# Simple Linear Regression Model

The model is

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \qquad i = 1\ 2\ 3\ 4 \dots n$$

Simple since there is only one $X_i$

It is linear in $\beta_0$ and $\beta_1$

$E(Y_i) = \beta_0 + \beta X_i$ is linear.

$Y = \beta_0 + \beta_1 X_1 + e^{\beta_2 X_2} + \varepsilon_i$ is not linear

$Y = \beta_0 + e^{\beta_1 X_1} + \varepsilon_i$ is not linear also.

- A simple linear regression model assumes that the following specification is true in the population

$$Y = \beta_0 + \beta_1 X_i + \varepsilon$$

- Where other unobserved factors determining $y$ are captured by the error term $\varepsilon$.

- To estimate the parameters $\beta_0$ & $\beta_1$ in the model Assumptions 1 - 5, Linearity, Identification, exogeneity, spherical error terms and data generation are required for Ordinary least squares. (OLS)

- Assumptions 1 - 6, ie, including normality are required to use Maximum likelihood to estimate $\beta_0$ & $\beta_1$ (MLE)

Assumptions summarized

i   $E(\varepsilon_i) = 0$ for all $i = 1\ 2 \dots n \Rightarrow E(Y_i) = \beta_0 + \beta_1 X_i$

ii  $Var(\varepsilon_i) = \sigma^2$ for all $i = 1\ 2 \dots n \Rightarrow Var(Y_i) = \sigma^2$ Constant

iii Covariance$(\varepsilon_i, \varepsilon_j) = 0$ for $i \neq j \Rightarrow Cov(Y_i, Y_j) = 0$ uncorrelatedness

iv  Normality $\varepsilon_i \sim N(0, \sigma^2)$, $\varepsilon_i$ and $Y_i$ are independent and uncorrelated.

## Ordinary least squares estimator

- Least squares Method does not require any distributional assumptions. (It does not require normality)

- MLE estimation requires the normality assumption (6)

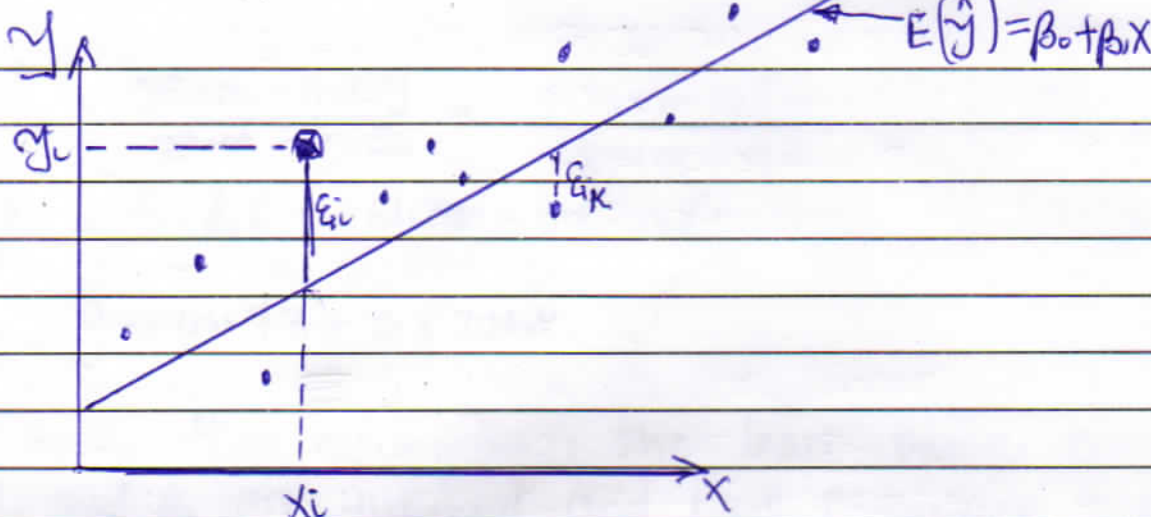- Ordinary least squares estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ minimizes the sum of squared residuals (SSR)

$$Q = \sum_{i=1}^{n} \hat{\varepsilon}_i^2 = \hat{\varepsilon}'\hat{\varepsilon} = \sum_{i=1}^{n}\left(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i\right)^2 = \sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2$$

$$E(\hat{y}) = \beta_0 + \beta_1 x$$

$\hat{y}_i$ estimates $E(y_i) = \beta_0 + \beta_1 x_i$, not $\beta_0 + \beta_1 x + \varepsilon_i$.

The minimum is acquired by

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

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

Simplify

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum x_i = \sum y_i$$

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

Solving for

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

OR

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

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

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

Note

$$S_{xx} = \sum_i (x_i - \bar{x})^2 = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum x)}{n}$$

$$S_{yy} = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

Example: Is a student performance in final exam (y) determined by continuous assessment (x)

| y | 95 | 80 | 0 | 0 | 79 | 77 | 72 | 66 | 98 | 90 | 0 | 95 | 35 | 50 | 72 | 55 | 75 | 66 | $\sum y =$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| x | 96 | 77 | 0 | 0 | 78 | 64 | 89 | 47 | 90 | 93 | 18 | 86 | 0 | 30 | 59 | 77 | 74 | 67 | $\sum x =$ |
| $x^2$ | | | | | | | | | | | | | | | | | | | $\sum x^2 =$ |
| $xy$ | | | | | | | | | | | | | | | | | | | $\sum xy =$ |
| $y^2$ | | | | | | | | | | | | | | | | | | | $\sum y^2 =$ |

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x}\bar{y}}{\sum x_i^2 = n\bar{x}} = \frac{81,195 - 18(58.056)(61.389)}{80199 - 18(58.056)^2} = 0.8726$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 61.389 - 0.8726(58.056) = 10.73$$

$$\hat{y} = 10.73 + 0.8726 X.$$

With the three assumptions the least squares estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased and have minimum variance among all linear unbiased estimators "MVUE"

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

$$E(\hat{\beta}_0) = \beta_0$$

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{\sum(x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}$$

$$Var(\hat{\beta}_0) = \sigma^2\left(\frac{1}{n} + \frac{\bar{x}}{\sum(x_i - \bar{x})^2}\right) = \sigma^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)$$

If $\sigma^2 = E(y_i - E(y_i))^2$ using $\hat{y}$ as estimate for $E(y_i)$ we estimate $\hat{\sigma}^2$ by

$$\delta^2 = \frac{\sum(y_i - \hat{y})^2}{n-2}.$$

Note

$$Var(y_i) = E(y_i - E(y_i))^2 = E(y_i - \beta_0 - \beta_1 x_i)^2 = E(\varepsilon_i^2) = \sigma^2$$

$$S^2 = \frac{\sum(y_i - \hat{y}_i)^2}{n-2} = \frac{\sum(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)}{n-2} = \frac{SSE}{n-2}$$

$$\hat{\varepsilon} = (y_i - \hat{y}_i) \text{ is refered to residuals}$$

$$E(S^2) = E\frac{SSE}{n-2} = \frac{(n-2)\sigma^2}{n-2} = \sigma^2 \quad \text{unbiased}$$

$\delta^2$ is an unbiased estimator of $\sigma^2$

$$SSE = \sum(y_i - \hat{y}_i)^2 = S_{yy} - \frac{S_{xy}^2}{S_{xx}} \qquad \delta^2 = \frac{SSE}{n-2}$$

Write on both sides of the paper

REG No.

Question.......... ④

Do not write
in either
margin

# Hypothesis Testing and Confidence Intervals for $\hat{\beta}_1$

$H_0 : \beta_1 = 0$ Vs $H_1 \; \beta_1 \neq 0$ :

$H_0 : \beta_1 = 0$ There is no linear relationship between $y$ and $x$

Assuming $\varepsilon \sim N(0, \sigma^2)$ or $y_i \sim N(\beta_0 + \beta_1 x_1, \sigma^2)$

(i) $\hat{\beta}_1$ is $N\left[\beta_1, \dfrac{\sigma^2}{S_{xx}}\right]$

(ii) $\dfrac{(n-2)s^2}{\sigma^2}$ is $\chi^2(n-2)$

(iii) $\hat{\beta}_1$ and $s^2$ are independent.

Test statistics for $H_0 : \beta_1 = 0$ Vs $\beta_1 \neq 0$

$$t = \frac{\hat{\beta}_1}{S/\sqrt{S_{xx}}} \sim t(n-2) \quad \text{If } \beta = 0$$

Note

$$t = \frac{\text{Estimate}}{\text{Standard error of estimate}} = \frac{\hat{\beta}_1}{Sd(\hat{\beta}_1)}$$

$$\text{Reject } H_0 \text{ if } |t| \geq t_{\alpha/2}(n-2)$$

Detailed Hypothesis

Case 1 : $H_0 : \beta_1 \leq 0$ Vs $H_1 : \beta_1 > 0$

Case 2 : $H_0 : \beta_1 \geq 0$ Vs $H_1 : \beta_1 < 0$

Case 3 : $H_0 : \beta_1 = 0$ Vs $H_1 : \beta_1 \neq 0$

Test statistic (T.S) $\quad t = \dfrac{\hat{\beta}_1 - 0}{S_e/\sqrt{S_{xx}}}$
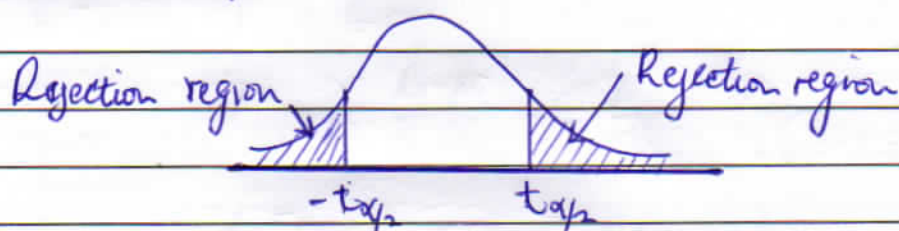
Rejection Region For d.f. $n-2$ and type 1 error $\alpha$

Case 1. Reject $H_0$ if $t > t_\alpha$

2. Reject $H_0$ if $t < -t_\alpha$

3. Reject $H_0$ if $|t| > t_{\alpha/2}$

For case 3: $H_0 : \beta_1 = 0$ Vs $H_1 : \beta_1 \neq 0$



$$t = \frac{\hat{\beta_1}}{S/\sqrt{S_{xx}}} = \frac{0.8726}{13.8547/139.753} = 8.8025$$

$t = 8.8025 > t_{0.025}(16) = 2.120$  Reject $H_0$.

## Confidence Interval for Slope $\hat{\beta_1}$

A $100(1-\alpha)\%$ Confidence interval for $\beta_1$ is

$$\hat{\beta_1} \pm t_{\alpha/2}(n-1) \frac{S}{\sqrt{S_{xx}}} \quad \text{Where } S = \frac{S_{yy} - \frac{S_{xy}^2}{S_{xx}}}{n-2}$$

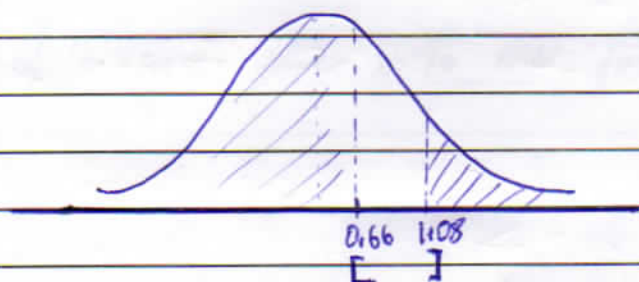$$\beta_1 \pm t_{\alpha/2}(n-1) \frac{S}{\sqrt{S_{xx}}}$$

A 95% Confidence interval for $\beta_1$ for the Marks example

$$\hat{\beta_1} \pm t_{0.025}^{(16)} \frac{S}{\sqrt{(x_i - \bar{x})^2}}$$

$$0.8726 \pm 2.120(0.09914)$$
$$0.8726 \pm 0.2102$$
$$(0.6624, 1.0828)$$



0.66  1.08
[    ]

# Confidence interval for Slope $\beta_1$

$$\left| \frac{\hat{\beta}_1 - \beta_1}{Se\sqrt{1/S_{xx}}} \right| < t_{\alpha/2} \quad \text{is a } (1-\alpha)100\% \text{ Confidence I.}$$

$$-t_{\alpha/2} < \frac{\hat{\beta}_1 - \beta_1}{Se\sqrt{1/S_{xx}}} < t_{\alpha/2}$$

$$-t_{\alpha/2} Se\sqrt{\frac{1}{S_{xx}}} < \hat{\beta}_1 - \beta_1 < t_{\alpha/2} Se\sqrt{1/S_{xx}}$$

$$\hat{\beta}_1 - t_{\alpha/2} Se\sqrt{\frac{1}{S_{xx}}} < \beta < \hat{\beta}_1 + t_{\alpha/2} Se\sqrt{1/S_{xx}} \quad \text{is the}$$

$(1-\alpha)100\%$ Confidence interval for $\beta_1$ where $t$ has $n-2$ df.

# F test for $H_0$ $\beta_1 = 0$ (Recall ANOVA)

$H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

Test Statistics $F = \dfrac{SS(Regression)/n-1}{SS\,Error/n-2} = \dfrac{MSReg}{MSE}$

Rejection Region: With $df_1 = 1$ $df_2 = n-2$, reject $H_0$ if $F > F_\alpha$

$$SSReg = \sum(\hat{y} - \bar{y})^2$$
$$SSE = \sum(y_i - \bar{y})^2$$

# Confidence Interval for intercept $\beta_0$

$$\sigma_{\hat{\beta}_0} = Se\sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

The $(1-\alpha)100\%$ Confidence interval for $\beta_0$ is given by

$$\beta_0 \pm t_{\alpha/2} Se\sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} \quad \text{for } t \text{ with } n-2 \text{ df.}$$

Q Determine the 95% CI for $\hat{\beta}_1$ for example 1.1.

Inference about the Correlation coefficient $\rho_{yx}$

The t statistic for $\beta_1$ can be expressed in terms of r as follows

$$t = \frac{\hat{\beta_1}}{\delta / \sqrt{\sum (x_i - \bar{x})^2}}$$

$$t = r_{yx} \frac{\sqrt{n-2}}{\sqrt{1 - r_{yx}^2}}$$

The sample correlation $r_{yx}$ is the basis for estimation and significance testing of the population correlation $\rho_{yx}$

Hypothesis

Case 1: $H_0 \; \rho_{yx} \leq 0$ vs. $H_1 \; \rho_{xy} > 0$

Case 2: $H_0 \; \rho_{yx} \geq 0$ vs. $H_a \; \rho_{yx} \leq 0$

Case 3: $H_0 \; \rho_{yx} = 0$ vs. $H_1 \; \rho_{yx} \neq 0$

Time Series T.S. $t = r_{yx} \dfrac{\sqrt{n-2}}{\sqrt{1 - r_{yx}^2}}$

R.R with $n-2$ df and type I error probability $\alpha$,

1. $t > t_\alpha$

2. $t < -t_\alpha$

3. $|t| > t_{\alpha/2}$

Check assumptions and draw conclusions

Exercise

| Soil PH : | 3.3 | 3.4 | 3.4 | 3.5 | 3.6 | 3.6 | 3.7 | 3.7 | 3.8 | 3.8 |
|---|---|---|---|---|---|---|---|---|---|---|
| Growth Red : | 17.78 | 21.59 | 23.84 | 15.13 | 23.45 | 20.87 | 17.78 | 20.09 | 17.78 | 12.46 |

$Q_1$ Examine the scatter plot and decide whether a straight line is reasonable model. Is regression significant?

$Q_{11}$ Identify the Least squares model estimates for $Y = \beta_0 + \beta_1 X + \epsilon$

$Q_{111}$ Predict the growth retardation for a soil PH of 4.0

# Simple Linear Regression.

**The Model:** The simple linear regression mode for $n$ observation can be written as

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \qquad i = 1,2,3,\dots n$$

Linear in $\beta_0$ and $\beta_1$

$$y = \beta_0 + \beta_1 x_i^2 + \varepsilon_i \quad \text{is Linear in } \beta_0 \text{ \& } \beta_1$$

$$y = \beta_0 + e^{\beta_1 x_i} + \varepsilon_i \quad \text{is not Linear}$$

Model
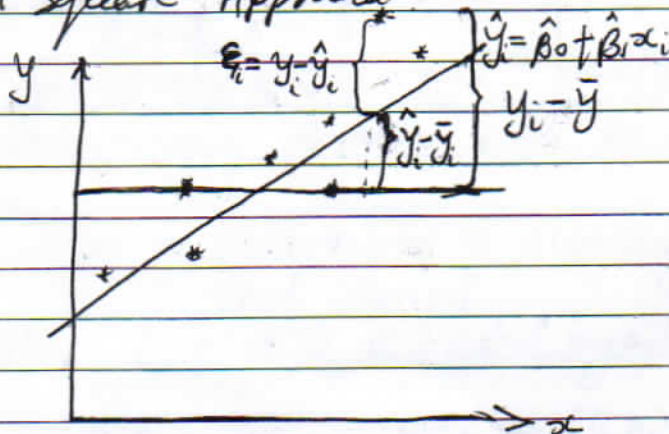$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \qquad i = 1,2\dots n$$

Assumption
(i)  $E(\varepsilon_i) = 0$ for all $i = 1,2\dots n$  or  $E(y_i) = \beta_0 + \beta_1 x_i$
(ii)  $Var(\varepsilon_i) = \sigma^2$ for all $i = 1,2\dots n$  or  $Var(y_i) = \sigma^2$ Homoscedestic
(iii)  $Cov(\varepsilon_i,\varepsilon_j) = 0$ for all $i \neq j$  or  $Cov(x_i, y_i) = 0$

$$Var(y_i) = E[y_i - Ey_i]^2 = E(y_i - \beta_0 - \beta_1 x_i)^2 = E(\varepsilon_i^2) = \sigma^2$$

## Estimation of $\beta_0$ $\beta_1$ and $\sigma^2$

### Least square Approach



$$\varepsilon_i = y_i - \hat{y}_i \qquad \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Total variation $= y_i - \bar{y}$
error (Residual) $= y_i - \hat{y}$
explained by regression $= \hat{y}_i - \bar{y}_i$

$$y_i - \hat{y}_i = (y_i - \bar{y}) - (\hat{y}_i - \bar{y}) \qquad \text{from the sum}$$

$$(y_i - \bar{y}) = (y_i - \hat{y}) + (\hat{y}_i - \bar{y}_i)$$

## Sum of squares

$$\sum \left(y_i - \bar{y}\right)^2 = \sum \left(\hat{y}_i - \bar{y}\right)^2 + \sum \left(y_i - \hat{y}_i\right)^2 + 2 \sum \left(\hat{y}_i - \bar{y}\right)\left(y_i - \hat{y}_i\right)$$

$$\left(\begin{array}{c}\text{Sum of squares} \\ \text{about mean} \\ \text{SST corrected}\end{array}\right) = \left(\begin{array}{c}\text{Sum of squares} \\ \text{due to regression} \\ \text{SSR}(b_1/b_0)\end{array}\right) + \left(\begin{array}{c}\text{Sum of square} \\ \text{about regression} \\ \text{Residual SS}\end{array}\right)$$

$$SST = SSR + SSE.$$

Therefore to test if regression is significant ie
$$H_0 : \beta_1 = 0 \quad Vs \quad H_1 \quad \beta_1 \neq 0$$
We have the anova.

## ANOVA Table.

| Source of variation | Degrees of freedom | Sum of Squares SS | Mean Square | F $F = \dfrac{MSReg}{S^2}$ |
|---|---|---|---|---|
| Due to regression | 1 | $\sum \left(\hat{y}_i - \bar{y}\right)^2$ | $MSreg$ | |
| About regression Residual | $n-2$ | $\sum \left(y_i - \hat{y}_i\right)^2$ | $\delta^2 = \dfrac{SSE}{n-2}$ | |
| Total corrected for $\bar{y}$ | $n-1$ | $\sum \left(y_i - \bar{y}\right)^2$ | | |

An accompanying p value will be available (given)

p-value = probability of observing a value larger than that observed.

In this case it is a 2 sided hypothesis

$$p\text{-value} = 2 * P\left(F \geq F_c(1, n-2)\right)$$

- If # $F_c$ is large reject $H_0$ $\beta_1=0$ : $\Rightarrow$ implying that most variation is caused by (explained by) regression.
- Reject $H_0$ if p-value is too small

## Coefficient of Determination

Coefficient of Determination is defined as

$$r^2 = \frac{SSR}{SST} = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

$SSR$ – Regression sum of squares $= \sum(\hat{y}_i - \bar{y})^2$
$SST$ – Total sum of squares $= \sum(y_i - \bar{y})^2$

$$SST = SSR + SSE$$

$$\sum(y_i - \bar{y})^2 = \sum_{i=1}^{n}(\hat{y}_i - \bar{y})^2 + \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

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

$$SSR = SST - SSE$$

$$r^2 = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$$

Proportion of variation in $y$ explained by the model or accounted by regression.

$$r = \frac{S_{xy}}{\sqrt{S_x^2 S_y^2}} = \frac{\sum(x_i - \bar{x})\sum(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}}$$

$$r^2 = \frac{SSR}{SST} = \frac{14873.0}{17944} = 0.8288$$

$$t = \frac{\hat{\beta}_1}{s/\sqrt{\sum_i(x_i - \bar{x})^2}} = \frac{\sqrt{n-2} \, r}{\sqrt{1 - r^2}}$$

The fit of the regression is "good" if the sum $\sum \hat{e}_i^2 \, (= SSE)$ is "small" ie the unexplained part of the variance of $y$ is small. If SSE is small then $R^2$ is close to 1

Coefficient of Determination. (Multiple Correlation Coefficient)

$$\gamma^2 = \frac{SSR}{SST} = \frac{\sum (\hat{y_i} - \bar{y})^2}{\sum (y_i - \bar{y})^2} \qquad \text{Where } SST = SSR + SSE$$

$$SST = \sum (y_i - \bar{y})^2 = \sum_i (\hat{y_i} - \bar{y})^2 + \sum_i (y_i - \hat{y})^2$$

$\gamma^2$ gives the proportion of variation in $y$ that is explained by the model or accounted for by regression on $x$.

Sample Correlation Coefficient $\gamma$ between $x$ and $y$

$$\gamma = \frac{S_{xy}}{\sqrt{S_x^2 S_y^2}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{[\sum (x_i - \bar{x})^2][\sum_{i=1} (y_i - \bar{y})^2]}}$$

For Example 1.1 $\quad \gamma^2 = \frac{SSR}{SST} = \frac{14873.0}{17244.3} = 0.8288$

$$\gamma = \sqrt{0.8288} = 0.91$$

Exercise

| Y | X | $x^2$ | $y^2$ | $xy$ |
|---|---|---|---|---|
| 25 | 10 | | | |
| 55 | 18 | | | |
| 50 | 25 | | | |
| 75 | 40 | | | |
| 110 | 50 | | | |
| 138 | 63 | | | |
| 90 | 42 | | | |
| 60 | 30 | | | |
| 10 | 5 | | | |
| 100 | 55 | | | |

(i) Determine the least square estimators $\hat{\beta_0}$ and $\hat{\beta_1}$

(ii) Determine $\gamma^2$ and explain how well your model fits the data

(iii) Use ANOVA to test whether regression is significant.