

# Probability & statistics

## ① Random variables

### (expectation)

$$E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx = \sum x \cdot f(x)$$

$$E(g(X)) = \int_{-\infty}^{\infty} g(x) \cdot f(x) dx = \sum g(x) \cdot f(x)$$

$$E(ax + b) = a E(X) + b$$

$$E(X+Y) = E(X) + E(Y)$$

### (variance)

$$\begin{aligned} V(X) &= E(X^2) - [E(X)]^2 \\ &= \sum (x_i - m_x)^2 f(x) \\ &= \int_{-\infty}^{\infty} (x - m_x)^2 f(x) dx \end{aligned}$$

$$V(ax + b) = a^2 V(X)$$

$$V(ax \pm by) = a^2 V(X) + b^2 V(Y) + 2ab \text{cov}(X, Y)$$

## ② Sampling

(bias)  $E(\hat{\theta}) - \theta \Rightarrow \text{bias} = 0 \text{ then unbiased}$

(variance)  $E(\hat{\theta}^2) - E(\hat{\theta})^2 \Rightarrow \text{smaller then more efficient}$

(MSE)  $\text{MSE} = \text{bias}^2 + \text{variance}$

(consistency)  $\Rightarrow \text{as } n \rightarrow \infty, \hat{\theta} \rightarrow \theta$

$\hat{\theta}$  is consistent iff  $\lim_{n \rightarrow \infty} P(\hat{\theta} - \theta = 0) = 1$

## ④ 2) Random variables

### (covariance)

$$\text{cov}(X, Y) = E(XY) - m_X m_Y$$

$$\text{cov}(ax + b, cy + d) = ac \text{cov}(X, Y)$$

$X, Y$  independent  $\rightarrow \text{cov}(X, Y) = 0$ , not converse

### (correlation)

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sqrt{V(X)} \sqrt{V(Y)}} = \frac{r_{xy}}{s_x s_y}$$

$$\hat{\rho}_{X,Y} = \frac{s_{xy}}{s_x s_y} = \frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum (x_i - \bar{x})^2 \cdot \frac{1}{n-1} \sum (y_i - \bar{y})^2}$$

## ③ distribution of $\bar{X}$

### (expectation & variance)

under SRS,  $X_i$  is i.i.d.  $\bar{X} = \frac{1}{n} \sum X_i$

$$E(\bar{X}) = m_x$$

$$V(\bar{X}) = \frac{\sigma^2}{n}$$

$$S_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$

$$SE(\bar{X}) = \frac{s_x}{\sqrt{n}}$$

### (central limit theorem)

As  $n \rightarrow \infty$ ,  $\bar{X} \rightarrow N(m_x, \frac{\sigma^2}{n})$

$n \geq 100$ ,  $\bar{X} \sim N(m_x, \frac{\sigma^2}{n})$  approximately

### (law of large numbers)

As  $n \rightarrow \infty$ ,  $P(\bar{X} - m_x = 0) \rightarrow 1$

$\hookrightarrow \bar{X} \rightarrow m_x$

### (hypothesis testing)

1.  $H_0, H_1$

2.  $\bar{X} \sim N(m_x, H_0, \frac{\sigma^2}{n})$  approximately,

$$\bar{X} \sim T_{n-1}(m_x, H_0, \frac{s_x^2}{n})$$

3. test statistic =  $\frac{\bar{X} - m_x, H_0}{\frac{\sigma}{\sqrt{n}}} \text{ or } \frac{\bar{X} - m_x, H_0}{\frac{s_x}{\sqrt{n}} t_{n-1}}$

4. calculate p-value

90%, 95%, 99%  
1.64, 1.96, 2.58

### (confidence interval)

$$\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \text{ or } \bar{X} \pm t_{n-1, \alpha/2} SE(\bar{X})$$

$\Rightarrow$  in the long run,  $1 - \alpha$  % of all samples will give an interval that contains the true  $m$

## linear regression

### ① OLS estimation of bivariate linear regression

#### (linear model)

$$E(Y) = E(E(Y|X=x))$$

$$E(Y_i|X_i) = \beta_0 + \beta_1 X_i$$

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

#### OLS estimation

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i, \quad \hat{u}_i = Y_i - \hat{Y}_i$$

$$\text{min. } \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2$$

$$\frac{\partial Q}{\partial \hat{\beta}_0} = -2 \sum Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i = 0$$

$$\frac{\partial Q}{\partial \hat{\beta}_1} = -2 \sum X_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0$$

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

$$\hat{\beta}_1 = \frac{S_{xy}}{S_x^2} = \frac{\frac{1}{n-1} \sum (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n-1} \sum (X_i - \bar{X})^2}$$

$$\frac{1}{n} \sum \hat{u}_i = 0$$

$$\frac{1}{n} \sum \hat{Y}_i = \bar{Y}, \quad \bar{Y} \text{ lies on line}$$

$$\sum u_i X_i = 0$$

$$\text{cov}(\hat{u} | X) = 0$$

### ④ studying vs: measures of fit

**R<sup>2</sup>** measures fraction of sample variance explained by model. mult-less & directionless.

$$Y_i = \hat{Y}_i + \hat{u}_i$$

predicted variance

$$R^2 = \frac{\text{explained sum of squares}}{\text{total sum of squares}} = \frac{\frac{1}{n-1} \sum (\hat{Y}_i - \bar{Y})^2}{\frac{1}{n-1} \sum (Y_i - \bar{Y})^2}$$

$$= 1 - \frac{\text{sum of squared residuals}}{\text{total sum of squares}} = \frac{\sum \hat{u}_i^2}{\sum (Y_i - \bar{Y})^2}$$

actual variance

fraction of variance not explained by model

$$= \frac{\text{var}(X, Y)}{\text{Var}(X) \text{Var}(Y)} = \frac{S_{xy}}{S_x S_y} \quad ESS + SSR = TSS$$

#### ② sampling distributions

$$E(\hat{\beta}_1) = \beta_1$$

(larger  $V(X)$ , smaller  $V(\hat{\beta}_1)$ )

$$V(\hat{\beta}_1) = \frac{1}{n} \frac{V((X_i - \bar{X}) u_i)}{V(X_i)^2} = \frac{\frac{1}{n-2} \sum (X_i - \bar{X})^2 \hat{u}_i^2}{\frac{1}{n} (X_i - \bar{X})^2}$$

#### (Assumptions)

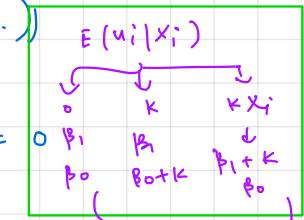
- $E(u_i | X = x) = 0$
- $(X_i, Y_i)$  are independently & identically distributed
- large outliers are rare

#### (bias)

$$E(\hat{\beta}_1 - \beta_1 | X_1, X_2, \dots, X_n) = E\left(\frac{\sum (X_i - \bar{X}) u_i}{\sum (X_i - \bar{X})^2} | X_1, X_2, \dots\right)$$

$$= \frac{\sum (X_i - \bar{X}) E(u_i | X_1, X_2, \dots)}{\sum (X_i - \bar{X})^2} \quad \text{I.I.D.}$$

$$= \frac{\sum (X_i - \bar{X}) E(u_i | X_i)}{\sum (X_i - \bar{X})^2}$$



#### (consistency)

$$E(\hat{\beta}_1) = \beta_1 \quad \& \quad V(\hat{\beta}_1) \propto \frac{1}{n} \Rightarrow \lim_{n \rightarrow \infty} \hat{\beta}_1 \rightarrow \beta_1$$

#### ③ homo & heteroskedasticity

$$\hookrightarrow V(u_i | X_i = x) = \frac{\sigma_u^2}{\text{constant}} \quad \text{v}(u_i | X_i = x) \neq k$$

iff homoskedastic:

$$V(\hat{\beta}_1) = \frac{1}{n} \frac{\sigma_u^2}{\sigma_x^2} = \frac{1}{n} \cdot \frac{\frac{1}{n-2} \sum \hat{u}_i^2}{\frac{1}{n} \sum (X_i - \bar{X})^2}$$

heteroskedastic robust and homoskedastic estimators of  $V(\hat{\beta}_1)$  converge as  $n \rightarrow \infty$

**standard error of regression** measure of spread of observations about regression line wrt to  $\gamma$   
ie. sample st. dev of OLS residuals

$$SER = \sqrt{\frac{1}{n-2} \sum (\hat{u}_i - \bar{u}_i)^2} = \sqrt{\frac{1}{n-2} \sum \hat{u}_i^2}$$

$\hookrightarrow = 0$  due to OLS

#### (root mean squared error)

$$RMSE = \sqrt{\frac{1}{n} \sum \hat{u}_i^2}$$

## ⑤ studying vs: hypothesis testing

$$1. H_0: \beta_1 = k, H_1: \beta_1 \neq k$$

2. compute  $SE(\hat{\beta}_1)$

$$\begin{aligned} SE(\hat{\beta}_1) &= \sqrt{\frac{\text{varc}(x_i - \bar{x})(u_i)}{\text{var}(x_i)^2}} \quad \text{scaling affects here} \\ &= \sqrt{\frac{\frac{1}{n-2} \sum (x_i - \bar{x})^2 u_i^2}{\left[ \frac{1}{n-2} \sum (x_i - \bar{x})^2 \right]}} \end{aligned}$$

$\Rightarrow$  take out  $k^2$  on var,  $k$  on  $\text{var}(\hat{\beta}_1)$

$\Rightarrow$  under  $H_0$ , since  $n \geq 100$ ,  
 $\hat{\beta}_1 \sim N(\beta_1, SE(\hat{\beta}_1))$  approximately

3. compute t-statistic

$$t = \frac{\hat{\beta}_1 - \beta_{1,H_0}}{SE(\hat{\beta}_1)} \sim \text{approx normal}$$

$$\begin{aligned} \hat{y}_i &= \hat{\beta}_1 x_i + \hat{\beta}_0 \quad \hat{u}_i = y_i - \hat{y}_i \\ \hat{\beta}_1 &= \frac{\hat{s}_{xy}}{\hat{s}_{xx}} = \frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

4. compute p-values

5. If two tailed, compute confidence interval

$$\hat{\beta}_1 \pm t_{n-1, \alpha/2} SE(\hat{\beta}_1) \approx \hat{\beta}_1 \pm z_{\alpha/2} \cdot SE(\hat{\beta}_1)$$

1.64, 1.96, 2.58  
90%, 95%, 99%

## ⑥ omitted variable bias

$$\begin{aligned} E(\hat{\beta}_1) &= \beta_1 + \underbrace{\beta_2 \rho_{xz} \cdot \frac{\sigma_z}{\sigma_x}}_{\text{omitted covariated in } \hat{u}_i} + E\left(\frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2}\right) \\ &= \beta_1 + E\left(\frac{\sum (x_i - \bar{x})(\beta_2 z_i + \eta_i)}{\sum (x_i - \bar{x})^2}\right) \\ &= \beta_1 + \underbrace{\rho_{xz} \cdot \frac{\sigma_u}{\sigma_x}}_{\text{coefficient}} \end{aligned}$$

## ⑦ unit change

$$\begin{aligned} Y \rightarrow aY &\Rightarrow \hat{\beta}_1 \rightarrow a\hat{\beta}_1, \hat{\beta}_0 \rightarrow a\hat{\beta}_0 \\ X \rightarrow bX &\Rightarrow \hat{\beta}_1 \rightarrow \frac{\hat{\beta}_1}{b}, \hat{\beta}_0 \text{ unchanged} \end{aligned}$$

## ⑧ interpretation of regression line

### (continuous variable)

$$\text{slope: } \frac{\Delta E(Y)}{\Delta X}$$

$$\text{intercept: } E(Y | X=0)$$

### (binary variable)

$\hat{\beta}_1$ : difference in  $E(Y)$  when  $X$  is present vs. not. (i.e. diff in avg.  $Y$ )

$$\hat{\beta}_0: E(Y | X=0)$$

### (predictions)

$x_i$  should be SRS from same distribution, and within range of data

### (when does it occur?)

- the omitted variable is correlated with the included regressor(s)
- the omitted variable is a determinant of the dependent variable

## multiple regression

### ① functional form & OLS

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \varepsilon_i$$

(OLS)

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots \Rightarrow E(Y | X_{1i}, X_{2i}, \dots)$$

$$\begin{aligned} \hat{u}_i &= Y_i - \hat{Y}_i \\ \delta &= \sum_{i=1}^n \hat{u}_i^2 \end{aligned} \quad \left\{ \begin{array}{l} \frac{\partial \delta}{\partial \beta_0} = 0 \\ \frac{\partial \delta}{\partial \beta_1} = 0 \end{array} \right\} \Rightarrow \text{solve}$$

### ② multicollinearity

(perfect multicollinearity)

one regressor is a perfect linear function of other regressors

$$\text{eg. } X_1 = kX_2$$

(imperfect multicollinearity)

regressors are said to exhibit imperfect multicollinearity when they are highly correlated

### ③ Assumptions of multi variable regression

1.  $E(u_i | X_{1i}, X_{2i}, X_{3i}, \dots) = 0$
2.  $(X_{1i}, X_{2i}, \dots)_j$  - each observation  $\rightarrow$  iid
3. larger outliers are not common
4. no perfect collinearity  $\rightarrow$  divide by 0

### ④ Estimator distribution

$\hookrightarrow$  same as bivariate case

(imperfect multicollinearity)

$\hookrightarrow$  when 2 or more regressors are highly correlated  $\rightarrow$  then during OLS estimation, in keeping all else constant  $\rightarrow V(X_j)$  very small.

$$V(\hat{\beta}_j) \propto \frac{1}{V(X_j)}$$

$\Rightarrow \hat{\beta}_j$  still unbiased & consistent, just larger  $SE(\hat{\beta}_j)$  since  $V(X_j)$  is smaller

### ⑤ Binary multiple regression

(dummy variable trap) if there are  $G$  binary categories and there are  $G$  regressors + intercept,  $\forall i$  due to perfect multicollinearity because  $\beta_0 \Leftrightarrow \beta_0(1) \Leftrightarrow X_{0,j} = 1$ ,  $1$  will always be scalar multiple of all other  $X_i$

$$\underbrace{(n-1 \text{ regressors} + \hat{\beta}_0)}_{\text{exclude } \hat{\beta}_0}$$

### ⑥ Hypothesis testing: F test

(joint hypothesis testing)

$$1. H_0: \underbrace{\beta_j = k_j}_{\text{restriction}} \wedge \beta_m = k_m \dots \stackrel{q}{\wedge} \text{restrictions}$$

$$H_1: \text{at least 1 of } q \text{ restrictions false}$$

2. calculate  $q$  t-statistics and sum them to get F test stat

$$F \sim \frac{\frac{X_m^2}{m}}{\frac{X_n^2}{n} + \hat{\beta}_j^2} \sim \chi^2_q \text{ approximately}$$

$$\text{t-stat} \quad \left| \sum_{j=1}^n \frac{\hat{\beta}_{j,0} - \hat{\beta}_j}{se(\hat{\beta}_j)} \right| \sim \frac{X_q}{q}$$

$$f = \frac{\frac{X_q^2}{q}}{\frac{X_n^2}{n}} \rightarrow \text{sample size}$$

(single restriction of multiple  $\hat{\beta}_j$ )  $\text{eg. } \hat{\beta}_1 = \hat{\beta}_2$

1. reorganize eqn so that  $(\hat{\beta}_1 - \hat{\beta}_2)$  becomes coefficient

2. t-test it

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$$

$$Y_i = \beta_0 + (\beta_1 - \beta_2) X_{1i} + \beta_2 (X_{1i} + X_{2i}) + \varepsilon_i$$

$$= \beta_0 + \gamma X_{1i} + \beta_2 V_i, \quad V = X_{1i} + X_{2i}$$

## ① measures of fit

$$SER = s_{\hat{y}} = \sqrt{s_y^2}$$

$$s_y^2 = \frac{1}{n-k-1} \sum (\hat{y}_i - \bar{y}_i)^2 = \frac{SSR}{n-k-1}$$

↑  
num. of non- $\beta_0$   
regressors

### adjusted $\bar{R}^2$

$$\begin{aligned} \bar{R}^2 &= 1 - \frac{n-1}{n-k-1} \frac{\sum \hat{u}_i^2}{\sum (y_i - \bar{y})^2} \\ &= 1 - \frac{n-1}{n-k-1} \frac{SSR}{TSS} = 1 - \frac{s_y^2}{s_{\hat{y}}^2} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \bar{R}^2 < R^2$$

## non-linear regression

### ① model & effect of $\Delta X$

$$y_j = f(\underbrace{\int(x_1, x_2, \dots)}_{\text{some non-linear functional form}} y_j) + u_j$$

### effect of $\Delta X$ on $y$

take the difference

$$\Delta \hat{y} = f(x_i + \Delta x_i, x_2 + \Delta x_2, \dots) - f(x_1, x_2, \dots)$$

### standard error of estimated effect

in the linear case,  $\hat{y} = \hat{\beta}_0 + \sum \hat{\beta}_i x_i$

$$\Rightarrow \Delta \hat{y} = \hat{\beta}_1 (\Delta x_i) \Rightarrow SE(\Delta \hat{y}) = SE(\hat{\beta}_1) \cdot \Delta x_i$$

in non-linear case, depends on functional form

$$\text{eg. } \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 \dots$$

$$\Rightarrow \Delta \hat{y} = \hat{\beta}_1 (\Delta x) + \hat{\beta}_2 ((x + \Delta x)^2 - x^2)$$

$$SE(\Delta \hat{y}) = SE(\hat{\beta}_1) \cdot \underline{\Delta x} + SE(\hat{\beta}_2) \underline{((x + \Delta x)^2 - x^2)}$$

### ② polynomial models

$$y_j = \beta_0 + \beta_1 x_{ij} + \beta_2 x_{ij}^2 \dots$$

$\Rightarrow$  we still treat as multiple regression, so multivariate OLS as per usual

$\Rightarrow$  when  $\Delta x \Rightarrow$  remember to sum effects of all  $\beta_j$  at that range difference

### ③ logarithmic models

#### 1) properties of logarithms

$$\ln\left(\frac{a}{b}\right) = \ln a - \ln b$$

$$\ln(ab) = \ln a + \ln b$$

$$\ln(a^b) = b \ln a$$

#### (logarithms & percentages)

$$\begin{aligned} \ln(x + \Delta x) - \ln(x) &= \ln\frac{x + \Delta x}{x} \\ &= \ln(1 + \frac{\Delta x}{x}) \\ &\approx \frac{\Delta x}{x} \text{ when } \Delta x \text{ is small} \end{aligned}$$

#### 2) Three log models

##### (linear-log model)

$$Y_j = \beta_0 + \sum \beta_i \ln(X_{ij}) + u_j$$

$$Y_i + \Delta Y_i = \beta_0 + \beta_1 \ln(X_i + \Delta X)$$

$$\Delta Y_i = \beta_1 [\ln(X_i + \Delta X) - \ln(X_i)] \approx \beta_1 \left(\frac{\Delta X_i}{X_i}\right) \approx \beta_1 (0.01) K$$

$\Rightarrow$  A  $K\%$  change in  $X$  predicts an average change in  $Y$  by  $\beta_1 \times K\%$ .

Case	Population Regression Model	Interpretation of $\beta_1$
I. linear-log	$Y_i = \beta_0 + \beta_1 \ln(X_i) + u_i$	A 1% change in $X$ is associated with a change in $Y$ of $0.01\beta_1$
II. log-linear	$\ln(Y_i) = \beta_0 + \beta_1 X_i + u_i$	A change in $X$ by one unit ( $\Delta X = 1$ ) is associated with a $100\beta_1\%$ change in $Y$
III. log-log	$\ln(Y_i) = \beta_0 + \beta_1 \ln(X_i) + u_i$	A 1% change in $X$ is associated with a $\beta_1\%$ change in $Y$

##### (log-linear model)

$$\ln(Y_j) = \beta_0 + \sum \beta_i X_{ij} + u_j$$

$$\ln(Y + \Delta Y) = \beta_0 + \beta_1 (X + \Delta X)$$

$$\ln(Y + \Delta Y) - \ln(Y) = \beta_1 \Delta X \approx \frac{\Delta Y}{Y}$$

$\Rightarrow$  A unit change in  $X$  is associated with a  $\beta_1\%$  change in  $Y$  on average.

##### (log-log model)

$$\ln(Y_j) = \beta_0 + \sum \beta_i \ln(X_{ij}) + u_j$$

$$\ln(Y + \Delta Y) = \beta_0 + \beta_1 \ln(X + \Delta X)$$

$$\ln(Y + \Delta Y) - \ln(Y) = \beta_1 [\ln(X + \Delta X) - \ln(X)]$$

$$\frac{\Delta Y}{Y} \approx \beta_1 \frac{\Delta X}{X}$$

$\Rightarrow$  a  $K\%$  change in  $X$  is associated with an average  $\beta_1 \cdot K\%$  change in  $Y$ .

##### (predicted values of $y$ when $Y$ is in log)

$\Rightarrow$  taking the exponent is biased

$$\text{eg. } \ln(y) = \beta_0 + \beta_1 x_i + u_i$$

$$\hookrightarrow y_i = e^{\beta_0 + \beta_1 x_i + u_i}$$

$$E(y_i | x_i) = E(e^{\beta_0 + \beta_1 x_i} | x_i) \cdot E(e^{u_i} | x_i)$$

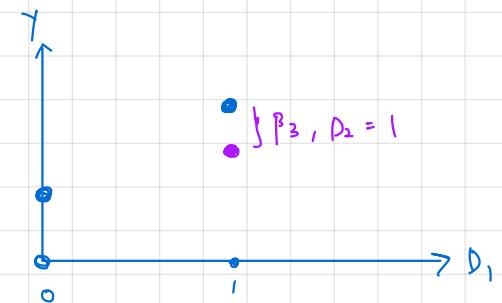
$\Rightarrow$  even if  $E(u_i | x_i) = 0$ ,  $E(e^{u_i} | x_i) \neq 1$  in general

scaled bias

#### ④ Interaction between variables

##### (binary-binary)

$$Y_i = \beta_0 + \beta_1 D_1 + \beta_2 D_2 + \beta_3 D_1 D_2 + u_i$$



$$E(Y_i | D_1 = 1, D_2 = d_2) = \beta_0 + \beta_1 + \beta_2 d_2 + \beta_3 d_2$$

$$E(Y_i | D_1 = 0, D_2 = d_2) = \beta_0 + \beta_2 d_2$$

Expected difference when  $D_1 = 1$  vs.  $D_1 = 0 \Rightarrow \beta_1 + \beta_3 d_2$

$\Rightarrow$  effect on  $Y$  of change in  $D_1$  now depends on  $D_2$

$\Rightarrow \beta_3$  is the difference in effect on  $Y$  of  $D_1$  when  $D_2 = 1$  vs.  $D_2 = 0$

##### (continuous-continuous)

$\hookrightarrow$  effectively finer granularity of binary-binary case

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \underbrace{\beta_3 X_1 X_2}_{\substack{\text{extra effect from} \\ \text{changing both } X_1 \text{ or } X_2}} \quad \hookrightarrow \text{effect of } X_1 \text{ holding } X_2 \text{ constant}$$

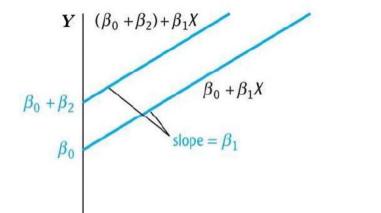
$$\frac{dy}{dX_1} = \beta_1 + \beta_3 X_2$$

$\Rightarrow \beta_3 \cdot X_2$  : the expected difference in effect of a unit change in  $X_1$  given  $X_2 = x_2$

##### (continuous-binary)

$\hookrightarrow$  1. test if they have diff. intercept i.e.  $\beta_2 = 0$

$\hookrightarrow$  inclusion of binary term allows for different offset



(a) Different intercepts, same slope

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i$$

$$D_i = 0$$

$$D_i = 1$$

$$y_i = \beta_0 + \beta_1 X_i + u_i \quad Y_i = (\beta_0 + \beta_2) + \beta_1 X_i + u_i$$

$\hookrightarrow$  inclusion of continuous-binary interaction allows the effect of change in  $X$  to be different (i.e. different gradient)  $\hookrightarrow$  2. test if they have diff gradient i.e.  $\beta_3 = 0$  (t-test)

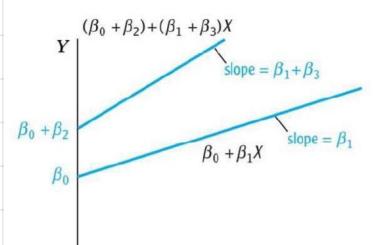
$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + \beta_3 X_i D_i + u_i$$

$$\text{if } D_i = 1$$

$$\text{if } D_i = 0$$

$$y_i = \beta_0 + (\beta_1 + \beta_3) X_i + \beta_2$$

$$Y_i = \beta_0 + \beta_1 X_i$$



(b) Different intercepts, different slopes

$\Rightarrow \beta_3$  - difference in effect of unit change in  $X$ , for observations where  $D_i = 1$  vs.  $D_i = 0$

## Assessing regression

### ① Validity of a study

(internal validity) for a study to be internally valid:

1. the estimator of the causal effect of interest should be unbiased and consistent
2. hypothesis tests have the correct significance level; standard errors are correctly estimated (hetero vs. homoscedasticity)

(external validity)

A study is externally valid if its inferences can be generalized to other populations & settings

### ② Threats to internal validity

Omitted variable bias  $E(u_i | x_i) \neq 0$

1. the omitted variable is correlated with the included regressor(s)
2. the omitted variable is a determinant of the dependent variable

$\Rightarrow$  solutions: include omitted variables, use panel data, instrumental variables or RCT

incorrect functional form

$\hookrightarrow$  bias can arise if the wrong functional form is used. we will observe trends in the residuals. solution: fit better models

measurement error

errors in  $x_i$

$\hookrightarrow$  in general,  $\hat{\beta}_i$  will be biased and inconsistent

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + u_i \\ &= \beta_0 + \beta_1 \tilde{x}_i + [\beta_1(x_i - \tilde{x}_i) + u_i] \end{aligned}$$

$\uparrow$   
measured data       $v_i$

errors in  $y_i$

$\hookrightarrow$  in general, no bias to  $\hat{\beta}_i$ , but  $\text{Var}(\hat{\beta}_i)$  will be larger

$$Y_i = \beta_0 + \beta_1 \tilde{x}_i + u_i + w_i$$

$\Rightarrow$  so long as measurement error on  $y_i$  is random,  $E(w_i | x_i) = 0$ , so  
 $E(v_i | x_i) = E(w_i | x_i) + E(u_i | x_i) = 0$

$$\Rightarrow \text{check corr. between } \tilde{x}_i \text{ and } v_i = \frac{\text{cov}(\tilde{x}_i, v_i)}{\sqrt{\text{cov}(\tilde{x}_i, \tilde{x}_i)} \sqrt{\text{cov}(v_i, v_i)}}$$

1.  $\text{cov}(\tilde{x}_i, v_i) = 0$  iff  $E(v_i | x_i) = 0$ , then  $\hat{\beta}_i$  unbiased & consistent

$$\begin{aligned} 2. \text{cov}(\tilde{x}_i, v_i) &= \text{cov}(\tilde{x}_i, \beta_1 x_i - \beta_1 \tilde{x}_i + u_i) \\ &= \text{cov}(\tilde{x}_i, \beta_1 x_i) - \text{cov}(\tilde{x}_i, \beta_1 \tilde{x}_i) + \underbrace{\text{cov}(\tilde{x}_i, u_i)}_{=0, \text{since } E(u_i | x_i) = 0} \\ &= \beta_1 (\text{cov}(\tilde{x}_i, x_i) - \text{cov}(\tilde{x}_i, \tilde{x}_i)) \\ &= \beta_1 (\sigma_x^2 - \tilde{x}_i^2) \end{aligned}$$

$\leftarrow \text{in general, } \neq 0, \text{ so } E(v_i | x_i) \neq 0$

### sample selection bias

in some cases, SRS is thwarted because the dependent variable is only observed for a restricted sample e.g. education or earnings — we only study the sub-group that is employed → we fail to study effect of education earnings had that group worked.

### Simultaneous causality bias

if  $x \rightarrow y$  but  $y \rightarrow x$ , then  $\hat{\beta}_1$  will be picking up the simultaneous effects  
 $\Rightarrow$  OLS estimator of  $\hat{\beta}_1$  will be biased and inconsistent

Suppose:  $y_i = \beta_0 + \beta_1 x_i + u_i \Rightarrow$  we are fitting this

$$x_i = \gamma_0 + \gamma_1 y_i + v_i$$

$$\begin{array}{c} \curvearrowleft \quad \curvearrowright \\ \gamma_1 > 0 \quad \gamma_1 < 0 \end{array}$$

$$u_i \uparrow, y_i \uparrow, x_i \uparrow \quad u_i \uparrow, y_i \uparrow, x_i \downarrow$$

$$P_{u_i, x_i} > 0 \quad P_{u_i, x_i} < 0$$

$$\underbrace{E(u_i | x_i)}_{\neq 0}$$

↳ solution: RCT

### ③ Threats to external validity

↳ whether one can generalize results found in a study to other populations & settings depends on how similar the populations & settings are

e.g. results from Canada → can apply to SG?

↳ practical significance, beyond statistical significance?

## Regression w/ panel data

### ① fixed effects regression

#### T=2, Take difference

$$y_{i1} = \alpha_i + \gamma_1 + \beta_1 x_{i1} + \varepsilon_{i1}$$

$$y_{i2} = \alpha_i + \gamma_2 + \beta_1 x_{i2} + \varepsilon_{i2}$$

$$y_{i2} - y_{i1} = \Delta y_{it}$$

$$\Delta y_{it} = \beta_0 + \beta_1 \Delta x_{it} + \Delta \varepsilon_{it}$$

$$(\gamma_2 - \gamma_1) + \beta_1 (x_{i2} - x_{i1}) + (\varepsilon_{i2} - \varepsilon_{i1})$$

$$y_{it} = \alpha_i + \gamma_t + \beta_1 x_{it} + \varepsilon_{it}$$

$$\bar{y}_i = \alpha_i - \gamma_1$$

$$\bar{\gamma}_t = \gamma_t - \gamma_1$$

#### using binary variables

↳ model has  $T \times n$  different intercepts  $(\gamma_t + \alpha_i)$  for each entity-period pair.

↳ to avoid perfect collinearity, we can model intercepts as  $(n-1)$  binary variables, one for each entity and another  $(T-1)$  binary variables, one for each period.

$$\beta_0 = (\alpha_1 + \gamma_1)$$

$$y_{it} = \beta_0 + \beta_1 x_{it} + \dots$$

$$+ \gamma_2 D_{2i} + \dots + \gamma_T D_{Ti}$$

$$+ f_2 T_{2t} + \dots + f_T T_{Tt}$$

offset for entity  $i$   
is  $i$  the  $n^{\text{th}}$  entity  
is  $t$  the  $T^{\text{th}}$  period  
offset for time  $T$

#### (de-)meaning

$$y_{it} = \alpha_i + \gamma_t + \beta_1 x_{it} + \varepsilon_{it}$$

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

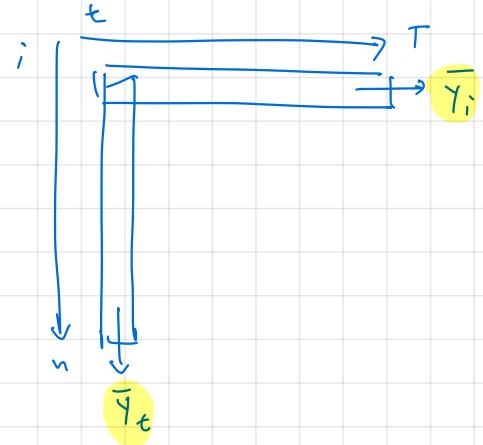
$$= \cancel{\alpha_i} + \cancel{\gamma_t} + \beta_1 x_{it} + \varepsilon_{it}$$

$$- (\cancel{\alpha_i} + \beta_1 \bar{x}_i + \bar{\varepsilon}_i) \quad \begin{matrix} \text{time effect lost} \\ \text{assuming avg. = 0} \end{matrix}$$

$$- (\cancel{\gamma_t} + \beta_1 \bar{x}_t + \bar{\varepsilon}_t) \quad \begin{matrix} \text{entity effect lost} \\ \text{assuming avg. = 0} \end{matrix}$$

$$= \beta_1 (x_{it} - \bar{x}_i - \bar{x}_t) + (\varepsilon_{it} - \bar{\varepsilon}_i - \bar{\varepsilon}_t)$$

$$\underbrace{x_{it}}_{\tilde{x}_{it}} \quad \underbrace{\varepsilon_{it}}_{\tilde{\varepsilon}_{it}}$$



$$\frac{1}{n} \sum_{i=1}^n \alpha_i = 0 \quad \frac{1}{T} \sum_{t=1}^T \alpha_i = \alpha_i$$

$$\frac{1}{n} \sum_{i=1}^n \gamma_t = \gamma_t \quad \frac{1}{T} \sum_{t=1}^T \gamma_t = 0$$

$$\bar{Y}_i = \frac{1}{T} \sum_{t=1}^T y_{it} = \beta_1 \frac{1}{T} \sum_{t=1}^T x_{it} + \frac{1}{T} \sum_{t=1}^T u_{it} + \frac{1}{T} \sum_{t=1}^T \alpha_i + \frac{1}{T} \sum_{t=1}^T \gamma_t$$

$\alpha_i$  ↗ estimated directly by mean if no other regressors

$$\bar{Y}_t = \frac{1}{n} \sum_{i=1}^n y_{it} = \beta_1 \frac{1}{n} \sum_{i=1}^n x_{it} + \frac{1}{n} \sum_{i=1}^n u_i + \frac{1}{n} \sum_{i=1}^n \alpha_i + \frac{1}{n} \sum_{i=1}^n \gamma_t$$

## ② estimation and estimator distribution

### OLS estimation and estimator distributions

$$\hat{\beta}_j^{\text{BeforeAfter}} = \frac{\sum_{i=1}^n (x_{i2} - \bar{x}_{i1})(y_{i2} - \bar{y}_{i1})}{\sum_{i=1}^n (x_{i2} - \bar{x}_{i1})^2}$$

$$\hat{\beta}_j^{\text{demean}} = \frac{\sum_{i=1}^n \sum_{t=1}^T \tilde{x}_{ijt} \tilde{y}_{it}}{\sum_{i=1}^n \sum_{t=1}^T \tilde{x}_{ijt}^2}$$

$$Q_x = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \tilde{x}_{it}^2 \quad n_i = \sqrt{\frac{1}{T} \sum_{t=1}^T \tilde{x}_{it}^2}$$

$$\sqrt{nT} (\hat{\beta}_j - \beta_j) \sim N(0, \frac{\text{Var}(u)}{Q_x^2}) \text{ approximately as } n \rightarrow \infty$$

$$E(\hat{\beta}_j) = \beta_j + \frac{\frac{1}{n} \cdot \frac{1}{T} \cdot \sum_{i=1}^n \sum_{t=1}^T \tilde{x}_{ijt} E(u_{it} | X_{ijt})}{\frac{1}{n} \cdot \frac{1}{T} \sum_{i=1}^n \sum_{t=1}^T \tilde{x}_{ijt}^2}$$

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\beta}_j) = \frac{1}{nT} \frac{\sum_{i=1}^n \tilde{u}_{ii}^2}{Q_x^2}$$

$$\begin{aligned} \text{Var}(n_i) &= \text{Var}\left(\sqrt{\frac{1}{T} \sum_{t=1}^T \tilde{x}_{it}^2} \tilde{u}_{it}\right) = \frac{1}{T} \text{Var}(\tilde{x}_{i1} \tilde{u}_{i1} + \tilde{x}_{i2} \tilde{u}_{i2} \dots) \\ &= \frac{1}{T} \left( \text{Var}(\tilde{x}_{i1} \tilde{u}_{i1}) + \text{Var}(\tilde{x}_{i2} \tilde{u}_{i2}) \dots \right. \\ &\quad \left. + 2 \text{cov}(\tilde{x}_{i1} \tilde{u}_{i1}, \tilde{x}_{i2} \tilde{u}_{i2}) \dots \right) \end{aligned}$$

### assumptions of fixed effects regression

1.  $E(u_{it} | X_{i1}, X_{i2} \dots) = 0$  for each entity,  $u_{it}$  has mean 0 given a specific entity effect and the history of  $X_j$  for that state. i.e. no OVB
2.  $(X_{i1}, X_{i2} \dots X_{iT})$  are i.i.d. entity tuples are independently distributed. i.e. collection of observations for a given entity has no bearing on that of another.
3.  $(X_{it}, Y_{it})$  have finite 4th moments. no extreme outliers
4. No perfect multicollinearity
5.  $\text{corr}(u_{it}, u_{is} | X_{ijt}, X_{jis}, x_i) = 0$  for  $t \neq s$  given some variable  $X_j$ , the error terms are uncorrelated over time within a state. i.e. no serial correlation  $\hookrightarrow$  else we don't need SE

## Heteroskedasticity robust vs. clustered SE

↳ recall that the heteroskedasticity robust estimator of variance is as:

$$\widehat{SE}(\hat{\beta}_j) = \sqrt{\text{Var}(\hat{\beta}_j)} \quad \text{Var}(\hat{\beta}_j) = \frac{1}{n} \sum \frac{\text{Var}((X_i - \bar{X}) u_i)}{\text{Var}(X_i)^2}$$

↳ above, we have derived that in fixed effects regression:

$$\begin{aligned} \lim_{n \rightarrow \infty} \text{Var}(\hat{\beta}_j) &= \frac{1}{nT} \frac{\sum u_i^2}{Q_x^2} \\ \text{Var}(u_i) &= \text{Var}\left(\frac{1}{T} \sum_{t=1}^T \tilde{x}_{it} \tilde{u}_{it}\right) = \frac{1}{T} \text{Var}(\tilde{x}_{i1} \tilde{u}_{i1} + \tilde{x}_{i2} \tilde{u}_{i2} \dots) \\ &= \frac{1}{T} \left( \text{Var}(\tilde{x}_{i1} \tilde{u}_{i1}) + \text{Var}(\tilde{x}_{i2} \tilde{u}_{i2}) \dots \right. \\ &\quad \left. + 2 \text{cov}(\tilde{x}_{i1} \tilde{u}_{i1}, \tilde{x}_{i2} \tilde{u}_{i2}) \dots \right) \end{aligned}$$

⇒ heteroskedasticity robust estimator does not take into account of correlations in time. So we need one that does.

## clustered SE

$$\text{SE}_{\text{clustered}}(\hat{\beta}) = \sqrt{\frac{1}{nT} \frac{\hat{s}_u^2}{Q_x^2}}$$

$$\hat{s}_u^2 = \frac{1}{n-1} \sum_{i=1}^n \left( \underbrace{\sum_{t=1}^T \tilde{x}_{it} \tilde{u}_{it}}_{\tilde{u}_i} - \underbrace{\frac{1}{n} \sum_{i=1}^n \tilde{u}_i}_{\bar{u}_i} \right)^2$$

works out to 0 assuming  $u_i$  and  $x_i$  uncorr.

clustering by state

⇒ i.e. for one  $i$ , avg. across all its time, because unitary likely dependent across time, but independent across states

## Regression w/ dependent variables

### ① binary regression

↪ interpret the regression as modelling the probability that the dependent variable = 1.

model  $f(\{x_j\})$  predicts  $E(Y | \{x_j\})$

$$E(Y | \{x_j\}) = P(Y=1 | \{x_j\}) (1) + P(Y=0 | \{x_j\}) (0)$$

### ② linear probability model

↪ use of a linear model to predict probability

(functional form)

$$E(Y | \{x_j\}) = P(Y=1 | \{x_j\}) = \beta_0 + \sum_j \beta_j x_j$$

(estimation) OLS

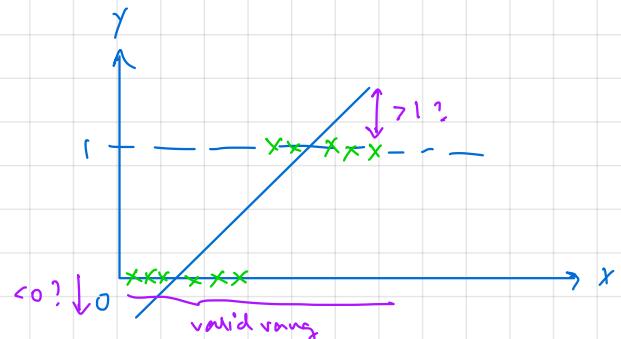
(estimator distributions) same as that of OLS

(inference) same as that of OLS

(shortcomings)

1. probabilities cannot exceed 1 or be negative.

2. effects of regressors on probability likely to be non-linear, have saturation, etc.



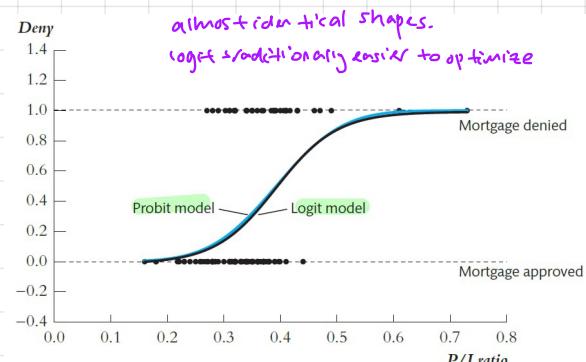
### ③ probit & logit regression

$$z = \beta_0 + \sum_j \beta_j x_j \quad (\text{probit model})$$

$$P(Y=1 | \{x_j\}) \quad \Phi(z), \quad \Phi \text{ is cumulative standard normal}$$

(logit model)

$$\frac{1}{1+e^{-z}}, \text{ logistic function}$$



## ④ estimation and inference of probit & logit regression

### i) non-linear least squares

↳ in probit & logit regression, the coefficients are non-linear with respect to the prediction, so ordinary least squares no longer works. But the idea of taking derivatives to find a closed form solution still applies

$$Q = \sum_{i=1}^n \left[ y_i - f(\hat{\beta}_0 + \sum_j \hat{\beta}_j x_{ij}) \right]^2 \quad L^2 \text{ loss}$$

$$\frac{\partial Q}{\partial \hat{\beta}_j} = 0 \Rightarrow \text{solve}$$

1. estimators are consistent

2. normally distributed in large samples

} but, MLE estimators have smaller variance and are thus more efficient

### ii) maximum likelihood estimation

(idea) assuming our data was collected SRS, iid, then can we estimate  $\beta_j$  by finding the coefficients that maximize the probability of observing this dataset?

1. probability of observing  $y_i | \{x_{ij}\}$ :  $p_i = f(\{ \beta_j \}, \{ x_{ij} \})$

2. we can model it as  $P(Y_i = y_i) = p_i^{y_i} (1-p_i)^{1-y_i}$

3. assuming iid,  $P(Y_1 = y_1, Y_2 = y_2, \dots) = \prod_{i=1}^n P(Y_i = y_i)$

4. convert product to log for convenience

$P(Y_1 = y_1, Y_2 = y_2, \dots) = \sum_{i=1}^n (y_i) \ln(p_i) + (1-y_i) \ln(1-p_i)$

5. optimise by gradient ascent / descent (if negative log, which STATA does)

### iii) MLE estimator distribution

↳  $\hat{\beta}_j^{\text{MLE}}$  is unbiased & consistent, and normally distributed in large samples

↳  $\hat{\beta}_j^{\text{MLE}}$  achieves minimum asymptotic variance among unbiased estimators, and is thus the most efficient estimator in binary dependent regression

⇒ note that  $\text{SE}(\hat{\beta}_j^{\text{MLE}})$  also needs to be robustly adjusted for heteroskedasticity

⇒ inference is identical to OLS for practical purposes

## ⑤ measures of fit

(R<sup>2</sup>) R<sup>2</sup> is a poor measure of fit for the LPM. This is also true for the probit and logit models. Ultimately, the data is inherently non-linear and limited; R<sup>2</sup>, which measures deviation from a continuous line is thus not a good metric.

[accuracy] Does not give a good indicator of quality of prediction  $\Rightarrow$  whatever P = 50% or 90% is the same.

[pseudo-R<sup>2</sup>] compares the likelihood of the estimated model ie.  $P(Y|b_j|X_{ij})$  to the likelihood when none of the X's are included as regressors — that is, probability  $p_i$  is tilted

$$\text{pseudo-}R^2 = 1 - \frac{\ln(L^{\text{max probit}})}{\ln(L^{\text{max Bernoulli}})}$$

conversion to log-likelihood

analogous to  $R^2 = 1 - \frac{ESS}{TSS}$ ,

$1 - \frac{P'}{P} \rightarrow$  explained by model  
 $P \rightarrow$  baseline p of occurrence

$$f^{\text{max probit}} = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}, \quad p_i = \sigma(\hat{b}_0 + \sum b_j x_{ij})$$

$$f^{\text{max Bernoulli}} = \prod_{i=1}^n \hat{p}^{y_i} (1-\hat{p})^{1-y_i}$$

fixed

## ⑥ comparing LPM, logit & probit models

$\hookrightarrow$  all three models are just approximations to the unknown population regression function  $E(Y|X_j) = P(Y=1|X_j)$

1. The LPM is easiest to use and interpret, but cannot capture the non-linear nature of the true function
2. The probit and logit models are identical for all intents and purposes, but as with most non-linear models, interpretation is trickier
3. In some datasets, all three models may be very similar

## IV regression

### ① instrumental variables

1. instrument relevance:  $\text{corr}(z_i, x_j) \neq 0 \Rightarrow z_i \text{ predicts } x_j$
2. instrument exogeneity:  $\text{corr}(z_i, u_i) = 0 \Rightarrow z_i \text{ does not introduce bias}$

### ② TSLS

↳ idea: to find the unbiased effect of  $X$  on  $Y$ , we induce unbiased changes in  $X$  through an exogenous variable  $Z$ .  $Z$  might also predict  $Y$ , though. So to estimate  $\hat{\beta}_j$ , we simply need to account for the effects of  $Z$  on  $Y$ .

$$x_i = \pi_0 + \pi_1 z_i + v_i \Leftrightarrow z_i = -\frac{\pi_0}{\pi_1} + \frac{1}{\pi_1} x_i + v_i$$

$$y_i = \gamma_0 + \gamma_1 z_i + w_i$$

$z_i$  predicts both  $x_i$  and  $y_i$

w/o OVB.

$$E(v_i | z_i) = 0$$

$$E(w_i | z_i) = 0$$

$$y_i = (\underbrace{\gamma_0 - \gamma_1 \frac{\pi_0}{\pi_1}}_{\beta_0} + \underbrace{\frac{\gamma_1}{\pi_1} x_i}_{\beta_1} + \underbrace{(w_i - \frac{1}{\pi_1} v_i)}_{u_i}) \Rightarrow E(u_i | x_i) = 0!$$

$$\beta_1 = \frac{\gamma_1}{\pi_1}$$

OLS estimation of  $\pi_0, \pi_1, \gamma_1, \gamma_0$ .

$$\left. \begin{array}{l} \gamma_1 = \frac{\text{cov}(y_i, z_i)}{\text{var}(z_i)} \\ \pi_1 = \frac{\text{cov}(x_i, z_i)}{\text{var}(z_i)} \end{array} \right\} \quad \boxed{\beta_1 = \frac{\gamma_1}{\pi_1} = \frac{\text{cov}(y_i, z_i)}{\text{cov}(x_i, z_i)}}$$

$$\hat{\beta}_j^{\text{TSLS}} = \frac{s_{yz}}{s_{xz}} = \frac{\cancel{\frac{1}{n} \sum (y_i - \bar{y})(z_i - \bar{z})}}{\cancel{\frac{1}{n} \sum (x_i - \bar{x})(z_i - \bar{z})}}$$

$$\hat{\beta}_j^{\text{TSLS}} \sim N(\beta_1, \frac{\text{Var}((z_i - \mu_z) v_i)}{n \cdot \text{cov}(x_i, z_i)^2})$$

(assumptions of IV regression)

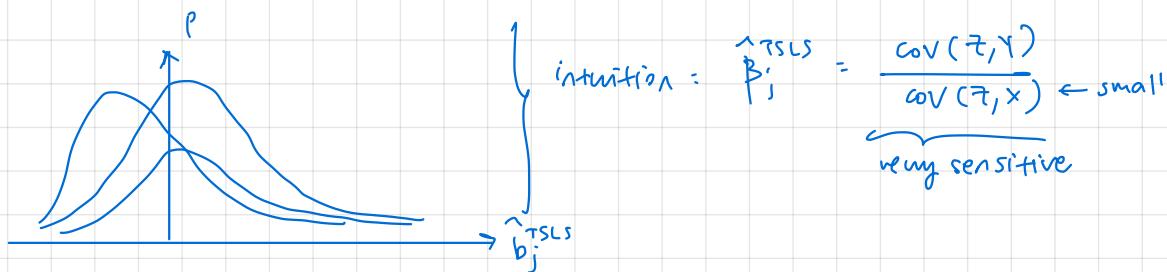
1.  $E(u_i | w_{i1}, \dots, w_{ir}) = 0$
2.  $(y_i, x_i, \dots, w_{ij}, \dots)$  is i.i.d.
3. no extreme outliers
4. IV  $z_1, \dots, z_k$  are valid

$$y_i = \underbrace{\sum_{k=1}^r \beta_k x_{ki}}_{\text{exogenous}} + \underbrace{\sum_{j=k+1}^r \beta_j w_{ji}}_{\text{endogenous}} + u_i$$

### ③ Checking instrument relevance

#### i) weak (irrelevant) instruments

- ↳ a way to think about instrument relevance wrt.  $\hat{\beta}_j$  is like sample size — higher corr.  $\rightarrow$  larger variation in  $X$  explained and triggered
  - ↳ Statistical inference using TSLS relies on  $\hat{\beta}_j^{\text{TSLS}}$  being approx. normal in large samples. But if instruments are weak, distribution is non-normal  $\rightarrow$  even in large samples, and distributed more like ratio of two normal variables
- $\Rightarrow$  not consistent, even if  $n \rightarrow \infty$



#### ii) checking instrument relevance

(idea) regress endogenous  $X$  on  $\tilde{z}_1, \dots, \tilde{z}_m$ , controlled (covariates) with  $w_1, \dots, w_n$ . Then do an F test that coefficients on  $\tilde{z}_1, \dots, \tilde{z}_m$  are zero. If rejection probability is low, instruments are weak

(rule of thumb) first stage F test,  $f \text{ stat} < 10 \Rightarrow$  weak instrument set

↳ bias of TSLS in weak instrument:  $E(\hat{\beta}_j^{\text{TSLS}}) - \beta_j \approx \frac{E(\hat{\beta}_j^{\text{OLS}}) \pi}{E(F) - 1}$

↳ if  $E(F) = 10$ , then  $\frac{E(\hat{\beta}_j^{\text{TSLS}}) - \beta_j}{E(\hat{\beta}_j^{\text{OLS}}) - \beta_j} \approx \frac{1}{10-1} \Rightarrow$  that is, bias of TSLS relative to bias of (bad, since simultaneous) OLS is about 10%.  $\rightarrow$  which is generally acceptable

$\Rightarrow$  if  $f > 10$ , relative bias is  $< 10-1$ , vice versa

#### ④ checking instrument validity

under, exactly, over identified

no. of IVs less than, equal to, more than endogenous variables  $X_j$

(checking instrument exogeneity: J test of overidentified instruments)

↳ idea: exogeneity of IVs mean that they are uncorrelated with  $u_i$ . It logically follows that they should be statistically uncorrelated in our estimate  $\hat{u}_i$

- true  $X$ , rather than  
predicted  $\hat{X}$ . Take note! Not  $\hat{X}$ .
1. estimate  $\hat{y}_i = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + w_{k+1} + \dots \Rightarrow$  get  $\hat{u}_i = \hat{y}_i - \hat{y}_i$
  2. regress  $\hat{u}_i$  on all supposedly exogenous variables  $\Rightarrow \hat{u}_i = f_1 z_1 + \dots + f_m z_m + \dots$
  3. compute F stat. that  $f_1, \dots, f_m$ , all coefficients on IV, are 0.
  4. compute J stat., that is  $J = m f$   
↓  
no. of IV      no. of IV - no. of X endogenous
  5. compute p-value.  $J \sim \chi^2(m-k \text{ d.o.f.}) \rightarrow$  reject  $\Rightarrow$  not exogenous

interpreting the J test

→ what it does say

- hypothesis about instrument exogeneity (all)

↓  
what it does not say

- which instruments are the non-exogenous ones → use theory / intuition to decide removal

effects of OVB in instrument (corr.  $z_i$  &  $w_i$ )

$$\hat{\beta}_1 = \frac{\gamma_1}{\pi_1} = \frac{\text{cov}(y_i, z_i)}{\text{cov}(x_i, z_i)} = \frac{\text{cov}(z_i, \beta_0 + \beta_1 x_i + \beta_2 w_i + u_i)}{\text{cov}(x_i, z_i)}$$

"true" functional form

$$= \frac{\beta_1 \text{cov}(z_i, x_i)}{\text{cov}(x_i, z_i)} \quad \text{ideally since all other } \text{cov}(z_i, \text{---}) = 0$$

⇒ but if  $z_i$  corr. to  $w_i$ , then

$$\hat{\beta}_1 = \frac{\beta_1 \text{cov}(z_i, x_i) + \beta_2 \text{cov}(z_i, w_i)}{\text{cov}(x_i, z_i)}$$

biased & inconsistent!

# Time series

## ① fundamental concept

(The  $j$ th lag)  $y_{t-j}$  is the  $j$ th lag. we use the lag operator  $L$  to denote:

$$Ly_t = y_{t-1} \Leftrightarrow L^j y_t = y_{t-j}$$

(Differences)  $\Delta y_t = y_t - y_{t-1}, \Delta = (1-L)$

$$\begin{aligned} \hookrightarrow \Delta \text{ can be iteratively applied. } \Delta^2 y_t &= (1-L)^2 y_t \\ &= (1 - 2L + L^2) y_t \\ &= y_t - 2y_{t-1} + y_{t-2} \end{aligned}$$

$\hookrightarrow$  note that  $\Delta^j y_t \neq y_t - y_j \Rightarrow$  substitute to  $L$ , then solve

→ autocovariance and autocorrelation

( $i$ th autocovariance) the covariance between  $y_t$  and  $y_{t+j}$ .

( $i$ th auto correlation) the correlation between  $y_t$  and  $y_{t+j}$ .

$$\rho_j = \text{corr}(y_t, y_{t+j}) = \frac{\text{cov}(y_t, y_{t+j})}{\sqrt{\text{var}(y_t)} \sqrt{\text{var}(y_{t+j})}}$$

## 3) types of autocorrelation

(geometric decay)

$$\rho(y_t, y_{t+k}) \approx c^k \text{ for some } c < 1$$

(slow decay) power law

$$\rho(y_t, y_{t+k}) \approx k^\alpha \text{ for some } \alpha < 0$$

## ② stationarity & ergodicity

(stationarity) a mathematical property that a stochastic process has when all the random variables of that process are identically distributed.

$\Rightarrow$  for any set  $x_t \dots x_{t+n}$  in a stochastic process, all  $x_t \dots x_{t+n}$  have the same probability distribution

$\Rightarrow$  the joint distribution of  $Y_t, Y_{t-1} \dots Y_{t-k}$  is identical to  $Y_{t-h}, Y_{t-h-1} \dots Y_{t-h-k}$  for any  $h$ . All corr. joint distributions are identical & time displacement invariant

The joint distributions  $(Y_{st1}, Y_{st2} \dots Y_{sT})$  do not depend on  $s$ , regardless of  $T$  (though may be different for different  $T$ )

(covariance) stationarity

a stochastic process exhibits covariance stationarity if its mean, variance and autocovariance do not depend on  $t$  and variables  $Y_t$  have finite variance.

$\text{cov}(X_t, X_{t+j})$  is independent of  $t$ , at most dependent on  $j$

$E(Y_{t+j}) = \mu_j$  is independent of  $t$ , at most dependent on  $j$

(stationarity up to order  $m$ )

$\hookrightarrow$  strict stationarity is hard in practice

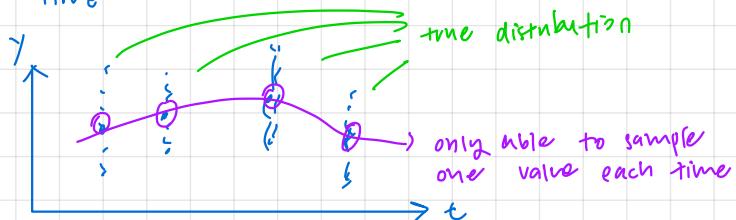
$\hookrightarrow$  often we simply require that stationarity holds up to a certain point in history

$\hookrightarrow$  that is,  $\forall h \in [0, H]$ ,  $\{Y_{t+h}\}_{t \in [h, T+h]}$ , all  $Y_{t+h}$  are identically distributed or have the same mean, variance &  $j$ th autocovariance in the case of weak (autocovariance) stationarity

(ergodicity) a stochastic process is said to be ergodic if its statistical properties can be deduced from a single, sufficiently long random sample of the process

(why are stationarity + ergodicity important?)

$\hookrightarrow$  with stochastic processes of time, it is by definition that there is no way to collect multiple data points for each variable, since that entails going back in time



$\hookrightarrow$  stationarity  $\Rightarrow$   $Y_t$  is identically distributed at all  $t$

$\hookrightarrow$  ergodicity  $\Rightarrow$  in enough time steps, we can estimate the stationary properties!

## ② Autoregression in stationary processes

### (autoregressive distributed lag model)

↪ suppose there is some stable (stationary) relationship between  $y_t$  and its previous  $p$  values, and also some other  $x$  and their previous  $k$  values

↪ that is, the population distribution of  $y_t$  follows

$$y_t = \beta_0 + \beta_1 y_{t-1} \dots + \beta_p y_{t-p} \\ + f_1 x_{t-1} \dots f_k x_{t-k} + u_t \quad \text{some error term}$$

if  $E(u_t | y_{t-1}, y_{t-2} \dots y_{t-p}, x_{t-1}, \dots x_{t-k}) = 0$ , then we can estimate

$\hat{y}_t$  in OLS regression such that  $E(\hat{y}_t) = E(y_t)$ .

ergodicity → we set up tuples of  $y_t, y_{t-1} \dots y_{t-k}$ , and use OLS to estimate  $\beta_0, \beta_1 \dots f_1, \dots$

stationarity → no variables dependent on time

### least squares assumptions

↪ analogous to those of OLS on cross sectional / panel data

1.  $E(u_t | y_{t-1}, y_{t-2} \dots x_{t-1}, x_{t-2} \dots) = 0$ . That is  $E(u_t | \text{regressors}) = 0$
2. All  $y_t, x_{1t}, x_{2t} \dots$  have a stationary distribution. (No dependence on time needs to be modelled)
3.  $(y_t, x_{1t}, x_{2t} \dots)$  and  $(y_{t+j}, x_{1,t+j}, x_{2,t+j} \dots)$  are independent as  $j$  gets large. That is, all regressors are ergodic.
4. No perfect multicollinearity or outliers.

### ③ forecasting errors

(forecast error)

$$FE = Y_{t+1} - \hat{Y}_t$$

(MSFE)

$$MSFE = E((Y_{t+1} - \hat{Y}_{t+1})^2)$$

(RMSE)

$$RMSE = \sqrt{MSFE}$$

(oracle forecast)

suppose  $\hat{\beta}_0, \hat{\beta}_1, \dots$  are estimated perfectly, with zero variance.

$$\text{then } \hat{Y}_t = \hat{\beta}_0 Y_{t-1} + \hat{\beta}_1 Y_{t-2} + \dots + u_t.$$

$$\text{then } E(\hat{Y}_t) = 0, \text{Var}(\hat{Y}_t) = \text{Var}(u_t).$$

$\Rightarrow$  This case is known as the oracle forecast, and this variance and  $MSFE = \text{Var}(u_t)$  for an oracle model.

$$1. \text{ observe } \text{Var}(x) = E(x^2) - E(x)^2$$

but in the general case, is bias variance trade off.

$$2. MSFE = E(Y_T - \hat{Y}_T)^2$$

$$MSE = \text{bias}^2 + \text{var} + \sigma_u^2$$

$$= E(u_T^2) + E((\hat{\beta}_0 - \beta_0) + (\hat{\beta}_1 - \beta_1) Y_{T-1} + \dots)^2$$

$$= E(u_T^2) \text{ for oracle forecast}$$

$$\text{observe that this is } \text{Var}(\hat{\beta}_0 + \dots) = \text{Var}(\hat{Y}_T)$$

$$= \text{Var}(u_T) \text{ since } E(u_T | Y_{T-1}, \dots) = 0$$

$\uparrow$  for the same reason that mean  $> 0$ ,  $E(x)^2 = 0$  if unbiased

3) Estimation of MSFE

$$MSFE = \underbrace{E(u_T^2)}_{\text{variance of random error}} + \underbrace{E((\hat{\beta}_0 - \beta_0) + (\hat{\beta}_1 - \beta_1) Y_{T-1} + \dots)^2)}_{\text{accumulated squared error from coefficients (variance of } \hat{Y})}$$

recall that ergodicity only maintains for long time horizon! How long?  $\Rightarrow$  could be biased

variance of random error

accumulated squared error from coefficients (variance of  $\hat{Y}$ )

(bias variance trade off)

$$E((Y_t + u_t - \hat{Y}_t)^2) = (Y_t - E(\hat{Y}_t))^2 + E(u_t^2) + E((\hat{Y}_t) - \hat{Y})^2$$

$$= \text{bias}^2 + \text{noise} + \text{variance of coefficients}$$

1. approximation by SER if no. of predictors small compared to sample size, contribution of estimation error will be small

$$\widehat{\text{MSFE}} \approx E(u_t^2) = s_u^2 \approx \text{SER}$$

⇒ assumes specific functional form above

2. estimation by final prediction error

↪ assuming errors are homoskedastic

$$\widehat{\text{MSFE}} = \frac{T+p+1}{T} s_{\hat{u}}^2 = \frac{T+p+1}{T-p-1} \frac{\text{SSR}}{T}$$

3. estimate by pseudo OOS

- ↪ the SER method ignores estimation error
  - ↪ the FPE method requires homoskedasticity
- ↪ & formulation also requires stationarity, since functional form of  $\hat{y}$  above based on stationarity
- ↪ pseudo OOS simulates "real time"

Pseudo out-of-sample forecasts are computed using the following steps:

1. Choose a number of observations,  $P$ , for which you will generate pseudo out-of-sample forecasts; for example,  $P$  might be 10% or 20% of the sample size. Let  $s = T - P$ .
2. Estimate the forecasting regression using the estimation sample—that is, using observations  $t = 1, \dots, s$ .
3. Compute the forecast for the first period beyond this shortened sample,  $s + 1$ ; call this  $\tilde{Y}_{s+1|s}$ .
4. Compute the forecast error,  $\tilde{u}_{s+1} = Y_{s+1} - \tilde{Y}_{s+1|s}$ .
5. Repeat steps 2 through 4 for the remaining periods,  $s = T - P + 1$  to  $T - 1$  (reestimate the regression for each period). The pseudo out-of-sample forecasts are  $\tilde{Y}_{s+1|s}, s = T - P, \dots, T - 1$ , and the pseudo out-of-sample forecast errors are  $\tilde{u}_{s+1}, s = T - P, \dots, T - 1$ .

$$\widehat{\text{MSFE}}_{\text{POOS}} = \frac{1}{P} \sum_{s=T-P+1}^T \tilde{u}_s^2$$

P samples for OOS testing, take avg. estimated errors

## Forecast interval

and heteroskedastic

In practice, it is convenient to assume  $u_t$  is normally distributed. Under assumptions of stationarity, forecast error is the sum of  $u_T$  and estimation error of coefficients.

$$E(\hat{Y}_T - \tilde{Y}_T) = E(u_T) + E((\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) Y_1 + \dots)$$

assuming ergodicity, second term is approx. normal by CLT, then if  $u_T$  is normal, forecast error is also normal

$$\hat{Y}_T \pm z_\alpha \widehat{\text{RMSE}} \leftarrow \text{normally distributed}$$

$\widehat{\text{RMSE}}_{\text{PFE}}$  or  $\widehat{\text{RMSE}}_{\text{oos}}$

## ④ Choosing order

F-stat similar to multivariate regression

### Bayesian information criterion

$$BIC(p) = \gamma_n \frac{SSR(p)}{T} + (p+1) \frac{\gamma_n T}{T} \quad \text{intuition}$$

$\rightarrow SSR \downarrow, BIC \downarrow$   
 $\rightarrow p \uparrow, BIC \uparrow$  (penalty)

$$SSR = \sum_{i=1}^N \hat{u}_i^2 \quad \text{sum of squared residuals}$$

Fix  $p$ th model and compute  $BIC(p)$ . Pick model w/ lowest  $BIC$ .

$\Rightarrow \hat{p}$  which minimizes  $BIC$  is a consistent estimator of true lag length, ie. length of lag in underlying population model

$\Rightarrow$  minimising  $BIC$  maximises the posterior probability of the model being the true model

### Akaike information criterion

but in forecasting, "true model" is irrelevant — all that matters is MSE

$$AIC(p) = \gamma_n \frac{SSR(p)}{T} + (p+1) \frac{2}{T}$$

$\rightarrow SSR \downarrow, AIC \downarrow$   
 $\rightarrow p \uparrow, AIC \uparrow$   
 $\rightarrow T \uparrow, AIC \downarrow$

penalty for more  $p$  smaller than  $BIC$

## ⑤ Non-stationarity

### (deterministic trends)

$y_t$  is a non-random function of time. That is, we can include  $t$  as a regressor in the model, and the relationship between  $y_t$  and  $t$  is stable

$$\text{eg. } y_t = \beta_0 + \alpha_1 t + \alpha_2 t^2 + \beta_1 y_{t-1} \dots$$

### (stochastic trends)

$y_t$  is a random function of time. That is, the way  $y_t$  fluctuates is a function of time. i.e.  $\text{Var}(y_t) = \gamma(t)$

$\Rightarrow$  when a trend is present, by definition,  $(y_t, y_{t-1} \dots)$  cannot be stationary. That is, the distribution of  $y_t$  must be different at each time step.

#### random walk model

special case of AR(1) where  $\beta_0 = 0$ ,  $\beta_1 = 1$

$$y_t = y_{t-1} + u_t, u_t \text{ is i.i.d.}$$

#### random walk with drift

$$y_t = \beta_0 + y_{t-1} + u_t, u_t \text{ is i.i.d.}$$

$$1. E(y_t) = \beta_0 + y_{t-1} + E(u_t)$$

$$2. \text{Var}(y_t) \text{ same} = t \text{Var}(u_t)$$

$$3. \text{cov}(y_t, y_{t-k}) = (t-k) \text{Var}(u_t), \text{ since } \text{Var}(\beta_0) = 0$$

$$2. y_t = y_{t-2} + u_{t-1} + u_t \dots$$

$$= y_{t-k} + \sum_{i=1}^k u_{t-i} = y_0 + \sum_{i=1}^t u_{t-i}$$

$$\text{Var}(y_t) = 0 + t \text{Var}(u_t) \text{ if homoskedasticity}$$

$$3. \text{cov}(y_t, y_{t-k}) = \text{cov}\left(y_0 + \sum_{i=1}^t u_i, y_0 + \sum_{i=1}^k u_i\right)$$

$$= (t-k) \text{Var}(u_i) \text{ since } \text{cov}(u_i, u_j) = 0 \text{ for } i \neq j$$

## ⑥ Testing non-stationarity

$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} \dots + u_t, u_t \text{ is i.i.d.}$$

$$= \beta_0 + \beta_1 (\beta_0 + \beta_1 y_{t-2} \dots) + \beta_2 (\beta_0 + \beta_1 y_{t-3} \dots) + u_t$$

roots are then  $\bar{\beta}$  that satisfy  $1 - \beta_1 \bar{\beta} - \beta_2 \bar{\beta}^2 \dots - \beta_p \bar{\beta}^p = 0$   
 $p$  roots. If all  $|\bar{\beta}| > 1$ , stationary.

$\Leftrightarrow$  sum of coefficients  $\beta_1 + \dots + \beta_p = 1$  if not stationary

### (Dickey Fuller test in the AR(1) model)

↪ main idea: not stationary  $\rightarrow \sum_{i=1}^p \beta_i = 1 \Rightarrow H_0: \text{not stationary, try to reject}$

$$Y_t = \beta_0 + \beta_1 Y_{t-1} + u_t$$

$$\begin{aligned} H_0: \beta_1 - 1 &= 0 \\ H_1: \beta_1 - 1 &< 0 \end{aligned} \quad \left. \begin{array}{l} \text{fit } Y_t - Y_{t-1} = \beta_0 + \beta_1 Y_{t-1} - Y_{t-1} + u_t \\ \Leftrightarrow \Delta Y_t = \beta_0 + (\beta_1 - 1) Y_{t-1} + u_t \end{array} \right\}$$

test  $\frac{f}{SE(f)}$ , where  $f$  is DF distributed against critical values

### (Augmented Dickey Fuller test for AR(p))

1. choose  $p$  using F-test / BIC / AIC

2. ADF test

$$\Delta Y_t = \beta_0 + f Y_{t-1} + \sum_{i=1}^{p-1} \gamma_i \Delta Y_{t-i}, \quad f = \sum_{i=1}^{p-1} \beta_i - 1$$

$$\begin{aligned} H_0: f &= \sum_{i=1}^p \beta_i - 1 = 0 \rightarrow \text{non-stationary} \\ H_1: f &< 0 \rightarrow \text{not non-stationary} \end{aligned} \quad \left. \begin{array}{l} \text{test against ADF} \\ \text{statistic distribution} \\ \text{by } \frac{f}{SE(f)} \text{ vs. } ADF_\alpha \end{array} \right\}$$

### (testing for stochastic & deterministic non-Horizontality)

↪ we extend the augmented DF test by including function of  $t$

$$\Delta Y_t = \beta_0 + \cancel{\alpha t} + f Y_{t-1} + \sum_{i=1}^{p-1} \gamma_i \Delta Y_{t-i}, \quad f = \sum_{i=1}^{p-1} \beta_i - 1$$

$$\begin{aligned} H_0: f &= \sum_{i=1}^p \beta_i - 1 = 0 \rightarrow \text{non-stationary} \\ H_1: f &< 0 \rightarrow \text{not non-stationary} \end{aligned} \quad \left. \begin{array}{l} \text{test against ADF} \\ \text{statistic distribution} \\ \text{by } \frac{f}{SE(f)} \text{ vs. } ADF_\alpha \end{array} \right\}$$