\beta} = \hat{\beta}_{MLE}$$

$$-\frac{1}{2\hat{\sigma}_{MLE}^2} + \frac{1}{2\hat{\sigma}_{MLE}^4} S(\hat{\beta}) = 0 \Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{\hat{\epsilon}' \hat{\epsilon}}{n} = \hat{\sigma}_{MLE}^2$$

# GEOMETRIC ASPECTS OF OLS

FOC of OLS is  $X'\hat{\epsilon} = 0$  (so residuals and regressors are orthogonal)

If our model has an intercept,  $X$  will have a column of 1s. Then:

- $\sum_{i=1}^n \hat{\epsilon}_i = 0$  (Least Squares residuals sum to 0)
- $\bar{y} = \hat{y}$
- $\bar{y} = \bar{x}\hat{\beta}$  (Regression "line" passes through the mean)

Residual maker  $M = I_n - X(X'X)^{-1}X'$   $M \times Z$  gives residuals of reg  $Z | X$   
 Projection matrix  $P = X(X'X)^{-1}X'$  Symmetric  
and idempotent

Estimator:  $\hat{\beta} = (X'X)^{-1}X'y$

Fitted values:  $\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y = Py$  •  $P$  projects  $y$  onto the space

Residual:  $\hat{\epsilon} = y - \hat{y} = y - X(X'X)^{-1}X'y = My = M\epsilon$  • Orthogonal to the projection space

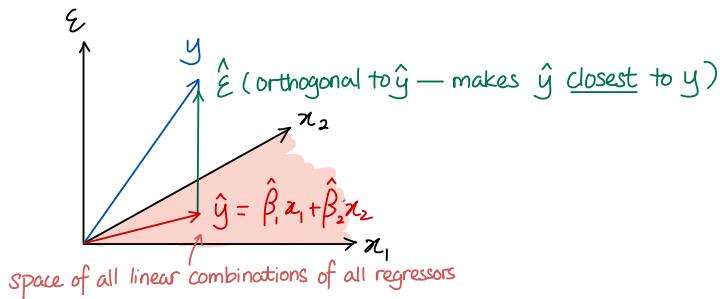
$PX = X$

$MX = 0$  (M orthogonal to X) regressing X on X gives no residuals

$PM = 0$

$\hat{y}'\hat{\epsilon} = 0$  (OLS residuals and fitted values are orthogonal)

$\hat{\epsilon}'\hat{\epsilon} = \epsilon'M\epsilon$   $E(\epsilon'M\epsilon) = E(\text{tr}(\epsilon'M\epsilon))$



OLS is a problem of finding the point in the space spanned by  $x_1, \dots, x_k$  that is "closest" to  $y$ . The solution is  $\hat{y}$ , which minimises  $\|y - \hat{y}\|$

Mathematical Intuition: **Projection Theorem.** Project  $y$  onto subspace of  $X$ .

# FINITE SAMPLE PROPERTIES OF $s^2 = \frac{\hat{\epsilon}'\hat{\epsilon}}{n-k} = \frac{\epsilon'M\epsilon}{n-k}$

**ESTIMATOR FOR  $\sigma^2$**   $s^2 = \frac{RSS}{n-k} = \frac{\hat{\epsilon}'\hat{\epsilon}}{n-k}$

$$\begin{aligned} RSS &= \hat{\epsilon}'\hat{\epsilon} = (My)'(My) = (M(X\beta + \epsilon))'(M(X\beta + \epsilon)) = (MX\beta + M\epsilon)'(MX\beta + M\epsilon) \\ &= (M\epsilon)'(M\epsilon) = \epsilon'M'M\epsilon = \epsilon'MM'\epsilon = \epsilon'M\epsilon \end{aligned}$$

**UNBIASEDNESS**  $E(s^2) = \sigma^2$  by A1-4

$$E(s^2) = E\left(\frac{\epsilon'M\epsilon}{n-k}\right)$$

If  $X$  is fixed

$$\begin{aligned} E(\epsilon'M\epsilon) &= E(\text{tr}(\epsilon'M\epsilon)) = E(\text{tr}(M\epsilon\epsilon')) = \text{tr}(E(M\epsilon\epsilon')) = \text{tr}(ME(\epsilon\epsilon')) = \text{tr}(M\sigma^2 I) \\ &= \sigma^2 \text{tr}(M) = \sigma^2(n-k) \quad \text{PSI Q2} \end{aligned}$$

If  $X$  is stochastic,

$$E(\epsilon'M\epsilon) = E(E(\text{tr}(\epsilon'M\epsilon)|X)) = E(\text{tr}(ME(\epsilon\epsilon'|X))) \stackrel{A4}{=} E(\text{tr}(M\sigma^2 I)) = \sigma^2 n - k$$

**SAMPLING DISTRIBUTION**  $\frac{(n-k)s^2}{\sigma^2} \sim \chi_{n-k}^2$  by A1-5

If  $X$  (and thus  $M$ ) is fixed,  $\frac{(n-k)s^2}{\sigma^2} = \frac{\epsilon'M\epsilon}{\sigma^2} \sim \chi_{n-k}^2$  A5  
 $N(0, \sigma^2 I)$  idempotent;  $n-k$   $\lambda_s$  with  $=1$  (PSI)

If  $X$  (and thus  $M$ ) is stochastic,  $\frac{(n-k)s^2}{\sigma^2} | X \sim \chi_{n-k}^2$

**ZERO COVARIANCE**  $\hat{\beta}$  and  $s^2$  are independent

Prove  $\text{Cov}(\hat{\beta}, \hat{\epsilon}) = 0 \stackrel{+A5}{\implies} \hat{\beta}$  and  $\hat{\epsilon}$  are independent  
 $\Rightarrow \hat{\beta}$  and any function of  $\hat{\epsilon}$  (ind.  $s^2$ ) are independent

$$\begin{aligned} \text{Cov}(\hat{\beta}, \hat{\epsilon}) &= E((\hat{\beta} - E(\hat{\beta}))(\hat{\epsilon} - E(\hat{\epsilon})))' = E((\hat{\beta} - \beta)(M\epsilon - E(M\epsilon)))' = E((X(X')^{-1}X'\epsilon)M')' \\ &\stackrel{\hat{\beta} \text{ unbiased (A1-3)}}{=} E((X(X')^{-1}X'\epsilon)M')' = \sigma^2(X(X')^{-1}X'M')' = \sigma^2(X'X)^{-1}(MX)' = 0 \end{aligned}$$

# REGRESSION ANATOMY

## PARTITION

Let  $y = X\beta + \varepsilon = \begin{pmatrix} X_1 & X_2 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \varepsilon$ ,  $k_1 + k_2 = k$

$$= X_1\beta_1 + X_2\beta_2 + \varepsilon$$

- Maybe we're only interested in a subset of parameters
- Maybe  $k_1 + k_2$  is too large and  $X'X$  too large to invert (e.g. fixed effects panel)

FOC (as usual):  $X'X\hat{\beta} = X'y \Rightarrow \begin{pmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1'y \\ X_2'y \end{pmatrix}$

$$\Rightarrow X_1'X_1\hat{\beta}_1 + X_1'X_2\hat{\beta}_2 = X_1'y$$

$$X_2'X_1\hat{\beta}_1 + X_2'X_2\hat{\beta}_2 = X_2'y$$

PS3  $\Rightarrow \begin{cases} \hat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y, \text{ with } M_2 = I - X_2(X_2'X_2)^{-1}X_2' \\ \hat{\beta}_2 = (X_2'M_1X_2)^{-1}X_2'M_1y, \text{ with } M_1 = I - X_1(X_1'X_1)^{-1}X_1' \end{cases}$  Consistent (2018 3b)

**BIVARIATE EXAMPLE**  $y = \beta_1 + \beta_2 x + \varepsilon$ , let  $X_1 = \begin{pmatrix} 1 & \dots & 1 \end{pmatrix} = \mathbf{1}_{n \times 1}$ ,  $X_2 = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} = x$

Then  $\hat{\beta}_2 = (x'Mx)^{-1}x'My$  where  $M = I_n - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}' = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}'$

$$x'Mx = (Mx)'(Mx) = \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{PS1 Q2}$$

$$x'My = (Mx)'(My) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\hat{\beta}_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{as expected}$$

**ORTHOGONAL REGRESSORS** e.g. MECE dummies (female/male, seasons)

$y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon$ ,  $X_1'X_2 = 0$  ( $X_1 \perp\!\!\!\perp X_2$ )  $\star$  This reg has no constants

$$\hat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y = (X_1'X_1)^{-1}X_1'y$$

- can just ignore the uninteresting set of dummies, regress  $y$  on  $X_1$  only
- But  $X_1'X_2 \neq 0$  in general! Ignoring causes OVB

**RESIDUAL-BASED MODEL**  $y = X_1\beta_1 + X_2\beta_2 + \varepsilon$

$$\Rightarrow M_2y = M_2X_1\beta_1 + M_2\varepsilon \quad (M_2X_2 = 0)$$

$$\Rightarrow y^* = X_1^*\beta_1 + M_2\varepsilon \quad \text{error}$$

$y^* = M_2y$  contains residuals from reg  $y$  on  $X_2$

$X_1^* = M_2X_1$  contains residuals from  $k_1$  regressions of  $X_1$  on  $X_2$ .



Essentially,  $M_2$  finds the transformation to the original regression that removes  $\beta_2$

If  $X_2 = \begin{pmatrix} 1 & \dots & 1 \end{pmatrix}$  (so  $\beta_2$  is just the constant),  $M_2y = y - \bar{y}$  PS)

**FRISCH-WAUGH-LOVELL THEOREM** We can get  $\hat{\beta}_1$  by performing OLS on the residual based model

- Furthermore, residuals from the residual based model  $M_2\hat{\varepsilon} = \hat{\varepsilon}$  (from full regression)! Identical RSS!
- Residual based regression removes the correlation that  $X_1$  and  $y$  have with  $X_2$  before estimating  $\beta_1$
- We can also regress  $y = X_1^*\beta_1 + e$  but residuals incorrect for obtaining  $s^2$

**SE OF RESIDUAL-BASED MODEL**  $\text{Var}(\hat{\beta}_1) = \sigma^2(X_1'M_2X_2)^{-1}$ ,  $\widehat{\text{Var}}(\hat{\beta}_1) = s^2(X_1'M_2X_2)^{-1}$

$$s^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k_1-k_2} \quad . \quad \hat{\varepsilon} = y - X_1\hat{\beta}_1 - X_2\hat{\beta}_2 = y^* - X_1^*\hat{\beta}_1 \quad \text{by FWL}$$

↑ still controlling for  $X_2$   
(implicitly/explicitly)

# DESEASONALISATION

time series/panel  
with seasonal fluctuations

If we have  $y = S\gamma + M_S X + \varepsilon$ , where  $S$  is a complete set of seasonal dummies  
 $X$  are economic variables

mutually orthogonal!

By FWL we can estimate  $S$  by  $M_S y = M_S X S + M_S \varepsilon$   
 $M_S y$  is residuals from regressing  $y$  on seasonal dummies, e.g.

$$y_t = \gamma_1 S_{1t} + \gamma_2 S_{2t} + \gamma_3 S_{3t} + \gamma_4 S_{4t} + \varepsilon_t \text{ where } S_{it} = \begin{cases} 1 & \text{if } t \text{ lies in the } j^{\text{th}} \text{ season} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Then } \hat{\varepsilon}_t = y_t - (\gamma_1 S_{1t} + \gamma_2 S_{2t} + \gamma_3 S_{3t} + \gamma_4 S_{4t})$$

Since  $S_{1t}, S_{2t}, S_{3t}, S_{4t}$  are mutually orthogonal, by FWL, we can estimate  $\gamma_j$  by  $y_t = \gamma_j S_{1t} + v_t$

$$\text{This means that } \hat{\gamma}_j = (S'_1 S_1)^{-1} S'_1 y = \frac{\sum S_{1t} y_t}{\sum S_{1t}^2} = \frac{\sum S_{1t} y_t}{\sum S_{1t}} \quad (\text{since } S_{1t} = 0 \text{ or } 1) \\ = \bar{y}_1, \text{ average of } y \text{ values falling in first season}$$

$$\text{By similar logic, } \hat{\varepsilon}_t = y_t - \bar{y}_1 S_{1t} - \bar{y}_2 S_{2t} - \bar{y}_3 S_{3t} - \bar{y}_4 S_{4t}$$

$$M_S y = \hat{\varepsilon}_t = \begin{cases} y_t - \bar{y}_1 & \text{if } t \text{ is first season} \\ y_t - \bar{y}_2 & \text{if } t \text{ is second season} \\ y_t - \bar{y}_3 & \text{if } t \text{ is third season} \\ y_t - \bar{y}_4 & \text{if } t \text{ is fourth season} \end{cases}$$

This deseasonalises the  $y$  variable by subtracting seasonal averages from  $y_t$  before further regressing on economic variables ( $M_S y = M_S X S + M_S \varepsilon$ )

- Many time series already deseasonalised before access

## PANEL

with unobserved heterogeneity

Individual heterogeneity (fixed effect)

$$y_{it} = \alpha_i + Z'_{it} \beta + \varepsilon_{it} \quad i = 1, \dots, n \quad k+n \text{ parameters} \\ \downarrow \text{too many!} \\ \begin{pmatrix} y_{11} \\ \vdots \\ y_{nT} \end{pmatrix} = \begin{pmatrix} \top & \top & \top & \cdots & \top \\ Z'_{11} & Z'_{12} & Z'_{13} & \cdots & Z'_{1T} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ Z'_{n1} & Z'_{n2} & Z'_{n3} & \cdots & Z'_{nT} \end{pmatrix} \begin{pmatrix} \beta \\ \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} + \varepsilon_{it}, \quad d_{j,it} = \begin{cases} 1 & \forall t \text{ if } i=j \\ 0 & \forall t \text{ otherwise} \end{cases} \\ = (Z D)(\alpha) + \varepsilon_{it} = Z \beta + D \alpha + \varepsilon_{it}$$

We want  $M_D y$ , the residual from reg  $y$  on  $D$ , to estimate  $\beta$  with  $M_D y = M_D Z \beta + M_D \varepsilon_{it}$

Since  $d_{j,it}$  are mutually orthogonal, by FWL we can estimate  $\alpha_j$  by  $y_{it} = d_{j,it} \alpha_j + v_{it}$

$$\text{Then } \alpha_j = (d_{j,it}' d_{j,it})^{-1} d_{j,it}' y = \frac{\sum d_{j,it} y_{it}}{\sum d_{j,it}^2} = \frac{\sum d_{j,it} y_{it}}{\sum d_{j,it}} \quad \text{since } d_{j,it} = 0, 1 \\ = \frac{\sum_{t=1}^T y_{jt}}{T} = \bar{y}_j \quad \text{individual } j \text{'s average over time} \\ \text{since } d_{j,it} y_{it} = \begin{cases} y_{it} & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \hat{v}_{it} = y_{it} - \bar{y}_j \quad \left. \right|_{1 \times k \text{ vectors}} \\ M_D y = \begin{pmatrix} \hat{v}_{11} \\ \vdots \\ \hat{v}_{nT} \end{pmatrix} = \begin{pmatrix} y_{11} - \bar{y}_1 \\ \vdots \\ y_{nT} - \bar{y}_n \end{pmatrix} \quad \text{Similarly, } M_D Z = ((z_{11} - \bar{z}_1) \cdots (z_{1T} - \bar{z}_1), \dots, (z_{n1} - \bar{z}_n) \cdots (z_{nT} - \bar{z}_n))$$

$$y_{it} - \bar{y}_j = (Z_{it} - \bar{z}_j)' \beta + M_D \varepsilon_{it}$$

# OVB

True (Long):  $y = X\beta + Z\delta + \varepsilon$        $E(\hat{\beta}) = \beta$        $\text{Var}(\hat{\beta}) = \sigma^2 I$   
 Misspecified (Short):  $y = X\beta + v$   
 assume  $X, Z$  fixed for simplicity

The short  $\tilde{\beta} = (X'X)^{-1}X'y$   
 $= (X'X)^{-1}X'(X\beta + Z\delta + \varepsilon)$   
 $= \beta + (X'X)^{-1}X'Z\delta + (X'X)^{-1}X'\varepsilon$  is biased!

$E(\tilde{\beta}) \neq \beta$

$$\begin{aligned} E(\tilde{\beta}) &= E(\beta + (X'X)^{-1}X'Z\delta + (X'X)^{-1}X'\varepsilon) \\ &= \beta + (X'X)^{-1}X'Z\delta + (X'X)^{-1}X'E(\varepsilon) \\ &= \beta + (X'X)^{-1}X'Z\delta \end{aligned}$$

bias!! rep indirect effect of  $X$  on  $y$  through its correlation with  $Z$   
 "capturing too much"

**CONDITIONS**

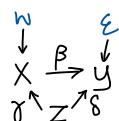
- $X'Z \neq 0$  ( $X$  affects  $Z$ )
- $\delta \neq 0$  ( $Z$  is relevant)

**SE IS BIASED**  $\text{Var}(\tilde{\beta}) = \sigma^2(X'X)^{-1}$  which is "correct", but how we estimate  $\sigma^2$  is wrong!

$\tilde{s}^2 = \frac{\hat{v}'\hat{v}}{n-k_x}$ ,  $\hat{v} = y - X\tilde{\beta}$  is not unbiased! se/t/F-stats based on  $\tilde{s}^2$  invalid!

$y = X\beta + v \Rightarrow \hat{v} = M_{xy} \Rightarrow \hat{v} = M_x(X\beta + Z\delta + \varepsilon) = M_xZ\delta + M_x\varepsilon$   
 Therefore,  $\hat{v}'\hat{v} = (M_xZ\delta + M_x\varepsilon)'(M_xZ\delta + M_x\varepsilon)$   
 $= \varepsilon'M_x\varepsilon + \delta'ZM_xZ\delta + 2\delta'Z'M_x\varepsilon$   
 $E(\hat{v}'\hat{v}) = E(\varepsilon'M_x\varepsilon) + E(\delta'ZM_xZ\delta) + 2\delta'Z'M_xE(\varepsilon)$  assuming  $X, Z$  fixed  
 $= \sigma^2(n-k_x) + E(\delta'ZM_xZ\delta)$  by A1  
  
 So  $E(\tilde{s}^2) = E\left(\frac{\hat{v}'\hat{v}}{n-k_x}\right) \neq \sigma^2$

**OVB VIOLATES A3**



True:  $y = X\beta + Z\delta + \varepsilon$        $E(\varepsilon) = 0$   
 $\Rightarrow v = Z\delta + \varepsilon$

Short:  $y = X\beta + v$

Auxiliary:  $X = Z\gamma + w$        $E(w) = 0$   
 Assume  $w \perp \varepsilon$ , so  $w'\varepsilon = 0$

If  $E(v|X) = 0$ , by LIE,  $E(Xv) = E(E(Xv|X)) = E(X'E(v|X)) = 0$  BUT

$$\begin{aligned} \text{However, } E(Xv) &= E[(Z\gamma + w)'(Z\delta + \varepsilon)] \\ &= E(\gamma'Z\delta + \gamma'Z\varepsilon + w'Z\delta + w'\varepsilon) \\ &= E(\gamma'Z\delta) + \gamma'Z E(\varepsilon) + E(w')Z\delta + E(w'\varepsilon) \\ &= E(\gamma'Z\delta) \\ &\neq 0 \end{aligned}$$

Thus,  $E(v|X) \neq 0$ . A3 is broken!

# INCLUDING IRRELEVANT VARIABLES

$$\text{True : } y = X\beta + \varepsilon$$

$$\text{Misspecified: } y = X\beta + Z\delta + v$$

Note:  $Z$  is irrelevant to  $y$ :

$\tilde{\beta} = (X'M_z X)^{-1} X'M_z y$  is still unbiased and  $\tilde{s}^2$  is unbiased, so inference is valid  
That said,  $\tilde{\beta}$  not the most efficient (power)

$\tilde{\beta}$  IS UNBIASED

$$\tilde{\beta} = (X'M_z X)^{-1} X'M_z (X\beta + v) = \beta + (X'M_z X)^{-1} X'M_z v$$

$$E(\tilde{\beta}) = E[\beta + (X'M_z X)^{-1} X'M_z v] \underset{X, Z \text{ fixed}}{=} \beta + (X'M_z X)^{-1} X'M_z E(v) = \beta$$

SE IS UNBIASED

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= \text{Var}(\beta + (X'M_z X)^{-1} X'M_z v) \\ &= (X'M_z X)^{-1} X'M_z \text{Var}(v) M_z' X (X'M_z X)^{-1} \underset{X, Z \text{ fixed}}{} \\ &= \sigma^2 (X'M_z X)^{-1} \checkmark \end{aligned}$$

$$S^2 = \frac{\hat{v}' \hat{v}}{N - k_x - k_z} \text{ is unbiased since}$$

$$\begin{aligned} \hat{v} &= M_{xz} y = M_{xz} (X\beta + \varepsilon) \underset{\text{AI}}{=} M_{xz} \varepsilon \\ \Rightarrow E(\hat{v}' \hat{v}) &= E(\varepsilon' M_{xz} \varepsilon) \underset{\substack{\text{usual} \\ \text{process}}}{=} \sigma^2 (N - k_x - k_z) \checkmark \end{aligned}$$

EFFICIENCY LOSS  $\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta}) \quad \sigma^2 (X'M_z X)^{-1} \geq \sigma^2 (X'X)^{-1}$

- Select relevant variables based on economic arguments, not statistical ones

To prove:  $(X'M_z X)^{-1} - (X'X)^{-1}$  is positive semidefinite

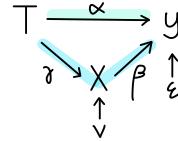
$\Updownarrow$

$X'X - X'M_z X$  is positive semidefinite

$$\begin{aligned} X'X - X'M_z X &= X'X - X'X + X'Z(Z'Z)^{-1}Z'X \\ &= DD' \text{ with } D = (Z'Z)^{-\frac{1}{2}} Z'X \text{ positive semidefinite} \underset{\text{PSI}}{} \end{aligned}$$

# MEDIATION

direct effect:  $\alpha$   
indirect effect:  $\gamma\beta$



Assume  
 $\epsilon \perp T$     $T \perp v$   
 $v \perp \epsilon$     $X \perp \epsilon$

$$\begin{aligned} y &= \delta + T\alpha + X\beta + \epsilon \\ X &= \gamma T + v \end{aligned}$$

GM assumptions are satisfied!  
e.g. A3:  $E(\epsilon|T, X) = 0$

If we exclude  $X$ :  $y = \delta + T\alpha + (\gamma + T\beta + v)\beta + \epsilon$   
 $= \delta + \gamma\beta + T(\alpha + \beta) + v\beta + \epsilon$

This means that the coefficient from reg  $y$  on  $T$  gives the total effect  $\alpha + \gamma\beta$

reg  $y$  on  $T, X$  gives  $\beta$   
 reg  $X$  on  $T$  gives  $\gamma$

- | Observing  $y, T, X$  allows us to differentiate direct/indirect effect
- | Observing  $y, T$  only gives total effect estimate

This means that if you include "controls" that are actually mediators, the resultant decrease in the coefficient on  $T$  just means the direct effect < total effect.  
Be careful about interpretation!

## SELECTING REGRESSORS

**GOODNESS OF FIT**  $R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$   $\leftarrow \hat{\epsilon} \hat{\epsilon}$   
 $\leftarrow \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$   $\leftarrow$  also called SST  $\star$  implies  $\hat{\epsilon} \hat{\epsilon} = \text{RSS} = \text{TSS}(1 - R^2)$

- Not suitable in helping to select regressors since additional regressors always  $\uparrow R^2$

**ADJUSTED R<sup>2</sup>**  $\bar{R}^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \frac{n-1}{n-k}$

- Allows for a tradeoff between fit and parsimony ( $1k \Rightarrow \sqrt{n-k} \Rightarrow \uparrow \bar{R}^2$ )

## SUMMARY

$$y = X\beta + Z\delta + \epsilon$$

True state of the world	Include Z	Don't Include Z
$\delta = 0$	Unbiased but higher variance	BLUE
$\delta \neq 0$	BLUE	Biased but lower variance

Unbiasedness is important but so is variance!

Don't just include all potential regressors! Think about the tradeoff

# HYPOTHESIS TESTING

under A1-5  
assume  $X$  fixed for simplicity

## SETUP

From A1-5, the exact sampling distribution of

$$\begin{aligned}\hat{\beta} &\sim N(\beta, \sigma^2(X'X)^{-1}) \\ \hat{\beta}_i &\sim N(\beta_i, \sigma^2 c_{ii}), \quad c_{ii} = [(X'X)^{-1}]_{ii} \\ (n-k)\sigma^{-2}s^2 &\sim \chi^2_{n-k}\end{aligned}$$

$\hat{\beta}$  and  $s^2$  are independent

## T-TEST (SINGLE PARAMETER)

If  $\sigma^2$  is known, we use  $Z = \frac{\hat{\beta}_i - \mu}{\sigma/\sqrt{c_{ii}}} \sim N(0, 1)$  under  $H_0$

Reject  $H_0$  (two-sided) if  $\left| \frac{\hat{\beta}_i - \mu}{\sigma/\sqrt{c_{ii}}} \right| > Z_{\frac{\alpha}{2}} \stackrel{\alpha=5\%}{=} 1.96$

If  $\sigma^2$  is unknown, we use  $t = \frac{\hat{\beta}_i - \mu}{s/\sqrt{c_{ii}}} = \frac{\hat{\beta}_i - \mu}{SE(\hat{\beta}_i)} \sim t_{n-k}$  under  $H_0$

Reject  $H_0$  (two-sided) if  $\left| \frac{\hat{\beta}_i - \mu}{SE(\hat{\beta}_i)} \right| > t_{n-k, \frac{\alpha}{2}}$

$t$  is derived from  $\frac{\hat{\beta}_i - \mu}{s/\sqrt{c_{ii}}} = \frac{\left( \frac{\hat{\beta}_i - \mu}{\sigma/\sqrt{c_{ii}}} \right) \leftarrow N(0, 1)}{\frac{(n-k)\sigma^{-2}s^2}{n-k} \downarrow \chi^2_{n-k}} \sim t_{n-k}$  under  $H_0$  as  $\hat{\beta}$  and  $s^2$  are independent

For small  $n$ , critical values are typically larger when  $\sigma^2$  is unknown

- t-distribution has fatter tails than  $N$

- "Reject the null later"

As  $n \rightarrow \infty$ ,  $t \rightarrow z$  so critical values are identical whether  $\sigma^2$  is known or not

- If you get a low  $t$  with high  $se(\hat{\beta})$  but  $\hat{\beta}$  far from 0, it's noisy/imprecise data (can't conclude)
- If you get a low  $t$  with low  $se(\hat{\beta})$  but  $\hat{\beta}$  close to 0, you're "not rejecting precisely" (likely no effect)

**P-VALUE**  $\Pr(|T| \geq \frac{\hat{\beta}_i - \mu}{SE(\hat{\beta}_i)})$

Lowest level of significance at which you want to reject  $H_0$

**CONFIDENCE INTERVAL**  $100(1-\alpha)\%$  CI for  $\beta_i$  ( $\sigma^2$  unknown)

$[\hat{\beta}_i - t_{n-k, \frac{\alpha}{2}} SE(\hat{\beta}_i), \hat{\beta}_i + t_{n-k, \frac{\alpha}{2}} SE(\hat{\beta}_i)]$  derive in PS4 Q3

such that  $\Pr(\hat{\beta}_i - t_{n-k, \frac{\alpha}{2}} SE(\hat{\beta}_i) < \text{true value } \beta_i < \hat{\beta}_i + t_{n-k, \frac{\alpha}{2}} SE(\hat{\beta}_i)) = \Pr(-t_{n-k, \frac{\alpha}{2}} < \frac{\hat{\beta}_i - \beta_i}{SE(\hat{\beta}_i)} < t_{n-k, \frac{\alpha}{2}}) = 1 - \alpha$

Reject  $H_0: \beta_i = \mu$  with sig. level  $\alpha$  with t-test  $\xrightarrow{\text{PS4Q3}} \mu \notin 100(1-\alpha)\%$  CI

**T-TEST (LINEAR RESTRICTION)**  $H_0: \beta_1 + \beta_2 + \dots + \beta_k = c$  or equivalently  $H_0: \beta = c$

$r'\hat{\beta} \sim N(c, \sigma^2 r'(X'X)^{-1}r)$  under  $H_0$

because  $\hat{\beta} \sim N(\beta + \text{Var}(\hat{\beta}))$   
 $\Rightarrow r'\hat{\beta} \sim N(r'\beta + c, r'\text{Var}(\hat{\beta})r)$  under  $H_0$

If  $\sigma^2$  is known, we use  $Z = \frac{r'\hat{\beta} - c}{\sigma \sqrt{r'(X'X)^{-1}r}} = \frac{r'\hat{\beta} - c}{SD(r'\hat{\beta} - c)} \sim N(0, 1)$  under  $H_0$

If  $\sigma^2$  is unknown, we use  $t = \frac{r'\hat{\beta} - c}{S\sqrt{r'(X'X)^{-1}r}} = \frac{r'\hat{\beta} - c}{SE(r'\hat{\beta} - c)} \sim t_{n-k}$  under  $H_0$

★  $SE(\hat{\beta}_1 + \hat{\beta}_2) \neq SE(\hat{\beta}_1) + SE(\hat{\beta}_2)$

Just reparameterise: rewrite model with  $\gamma = \beta_1 + \beta_2 - c$ , test  $H_0: \gamma = 0$  PS5 Q1b  
 Allows us to just use the "standard OLS package"

## STEPS

- ① Define  $H_0, H_1$
- ② Calculate test statistic and give its distribution under  $H_0$
- ③ Give sig. level and critical values
- ④ State rejection rule and interpret findings
- ⑤ Clearly indicate assumptions for test to be valid (A1-5)

## IMPORTANT

Single hypothesis tests need to be followed by joint hypothesis test

If near multicollinear, marginal contribution of each regressors may be small (low t) when added last but they may be jointly significant (high F/W)

## POWER

$$\begin{aligned} \text{Power} &= 1 - \Pr(\text{Type II error}) \\ &= \Pr(H_0 \text{ rejected} \mid H_0 \text{ true}) \end{aligned}$$

e.g.  $H_0: \beta = \beta_1$  vs  $H_1: \beta \neq \beta_1$

$$\begin{aligned} \text{Power of this test when } \beta = \beta_1 \text{ is } &\Pr\left(\left|\frac{\hat{\beta}}{\sigma}\right| \geq Z_{\alpha/2} \mid \beta = \beta_1\right) \\ &= \Pr\left(\frac{\hat{\beta}}{\sigma} \leq -Z_{\alpha/2} \text{ or } \frac{\hat{\beta}}{\sigma} \geq Z_{\alpha/2} \mid \beta = \beta_1\right) \\ &= \Pr\left(\frac{\hat{\beta} - \beta_1}{\sigma} \leq -Z_{\alpha/2} - \frac{\beta_1}{\sigma}\right) + \Pr\left(\frac{\hat{\beta} - \beta_1}{\sigma} \geq Z_{\alpha/2} - \frac{\beta_1}{\sigma}\right) \end{aligned}$$

$\uparrow \beta_1 \Rightarrow \uparrow \text{Power} \Rightarrow \downarrow \Pr(\text{Type II error})$

$\uparrow \alpha \Rightarrow \uparrow \text{Power} \Rightarrow \downarrow \Pr(\text{Type II error})$   
 $= \uparrow \Pr(\text{Type I error})$

To ensure test is powerful,

# HYPOTHESIS TESTING - MULTIPLE LIN RESTRICTIONS

**F-TEST**  $H_0: R\beta = C$  R is a  $J \times k$  full rank matrix,  $J \leq k < n$

e.g. let  $k=4$ ,  $\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix}$

$R = (0 1 -1 0)$   $C = (0)$  then  $H_0: \beta_2 = \beta_3$   
 ↪ single linear restriction

$R = (0 1 0 0)$   $C = (2)$  then  $H_0: \beta_2 = 2$  and  $\beta_3 = 1$

$R = (0 0 1 0)$   $C = (0)$  then  $H_0: \beta_2 = \beta_3 = \beta_4 = 0$

$R = (0 1 1 0)$   $C = (0)$  then  $H_0: \beta_2 = \beta_3 = 2\beta_4 = 0$

NOT  $(0 1 -1 0)$  which doesn't have full row rank

## WALD TEST

One of two hypothesis testing principles (the other is RRSS/URSS)

Do the estimates come reasonably close to satisfying the restrictions implied by the hypothesis?

$H_0: R\beta = C$  vs  $H_1: R\beta \neq C$

We check if discrepancy of  $R(\hat{\beta} - C)$  from 0 is statistically significant or just due to sampling error

**DISCREPANCY VECTOR**  $d = R\hat{\beta} - C$ .

Typically use the quadratic form  $d'[\text{Var}(d)]^{-1}d$  that is a scalar measure of the size of vector

If  $\sigma^2$  is known,

Under  $H_0$ ,  $d \sim N(0, \text{Var}(d))$ ,  $\text{Var}(d) = \sigma^2 R(X'X)^{-1}R'$   
 ↪ must have full rr  
 for Var to be non-singular

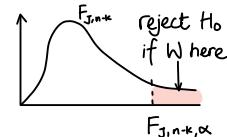
$$\frac{(R\hat{\beta} - C)[R(X'X)^{-1}R]^{-1}(R\hat{\beta} - C)}{\sigma^2} = d'[\text{Var}(d)]^{-1}d \sim \chi^2_J$$

PS2 Q3

If  $\sigma^2$  is unknown,

Under  $H_0$ ,  $W = \frac{(R\hat{\beta} - C)[R(X'X)^{-1}R]^{-1}(R\hat{\beta} - C)}{S^2 J} = \frac{d'[\widehat{\text{Var}}(d)]^{-1}d}{J} \sim F_{J, n-k}$

Then we reject  $H_0$  if  $W > F_{J, n-k, \alpha}$   
 ↪ not  $\frac{\alpha}{2}$ !! Squaring d makes one-sided tests impossible anyway



**WHY KEEP T-TEST?** When there is only a single linear restriction such that  $d$  is scalar,

Under  $H_0$ ,  $d'[\text{Var}(d)]^{-1}d = \left(\frac{d}{\text{SE}(d)}\right)^2 \sim \chi^2_1$  if  $\sigma^2$  is known

$d'[\widehat{\text{Var}}(d)]^{-1}d = \left(\frac{d}{\text{SE}(d)}\right)^2 \sim t_{n-k}^2 = F_{1, n-k}$

By taking squares, we hide the direction of our violations and makes one-sided tests impossible.

# LOSS OF FIT (COMPARE RRSS / URSS)

$$H_0: R\beta = c \quad \text{vs} \quad H_1: R\beta \neq c$$

Does imposing the restrictions in  $H_0$  lead to a significant loss in fit?

① Regress  $y$  on  $X$  to get the **unconstrained least squares estimator**  $\hat{\beta}$ .

$$\text{② Calculate } URSS = \hat{\epsilon}'\hat{\epsilon} = (y - X\hat{\beta})'(y - X\hat{\beta})$$

③ Regress  $y$  on  $X$  subject to the constraint  $R\beta = c$  (minimise  $(y - X\beta)'(y - X\beta)$  st  $R\beta = c$ ) to get **constrained least squares estimator**  $\beta^*$

$$\text{④ Calculate } RRSS = \epsilon^*\epsilon^* = (y - X\beta^*)'(y - X\beta^*)$$

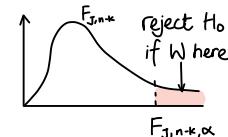
$$\text{⑤ } F = \frac{RRSS - URSS}{URSS} \left( \frac{n-k}{J} \right) \sim F_{J, n-k}$$

under  $H_0$ . also written as  $\left( \frac{R^2_{\text{unconstr}} - R^2_{\text{constr}}}{1 - R^2_{\text{constr}}} \right) \left( \frac{n-k}{J} \right) \sim F_{J, n-k}$

- $F$  is a ratio measuring % increase in residual variance (loss in fit) due to imposing  $H_0$ .

⑥ Find the **critical value**  $K^*$  st  $\alpha\%$  of the area under an  $F_{J, n-k}$  distribution lies to the right of  $K^*$

⑦ Reject  $H_0$  if  $F > K^*$



**SIMPLE EXAMPLE - "NOTHING IS HAPPENING HERE"** except for the constant

$$y = \beta_0 + \beta_1 x_{1,1} + \dots + \beta_k x_{k,1} + \epsilon_i \quad H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0 \quad \text{vs} \quad H_1: (\beta_1, \dots, \beta_k) \neq (0, \dots, 0)$$

$k-1$  restrictions

$$URSS = \sum_{i=1}^n (y_i - x_i \hat{\beta})^2 = RSS$$

$$\text{Restricted model is } y_i = \beta_0 + \epsilon_i \Rightarrow \beta_0^* = \bar{y} \Rightarrow RRSS = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \bar{y})^2 = TSS$$

$$F = \left( \frac{TSS - RSS}{RSS} \right) \left( \frac{n-k}{k-1} \right) = \frac{R^2}{1-R^2} \left( \frac{n-k}{k-1} \right) \sim F_{k-1, n-k} \text{ under } H_0 \quad R^2 = 1 - \frac{RSS}{TSS}$$

- More likely to reject  $H_0$  when  $R^2$  is large! Regressors more relevant!

**FINDING RRSS** Find  $\beta$  which minimises  $(y - X\beta)'(y - X\beta)$  subject to  $R\beta = c$

Other than plugging in the restrictions (see above), can also do Lagrangian

$$L(b, \lambda) = (y - Xb)'(y - Xb) - 2\lambda'(Rb - c)$$

just helps simplify (actual value of  $\lambda$  unimportant)

$$\begin{aligned} \text{FOCs: } & \begin{cases} -2X'y + 2X'X\beta^* - 2R'\lambda^* = 0 \\ -2(R\beta^* - c) = 0 \end{cases} \\ & \Rightarrow \beta^* = (X'X)^{-1}X'y + (X'X)^{-1}R'\lambda^* \\ & R\beta^* = R\hat{\beta} + R(X'X)^{-1}R'\lambda^* \\ & c = R\hat{\beta} + R(X'X)^{-1}R'\lambda^* \end{aligned}$$

$$\begin{aligned} \lambda^* &= (R(X'X)^{-1}R')^{-1}(c - R\hat{\beta}) && \text{if } R\hat{\beta} = c, \lambda^* = 0 \\ \beta^* &= \hat{\beta} - (X'X)^{-1}R'(R(X'X)^{-1}R')^{-1}(c - R\hat{\beta}) && \text{if } R\hat{\beta} = c, \beta^* = \hat{\beta} \end{aligned}$$

$\lambda^*$  is "shadow prices" of imposing the restrictions — high is bad

# LOSS OF FIT = WALD TEST

$$F = W \equiv \frac{d'(\widehat{\text{Var}}(d))^{-1}d}{J} \sim F_{J, n-k} \text{ under } H_0$$

$$\begin{aligned} d &= R\hat{\beta} - c \\ \text{Var}(d) &= \sigma^2 R(X'X)^{-1} R' \\ \widehat{\text{Var}}(d) &= S^2 R(X'X)^{-1} R' \end{aligned}$$

Allows us to conclude  $F = \frac{\text{RRSS} - \text{URSS}}{\text{URSS}} \left( \frac{n-k}{J} \right) \sim F_{J, n-k} \text{ under } H_0$

## PROOF

$$\begin{aligned} \frac{\text{RRSS} - \text{URSS}}{\text{URSS}} \left( \frac{n-k}{J} \right) &= \frac{\text{RRSS} - \text{URSS}}{S^2 J} \\ &= \frac{(y - X\beta^*)'(y - X\beta^*) - (y - X\hat{\beta})'(y - X\hat{\beta})}{S^2 J} \\ &= \frac{(y - X\hat{\beta} + D)'(y - X\hat{\beta} + D) - (y - X\hat{\beta})'(y - X\hat{\beta})}{S^2 J}, \quad D = X(X'X)^{-1}R'(R(X'X)^{-1}R')^{-1}(c - R\hat{\beta}) \\ &= \frac{(y - X\hat{\beta})'(y - X\hat{\beta}) + D'D + (y - X\hat{\beta})'D + D'(y - X\hat{\beta}) - (y - X\hat{\beta})'(y - X\hat{\beta})}{S^2 J} \\ &\text{since } (y - X\hat{\beta})'D = \hat{\epsilon}'X(X'X)^{-1}R'(R(X'X)^{-1}R')^{-1}(c - R\hat{\beta}) = 0 \\ &= \frac{D'D}{S^2 J} \\ &= \frac{d'[\widehat{\text{Var}}(d)]^{-1}d}{J} \\ &= W \end{aligned}$$

Note: Wald test is preferred to comparing  $\frac{\text{RRSS}}{\text{URSS}}$  (or equivalently  $\frac{R_r^2}{R_u^2}$ ) when there is **heteroscedasticity** (can use RSE for  $\widehat{\text{Var}}(d)$ )

$$t^2 = F$$

$$F = \frac{\text{RRSS} - \text{URSS}}{\text{URSS}} \left( \frac{n-k}{J} \right) = \frac{d'(\widehat{\text{Var}}(d))^{-1}d}{J} \stackrel{\text{from above}}{=} \left( \frac{d}{\text{SE}(d)} \right)^2 \stackrel{\text{with only 1 restriction } (J=1)}{=} t^2$$

# ASYMPTOTIC THEORY

In the absence of a finite (exact) sampling distribution, we use estimators that are consistent (convergence in probability) and has a limiting distribution (convergence in distribution)

asymptotically unbiased

enables hypothesis testing

"density of estimator converges to a spike at true value"

**CONSISTENT**  $\hat{\theta}$  is consistent if, as  $n \rightarrow \infty$ ,  $\text{plim } \hat{\theta} = \theta$  or  $\hat{\theta} \xrightarrow{P} \theta$

Formal definition:

Let  $X_n$  be a random variable indexed by the size of a sample

$X_n$  converges in probability to  $X$  ( $X_n \xrightarrow{P} X$ ) if  $\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0 \quad \forall \epsilon > 0$

## PROVING CONSISTENCY

### SUFFICIENT CONDITIONS

①  $\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta$  (Asymptotic unbiasedness)

②  $\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}) = 0$  requires  $\text{Var}(\hat{\theta})$  to be finite. Stronger than needed

① + ②  $\Rightarrow$  Chebychev's Inequality  $\Rightarrow$  Consistency

$$P(|\hat{\theta} - \theta| > \epsilon) \leq \frac{E[(\hat{\theta} - \theta)^2]}{\epsilon^2} = \frac{\text{MSE}(\hat{\theta})}{\epsilon^2} \rightarrow 0$$

- Unbiased  $\Rightarrow$  asymptotically unbiased (common sense)

## KHINCHINE'S WEAK LAW OF LARGE NUMBERS

If  $X_1, \dots, X_n$  are random samples (i.i.d.) with finite mean  $\mu$ .

$$\text{plim}(\bar{X}) = \text{plim} \frac{1}{n} \sum_{i=1}^n X_i = E(X_i) = \mu$$

no need to assume  $\text{Var}(X_i)$  is finite, or derive  $\text{Var}(\bar{X})$

- Sample averages converge to their population counterparts
- Can use to prove others,  $\text{plim}_{\substack{X_i \text{ iid} \\ X_i \text{ iid}}} \frac{1}{n} \sum X_i^2 = E(X_i^2)$ ,  $\text{plim}_{\substack{X_i \text{ iid} \\ X_i \text{ iid}}} \frac{1}{n} \sum X_i \epsilon_i = E(X_i \epsilon_i) = 0$ ! with Slutsky theorem
- Alt. LLNs exist for  $X_i, \epsilon_i$  not iid
- Approach: Write equations of estimators as summation, multiply & divide by  $n$  and apply WLLN

## SLUTSKY THEOREM

$$\text{plim } g(Y_n) = g(\text{plim}(Y_n))$$

$\forall$  continuous  $g(\cdot)$

- addition
- multiplication
- matrix multiplication
- inverse
- into elements of a matrix

## OTHER TOOLS

$$\text{plim}(c) = c \quad \bullet \text{ If something isn't a function of } n, \text{ it is equal to its plim}$$

$$\text{plim}_{n \rightarrow \infty} (X_n) = \lim_{n \rightarrow \infty} (X_n) \text{ if } X_n \text{ is not random} \quad \bullet \text{ e.g. } \text{plim}_{n \rightarrow \infty} \frac{n}{n-k} = \lim_{n \rightarrow \infty} \frac{n}{n-k} \text{ if } k \text{ is not changing}$$

# PROOFS OF CONSISTENCY

using Slutsky, WLLN

## SAMPLE MEAN IS CONSISTENT

If  $X_1, \dots, X_n$  is a random sample (i.i.d.),  $X_i \sim (\mu, \sigma^2)$ .

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  is a consistent estimator of  $\mu$

sufficient conditions

Method ①  $E(\bar{X}) = \mu$ ,  $\text{Var}(\bar{X}) = \frac{\sigma^2}{n} \Rightarrow \lim_{n \rightarrow \infty} E(\bar{X}) = \mu$ ,  $\lim_{n \rightarrow \infty} \text{Var}(\bar{X}) = 0 \Rightarrow \text{plim} \left( \frac{1}{n} \sum_{i=1}^n X_i \right) = \mu$  assumes  $\sigma^2$  is finite

Method ②  $E(X_i)$  is finite, so apply WLLN directly:  $\text{plim } \bar{X} = E(X_i) = \mu$ . Consistent! No need  $\text{Var}(X_i) = \sigma^2 < \infty$   
No need to derive  $\text{Var}(\bar{X})$

## $\bar{X}^2$ IS CONSISTENT

$X_1, \dots, X_n$  is a random sample (i.i.d.) with finite mean  $\mu$

sufficient conditions

Method ①  $\text{Var}(\bar{X}^2) = E(\bar{X}^2) - (E(\bar{X}))^2 \Rightarrow E(\bar{X}^2) = \frac{\sigma^2}{n} + \mu^2 \xrightarrow{n \rightarrow \infty} \mu^2$  biased but asymptotically unbiased

$\text{Var}(\bar{X}^2) = \dots$ ? too hard.

Method ②  $\xrightarrow{\text{WLLN, Slutsky}} \text{plim } \bar{X} = \mu \xrightarrow{\text{Slutsky}} \text{plim } \bar{X} = \mu^2$  ✓

## $\widehat{\text{Cov}}(x_i, \varepsilon_i)$ IS CONSISTENT

$\text{plim } \widehat{\text{Cov}}(x_i, \varepsilon_i) = \text{plim} \left( \frac{1}{n} \sum_{i=1}^n x_i \varepsilon_i - \bar{x} \bar{\varepsilon} \right) \xrightarrow{\text{Slutsky}} \text{plim } \bar{x}_i \bar{\varepsilon}_i - \text{plim } \bar{x} \text{ plim } \bar{\varepsilon} = E(x_i \varepsilon_i) - E(x_i)E(\varepsilon_i) = \text{Cov}(x_i, \varepsilon_i)$

## $\hat{\beta}$ IS CONSISTENT - BIVARIATE

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

sufficient conditions

Method ①  $E(\hat{\beta}) = \beta$  and  $\text{Var}(\hat{\beta}|X) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} \xrightarrow{n \rightarrow \infty} 0$  as  $n \rightarrow \infty$  excitation condition ✓

Method ②  $\text{plim}(\hat{\beta}) = \text{plim}(\beta + \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}) = \beta + \frac{\text{plim} \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})}{\text{plim} \frac{1}{n} \sum (x_i - \bar{x})^2} = \beta + \frac{\text{plim} \widehat{\text{Cov}}(x_i, \varepsilon_i)}{\text{plim} \text{Var}(x_i)} \xrightarrow{\text{WLLN}} \beta + \frac{\text{Cov}(x_i, \varepsilon_i)}{\text{Var}(x_i)} = \beta$  assume  $(x_i, \varepsilon_i)$  iid

## $\hat{\beta}$ IS CONSISTENT - MULTIVARIATE

$$\hat{\beta} = (X'X)^{-1} X' y = \beta + (X'X)^{-1} X' \varepsilon = \beta + \left( \frac{X'X}{n} \right)^{-1} \frac{X' \varepsilon}{n} \quad \text{write in terms of averages}$$

$$\text{In bivariate case, } \frac{X'X}{n} = \begin{pmatrix} 1 & \sum_{i=1}^n \frac{x_i}{n} \\ \sum_{i=1}^n \frac{x_i}{n} & \sum_{i=1}^n \frac{x_i^2}{n} \end{pmatrix} \quad \text{PSI Q3}$$

\*  $E(x_i^2)$  must exist  
( $E(x_i)$  exists by existence of moments)

$$\text{plim} \left( \frac{X'X}{n} \right) \xrightarrow{\text{Slutsky}} \begin{pmatrix} 1 & \text{plim} \frac{1}{n} \sum x_i \\ \text{plim} \frac{1}{n} \sum x_i & \text{plim} \frac{1}{n} \sum x_i^2 \end{pmatrix} \xrightarrow{\text{WLLN}} \begin{pmatrix} 1 & E(x_i) \\ E(x_i) & E(x_i^2) \end{pmatrix} = D \quad (\text{assume } D^{-1} \text{ exists})$$

$$\text{plim} \left( \frac{X' \varepsilon}{n} \right) = \text{plim} \frac{1}{n} \sum_{i=1}^n x_i \varepsilon_i \xrightarrow{\text{WLLN}} E(x_i \varepsilon_i) \stackrel{\text{def}}{=} E[E(x_i \varepsilon_i | X)] = E[x_i E(\varepsilon_i | X)] \stackrel{\text{A3}}{=} 0$$

$$\text{plim } \hat{\beta} = \underset{\text{Slutsky}}{\beta} + \left( \text{plim} \frac{X'X}{n} \right)^{-1} \text{plim} \frac{X' \varepsilon}{n} = \beta$$

$$S^2 \text{ IS CONSISTENT} \quad S^2 = \frac{\hat{\varepsilon}' \hat{\varepsilon}}{n-k} = \frac{\varepsilon' M \varepsilon}{n-k} = \frac{n}{n-k} \frac{\varepsilon' M \varepsilon}{n} = \frac{n}{n-k} \left[ \frac{\varepsilon' \varepsilon}{n} - \frac{\varepsilon' X}{n} \left( \frac{X' X}{n} \right)^{-1} \frac{X' \varepsilon}{n} \right] \xrightarrow{\text{WLLN}} \sigma^2 \quad \text{PS5 extra}$$

For proving biased/inconsistent due to heteroscedasticity, see 2018 2c

# ASYMPTOTIC DISTRIBUTION

The asymptotic property of **convergence in distribution** enables hypothesis testing of an estimator without knowing its finite sample distribution — gives us a distribution we can use to approximate the estimator's distribution arbitrarily well for sufficiently large sample.

Formal definition: The sequence  $Z_n$  with CDF  $F_{Z_n}(z)$  converges in distribution to a random variable  $Z$  with CDF  $F_Z(z)$  if  $\lim_{n \rightarrow \infty} |F_{Z_n}(z) - F_Z(z)| = 0$  at all points of continuity of  $F_Z(z)$

Intuitively: Distribution of  $Z_n \rightarrow$  Distribution of  $Z$  as  $n \rightarrow \infty$

Example:  $t_{n-k} \xrightarrow{n \rightarrow \infty} N(0,1)$

**CENTRAL LIMIT THEOREM** If  $X_1, \dots, X_n$  is a random sample (i.i.d) with finite mean  $\mu$  and finite variance  $\sigma^2$ ,

$$\text{Standardised } \sqrt{n}(\bar{X} - \mu) \equiv \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - E(X_i)) \xrightarrow{d} N(0, \sigma^2)$$

" $\sqrt{n}(\bar{X} - \mu)$  has a  $N(0, \sigma^2)$  limiting distribution"  
 $\bar{X}$  does not depend on sample size  $n$

- If  $X_i$  normally distributed,  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$  exact
- If  $X_i$  unknown distribution,  $\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2)$  by CLT  
 • so  $\bar{X} \stackrel{\sim}{\sim} N(\mu, \frac{\sigma^2}{n})$  approximately

## OLS ESTIMATOR

Even if we don't have A5 ( $\varepsilon_i \sim N(0, \sigma^2)$ ), we can still rely on A1-4 and CLT to conclude  $\hat{\beta} | X \stackrel{\sim}{\sim} N(\beta, \sigma^2(X'X)^{-1})$  if  $\varepsilon_i$  is iid,  $E(\varepsilon_i) = 0$ ,  $\text{Var}(\varepsilon_i) = \sigma^2$  finite

## APPLICATION TO H-TEST

### SINGLE LINEAR RESTRICTION

$$H_0: c'\beta = \gamma$$

Under finite sample test,  $t = \frac{(c'\hat{\beta} - \gamma)}{\text{SE}(c'\hat{\beta} - \gamma)} \sim t_{n-k}$  by A1-5 under  $H_0$

Under asymptotic test,  $Z = \frac{(c'\hat{\beta} - \gamma)}{\text{SE}(c'\hat{\beta} - \gamma)} \stackrel{\sim}{\sim} N(0, 1)$  by A1-4 and CLT under  $H_0$

- Doesn't matter if we know  $\sigma^2$  or not —  $S^2$  is a consistent estimator of  $\sigma^2$

### MULTIPLE LINEAR RESTRICTION

$$H_0: R\beta = c \quad d = R\hat{\beta} - c$$

Under finite sample test,  $\frac{d'(\widehat{\text{Var}}(d))^{-1}d}{J} \sim F_{J, n-k}$  by A1-5 under  $H_0$   
 $d \sim N(0, \sigma^2 R(X'X)^{-1} R')$

Under asymptotic test,  $d'(\widehat{\text{Var}}(d))^{-1}d \stackrel{\sim}{\sim} \chi_J^2$  by A1-4 and CLT under  $H_0$   
 $d \stackrel{\sim}{\sim} N(0, \sigma^2 R(X'X)^{-1} R')$

- Doesn't matter if we know  $\sigma^2$  or not —  $S^2$  is a consistent estimator of  $\sigma^2$
- Asymptotically, no need to account for imprecision of  $S^2$  with F-dist
- t and F test are **exact tests** that rely on the assumption of normality
- Z and  $\chi^2$ -test don't — use CLT

# TIME-SERIES DATA

alt. title: why A3 doesn't work

$$\{(y_t, x_{t1}, \dots, x_{tk}) : t=1, \dots, T\} \text{ where } T \text{ is sample size}$$

OLS estimator biased! Need large samples

Two assumptions needed for LLN, CLT and standard statistical inference to work

- Stationarity
- Weak dependence

## MODELS USING TS DATA

$$\boxed{\text{STATIC}} \quad y_t = \beta_0 + \beta_1 x_t + \varepsilon_t, \quad t=1, \dots, T \quad \text{e.g. Philips curve}$$

- Models a contemporaneous relationship between  $y$  and  $x$  ( $y_t$  doesn't depend on  $y_{t-1}$ )
- Assumes a change in  $x_t$  at time  $t$  has immediate effect on  $y_t$
- If A1-4, OLS is BLUE
  - But there is often autocorrelation -  $\text{Cov}(\varepsilon_t, \varepsilon_s) \neq 0$  for some  $t \neq s$
  - But as long as A3  $E(\varepsilon|X)=0$ , OLS still unbiased, and we can use RSE

$$\boxed{\text{FINITE DISTRIBUTED LAG MODELS}} \quad y_t = \alpha + \gamma_t \gamma_1 + \gamma_{t-1} \gamma_2 + \varepsilon_t \quad \text{e.g. min W and UNt}$$

$\gamma_1, \gamma_2$  are distributed lag coefficients

- $\gamma_1$  is the contemporaneous effect - impact propensity
- $\gamma_1 + \gamma_2$  is the long run propensity - when  $x$  is changed permanently to a fixed value

$$X = \begin{pmatrix} 1 & T & T & T \\ 1 & x_t & x_{t-1} & x_{t-2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \text{income} & \text{income} & \text{income} \\ 2020 & 2019 & 2018 \end{pmatrix} \quad \text{If } x_t \text{ changes slowly over time (} x_t \text{ highly correlated to } x_{t-1}\text{), near multicollinearity may make estimates imprecise, but we can still estimate LR effect precisely}$$

## WEAKER EXOGENITY

But for OLS to be unbiased, must assume strict exogeneity  $E(\varepsilon_t | X) = E(\varepsilon_t | x_1, \dots, x_T) = 0$

- $\varepsilon_t$  uncorrelated with all past, current and future  $x$ -values
- In reality  $x_{t+1}$  may respond to  $\varepsilon_t$  e.g. min W responding to past employment shocks. PS6 Q1
- More reasonable to assume  $E(\varepsilon_t | x_t, x_{t-1}, \dots) = 0$
- Makes OLS of FDL/static models biased, but consistent under stationarity and weak dependence

$$\text{Aim: } \text{plim } \hat{\beta} = \beta + [\text{plim}(\frac{XX}{T})]^{-1} \text{plim}(\frac{X\varepsilon}{T}) = \beta$$

Assume  $\text{plim}(\frac{XX}{T}) = D$  and  $D^{-1}$  exists. PS6

$$\text{plim} \frac{X'\varepsilon}{T} = \begin{pmatrix} \text{plim} \frac{\frac{T}{n} \varepsilon_t}{T} \\ \text{plim} \frac{\frac{T}{n} \varepsilon_t x_t}{T} \\ \text{plim} \frac{\frac{T}{n} \varepsilon_t x_{t-1}}{T} \\ \text{plim} \frac{\frac{T}{n} \varepsilon_t x_{t-2}}{T} \end{pmatrix} \stackrel{\text{WUN}}{=} 0 \quad \text{if } E(\varepsilon_t | x_t, x_{t-1}, \dots) = 0$$

□

## DYNAMIC: AUTOREGRESSIVE DISTRIBUTED LAG MODELS

$$y_t = \alpha + \phi y_{t-1} + x_t \gamma_1 + x_{t-1} \gamma_2 + \varepsilon_t, |\phi| < 1, \quad E(\varepsilon_t | x_t, x_{t-1}, \dots, y_{t-1}, y_{t-2}, \dots) = 0 \quad (\text{contemporaneous exogeneity})$$

\*predetermined\*

- Omitting  $y_{t-1}$  if  $y_{t-1}$  is correlated with  $x_t$  causes OVB ( $\gamma_1$  captures additionally the effects of  $y_{t-1}$ ); control for it!!
- $\gamma_1$  is the contemporaneous effect
- $\frac{\gamma_1 + \gamma_2}{1-\phi}$  is the long run propensity in equilibrium  $(x^*, y^*)$ ,  $\varepsilon_t = E(\varepsilon_t) = 0$  so  $y^* = \alpha + \phi y^* + (\gamma_1 + \gamma_2) x^* \Rightarrow y^* = \frac{\alpha}{1-\phi} + \frac{\gamma_1 + \gamma_2}{1-\phi} x^*$

$$X = \begin{pmatrix} 1 & T & T & T \\ 1 & y_{t-1} & x_t & x_{t-1} & x_{t-2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \text{income} & \text{income} & \text{income} & \text{income} \\ 2020 & 2020 & 2019 & 2018 \end{pmatrix} \quad \beta = \begin{pmatrix} \alpha \\ \phi \\ \gamma_1 \\ \gamma_2 \end{pmatrix}$$

$$E(\varepsilon_t | X) = E(\varepsilon_t | x_1, \dots, x_T, y_1, \dots, y_{t-1}) = 0 \text{ cannot be true}$$

Since  $\text{Cov}(\varepsilon_t, y_t) \neq 0$  by definition.  $\hat{\beta}_{OLS}$  biased! PS6 Q2

- With contemporaneous exogeneity,  $\hat{\beta}_{OLS}$  consistent under stationarity and weak dependence Some proof as above
- Not consistent under  $\varepsilon_t$  autocorrelation

# STATIONARITY

replaces identically distributed

Joint dist of  $(x_{t_1}, \dots, x_{t_m})$  identical to that of  $(x_{t+h_1}, \dots, x_{t+m})$   $\forall h$   
 • Actually stronger than identically distributed ( $id \not\Rightarrow Cov\ stat$ )

**DEFINITION** A stochastic process  $\{x_t : t=1, 2, \dots\}$  is covariance stationary if

- ①  $E(x_t)$  is finite and constant (i.e. same  $\mu_t$ )
  - ②  $Var(x_t)$  is finite and constant (i.e. same  $\sigma^2_t$ )
  - ③  $Cov(x_t, x_{t+h})$  is finite and depends only on distance in time,  $h$  (i.e. same  $\rho_h$  holding  $h$  constant)
- Implies the first two moments must exist and not change over time
  - Can allow for deterministic trending behaviour that do not affect dependence of process over time
    - Just detrend it first!

## DETERMINISTIC TREND

$$y_t = \alpha + \beta t + \varepsilon_t, t=1, \dots, T$$

$$E(y_t) = \alpha + \beta t \text{ changes with } t$$

$$Var(y_t) = Var(\varepsilon_t)$$

$$Cov(y_t, y_{t+h}) = Cov(\varepsilon_t, \varepsilon_{t+h})$$



If  $\varepsilon_t$  is covariance stationary, the detrended  $y_t - \beta t = \alpha + \varepsilon_t$  is covariance stationary

## WEAK DEPENDENCE

replaces independence in TS data

Weak dependence requires  $\text{Corr}(x_t, x_{t+h}) = \frac{\text{Cov}(x_t, x_{t+h})}{\sqrt{\text{Var}(x_t)\text{Var}(x_{t+h})}} \rightarrow 0$  as  $h \rightarrow \infty$  (asymptotic independence)

**AUTOCORRELATION FUNCTION**  $\rho(h) = \text{Corr}(x_t, x_{t+h})$  (Sample analogue is the correlogram)

- $\rho(0) = \text{Corr}(x_t, x_t) = 1$

## COMMON DEPENDENCE PROCESSES

**AUTOREGRESSIVE PROCESS** of order 1 AR(1)

$$y_t = \phi y_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim i.i.d.(0, \sigma^2)$$

white noise

Stationarity requires  $|\phi| < 1$

$> 1$ : no finite mean var  
 $= 1$ : random walk

- By recursive substitution,  $y_t = \varepsilon_t + \phi \varepsilon_{t-1} + \phi^2 \varepsilon_{t-2} + \dots$  as  $\phi^s \rightarrow 0$
- $\text{Var}(y_t) = \text{Var}(\varepsilon_t) + \phi^2 \text{Var}(\varepsilon_{t-1}) + \phi^4 \text{Var}(\varepsilon_{t-2}) + \dots = (1 + \phi^2 + \phi^4 + \phi^6 + \dots) \sigma^2$  converges iff  $|\phi| < 1$ 
  - if an intercept is added, mean will diverge too
  - $\rho(h) = \phi^h$ : a shock affects all future obs with decreasing (geometric) effect
  - $E(y_t) = 0$ ,  $\text{Var}(y_t) = \frac{\sigma^2}{1-\phi^2}$ ,  $\text{Cov}(y_t, y_{t+h}) = \phi^h$

**MOVING AVERAGE PROCESS** of order 1 MA(1)

$$y_t = \varepsilon_t + \theta \varepsilon_{t-1}, \quad \varepsilon_t \sim i.i.d.(0, \sigma^2)$$

- Always stationary
- $\rho(0) = 1$ ,  $\rho(1) = \frac{\theta}{1+\theta^2}$ ,  $\rho(h) = 0 \quad \forall h = 2, 3, \dots$  Shock only affects  $y_t$  in 2 periods
- $E(y_t) = 0$ ,  $\text{Var}(y_t) = \sigma^2(1+\theta^2)$ ,  $\text{Cov}(y_{t+1}, y_t) = \theta \sigma$ ,  $\text{Cov}(y_{t+h}, y_t) = 0 \quad \forall h > 1$  PSB extra  $\theta$

**AUTOREGRESSIVE MOVING AVERAGE PROCESS**

ARMA(1,1)

$$y_t = \phi y_{t-1} + \varepsilon_t + \theta \varepsilon_{t-1}, \quad \varepsilon_t \sim i.i.d.(0, \sigma^2)$$

# HETEROSCEDASTICITY & AUTOCORRELATION

## DEFINITION

A1  $y = x\beta + \varepsilon$ ,  $E(\varepsilon) = 0$

A2 No perfect multicollinearity

A3  $E(\varepsilon|x) = 0$

A4\* General covariance matrix  $\text{Var}(\varepsilon|x) \stackrel{\text{A3}}{=} E(\varepsilon\varepsilon'|x) = \Sigma$

$\Sigma$  is a symmetric positive semidefinite matrix

often written as  $\Sigma = \sigma^2 \Omega$   
↑  
unknown scaling parameter

## HETEROSCEDASTICITY

$$\text{Cov}(\varepsilon_i, \varepsilon_j | X) = \begin{cases} \sigma^2 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases} \quad \Sigma = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{pmatrix}$$

## AUTOCORRELATION

$$\text{Cov}(\varepsilon_t, \varepsilon_s | X) \neq 0 \text{ for some } t \neq s$$

e.g. AR(1): if  $\varepsilon_t = \rho \varepsilon_{t-1} + v_t$ ,  $|\rho| < 1$ ,  $v_t \sim \text{iid}(0, \sigma^2)$

$$\Sigma = \frac{\sigma^2}{1-\rho^2} \begin{pmatrix} 1 & \rho & \rho^2 & \rho^3 & \dots & \rho^{n-1} \\ \rho & 1 & & & & \\ \rho^2 & & 1 & & & \\ \rho^3 & & & 1 & & \\ \vdots & & & & \ddots & \\ \rho^{n-1} & & & & & 1 \end{pmatrix} \quad \text{"fading memory"} \quad \text{homoscedastic, but autocorrelated by cov stationarity}$$

## ISSUE

$\hat{\beta}_{OLS}$  not BLUE!

A1-3  
Unbiased ✓ Consistent ✓ Linear ✓ But efficiency loss ✗

$$\text{Var}(\hat{\beta}|X) = (X'X)^{-1} X' \Sigma X (X'X)^{-1} \text{ given A.4*}$$

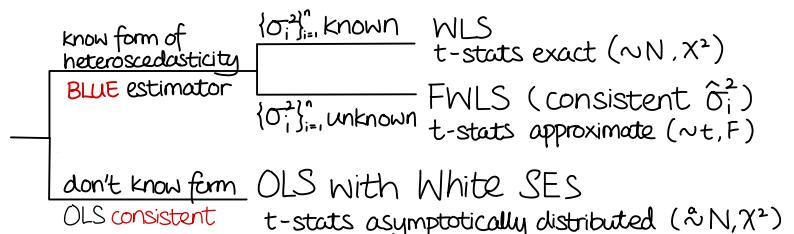
$$= \sigma^2 (X'X)^{-1} X' \Omega X (X'X)^{-1} \quad \text{Standard t- and F-tests invalid}$$

## SOLUTION

- ① Stick with  $\hat{\beta}_{OLS}$  but use <sup>RSE</sup> **correct** estimator of covariance matrix (heteroscedasticity-robust inference)
  - ✓ Don't have to specify form of heteroscedasticity/serial correlation
  - ✗ Lose efficiency
- ② Derive alternative estimator  $\hat{\beta}_{ALS}$  that is BLUE
  - ✓ Regain efficiency (asymptotically)
  - ✗ Must specify form of heteroscedasticity/serial correlation

## APPROACHES

Heteroscedasticity only  
 $\text{Var}(\varepsilon_i|x) = \sigma_i^2 \neq \sigma^2$ ,  $\text{Cov}(\varepsilon_i, \varepsilon_j|x) = 0$



Autocorrelation too

OLS consistent

OLS with HAC/Newey-West SEs  
t-stats asymptotically distributed (~N, X²)

Note: Should only use RSEs when there's evidence of heteroscedasticity/autocorrelation

# HETEROSCEDASTICITY-ROBUST INFERENCE

When there is heteroscedasticity but no autocorrelation,

$$\text{Var}(\hat{\beta}|X) = \left( \sum_{i=1}^n x_i x_i' \right)^{-1} \left( \sum_{i=1}^n \sigma_i^2 x_i x_i' \right) \left( \sum_{i=1}^n x_i x_i' \right)^{-1} \quad x_i' = (x_{j1} \dots x_{jk})$$

When  $\{\sigma_i^2\}_{i=1}^n$  is unknown

$$\text{Var}(\hat{\beta}|X) = \left( \sum_{i=1}^n x_i x_i' \right)^{-1} \left( \sum_{i=1}^n \hat{\varepsilon}_i^2 x_i x_i' \right) \left( \sum_{i=1}^n x_i x_i' \right)^{-1} \quad \text{HCSE/White SE}$$

- ★ Consistent but only valid asymptotically
- ★ Does not imply  $\hat{\varepsilon}_i^2$  is a consistent estimator of  $\sigma_i^2$   
but that  $\text{plim} \left( \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 x_i x_i' - \frac{1}{n} \sum_{i=1}^n \sigma_i^2 x_i x_i' \right)$  (averages)

## WEIGHTED LEAST SQUARES

With weights  $r_i = \frac{1}{\sigma_i^2}$

Satisfies A1-4  $\frac{y_i}{\sigma_i} = \left( \frac{x_i}{\sigma_i} \right)' \beta + \frac{\varepsilon_i}{\sigma_i}$  Same as having  $R = \begin{pmatrix} \frac{1}{\sigma_1} & \cdots & \frac{1}{\sigma_n} \\ 0 & \cdots & 0 \end{pmatrix}$

e.g. A3:  $\text{Var}\left(\frac{\varepsilon_i}{\sigma_i}|X\right) = \frac{1}{\sigma_i^2} \text{Var}(\varepsilon_i|X) = 1$  (constant)

$$\hat{\beta} = \left[ \frac{\sum_{i=1}^n x_i x_i'}{\sum_{i=1}^n \sigma_i^2} \right]^{-1} \sum_{i=1}^n \frac{x_i y_i}{\sigma_i^2}$$

can also replace  $\sigma_i$  with weights  $w_i$  s.t.  $\sigma_i^2 = \sigma^2 w_i$   
if we know  $\sigma_i^2$ , just do standard t/F-tests

## FEASIBLE WLS

$$= \text{Var}(\varepsilon_i|x_i) \stackrel{\text{A3}}{=} E(\varepsilon_i^2|x_i)$$

Don't know  $\sigma_i^2$ ? Parameterise  $\sigma_i^2$  to get  $\hat{\sigma}_i^2$

$$\hat{\beta}_{\text{FALS}} = \left[ \frac{\sum_{i=1}^n x_i x_i'}{\sum_{i=1}^n \hat{\sigma}_i^2} \right]^{-1} \sum_{i=1}^n \frac{x_i y_i}{\hat{\sigma}_i^2}$$

One way to parameterise  $\sigma_i^2$

**LINEAR HETEROSCEDASTICITY**  $\sigma_i^2 = \delta_0 + z_i \delta_1 = E(\varepsilon_i^2|x_i)$   $z_i$  is some function of  $x_i$   
 $\delta_1$  can be negative !!

$\hat{\varepsilon}_i^2$  is a good guess of  $\varepsilon_i^2$ , which on average is  $\sigma_i^2$

Then we find consistent  $\hat{\delta}_0, \hat{\delta}_1$

**DERIVING**  $\varepsilon_i^2 = \delta_0 + z_i \delta_1 + v_i$

$$E(\varepsilon_i^2|x_i) = \delta_0 + z_i \delta_1 + E(v_i|x_i) = \delta_0 + z_i \delta_1 \Rightarrow \text{assume GM} \Rightarrow \text{Regress } \hat{\varepsilon}_i \text{ on } Z \text{ to get BLUE } \hat{\delta}_0, \hat{\delta}_1$$

**TESTING**  $H_0: \delta_0 = \delta_1 = 0$  vs  $H_1: \neg H_0$  asymptotic bc we're using  $\hat{\varepsilon}_i^2$ , not  $\varepsilon_i^2$

$z_i$  is some function of  $x_i$

**EXPONENTIAL HETEROSCEDASTICITY**  $\sigma_i^2 = e^{\delta_0 + z_i \delta_1} = \sigma^2 e^{z_i \delta_1}$   $\sigma^2 = e^{\delta_0}$  is constant of proportionality

Regress  $\log \hat{\varepsilon}_i = e^{\delta_0 + z_i \delta_1} + e_i$

# GENERALISED LINEAR MODEL

LM but relaxes A4  $\text{Var}(\varepsilon|X) = \sigma^2 I$

Find square, invertible matrix  $R$  s.t.  $R'R = \Omega^{-1}$ . Such an  $R$  exist  
 Transform model  $Ry = RX\beta + R\varepsilon$  ( $y^* = X^*\beta + \varepsilon^*$ )  
 Then A1-4 are satisfied  $\Rightarrow$  OLS on transformed model is BLUE

quadratic form suitably weighted by precision

**MINIMISE GENERALISED SUM OF SQUARES**

$$\hat{\beta}_{OLS} = \arg \min_b S(b), S(b) = (y - Xb)' \Omega^{-1} (y - Xb)$$

$$\text{can also use } (y - Xb)' \Sigma^{-1} (y - Xb) = \varepsilon' \text{Var}(\varepsilon|X)^{-1} \varepsilon$$

Accounts for the fact that some observations are associated with higher variability than others by giving their discrepancies (residuals) less weight

When  $S(b)$  is at minimum,

$$\boxed{\frac{\partial S(b)}{\partial b} = \left( \begin{array}{c} \frac{\partial S(b)}{\partial b_1} \\ \vdots \\ \frac{\partial S(b)}{\partial b_k} \end{array} \right) = 2y' \Omega^{-1} X - 2X' \Omega^{-1} X b = 0 \Rightarrow \hat{\beta} = (X' X)^{-1} X' y}$$

$$\hat{\beta}_{GLS} = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1} y$$

but typically infeasible as we don't know  $\Omega$   
 replace with  $\hat{\Omega}$  (feasible GLS)

$$\text{or } \hat{\beta}_{GLS} = (X^* X^*)^{-1} X^* y^*$$

**PROPERTIES OF  $\hat{\beta}_{GLS}$**  Derive from transformed model  $y^* = X\beta + \varepsilon^*$ ,  $\begin{cases} E(\varepsilon^*|X^*) = 0 \\ \text{Var}(\varepsilon^*|X^*) = 0 \end{cases}$

$$E(\hat{\beta}_{GLS}) = \beta$$

$$\text{Var}(\hat{\beta}_{GLS}|X) = \sigma^2 (X^* X^*)^{-1} = \sigma^2 (X' \Omega^{-1} X)^{-1}$$

**PROPERTIES OF  $\hat{s}_{GLS}^2$**   $\hat{s}_{GLS}^2 = \frac{\hat{\varepsilon}^* \hat{\varepsilon}}{N-k} = \frac{(y - \hat{X}\hat{\beta}_{GLS})' \Omega^{-1} (y - \hat{X}\hat{\beta}_{GLS})}{N-k}$

Unbiased:  $\hat{s}_{GLS}^2 = \frac{\hat{\varepsilon}' M \hat{\varepsilon}}{N-k}$ ,  $M = I - X^*(X^* X^*)^{-1} X^*$ , proceed

If A5:  $\varepsilon|X \sim N(0, \sigma^2 \Omega)$ , then  $\hat{\beta}_{GLS}|X \sim N(\beta, \sigma^2 (X' \Omega^{-1} X)^{-1})$  and  $\hat{\beta}_{GLS} \perp \hat{s}_{GLS}^2$

## FEASIBLE GLS

$$\hat{\beta}_{FGLS} = (X' \hat{\Omega}^{-1} X)^{-1} X' \hat{\Omega}^{-1} y$$

$$\text{e.g. } \hat{\Omega}(\theta) = \frac{\sigma^2}{1-\theta^2} \begin{pmatrix} 1 & \theta & \theta^2 & \theta^3 & \dots & \theta^{n-1} \\ \theta & 1 & & & & \\ \theta^2 & & 1 & & & \\ \theta^3 & & & 1 & & \\ \vdots & & & & \ddots & \\ \theta^{n-1} & & & & & 1 \end{pmatrix}$$

put structure on

To derive a consistent  $\hat{\Omega}$ , we must parameterise  $\Omega$  in terms of a finite-dimensional parameter vector  $\theta$   
 $\Omega = \Omega(\theta)$ . Then use  $\hat{\varepsilon} = y - \hat{X}\hat{\beta}$  to derive  $\hat{\theta}$ , s.t.  $\text{plim}(\hat{\theta}) = \theta$ , so  $\text{plim}(\hat{\Omega}) = \Omega(\text{plim} \hat{\theta}) = \Omega(\theta)$

e.g. if  $\Omega$  is diagonal,  $\Omega_{ii} = z_i \gamma$   $\sim$  unobserved

1. Regress  $y$  on  $X$  to get  $\hat{\varepsilon} = y - \hat{X}\hat{\beta}_{OLS}$

2. Regress  $\hat{\varepsilon}_i^2$  on  $z_i$  to get  $\hat{\gamma}$  (since  $E(\varepsilon_i^2) = \text{Var}(\varepsilon_i) + [E(\varepsilon_i)]^2 = \text{Var}(\varepsilon_i) = z_i \gamma$ )

3. Thus  $\hat{\Omega}_{ii} = z_i \hat{\gamma}$ ,  $\hat{\Omega} = \begin{pmatrix} \hat{\Omega}_{11} & & \\ & \ddots & \\ & & \hat{\Omega}_{nn} \end{pmatrix}$

**$\hat{\beta}_{FGLS}$  ASYPTOTIC PROPERTIES**

$\hat{\beta}_{FGLS}$  is usually biased, and isn't even linear ( $\hat{\varepsilon}$ , which is a function of  $y$ , is used to derive  $\hat{\Omega}$ , so  $y$  enters  $>$  once)  $\hat{\Omega}$   
 so  $\hat{\beta}_{FGLS}$  is not BLUE; shouldn't be used with small  $N$  too

Asymptotic normal:  $\hat{\beta}_{FGLS}|X \stackrel{\text{d}}{\sim} N(\beta, \sigma^2 (X' \Omega^{-1} X)^{-1})$

- We have to use this property even if we assume/know  $\varepsilon|X \sim N(0, \sigma^2 \Omega)$  because we're also using  $\hat{\Omega}$
- Only asymptotically does  $\hat{\beta}_{FGLS}$  inherit the desirability of  $\hat{\beta}_{GLS}$

# TEST HETROSCEASTICITY

## WALD TEST (FGLS) SETTING

Test for heteroscedasticity of specific form

$$\begin{cases} H_0: \sigma_i^2 = \sigma^2 \forall i \\ H_A: \sigma_i^2 = \delta_0 + z_i' \delta_1 \text{ or } \sigma_i^2 = e^{\delta_0} e^{z_i' \delta_1} \end{cases} \Leftrightarrow \begin{cases} H_0: \delta_1 = 0 \\ H_A: \text{at least one } \delta_i \text{ is non-zero} \end{cases}$$

$\Leftrightarrow$  Test significance of the regression  $\hat{\varepsilon}_i^2 = \delta_0 + z_i' \delta_1 + v_i$   
or  $\ln(\hat{\varepsilon}_i^2) = \delta_0 + z_i' \delta_1 + v_i$

$\Leftrightarrow$  Read F-stat, or (more accurately since we don't have  $\varepsilon_i^2$ )

calculate  $\hat{\delta}_1' (\widehat{\text{Var}}(\hat{\delta}_1))^{-1} \hat{\delta}_1 \stackrel{\text{asymptotically!}}{\sim} X_p^2$ , where  $p = \dim(\delta_1)$

## BREUSH-PAGAN TEST

An example of a Lagrange Multiplier test

$$\begin{cases} H_0: \sigma_i^2 = \sigma^2 \forall i \\ H_A: \sigma^2 h(\delta_0 + z_i' \delta_1) \text{ unknown, continuously differentiable} \end{cases}$$

sample size

$$L = nR^2 \stackrel{\downarrow}{\sim} X_p^2 \text{ under } H_0 \text{ where } p = \dim(\delta_1)$$

Reject  $H_0$  if  $nR^2 > X_{p,\alpha}^2$

$\downarrow$  OLS residual

$R^2$  is goodness-of-fit from auxiliary regression  $\hat{\varepsilon}_i^2 = \delta_0^* + z_i' \delta_1^* + v_i$

- High  $R^2 \Rightarrow z_i$  and  $\hat{\varepsilon}_i$  very related  $\Rightarrow$  heteroscedasticity
- Cannot estimate  $\delta_1$  and do Wald  $\hat{\delta}_1' (\widehat{\text{Var}}(\hat{\delta}_1))^{-1} \hat{\delta}_1$ , since form is not explicit
- Also cannot do F test on  $\delta_1^*$  since form is not explicit

## WHITE TEST

No idea of potential form of heteroscedasticity

$$\begin{cases} H_0: \sigma_i^2 = \sigma^2 \forall i \\ H_A: \neg H_0 \end{cases} \Leftrightarrow \begin{cases} H_0: E(v_i^2 | x_i) = \sigma^2 \forall i \text{ if homoscedastic, then } x_i \text{/functions of } x_i \text{ shouldn't explain } E(v_i^2 | x_i) \\ H_A: \neg H_0 \end{cases}$$

Run auxiliary regression  $\hat{\varepsilon}_i^2 = \gamma_0 + z_i' \gamma_1 + v_i$

Test  $\gamma_1 = 0$ , compute  $R^2$ .  $nR^2 \stackrel{\downarrow}{\sim} X_p^2$  where  $p = \dim(\gamma_1)$

Reject  $H_0$  if  $nR^2 > X_{p,\alpha}^2$

- $z_i$  may include some/all variables in  $x_i$  and other variables that depends on  $x_i$
- White "use  $x_i$ , and set of all unique squares and cross products of variables in  $x_i$ "
  - e.g.  $x_{i2}, x_{i3}, x_{i2}^2, x_{i3}^2, x_{i2}x_{i3}$
  - But cannot use  $x_{i2}^2$  if  $x_{i2}$  is a dummy

# AUTOCORRELATION - ROBUST INFERENCE

$$\text{Var}(\hat{\beta}|X) = (X'X)^{-1} X' \Sigma X (X'X)^{-1}$$

**USE OLS WITH RSE** Only when OLS is consistent

- At least contemporaneous exogeneity  $E(\varepsilon_t|x_t)=0$
- Not useful with lagged y-values - use Newey-West/HAC RSE instead (consistent)
  - When  $\text{Cov}(\varepsilon_t, \varepsilon_{t+s}|X) = \gamma_s$ , with  $\gamma_s=0$  when  $s > \text{lag}$
  - HAC also allows lag to  $\uparrow$  as  $T \uparrow$
  - Stata: newey y  $x_1 - x_k$ , lag (insert no here)  
R can auto-det best lag

Hoodridge (v technical)

## TEST FOR AUTOCORRELATION

Must specify autocor. model, e.g. AR(1)  $\varepsilon_t = \rho \varepsilon_{t-1} + v_t$

$$H_0: \rho = 0 \quad \text{vs} \quad H_A: \rho \neq 0 \quad \text{or} \quad \rho > 0$$

Use  $\hat{\varepsilon}_t$  since we can't observe  $\{\varepsilon_t\}$

**UNDER STRICT EXOGENEITY**

- ① Estimate  $y_t = \beta_0 + \beta_1 x_{1t} + \dots + \beta_k x_{kt} + u_t$  by OLS, save  $\{\hat{u}_t\}$
  - ② Run AR(1)  $\hat{u}_t$  on  $\hat{u}_{t-1}$
  - ③  $\frac{\hat{\rho}}{\text{SE}(\hat{\rho})} \sim N(0,1)$  Asymptotic t-test as we use  $\hat{u}_t$ , not  $u_t$
- ★ Can also add lags - regress  $\hat{u}_t$  on  $\hat{u}_{t-1}, \hat{u}_{t-2}$ , test for joint significance

**UNDER CONTEMPORANEOUS EXOGENEITY**

Same as above but for ②

Run OLS of  $\hat{u}_t$  on  $\hat{u}_{t-1}, x_{1t}, \dots, x_{kt}$   $\leftarrow$  controls

- Account for the fact that  $u_{t-1}$  might be correlated with  $x_{1t}, \dots, x_{kt}$  not strictly exogenous
- Use if  $x_{ij}$ 's contain lagged y

might respond to past shocks

**NOTE** Not tested

We note that  $\begin{cases} y_t = \alpha_0 + \beta x_t + \varepsilon_t \\ \varepsilon_t = \rho \varepsilon_{t-1} + u_t \end{cases}$   $|\rho| < 1$ ,  $u_t \sim (0, \sigma^2)$  i.i.d. is GM error

Transform model to regain efficiency:  $y_t - \rho y_{t-1} = (1-\rho)\alpha + \beta x_t - \rho \beta x_{t-1} + \varepsilon_t - \rho \varepsilon_{t-1}$   
 $= (1-\rho)\alpha + \beta(x_t - \rho x_{t-1}) + u_t$   $\text{IS GM!!}$   
 new regressor

- We dk  $\rho$ , so make it feasible by regressing on residuals
- Will lose first obs

# ENDOGENEITY

Correlation between errors and regressors  
 $E(\varepsilon|X) \neq 0$

OVB  
 lagged dep var when errors exhibit dep  
 regressor measurement error  
 simultaneity

## CONSEQUENCE FOR OLS

biased  $E(\hat{\beta}|X) = \beta + (XX)^{-1}X'E(\varepsilon|X) \neq \beta$   
 inconsistent  $\text{plim } \hat{\beta} = \beta + \text{plim} \left( \frac{XX}{n} \right)^{-1} \text{plim} \left( \frac{X\varepsilon}{n} \right) \neq \beta$   
 $\downarrow$   
 $= E(x_i \varepsilon_i)$  by LLN

Not surprising as FOC of OLS require  $X'\hat{\varepsilon} = \sum_{i=1}^n x_i \hat{\varepsilon}_i = 0$ , which is proportional to  $\frac{1}{n} \sum_{i=1}^n x_i \hat{\varepsilon}_i = 0$ .  
 sample analogue of  $E(x_i \varepsilon_i) = 0$ . If  $E(x_i \varepsilon_i) \neq 0$  then imposing orthogonality is inappropriate

## LAGGED DEP VARS & SERIALLY CORREL ERRORS

$$y_t = \beta_1 + \beta_2 x_t + \beta_3 y_{t-1} + \varepsilon_t, \quad |\beta_3| < 1, \quad t=1, \dots, T$$

$$\varepsilon_t = \rho \varepsilon_{t-1} + v_t, \quad |\rho| < 1 \quad v_t \stackrel{iid}{\sim} (0, \sigma_v^2)$$

$\text{Cov}(y_{t-1}, \varepsilon_t) = E(y_{t-1} \varepsilon_t) \neq 0$  so  $y_{t-1}$  is endogenous

$$\begin{aligned} \text{Cov}(y_{t-1}, \varepsilon_t) &= \text{Cov}(\rho \varepsilon_{t-1} + v_t, \beta_1 + \beta_2 x_{t-1} + \beta_3 y_{t-2} + \varepsilon_{t-1}) \\ &= \rho \beta_3 \text{Cov}(\varepsilon_{t-1}, x_{t-1}) + \rho \beta_3 \text{Cov}(\varepsilon_{t-1}, y_{t-2}) + \rho \text{Var}(\varepsilon_{t-1}) \\ &= \rho \beta_3 \text{Cov}(\varepsilon_t, y_{t-1}) + \rho \sigma^2 \end{aligned}$$

$\text{Cov}_{\text{stat}}$

$$\Rightarrow \text{Cov}(\varepsilon_t, y_{t-1}) = \frac{\rho \sigma^2}{1 - \rho \beta_3} \neq 0 \quad \text{assuming covariance stationarity } \rho \beta_3 \neq 1$$

$$\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3 \text{ inconsistent: let } X = \begin{pmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_{t-1} & y_t & y_{t-2} \end{pmatrix} \text{ then } \text{plim } \frac{X\varepsilon}{T} = \text{plim} \begin{pmatrix} \frac{1}{T} \sum x_t \\ \frac{1}{T} \sum x_t \varepsilon_t \\ \frac{1}{T} \sum y_{t-1} \varepsilon_t \end{pmatrix} \stackrel{\text{LLN}}{=} \begin{pmatrix} E(\varepsilon_{t-1}) \\ E(\varepsilon_t x_t) \\ E(\varepsilon_t y_t) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \text{Cov}(y_{t-1}, \varepsilon_t) \end{pmatrix}$$

## MEASUREMENT ERROR

True model:  $y_i = x_i^* \beta + \varepsilon_i$

Mismeasurement:  $x_i = x_i^* + u_i$  composite error

Mis measured model:  $y_i = x_i^* \beta + v_i = x_i^* \beta + \varepsilon_i - u_i \beta$

## CLASSICAL MEASUREMENT ERROR ASSUMPTION

true value of regressor being measured  
 $u_i$  is independent of  $x_i^*$  and  $\varepsilon_i^*$  ( $u_i \sim i.i.d (0, \sigma_u^2)$ )

$\hat{\beta}$  inconsistent as  $\text{plim } \hat{\beta} = \beta + (\text{plim}(\frac{XX}{n}))^{-1} \text{plim}(\frac{Xv}{n}) \stackrel{\text{LLN}}{=} \beta + E(x_i x_i')^{-1} E(x_i v_i) \neq \beta$  because

$$\begin{aligned} E(x_i v_i) &= E((x_i^* + u_i)(\varepsilon_i - u_i \beta)) \\ &= E(x_i^* \varepsilon_i) + E(u_i \varepsilon_i) - E(x_i^* u_i) \beta - E(u_i u_i) \beta \\ &= -E(u_i u_i) \beta \neq 0 \quad \text{by GM/assume } x_i^* \text{ fixed} \end{aligned}$$

$$E(x_i x_i') = E(x_i^* x_i^*) + E(u_i u_i')$$

In general all parameters will be affected by measurement error in one/more regressors

Attenuation bias: when just one regressor has CME, estimated parameters will be closer to 0 PS8 Q1

**SIMULTANEITY**  $X$  jointly determined with  $y$  in the same economic model

**STRUCTURAL FORM** gives relationships of interest (behavioural relations)

$$\begin{cases} \text{Demand} & q_t = \alpha + \alpha_2 p_t + \alpha_3 m_t + u_{1t} \\ \text{Supply} & q_t = \beta_1 + \beta_2 p_t + \beta_3 c_t + u_{2t} \end{cases} \quad \begin{matrix} \text{exogenous/predetermined - independent of } \varepsilon_t \\ \text{usually correlated} \end{matrix}$$

If  $u_{1t}, u_{2t}$  correlated then  $\text{Cov}(u_{1t}, u_{2t}) = \sigma_{12}$

$$\begin{cases} C_t = \alpha + \beta Y_t + \varepsilon_t, E(\varepsilon_t) = 0 \\ Y_t = C_t + I_t \end{cases}$$

OLS  $(\alpha, \beta)$  inconsistent since  $\text{Cov}(Y_t, \varepsilon_t) = \frac{1}{1-\beta} \sigma_\varepsilon^2 \neq 0$

**REDUCED FORM** expresses endo var in terms of exo vars and errors

$$\begin{cases} Y_t = \frac{\alpha}{1-\beta} + \frac{1}{1-\beta} I_t + \frac{1}{1-\beta} \varepsilon_t, \beta \neq 1 \\ P_t = \pi_{11} + \pi_{12} m_t + \pi_{13} C_t + V_{1t} \\ q_t = \pi_{21} + \pi_{22} m_t + \pi_{23} C_t + V_{2t} \end{cases} \quad \leftarrow \text{fn of } u_{1t} \& u_{2t}$$

## IV ESTIMATION

$$y = X\beta + \varepsilon \quad \text{Cov}(X, \varepsilon) \neq 0$$

### ASSUMPTIONS

**Validity**:  $\text{Cov}(z_t, \varepsilon_t) = E(z_t \varepsilon_t) = 0$  not correlated with  $\varepsilon$ . Implied by  $E(\varepsilon_t | z_t) = 0$

**Relevance**:  $E(z_t' z_t)$  has rank  $k \Rightarrow$  if square, invertible.  $Z$  is correlated with  $X$

- # instruments  $\geq$  # regressors (**order condition**), if not rank  $< k$
- Instruments cannot have a direct effect on  $y$  for full rank

Our IV estimator  $\hat{\beta}_{IV}$  is given by the sample analogue of the moment condition  $E(z_i \varepsilon_i) = 0$  validity

$$\frac{1}{n} Z' \hat{\varepsilon}^{IV} = \frac{1}{n} Z' (y - X \hat{\beta}_{IV}) = 0$$

If  $Z'X$  is square and invertible (**just-identified**)

$$\hat{\beta}_{IV} = (Z'X)^{-1} Z' y \quad \text{if all regressors are exogenous, just take } Z = X$$

need to have enough instruments (order condition)

$$(\text{Calculus - bivariate}) \quad \hat{\beta}_{1,IV} = \beta_1 + \sum_{i=1}^n d_i u_i, \quad d_i = \frac{z_i - \bar{z}}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})} \leftarrow \neq 0 \text{ by relevance}$$

## BIASED BUT CONSISTENT

**CALCULUS**  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$  with  $E(\varepsilon_i) = 0$  but  $E(\varepsilon_i x_i) \neq 0$

$$\hat{\beta}_{I,IV} = \beta_1 + \sum_{i=1}^n d_i \varepsilon_i.$$

$$d_i = \frac{z_i - \bar{z}}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})}$$

$\uparrow$   
 $\neq 0$  by relevance

LIE  $E(\hat{\beta}_{I,IV}|x,z) = \beta_1 + \sum_{i=1}^n d_i E(\varepsilon_i|x,z)$  correlated..  $\neq \beta$  is biased

$$\text{plim } \hat{\beta}_{IV} = \beta + \frac{\text{plim} (\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(\varepsilon_i - \bar{\varepsilon}))}{\text{plim} (\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x}))} \stackrel{\text{LLN}}{=} \beta + \frac{\text{Cov}(z_i, \varepsilon_i)}{\text{Cov}(z_i, x_i)} \stackrel{\substack{=0 \text{ validity} \\ \neq 0 \text{ relevance}}}{=} \beta \text{ IS consistent}$$

$$\begin{aligned} \text{Var}(\hat{\beta}_{I,IV}|x,z) &= \text{Var}(\beta_1 + \sum_{i=1}^n d_i \varepsilon_i | x,z) = \sum_{i=1}^n d_i^2 \text{Var}(\varepsilon_i | x,z) \text{ assuming independence} \\ &= \sigma^2 \sum_{i=1}^n d_i^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \frac{1}{R_{zx}^2} = \frac{\sigma^2}{n \text{Var}(x_i)} \frac{1}{R_{zx}^2} \text{ assuming conditional heteroscedasticity} \end{aligned}$$

For more precision:  $\downarrow n \downarrow \sigma^2 \uparrow$  regressor variability  $\uparrow$  strong first stage effect

Use exogenous variables as instruments for themselves!

## LINEAR ALGEBRA

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

LIE  $E(\hat{\beta}_{IV}|X,Z) = \beta + (ZX)^{-1} Z'E(\varepsilon|X,Z) \neq \beta$  is biased

$$\text{plim } \hat{\beta}_{IV} = \beta + (\text{plim } \frac{1}{n} Z'X)^{-1} \text{plim } \frac{1}{n} Z'\varepsilon = \beta \text{ IS consistent}$$

- $\text{plim } \frac{1}{n} Z'X \stackrel{\text{iid (LLN)}}{=} M_{ZX} = E(Z_i X'_i)$  invertible by relevance and exclusion
- $\text{plim } \frac{1}{n} Z'\varepsilon \stackrel{\text{LLN}}{=} E(Z_i \varepsilon_i) = 0$  by validity

$$\text{Var}(\hat{\beta}_{IV}) = \sigma^2 (Z'X)^{-1} Z'Z (X'Z)^{-1}$$

- $\frac{X'Z}{n}$  could be viewed as the matrix of sample covariances between  $Z$  and  $X$
- Weak first stage will give  $(X'Z)^{-1}$  huge elements and give  $\hat{\beta}_{IV}$  MSE

By CLT,  $\tilde{\beta}_{IV} \approx N(\beta, \sigma^2 (Z'X)^{-1} Z'Z (X'Z)^{-1})$

Assuming conditional A4,  $\text{Var}(\varepsilon|Z) = \sigma^2 I$

$$S_{IV}^2 = \frac{\text{RSS}}{n-k} = \frac{(y - X\hat{\beta}_{IV})'(y - X\hat{\beta}_{IV})}{n-k}$$

## EXCLUSION RESTRICTION

If  $d_t$  has a direct effect on  $y_t$  in the true model:  $y_t = \alpha + \beta w_t + \gamma d_t + \varepsilon_t$ ,  $\text{Cov}(w_t, \varepsilon_t) = E(w_t \varepsilon_t) \neq 0$

Let  $x_t = \begin{pmatrix} 1 \\ d_t \\ w_t \end{pmatrix}$ . If we use the instrument  $z_t = \begin{pmatrix} 1 \\ d_t \\ d_t \end{pmatrix}$ , then  $E(z_t x_t') = \begin{pmatrix} 1 & E(d_t) & E(d_t) \\ E(d_t) & E(d_t^2) & E(d_t^2) \\ E(w_t) & E(w_t d_t) & E(w_t d_t) \end{pmatrix}$  rank < 3

- Thus instruments are variables that don't enter the regression equation itself
- In simultaneous equation models, instruments often come from other equations in the system
- Lags of the included exogenous variables may provide suitable lagged dependent variables

## OPTIMAL INSTRUMENT

When overidentified ( $\# \text{instrument} > \# \text{regressors}$ )

Best choice  $Z^{opt}$  is obtained by using the fitted values of an OLS regression of the columns of  $X$  on all instruments  $Z^{opt} = \hat{X} = Z(Z'Z)^{-1} Z'X = P_Z X$

↑ projects all endo vars to all exo vars

- Lin. comb. of all instruments ( $\checkmark$  valid) with the highest correlation ( $\checkmark$  relevant)

\* if var in  $X$  is exogenous (i.e. it is its own inst) then  $\hat{X}$  perfectly predicts (just returns the var)

- Basically regress  $X = Z\hat{\beta} + w$  to get  $\hat{X} = Z\hat{\beta}$ .

$$\text{Then } \hat{\beta}_{IV} = (Z^{opt'} X)^{-1} Z^{opt'} y = (X' P_Z X)^{-1} X' P_Z y \stackrel{\text{PSL}}{=} (\hat{X}' \hat{X})^{-1} \hat{X}' y$$

## IV FOR LAGGED Y WITH AUTOCORREL

$$y_t = \beta_1 + \beta_2 x_t + \beta_3 y_{t-1} + \varepsilon_t, \quad |\beta_3| < 1, \quad t=1, \dots, T$$

$$\varepsilon_t = \rho \varepsilon_{t-1} + v_t, \quad |\rho| < 1 \quad v_t \stackrel{iid}{\sim} (0, \sigma_v^2)$$

Instrument  $\begin{pmatrix} 1 \\ x_t \\ y_{t-1} \end{pmatrix}$  with  $\begin{pmatrix} 1 \\ x_t \\ x_{t-1} \end{pmatrix}$   $\hat{x} = \begin{pmatrix} 1 \\ x_t \\ \hat{y}_{t-1} \end{pmatrix}$

Can also use more lagged values  $\begin{pmatrix} 1 \\ x_t \\ x_{t-1} \\ x_{t-2} \\ \vdots \end{pmatrix}$  to ↑ correlation, but also ↑ bias

Note: What if supply/demand shifters don't exist

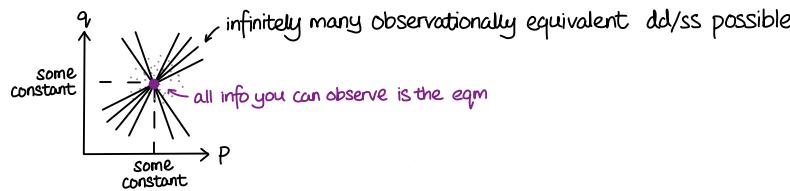
$$\begin{cases} \text{Demand} & q_{it} = \alpha_1 + \alpha_2 p_{it} + \varepsilon_{it} \quad \alpha_2 \neq \beta_2 \\ \text{Supply} & q_{it} = \beta_1 + \beta_2 p_{it} + \varepsilon_{it} \end{cases}$$

Observable data:  $\{(p_t, q_t)\}_{t=1}^T$

When  $D=S$ ,

$$p_t = \frac{\beta_1 - \alpha_1}{\alpha_2 - \beta_2} + \frac{\varepsilon_{2t} - \varepsilon_{1t}}{\alpha_2 - \beta_2} \quad \text{just noise}$$

$$q_t = \alpha_1 + \alpha_2 \frac{\beta_1 - \alpha_1}{\alpha_2 - \beta_2} + \alpha_2 \left( \frac{\varepsilon_{2t} - \varepsilon_{1t}}{\alpha_2 - \beta_2} \right) + \varepsilon_{it} \quad \text{just noise}$$



Parameters under-identified! Can't estimate them

## 2SLS in over-identification case

- ① Compute  $\hat{X}$ , the  $j^{th}$  column of which contains the fitted values of an OLS regression of the  $j^{th}$  column (endogenous regressor) of X on all instruments Z.  $\hat{X} = Z(Z'Z)^{-1}Z'X = P_Z X$   
i.e. regress reduced form of X
- ② Perform OLS of  $y$  on  $\hat{X}$ .  $\hat{\beta}_{IV} = (\hat{X}'\hat{X})^{-1}\hat{X}'y$

★ Correct the SEs (Stata will use an inconsistent  $S^2 = \frac{(y - \hat{X}\hat{\beta}_{IV})'(y - \hat{X}\hat{\beta}_{IV})}{n-k}$  instead of  $S^2 = \frac{(y - X\hat{\beta}_{IV})'(y - X\hat{\beta}_{IV})}{n-k}$ )

## TEST RELEVANCE OF INSTRUMENTS

Obtain reduced form and do t-test (one restriction) or F-test (joint) on coefficient estimates of instruments

# TEST FOR ENDOGENEITY

Given  $y_1 = \beta_0 + \beta_1 y_2 + \beta_2 x_i + u_1, E(u_1) = E(x_i u_1) = 0$

Question: is  $y_2$  and  $u_1$  correlated or not?

$$H_0: \text{Cov}(y_2, u_1) = 0 \quad (y_2 \text{ exogenous})$$

- $\hat{\alpha}_{OLS}$  consistent, efficient (BLUE)
- $\hat{\alpha}_{IV}$  consistent too (but not needed)

$$H_1: \text{Cov}(y_2, u_1) \neq 0 \quad (y_2 \text{ endogenous})$$

- only  $\hat{\alpha}_{IV}$  consistent

## HAUSMAN TEST

Under  $H_0$ ,  $\text{plim}(\hat{\beta}_{IV} - \hat{\beta}_{OLS}) = 0$  while under  $H_1$ ,  $\text{plim}(\hat{\beta}_{IV} - \hat{\beta}_{OLS}) \neq 0$

Furthermore,  $\text{Var}(\hat{\beta}_{IV} - \hat{\beta}_{OLS}) = \text{Var}(\hat{\beta}_{IV}) - \text{Var}(\hat{\beta}_{OLS})$

positive definite as  $\hat{\beta}_{OLS}$  is efficient (has minimum variance) under  $H_0$

Thus, we should compare if the difference is close to 0

In the scalar setting,  $\text{SE}(\hat{\beta}_{IV} - \hat{\beta}_{OLS}) = \sqrt{\text{Var}(\hat{\beta}_{IV}) - \text{Var}(\hat{\beta}_{OLS})} = \sqrt{\text{SE}(\hat{\beta}_{IV})^2 - \text{SE}(\hat{\beta}_{OLS})^2}$

Thus, t-stat is

scalar  $t = \frac{\hat{\beta}_{IV} - \hat{\beta}_{OLS}}{\text{SE}(\hat{\beta}_{IV} - \hat{\beta}_{OLS})} \stackrel{d}{\sim} N(0,1)$ . Reject if  $|t| > z_{\frac{\alpha}{2}}$

vector  $(\hat{\beta}_{IV} - \hat{\beta}_{OLS})' [\text{Var}(\hat{\beta}_{IV}) - \text{Var}(\hat{\beta}_{OLS})]^{-1} (\hat{\beta}_{IV} - \hat{\beta}_{OLS}) \stackrel{d}{\sim} \chi^2_{\dim(\beta)}$  \*  $\chi^2_1$  is quadratic form of  $N(0,1)$

## SIMPLE REGRESSION-BASED TEST

Obtain  $\hat{y}_2$  from fitting  $y_2 = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + \pi_3 x_i + v_2$ , where  $z_1, z_2$  are instruments for  $y_2$

- $\hat{y}_2$  is a "good" regressor where  $y_2$ 's endogeneity is "washed out"
- $\hat{v}_2 = y_2 - \hat{y}_2$  must contain the "bad" component of  $y_2$ 
  - $\text{Cov}(y_2, u_1) = \text{Cov}(\hat{y}_2 + \hat{v}_2, u_1) = \text{Cov}(\hat{y}_2, u_1) + \text{Cov}(\hat{v}_2, u_1)$
- Add  $\hat{v}_2$  to original regression

$$y_1 = \alpha_0 + \alpha_1 y_2 + \alpha_2 x_i + \hat{s} \hat{v}_2 + \text{error}$$

$$\text{Test } H_0: \hat{s} = 0$$

$$\text{Test statistic: } \frac{\hat{s}}{\text{SE}(\hat{s})} \sim N(0,1)$$

Reject $H_0$	Need to "control" for endogeneity of $y_2$ by including $\hat{v}_2$ Resulting parameter estimates $\hat{\alpha}$ will be same as when obtained through 2SLS
--------------	--

Don't reject $H_0$	Can just use original model w/o $\hat{v}_2$ to estimate $\alpha$ Then should just do OLS
--------------------	---

Note: if  $\geq 2$  endogenous variables  $y_1, y_2$ , decompose both, put both residuals back in original regressor and do asymptotic F test on joint hypothesis ( $H_0: S_1 = 0$  and  $S_2 = 0$ )

# MLE

Assume random variables  $y_1, \dots, y_n$  are distributed independently

$$L(\theta; y_1, \dots, y_n) = f(y_1, \dots, y_n; \theta) \stackrel{\text{independence}}{=} \prod_{i=1}^n f(y_i; \theta)$$

$$\text{Take log-likelihood } \log L(\theta; y_1, \dots, y_n) = \sum_{i=1}^n \log L_i(\theta)$$

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \log L(\theta; y_1, \dots, y_n) \equiv \underset{\theta \in \Theta}{\operatorname{argmax}} \log L(\theta; y)$$

$$\text{FOC: } \frac{\partial \log L(\theta; y)}{\partial \theta} \Big|_{\hat{\theta}_{MLE}} = 0$$

$$\text{SOC (to verify maximum): } \frac{\partial^2 \log L(\theta; y)}{\partial \theta \partial \theta'} \Big|_{\hat{\theta}_{MLE}} \text{ is negative definite matrix}$$

Usually hard to analytically determine (solve for)  $\hat{\theta}_{MLE}$

Need to do gradient ascent to find local (global if globally concave) maximum

## TERMS

Score vector/gradient is a  $k \times 1$  vector of first derivatives  $S(\theta) = \sum_{i=1}^n \frac{\partial \log L_i(\theta)}{\partial \theta} = \sum_{i=1}^n S_i(\theta)$  evaluated at the true  $\theta$

Hessian  $H$  is  $k \times k$  matrix of second derivatives  $H(\theta) = \sum_{i=1}^n \frac{\partial^2 \log L_i(\theta)}{\partial \theta \partial \theta'} = \sum_{i=1}^n H_i(\theta) = \begin{pmatrix} \frac{\partial^2 \log L(\theta)}{\partial \theta_1 \partial \theta_1} & \cdots & \frac{\partial^2 \log L(\theta)}{\partial \theta_1 \partial \theta_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 \log L(\theta)}{\partial \theta_k \partial \theta_1} & \cdots & \frac{\partial^2 \log L(\theta)}{\partial \theta_k \partial \theta_k} \end{pmatrix}$  is symmetric

Information matrix ( $k \times k$ )  $I(\theta) = -E(H(\theta))$  negative definite at  $\hat{\theta}_{MLE}$

## PROPERTIES OF MLE ESTIMATOR

Assuming certain regularity conditions are met

**Consistent** but may be biased

**Asymptotically Normal** enables hypothesis testing

- $\hat{\theta} \stackrel{\text{d}}{\sim} N(\theta, I(\theta)^{-1})$ , where  $I(\theta) = -E(H(\theta))$

**Asymptotically Efficient**, so estimates are precise and tests using them are powerful

- $I(\theta)^{-1}$  provides a lower bound on the asymptotic covariance matrix for any consistent, asymptotically normal estimator for  $\theta$  (Cramer-Rao lower bound)

★ The asymptotic variance can be estimated by  $\widehat{AVar}(\hat{\theta}) = -H(\hat{\theta})^{-1}$  or  $\widehat{AVar}(\hat{\theta}) = \left( \sum_{i=1}^n S_i(\hat{\theta}) S_i(\hat{\theta})' \right)^{-1}$

$H(\hat{\theta})$  is symmetric  
better properties in small samples  
outer product of first derivatives easier than computing  $H$

## MLE: BINARY CHOICE

$$f(y) = \begin{cases} \theta^y (1-\theta)^{1-y} & y=0,1 \\ 0 & \text{otherwise} \end{cases}$$

$$L(\theta; Y) = f(Y_1, \dots, Y_n; \theta) = \prod_{i=1}^n f(Y_i) \text{ assuming independence} = \theta^{\sum_{i=1}^n Y_i} (1-\theta)^{n-\sum_{i=1}^n Y_i}$$

$$\log L(\theta; Y) = \sum_{i=1}^n Y_i \log(\theta) + (n - \sum_{i=1}^n Y_i) \log(1-\theta)$$

$$\text{FOC: } \frac{\partial \log L}{\partial \theta} = \frac{\sum_{i=1}^n Y_i}{\theta} = \frac{n - \sum_{i=1}^n Y_i}{1-\theta} = 0 \Rightarrow \hat{\theta}_{MLE} = \frac{\sum_{i=1}^n Y_i}{n} = \bar{Y}$$

$$\text{Consistent: } \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{\text{law}} E(Y_i) = 0$$

$$\text{Unbiased: } E(\hat{\theta}_{MLE}) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \frac{1}{n} n\theta = \theta$$

$$\text{Asymptotically normal } \hat{\theta}_{MLE} \stackrel{\text{d}}{\sim} N(\theta, \frac{\theta(1-\theta)}{n}) \text{ so } SE(\hat{\theta}_{MLE}) = \sqrt{\frac{\hat{\theta}_{MLE}(1-\hat{\theta}_{MLE})}{n}}$$

To show: obtain  $I(\theta)$ , or look at properties of  $\bar{Y}$

# MLE: LINEAR REGRESSION

$$y_i = \mathbf{x}_i' \beta + u_i, \quad u_i | \mathbf{x}_i \sim \text{iid } N(0, \sigma^2), \quad i=1, \dots, n$$

$$y_i \sim N(\mathbf{x}_i' \beta, \sigma^2)$$

assume  $\mathbf{X}$  fixed  
(if not use conditional)

$$f(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i' \beta)^2}$$

$$L(\beta, \sigma^2) = f(y_1, \dots, y_n) \stackrel{\text{indep}}{=} \prod_{i=1}^n f(y_i) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2}$$

$$\log L(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta)$$

$$\text{FOCs: } \begin{cases} \frac{\partial \log L}{\partial \beta} = -\frac{1}{2\sigma^2} (-2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\beta) \\ \frac{\partial \log L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \end{cases}$$

$$\hat{\beta}_{\text{MLE}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{\mathbf{u}'\mathbf{u}}{n}, \quad \mathbf{u} = \mathbf{y} - \mathbf{X}\hat{\beta}_{\text{MLE}}$$

$$\begin{pmatrix} \hat{\beta}_{\text{MLE}} \\ \hat{\sigma}_{\text{MLE}}^2 \end{pmatrix} \stackrel{\text{?}}{\sim} N\left(\begin{pmatrix} \beta \\ \sigma^2 \end{pmatrix}, \begin{pmatrix} \sigma^2(\mathbf{X}'\mathbf{X})^{-1} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix}\right)$$

Calculate Hessian (W10 slides) to get  $\mathbb{I} \Rightarrow \mathbb{I}^{-1}$

Zero covariance + normality  $\Rightarrow$  Independence

$$\hat{\beta}_{\text{MLE}} \stackrel{\text{?}}{\sim} N(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}) \quad \text{OLS derivation tells us this is in fact true for any } n$$

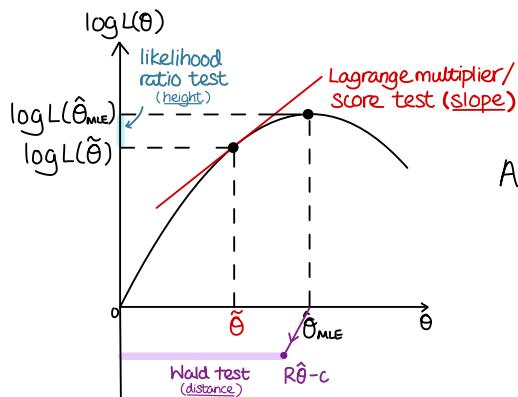
$$\hat{\sigma}_{\text{MLE}}^2 \stackrel{\text{?}}{\sim} N(\sigma^2, \frac{2\sigma^4}{n})$$

# TRINITY OF CLASSICAL TESTING

$$H_0: R\theta = c, \theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_k \end{pmatrix}$$

$\hat{\theta}$  is the **unrestricted** (in this case MLE) estimator of  $\theta$

$\tilde{\theta}$  is the **restricted** (in this case MLE) estimator of  $\theta$ , where  $R\tilde{\theta} = c$



All three tests are asymptotically equivalent

- Often chosen based on computational ease
- May have different finite (small) sample properties

## WALD TEST

Is the discrepancy  $d = R\hat{\beta} - c$  close to 0?

- Uses only  $\hat{\theta}$ , not  $\theta$
- Compared to Wald test in lin. reg. case:
  - Need to assume MLE regularity conditions
  - Need to use  $\hat{\beta}_{MLE} \stackrel{\text{asymptotic}}{\sim} N(\beta, V)$ , and  $V$  obtained differently, using  $(I(\beta))^{-1}$
  - Test statistic valid only asymptotically, as we don't have A1-5

## TEST STATISTIC

$$\text{Single linear restriction: } Z = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \xrightarrow{d} N(0,1) \text{ under } H_0$$

Called z-test (asymptotic) rather than t-test (finite)  
since  $\hat{\theta}_{MLE}$  are usually asymptotically normal

$$\text{Joint linear restriction: } W = d' (\widehat{\text{Var}}(d))^{-1} d \xrightarrow{d} \chi_r^2 \text{ under } H_0$$

$\widehat{\text{Var}}(d)$  or  $\text{Var}(d)$  don't matter. Both have same limiting distribution

- Since  $d$  is a vector, we want to check length of vector  $d'd$
- We don't know how  $d'd$  is distributed,  
but we know how the Mahalanobis distance  $d'(\widehat{\text{Var}}(d))^{-1}d$  is distributed

## LIKELIHOOD RATIO TEST

Does imposing restriction on estimator lead to a significant loss of fit similar to RRSS/URSS

- Is  $\log L(\hat{\theta}) - \log(\tilde{\theta})$  significantly different from 0
- Easy to apply, but need both  $\hat{\theta}, \tilde{\theta}$

## TEST STATISTIC

$$\begin{array}{c} \text{unrestricted} \quad \text{restricted} \\ \downarrow \quad \downarrow \\ 2(\log L(\hat{\theta}) - \log L(\tilde{\theta})) \xrightarrow{d} \chi_r^2 \text{ under } H_0 \end{array}$$

↑  
no. of restrictions

If  $H_0$  is true, likelihood ratio  $\frac{\log L(\tilde{\theta})}{\log L(\hat{\theta})} \approx 1 \Leftrightarrow \log L(\tilde{\theta}) - \log L(\hat{\theta}) \approx 0$

## LAGRANGE MULTIPLIER TEST

$$\max \log L(\theta) \text{ s.t. } R\theta = c$$

↓  
imposes restriction!

$$\mathcal{L} = \log L(\theta) + \lambda'(R\theta - c)$$

- $\lambda$  are **shadow prices**: how expensive it is to impose restrictions
- Reject restrictions ( $H_0$ ) if shadow prices are too high (far from 0)
- Uses only  $\tilde{\theta}$ , not  $\hat{\theta}$

## TEST STATISTIC

$$\tilde{\lambda}' (\widehat{\text{Var}}(\tilde{\lambda}))^{-1} \tilde{\lambda} \xrightarrow{d} \chi_r^2 \quad \text{under } H_0$$

or equivalently,

$$S(\tilde{\theta})' (\widehat{\text{Var}}(S(\tilde{\theta})))^{-1} S(\tilde{\theta}) \xrightarrow{d} \chi_r^2 \quad \text{under } H_0$$

$$S(\tilde{\theta}) = \frac{\partial \log L(\theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial \log L_i(\theta)}{\partial \theta} \Big|_{\tilde{\theta}} = \sum_{i=1}^n S_i(\tilde{\theta})$$

the latter can be rewritten as  $\left[ \sum_{i=1}^n S_i(\tilde{\theta}) \right]' \left[ \sum_{i=1}^n S_i(\tilde{\theta}) S_i(\tilde{\theta})' \right]^{-1} \left[ \sum_{i=1}^n S_i(\tilde{\theta}) \right]$  because

$$\text{Var}\left(\frac{\partial \log L(\theta)}{\partial \theta}\right) = E \left\{ \left[ \frac{\partial \log L}{\partial \theta} - E\left(\frac{\partial \log L}{\partial \theta}\right) \right] \left[ \frac{\partial \log L}{\partial \theta} - E\left(\frac{\partial \log L}{\partial \theta}\right) \right]' \right\} = E \left[ \left( \frac{\partial \log L}{\partial \theta} \right) \left( \frac{\partial \log L}{\partial \theta} \right)' \right]$$

$= 0$  at  $\tilde{\theta}_{MLE}$

which can be estimated by  $\sum_{i=1}^n S_i(\tilde{\theta}) S_i(\tilde{\theta})'$

notation is different  
from MLE chapter

→ This can be obtained by computing the  $nR^2$  of regressing  $l_i = S_i(\tilde{\theta})' \gamma + v_i$ ,  $i=1, \dots, n$  where  $S_i(\tilde{\theta}) = \begin{pmatrix} 1 & S_i(\tilde{\theta})' \\ \vdots & \vdots \\ 1 & S_n(\tilde{\theta})' \end{pmatrix}$

or equivalently  $\tau = S\hat{\gamma} + v$  in matrix form, where  $S = \begin{pmatrix} 1 & S_1(\tilde{\theta})' \\ \vdots & \vdots \\ 1 & S_n(\tilde{\theta})' \end{pmatrix}$ . This gives  $\hat{\gamma} = (SS')^{-1}Sv$ .

It works because  $nR^2 = n \frac{ESS}{TSS} = n \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} = n \frac{\hat{y}'\hat{y}}{n} = (S\hat{\gamma})'(S\hat{\gamma}) = \hat{\gamma}'S'S\hat{\gamma} = \hat{\gamma}'S(S'S)^{-1}S\hat{\gamma}$

is equivalent to  $\left[ \sum_{i=1}^n S_i(\tilde{\theta}) \right]' \left[ \sum_{i=1}^n S_i(\tilde{\theta}) S_i(\tilde{\theta})' \right]^{-1} \left[ \sum_{i=1}^n S_i(\tilde{\theta}) \right]$  instead of the usual  $\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$  which equals 0  
i.e.  $R^2$  is uncentred

# BINARY CHOICE: LPM/NLLS

**LINEAR PROBABILITY MODEL**  $y_i = \mathbf{x}_i' \beta + \varepsilon_i$ ,  $E(\varepsilon_i | \mathbf{x}_i) = 0$  where  $y_i \in \{0, 1\}$

Benefits:

- Since  $E(\varepsilon_i | \mathbf{x}_i) = 0$  and  $y_i$  binary,  $E(y_i | \mathbf{x}_i) = 1 \times \Pr(y_i=1 | \mathbf{x}_i) + 0 \times \Pr(y_i=0 | \mathbf{x}_i) = \Pr(y_i=1 | \mathbf{x}_i) = \mathbf{x}_i' \beta$
- $\Pr(y_i=1 | \mathbf{x}_i)$  linear in the parameters; easy to interpret  $\beta_j$
- $\beta_j$ : How does the  $j^{\text{th}}$  explanatory variable affect probability that  $y=1$ , CP
- Constant marginal effects

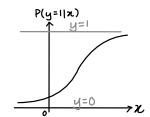
Drawback:

- $\mathbf{x}_i' \hat{\beta}$  is not necessarily constrained to lie in  $[0, 1]$
- Error term  $\varepsilon_i = y_i - \mathbf{x}_i' \beta$  can only take 2 values  $1 - \mathbf{x}_i' \beta$  or  $-\mathbf{x}_i' \beta$ ; highly non-normal
- Heteroscedasticity.  $\text{Var}(\varepsilon_i | \mathbf{x}_i) = \text{Var}(y_i | \mathbf{x}_i) = \Pr(y_i=1 | \mathbf{x}_i)(1 - \Pr(y_i=1 | \mathbf{x}_i)) = \mathbf{x}_i' \beta(1 - \mathbf{x}_i' \beta)$  is not constant
  - Must use RSE

$$\begin{aligned}\text{Var}(y) &= E(y^2) - [E(y)]^2 \\ E(y) &= E(y^2) = \Pr(y=1)\end{aligned}$$

Thus, we may decide to model  $P(y=1 | \mathbf{x}) = F(\mathbf{x}' \beta)$ , where  $F(\cdot)$  is the CDF

Usually  $F(\mathbf{x}' \beta) = \Phi(\mathbf{x}' \beta) = \int_{-\infty}^{\mathbf{x}' \beta} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2} du$  probit or  $\Lambda(\mathbf{x}' \beta) = \frac{e^{\mathbf{x}' \beta}}{1+e^{\mathbf{x}' \beta}}$  logit



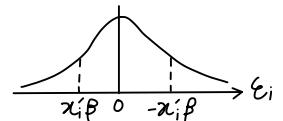
CDF of logistic RV

**WHY USE CDF?** Intuitive explanation using behavioural model

Let  $w_i$  be the excess utility of working (vs not working) s.t.  $w_i = \mathbf{x}_i' \beta + \varepsilon_i$ .

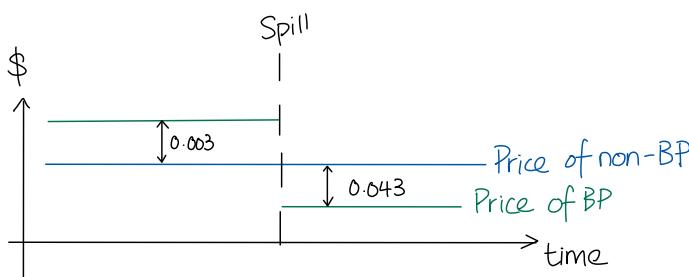
We don't observe  $w_i$  but  $y_i = 1$  if  $w_i \geq 0$ ,  $y_i = 0$  if  $w_i < 0$  therefore

$$\begin{aligned}\Pr(y_i=1 | \mathbf{x}_i) &= \Pr(w_i > 0 | \mathbf{x}_i) = \Pr(\mathbf{x}_i' \beta + \varepsilon_i > 0 | \mathbf{x}_i) \\ &= \Pr(\varepsilon_i > -\mathbf{x}_i' \beta | \mathbf{x}_i) \\ &= \Pr(\varepsilon_i < \mathbf{x}_i' \beta | \mathbf{x}_i) \text{ given } \varepsilon_i \text{'s distribution (if symmetric)} \\ &= F_\varepsilon(\mathbf{x}_i' \beta) \text{ depends on pdf of } \varepsilon\end{aligned}$$

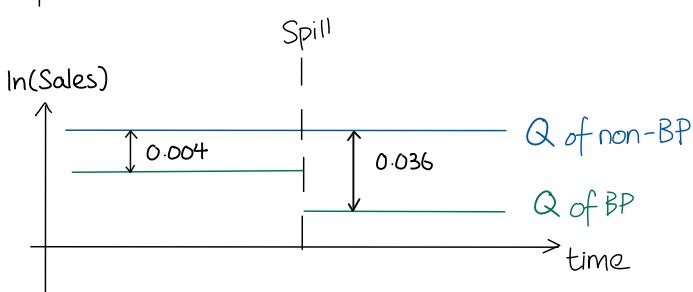


## NON-LINEAR REGRESSION

- We can run  $y_i = F(\mathbf{x}_i' \beta) + \varepsilon_i$  to get a non-linear least squares estimator  $\hat{\beta}_{NLS} = \underset{b}{\operatorname{argmin}} \sum_{i=1}^n (y_i - F(\mathbf{x}_i' b))^2$  usually no explicit form
- But doesn't solve heteroscedasticity (use RSEs). So next use MLE: asymptotically efficient



$$\Delta P = P_{\text{after}} - P_{\text{before}} = -0.046$$



$$\begin{aligned}\Delta \ln Q &= e^{\ln Q_{\text{after}}} - e^{\ln Q_{\text{before}}} = \frac{1}{Q} \\ \Delta \ln Q &= \beta \\ \Delta \ln P &= \beta\end{aligned}$$

# BINARY CHOICE: MLE

$$y_i: \text{Binary / Bernoulli RV} \quad f(y_i) = \begin{cases} \theta_i^{y_i} (1-\theta_i)^{1-y_i} \\ 0 \end{cases}$$

$\theta_i = \Pr(y_i=1|x_i)$  should differ from individual to individual

We can specify (model) a relation between  $\theta_i$  and  $x_i$  using a CDF

$$\theta_i = F(x_i'\beta) = \Phi(x_i'\beta) \text{ probit or } \frac{e^{x_i'\beta}}{1+e^{x_i'\beta}} \text{ logit}$$

## PROBIT

$$L(\beta) = \prod_{i=1}^n \Phi(x_i'\beta)^{y_i} (1 - \Phi(x_i'\beta))^{1-y_i}$$

$$\log L(\beta) = \sum_{i=1}^n y_i \log \Phi(x_i'\beta) + (1-y_i) \log [1 - \Phi(x_i'\beta)]$$

$$\text{FOC: } \left. \frac{\partial \ln L(\beta)}{\partial \beta} \right|_{\hat{\beta}_{MLE}} = \sum_{i=1}^n \hat{\epsilon}_i^g x_i = 0 \quad \text{PSQ3}$$

$$\text{where generalised residuals } \hat{\epsilon}_i^g = \frac{y_i - \Phi(x_i'\hat{\beta}_{MLE})}{\Phi(x_i'\hat{\beta}_{MLE})[1 - \Phi(x_i'\hat{\beta}_{MLE})]} \downarrow \Phi(x_i'\hat{\beta}_{MLE}) \text{ PDF}$$

However, FOC gives no explicit form for  $\hat{\beta}_{MLE}$  (may need gradient ascent)

↑ "the parameter iteratively converging...  
globally concave obj fn..."

## INTERPRET

Predicted probability:  $\Pr(\widehat{y_i}=1|x_i) = \Phi(x_i'\hat{\beta}_{MLE})$

Intuitively,  $\hat{\beta}_{MLE}$  is chosen s.t. individuals with  $y_i=1$  have high  $\Pr(\widehat{y_i}=1|x_i) = \Phi(x_i'\hat{\beta}_{MLE})$   
individuals with  $y_i=0$  have high  $\Pr(\widehat{y_i}=0|x_i) = 1 - \Phi(x_i'\hat{\beta}_{MLE})$

Marginal effect (if  $x_j$  is continuous):  $\frac{\partial \Phi(x'\beta)}{\partial x_j} = \phi(x'\beta)\beta_j \quad \text{not } \beta_j \text{ (unlike OLS)}$

- Need "ceteris paribus"
- Not constant: Depends on individual characteristics  $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix}$ , though sign of  $\beta_j$  still informative
- May want to evaluate at the means  $\phi(\bar{x}'\beta)\beta_m$ . (partial effect on average)
  - That said, individual  $\bar{x}$  may not exist, esp for dummy variable ("half female")
  - Stata command: margins, dydx (\*) atmeans
- May want to average over all individuals  $\frac{1}{n} \sum_i \phi(x_i'\beta)\beta_j$  (partial average treatment effect)
  - "On average, increasing  $x_j$  by 1 increases probability of  $y_i=1$  by ..., ceteris paribus"
  - Stata command: margins, dydx (\*)

Marginal effect (if  $x_j$  is a dummy):  $\frac{\Delta \Pr(y_i=1|x)}{\Delta x_j} = \Phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots) - \Phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots)$   
 $= \Pr(y_i=1|x_1=1, x_2, \dots) - \Pr(y_i=1|x_1=0, x_2, \dots)$

$$x + \varepsilon$$

# BINARY CHOICE: MLE TEST

Same as generic cases but with additional pointers for MLE binary choice

$$\text{example } H_0: \beta_j = 0$$

## WALD TEST

$$Z = \frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)} \stackrel{a}{\sim} N(0, 1) \text{ under } H_0, \text{ reject if } |z| > z_{\frac{\alpha}{2}}$$

- $\text{SE}(\hat{\beta}_j)$  calculated by taking square root of the  $j,j^{\text{th}}$  element of  $I(\hat{\beta})^{-1} = \left(-\frac{\partial^2 \ln L(\beta)}{\partial \beta \partial \beta}\Big|_{\hat{\beta}_{\text{MLE}}}\right)^{-1}$   
might be  $j+1, j+1$  if intercept  $\beta_0$
- Recall: Under suitable regularity conditions,  $\hat{\beta}_{\text{MLE}} \stackrel{a}{\sim} N(\beta, I(\beta)^{-1})$  so  $\hat{\beta}_j \stackrel{a}{\sim} N(\beta_j, [I(\hat{\beta})]_{jj})$
- Not using t-dist to account for imprecision as this test is only valid asymptotically  
and asymptotically it doesn't matter if we know the true SE or use a consistent estimator for it

## LIKELIHOOD RATIO TEST

$$2(\log L^U - \log L^R) \xrightarrow{a} \chi^2_k \text{ under } H_0, \quad \text{reject if } > \chi^2_{1,\alpha}$$

↑ no. of restrictions  
unconstrained probit/logit

"difference in dimensionality" (df of unrestricted model - df of restricted)

Stata commands (for testing  $\beta_j = 1$ ):

constr 1 variable\_j=1 set 1 restriction

probit depvar variable\_1 variable\_2 variable\_j ..., constr(1) estimate probit n constraint

## LAGRANGE MULTIPLIER TEST

$$nR^2 \sim \chi^2_k \text{ under } H_0 \quad \text{reject if } > \chi^2_{1,\alpha}$$

Uncentred  $R^2$  obtained from regressing  $l_i = (\hat{\epsilon}_i^G)\gamma_1 + (\hat{\epsilon}_i^G x_{i1})\gamma_2 + \dots + v_i$

$$nR^2 = n \frac{\text{ESS}}{\text{TSS}} = n \frac{\sum_i \hat{y}_i^2}{\sum_i y_i^2} = n \frac{\hat{y}'\hat{y}}{n} = (S\hat{\beta})'(S\hat{\beta}) = \hat{\beta}'S'S\hat{\beta} = \hat{\beta}'S(S')^{-1}S\hat{\beta}$$

Score vector in probit

In Stata probit output, joint significance of the slopes is automatically given as "LR chi2(k)"  
where k is the number of regressions in the top right hand corner

Test statistic will be small/derivative will be close to zero if the restrictions imposed were true. It is an indication that imposing the restrictions on the model is not unreasonable or costly.

## CHOW TEST

rural/urban, male/female

"Are slopes/intercepts the same for a different dataset/Segment of the population?"

Conduct likelihood ratio test.  $H_0: \beta_j^{\text{rural}} = \beta_j^{\text{urban}} \forall j$  vs  $H_1: \text{at least one } \beta_j^{\text{rural}} \neq \beta_j^{\text{urban}}$

Restricted model:

probit y x

Then take the log likelihood directly from the Stata output to get  $\log L^R$

Unrestricted model:

probit y x if urban == 1

probit y x if urban == 0

Then take the log likelihood from both outputs and add them up to get  $\log L^U$

$$2(\log L^U - \log L^R) \xrightarrow{a} \chi^2_r$$

↑ number of variables (don't forget intercept)

"difference in dimensionality" (df of unrestricted model - df of restricted =  $2k - k = k$ )

# COUNT VARIABLE

$y_i = 0, 1, 2, \dots$  (nonnegative integer values)

$$f(y_i) = \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!}, \quad y_i = 0, 1, 2, 3, \dots$$

where  $\lambda_i = E(y_i) = \text{Var}(y_i)$   
 $\uparrow \lambda_i \text{ not constant. dep on } x_i$

- OLS has similar drawbacks as before (negative  $E(y_i|x_i)$ , heteroscedasticity)
- Weighted NLLS more efficient than NLLS due to heteroscedasticity; MLE even more efficient

We can specify (model) a relation between  $\lambda_i$  and  $x_i$  using an exponential (to ensure  $\lambda_i \geq 0$ )

$$\lambda_i = E(y_i|x_i) = e^{x_i \beta}$$

## MLE ESTIMATION

$$\ln L(\beta) \stackrel{\text{indep}}{=} \sum_{i=1}^n \ln \left( \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} \right) = \sum_{i=1}^n \left[ -\lambda_i + y_i \ln \lambda_i - \ln(y_i!) \right] = \sum_{i=1}^n \left[ -x_i \beta + y_i \ln x_i \beta - \ln(y_i!) \right]$$

FOC:  $\frac{\partial \ln L(\beta)}{\partial \beta} \Big|_{\hat{\beta}_{MLE}} = \sum_{i=1}^n \hat{\epsilon}_i^R x_i = 0$  gives no explicit form for  $\hat{\beta}_{MLE}$  (may need gradient ascent)

Where residuals  $\hat{\epsilon}_i^R = y_i - e^{-x_i \hat{\beta}_{MLE}}$

- $SE(\hat{\beta}_j)$  calculated by taking square root of the  $j, j^{\text{th}}$  element of  $I(\hat{\beta})^{-1} = \left( -\frac{\partial^2 \ln L(\beta)}{\partial \beta \partial \beta} \Big|_{\hat{\beta}_{MLE}} \right)^{-1}$   
 $\uparrow$  might be  $j+1, j+1$  if intercept  $\beta_0$
- Recall: Under suitable regularity conditions,  $\hat{\beta}_{MLE} \stackrel{\text{d}}{\sim} N(\beta, I(\beta)^{-1})$  so  $\hat{\beta}_j \stackrel{\text{d}}{\sim} N(\beta_j, [I(\hat{\beta})^{-1}]_{jj})$

$$\text{Marginal effect: } \frac{\partial E(y|x)}{\partial x_k} = \frac{\partial e^{x \beta}}{\partial x_k} = \beta_k e^{x \beta}$$

"How does expected  $y_i$  change due to  $x_k$ , CP"