

$Cov(X, Y) = E((X - \mu_X)(Y - \mu_Y))$      $Cov > 0$ : on average  $X, Y$  above (below) mean together.

$Cov < 0$ : on average, when  $X$  is above (below) mean,  $Y$  is below (above).

$(X, Y)$ : independent  $\implies Cov(X, Y) = 0$      $Cov(X, Y) = 0 \not\Rightarrow (X, Y)$ : independent

$Cov((a_1X + b_1), (a_2Y + b_2)) = a_1a_2Cov(X, Y)$      $Var(aX + bY) = a^2Var(X) + b^2Var(Y) + 2abCov(X, Y)$

$Corr(X, Y) = \rho_{XY} = \frac{Cov(X, Y)}{\sigma_X\sigma_Y}$      $-1 \leq \rho \leq 1$      $(X, Y)$ : independent  $\implies Corr(X, Y) = 0$

$Corr(X, Y) = 0 \not\Rightarrow (X, Y)$ : independent     $E(X) = \sum_{i=1}^k x_i f_i(x_i)$  [or]  $\int_{-\infty}^{\infty} xf(x)dx$

$E(g(x)) = \sum_{i=1}^k g(x_i) f_i(x_i)$  [or]  $\int_{-\infty}^{\infty} g(x)xf(x)dx$      $E(c) = c$ ,     $E(aX + b) = aE(X) + b$ ,

$E(X + Y) = E(X) + E(Y)$      $Var(X) = E(X - \mu)^2 = \sigma^2 = E(X^2) - (E(X))^2$      $Var(c) = 0$ ,     $Var(aX +$

$b) = a^2Var(X)$      $Z = \frac{X-\mu}{\sigma} \sim N(0,1)$      $\Phi(c) = z$  - table reading for number  $c$

$P(Z > z) = 1 - \Phi(z)$ ,     $P(Z < -z) = P(Z > z)$ ,     $P(a \leq Z \leq b) = \Phi(b) - \Phi(a)$

$P(|Z| > c) = P(Z > c) + P(Z < -c) = 2(1 - \Phi(c))$      $E(Y|X = x) = \sum_j y_j f_{Y|X}(y_j|x)$  or  $\int y_i f_{Y|X}(y_i|x) dy$

Simple linear regression model.  $y = \beta_0 + \beta_1 x + u$  assuming  $E\{u|x\} = 0 \implies E\{y|x\} = \beta_0 +$

$\beta_1 x$      $E(u) = 0 \implies E(y - \beta_0 - \beta_1 x) = 0$      $Cov(x, u) = E(xu) = 0 \implies E(x(y - \beta_0 - \beta_1 x)) = 0$

$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$      $\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i(y_i - \bar{y})}{\sum_{i=1}^n x_i(x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$      $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x$      $\hat{u}_i = y_i - \hat{y}_i$

$\sum_{i=1}^n \hat{u}_i = 0$      $\sum_{i=1}^n x_i \hat{u}_i = 0$     The point  $(\bar{x}, \bar{y})$  is always on the regression line.

$bias(\hat{\beta}_1) = E(\hat{\beta}_1) - \beta_1$      $E(\hat{\beta}_1) = \beta_1 + E\left(\frac{\sum u_i x_i}{\sum x_i^2}\right) = \beta_1 + \frac{cov(u, x)}{var(x)} = \beta_1 + 0$

$\therefore bias(\hat{\beta}_1) = \beta_1 - \beta_1 = 0$  ■

$SST =$  total sum of squares = total variation of  $y$  about its mean  $\sum (y_i - \bar{y})^2$

$SST_x =$  total variation of  $x$  about its mean  $\sum (x_i - \bar{x})^2$      $SSR = \sum (y_i - \hat{y}_i)^2$

$SSE =$  explained sum of squares = predicted variation of  $y$  about  $\bar{y} = \sum (\hat{y}_i - \bar{y})^2$

$R^2 = \frac{SSE}{SST} = 1 - \frac{SSR}{SST}$     shows fraction of variation in  $y$  explained by  $x$

$E(\hat{\beta}_0) = \beta_0$      $E(\hat{\beta}_1) = \beta_1$      $Var(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{nSST_x}$      $Var(\hat{\beta}_1) = \frac{\sigma^2}{SST_x}$

$[SST_x = \sum (x_i - \bar{x})^2]$  Unbiased estimator of population variance:  $\hat{\sigma}^2 = \frac{1}{n-2} \sum \hat{u}_i^2 = \frac{SSR}{n-2}$

$SER = \hat{\sigma} = \sqrt{\hat{\sigma}^2}$      $SE(\hat{\beta}_0) = \sqrt{Var(\hat{\beta}_0)} = \sqrt{\frac{\sigma^2 \sum x_i^2}{nSST_x}}$      $SE(\hat{\beta}_1) = \sqrt{Var(\hat{\beta}_1)} = \sqrt{\frac{\sigma^2}{SST_x}}$

modify dependent:  $cy = c\beta_0 + c\beta_1 x$     modify independent:  $y = \beta_0 + \left(\frac{\beta_1}{k}\right)(kx)$

OVB:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$  but we estimate  $y = \widetilde{\beta}_0 + \widetilde{\beta}_1 x_1$  and  $\widetilde{x}_2 = \widetilde{\delta}_0 + \widetilde{\delta}_1 x_1$  then  $\widetilde{\beta}_1 = \widehat{\beta}_1 + \widehat{\beta}_2 \widetilde{\delta}_1$  then  $E(\widetilde{\beta}_1) = \beta_1 + \beta_2 \widetilde{\delta}_1$  Including irrelevant vars maintains unbiasedness, but precision of estimates is bad.

Multicollinearity: correlation among expl vars leads to bias in coefficients

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{SST_j(1-R_j^2)} \quad \text{where } SST_j = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \quad \text{and}$$

test statistic:  $t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$     null hypothesis  $H_0: \beta_1 = 0$     alt  $H_A: \beta_1 < 0, \beta_1 > 0, \beta_1 \neq 0$

$Y$  – Variables: Quantitative (Cardinal or Ordinal), or Qualitative (Ordered or Unordered)

Regression: Cardinal (OLS), Ordinal (try to make cardinal, if not, ordered logit),

Categorical 2 groups (logit), Unordered > 2 groups

(multinomial logit), Ordered > 2 groups (ordered logit)

$X$  causes  $Y$ : causality;  $Y$  causes  $X$ : reverse causality; both: simultaneous causality

Sources of Bias: OVB, Measurement error, Sample Selection Bias, Misspecification.

Overstated:  $|\hat{\beta}_1| > \hat{\beta}_1$     Understated:  $|\hat{\beta}_1| < \hat{\beta}_1$

OVB:  $Y = X + Z$ , but  $Z$  omitted. Biased when  $\text{Corr}(X, Z) \neq 0$  and  $\text{Corr}(Y, Z) \neq 0$  (hold  $X$  constant)

	$\hat{\beta}_1 > 0$		$\hat{\beta}_1 < 0$	
	$\text{Corr}(Y, Z) > 0$	$\text{Corr}(Y, Z) < 0$	$\text{Corr}(Y, Z) > 0$	$\text{Corr}(Y, Z) < 0$
$\text{Corr}(X, Z) > 0$	Overstated	Understated	Understated	Overstated
$\text{Corr}(X, Z) < 0$	Understated	Overstated	Overstated	Understated

Error in  $X$ : Attenuation (slope understated), Imprecision (large std errors)

Errors in  $Y$ : None (slope unbiased), Imprecision (large std errors)

Coefficient interpretation with quadratic  $X$  (...  $\beta_1 X + \beta_2 X^2$  ...):

$\beta_1 > 0, \beta_2 > 0$ : incr at incr rate     $\beta_1 > 0, \beta_2 < 0$ : incr at decr rate

$\beta_1 < 0, \beta_2 > 0$ : decr at incr rate     $\beta_1 < 0, \beta_2 < 0$ : decr at decr rate

lin – lin:  $y = \beta_0 + \beta_1 x$      $\beta_1 =$  unit change in  $y$  for a 1 unit change in  $x$

lin – log:  $y = \beta_0 + \beta_1 \ln x$      $\frac{\beta_1}{100} =$  unit change in  $y$  for a 1% change in  $x$

log – lin:  $\ln y = \beta_0 + \beta_1 x$      $100\beta_1 =$  % change in  $y$  for a 1 unit change in  $x$

log – log:  $\ln y = \beta_0 + \beta_1 \ln x$      $\beta_1 =$  % change in  $y$  for a 1% change in  $x$

Difference in Differences:  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 (X_1 * X_2) + \varepsilon$

	$X_2 = 0$	$X_2 = 1$	Difference
$X_1 = 1$	$\beta_0 + \beta_1$	$\beta_0 + \beta_1 + \beta_2 + \beta_3$	$\beta_2 + \beta_3$
$X_1 = 0$	$\beta_0$	$\beta_0 + \beta_2$	$\beta_2$
	Difference in Differences:		$\beta_3$

Good instrument: Relevant ( $\text{Corr}(X, Z) \neq 0$ ) and Exogenous ( $\text{Corr}(Z, \varepsilon) = 0$ )

$P\text{value} = .02 \Rightarrow$  in one out of fifty samples, we would observe result if  $H_0$ : true.