## **Revision Notes for ETC2410 – Sample**

## 3. Linear Regression Model

## Lecture #3 – The Linear Regression Model

- <u>Simple LRM</u>.  $\mathbf{E}(\mathbf{y}_i|\mathbf{x}_i) = \beta_0 + \beta_1\mathbf{x}_i + \mathbf{u}_i$ . (#1) **CONDITIONAL DISTRIBUTION** model Y with a mean and variance for each level of the predictor X; (#2) **LINEARITY** that  $\mathbf{E}(\mathbf{y}_i|\mathbf{x}_i) = \beta_0 + \beta_1\mathbf{x}_i$  for each observation i. (#3) **ZERO CONDITIONAL MEAN**  $\mathbf{E}(\mathbf{u}|\mathbf{x}) = 0$ . Deviations occur, but on average = 0
  - MARGINAL EFFECT:  $\beta_1$ ; 1-UNIT INCREASE in X  $\rightarrow$  E(Y) rises by  $\beta_1$ , all other factors

constant. Can take the <u>DERIVATIVE</u> of y with respect to x:  $\frac{\partial y}{\partial x} = \beta_1$ . Also,  $\beta_1 = Cov / Var$ .

- From the linearity assumption, follows that:  $E(y_i|x_i+1) E(y_i|x_i) = \beta_1$ .
- **INTERCEPT:**  $\beta_0$ ; <u>EXPECTED VALUE</u> of the DV (Y) when all predictor variables (X) = **ZERO**.  $\hat{\beta}_0$  = the predicted value of  $y_i$  when  $x_i = 0$ . This is because- **E(y|x=0) = \beta\_0 + \beta\_1(0)**.
- **RESIDUALS**: **u**<sub>i</sub>. captures **UNEXPLAINED VARIATION**; factors affecting the DV <u>NOT</u> <u>CAPTURED</u> by the model's predictors (xi). The residual is:  $\hat{\mathbf{u}} = \hat{\mathbf{Y}} - \mathbf{Y} = \mathbf{y}_i - \mathbf{E}(\mathbf{y}_i | \mathbf{x}_i)$ .
- <u>OLS estimator</u>. **OLS ESTIMATOR**: <u>ESTIMATES REGRESSION COEFFICIENTS</u> ( $\beta_0$  and  $\beta_1$ ) by **MINIMIZING** the <u>SUM OF SQUARED DIFFERENCES</u> b/w **OBSERVED** and **PREDICTED** values.
  - Aim of OLS estimator: use OBSERVED DATA to estimate the <u>UNKNOWNS</u>  $\beta_0$  and  $\beta_1$ . Control the fit of the model. There is sampling variability, as these are random variables.
  - Method of the OLS: (#1) calculate PREDICTION ERRORS-  $\hat{\mathbf{u}}_{i} = \mathbf{y}_{i} \mathbf{b}_{0} \mathbf{b}_{1}\mathbf{x}_{i}$ . (#2) substitute into the SSR Eq. SSR( $\mathbf{b}_{0}, \mathbf{b}_{1}$ ) =  $\sum_{i=1}^{n} \hat{\mathbf{u}}_{i} = \sum_{i=1}^{n} (\mathbf{y}_{i} - \mathbf{b}_{0} - \mathbf{b}_{1}\mathbf{x}_{1})^{2}$ . (#3) To MINIMISE SSR, set the derivative = 0 with respect to each parameter ( $\beta_{0}$  and  $\beta_{1}$ ). This gives TWO equations to solve for TWO unknowns.
    - gives TWO equations to solve for TWO unknowns. • Formula for b1:  $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\widehat{cov(x,y)}}{Var(x)} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2}$ . Formula for b0:  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .

<u>Properties of data that hold during OLS</u>. ( $y_i = \beta_0 + \beta_1 x_i + u_i$ , i = 1, ..., n by the method of OLS):

- (P1):  $\sum_{i=1}^{n} \hat{u}_{i} = 0$ . Residuals sum to zero. On avg, model is right-balance in middle of the data.
- (P2): ∑<sub>i=1</sub><sup>n</sup> x<sub>i</sub>û<sub>i</sub> = 0. Errors are <u>NOT</u> SYSTEMATICALLY RELATED to predictor. Errors are EVENLY SPREAD across the x-axis, model captured all linear information in the predictor.
   For P2: the vector UHAT is ORTHOGONAL to columns of X- x'u = 0.
- (P3):  $\overline{y} = \overline{y^{\wedge}}$ , where  $\overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$  and  $\overline{y^{\wedge}} = \frac{1}{n} \sum_{i=1}^{n} \hat{y}_i$ . Avg of actual values = avg of the predicted.
- <u>R<sup>2</sup></u>. **R-SQUARED** (COEFFICIENT OF DETERMINATION): measure of GOODNESS OF FIT.
  <u>PROPORTION</u> of SAMPLE VARIATION in y explained by x. R<sup>2</sup> = SSE / SST = 1 SSR/SST.
  - SST: SST → total sample variation in y (total variation of observed data around its mean).
    SST does not consider a model. SST = SSE + SSR. DECOMPOSED into SSR and SSE:
    - SSE: part of SST explained by IVs in model.  $\hat{\beta}_0 + \hat{\beta}_1 x \rightarrow SSE$ , explaining Var(YHAT). Measures variation in <u>PREDICTED VALUES</u> around the mean. :  $\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$ .
    - SSR: the part of SST NOT explained by the model. UHAT  $\rightarrow$  SSR, accounting for variation not explained by predictors.  $\hat{u}: \sum_{i=1}^{n} \hat{u}_{i}^{2}$ .
- <u>Multiple LRM</u>. The conditional mean of y depends on "k" explanatory variables. For the general MLR, it follows that: E(y<sub>i</sub> | x<sub>i1</sub> + 1, x<sub>i2</sub>, ..., x<sub>ik</sub>) E(y<sub>i</sub> | x<sub>i1</sub>, x<sub>i2</sub>, ..., x<sub>ik</sub>) = β<sub>1</sub>.
  - Coefficient interpretation: b1 measures average  $\Delta$  in Y in response to 1-unit  $\Delta$  in X, HOLDING VALUES OF ALL OTHER PREDICTORS CONSTANT. EXPLICITLY account.

- <u>MLR in matrix notation</u>. (#1) n observations = n equations with the unknown coefficients. (#2) Formulate in matrix algebra – stack equations for each n. (#3) write the model as  $y = X\beta + u$ .
  - $\circ \quad \text{Vectors-(\#1) y \& u = n x 1 (stack observations from 1 to n). (\#2) coefficient V \beta (k+1) x 1.}$
  - X is a regressor matrix, with (#1) a **column of 1's** for b0, (#2) **k+1** columns and **n** rows.
- OLS in Matrix Form.  $\hat{\beta} = (X'X)^{-1}X'y$  is the OLS estimator of  $\beta$  in matrix notation.
  - Define a <u>VECTOR</u> of **OLS estimates** for  $\beta$ . (#1) vector of OLS **FITTED VALUES** is  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ . (#2) vector of OLS **RESIDUALS** is  $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ . (#3) So, the **SSR** equals  $\sum_{i=1}^{n} \hat{\mathbf{u}}_{i}^{2} = \hat{\mathbf{u}}'\hat{\mathbf{u}} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ .
  - X'X <u>MUST BE</u> an **INVERTIBLE MATRIX**  $\rightarrow$  columns of X must be **LINEARLY INDEPENDENT**.
    - **LINEAR INDEPENDENCE**: no columns in the matrix can be built by adding or scaling others (linear combination). Can't recreate one column by mixing others.
      - Interpreting linear dependence: a regressor is not adding any <u>NEW</u> info.