

743- Regression and Time Series

Mamikon S. Ginovyan

Analysis of Variance and Correlation

Goodness of Fit

Regression residuals can provide a useful measure of the fit between the estimated regression line and the data.

- A **good regression equation** is one which **allows explain a large proportion of the variance** of Y .
- **Large residuals imply a poor fit**, while **small residuals imply a good fit**.

Thus, residuals (more precisely, the **SSE**) can be used as a measure of goodness of fit.

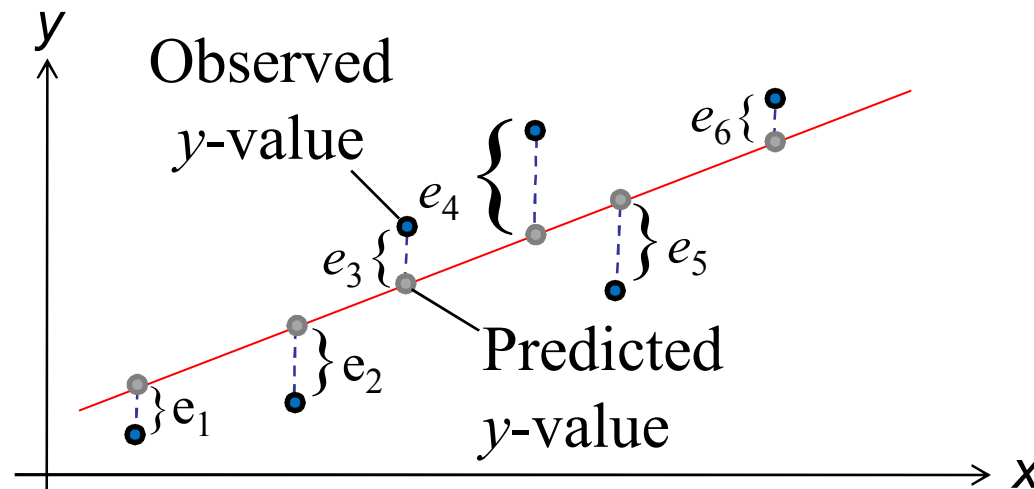
Goodness of Fit

Recall Residuals

- The **difference** between the **observed y-value** and the **predicted y-value** for a given **x-value** on the line.

For a given x -value,

$$e_i = (\text{observed } y\text{-value}) - (\text{predicted } y\text{-value})$$



Goodness of Fit

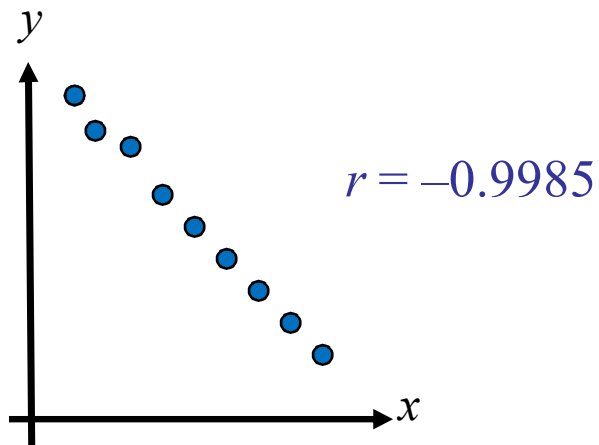
Note: **SSE** can be used as a measure of goodness of fit (explaining y-variation).

(a) Perfect Fit	(b) Moderate Fit	(c) Poor Fit
All variation explained.	Most variation explained.	Regression line does not fit the data, little variation explained.
<i>SSE=0</i>	<i>SSE</i> is small	<i>SSE</i> is large

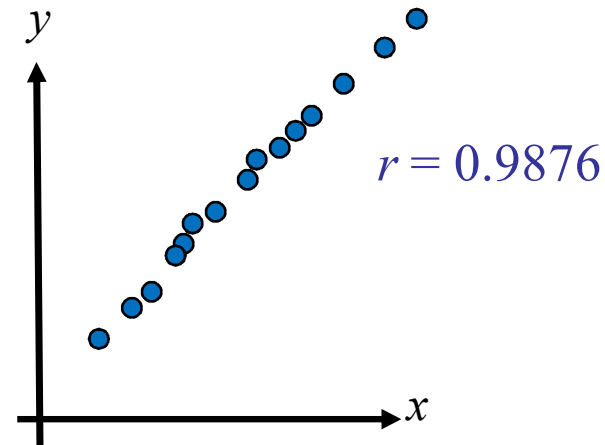
Goodness of Fit

✓ Example 1.

Consider the following **scatter plots** (explaining y-variation).



Strong negative correlation



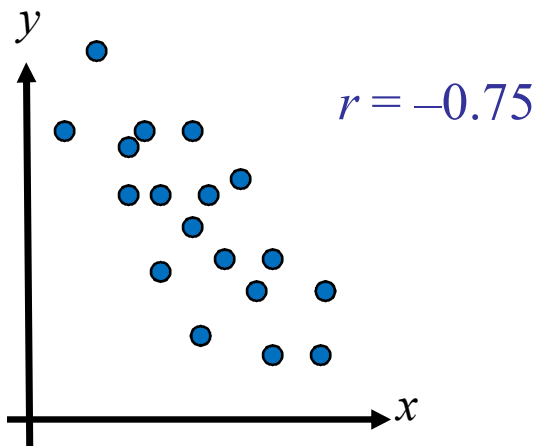
Strong positive correlation

(a) Perfect Fit

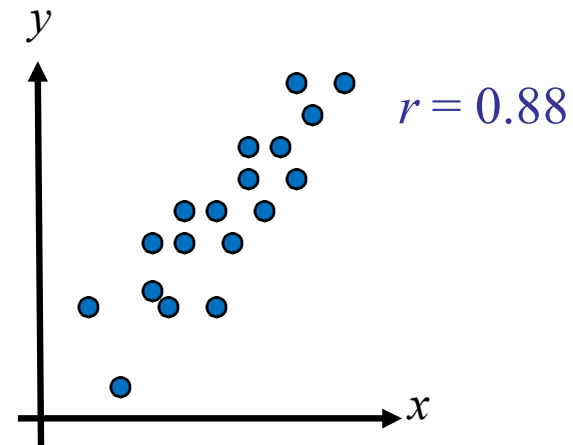
All variation explained.

$$SSE=0$$

Goodness of Fit



Strong negative correlation



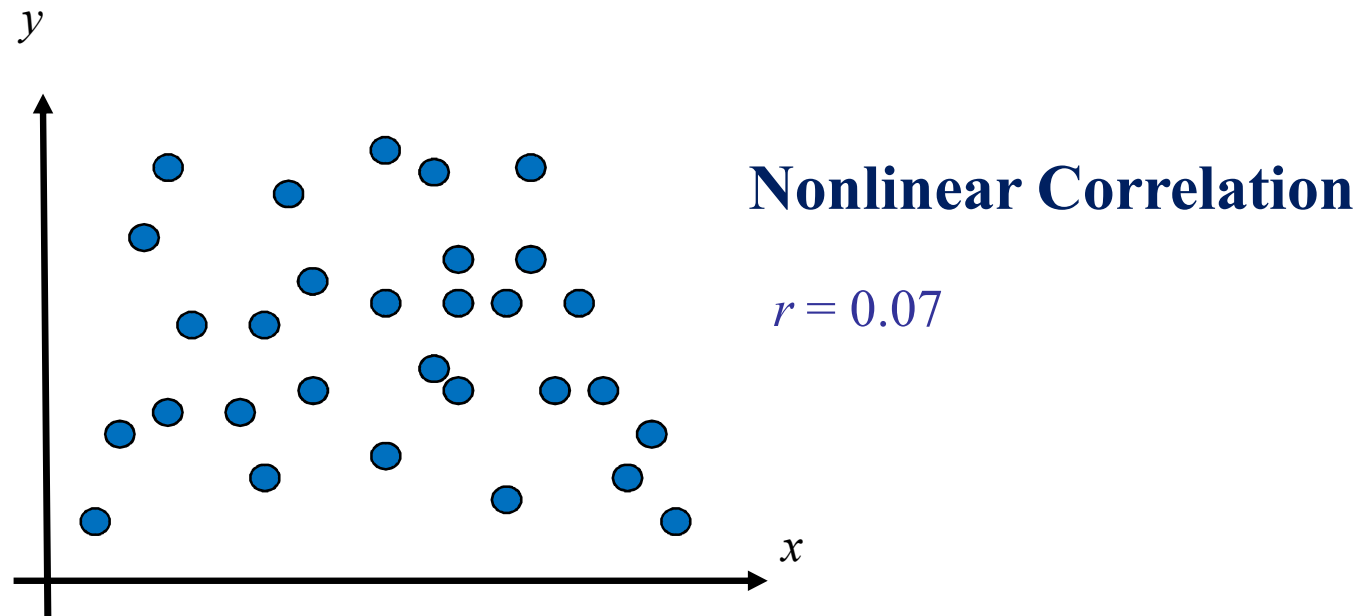
Strong positive correlation

(b) Moderate Fit

Most variation explained.

***SSE* is small**

Goodness of Fit



(c) Poor Fit

**Regression line does not fit the data,
little variation explained.**

***SSE* is large**

Goodness of Fit

A **quantitative** measure of the **total amount of variation** in observed y -values is given by the **total sum of squares**:

$$\text{Variation } (y) = SST = S = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (1)$$

the sum of squared deviation about the sample mean of the observed y -values.

Let $\hat{y} = \hat{a} + \hat{b}x$ be the **estimated regression line**, then for all observations we have

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \quad (2)$$

Goodness of Fit

From (1) and (2) we obtain the following

ANOVA Identity for Regression

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= SSE + SSR. \end{aligned} \quad (3)$$

(the cross-product is equal to 0), where

Goodness of Fit

$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ is the **Error Sum of squares**
(equivalently, the **Residual Sum of Squares**),

$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ is the **Regression Sum of Squares**,

$SST = \sum_{i=1}^n (Y_i - \bar{y})^2$ is the **Total Sum of Squares**.

Goodness of Fit

Observe that

- SSE measures the unexplained variation of the data, while
- SSR measures variation explained by the linear relationship, that is,
SSR is the amount of total variation that is explained by the model.

Observe that (by (3))

$$SSE < SST$$

(**unless** the horizontal line $\hat{Y} = \bar{Y}$, $\hat{b} = 0$, is the **least squares line**).

Goodness of Fit

The **ratio**

$$\frac{SSE}{SST}$$

is the proportion of total variation that cannot be explained by the simple regression model, and

$$1 - \frac{SSE}{SST} = \frac{SSR}{SST}$$

is the proportion of the observed y -variation explained by the model.

This leads to the following definition.

Coefficient of Determination

Definition.

The number

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

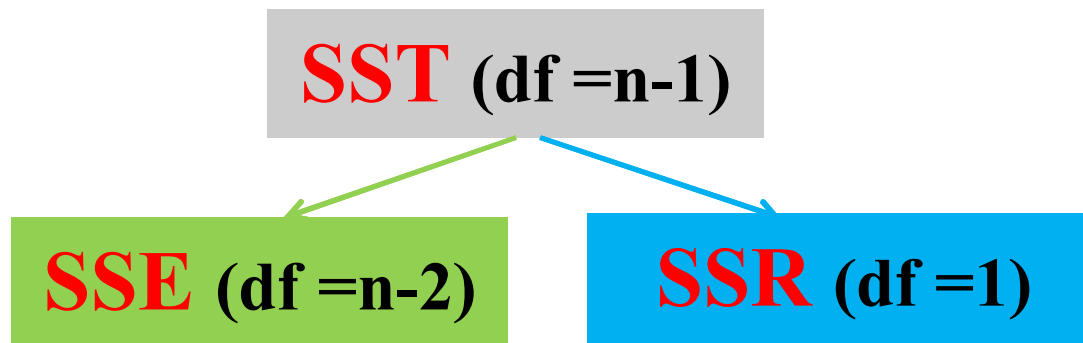
is called the **coefficient of determination**, and describes the proportion of observed y -variation that can be explained by the **simple linear regression model** (attributed to an approximate linear relationship between y and x).

ANOVA

It is occasionally useful to summarize the breakdown of the y - variation in terms of an ANOVA.

Taking into account the **ANOVA Identity**:

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= SSE + SSR, \end{aligned}$$



we can construct the following ANOVA table, which is also useful for hypothesis testing problem.

ANOVA

ANOVA Table for Simple Linear Regression

Source of Variation	df	Sum of Squares	Mean Squares	F
Regression	1	SSR	$MSR = SSR / 1$	$\frac{SSR}{SSE / (n - 2)}$
Error	$n - 2$	SSE	$MSE = s^2 = \frac{SSE}{n - 2}$	
Total	$n - 1$	SST		

ANOVA

Remark.

The hypothesis

$$H_0 : b = 0 \quad \text{vs.} \quad H_a : b \neq 0$$

can be tested using the **ANOVA Table** and the **Decision Rule:**

Reject H_0 if $F \geq F_{a,1,(n-2)}$.

Thus, the F -test gives exactly the same result as the **model utility t -test** because

$$T^2 = F \quad \text{and} \quad t_{a/2,(n-2)}^2 = F_{a,1,(n-2)}.$$

Correlation vs. Coefficient of Determination

Recall that the correlation $r(X, Y)$ is defined by

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \cdot \sqrt{\text{Var}(Y)}}, \quad \text{and}$$

$$r_{XY} = \hat{r}_{XY} = \frac{S_{XY}}{S_X S_Y} = \frac{\frac{1}{n-1} \sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\frac{1}{n-1} \sum (X_i - \bar{X})^2} \cdot \sqrt{\frac{1}{n-1} \sum (Y_i - \bar{Y})^2}}, \quad (1)$$

$$S_X^2 = S_{XX}$$

is the **point estimator** for $r(X, Y)$, and

$$r_{xy} = \hat{r}_{xy} = \frac{S_{xy}}{S_x S_y} \quad \text{is the } \mathbf{\underline{point estimate}} \text{ for } r(X, Y).$$

Correlation vs. Coefficient of Determination

Remark.

Taking into account that

$$\hat{b} = \frac{S_{XY}}{S_{xx}} \quad (S_X^2 = S_{XX}). \quad (2)$$

we obtain the following useful equality.

$$r_{XY} = \hat{r}_{XY} = \frac{S_{XY}}{S_X S_Y} = \frac{S_{XY}}{S_{xx}} \cdot \frac{S_x}{S_Y} = \hat{b} \cdot \frac{S_x}{S_Y}. \quad (3)$$

Correlation vs. Coefficient of Determination

Recall also that,
for any two RV's X and Y ,

$$(a) \quad |r(X, Y)| \leq 1$$

$$(b) \quad |r(X, Y)| = 1 \Leftrightarrow Y = aX + b.$$

Interpreting $r_{XY} = \hat{r}_{XY}$.

The above result is not sufficient to provide a useful interpretation of $r_{XY} = \hat{r}_{XY}$.

Correlation vs. Coefficient of Determination

For example, what does it mean to say that the sample correlation coefficient is .73, .55 or .51?

One way to answer such a question focuses on the **square** of $r_{XY} = \hat{r}_{XY}$ rather than $r_{XY} = \hat{r}_{XY}$ on itself.

Now we show that

$$r_{XY}^2 = \hat{r}_{XY}^2 = R^2, \quad (4)$$

where $R^2 = \frac{SSR}{SST}$ is the **coefficient of determination**.

Correlation vs. Coefficient of Determination

Indeed, first observe that

$$\begin{aligned}SSE &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}x_i)^2 \quad (\text{since } \hat{a} = \bar{Y} - \hat{b}\bar{x}) \\&= \sum_{i=1}^n (Y_i - \bar{Y} + \hat{b}\bar{x} - \hat{b}x_i)^2 = \sum_{i=1}^n \left[(Y_i - \bar{Y}) - \hat{b}(x_i - \bar{x}) \right]^2 \\&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + \hat{b}^2 \sum_{i=1}^n (x_i - \bar{x})^2 - 2\hat{b} \sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x}) \\&= S_{YY} + \hat{b}^2 S_{XX} - 2\hat{b} S_{XY} \quad (\text{since } S_{XY} = \hat{b} S_{XX}) \\&= S_{YY} + \hat{b}^2 S_{XX} - 2\hat{b}^2 S_{XX} = S_{YY} - \hat{b}^2 S_{XX}.\end{aligned}$$

Correlation vs. Coefficient of Determination

Thus

$$SSE = S_{YY} - \hat{b}^2 S_{XX}. \quad (5)$$

Now by definition of R^2 and (5) we have

$$\begin{aligned} R^2 &= 1 - \frac{SSE}{SST} = 1 - \frac{S_{YY} - \hat{b}^2 S_{XX}}{S_{YY}} \\ &= 1 - 1 + \hat{b}^2 \frac{S_{XX}}{S_{YY}} = \hat{b}^2 \frac{S_{XX}}{S_{YY}} \quad (by(3)) \\ &= \hat{r}_{XY}^2 = r_{XY}^2. \end{aligned}$$

Correlation vs. Coefficient of Determination

Remark 1.

The equality (4) we proved under **linear relationship** between X and Y . **In general**, it is not true.

More precisely, if X and Y are in **non-linear relationship**, then

$$R^2 \geq r_{XY}^2, \text{ and the difference } R^2 - r_{XY}^2$$

is a measure of non-linearity.

Correlation vs. Coefficient of Determination

Remark 2.

The definition of the coefficient of determination $R = R_{YX}$ is based on the **Conditional Variance Formula**:

$$\text{Var}(Y) = \text{Var}(E[Y|X]) + E[(\text{Var}(Y|X))],$$

and is defined to be

$$R_{YX}^2 = 1 - \frac{E[\text{Var}(Y|X)]}{\text{Var}(Y)}.$$

Correlation vs. Coefficient of Determination

Properties of the Coefficient of Determination R_{YX}^2

1. $0 \leq R_{YX}^2 \leq 1$, and $R_{YX}^2 = 1$ if and only if

$$E[Y - E(Y|X = x)]^2 = 0 \Leftrightarrow P(Y = E[Y|X = x]) = 1,$$

that is, there is a **functional relationship** between X and Y .

2. In general, R_{YX}^2 is **not symmetric**, that is, $R_{YX}^2 \neq R_{XY}^2$.

3. $R_{YX}^2 = r_{YX}^2$ if $Y = aX + b$, otherwise $R_{YX}^2 > r_{YX}^2$.

4. If $E[Y|X = x]$ is **independent** of x , then $R_{YX} = 0$.

This is, the case when the RV's X and Y are **independent**.

Hypothesis Test for Correlation

- To test the null hypothesis of **no linear association**:

$$H_0: \rho = 0$$

we use the **test statistic**:

$$T = \frac{r}{\sqrt{\frac{1 - r^2}{n - 2}}} = \frac{r \sqrt{n - 2}}{\sqrt{1 - r^2}} \sim t_{(n - 2)}$$

which follows the **Student's t-distribution** with $(n - 2)$ degrees of freedom:

Remark: Hypothesis Test for Correlation

- The test of the hypothesis $H_0: \rho = 0$ of **no linear association** may also be based on the **Fisher statistic**:

$$z = \frac{Z - m_Z}{s_Z} = (Z - m_Z) \sqrt{n-3} \sim N(0,1),$$

where

$Z = \frac{1}{2} \cdot \ln \frac{1+r}{1-r}$ is the **Fisher Z- transformation**,

$$m_Z = \frac{1}{2} \cdot \ln \frac{1+r}{1-r}; \quad s_Z = \frac{1}{\sqrt{n-3}} \quad \text{and} \quad z(\text{obs}) = (Z(\text{obs}) - m_Z) \sqrt{n-3}.$$

Observe that under $H_0: \rho = 0$, we have $m_Z = (.5) \cdot \ln 1 = 0$.

Example: Hypothesis Test for Correlation

Age and Price of Orions. The data on age and price for a sample of **11** of Orions are given in the following table.

Age (yr) x	5	4	6	5	5	5	6	6	2	7	7
Price (\$100) y	85	103	70	82	89	98	66	95	169	70	48

At the **5%** significance level, do the data provide sufficient evidence to conclude that age and price of Orions are **negatively linearly correlated?**

Example: Hypothesis Test for Correlation

Solution:

Step 1. State the null and alternative hypotheses.

Let ρ denote the population **linear correlation coefficient** for the variables age and price of Orions. Then

$H_0: \rho = 0$ (age and price are linearly **uncorrelated**)

$H_1: \rho < 0$ (age and price are linearly **correlated**)

Step 2. Compute the **observed value** of test statistic:

$$t_0 = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}} = \frac{-0.924}{\sqrt{\frac{1-(-0.924)^2}{11-2}}} = -7.249.$$

Example: Hypothesis Test for Correlation

Calculations: We have

$$\sum_{i=1}^{11} x_i = 58; \quad \sum_{i=1}^{11} y_i = 975; \quad \sum_{i=1}^{11} x_i y_i = 4732;$$

$$\sum_{i=1}^{11} x_i^2 = 326; \quad \sum_{i=1}^{11} y_i^2 = 96129;$$

$$r = r_{XY} = \hat{r}_{XY} = \frac{S_{XY}}{S_X S_Y} = \frac{SS_{xy}}{\sqrt{SS_{xx} SS_{yy}}} = -0.924.$$

Example: Hypothesis Test for Correlation

Step 4. Decide on the significance level: $\alpha = .05$.

The **critical value** for $df = n-2 = 11-2 = 9$ and $\alpha = .05$ is

$$-t_{\alpha, n-2} = -t_{.05, 9} = -1.833,$$

The **critical region** is: $CR = \{T < -1.833\}$.

Step 5. Since $t_0 = -7.249 < -1.833 = t_{\alpha}$, we **reject** $H_0: \rho = 0$, and conclude that the test results are **statistically significant** at the **5%** level of significance.

Inferences about the mean

Inferences about the mean $m_{Y|X} = E[Y|X = x]$.

In addition to statistical inference about the regression parameters (*unknown*) a, b and S^2 , we will sometimes find it also helpful to draw inferences about the regression line

$$m_{Y|x_0} = E[Y|X = x_0] = a + b x_0, \quad (1)$$

where x_0 is a specified value of **independent** variable x .

Inferences about the mean

Once the point estimators \hat{a} and \hat{b} have been calculated,

$$\hat{Y} = \hat{a} + \hat{b}x_0 \quad (2)$$

can be regarded either as:

- a **point estimate** of $m_{Y|x_0}$, or
- a **prediction** of the Y value that will result from a single observation made when $x = x_0$.

Inferences about the mean

Observe that the point estimate or prediction by itself gives no information concerning how precisely $m_{Y|x_0}$ has been **estimated** or Y has been **predicted**.

This can be remedied by developing CI's for $m_{Y|x_0}$ and a **prediction interval (PI)** for a single Y value.

Thus, as a point estimator for **regression line** $m_{Y|x_0}$ given by (1) we consider the statistic \hat{Y} given by (2).

Inferences about the mean

Properties of point estimator $\hat{Y} = \hat{a} + \hat{b} x_0$.

1. \hat{Y} is an **unbiased estimator** for $m_{Y|x_0}$, that is,

$$E[\hat{Y}] = a + b x_0. \quad (3)$$

Indeed, since \hat{a} and \hat{b} are unbiased estimators for a and b , respectively, we have

$$\begin{aligned} E[\hat{Y}] &= E[\hat{a} + \hat{b} x_0] \\ &= E[\hat{a}] + x_0 E[\hat{b}] = a + b x_0. \end{aligned}$$

Properties of Point Estimator

2. The **variance** $Var(\hat{Y})$ of \hat{Y} is given by

$$s_{\hat{Y}}^2 = Var(\hat{Y}) = s^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]. \quad (4)$$

2'. The **estimated variance** $s_{\hat{Y}}^2$ is given by

$$s_{\hat{Y}}^2 = Var(\hat{Y}) = s^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right], \quad \text{where}$$

$$S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2. \quad (5)$$

Properties of Point Estimator

Proof of (4).

We have

$$\begin{aligned} \text{Var}(\hat{Y}) &= \text{Var}(\hat{a} + \hat{b}x_0) \quad (\text{since } \hat{a} = \bar{Y} - \hat{b}\bar{x}) \\ &= \text{Var}(\bar{Y} - \hat{b}\bar{x} + \hat{b}x_0) \\ &= \text{Var}(\bar{Y} + \hat{b}(x_0 - \bar{x})) \\ &\quad (\text{since } \bar{Y} \text{ and } \hat{b} \text{ are independent}) \\ &= \text{Var}(\bar{Y}) + \text{Var}[\hat{b}(x_0 - \bar{x})] \\ &= \text{Var}(\bar{Y}) + (x_0 - \bar{x})^2 \text{Var}(\hat{b}) \\ &= \frac{s^2}{n} + (x_0 - \bar{x})^2 \cdot \frac{s^2}{s_{xx}} = s^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}} \right]. \end{aligned}$$

Properties of Point Estimator

3. \hat{Y} has a **normal distribution**: $\hat{Y} \sim N(a + b x_0, s_{\hat{Y}}^2)$.

Proof.

Taking into account that

$$S_{xY} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) = \sum_{i=1}^n (x_i - \bar{x})Y_i, \quad (6)$$

$$\begin{aligned} (\text{since } \sum_{i=1}^n (x_i - \bar{x})\bar{Y} &= \bar{Y} \sum_{i=1}^n (x_i - \bar{x}) \\ &= \bar{Y} (n\bar{x} - n\bar{x}) = 0.) \end{aligned}$$

we obtain,

Properties of Point Estimator

$$\hat{Y} = \hat{a} + \hat{b}x_0 = \bar{Y} + \hat{b}(x_0 - \bar{x}) \quad (\text{since } \hat{a} = \bar{Y} - \hat{b}\bar{x})$$

$$= \bar{Y} + \frac{S_{xY}}{S_{xx}}(x_0 - \bar{x}) \quad (\text{since } \hat{b} = \frac{S_{xY}}{S_{xx}})$$

$$= \frac{1}{n} \sum_{i=1}^n Y_i + \sum_{i=1}^n (x_i - \bar{x}) Y_i \frac{(x_0 - \bar{x})}{S_{xx}} \quad (\text{by (6)})$$

$$= \sum_{i=1}^n \left[\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}} \right] Y_i = \sum_{i=1}^n a_i Y_i. \quad (7)$$

Properties of Point Estimator

Thus, $\hat{Y} = \sum_{i=1}^n a_i Y_i$, where the coefficients $a_i, i = \overline{1, n}$

involve the x_i 's and x_0 , all of which are **fixed**.

Because Y_i are **normally distributed** and **independent**,
 \hat{Y} as a **linear combination** of Y_i 's will also be **normally distributed**.

A Confidence Interval for regression line

A Confidence Interval Regression line:

$$m_{Y|x_0} = E[Y | X = x_0].$$

Using **Properties 1-3** we can state

Theorem 1. The statistic

$$T = \frac{\hat{Y} - (a + b x_0)}{s_{\hat{Y}}} \sim t_{(n-2)} \quad (8)$$

has a **Student t -distribution** with $(n - 2)$ df.

A Confidence Interval for regression line

Proof.

The statistic T can be represented as follows

$$T = \frac{\hat{Y} - (a + b x_0)}{s_{\hat{Y}}} \cdot \left[\sqrt{\frac{(n-2)s^2}{s^2}} \right]^{-1} \cdot (n-2).$$

Since $\frac{\hat{Y} - (a + b x_0)}{s_{\hat{Y}}} \sim N(0,1)$ and

$$\frac{(n-2)s^2}{s^2} \sim c^2(n-2),$$

the result follows from the **definition** of **t -distribution**.

A Confidence Interval for regression line

Now using **standard t -procedure** with $(n - 2)$ df, we can write, for given a ($0 \leq a \leq 1$)

$$P(-t_{a/2, (n-2)} \leq T \leq t_{a/2, (n-2)}) = 1 - a. \quad (9)$$

Substituting, T from (8) into (9) we obtain

$$P(-t_{a/2, (n-2)} \leq \frac{\hat{Y} - m_{Y|x_0}}{s_{\hat{Y}}} \leq t_{a/2, (n-2)}) = 1 - a.$$

Solving the inside inequalities for $m_{Y|x_0}$ we get the **desired CI** for $m_{Y|x_0}$.

Summing up we can state the following result.

A Confidence Interval for regression line

Theorem 2.

For given α ($0 \leq \alpha \leq 1$) a $100(1 - \alpha)\%$ **CI** for **regression line**

$$m_{Y|x_0} = E[Y | X = x_0]$$

is the interval

$$\hat{Y} \pm t_{\alpha/2, (n-2)} \cdot s_{\hat{Y}}, \quad (10)$$

where $s_{\hat{Y}}$ is given by

$$s_{\hat{Y}}^2 = \text{Var}(\hat{Y}) = s^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right].$$

A Prediction Interval for a Future Value of Y

Similar to **CI** (10) for $m_{Y|x_0}$, we frequently wish to obtain an **interval of plausible values for the value** of Y associated with some future observation when the independent variable x has value x_0 , that is,

- the value of Y of a single future observation to be recorded at some given level of $x = x_0$.

Prediction Interval for a Future Value of Y

Remark 1.

A **CI** refers to a **parameter** (= population characteristic), whose value is fixed but unknown to us.

In contrast, a **future value** of Y is not a parameter but instead is a **RV**.

For this reason we refer to an interval of plausible values for a future Y as a **Prediction Interval (PI)** rather than **CI**.

Prediction Interval for a Future Value of Y

- For the **CI** we use the **error of estimation**:

$$(a + b x_0) - (\hat{a} + \hat{b} x_0) =$$

= a **difference** between
a **fixed (but unknown) quantity** and
a **RV**.

- The **error of prediction** is

$$Error = Y - \hat{Y} = Y - (\hat{a} + \hat{b} x_0)$$

$$= (a + b x_0 + e) - (\hat{a} + \hat{b} x_0)$$

= a **difference** between **two RV's**.

Prediction Interval for a Future Value of Y

With the additional random term e , there is **more uncertainty in prediction** than in **estimation**, so a **PI** will be **wider** than a **CI**.

Let $(x_1, Y_1), (x_2, Y_2), \dots, (x_n, Y_n)$ be a set of n observations that satisfy the **Model Assumptions**, and let (x_0, Y) be a hypothetical **future observation**, where Y is **independent** of the Y_i 's, $i = 1, \dots, n$.

Definition: A **prediction interval** is a **range** of numbers that contains Y with a specified **probability** $= (1 - \alpha)$.

Prediction Interval for a Future Value of Y

Consider the difference $Y - \hat{Y} = \text{error of prediction.}$

- For the expectation of this error we have

$$E[Y - \hat{Y}] = E[Y] - E[\hat{Y}] = (a + b x_0) - (a + b x_0) = 0 \quad (1)$$

- For variance of the $Y - \hat{Y}$ we have

$$Var(Y - \hat{Y}) = Var(Y) + Var(\hat{Y})$$

(since \hat{Y} and Y are **independent**)

$$\begin{aligned} &= s^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] + s^2 \\ &= s^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]. \end{aligned} \quad (12)$$

Prediction Interval for a Future Value of Y

As in the **CI case**, we can show that the RV

$$Z = \frac{Y - \hat{Y}}{\sqrt{\text{Var}(\hat{Y} - Y)}} \sim N(0,1).$$

Therefore the RV

$$T = \frac{Y - \hat{Y}}{s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t_{(n-2)} \quad (13)$$

has a **Student t -distribution** with $(n-2)$ df.

Prediction Interval for a Future Value of Y

Therefore, for given a ($0 \leq a \leq 1$),

$$P(-t_{a/2, (n-2)} \leq T \leq t_{a/2, (n-2)}) = 1 - a. \quad (14)$$

Substituting T from (13) into (14), and solving the inside inequalities for Y we get the following result.

Prediction Interval for a Future Value of Y

Theorem 3 (PI for Y).

For given α ($0 \leq \alpha \leq 1$) a $100(1 - \alpha)\%$ **PI** for a future observation Y when $\mathbf{x} = \mathbf{x}_0$ is the interval

$$\begin{aligned} PI &= (\hat{a} + \hat{b}x_0) \pm t_{\alpha/2, (n-2)} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \\ &= (\hat{a} + \hat{b}x_0) \pm t_{\alpha/2, (n-2)} \cdot \sqrt{s^2 + s_{\hat{Y}}^2} \\ &= \hat{Y} \pm t_{\alpha/2, (n-2)} \cdot \sqrt{s^2 + s_{\hat{Y}}^2}. \end{aligned}$$

Prediction Interval for a Future Value of Y

Remark.

Notice that the **length** of a **CI** for $E[Y]$ when $x = x_0$ is given by

$$2 t_{\alpha/2, (n-2)} \cdot s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}, \quad (15)$$

whereas the **length** of a **PI** for an actual value of Y when $x = x_0$ is given by

$$2 t_{\alpha/2, (n-2)} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}. \quad (16)$$

Thus, we observe that **PI** for the actual value of Y is longer than **CI**'s for expectation $E[Y]$ if both are determined for the same value of $x = x_0$.

Testing the Equality of Two Slopes

We often are interested in comparison of two linear XY - relationships. In such situations we test the hypothesis

$$H_0 : b_1 = b_2,$$

where b_1 and b_2 are the true slopes associated with the two regressions.

If the data points taken on the two regressions are all independent, a **two-sample t -test** can be set up based on the properties proved in **Theorem 2 and 3** (for \hat{a} and \hat{b}).

The following result is true.

Testing the Equality of Two Slopes

Theorem 1.

Let $(x_1, Y_1), (x_2, Y_2), \mathbf{L}, (x_n, Y_n)$ and $(x_1^*, Y_1^*), (x_2^*, Y_2^*), \mathbf{L}, (x_n^*, Y_n^*)$ be two **independent** observations, each satisfying the **Model**

Assumptions:

$$E[Y | x_0] = a_1 + b_1 x_0, \quad E[Y^* | x^*] = a_2 + b_2 x^*.$$

(a) Let
$$T = \frac{(\hat{b}_1 - \hat{b}_2) - (b_1 - b_2)}{s \sqrt{\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} + \frac{1}{\sum_{i=1}^m (x_i^* - \bar{x}^*)^2}}},$$

where

Theorem 1.

$$s = \frac{1}{\sqrt{n+m-4}} \sqrt{\sum_{i=1}^n [Y_i - (\hat{a}_1 + \hat{b}_1 x_i)]^2 + \sum_{i=1}^m [Y_i^* - (\hat{a}_2 + \hat{b}_2 x_i^*)]^2}.$$

Then $T \sim t_{(n+m-4)}$ = ***t*-distribution** with ($n + m - 4$) df.

Theorem 1.

(b) To test the hypothesis

$$H_0 : b_1 = b_2 \quad vs. \quad H_1 : b_1 \neq b_2$$

at the α level of significance,

Reject H_0 if either $T_{obs} \leq -t_{\alpha/2, (n+m-4)}$ or $T_{obs} \geq t_{\alpha/2, (n+m-4)}$,

where

$$T_{obs} = \frac{(\hat{b}_1 - \hat{b}_2)}{s \sqrt{\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} + \frac{1}{\sum_{i=1}^m (x_i^* - \bar{x}^*)^2}}}.$$

Theorem 1.

Remark.

One-sided tests are defined in the usual way by replacing the critical values

$$\pm t_{\alpha/2, (n+m-4)}$$

by $t_{\alpha/2, (n+m-4)}$ for **upper-tailed** test $(H_1 : b_1 > b_2)$, or

by $-t_{\alpha, (n+m-4)}$ for **upper-tailed** test $(H_1 : b_1 < b_2)$.

Summary Example

Suppose we have a data set consisting of $n = 5$ points. The data set is given by the following table.

x	-2	-1	0	1	2
y	0	0	1	1	3

- (i) Use the method of least squares to fit a straight line to given points.
- (ii) Find the variance of the estimates \hat{a} and \hat{b} obtained in part (i).
- (iii) Estimate s^2 from the data.
- (iv) Calculate a 95% CI for the slope parameter b .

Summary Example

- (v) Do the data present sufficient evidence to indicate that the slope differs from 0? **(Model Utility Test)**.

Test the hypothesis at $\alpha = .05$ level of significance using

- (a) Critical-value method
- (b) P -value method
- (c) CI-method.

- (vi) Find a 90% CI for the regression line $m_{y|x_0} = E[Y|X = x_0]$, when $x_0 = 1$.

- (vii) Suppose that the experiment that generated the data in Part (i) is to be run again with $x = x_0 = 2$.

Predict the particular value of Y with $\alpha = .1$, that is, construct 90% PI for $m_{Y|x=2} = E[Y|x = 2]$.

Summary Example

Solution.

(i) Use the method of least squares to fit a straight line to given points.

First we construct the following **Calculation Table**.

x_i	y_i	$x_i y_i$	x_i^2	y_i^2
-2	0	0	4	0
-1	0	0	1	0
0	1	0	0	1
1	1	1	1	1
2	3	6	4	9
$\sum_{i=1}^n x_i = 0$	$\sum_{i=1}^n y_i = 5$	$\sum_{i=1}^n x_i y_i = 7$	$\sum_{i=1}^n x_i^2 = 10$	$\sum_{i=1}^n y_i^2 = 11$

Summary Example-Solution

Compute

$$S_{xy} = \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i y_i = 7 - \frac{1}{5}(0)(5) = 7,$$

$$S_{xx} = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 = 10 - \frac{1}{5}(0)^2 = 10.$$

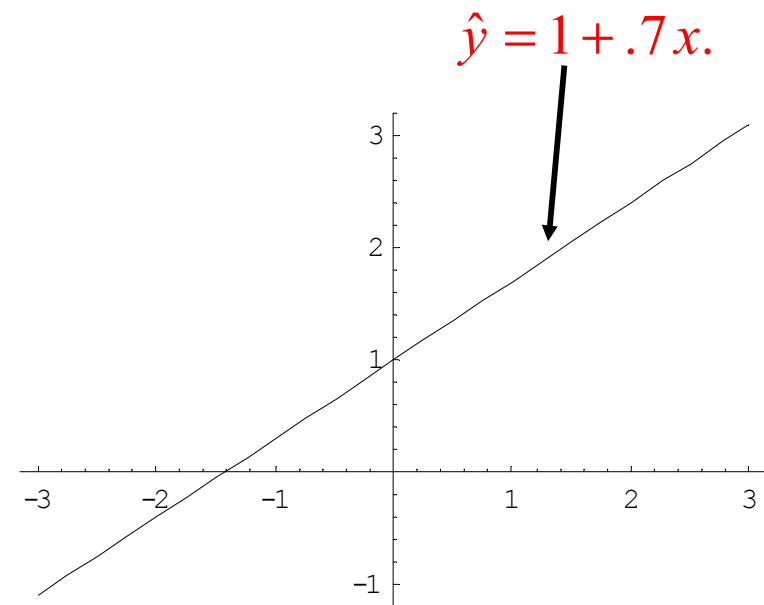
Thus,

$$\hat{b} = \frac{S_{xy}}{S_{xx}} = \frac{7}{10} = .7,$$

$$\hat{a} = \bar{y} - \hat{b} \bar{x} = \frac{5}{5} - (.7)(0) = 1.$$

Therefore the **fitted line** is:

$$\hat{y} = \hat{a} + \hat{b}x = 1 + .7x.$$



Summary Example-Solution

(ii) Find the variance of the estimates \hat{a} and \hat{b} obtained in Part (i).

Solution. We have

$$Var(\hat{a}) = \frac{s^2 \sum_{i=1}^n x_i^2}{nS_{xx}} = \frac{s^2(10)}{5(10)} = \frac{1}{5}s^2$$

and

$$Var(\hat{b}) = \frac{s^2}{S_{xx}} = \frac{s^2}{10}.$$

Remark. Notice that in this example $\sum x_i = 0$. Hence $Cov(\hat{a}, \hat{b}) = 0$, since

$$Cov(\hat{a}, \hat{b}) = -\frac{\bar{x}s^2}{S_{xx}}, \quad \text{where } s^2 = s_{LS}^2, \quad \text{and } \bar{x} = \frac{1}{n} \sum x_i = 0.$$

Summary Example-Solution

(iii) Estimate s^2 from the data.

Solution. We have

$$\sum_{i=1}^5 y_i = 5, \quad \sum_{i=1}^5 y_i^2 = 11, \quad S_{xy} = 7, \quad \hat{b} = .7, \quad \bar{y} = 1.$$

Hence

$$S_{yy} = \sum_{i=1}^5 (y_i - \bar{y})^2 = \sum_{i=1}^5 y_i^2 - 5(\bar{y}^2) = 11 - 5(1)^2 = 6,$$

and

$$SSE = S_{yy} - \hat{b}S_{xy} = 6 - (.7)(7) = 1.1.$$

Therefore

$$s^2 = \frac{SSE}{n-2} = \frac{1.1}{5-2} = \frac{1.1}{3} = .367.$$

Summary Example-Solution

(iv) Calculate a 95% CI for the slope parameter b .

Solution.

We have that for $0 \leq a \leq 1$, a $100(1-a)\%$ CI for b is the interval

$$\hat{b} \pm t_{a/2, (n-2)} s_{\hat{b}} = \hat{b} \pm t_{a/2, (n-2)} \frac{s}{\sqrt{S_{xx}}} \quad (1)$$

Now we compute the quantities in (1).

Since

$$n = 5 \Rightarrow n - 2 = 3, a = .05 \Rightarrow \frac{a}{2} = .025.$$

So from t -table we find

$$t_{a/2, (n-2)} = t_{.025, 3} = 3.182. \quad (2)$$

Summary Example-Solution

From Part (iii) we have

$$s = \sqrt{\frac{SSE}{n-2}} = \sqrt{.367} = .606. \quad (3)$$

From Part (i) we have

$$\hat{b} = .7 \quad \text{and} \quad S_{xx} = 10. \quad (4)$$

Substituting (2) - (4) into (1) we get the desired 95% CI for β

$$\begin{aligned} & \hat{b} \pm t_{\alpha/2, (n-2)} \frac{s}{\sqrt{S_{xx}}} \\ & = .7 \pm (3.182)(.606)\sqrt{.1} = .7 \pm .61 = (.09, 1.31). \end{aligned}$$

Summary Example-Solution

Remark. If we wish to estimate β correct to within .15 unit, it is obvious that the CI is too wide and that the sample size ($n=5$) must be increased.

(v) Do the data present sufficient evidence to indicate that the slope differs from 0? (Model Utility Test).

Test the hypothesis at $\alpha=.05$ level of significance using

- (a) Critical-value method**
- (b) P -value method**
- (c) CI-method.**

Summary Example-Solution

Solution.

(a) **Step 1.** Set up the hypothesis:

$$H_0 : b = 0 \quad vs. \quad H_1 : b \neq 0. \quad (b_0 = 0)$$

Step 2. Specify test statistic and its distribution under H_0

$$TS = T = \frac{\hat{b} - b}{s_{\hat{b}}}.$$

under

$$H_0 : b = 0, T = \frac{\hat{b} - 0}{s_{\hat{b}}} = \frac{\hat{b}}{s_{\hat{b}}} \sim t_{(n-2)}.$$

Summary Example-Solution

Step 3. Compute the observed value of TS:

$$T_0 = T_{ob} = \frac{\hat{b} - 0}{s_{\hat{b}}} = \frac{.7 - 0}{.192} = 3.65,$$

Since

$$s_{\hat{b}} = \frac{s}{\sqrt{S_{xx}}} = (.606)(\sqrt{.1}) = .192.$$

Step 4. Compute the Critical Values and set up the Rejection Region (RR).

Since the test is two-sided, the Critical Values are $\pm t_{\alpha/2} = \pm t_{\alpha/2, (n-2)}$.

For $n = 5$ and $\alpha = .05$ from t -table we find

$$\pm t_{\alpha/2, (n-2)} = \pm t_{.025, 3} = \pm 3.182.$$

So the **RR** is: $RR = \{t : |T| \geq 3.182\}.$

Summary Example-Solution

Step 5. Decision (Use **Decision Rule** based on CV method).

Since **the observed value** of TS: $T_{ob} = 3.65$
falls the RR ($3.65 > 3.182!$),

we **reject** H_0 at $\alpha = .05$ level of significance.

Step 6. Conclusion (The answer of question)

Yes, the data provide sufficient evidence to indicate that the slope β differs from 0.

Summary Example-Solution

(b) P-value approach

Steps 1-3 are the same as in CV approach.

Step 4' Compute the P -value corresponding to the observed value of $TS : T_{ob} = 3.65$.

Since the test is two-sided the P -value is

$$P\text{-value} = 2.P(T > |T_{ob}|) = 2.P(T > 3.65),$$

where $T \sim t_{(n-2)} = t_3$.

Using t -table we find that

$$.01 < P(T > 3.65) < .025.$$

Therefore

$$.02 < P\text{-value} < .05.$$

Summary Example-Solution

Remark.

In contrast to Z -table, t -table is not reach, by this reason we can find only limits for P -values.

Step 5' Decision

(Use **Decision Rule** based on P -value approach)

Since given significance level $\alpha = .05 > P - value$, we **reject** H_0 .

Step 6' (**Conclusion**) is the same as **Step 6** in CV approach.

Summary Example-Solution

(c) CI-approach

Step1. Compute $100(1 - \alpha)\% = 100(1 - .05)\% = 95\%$ CI for b .
In Part (iv) we have constructed this **CI**:

$$\hat{b} \pm t_{\alpha/2, (n-2)} \cdot s_{\hat{b}} = (.09, 1.31).$$

Step 2. Decision (Use **Decision Rule** based on CI approach)
Since the **specified value** $b_0 = 0$ does not fall the **95% CI**,
we **reject** H_0 .

Summary Example-Solution

**(vi) Find a 90% CI for the regression line $m_{y|x_0} = E[Y|X = x_0]$,
when $x_0 = 1$.**

Solution.

We know that a $100(1-a)\%$ CI for $m_{y|x_0}$ is

$$\hat{Y} \pm t_{a/2, (n-2)} \cdot s_{\hat{Y}} = \hat{a} + \hat{b}x_0 \pm t_{a/2, (n-2)} \cdot s_{\hat{Y}}, \quad (1)$$

where

$$s_{\hat{Y}} = s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}.$$

Summary Example-Solution

We have $\hat{a} = 1, \hat{b} = .7, s = .606$. (2)

For $n = 5, \bar{x} = 0, S_{xx} = 10$, we find

$$\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} = \frac{1}{5} + \frac{(1-0)^2}{10} = .3. \quad (3)$$

For $n = 5, a = .10, \frac{a}{2} = .05$, from t -table we find

$$t_{a/2, (n-2)} = t_{.05, 3} = 2.353. \quad (4)$$

Summary Example-Solution

Substituting (2) - (4) into (1) we find the desired 90% CI for

$$m_{y|x=1} = E[Y|x=1]:$$

$$\begin{aligned} & \hat{a} + \hat{b}x_0 \pm t_{\alpha/2, (n-2)} \cdot s_{\hat{Y}} \\ &= [1 + (.7)(1)] \pm (2.353)(.606)\sqrt{.3} \\ &= 1.7 \pm .781 = (.919, 2.481). \end{aligned}$$

Thus, a **90% CI** for $m_{y|x=1} = E[Y|x=1]$ is the interval **(.919, 2.481)**.

Summary Example-Solution

Interpretation of CI.

We are 90% confident that, when the independent variable takes on the value $x=1$, the conditional mean value $m_{y|x=1} = E[Y|x=1]$ of the dependent variable Y is between **.918** and **2.481**.

Obviously, the CI is very wide, but it is not surprising, because it is based on very small ($n=5$) sample size.

Summary Example-Solution

(vii) Suppose that the experiment that generated the data in Part (i) is to be run again with $x = x_0 = 2$. Predict the particular value of Y with $\alpha = .1$, that is, construct 90% PI for $m_{Y|x=2} = E[Y|x = 2]$.

Solution.

We know that for given a ($0 \leq a \leq 1$), a $100(1 - a)\%$

PI for $m_{Y|x_0}$ is

$$\hat{Y} \pm t_{a/2, (n-2)} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}. \quad (1)$$

Summary Example-Solution

From Part (i), we have

$$\hat{a} = 1, \text{ and } \hat{b} = .7.$$

Hence the **predicted value** of Y with $x=2$ is

$$\hat{Y} = \hat{a} + \hat{b}x_0 = 1 + (.7)(2) = 1 + 1.4 = 2.4. \quad (2)$$

From Part (vi), we have for $n = 5, a = .10,$

$$t_{a/2, (n-2)} = t_{.05, 3} = 2.353 \quad (3)$$

From Part (iv), we have $s = .606,$

and from Part (i) $S_{xx} = 10$ and $\bar{x} = 0.$

Summary Example-Solution

Hence for $x_0 = 2$,

$$1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} = 1 + \frac{1}{5} + \frac{(2-0)^2}{10} = 1.6. \quad (4)$$

Substituting (2) - (4) into (1) we obtain

$$\begin{aligned} & \hat{Y} \pm t_{\alpha/2, (n-2)} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \\ &= 2.4 \pm (2.353)(.606)\sqrt{1.6} = 2.4 \pm 1.804 = (.596, 4.204). \end{aligned}$$

Thus, a 90% PI for $m_{Y|x=2}$ is the interval **(.596, 4.204)**.

Summary Example-Solution

Remark.

If we construct a **90% PI** for $m_{Y|x=1}$ (instead of $x=2$), we obtain

$$PI = 1.7 \pm 1.63 = (.07, 3.33).$$

Comparing this **PI** with the **90% CI** for $m_{Y|x=1}$ obtained in Part (iv):

$$CI = 1.7 \pm .781 = (.919, 2.481),$$

we found that **PI** is **wider** than **CI**.