

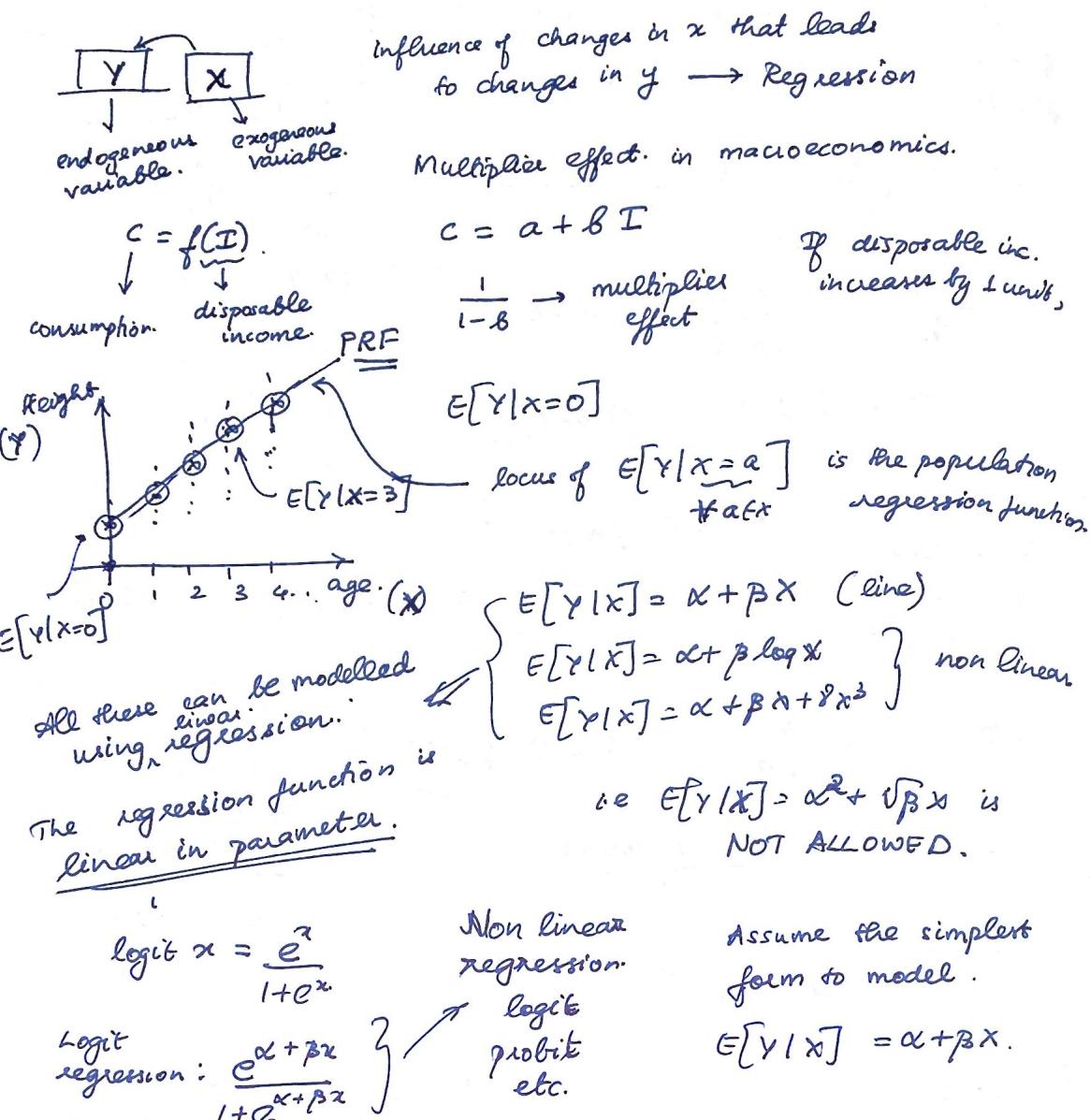
Linear Regression

Statistics and Econometrics.

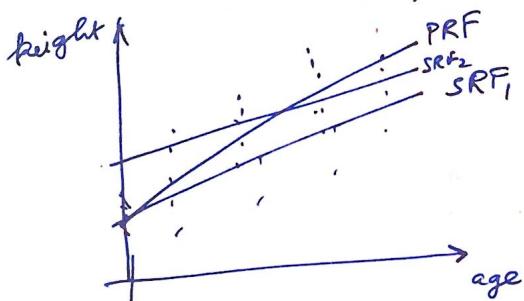
Linear Regression - Fundamentals

- Linear Regression estimator \hat{Y} is makes an approximate estimate of $E(Y|X)$ given that Y is the dependent and X is the independent variable.
- Population Regression function and Sample Regression function.
- Linear Regression is linear in parameter. Not necessarily in variables.
- Stochastic and deterministic terms (regressors) in the model.
- Minimization of squares of residuals and the OLS criterion. Normal equations.
- Properties of OLS Regression. Goodness of fit measure.
- Classical Linear Regression Model (CLRM) Assumptions. Constant mean for errors, Homoscedasticity, No autocorrelation or serial correlation, etc. Problems of endogeneity. Omitted variable bias, Reverse causality.
- Unbiasedness and Variance properties of linear regression (referred to as BLUE properties).
- Endogeneity leads to OLS estimators being biased. Not autocorrelation nor heteroscedasticity.
- Hypothesis testing in regression models, Interval estimation, Q-statistic, Standard errors, Confidence interval, Life expectancy - GDP example.
- Linear regression model with 2 independent variables. Derivations for estimators, coefficients, R^2 , and unbiasedness and variance of estimators.

Linear Regression - basics



Now each one takes a different sample and does the same except → obtains the SRF sample regression function



$$\hat{y} = \hat{a} + \hat{b}x$$

predicted \hat{y} →
it is an estimate of $E[Y|x]$

$$\hat{y} = \hat{a} + \hat{b}x$$

$$\hat{y} \approx E[Y|x]$$

treating this in a deterministic format.
Because x influences \hat{y} is known & fixed.

But there could be other factors too that influence Y (height) not just x (age).

Stochastic and deterministic terms in regression model

The actual height can be explained by using stochastic terms to model height.

$$Y_i = \alpha + \beta_i x_i + u_i$$

deterministic term
randomness term.
Because of inherent randomness in Human behavior.

captures all the remaining effects that are not captured by x_i

We introduce this because,

$$Y_i = \alpha + \beta x_i + u_i$$

deterministic regressor
stochastic regressor.

PRF.

$$SRF: \hat{Y}_i = \hat{\alpha} + \hat{\beta} x_i + e_i$$

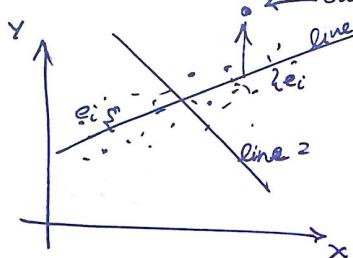
e_i 's and u_i 's will not be the same.

You can't reduce u_i 's but you can reduce the e_i 's

error terms.
"residuals"

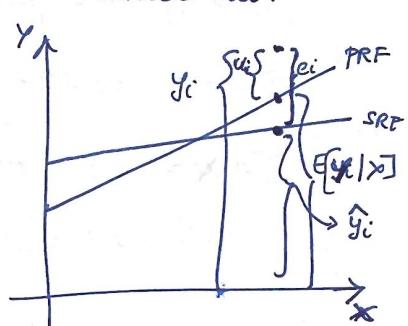
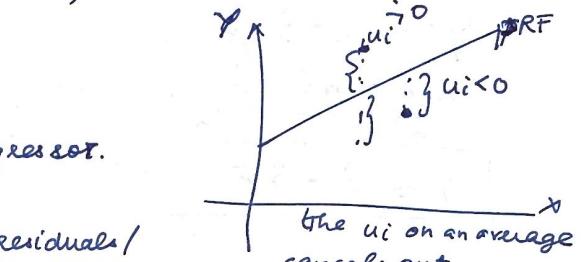
residuals /

inherent randomness in the systems



regressors → means

Always influenced by outliers. Hence outliers significantly affect regression.



Why line 1 & not line 2?

→ minimize error

→ $\min e_i^2$

$e_i^1 \rightarrow$ +ve & -ve cancel out.
 $e_i^3 \rightarrow$ same problem.

$$\min \sum (y_i - \hat{\alpha} - \hat{\beta} x_i)^2$$

The $\hat{\alpha}$ and $\hat{\beta}$ obtained by minimizing are called ordinary least square estimates.

Linear : Use OLS

Non linear : Use MLE

why? Gaussian Markov.

: Use MOM

$$\min \sum (y_i - \hat{\alpha} - \hat{\beta} x_i)^2$$

$$\frac{\partial^2 \sum e_i^2}{\partial \hat{\alpha}^2} = 0 \Rightarrow 2 \sum (y_i - \hat{\alpha} - \hat{\beta} x_i)(-1) = 0$$

$$\frac{\partial^2 \sum e_i^2}{\partial \hat{\beta}^2} = 0 \Rightarrow 2 \sum (y_i - \hat{\alpha} - \hat{\beta} x_i)(-x_i) = 0$$

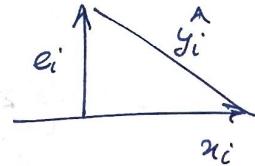
Note : You talk of distribution only for u_i ; not for e_i 's

Normal Equations

$$\begin{aligned}\sum (y_i - \hat{\alpha} - \hat{\beta} x_i) &= 0 \\ \sum (y_i - \hat{\alpha} - \hat{\beta} x_i) x_i &= 0\end{aligned} \Rightarrow \boxed{\begin{array}{l}\sum e_i = 0 \\ \sum e_i x_i = 0\end{array}}$$

e_i 's are always perpendicular to x_i 's.

$$\sum e_i = 0$$



$$\sum y_i - n\hat{\alpha} - \hat{\beta} \sum x_i = 0$$

$$\frac{\sum y_i}{n} = \frac{n\hat{\alpha}}{n} + \hat{\beta} \frac{\sum x_i}{n}$$

$$\boxed{\bar{y} = \hat{\alpha} + \hat{\beta} \bar{x}}$$

line passing through \bar{x} and \bar{y}

\therefore Any line passing through \bar{x} and \bar{y} will have $\boxed{\sum e_i = 0} \rightarrow \bar{e} = 0$

$\sum e_i = 0 \Rightarrow$ SRF always passes through \bar{x} and \bar{y}
sample covariance. sample means // not population means.

$$\sum e_i x_i = 0 \Rightarrow \text{cov}(e, x) = 0$$

$$\text{cov}(e, x) = \frac{\sum (x_i - \bar{x})(e_i - \bar{e})}{n} = 0 \quad \text{when } \bar{e} = 0;$$

$$= \frac{\sum (x_i - \bar{x})e_i}{n} = \frac{\sum x_i e_i - \bar{x} \sum e_i}{n} = 0$$

Properties of OLS Regression

① $\sum e_i = 0 \Rightarrow$ SRF always passes through \bar{x}, \bar{y} .

$$\Rightarrow \bar{y} = \hat{y}$$

$$y = \hat{\alpha} + \hat{\beta} x_i + e_i$$

$$\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$$

$$\sum_n y = \frac{n\hat{\alpha}}{n} + \hat{\beta} \frac{\sum x_i}{n} + \sum e_i^0$$

$$\sum_n \hat{y}_i = \frac{n\hat{\alpha}}{n} + \hat{\beta} \frac{\sum x_i}{n}$$

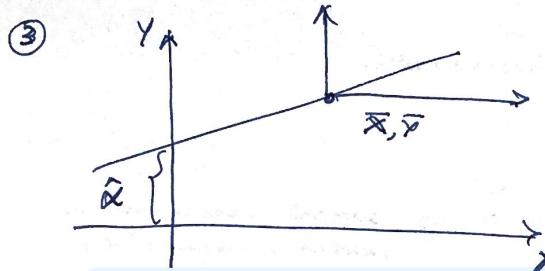
② $\sum x_i e_i = 0 \Rightarrow x$ and e are not correlated.

\hat{y} is a linear function of x .

$$\boxed{\text{Cov}(x, e) = 0} \Rightarrow \text{Cov}(\hat{y}, e) = 0$$

$$\text{cov}(\hat{y}, e) = \text{cov}(\hat{\alpha} + \hat{\beta} x, e) = \text{cov}(\hat{\alpha}^0, e)$$

$$\boxed{\text{Cov}(\hat{y}, e) = 0} \quad + \hat{\beta} \text{cov}(x, e) = 0$$



This shift of origin leads to disappearance of α . $\alpha = 0$ but slope $\hat{\beta}$ remains the same.

$$y_i = \hat{\alpha} + \hat{\beta} x_i + e_i$$

$$\bar{Y} = \hat{\alpha} + \hat{\beta} \bar{X}$$

$$y_i - \bar{Y} = \hat{\beta} (x_i - \bar{x}) + e_i$$

↓ shifting form of origin
to (\bar{x}, \bar{y})

SRF in deviation form:

$$y_i = \hat{\beta} x_i + e_i$$

SRF in the raw form

subtract averages.
make $\hat{\alpha} = 0$.

$$\min \sum e_i^2 = \sum (y_i - \hat{\beta} x_i)^2$$

$$\frac{\partial \sum e_i^2}{\partial \beta} = 2 \sum (y_i - \hat{\beta} x_i)(-x_i) = 0 \Rightarrow \sum y_i x_i = \hat{\beta} \sum x_i^2$$

$$\sum e_i = 0 \quad \sum y_i = n \hat{\alpha} + \hat{\beta} \sum x_i \Rightarrow \frac{\sum y_i - \hat{\beta} \sum x_i}{n} = 0$$

$$\sum e_i x_i = \hat{\alpha} \sum x_i + \hat{\beta} \sum x_i^2$$

$$\sum y_i x_i = \frac{\sum y_i - \hat{\beta} \sum x_i}{n} + \hat{\beta} \sum x_i^2$$

$$\sum y_i x_i = \sum y_i + \hat{\beta} (n \sum x_i^2 - \sum x_i) \quad \hat{\beta} = \frac{n \sum y_i x_i - \sum y_i \sum x_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{\beta} = \frac{\sum y_i x_i / n}{\sum x_i^2 / n} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

④ Goodness-of-fit

$$y_i = \hat{\beta} x_i + e_i$$

$$\sum y_i^2 = \sum (\hat{\beta} x_i + e_i)^2$$

$$\sum (y_i - \hat{y}_i)^2 = \hat{\beta}^2 \sum x_i^2$$

$$\sum y_i^2 = \underbrace{\hat{\beta}^2 \sum x_i^2}_{\text{Total}} + \underbrace{\sum e_i^2}_{\text{Explained}} + \underbrace{2 \hat{\beta} \sum x_i e_i}_{\text{Residual}}$$

R^2 is a factor only in OLS method.
Not for MLE etc.

$$R^2 = \frac{ESS}{TSS} \times 100 = \text{%.}$$

$$R^2 = 0.95 \leftarrow 95\%$$

Variation in y
can be explained
by x.

CLRM Assumptions

$$\begin{aligned}
 \text{PRF} &\Rightarrow Y_i = \beta_1 + \beta_2 x_i + u_i & u_i &\leftarrow \text{random variable} \\
 \text{SRF} &\Rightarrow Y_i = \hat{\beta}_1 + \hat{\beta}_2 x_i + e_i \\
 \min \sum e_i^2 &\quad \frac{\partial \sum e_i^2}{\partial \hat{\beta}_1} = 0 \Rightarrow \sum e_i = 0 \Rightarrow \text{sample regression line} \\
 &\quad \text{passes through } \bar{x} \text{ and } \bar{Y} \\
 &\quad \frac{\partial \sum e_i^2}{\partial \hat{\beta}_2} = 0 \Rightarrow \sum e_i x_i = 0 \Rightarrow \text{cov}(x, e) = 0 \\
 &\quad \text{cov}(\hat{Y}, e) = 0 \\
 Y_i &= \hat{\beta}_1 + \hat{\beta}_2 x_i + e_i \\
 \bar{Y} &= \hat{\beta}_1 + \hat{\beta}_2 \bar{x} \\
 Y_i - \bar{Y} &= \hat{\beta}_2 (x_i - \bar{x}) + e_i \\
 Y_i &= \hat{\beta}_2 x_i + e_i \\
 \sum y_i^2 &= \hat{\beta}_2^2 \sum x_i^2 + \sum e_i^2 \\
 \text{TSS} &\quad \text{ESS} \qquad \text{RSS}
 \end{aligned}$$

Numerical properties of OLS.

$u_i \leftarrow \text{random variable}$
 $y_i \rightarrow$
 $\hat{\beta}_2 \rightarrow$

Statistical Properties of OLS
CLRM assumptions

① $E[u_i | x_i] = 0$

The average of u_i 's for a given x_i is zero. \Rightarrow consequence of randomness of u_i 's

Otherwise \rightarrow systematic over/under estimation.

$$E[u_i | x_i] = 0 \Rightarrow E(u_i) = 0$$

② $E[u_i^2 | x_i] = \sigma^2$ (Homoscedasticity)

Heteroscedasticity problem:

- ✓ variance of u_i changes with x_i
- common with cross section data.

$$\text{Homoscedasticity} \Rightarrow E[u_i^2 | x_i] = \sigma^2 \Rightarrow E[E[u_i^2 | x_i]] =$$

$$\text{Heteroscedasticity} \Rightarrow E[u_i^2 | x_i] = \sigma_i^2 = E[u_i^2] = \sigma^2$$

$$\text{Var}[u_i] = \{E[u_i^2] - (E[u_i])^2\} = \sigma^2 - 0 = \sigma^2$$

③ $E[u_i u_j | x_i x_j] = 0$ \leftarrow mostly untrue for time series data. Eg: GDP \rightarrow always correlated with previous years data

$$\text{Cov}(u_i, u_j) = E[u_i u_j] - E[u_i] E[u_j]$$

\downarrow \uparrow
 $\circ \quad \circ$
(no correlation)

(→ doesn't change overnight?
persistence to shock?)

$$\text{Cov}(u_i, u_j) \neq 0$$

$E[u_t, u_{t-1}] = 0 \leftarrow$ a very strong assumption to make in time series data.

$E[u_t, u_{t-4}] = 0 \leftarrow$ the data is not correlated for 4 lags
 \nwarrow lag: difference between t & $t-4$

Individual outcome is not correlated with another individuals' outcome.

$E[u_i u_j | x_i, x_j] = 0 \Rightarrow$ no auto correlation / serial correlation

cluster level correlation \Rightarrow stratified sampling outcomes to a particular data are correlated \Rightarrow school data

(4) $E[u_i x_i] = 0 \quad \text{cov}(u_i, x_i) = 0$ $\xrightarrow{\text{School}}$
 $wage_i = \alpha + \beta x_i + u_i$
 \nwarrow therefore goes into u_i
since ability also determines the number of years of schooling one has attended \Rightarrow problem of endogeneity

Return to education estimation
 $x_i \rightarrow$ schooling.
ability is also a factor for wage
 \downarrow can't be measured & can't be included as a factor

- Problem of endogeneity \Rightarrow
 - omitted variable bias
 - reverse causality
 - simultaneity.
- Problem of feedback:

$x_i \xrightarrow{\text{causes}} y_i$ but also,
 $y_i \longrightarrow x_i$

Unbiasedness and Variance Properties of Linear Regression

$$\hat{\beta}_2 = \frac{\sum x_i y_i}{\sum x_i^2} \quad w_i = \frac{x_0}{\sum x_i^2} \quad \sum w_i = \frac{\sum x_i}{\sum x_i^2} = \frac{\sum x_i}{\sum x_i^2} = \frac{\sum (x_i - \bar{x})}{\sum x_i^2}$$

$$\therefore \sum w_i = 0$$

$$\sum w_i x_i = \sum \frac{x_i}{\sum x_i^2} x_i = \frac{\sum x_i^2}{\sum x_i^2} \Rightarrow \sum w_i x_i = 1$$

$$\sum w_i^2 = \sum \left(\frac{x_i}{\sum x_i^2} \right)^2 = \frac{\sum x_i^2}{(\sum x_i^2)^2} = \frac{1}{\sum x_i^2} \therefore \sum w_i^2 = \frac{1}{\sum x_i^2}$$

$$\sum w_i x_i = \sum w_i (x_i - \bar{x}) = \sum w_i x_i - \sum w_i \bar{x} = \sum w_i x_i$$

$$\therefore \sum w_i x_i = \sum w_i x_i = 1.$$

Note

$$w_i = \frac{x_i}{\sum x_i^2}$$

$$\sum w_i = 0$$

$$\sum w_i x_i = 1$$

$$\sum w_i^2 = \frac{1}{\sum x_i^2}$$

$$\sum w_i x_i = 1$$

Unbiased if: $E[\hat{\beta}_1] = \beta_1, E[\hat{\beta}_2] = \beta_2$

$$E[\hat{\beta}_2] = E\left[\frac{\sum x_i y_i}{\sum x_i^2}\right] = E\left[\sum \frac{x_i}{\sum x_i^2} y_i\right] = E\left[\sum w_i y_i\right]$$

$$\sum w_i y_i = \sum w_i (Y_i - \bar{Y}) = \sum w_i Y_i - \sum w_i \bar{Y}$$

$$E[\hat{\beta}_2] = E\left[\sum w_i Y_i\right] = E\left[\sum w_i (\beta_1 + \beta_2 X_i + u_i)\right]$$

At the same time this is WRONG!!

$$\hat{\beta}_2 = \sum w_i Y_i = \sum w_i [\beta_1 + \beta_2 X_i + u_i] \leftarrow \text{WRONG}$$

Why use β_2 & not $\hat{\beta}_2$? \nearrow here?
Because expectations are for populations and not for samples.
 \bar{Y} - sample average
 $E[Y]$ - population average

$$E[\hat{\beta}_2] = E[\beta_2 + \sum w_i \hat{x}_i^0 + \sum w_i u_i]$$

$$\sum w_i x_i^0 = \sum w_i x_i = 1$$

$$E[\hat{\beta}_2] = E[\beta_2 + \sum w_i u_i]$$

$u_i \rightarrow$ random variable
 $y_i \rightarrow$ function of u_i . Hence R.V.

$$E[\hat{\beta}_2] = E[\beta_2] + E[\sum w_i u_i]$$

$x_i \rightarrow$ are fixed. Not R.V.

$$E[\hat{\beta}_2] = \beta_2 + \sum w_i E[u_i^0]$$

$w_i \rightarrow$ function of u_i . Not R.V. \downarrow
 Because of C.R.M assumption. Expectation operator does NOT act on them

$$E[\hat{\beta}_2] = E[\bar{y} - \hat{\beta}_1 \bar{x}]$$

$$E[\hat{\beta}_1] = E[(\beta_1 + \beta_2 \bar{x} + \bar{u}) - \hat{\beta}_2 \bar{x}]$$

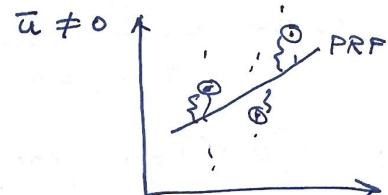
$$E[\hat{\beta}_1] = E[\beta_1 - (\hat{\beta}_2 - \beta_2) \bar{x} + \bar{u}]$$

$$E[\hat{\beta}_1] = \beta_1 - \bar{x} E[\hat{\beta}_2^0 - \beta_2] + E[\bar{u}]$$

$$\therefore E[\hat{\beta}_2] = \beta_2$$

$$E[\bar{u}] = E\left[\frac{\sum u_i}{n}\right] = \frac{\sum E[u_i]}{n}$$

$$\therefore E[\bar{u}] = 0$$



$$\text{population } E[u_i] = 0$$

$$\text{sample: } \bar{u} \neq 0$$

$$\text{OLS min. } \Rightarrow \bar{e} = 0$$

you take points corresp to the sample above.

$$\boxed{E[\hat{\beta}_1] = \beta_1} \quad \boxed{E[\hat{\beta}_2] = \beta_2}$$

OLS estimators are unbiased.

OLS estimators are still unbiased even if you have autocorrelation or heteroscedasticity.

OLS estimator doesn't work if $E[u_i] \neq 0$

\therefore Endogeneity problem \Rightarrow OLS is biased.

$$\because \hat{\beta}_2 = \sum w_i y_i$$

$$\begin{aligned} \text{var}[\hat{\beta}_2] &= E[\hat{\beta}_2 - \beta_2]^2 = E[\sum w_i (\beta_1 + \beta_2 x_i + u_i) - \beta_2]^2 \\ &= E[\beta_1 \sum w_i + \beta_2 \sum w_i x_i + \sum w_i u_i - \beta_2]^2 \\ &= E[\beta_1 \sum w_i^0 + \beta_2 \sum w_i^0 x_i + \sum w_i u_i - \beta_2]^2 \\ &= E[\beta_2 + \sum w_i u_i - \beta_2]^2 = E\left[\sum_{i=1}^n w_i u_i\right]^2 \\ &= E[w_1 u_1 + \dots + w_n u_n]^2 = E[w_1^2 u_1^2 + \dots + w_n^2 u_n^2 \\ &\quad + 2w_1 w_2 u_1 u_2 + \dots] \\ &= E\left[\sum w_i^2 u_i^2 + \sum_{i \neq j} \sum w_i w_j u_i u_j\right]. \end{aligned}$$

$$\therefore \sum w_i^2 E[u_i^2] + \sum_{i \neq j} \sum w_i w_j E[u_i u_j].$$

$\downarrow \sigma^2$ C.R.M assumption $\downarrow 0$

$$\text{var}[\hat{\beta}_2] = \sum w_i^2 \sigma^2$$

$$\therefore \text{var}[\hat{\beta}_2] = \frac{\sigma^2}{\sum x_i^2}$$

$$\begin{aligned}
\text{Var}[\hat{\beta}_1] &= E[(\hat{\beta}_1 - \beta_1)^2] = E[(\bar{Y} - (\hat{\beta}_2 \bar{X} + \beta_0) - \beta_1)^2] \\
&= E[(\beta_0 + \beta_2 \bar{X} + \bar{u}) - (\hat{\beta}_2 \bar{X} + \beta_0)]^2 \\
&= E[(\beta_0 - (\hat{\beta}_2 - \beta_2) \bar{X} + \bar{u} - \beta_1)^2] \\
&= E[(\hat{\beta}_2 - \beta_2)^2 \bar{X}^2 + \bar{u}^2 - 2(\hat{\beta}_2 - \beta_2) \bar{X} \bar{u}] \\
&= \bar{X}^2 \text{Var}[\hat{\beta}_2] + E[\bar{u}^2] - 2\bar{X} E[(\hat{\beta}_2 - \beta_2) \bar{u}] \quad \text{Var}[u] = \sigma^2 \\
&= \frac{\bar{X}^2 \sigma^2}{\sum x_i^2} + E[\bar{u}^2] - 2\bar{X} E[(\hat{\beta}_2 - \beta_2) \bar{u}] \quad \text{Var}[\bar{u}] = \frac{\sigma^2}{n}.
\end{aligned}$$

$$\begin{aligned}
E[(\hat{\beta}_2 - \beta_2) \bar{u}] &= E[(\hat{\beta}_2 - \beta_2) \frac{\sum u_i}{n}] = \frac{1}{n} E[\sum w_i u_i \sum u_i] \\
&= \frac{1}{n} E[\sum w_i u_i^2 + \sum_{i \neq j} (w_i + w_j) u_i u_j] \\
&= \frac{1}{n} \sum w_i E[u_i^2] + \underbrace{\sum_{i \neq j} (w_i + w_j)}_{\sigma^2} E[u_i u_j] = 0 \\
&= \frac{1}{n} \sum w_i \sigma^2 = 0.
\end{aligned}$$

$$\begin{aligned}
E[\bar{u}^2] &= E\left[\frac{\sum u_i}{n}\right]^2 = \frac{1}{n^2} E[\sum u_i^2] \\
&= \frac{1}{n^2} E[\sum u_i^2 + \sum_{i \neq j} u_i u_j] \\
&= \frac{1}{n^2} [\sum u^2 + 0] = \frac{n \sigma^2}{n^2} = \frac{\sigma^2}{n}
\end{aligned}$$

$$\therefore \text{Var}[\hat{\beta}_1] = \frac{\bar{X}^2 \sigma^2}{\sum x_i^2} + \frac{\sigma^2}{n} = \sigma^2 \left[\frac{1}{n} + \frac{(\sum x_i)^2}{n \sum x_i^2} \right]$$

Hypothesis testing in regression models

$$\begin{aligned} E[u_i | x_i] &= 0 \\ E[u_i^2 | x_i] &= \sigma^2 \\ E[u_i u_j | x_i x_j] &= 0 \\ E[u_i x_i] &= 0 \end{aligned}$$

$$E[\hat{\beta}_2] = \beta_2 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{unbiasedness.}$$

$$E[\hat{\beta}_1] = \beta_1$$

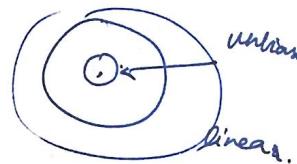
$$\text{var}[\hat{\beta}_2] = \frac{\sigma^2}{\sum x_i^2}$$

$$\text{var}[\hat{\beta}_1] = \sigma^2 \left[\frac{1}{n} - \frac{\bar{x}}{\sum x_i^2} \right]$$

OLS is the best estimator

only among linear unbiased estimators.

statistical questions :
 ① Estimation ? $\xrightarrow{\text{PT}}$
 ② Inference ? $\xrightarrow{\text{Interval + Hyp. Testing.}}$



Interval Estimation / Hypothesis Testing.

$$u_i \sim N(0, \sigma^2)$$

$\xrightarrow{\text{assumption.}}$

$$y_i = \beta_1 + \beta_2 x_i + u_i$$

$$y_i \sim N(\beta_1 + \beta_2 x_i, \sigma^2) \quad \text{why? because you treat the } x_i \text{'s as fixed.}$$

$$w_i = \frac{x_i}{\sum x_i^2}$$

$$\hat{\beta}_2 = \sum w_i u_i$$

$$\hat{\beta}_2 \sim N(\beta_2, \frac{\sigma^2}{\sum x_i^2})$$

You make assumption only on unknowns, viz, u_i 's.

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$$

$$\hat{\beta}_1 \sim N\left(\beta_1, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}}{\sum x_i^2} \right) \right)$$

$\hat{\beta}_2$ here.

Pivot \rightarrow should be a function of the parameters that you are trying to estimate, but the distribution should be independent of the parameters.

$$Q = \frac{\hat{\beta}_2 - \beta_2}{\sqrt{\sigma^2 / \sum x_i^2}} \sim N(0, 1).$$

$\sigma^2 \rightarrow \text{var of } u.$
 $\& u \text{ is not observed.}$

u is not observed in the sample. Therefore you can't find 's' i.e. the sample variance of u .

Thus, you instead use the counterpart of u in the sample, which is, i.e.

$$\hat{\sigma}^2 = \frac{\sum (e_i - \bar{e})^2}{n-1} \quad \text{But } \bar{e} = 0 \quad \text{OLS property.}$$

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-1}$$

n random. You get e_i by doing:
 $e_i = y_i - (\hat{\beta}_1 + \hat{\beta}_2 x)$

To get e_i you have to estimate $\hat{\beta}_1$ & $\hat{\beta}_2$.

$\hat{\beta}_1$ & $\hat{\beta}_2 \Rightarrow$ two parameters are estimated
 in regression.
 \therefore You lose 2 degrees of freedom.

$$\therefore \hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-2}$$

$$\frac{Z}{\sqrt{n/m}} \sim t(m).$$

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1)$$

$$\therefore \frac{\hat{\beta}_2 - \beta_2}{\sqrt{\hat{\sigma}^2 / \sum x_i^2}} \sim N(0, 1)$$

$$\frac{(n-2) \hat{\sigma}^2}{\sigma^2} \sim \chi^2(n-2)$$

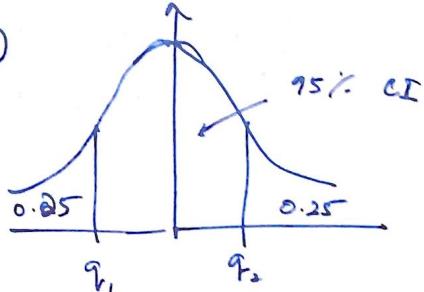
$$\therefore \frac{\frac{\hat{\beta}_2 - \beta_2}{\sqrt{\hat{\sigma}^2 / \sum x_i^2}}}{\frac{\sqrt{(n-2) \hat{\sigma}^2}}{\sigma^2 (n-2)}} \sim t(n-2)$$

$$\boxed{\frac{\hat{\beta}_2 - \beta_2}{\sqrt{\hat{\sigma}^2 / \sum x_i^2}} \sim t(n-2)}$$

means std dev.
 called SE because it
 is an estimated
 value.

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}_1^2 \left[\frac{1}{m} + \frac{\bar{x}^2}{\sum x_i^2} \right]}} \sim t(n-2)$$

$$P[q_1 < Q < q_2] = 0.95$$



$$P\left[q_1 < \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} < q_2\right] = 0.95$$

$$P\left[q_1 SE(\hat{\beta}_1) < \hat{\beta}_1 - \beta_1 < q_2 SE(\hat{\beta}_1)\right] = 0.95$$

Symm. dist.

$$P[-q_1 SE(\hat{\beta}_1) < \hat{\beta}_1 - \beta_1 < q_1 SE(\hat{\beta}_1)] = 0.95$$

$CI = \hat{\beta}_1 \pm q_1 SE(\hat{\beta}_1)$
$CI = \hat{\beta}_2 \pm q_1 SE(\hat{\beta}_2)$

"q" value comes from
 the t-table at
 $(n-2)$ dof.

If u_i takes ∞ binary values $\{0, 1\}$, the normality assumption cannot be made.

- You go for different functional forms (Now?)

- Sometimes changing the functional forms to log often makes a skewed distribution \rightarrow normal.

Hypothesis testing.

- Question: Does x have any effect on y ?

~~On slope~~ If not the corresponding β should be zero.

$H_0: \beta_2 = 0$	$\Rightarrow x$ has NO EFFECT on y .
$H_1: \beta_2 \neq 0$	

$y_i = f(L, K)$. Assume CRS. Log func. form

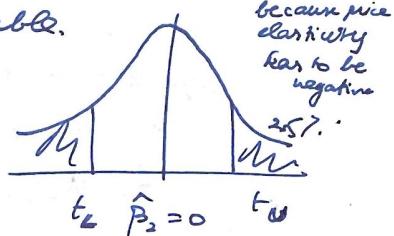
$$y_i = \beta_0 + \beta_1 L^\alpha K^\beta, \quad \alpha + \beta = 1$$

$t_{\text{crit}} \leftarrow 5\% \text{ of } 1\% \rightarrow$ got from the table.

$$t_{\text{calc}} = \frac{\hat{\beta}_2 - \beta_2}{\text{SE}(\hat{\beta}_2)} \stackrel{H_0}{=} \frac{\alpha}{\text{SE}(\hat{\beta}_2)}$$

If t_{calc} falls within the acceptance region, accept H_0 .

Demand
Elasticity
Chkd .4
Corporates want
to know this
 $H_1: \beta_2 < 0$
because price
elasticity has to be
negative



~~On Intercept~~

$$H_0: \beta_1 = 0$$

$$y_i = \beta_1 + \beta_2 x_i + u_i$$

$$\bar{Y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x}$$

$$\hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{x}$$

If $\beta_1 = 0 \rightarrow$ test for \bar{Y} .

$$R^2 = \frac{ESS}{TSS}$$

OLS property.

$$\sum y_i^2 = \hat{\beta}_2^2 \sum x_i^2 + \sum u_i^2.$$

$$TSS = ESS + RSS$$

If you don't have the intercept, this formula is violated.
 \rightarrow all positive numbers

t_{act}

$$t_{\text{calc}} = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)}$$

Say, $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

- ⇒ Check for significance of the parameter.
- ⇒ Statistical significance matters more than the values of parameters.
- ⇒ Although the eqn suggest $\Delta x \Rightarrow \Delta y$. Doesn't mean anything ~~actual~~ unless we know the statistical significance.

Say, $y_i = 4 + 0.17x_i + \epsilon_i$

- ⇒ Doesn't mean that x_i has a small influence on y_i .
- ⇒ It depends rather on the significance (again). This is because the values are governed by the ~~number~~ units in which y_i & x_i are measured.

Say, $y_i = 4 + 0.009x_i + \epsilon_i$

- ✗ $\$4 \text{ GDP} \Rightarrow \text{life expectancy increases by } 0.009 \text{ years}$
- ✓ $1000 \$ \text{ GDP} \Rightarrow \text{life expectancy increases by } 9 \text{ years}$

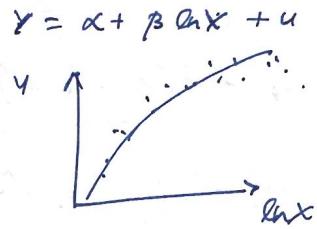
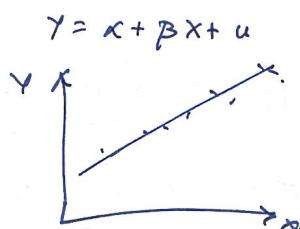
Say, $y_i = \text{Inc} + 0.009x_i + \epsilon_i$
doesn't matter

- There are very few cases where intercept actually matters. However you can't eliminate the intercept from regression.

If you take the GDP of the entire population of all countries & regress on life expectancy

→ Is it the PRF or SRF data.

→ Ans \Rightarrow SRF because,



$Y \rightarrow$ life expectancy in years

$x \rightarrow \ln(\text{GDP per capita})$.

Why non linear? → At lower levels of GDP per capita - adding income increases life exp significantly.
 $Y = \alpha + \beta \ln x + u$. However at higher levels of GDP, a saturation is reached.

$$\frac{\partial Y}{\partial x} = \beta \frac{1}{x}$$

absolute change

$$\frac{\partial Y}{\partial x} \times 100 = \frac{\beta}{100}$$

relative change.

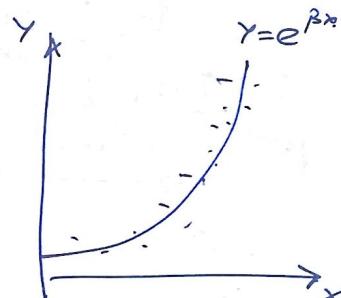
i.e. $\Delta \ln x$ leads to $\beta/100$ units change in Y .

$$\ln Y = \alpha + \beta x + u$$

$$\frac{1}{Y} \frac{\partial Y}{\partial x} = \beta$$

$$\frac{\partial Y / Y \times 100}{\partial x} = \beta \times 100$$

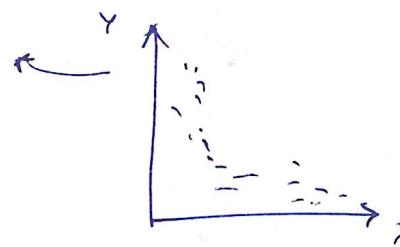
1 unit $\Delta \ln x$ leads $\beta \times 100\%$ change in Y .



$$\ln Y = \alpha + \beta \ln x + u$$

$$\frac{1}{Y} \frac{\partial Y}{\partial x} = \beta \frac{1}{x}$$

$$\frac{\partial Y / Y \times 100}{\partial x / x \times 100} = \beta$$



On $y = \beta \ln x$.

$$e^{\ln y} = e^{\beta \ln x}$$

$$y = e^\beta + e^{\ln x}$$

$$y = e^\beta + x ?$$

You can compare R^2 between these

$$\begin{cases} Y = \alpha + \beta x + u \\ Y = \alpha + \beta \ln x + u. \end{cases}$$

$$Y = \alpha + \beta x + \gamma z + \delta (xz) + u$$

$$\frac{\partial Y}{\partial x} = \beta + \gamma z \leftarrow$$

nonlinear effect of x on Y .

effect depends on the value of z .

Linear Regression with 2 independent variables

$$PRF \quad Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u_i$$

$$\sum e_i^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})^2$$

$$SRF \quad Y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + e_i$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_1} = 0 \Rightarrow 2 \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-1) = 0 \Rightarrow \boxed{\sum e_i = 0}$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_2} = 0 \Rightarrow 2 \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-x_{2i}) = 0 \Rightarrow \boxed{\sum e_i x_{2i} = 0}$$

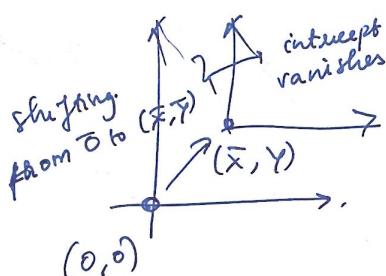
$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_3} = 0 \Rightarrow 2 \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-x_{3i}) = 0 \Rightarrow \boxed{\sum e_i x_{3i} = 0}$$

3 normal equations \rightarrow

$$\begin{cases} \sum e_i = 0 \\ \sum e_i x_{2i} = 0 \\ \sum e_i x_{3i} = 0 \end{cases} \quad \Rightarrow \quad \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \hat{\beta}_2 \bar{x}_2 + \hat{\beta}_3 \bar{x}_3$$

3dim ~
Regression plane

k dim \Rightarrow hyperplane.



$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + e_i$$

$$\bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \hat{\beta}_2 \bar{x}_2 + \hat{\beta}_3 \bar{x}_3$$

$$(\bar{Y} - y_i) = \hat{\beta}_2 (x_{2i} - \bar{x}_2) + \hat{\beta}_3 (x_{3i} - \bar{x}_3) + e_i$$

$$\boxed{y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + e_i}$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_2} = 0 \Rightarrow \sum y_i x_{2i} = \hat{\beta}_2 \sum x_{2i}^2 + \hat{\beta}_3 \sum x_{2i} x_{3i}$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_3} = 0 \Rightarrow \sum y_i x_{3i} = \hat{\beta}_2 \sum x_{2i} x_{3i} + \hat{\beta}_3 \sum x_{3i}^2$$

$$\sum e_i^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})^2$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_2} = 2 \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-x_{2i}) = 0$$

$$\cancel{\sum} (x_{2i} y_i - \hat{\beta}_1 x_{1i}^2 - \hat{\beta}_2 x_{2i} x_{3i}) = 0$$

$$\Rightarrow \sum x_{2i} y_i = \cancel{\sum} \hat{\beta}_1 x_{1i}^2 + \sum \hat{\beta}_2 x_{2i} x_{3i}$$

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_3} = 0 \Rightarrow 2 \sum (x_{3i} y_i - \hat{\beta}_1 x_{1i}^2 - \hat{\beta}_2 x_{2i} x_{3i}) = 0$$

$$\begin{aligned}\sum y_i x_{2i} &= \hat{\beta}_2 \sum x_{2i}^2 + \hat{\beta}_3 \sum x_{2i} x_{3i} \\ \sum y_i x_{3i} &= \hat{\beta}_2 \sum x_{2i} x_{3i} + \hat{\beta}_3 \sum x_{3i}^2\end{aligned}$$

$y = \alpha(x) + \beta(x)$
 $\alpha = \alpha(x) + \beta(x)$

$$\hat{\beta}_2 = \frac{\begin{vmatrix} \sum y_i x_{2i} & \sum x_{2i} x_{3i} \\ \sum y_i x_{3i} & \sum x_{3i}^2 \end{vmatrix}}{\begin{vmatrix} \sum x_{2i}^2 & \sum x_{2i} x_{3i} \\ \sum x_{2i} x_{3i} & \sum x_{3i}^2 \end{vmatrix}} = \frac{\sum y_i x_{2i} \sum x_{3i}^2 - \sum x_{2i} x_{3i} \sum y_i x_{3i}}{\sum x_{2i}^2 \sum x_{3i}^2 - (\sum x_{2i} x_{3i})^2}$$

$$\hat{\beta}_2 = \frac{\sum y_i x_{2i} \sum x_{3i}^2 - \sum x_{2i} \sum x_{3i} \sum y_i x_{3i}}{\sum x_{2i}^2 \sum x_{3i}^2 - (\sum x_{2i} x_{3i})^2}$$

$$\hat{\beta}_2 = \frac{\frac{\sum y_i x_{2i} \sum x_{3i}^2}{\sum x_{2i}^2 \sum x_{3i}^2} - \frac{\sum x_{2i} \sum x_{3i} (\sum y_i x_{3i})}{\sum x_{2i}^2 \sum x_{3i}^2}}{\frac{\sum x_{2i}^2 \sum x_{3i}^2}{\sum x_{2i}^2 \sum x_{3i}^2} - \frac{\sum x_{2i} \sum x_{3i}^2}{\sum x_{2i}^2 \sum x_{3i}^2}}$$

$$\hat{\beta}_2 = \frac{\ell_{12} - \ell_{13} \ell_{32}}{1 - \ell_{32} \ell_{23}}$$

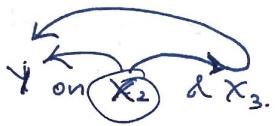
$\ell_{12} \rightarrow \text{simple regression of } y \text{ on } x_2.$

$$\ell_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{\sum (x_i - \bar{x})^2}{n}} \sqrt{\frac{\sum (y_i - \bar{y})^2}{n}}} \leftarrow \text{cov} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}} \leftarrow \text{std dev.}$$

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} \quad \hat{\ell} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}} = \frac{\sqrt{\sum y_i^2}}{\sqrt{\sum x_i^2}}$$

$$\hat{\sigma} = \frac{\rho \sqrt{\sum y_i^2}}{\sqrt{\sum x_i^2}} = \frac{\rho s_y}{s_x}$$

$$\hat{\beta}_2 = \frac{b_{12} - b_{13} b_{32}}{1 - b_{32} b_{23}} = b_{12} \text{ if } \rho_{x_2, x_3} = 0$$



If X_2 & X_3 are not correlated,
then 2 var regression & a 3 var
regression will give the same
coefficients.

2 var reg \Rightarrow both effects taken
into account

3 var reg \Rightarrow removing the effect of X_3 on X_2 .
i.e. the $b_{13} b_{32}$ term.

assuming if all the effects are true.

Then $b_{12} > \hat{\beta}_2 \Rightarrow$ you overestimate by
performing simple regression.

$$\hat{\beta}_2 = \frac{b_{12} - b_{13} b_{32}}{1 - b_{32} b_{23}} = \frac{b_{12} \frac{s_1}{s_2} - \rho_{13} \frac{s_1}{s_3} \rho_{32} \frac{s_3}{s_2}}{1 - \rho_{32} \rho_{23} \frac{s_3}{s_2} \cdot \frac{s_2}{s_3}}$$

$$\hat{\beta}_2 = \frac{\rho_{12} \frac{s_1}{s_2} - \rho_{13} \rho_{32} \frac{s_1}{s_2}}{1 - \rho_{32} \rho_{23}}$$

$$\hat{\beta}_2 = \frac{s_1}{s_2} \left[\frac{\rho_{12} - \rho_{13} \rho_{32}}{1 - \rho_{32} \rho_{23}} \right] = \left(\frac{\rho_{12} - \rho_{13} \rho_{23}}{1 - \rho_{23}^2} \right) \frac{s_1}{s_2}$$

If corr x_2, x_3 is high $(1 - \rho_{23}^2)$ is very small,
it just inflates the effect much more

Regress Y on X_3

\rightarrow take residuals $\rightarrow e_{13i}$

Regress Y on X_2

\rightarrow take residuals $\rightarrow e_{23i}$

Regress e_{13i} on e_{23i}

$$\text{slope} = \frac{\sum e_{13i} e_{23i}}{\sum e_{23i}^2} = \frac{\sum (y_i - b_{13} x_{3i}) (e_{23i})}{\sum e_{23i}^2}$$

$$\text{slope} = \frac{\sum y_i e_{23i} - \sum b_{13} x_{3i} e_{23i}}{\sum e_{23i}^2} \rightarrow 0 \text{ because of the normal equation.}$$

$$\begin{aligned}
 \text{slope} &= \frac{\sum y_i e_{23i}}{\sum e_{23i}^2} = \frac{\sum y_i (x_{2i} - b_{23} x_{3i})}{\sum (x_{2i} - b_{23} x_{3i})^2} \\
 &= \frac{\sum y_i x_{2i} - b_{23} \sum y_i x_{3i}}{\sum x_{2i}^2 + b_{23}^2 \sum x_{3i}^2 - 2 b_{23} \sum x_{2i} x_{3i}} \\
 &= \frac{\sum y_i x_{2i} - \frac{\sum x_{2i} x_{3i}}{\sum x_{3i}^2} \sum y_i x_{3i}}{\sum x_{2i}^2 + \left(\frac{\sum x_{2i} x_{3i}}{\sum x_{3i}^2} \right)^2 \sum x_{3i}^2 - 2 \frac{\sum x_{2i} x_{3i} \sum x_{2i} x_{3i}}{\sum x_{3i}^2}} \\
 &= \frac{\sum y_i x_{2i} \sum x_{3i}^2 - \sum x_{2i} x_{3i} \sum y_i x_{3i}}{\sum x_{3i}^2} \\
 &\quad \frac{\sum x_{2i}^2 \sum x_{3i}^2 + (\sum x_{2i} x_{3i})^2 - 2(\sum x_{2i} x_{3i})^2}{\sum x_{3i}^2}
 \end{aligned}$$

$$\hat{\beta}_{12,3} = \frac{\sum y_i x_{2i} \sum x_{3i}^2 - \sum x_{2i} x_{3i} \sum y_i x_{3i}}{\sum x_{2i}^2 \sum x_{3i}^2 + (\sum x_{2i} x_{3i})^2}$$

<sup>neg of 1 on 2
keeping 3 const</sup>

$$\text{Var}(\hat{\beta}_2) = \frac{\sigma^2}{\sum x_i^2} \quad \text{Var}(\hat{\beta}_r) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}}{\sum x_i^2} \right]$$

$$\begin{aligned}
 \text{Var}(E[\hat{\beta}_{12,3}]) &= E \left[\frac{\sum e_{13i} e_{23i}}{\sum e_{23i}^2} \right] = E \left[\sum (y_i - b_{13} x_{3i}) e_{23i} \right] \xrightarrow{\text{normal equations}} \\
 &= E \left[\frac{\sum (y_i - b_{13} x_{3i})^0 (e_{23i})}{\sum e_{23i}^2} \right] \\
 &= E \left[\frac{\sum y_i e_{23i}}{\sum e_{23i}^2} \right] = E \left[\frac{\sum (\beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i) e_{23i}}{\sum e_{23i}^2} \right] \\
 &= E \left[\frac{\beta_1 \sum e_{23i}^0 + \beta_2 \sum x_{2i} e_{23i} + \beta_3 \sum x_{3i} e_{23i} + \sum u_i e_{23i}}{\sum e_{23i}^2} \right]
 \end{aligned}$$

you can change
one of the $y_i \rightarrow y_i^0$
or $x_i \rightarrow x_i^0$

$$E[\hat{\beta}_{12.3}] = E\left[\frac{\beta_1 \cdot 0 + \beta_2 \sum x_{2i} e_{23i} + \beta_3 \cdot 0 + \sum u_i e_{23i}}{\sum e_{23i}^2} \right]$$

$$E[\hat{\beta}_{12.3}] = E\left[\beta_2 \frac{\sum x_{2i} e_{23i}}{\sum e_{23i}^2} + \frac{\sum u_i e_{23i}}{\sum e_{23i}^2} \right]$$

$$E[\hat{\beta}_{12.3}] = \beta_2 \frac{\sum x_{2i} e_{23i}}{\sum e_{23i}^2} + \frac{\sum e_{23i} E[u_i]}{\sum e_{23i}^2} \quad \downarrow 0 \text{ by CLRM.}$$

$$E[\hat{\beta}_{12.3}] = \beta_2 \frac{\sum x_{2i} e_{23i}}{\sum e_{23i}^2} \quad \text{on } X$$

$$\sum y_i e_i = \frac{\sum (\hat{y}_i + e_i) e_i}{\sum e_i^2}$$

$$E[\hat{\beta}_{12.3}] = \beta_2 \frac{\sum (\hat{x}_{2i} + e_{23i}) e_{23i}}{\sum e_{23i}^2} = \frac{\sum \hat{y}_i e_i + \sum e_i^2}{\sum e_i^2}$$

$$E[\hat{\beta}_{12.3}] = \beta_2 \frac{\sum \hat{x}_{2i} e_{23i} + \sum e_{23i}^2}{\sum e_{23i}^2}$$

$$E[\hat{\beta}_{12.3}] = \beta_2 \cdot 1 \rightarrow \boxed{E[\hat{\beta}_{12.3}] = \beta_2}$$

unbiasedness

$$\text{Var}[\hat{\beta}_{12.3}] = E\left[\frac{\sum e_{13} e_{23i}}{\sum e_{23i}^2} - \beta_2 \right]^2$$

$$= E\left[\frac{\sum (y_i - \beta_{13} x_{3i}) e_{23i}}{\sum e_{23i}^2} - \beta_2 \right]^2$$

$$= E\left[\frac{\sum y_i e_{23i}}{\sum e_{23i}^2} - \beta_2 \right]^2$$

$$= E\left[\frac{\sum (\beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i) e_{23i}}{\sum e_{23i}^2} - \beta_2 \right]^2$$

$$= E\left[\frac{\sum \beta_1 \cancel{\sum e_{23i}} + \beta_2 \sum x_{2i} e_{23i} + \beta_3 \cancel{\sum x_{3i} e_{23i}} + \cancel{\sum u_i e_{23i}}}{\sum e_{23i}^2} - \beta_2 \right]$$

$$\begin{aligned}
 \text{var}[\hat{\beta}_{12,3}] &= E \left[\beta_2 \frac{\sum x_{2i} e_{23i}^2}{\sum e_{23i}^2} + \frac{\sum u_i e_{23i}}{\sum e_{23i}^2} - \beta_2 \right]^2 \\
 \text{var}[\hat{\beta}_{12,3}] &= E \left[\frac{\sum u_i e_{23i}}{\sum e_{23i}^2} \right]^2 = \frac{1}{(\sum e_{23i}^2)^2} E \left[\sum u_i e_{23i} \right]^2 \\
 &= \frac{1}{(\sum e_{23i}^2)^2} E \left[\sum u_i^2 e_{23i}^2 + \sum_{i \neq j} u_i u_j e_{23i} e_{23j} \right] \\
 &= \frac{1}{(\sum e_{23i}^2)^2} \left[\sum e_{23i}^2 E[u_i^2] + \sum_{i \neq j} E[u_i u_j] \right] \\
 &= \frac{\sigma^2 \sum e_{23i}^2}{(\sum e_{23i}^2)^2} = \frac{\sigma^2}{\sum e_{23i}^2} = \frac{\sigma^2}{(1-R_{23}^2) \sum x_{2i}^2}
 \end{aligned}$$

If $R^2 = R^2$ in 2nd reg. $\rightarrow R^2 \uparrow$.
 $y \rightarrow x_2, x_3$ $1 - R^2$ is very low \rightarrow stdev is very high
 $x_2 \& x_3$ are corr. $\Rightarrow t\text{-stat} \downarrow$ is very less.
 \therefore you would wrongly conclude that..
 $x_2 \& x_3$ have no effect on y .
 But R^2 still is high \rightarrow problem of MULTICOLLINEARITY.

In 2nd regression : $\text{Var} \Rightarrow \frac{\sigma^2}{(1-R_{2,345\dots k}^2) \sum x_{2j}^2}$

Multicollinearity

- not a problem of the population but of the sample.
- SE is high $\Rightarrow \beta$ is low.
- Remedial measure?

y on K & L large firms $\underbrace{high K \& high L}$
 $K \& L$ are correlated.

You end up concluding that $K \& L$ has no effect on y .

Reg. $\frac{Y}{L}$ on $\frac{K}{L}$. But forecast will be correct because R^2 is high.

BUILD YOUR FUNDAMENTALS!