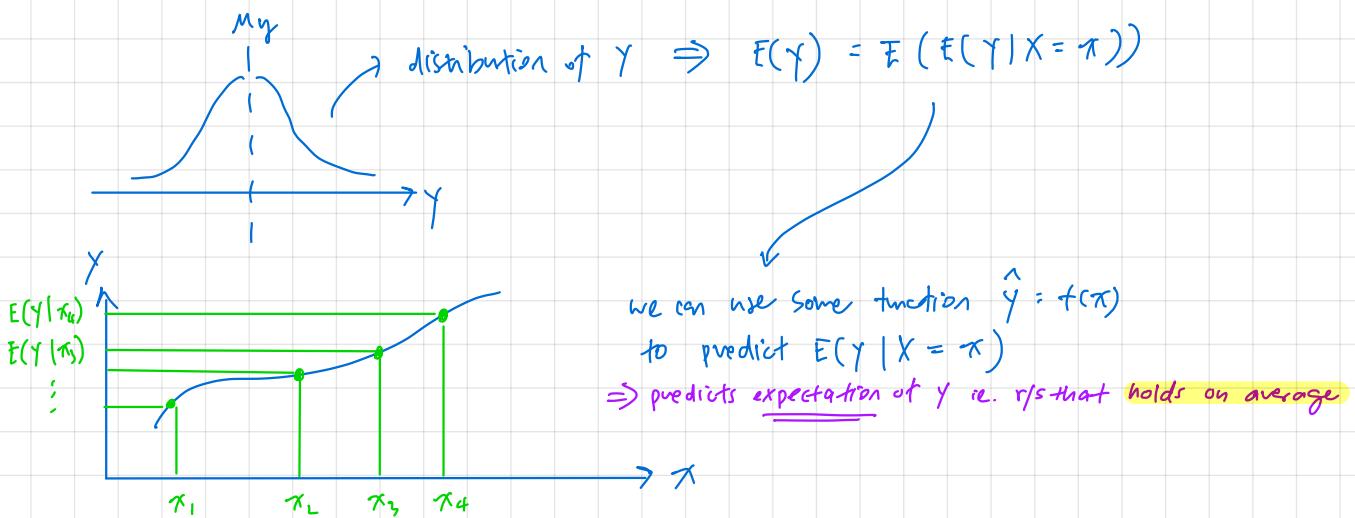


bivariate linear regression

① Linear regression model

- i) introducing another variable \bar{u} in iterated expectation
 - b) suppose we believe some RV Y is related to some RV X .
i.e. for a given value of x , there is some distribution of y .



- 2) using a linear model

- ↳ suppose we think there is some linear population r/s between Y and X

$$E(Y|X) = \beta_0 + \beta_1 X \Rightarrow E(Y_i|X_i) = \beta_0 + \beta_1 X_i \text{ for some observation } i$$

- ↳ X is not the sole determinant of Y , so for each X, Y observation we can abstract out these effects as some random error u

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

dependent variable independent variable

coefficients

- population error
- ↳ incorporates factors responsible for difference between prediction & actual values

② ordinary least squares estimation

- ↳ given some model $f(x)$, we can find the values of coefficients by solving simultaneously the closed form solutions describing the minimisation of sum of squared difference between model & observations

Lg intuition for using SSE: square penalizes large deviations and punishes all deviations, positive or negative

Estimating β_0 and β_1 using OLS estimation

1. we posit that $y_i = \beta_1 x_i + \beta_0 + u_i$

1.1 so we can define estimators $\hat{y}_i = \hat{\beta}_1 x_i + \hat{\beta}_0$, $\hat{u}_i = y_i - \hat{y}_i$

1.2 To perform OLS estimation, we minimize $(Y - \hat{Y})^2$ on aggregate by summing.

so our objective function $Q = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$

2. We can use partial differentiation to derive the relevant equations.

$$Q = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

$$\frac{\partial Q}{\partial \hat{\beta}_0} = -2 \sum y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i = 0 \quad \sum y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum x_i = 0$$

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

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

$$\sum y_i x_i - \hat{\beta}_0 \sum x_i - \hat{\beta}_1 \sum x_i^2 = 0$$

$$\hat{\beta}_1 = \frac{\sum y_i x_i - \hat{\beta}_0 \sum x_i}{\sum x_i^2}$$

$$= \frac{\sum y_i x_i - (\bar{Y} - \hat{\beta}_1 \bar{X}) \sum x_i}{\sum x_i^2} = \frac{\sum y_i x_i - \bar{Y} \sum x_i}{\sum x_i^2 - \bar{X} \sum x_i} \quad \text{re-express as } \bar{X}$$

$$= \frac{\sum y_i x_i - \bar{Y} (n \bar{X})}{\sum x_i^2 - \bar{X} (n \bar{X})} \cdot \frac{1/n}{1/n} = \frac{E(XY) - E(X)E(Y)}{E(X^2) - E(X)^2}$$

$$\hat{\beta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = \frac{S_{xy}}{S_x^2}$$

$$= \frac{\frac{1}{n-1} \sum (x_i - \bar{X})(y_i - \bar{Y})}{\frac{1}{n-1} \sum (x_i - \bar{X})^2}$$

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

$$\frac{1}{n} \sum \hat{y}_i = \bar{Y}$$

$$\sum \hat{u}_i x_i = 0$$

$$\text{Cov}(\hat{u}, X) = 0$$

$$TSS = SSR + ESS$$

③ Assumptions & sampling distributions of $\hat{\beta}_1$ and $\hat{\beta}_0$

i) Assumptions

1. given X_i , the conditional distribution of the population error term has a mean of 0.

$$E(u_i | X_i = x) = 0$$

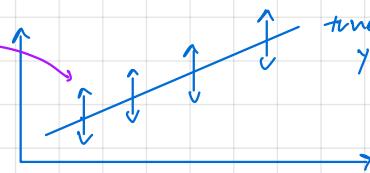


$$\hat{P}_{X_i, u_i} = 0 \text{, but not necessarily converse}$$

more importantly, $P_{X_i, u_i} = 0$ for $\hat{\beta}_1$

- the conditional distribution of the population error term has a mean of 0.

true pop. line
 $y = \beta_1 x + \beta_0$



2. (X_i, Y_i) are independently and identically distributed

3. large outliers are rare

ii) mean and variance of $\hat{\beta}_1$

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

$$\left(\sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2 \right) \xrightarrow{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \xrightarrow{n \rightarrow \infty} N(\beta_1, \frac{\text{Var}(X - \bar{X})n}{\text{Var}(X)^2})$$

$$\text{Var}(\hat{\beta}_1) = \frac{1}{n}$$

$$\frac{\text{Var}((X_i - \bar{X})u_i)}{\text{Var}(X_i)^2} \xrightarrow{n \geq 100 \text{ normally sufficient.}}$$

scaling affects here

↳ Because $\text{Var}(\hat{\beta}_1) \propto \frac{1}{n}$ and $E(\hat{\beta}_1) = \beta_1$, $\hat{\beta}_1$ is a consistent estimator.

$$\lim_{n \rightarrow \infty} \hat{\beta}_1 = \beta_1$$

↳ The larger $\text{Var}(X)$, the smaller $\text{Var}(\hat{\beta}_1)$
 \Rightarrow intuition: longer horizontal, better line



(proof that $\hat{\beta}_1$ is unbiased given OLS assumptions)

$$1. \text{ From partial derivative: } \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\left. \begin{aligned} 1.1 \text{ Also, } y_i &= \beta_0 + \beta_1 x_i + u_i \\ \bar{y} &= \beta_0 + \beta_1 \bar{x} + \bar{u} \end{aligned} \right\} \rightarrow y_i - \bar{y} = \beta_1(x_i - \bar{x}) + (u_i - \bar{u})$$

$$1.2 \text{ so } \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(\beta_1(x_i - \bar{x}) + (u_i - \bar{u}))}{\sum (x_i - \bar{x})^2}$$

$$1.3 \text{ so } \hat{\beta}_1 - \beta_1 = \frac{\sum (x_i - \bar{x})(u_i - \bar{u})}{\sum (x_i - \bar{x})^2}$$

$$\begin{aligned} 1.4 \quad \hat{\beta}_1 - \beta_1 &= \frac{\sum (x_i - \bar{x})u_i - \cancel{\sum (x_i - \bar{x})\bar{u}}}{\sum (x_i - \bar{x})^2} \xrightarrow{\text{cancel}} = \sum x_i u_i = 0 \\ &= \frac{\sum (x_i - \bar{x})u_i}{\sum (x_i - \bar{x})^2} \end{aligned}$$

$$2. \text{ theorem: } E(Y) = E(E(Y|X))$$

$$2.1 \quad E(\hat{\beta}_1 - \beta_1) = E(E(\hat{\beta}_1 - \beta_1 | X_1, X_2, \dots, X_n))$$

↓
estimation conditioned
on observations

$$= E\left(E\left(\frac{\sum (x_i - \bar{x})u_i}{\sum (x_i - \bar{x})^2} | X_1, \dots\right)\right)$$

↓
conditioned on X_1 , is ~~itself~~

$$= E\left(\frac{1}{\sum (x_i - \bar{x})^2} E\left(\sum (x_i - \bar{x})u_i | X_1, \dots, X_n\right)\right)$$

$$\begin{aligned} &E\left(\frac{1}{\sum (x_i - \bar{x})^2} E\left(\sum (x_i - \bar{x})u_i | X_1, \dots\right)\right) \leftarrow \text{conditioned on } X \text{ is } \cancel{\text{itself}} \\ &= E\left(\frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \cdot k\right) \end{aligned}$$

↓
is a constant
 \Rightarrow take out

↓
linear

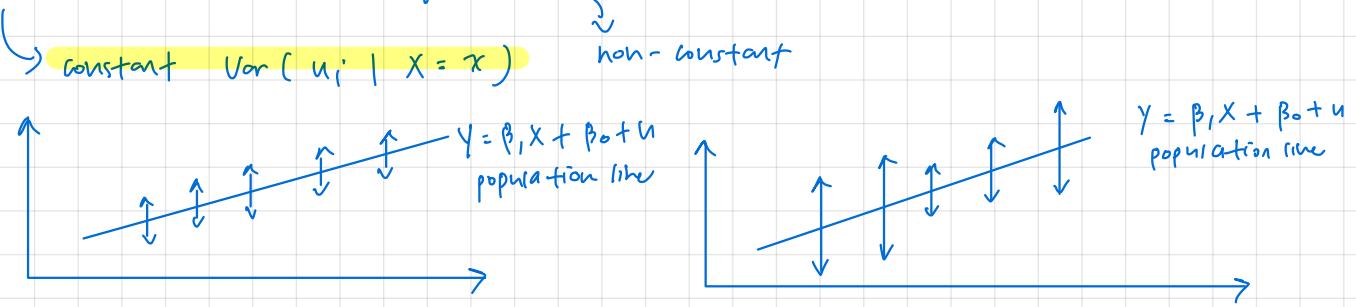
$$E(u_i | X_1, \dots, X_n) = kx_i$$

$$k \frac{\sum (x_i - \bar{x})x_i - (\bar{x} - \bar{x})\bar{x} + (\bar{x} - \bar{x})\bar{x}}{\sum (x_i - \bar{x})^2} = k + \frac{\sum (x_i - \bar{x})\bar{x}}{\sum (x_i - \bar{x})^2}$$

$$= k$$

↳ constant bias

3) homo- and heteroskedasticity



↳ whether homo- or heteroskedastic, need to consider underlying relationships and behaviours being modelled (e.g. income vs. years of education \Rightarrow older likely to be more spread)

↳ implications on computation of variance

regardless : $\text{Var}(\hat{\beta}_1) = \frac{1}{n} \frac{\text{Var}((X_i - \bar{X})u_i)}{\text{Var}(X_i)^2} \Rightarrow \text{SE}(\hat{\beta}_1) = \sqrt{\text{Var}(\hat{\beta}_1)}$

if homoskedastic, then $\text{Var}(\hat{\beta}_1)$ simplifies :

$$\text{Var}(\hat{\beta}_1) = \frac{s_n^2}{S_{xx}^2} \cdot \frac{1}{n} \frac{1}{n-2} \sum \hat{u}_i^2$$

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

↳ robust to heteroskedasticity



→ if homoskedastic, estimates will converge as $n \rightarrow \infty$

④ Hypothesis testing

↳ using OLS estimators to estimate β_1 and β_0 — coefficients that describe a relationship, we can test if given the data, their magnitudes are statistically different from 0 (i.e. no relationship) or some other value

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

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

2. compute b_1 and se

3. compute t -statistic

4. compute p-value

5. If two tailed, compute confidence interval

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)} \sim \text{approx normal}$$

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

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

$$\text{se}(\hat{\beta}_1) = \sqrt{\frac{\text{Var}((x_i - \bar{x})u_i)}{\text{Var}(x_i)^2}}$$

$$\hat{b}_1 = \frac{s_{xy}}{S_{xx}} = \frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

scaling affects here
take out k^2 on var,
 k on $\text{se}(\hat{\beta}_1)$

$$= \sqrt{\frac{\frac{1}{n-2} \sum (x_i - \bar{x})^2 \hat{u}_i^2}{[\frac{1}{n-2} \sum (x_i - \bar{x})]^2}}$$

$$m = \pm t_{n-1, \alpha/2} \cdot \text{se}$$

5) measures of fit

1) R^2 value

↳ measures fraction of variance of Y that is explained by X . unitless. $0 \leq R^2 \leq 1$

↳ in the case of a single regressor, R^2 is the square of the sample correlation between X & Y

⇒ perfect correlation ($R=1$) means all points lie exactly on the line

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

$$R^2 = \frac{\text{explained sum of squares}}{\text{total sum of squares}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

predicted variance given x_i

actual observed variance

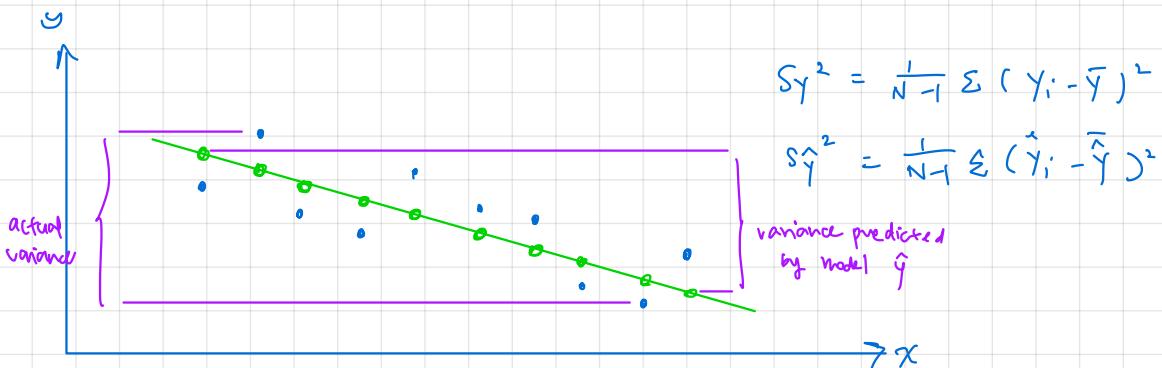
$$= 1 - \frac{\text{sum of squared residuals}}{\text{total sum of squares}} = \frac{\sum \hat{u}_i^2}{\sum (y_i - \bar{y})^2}$$

fraction of variance not explained by regression

↳ explained sum of squares + sum of squared residuals = total sum of squares

↳ intuition: there are two sources of variation in Y : X and some other unmodelled factors.

If $R^2=0$, then the regression of Y on X explains none of the variation of Y and $\text{ESS} = 0$. If $R^2=1$, the regression explains all variation of Y and $\text{ESS} = \text{TSS}$



2) standard error of regression (SER)

↳ measure of the spread of observations about regression (i.e. wrt \hat{y} (in units of y))

⇒ sample std. deviation of OLS residuals

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

$\bar{u} = \frac{1}{n} \sum \hat{u}_i = 0$ by definition of OLS estimation

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

both can normalise wrt. \bar{y}

⑥ interpretation of regression line

i) slope, intercept, predictions

$$\downarrow \quad \downarrow \hat{y} \text{ when } x=0$$

ΔY for every unit change in x

ii) changing units

iii) regression when X is a binary variable

predicting outside of modelled range is not a good idea as it is uninformed

$$\boxed{Y \rightarrow aY \Rightarrow \hat{\beta}_1 \rightarrow a\hat{\beta}_1, \hat{\beta}_0 \rightarrow a\hat{\beta}_0}$$

$$X \rightarrow bX \Rightarrow \hat{\beta}_1 \rightarrow \frac{\hat{\beta}_1}{b}, \hat{\beta}_0 \text{ unchanged}$$

$\hookrightarrow \hat{\beta}_1$ is not so much a slope, but a **translation** brought about by presence of X

$$Y_i = \beta_1 x_i + \beta_0 + u_i \quad \begin{cases} x_i = 0 & Y_i = \beta_0 + u_i \\ x_i = 1 & Y_i = \beta_1 + \beta_0 + u_i \end{cases}$$

\hat{Y} in this case gives the mean value when X is either value.

Report results

$$\hat{Y} = b_0 + b_1 X, R^2 = \underline{\quad}, SER = \underline{\quad}$$

$$(se_{b_0}) \quad (se_{b_1})$$

⑦ Prediction using OLS

(x, Y) are randomly drawn from the same population as (x_i, Y_i)

⑧ omitted variable bias

$$\hat{\beta}_1 = \beta_1 + \frac{\frac{1}{n} \sum (\pi_i - \bar{\pi})(u_i - \bar{u})}{\frac{1}{n} \sum (\pi_i - \bar{\pi})^2} = \beta_1 + \frac{s_{\pi u}}{s_{\pi}^2}$$

if

$$= \beta_1 + \frac{s_{\pi u}}{s_{\pi} \cdot s_{\pi}} \cdot \frac{s_u}{s_u}$$

$$= \beta_1 + \frac{s_{\pi u}}{s_{\pi} s_u} \cdot \frac{s_u}{s_{\pi}}$$

$$E(\hat{\beta}_1) = \beta_1 + \rho_{xy} - \frac{\tau_y}{\tau_x}$$

intuition: bias arises from the violation of $E(u_i | X_i) = 0$
is fundamentally because u_i and X_i are correlated

$E(u_i | X_i) = 0 \rightarrow \rho_{u_i X_i} = 0$, but converse untrue.

$E(u_i | X_i) = k \rightarrow E(\hat{\beta}_1) = \beta_1 + k, E(\hat{\beta}_0) = \beta_0 + k$

(intuition behind ρ_{xy} and omitted variable bias)

1. suppose we have an omitted variable π_i and coefficient β_2 .

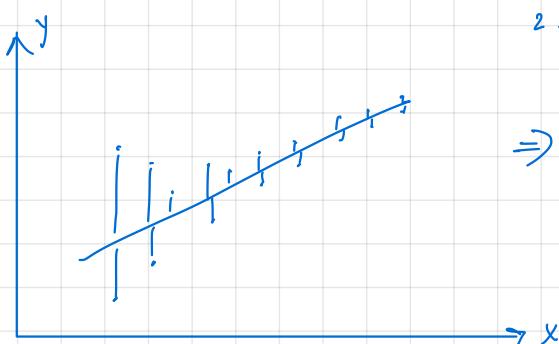
$$y_i = \beta_0 + \beta_1 x_i + (\beta_2 \pi_i + \varepsilon_i)$$

\nearrow true error
 \searrow omitted variable

1.1 we have erroneously modelled it as $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, so any error due to a change in π is encapsulated in our \hat{u}_i when we fit the model.

2. suppose β_2 is negative. suppose also that x_i and π_i are positively correlated. suppose also that β_1 is positive

2.1 so as x_i increases, π_i increases, and because $\beta_1 > 0$ but $\beta_2 < 0$ the correlation between



1. suppose $\rho_{xy} < 0$. As x increases, u decreases.

2. then $E(\hat{\beta}_1) < \beta_1$ - ie. flatter or more negative than it really is

\Rightarrow intuition:

ρ_{xy} arises because

ρ

Notes

1. wait longer \rightarrow rate note of whether it is applied to Y or $X \rightarrow \hat{\beta}_i$, compensate
 \downarrow scale

Multiple regression

① omitted variable bias

In a basic bivariate model, we assume that y is determined by a single variable x . But there might be other variables that also determine y that are correlated with x , muddling our study of the relationship between y and x .

(when does it occur?)

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

(why do omitted variables cause bias?)

1. let us simplify the expression of $\hat{\beta}_1$:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (\bar{x}_i - \bar{x}) y_i - (\bar{x}_i - \bar{x}) \bar{y}}{\sum (\bar{x}_i - \bar{x}) \bar{x}_i - (\bar{x}_i - \bar{x}) \bar{x}}$$

2. suppose the true model is $y = \beta_0 + \beta_1 x + \underbrace{\beta_2 z}_{\text{omitted variable}} + \varepsilon$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(\beta_0 + \beta_1 x_i + \beta_2 z_i + u_i)}{\sum (x_i - \bar{x}) x_i}$$

$$= \beta_1 + \beta_2 \frac{\sum (x_i - \bar{x}) z_i}{\sum (x_i - \bar{x}) x_i} + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x}) x_i}$$

$$= \beta_1 + \beta_2 \frac{\sum (x_i - \bar{x})(z_i - \bar{z})}{\sum (x_i - \bar{x})^2} + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2}$$



$$\frac{s_{xz}}{s_x^2} = \frac{s_{xz} \cdot s_z}{s_x \cdot s_x \cdot s_z} = \frac{s_{xz}}{s_x s_z} \cdot \frac{s_z}{s_x}$$

3. Taking expectation, the omitted variable causes bias in the estimation of $\hat{\beta}_1$:

$$E(\hat{\beta}_1) = \beta_1 + \beta_2 \rho_{xz} \frac{s_z}{s_x} + E\left(\frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} E(u_i | x_i)\right)$$

(how do we observe omitted variable bias in our model?)

1. let us simplify the expression of $\hat{\beta}_1$.

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (\bar{x}_i - \bar{x}) y_i - (\bar{x}_i - \bar{x}) \bar{y}}{\sum (\bar{x}_i - \bar{x})^2}$$

2. suppose the true model is $y = \beta_0 + \beta_1 x + \underbrace{\beta_2 z_i}_{\text{omitted variable}} + \varepsilon$

$$\begin{aligned} 2.1 \text{ Then } \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(\beta_0 + \beta_1 x_i + \beta_2 z_i + u_i)}{\sum (x_i - \bar{x})^2} \\ &= \beta_1 + \frac{\sum (x_i - \bar{x})(\beta_2 z_i + u_i)}{\sum (x_i - \bar{x})^2} \end{aligned}$$

2.2 if we consider that in the model we do not consider z_i , then the $\beta_2 z_i + u_i$ term are represented as a single estimator of

$$\text{more, } \tilde{u} = \beta_2 z_i + u_i$$

$$\text{so } \hat{\beta}_1 = \beta_1 + \frac{\sum (x_i - \bar{x}) \tilde{u}_i}{\sum (x_i - \bar{x})^2} = \beta_1 + \frac{\sum (x_i - \bar{x})(\tilde{u}_i - \bar{\tilde{u}}_i)}{\sum (x_i - \bar{x})^2}$$

3. We can express this in terms of correlation between x_i and \tilde{u}_i .

$$\begin{aligned} 3.1 \quad \beta_1 + \frac{s_{xu}}{s_x^2} &= \beta_1 + \frac{s_{x\tilde{u}}}{s_x^2} \cdot \frac{s_{\tilde{u}}}{s_u} \\ &= \beta_1 + \frac{s_{x\tilde{u}}}{s_u s_x} \cdot \frac{s_{\tilde{u}}}{s_x} \\ &= \beta_1 + r_{x\tilde{u}} \cdot \frac{s_{\tilde{u}}}{s_x} \end{aligned}$$

3.2 Taking expectation, we observe that the correlation between x and \tilde{u} reflects the omitted variable bias.

$$E(\hat{\beta}_1) = \beta_1 + r_{x\tilde{u}} \cdot \frac{s_{\tilde{u}}}{s_x}$$

(2) multiple regression model

↳ functional form: we posit that $y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots + \beta_k x_k + \varepsilon_i$

interpretation

↳ the model describes the relationship between y and regressors $x_1 \dots x_k$ on average.

$$\text{i.e. } E(y | x_1, x_2 \dots x_k) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

↳ the coefficient β_j is the effect on y of a unit change in x_j , holding all other regressors constant

OLS estimator

↳ just as in the bivariate case, we try to minimize the total sum of squared errors.

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \dots + \hat{\beta}_k x_k$$

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

$$S = \sum_{i=1}^n \hat{u}_i^2$$

$$\begin{aligned}\frac{\partial S}{\partial \beta_0} &= 0 \\ \frac{\partial S}{\partial \beta_1} &= 0 \\ &\vdots \\ \frac{\partial S}{\partial \beta_k} &= 0\end{aligned}$$

⇒ solve in linear algebra and calculus to get closed form solutions

(3) concepts of multicollinearity

1) perfect multicollinearity

↳ regressors are said to exhibit perfect multicollinearity when one of the regressors is a perfect linear function of other regressors. e.g. $X_1 = kX_2$

2) imperfect collinearity

↳ independent concept from perfect collinearity

↳ regressors are said to exhibit imperfect multicollinearity when they are highly (but not perfectly) correlated

④ Sampling distributions in multivariable regression

i) Least-squares assumptions

(the standard three)

1. $E(u_i | X_{1i}, X_2i, \dots, X_{ki}) = 0$

2. $(X_{1i}, X_{2i}, \dots, X_{ki}), i \in \{1, \dots, n\}$ are identically & independently distributed

3. Large outliers are unlikely

(no perfect collinearity) \rightarrow if present, likely a logical error in dataset / choice of variables

\hookrightarrow when perfect multicollinearity is present, then it's as if you have >1 instance of a given variable. e.g. $X_2 = X_1 \times 100$ (percentages)

\hookrightarrow in trying to estimate β_j , we intuitively hold all other X constant. Given perfect multicollinearity, it is as if trying to simultaneously hold constant & vary a variable

\hookrightarrow mathematically, it leads to a divide by zero error in the OLS estimator

ii) Estimator distributions

In general, if LSA assumptions hold, sampler is large

$\hookrightarrow \hat{\beta}_j$ is unbiased & consistent: $E(\hat{\beta}_j) = \beta_j$, $\text{Var}(\hat{\beta}_j) \propto \frac{1}{n}$

$\hookrightarrow \hat{\beta}_j$ are well approximated by multivariate normal distribution

$$\hat{\beta}_j \sim N(\beta_j, \text{Var}(\hat{\beta}_j)) \text{ approximately.}$$

\hookrightarrow because we don't know $\text{Var}(\hat{\beta}_j)$, we estimate with $\text{SE}(\hat{\beta}_j)$

$$\hat{\beta}_j \sim t_{n-1}(\beta_j, \text{SE}(\hat{\beta}_j))$$

(imperfect multicollinearity)

\hookrightarrow when two or more regressors are highly correlated, then sampling variance of these variables become very large.

\hookrightarrow intuition: high correlation, move together. In OLS estimation, holding all other variables constant, those regressors have a very small independent spread

$$\Rightarrow \text{Var}(\hat{\beta}_j) \propto \frac{1}{\text{Var}(X_j)} \Rightarrow \text{high sampling variance}$$

$\hookrightarrow \hat{\beta}_j$ will still be unbiased and consistent, though

⑤ Binary multivariate regression

1) dummy variable trap

↳ if there are G_1 binary variables, if there's an intercept in the regression, and all G_1 regressors are included, it will fail because of perfect collinearity

(intuition)

↳ imagine that β_0 is coefficient of a variable $X_0 = 1$.

↳ if many variables are mutually exclusive $\Rightarrow X_j = 1 - \underbrace{X_1 - X_2 - \dots}_{\text{if any is } 1, \text{ then } 0}$

\Rightarrow perfect collinearity between X_0 and X_j

\Rightarrow solution: exclude β_0 , or regress only $G-1$ variables

2) Interpretation of coefficients

1. Suppose K categories, inclusion of intercept & removal of 1 regressor

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 \dots \hat{\beta}_{K-1} X_{K-1i} + \hat{u}_i$$

↳ $E(\hat{y} | X_1=0, X_2=0 \dots) = \hat{\beta}_0 \Rightarrow \hat{\beta}_0$ estimates expected value of y given K^{th} category

↳ $E(\hat{y} | X_1=0, \dots X_j=1, \dots) = \hat{\beta}_0 + \hat{\beta}_j \Rightarrow \hat{\beta}_j$ estimates expected difference in y between j^{th} & K^{th} (base) categories

2. Suppose exclusion of intercept and K regressors

$$\hat{Y}_i = \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} \dots \hat{\beta}_K X_{Ki} + \hat{u}_i$$

↳ $E(\hat{y}_i | X_1=0 \dots X_j=1) = \hat{\beta}_j \Rightarrow \hat{\beta}_j$ estimates expected value of y given one of these categories

⑥ Hypothesis testing

i) single variable tests

same as bivariate case: $t_{n-1} = \frac{\hat{b}_j - b_0}{se(\hat{b}_j)}$

ii) joint hypothesis testing

↳ joint hypothesis: a hypothesis that imposes two or more restrictions on the regression coefficients

$$H_0: \beta_j = k_j \wedge \underbrace{\beta_m = k_m, \dots}_{1 \text{ restriction}} \text{ for a total of } q \text{ restrictions}$$

H_1 : at least 1 of q restrictions do not hold

(why do q t-tests not work?)

↳ if we did q t-tests, we fail to model a distribution that considers simultaneous restrictions \Rightarrow akin to binomial distribution rather than modelling from a single joint distribution

(The F statistic)

↳ we need some statistic with a known distribution that models the joint probability distribution of k random variables, that could be correlated with one another

\Rightarrow we can use the F-statistic

1. The F statistic models the ratio of two χ^2 variables.

$$F \sim \frac{\frac{\chi_m^2}{m}}{\frac{\chi_n^2}{n}}$$

2. we can cleverly aggregate our \hat{p}_j distributions as a single χ^2 random variable.

2.1 $\hat{p}_j \rightarrow N(\beta_j, \text{Var}(\hat{p}_j))$ as $n \rightarrow \infty$

2.2 if we square them: $\hat{p}_j^2 \sim \chi_1^2$ approximately

2.3 if we sum them and divide by number of restrictions q : $\frac{\chi_q^2}{q}$

2.4 by CLT: $\frac{\chi_m^2}{m} \rightarrow 1$ as $m \rightarrow \infty$. so $\frac{\chi^2}{q} = \frac{\frac{\chi_q^2}{q}}{\frac{\chi_m^2}{m}}$, $m \rightarrow \infty$

\Rightarrow F statistic!

The F-test

1. H_0, H_1, q joint hypotheses

2. calculate q t-statistics.

$$3. f = \frac{\sum t_i^2}{q} / \frac{\chi_n^2}{n} \quad \text{normally we assume } n \text{ is large } \approx \infty$$

4. calculate p-value

} note that even if $q=1$, the F-test is valid. But by convention, t test for single variable, else F test.

3) single restrictions involving multiple coefficients eg. theory suggests $\beta_1 = \beta_2$

↳ not a joint hypothesis test in multiple restriction.

↳ we can do an F test (abstacted by package)

↳ we can also transform the problem into a single variable test by forming desired r/s as a coefficient eg. $H_0: \beta_1 = \beta_2$

1. suppose we believe $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$

$$\begin{aligned} 2. \text{ then we can rewrite as } y_i &= \beta_0 + (\beta_1 - \beta_2)x_{1i} + \beta_2(x_{1i} + x_{2i}) + \varepsilon_i \\ &= \beta_0 + \gamma x_{1i} + \beta_2 v_i, \quad v = x_{1i} + x_{2i} \end{aligned}$$

3. then we can regress and test that $H_0: \gamma = \beta_1 - \beta_2 = 0$
using a t-test.

⑦ measures of fit

(standard error of regression)

↳ estimates standard deviation of error term $\epsilon_i \Rightarrow$ spread of y_i about \hat{y}_i

$$SER = s_{\hat{y}} = \sqrt{s_{\hat{u}}^2}, s_{\hat{u}}^2 = \frac{1}{n-k-1} \sum (\hat{u}_i - \bar{\hat{u}})^2 = \frac{1}{n-k-1} \sum_{i=1}^n \hat{u}_i^2 = \frac{SSR}{n-k-1}$$

$\bar{\hat{u}} = \frac{1}{n} \sum \hat{u}_i = 0$ no. of regressors

R^L

↳ estimates fraction of sample variance of y_i predicted by regressors.

$$R^L = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{ESS}{TSS} = 1 - \frac{SSR}{TSS}$$

\Rightarrow but R^2 never decreases when you increase no. of regressors, so it inflates how well model is fitting the data

↳ intuition = OLS estimator minimises TSS.

$$\text{i.e. } \min \sum (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots)^2$$

anytime you add a new \hat{b}_j , so long as $b_j \neq 0$, by definition TSS will fall, since minimisation.

\Rightarrow we need a metric that punishes model complexity

(Adjusted R^2)

↳ R^2 that deflates as more regressors are included

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

↑
no. of regressors

↳ things to note about \bar{R}^2 :

1. $\bar{R}^2 < R^2$, because $\frac{n-1}{n-k-1} > 1$

2. Adding a regressor has 2 effects on \bar{R}^2 . \rightarrow SSR \downarrow due to complexity, $\uparrow \bar{R}^2$
 $\frac{n-1}{n-k-1} \downarrow$, $\downarrow \bar{R}^2$

3. \bar{R}^2 can be negative. occurs when regressors reduce SSR by such a small amount that it fails to offset $\frac{n-1}{n-k-1}$

\Rightarrow rule of thumb: negative, fails to explain any variation

RMSE

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

Notes on R^2

1. what R^2 tells you

↳ how much of variability is explained by model

↳ high $R^2 \rightarrow$ regressors are good at predicting y_i

2. what R^2 does not tell you

↳ does not say anything about P_{UX} and omitted variable bias

↳ does not say anything about causality \rightarrow RCTs and hypothesis tests

↳ does not say anything about statistical significance \rightarrow hypothesis tests

③ model design

↳ when estimating a causal effect, we want to include variables of interest and control variables (to avoid omitted variable bias)

↳ all we need is that estimators of coefficients of variables of interest are unbiased and consistent. It does not matter if estimators of coefficients for control variables are biased

relaxing LSAT #1

↳ to study variables of interest (i.e. $\hat{\beta}_1$; unbiased, consistent), we simply need to ensure that u_i is uncorrelated with them.

↳ suppose we have $X_1, X_2 \dots X_K$ variables of interest & $w_1, w_2 \dots w_j$ control variables.

\Rightarrow To ensure $\hat{\beta}_{X_1} \dots \hat{\beta}_{X_K}$ unbiased, we simply need $E(u_i | X_1 \dots X_K, w_1 \dots w_j) = E(u_i | w_1 \dots w_j)$

A general approach

1. specify a base specification (set of variable(s))
 - ↳ should contain the variables of interest & a set of control variables suggested by judgement or theory
2. specify a range of plausible alternative specifications
 - ↳ other possible control variables you can use to observe if your coefficients change when you include them (they shouldn't, at least within CI across diff. specifications)
 - ↳ can also be alternative proxies for control variables
 - e.g. % on income support / yes / no
3. run regressions
 - ↳ pay attention to whether estimates of coefficients of interest, fitting metrics, numerically similar across the different specifications
 - yes ↳ estimates are generally reliable
 - no ↳ probably suffers from omitted variable bias

Non-linear regression

① non-linear models

↳ In general, we can model the average relationship between variables as some generic function and use different methods to estimate them

$$Y = f(\{x_{ij}\}) + u_i$$

$$\hat{f}(Y | \{x_{ij}\}) = f(\{x_{ij}\})$$

Effect on Y of a change in X

↳ In linear specifications, coefficients are intuitively interpretable as $\frac{\Delta Y}{\Delta X_i}$.

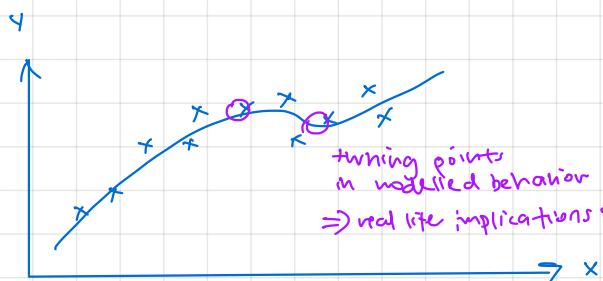
↳ In non-linear specifications, the interpretation depends on functional form

↳ More generally, when we fit non-linear models, we should consider the expected change in Y given some change in X_i from one value to another

$$\Delta \hat{Y} = \hat{f}(x_1 + \Delta x_1, x_2, x_3, \dots) - \hat{f}(x_1, x_2, \dots)$$

⇒ Our fitted function estimates expected change at each value of independent variable

↳ Shape, gradients & predictions at ranges of X are informative about behavior trying to study - individual coefficients less so



Standard errors of estimated effects

↳ In the linear case, $\hat{Y} = \hat{\beta}_0 + \sum \hat{\beta}_i x_i \dots$

$$\Rightarrow \Delta \hat{Y} = \hat{\beta}_i (\Delta x_i) \Rightarrow SE(\Delta \hat{Y}) = SE(\hat{\beta}_i)$$

↳ In the non-linear case, it depends on the functional form.

$$\text{eg. } \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 \Rightarrow \Delta \hat{Y} = \hat{\beta}_1 (\Delta x) + \hat{\beta}_2 ((x + \Delta x)^2 - x^2)$$

$$\Rightarrow SE(\Delta \hat{Y}) = SE(\underbrace{\Delta x(\hat{\beta}_1) + ((x + \Delta x)^2 - x^2)(\hat{\beta}_2)}_{\text{depends on functional form}})$$

② A general approach to fitting non-linear models

1. use economic theory to identify a shape between variables \Rightarrow choose functional form(s)
eg. gradient, turning points, plateaus, etc.
2. estimate parameters using OLS or some other techniques
 \hookrightarrow note that in many cases we do a non-linear transformation on the input variables, and exponentiating them allows us to use multi-variable OLS estimation on them
3. Evaluate the model

3.1 F-test: non-linear vs. linear model

3.2 T-test: determine degree of non-linearity

3.3 graph fit: see how well it fits the data

residuals

data about line

4. compare statistically significant models

\hookrightarrow when using \bar{R}^2 to compare models, make sure the basis is same eg. $\ln(y) - \ln(x)$ vs. $\ln(y) - x$

③ Polynomial models

1) functional form and OLS estimation

$\ln(y) - x$ } comparing
vs. $y - x$ } diff things!

$\Rightarrow \bar{R}^2$'s capture of variation not comparable

$$y_j = \beta_0 + \beta_{1j} x_{ij} + \beta_{2j} x_{ij}^2 + \dots + \beta_{rj} x_{ij}^r + u_j$$

to the r^{th} degree

\hookrightarrow observe that polynomial regression is identical to multiple regression, except that instead of independent variables, the variables are powers
 \Rightarrow we can run OLS estimation as in multiple regression

\Rightarrow interpret as we do for non-linear models

2) evaluating and choosing degrees

1. based on theory, pick a maximum value of r and fit a model

2. test for non-linearity: if all X_i^r , $r > 1$ terms are needed using F-test

3. use a t-test to test X^r if single variable, use F-test to test X_i^r .
if you do not reject it, remove power r terms and run the regression
 $\tilde{n} = r - 1$

\hookrightarrow remove to avoid imperfect multicollinearity, since statistically insignificant terms will move closely w/ others to fit noise

4. compare all statistically significant models w/ \bar{R}^2

④ 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)

$$y_j = \beta_0 + \sum_i \beta_i \ln(x_{ij}) + u_j$$

$$y_i + \Delta y_i = \beta_0 + \beta_1 \ln(x_i + \Delta x)$$

$$\Delta y_i = \beta_1 [\ln(x_i + \Delta x) - \ln(x)] \approx \beta_1 \left(\frac{\Delta x_i}{x_i} \right) \approx \beta_1 (0.01) K$$

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

Case	Population Regression Model	Interpretation of β_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)

$$\ln(y_j) = \beta_0 + \sum_i \beta_i x_{ij} + u_j$$

$$\ln(y + \Delta y) = \beta_0 + \beta_1 (x + \Delta x)$$

$$\ln(y + \Delta y) - \ln(y) = \beta_1 \Delta x \approx \frac{\Delta y}{y}$$

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

(log-log model)

$$\ln(y_j) = \beta_0 + \sum_i \beta_i \ln(x_{ij}) + u_j$$

$$\ln(y + \Delta y) = \beta_0 + \beta_1 \ln(x + \Delta x)$$

$$\ln(y + \Delta y) - \ln(y) = \beta_1 [\ln(x + \Delta x) - \ln(x)]$$

$$\frac{\Delta y}{y} \approx \beta_1 \frac{\Delta x}{x}$$

\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

3) fitting logarithmic models

↪ observe that as in the polynomial case, we are simply doing a transformation of the variables, and they effectively take the form of multiple regression \Rightarrow OLS estimation

⑤ functions that are non-linear in their parameters

↪ so far, we have expressed non-linearity by transforming the input variables, but the estimated parameters β_i are still linear. we like them because they are conveniently decomposable into a linear function and β_i can be estimated by OLS

\Rightarrow how about non-linear parameters?

(general) functions that are non-linear in parameters

$$y_i = f(x_1, x_2, \dots, \underbrace{\beta_0, \beta_1, \dots}_{\text{parameters also} \atop \text{inside}}) + u_i$$

(non-linear OLS estimation)

$$Q = \sum [y_i - f(x_1, x_2, \dots, b_0, \dots, b_1)]^2 \Rightarrow \text{minimize}$$

↪ in general, if there is no closed form solution, then we estimate by numerical solving (guessing until minimum) \Rightarrow like gradient descent

↪ if LSEs hold, the estimators of $\hat{\beta}_i$ are consistent and normally distributed with a large dataset, so we can use t-tests to test their significance

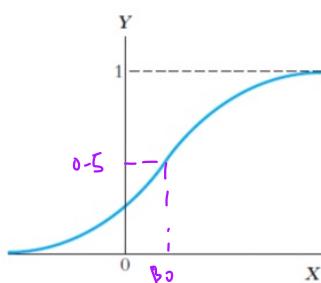
(logistic regression)

$$y_j = \frac{1}{1 + e^{-(\beta_0 + \sum \beta_i x_{ij})}} + u_i$$

↑
shifts laterally

↪ properties:

1. $y \rightarrow 1$ as $x \rightarrow \infty$
2. $y \rightarrow 0$ as $x \rightarrow -\infty$
3. gradient steepest at $x = 0$
3. as $x \rightarrow \infty$ or $-\infty$, $\frac{dy}{dx} \rightarrow 0$



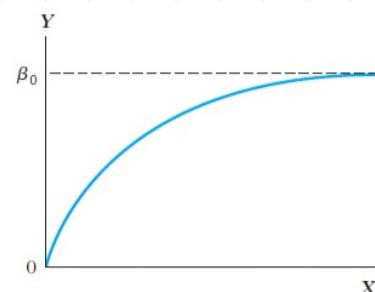
(negative exponential regression)

$$y_j = \beta_0 [1 - e^{-\beta_1 (x_i - \beta_2)}]$$

↑
two parameters per regressor

↪ properties:

1. $y \rightarrow \beta_0$ as $x \rightarrow \infty$
2. $\frac{dy}{dx} > 0$, $\frac{d^2y}{dx^2} < 0$
3. undefined for negative values of x



⑥ Interaction between variables

↳ suppose variables modify each others' effects on y at different values

e.g. at high student - teacher ratio, effects of SES more pronounced

⇒ how can we model this? we can introduce interaction terms and estimate their coefficients between variables we think will interact

↳ is there an interaction?

⇒ we do a t-test on the interaction terms' coefficient(s)

(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. $D_1 = 0 \Rightarrow \beta_1 + \beta_3 d_2$

⇒ effect on Y of change in D_1 now depends on D_2

⇒ β_3 is the difference in effect on Y of D_1 when $D_2 = 1$ vs. $D_2 = 0$

(continuous-continuous)

↳ effectively finer granularity of binary-binary case

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2$$

extra effect from changing both X_1 and X_2
↳ effect of X_1 holding X_2 constant

$$\frac{dy}{dx_1} = \beta_1 + \beta_3 X_2$$

⇒ $\beta_3 X_2$: the expected difference in effect of a unit change in X_1 , given $X_2 = x_L$

Continuous - binary

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

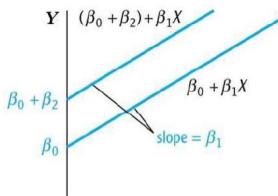
↳ inclusion of binary term allows for different offset

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 D_i$$

$$D_i = 0$$

$$D_i = 1$$

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



(a) Different intercepts, same slope

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

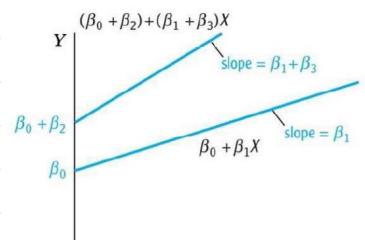
$$y_i = \beta_0 + \beta_1 x_i + \beta_2 D_i + \beta_3 x_i D_i + u_i$$

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

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

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

$$y_i = \beta_0 + \beta_1 x_i$$



(b) Different intercepts, different slopes

⇒ β_3 - difference in effect of unit change in X , for observations where $D_i = 1$ vs. $D_i = 0$

Notes

1. percentages in decimal. But if % change is large, lg approximation will be bad

2.

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. homoskedasticity)

(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

⇒ solutions: include omitted variables, use panel data, instrumental variables or RCT

(incorrect functional form)

↳ 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

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

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

$$= \beta_0 + \beta_1 \tilde{x}_i + [\beta_1(x_i - \tilde{x}_i) + u_i]$$

\uparrow \uparrow
 errors in x_i v_i

↳ errors in y_i

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

$$Y_i = \beta_0 + \beta_1 x_i + u_i + w_i$$

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

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

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

$$2. \text{cov}(\tilde{x}_i, v_i) = \text{cov}(\tilde{x}_i, \beta_1 x_i - \beta_1 \tilde{x}_i + u_i)$$

$$= \text{cov}(\tilde{x}_i, \beta_1 x_i) - \text{cov}(\tilde{x}_i, \beta_1 \tilde{x}_i) + \cancel{\text{cov}(\tilde{x}_i, u_i)}$$

$$= \underbrace{\beta_1 (\text{cov}(\tilde{x}_i, \beta_1 x_i) - \text{std}^2(\tilde{x}_i))}_{\text{in general, } \neq 0}$$

$= 0$, since $E(u_i | x_i) = 0$

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

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

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

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

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

$$\rho_{u_i, x_i} > 0 \quad \rho_{u_i, y_i} < 0$$

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

more formally

$$\text{cov}(x_i, u_i) = 0 \text{ iff. } \mathbb{E}(u_i | x_i) = 0$$

$$1. \text{ cov}(x_i, u_i) = \text{cov}(\gamma_0 + \gamma_1 y_i + v_i, u_i)$$

$$= \cancel{\text{cov}(\gamma_0, u_i)} + \text{cov}(\gamma_1 y_i, u_i) + \cancel{\text{cov}(v_i, u_i)}$$

$$= \text{cov}(\gamma_1 (\beta_0 + \beta_1 x_i + u_i), u_i)$$

$$= \cancel{\text{cov}(\gamma_1 \beta_0, u_i)} + \text{cov}(\gamma_1 \beta_1 x_i, u_i) + \text{cov}(\gamma_1 u_i, u_i)$$

$$= \gamma_1 \beta_1 \text{cov}(x_i, u_i) + \gamma_1 \sigma_u^2$$

$$(1 - \gamma_1 \beta_1) \text{cov}(x_i, u_i) = \gamma_1 \sigma_u^2 \Rightarrow \text{cov}(x_i, u_i) = \frac{\gamma_1 \sigma_u^2}{1 - \gamma_1 \beta_1} \neq 0$$

↳ solution: RCT

$$\Rightarrow \mathbb{E}(u_i | x_i) \neq 0$$

③ Threats to external validity

↳ whether one can generalise 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?

Notes

1. comment — statistics , practical significance , bias $\left\{ \begin{array}{l} \text{OVB} \\ \text{simultaneous causality} \\ \text{errors in measurement} \end{array} \right.$

Regression with panel data

① panel data and formulation

(panel data) data for n different entities observed at T time periods.
 notation as $(\{x_{it}\}_{j \in L, i, v}, y_{it})$
 $\underbrace{\quad}_{v \text{ variables of entity } i \text{ at time } t}$

b) a balanced panel has all observations — variables are observed for each entity and each time period.

Regression and OVB

↪ suppose we want to study the causal relationship between y and $\{x_j\}_{j \in L, v}$, and we have a set of panel data.

↪ simplistically, we can model as $y_{it} = \beta_0 + \beta_1 x_{it} + \dots + \varepsilon$

↪ but it is conceivable to think that there is omitted variable bias for unobserved variables. The dimensions of time and entity mean that we could split these variables.

$$y_{it} = \beta_0 + \beta_1 x_{it} + \dots + A_i z_i + T w_t + \varepsilon$$

variable of interest
 changing in space and time
 omitted variable that
 only changes across time
 omitted variable
 that only changes across
 entities (e.g. gender)

y. GST rate

(entity and time fixed effects)

↪ variables in panel data inherently vary across entities and time. In studying whether some set of X affects y , we can simply abstract out the effects of "irrelevant" effects specific to each entity and time period as fired effects.

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

/ /
 entity specific time specific
 consilive base line

② fixed effects regression

$T=2$, Take difference

$$Y_{i1} = \alpha_i + \gamma_1 + \beta_1 X_{i1} + \varepsilon_{i1}$$

$$Y_{i2} = \alpha_i + \gamma_2 + \beta_1 X_{i2} + \varepsilon_{i2}$$

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

$$\Delta Y_{it} = \beta_0 + \beta_1 \Delta X_{it} + \Delta \varepsilon_{it}$$

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

$$Y_{it} = \alpha_i + \gamma_t + \beta_1 X_{it} + \varepsilon_{it}$$

$$Y_i = \alpha_i - \alpha_1$$

$$\gamma_t = \gamma_t - \gamma_1$$

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.

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

$$Y_{it} = \beta_0 + \beta_1 X_{it} + \dots \quad \text{offset for entity } i$$

$$+ \gamma_2 D_{2i} + \dots \quad \begin{cases} \gamma_2 \\ D_{2i} \end{cases} \quad \text{is } i \text{ the } n^{\text{th}} \text{ entity}$$

$$+ \gamma_3 D_{3i} + \dots \quad \begin{cases} \gamma_3 \\ D_{3i} \end{cases} \quad \text{is } t \text{ the } T^{\text{th}} \text{ period}$$

$$+ \dots \quad \text{offset for time } T$$

(de-meaning)

$$Y_{it} = \alpha_i + \gamma_t + \beta_1 X_{it} + \varepsilon_{it}$$

$$\tilde{Y}_{it} = Y_{it} - \bar{Y}_i - \bar{Y}_t$$

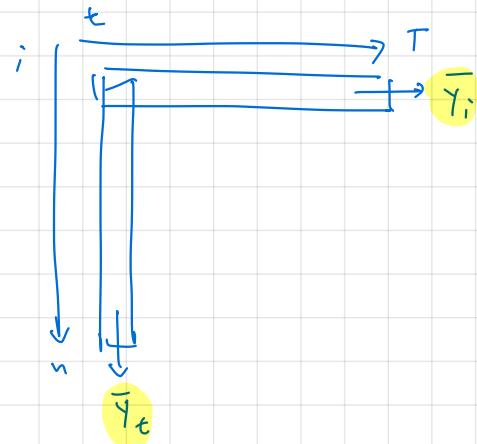
$$= \cancel{\alpha_i} + \cancel{\gamma_t} + \beta_1 X_{it} + \varepsilon_{it}$$

- $(\cancel{\alpha_i} + \beta_1 \bar{X}_i + \bar{\varepsilon}_i)$ ↗ time effect lost, assuming avg. = 0
- $(\cancel{\gamma_t} + \beta_1 \bar{X}_t + \bar{\varepsilon}_t)$ ↗ entity effect lost, assuming avg. = 0

$$= \beta_1 (\bar{X}_{it} - \bar{X}_i - \bar{X}_t) + (\bar{\varepsilon}_{it} - \bar{\varepsilon}_i - \bar{\varepsilon}_t)$$

$$\underbrace{\bar{X}_{it}}_{\tilde{X}_{it}}$$

$$\underbrace{\bar{\varepsilon}_{it}}_{\tilde{\varepsilon}_{it}}$$



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

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

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

estimated directly by mean if
no other regressors

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

③ estimation and estimator distribution

OLS estimation and estimator distributions

↪ we formulate the model as one of the three forms in fixed effects regression and use OLS estimation to estimate $\hat{\beta}_j$.

↪ Here, we describe the fixed effects estimator $\hat{\beta}_j$ is the OLS estimator obtained using demeaning.

↪ recall that for cross sectional data:

$$Q = \sum_{i=1}^n (y_i - \bar{y}_i)^2 = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_j x_j + \dots))^2$$

$$\frac{\partial Q}{\partial \hat{\beta}_0} = 0, \quad \frac{\partial Q}{\partial \hat{\beta}_j} = 0$$

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

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

replace x_{ij} by $\tilde{x}_{ijt} = x_{ijt} - \bar{x}_{ij} - \bar{x}_{jt}$
 y_i by $\tilde{y}_{it} = y_{it} - \bar{y}_i - \bar{y}_t$

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

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

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

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

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

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

Serial correlation

The most common form of serial correlation is first-order serial correlation, in which the error in time t is correlated to the previous ($t-1$) period error.

$$\varepsilon_t = \rho \varepsilon_{t-1} + u_t, \quad |\rho| < 1$$

$\rho > 1$: positive autocorrelation \rightarrow error terms tend to have same sign in related periods

$\rho < 1$: negative autocorrelation \rightarrow error terms tend to have different sign in related periods

Heteroskedasticity robust vs. clustered SE

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

$$\widehat{SE}(\hat{\beta}_1) = \sqrt{\text{Var}(\hat{\beta}_1)} \quad \text{Var}(\hat{\beta}_1) = \frac{1}{n} \sum \frac{\text{Var}((X_i - \bar{X}) u_i)}{\text{Var}(X_i)^2}$$

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

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

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

Clustered SE

$$\text{SE}_{\text{clustered}}(\hat{\beta}) = \sqrt{\frac{1}{nT} \frac{\sum u_i^2}{Q_X^2}}$$

$$\hat{s}_{\text{h.c.}}^2 = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{1}{T} \sum_{t=1}^T \tilde{x}_{it} \tilde{u}_{it} - \left(\frac{1}{n} \sum_{i=1}^n \tilde{x}_{i1} \tilde{u}_{i1} \right) \right)^2$$

works out to 0 assuming u_i and x_i uncorr.

clustering by state

\Rightarrow i.e. for one i , avg. across all its time, because unit-level dependent over time, but independent across states

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
↳ else we distort SE
↳ if holds, can use heteroskedasticity robust SE
↳ also len if just 2 time periods

Regression with limited dependent variable

① binary dependent variables

i) limited dependent variables dependent variables with limited range - e.g. binary

ii) regression with binary variables

↳ what does it mean to fit a line to a dependent variable that can only take on only two values, 0 and 1? \Rightarrow conditional probability!

↳ 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

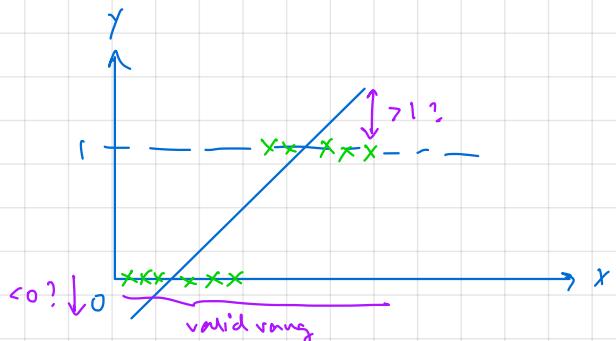
estimator distributions same as that of OLS

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.



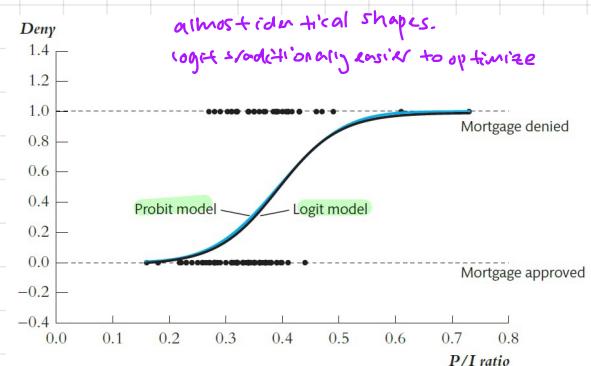
③ probit & logit regression

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

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

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 β_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

(R²) R² 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; R², 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-R²] 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

Notes

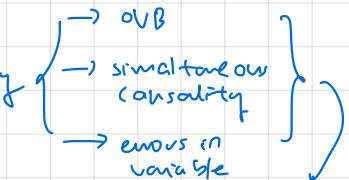
1. linear vs. probit — talk about saturation / bias, nonsensical predictions, qualitative difference in predictions across a range

Instrumental variables regression

① Instrumental variable regression

1) linear regression, endo- and exogeneity

↳ in linear regression, there are three big threats to internal validity



(endogenous variables) $\text{corr.}(X_{ij}, u_i) \neq 0 \Rightarrow E(u_i | X_{ij}) \neq 0$

(exogenous variables) $\text{corr.}(X_{ij}, u_i) = 0 \quad E(u_i | X_{ij}) \neq 0 \leftarrow X_j \text{ is correlated with } u$

↳ multiple regression & panel data regression help \bar{u} OVB \Rightarrow but how do we fix other biases?

2) idea behind IV regression

↳ corr. between X_{ij} and u_i is problematic because the estimator of $\hat{\beta}_j$ becomes biased, and intuitively so because minimising L^2 loss wouldn't be "representative".

↳ but, if we could replace X_{ij} with the "component" of X_{ij} not corr. to u_i , that would work! After all, it is still X_{ij} changing, just that the reason it is changing is entirely unrelated to $u_i \Rightarrow$ so $E(u_i | X_{ij}) = 0$

\Rightarrow how do we isolate that "component" of X_{ij} ? \Rightarrow use other variables unrelated to u_i to predict X_{ij} , then regress y on \hat{X}_j !

those other variables are precisely IVs

3) instrumental variables

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

② Two stage least-squares estimator

1) model formulation

↳ 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 .

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

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

z_i predicts both x_i and y_i

w/o OVB.

$$\mathbb{E}(v_i | z_i) = 0$$

$$\mathbb{E}(w_i | z_i) = 0$$

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

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

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

$$\begin{aligned} \hat{\gamma}_1 &= \frac{\text{cov}(y_i, z_i)}{\text{var}(z_i)} \\ \hat{\pi}_1 &= \frac{\text{cov}(x_i, z_i)}{\text{var}(z_i)} \end{aligned}$$

$$\hat{\beta}_1 = \frac{\hat{\gamma}_1}{\hat{\pi}_1} = \frac{\text{cov}(y_i, z_i)}{\text{cov}(x_i, z_i)}$$

2) estimator distribution

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

$$\hat{\beta}_1^{\text{TSLS}} \sim N(\beta_1, \frac{\text{Var}((z_i - \mu_z) u_i)}{n \cdot \text{cov}(x_i, z_i)^2})$$

$\hat{\beta}_j^{\text{TSLS}}$ is unbiased if z exogenous

$$\begin{aligned} s_{zy} &= \frac{1}{n-1} \sum (z_i - \bar{z})(y_i - \bar{y}) \\ &= \frac{1}{n-1} \sum (z_i - \bar{z}) (\beta_1 (x_i - \bar{x}) + u_i) \\ &= \beta_1 s_{zx} + \sum (z_i - \bar{z})(u_i - \bar{u}) \end{aligned}$$

$$\hat{\beta}_1^{\text{TSLS}} = \frac{s_{zy}}{s_{zx}} = \beta_1 + \frac{\frac{1}{n-1} \sum (z_i - \bar{z}) u_i}{\frac{1}{n-1} \sum (z_i - \bar{z})(x_i - \bar{x})} \rightarrow s_{zu} \Rightarrow \mathbb{E}(s_{zu}) = \text{cov}(z, u) = 0 \text{ if } z \text{ is exogenous}$$

$\hat{\beta}_j^{\text{TSLS}}$ approximately normal \bar{w} known variance in large samples

1. $\hat{\beta}_1^{\text{TSLS}} = \beta_1 + \frac{\frac{1}{n-1} \sum (z_i - \bar{z}) u_i}{\frac{1}{n-1} \sum (z_i - \bar{z})(x_i - \bar{x})} \stackrel{s_{zu}}{\sim} \frac{\frac{1}{n} \sum (z_i - \bar{z}) u_i - \text{cov}(z, u)}{\frac{1}{n} \sum (z_i - \bar{z})(x_i - \bar{x}) - \text{cov}(z, x)} \rightarrow \text{exogeneity assumption}$
2. since $\frac{1}{n} \sum (z_i - \bar{z}) u_i$ is the mean of $(z - \bar{z}) u_i$, and because $\mathbb{E}(u_i | z_i) = 0$ we know $\mathbb{E}(\frac{1}{n} \sum (z_i - \bar{z}) u_i) = 0$, $\text{Var}(\frac{1}{n} \sum (z_i - \bar{z}) u_i) = \frac{\text{cov}((z_i - \bar{z}) u_i)}{n}$
3. Also, rearranging, we have $(\hat{\beta}_1 - \beta_1) \text{cov}(z, x) = \frac{1}{n} \sum (z_i - \bar{z}) u_i$ as $n \rightarrow \infty$.
4. By CLT, as $n \rightarrow \infty$, $\frac{1}{n} \sum (z_i - \bar{z}) u_i \sim N(0, \frac{\text{Var}((z_i - \bar{z}) u_i)}{n})$ approximately.
5. Since $\text{cov}(z, x)$ is a constant, as $n \rightarrow \infty$, $\hat{\beta}_1 \sim N(\beta_1, \frac{\text{Var}((z_i - \bar{z}) u_i)}{n \cdot \text{cov}(z, x)^2})$ scaling the variance

3) estimation and inference in TSLS

(one shot) "two stage"

$$\hat{\beta}_1^{\text{TSLS}} = \frac{S_{yz}}{S_{xz}}$$

$\hat{SE}(\hat{\beta}_1^{\text{TSLS}})$, robust is ok
to use directly

(preferred)

\hat{s}_e must be heteroskedasticity robust, still!

$$\hat{x}_i = \hat{\pi}_0 + \hat{\pi}_1 z_i$$

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

unbiased but standard errors will be wrong since they do not account for 2 stage

4) assumptions of IV regression

$$Y_i = \beta_0 + \sum_{i=1}^k \beta_i X_{ki} + \sum_{j=k+1}^r \beta_j W_{ji} + u_i$$

$$1. E(u_i | W_{i1}, \dots, W_{ir}) = 0$$

$$2. (Y_i, X_{i1}, \dots, W_{ir}) \text{ is i.i.d.}$$

3. no extreme outliers

4. IV Z_1, \dots, Z_k are valid

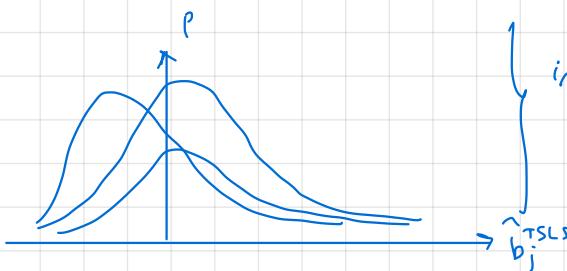
③ checking instrument strength

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 $n \rightarrow \infty$



intuition: $\hat{\beta}_j^{\text{TSLS}} = \frac{\text{cov}(z, y)}{\text{cov}(z, x)}$ ← small
very sensitive

2) checking instrument relevance

idea regress endogenous X on Z_1, \dots, Z_m , controlled covariates W_1, \dots, W_n .
 then do an F test that coefficients on Z_1, \dots, 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

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

\hookrightarrow 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%, which is generally acceptable

\Rightarrow if $F > 10$, relative bias is $< 10\%$, vice versa

3) dealing with weak instruments

\hookrightarrow if multiple instruments, some probably weaker than others — drop weaker, retest until satisfactory subset

④ checking instrument validity

under, exactly, over identified no. of IVs less than, equal to, more than endogenous variables X

Checking instrument exogeneity: J test of overidentified instruments

\hookrightarrow idea: exogeneity of IVs mean that they are uncorrelated with u_i . It logically follows that they should be statistically uncorrelated in our estimate of

$$u_i, \hat{u}_i = y_i - (\hat{\beta}_0^{\text{TSLS}} + \hat{\beta}_1^{\text{TSLS}} x_1 + \dots + \hat{\beta}_{K+1}^{\text{TSLS}} w_{K+1} + \dots)$$

↑ true x , rather than predicted \hat{x} . Take note! Not \hat{x} .

$$1. \text{ estimate } \hat{y}_i = \beta_0 + \beta_1 x_1 + \dots + \beta_{K+1} w_{K+1} + \dots \Rightarrow \text{get } \hat{u}_i = y_i - \hat{y}_i$$

$$2. \text{ regress } \hat{u}_i \text{ on all supposedly exogenous variables} \Rightarrow \hat{u}_i = f_1 z_1 + \dots + f_m z_m + \dots$$

3. compute F stat. that f_1, \dots, f_m , all coefficients on IV, are 0.

4. compute J stat., that is

$$J = mf$$

\downarrow
no. of IV

no. of IV - no. of X endogenous

5. compute p-value. $J \sim \chi^2(m-k)$ → reject \Rightarrow not exogenous

how the J test works

we can derive a relationship using the above assumptions and equations that $J \sim \chi^2(m-k \text{ dof}) \Rightarrow$ reject using p-value.

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

⑤ choosing instrumental variables

1. economic theory

2. random phenomena → cause variation in one, but not other (e.g. birthdate)

Notes

$$1. \text{ recall } \hat{\beta}_1 = \frac{\gamma_1}{\pi_1} \Rightarrow \frac{\text{cov}(\gamma_i, z_i)}{\text{cov}(x_i, z_i)} = \frac{\text{cov}(z_i, \beta_0 + \beta_1 x_i + \beta_2 w_i + u_i)}{\text{cov}(x_i, z_i)}$$

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

⇒ 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 regression (forecasting)

① fundamental concepts

1) fundamental operators

(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!

2) autocovariance and autocorrelation

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

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

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

3) types of autocorrelation

(geometric decay)

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

(slow decay) power law

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

② stationarity of processes

(stochastic processes)

\hookrightarrow traditionally defined as a collection of random variables indexed by some $t \in T$

$$\{x_t\}_{t \in T} \Rightarrow \text{re. an indexed set of RVs}$$

\hookrightarrow time series processes can be thought of as stochastic processes

(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 $Y_{t-n}, Y_{t-n-1} \dots Y_{t-n-k}$ for any n . All corr. joint distributions are identical & time displacement invariant

The joint distributions $(Y_{t+1}, Y_{t+2} \dots Y_{t+T})$ do not depend on t , 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

\Rightarrow for I(3), we set:

$$\text{cov}(Y_t, Y_{t+j}) = \sigma^2 \delta_j \Leftrightarrow \text{Var}(Y_t) = \sigma^2 \quad (j=0)$$

$$E(Y_t) = E(Y_{t+j}) = \mu \delta_j$$

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

③ ergodicity

(definition)

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

↳ suppose autocovariance stationary process $\{Y_{t+h}\}_{h \in [0, T]}$ has mean M .

$$\text{then } E(Y_t) = E(Y_{t+1}) = E(Y_{t+h}) = M$$

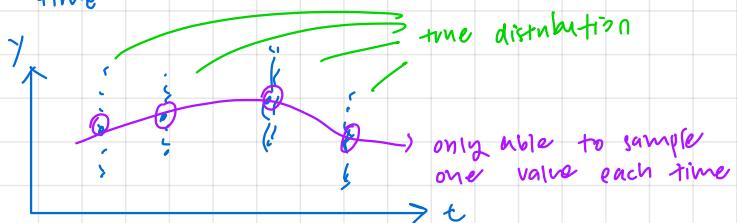
↳ then if $\frac{1}{T} \int_0^T \{Y_{t+h}\}_{h \in [0, T]} dt \rightarrow M$ as $T \rightarrow \infty$, we say it is ergodic

for the 1st moment (mean)
 ↗ observe that this is only possible
 if faraway variables are essentially independent

\Rightarrow If the statistical property is ergodic, then the long-horizon forecast converges to the unconditional value of property

④ why are stationarity & ergodicity important?

↳ 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



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

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

\Rightarrow can be simple like mean \rightarrow then $\hat{m} = \frac{1}{N} \sum_{i=1}^N Y_{t,n} = \frac{1}{T} \sum_{i=1}^T Y_t$

\Rightarrow can also be a more complicated recurrence relationship like

$Y_t = \beta_0 + \beta_1 Y_{t-1}$, and β_0 and β_1 can be estimated in enough time steps of data

⑤ Autoregression on stationary & ergodic time series

(autoregression)

An autoregression is a regression model in which Y_t is regressed against its own lagged values

↳ the number of lagged values is called the order of the regression

autoregressive distributed lag model

- ↪ suppose there is some stable (stationary) relationship between y_t and its previous p values, and also other 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 via 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 get large. That is, all regressors are ergodic.
4. No perfect multicollinearity or outliers.

⑥ forecasting errors

i) forecast error, MSFE & RMSFE

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

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

$$RMSFE = \sqrt{MSFE}$$

ii) 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

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

variance of random error

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

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

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

$$\hat{MSFE} \approx E(u_T^2) = \sigma_u^2 \approx SER$$

2. estimation by final prediction error

\hookrightarrow assuming errors are homoskedastic

$$\hat{MSFE} = \frac{T+p+1}{T} s_u^2 = \frac{T+p+1}{T-p-1} \frac{SSR}{T}$$

3. estimate by pseudo OOS

- ↳ the SER method ignores estimation error
- ↳ the FPE method requires homoskedasticity
- ↳ & formulation also requires stationarity, since functional form \hat{Y} above based on stationary
- ↳ 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 $\tilde{u}_{s+1}, s = T - P, \dots, T - 1$.

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

P samples for OOS testing, take avg. estimated error

⑦ forecast uncertainty

(forecast intervals) forecast intervals are like a confidence interval except that it pertains to a forecast.

- ↳ A key difference, though, is that CI are justified by CLT, so we can use the appropriate CI.
- ↳ in time series forecasting, we need to estimate the distribution of u_t or make assumptions about u_t in order to compute forecast intervals.
- \Rightarrow 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.

$$E(Y_t - \hat{Y}_t) = E(u_t) + E((\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) Y_t + \dots)$$

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

\Rightarrow in which case we formulate

$$\hat{y}_t \pm \hat{\sigma}_{\epsilon} \hat{RMSE} \leftarrow \begin{array}{l} \text{normally distributed} \\ \text{RMSE PFE or RMSE oos} \end{array}$$

(fan chart)

Forecast distributions are frequently portrayed graphically as a fan chart



⑧ Determining the order of autoregression

- ↳ how do we decide on autoregressors? unlike multiple regression on panel data, where the objective was unbiased estimates, here we simply care about MSFE
- ↳ how do we decide/ trade-off between exploiting autocorrelation vs. estimating too many and introducing more error into our MSFE?

F-stat similar to multivariate regression

(Bayesian information criterion)

$$BIC(p) = \ln \frac{SSR(p)}{T} + (p+1) \frac{\ln T}{T} \quad \begin{array}{l} \text{intuition:} \\ \text{SSR} \downarrow, BIC \downarrow \\ p \uparrow, BIC \uparrow \text{ (penalty)} \end{array}$$

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

→ fit p th model and compute $BIC(p)$. Pick model \bar{p} 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

$$AIC(p) = \ln \frac{SSR(p)}{T} + (p+1) \frac{2}{T} \quad \begin{array}{l} \text{SSR} \downarrow, AIC \downarrow \\ p \uparrow, AIC \uparrow \\ T \uparrow, AIC \downarrow \end{array}$$

- ↳ penalty for more p smaller than BIC

\Rightarrow AIC estimates kL divergence between model distribution of \hat{Y}_t and data distribution Y_t \Rightarrow how much info lost using the p^{th} model

\Rightarrow not consistent like BIC, but efficient — if true model not considered, picks best approximation w probability 1

(Data for comparing different models)

- make sure same amt. of data for each
- i.e. if necessarily decreases no. of samples

⑨ Non-stationarity (trends in time)

) trends

(definition) a persistent, long-term movement of a variable over time

(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

$$\text{e.g. } Y_t = \beta_0 + \alpha_1 t + \alpha_2 t^2 + \beta_1 Y_{t-1} \dots$$

(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

random walk with drift

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

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

$$1. E(Y_t) = Y_{t-1} + E(u_t), \text{ since } Y_{t-1} \text{ is realized}$$

$$2. Y_t = Y_{t-2} + u_{t-1} + u_t \dots$$

$$= Y_{t-k} + \sum_{i=1}^k u_{t-i} = Y_0 + \sum_{i=1}^t u_{t-i}.$$

$$\text{Var}(Y_t) = 0 + t \text{Var}(u_t) \text{ if homoskedastic}$$

$$3. \text{cov}(Y_t, Y_{t-k}) = \text{cov}\left(Y_0 + \sum_{i=1}^{t-k} u_i, Y_0 + \sum_{i=1}^{t-k} u_i\right)$$

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

$$1. E(Y_t) = \beta_0 + Y_{t-1} + E(u_t)$$

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

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

Since u_i is i.i.d, u_i in different periods are independent.

so $\text{cov}(u_i, u_j)$, $i \neq j = 0$

Then $\text{cov}\left(\sum_{i=1}^t u_i, \sum_{i=1}^k u_{t-i}\right)$ resolves to $\sum_{i=1}^{t-k} \text{var}(u_i) = (t-k) \text{var}(u_i)$

(

Since $\text{Var}(Y_t)$ and $\text{Var}(Y_t, Y_{t+k})$ depend on t , stochastic trend

2) testing for non-stationarity of Y_t

stationarity of Y_t

1. if Y_t follows an AR(p) model, Y_t is stationary if its roots are all greater than 1 in absolute value.

$$1.1 \quad Y_t = \beta_0 + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} \dots + u_t, \quad u_t \text{ is i.i.d.}$$

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

$$1.2 \quad Y_t = \text{some function of } u_t.$$

2. recall that a stochastic process is autocovariance stationary if $E(Y_T)$, $\text{cov}(Y_T, Y_{T-k})$ and $\text{Var}(Y_T) = \text{cov}(Y_T, Y_{T-0})$ is constant.

$$2.1 \quad \text{cov}(Y_T, Y_{T-k}) = \text{some function of } u_T.$$

$$2.2 \quad \text{Var}(Y_T) = \text{some function of } u_T$$

$$2.3 \quad E(Y_T) = \text{some function of } u_T.$$

2.4 For 2.1, 2.2, 2.3 to be constant, roots $z_1 \dots z_p$ must be greater than 1.

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

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

roots are then z that satisfy $1 - \beta_1 z - \beta_2 z^2 \dots - \beta_p z^p = 0$
p roots. If all $|z| > 1$, 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 $\rightarrow \sum_{i=1}^p \beta_i = 1 \Rightarrow H_0: \text{not stationary, try to reject}$

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

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

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

Extending the Dickey Fuller test to AR(p)

1. For $p=2$, observe that:

$$\begin{aligned} Y_t &= \beta_0 + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} \\ &= \beta_0 + (\beta_1 + \beta_2) Y_{t-1} - \beta_2 Y_{t-1} + \beta_2 Y_{t-2} \\ &= \beta_0 + (\beta_1 + \beta_2) Y_{t-1} - \beta_2 \Delta Y_{t-1} \end{aligned}$$

2. For $p=3$, observe that:

$$\begin{aligned} Y_t &= \beta_0 + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \beta_3 Y_{t-3} \\ &= \beta_0 + (\beta_1 + \beta_2 + \beta_3) Y_{t-1} - (\beta_2 + \beta_3) (Y_{t-1} - Y_{t-2}) - \beta_3 (Y_{t-2} - Y_{t-3}) \\ &= \beta_0 + (\beta_1 + \beta_2 + \beta_3) Y_{t-1} - (\beta_2 + \beta_3) \Delta Y_{t-1} - \beta_3 \Delta Y_{t-2} \end{aligned}$$

3. By induction, we can show that if we regress

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

3.1 This is the augmented Dickey Fuller test.

$$\begin{aligned} H_0: f &= \sum_{i=1}^p \beta_i - 1 = 0 \rightarrow \text{non-stationary} \\ H_1: f &< 0 \rightarrow \text{not non-stationary} \end{aligned} \quad \left. \begin{array}{l} \text{test against ADF} \\ \text{statistic distribution} \\ \text{by } \frac{f}{SE(f)} \text{ vs. } ADF_2 \end{array} \right.$$

4. We choose p by the same ways.

3) testing for stochastic & deterministic non-stationarity

↳ we extend the augmented DF test by including function of t

$$\Delta Y_t = \beta_0 + \text{include time terms} + \sum_{i=1}^{p-1} Y_i \Delta Y_{t-i}, f = \sum_{i=1}^{p-1} \beta_i - 1$$

$$H_0: f = \sum_{i=1}^p \beta_i - 1 = 0 \rightarrow \text{non-stationary}$$

$$H_1: f < 0 \rightarrow \text{not non-stationary}$$

} test against ADF statistic distribution by $\frac{f}{SE(f)}$ vs. ADF_{α}

(10) problems caused by stochastic trends

i) bias & non-normality of t-statistic

↳ t-stat of OLS estimates of β_i in AR(p) models will be non-normal and biased toward 0, even in large samples \Rightarrow invalid hypothesis testing

ii) spurious regressions

↳ stochastic trends can lead to two time series e.g. Y_t, X_t to appear related when they are not, if both are stochastic

(11) Dealing with non-stationarity

(1) Detrending the data

↳ observe that in random walk model, stochasticity arises from influence of previous step \Rightarrow de-trend by differencing

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

$$\Delta Y_t = \beta_0 + u_t \Rightarrow \text{stationary about } \beta_0$$

