

Simple Linear Regression

Find the Least-Squares Line.

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

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

Computational Formulas for Sum of Squares

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \hat{\beta}_0 \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i$$

$$SSTO = SSTO - SSK$$

$$SSR = \hat{\beta}_1^2 \cdot \left[\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n} \right] = \hat{\beta}_1^2 \cdot \sum_{i=1}^n (x_i - \bar{x})^2$$

$$SSTO = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{(\sum y_i)^2}{n}$$

Coefficient of Determination

$$R^2 = \frac{SSR}{SSTO}$$

Note: $0 \leq R^2 \leq 1$

"large", "moderate", "low".

Coefficient of Correlation

$$r = \sqrt{R^2} \text{ if } \hat{\beta}_1 > 0;$$

$$\text{or } -\sqrt{R^2} \text{ if } \hat{\beta}_1 < 0$$

$$r = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{\sqrt{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \cdot \sqrt{n \sum_{i=1}^n y_i^2 - (\sum_{i=1}^n y_i)^2}}$$

Note: $-1 \leq r \leq 1$ and $|r| \geq R$

If $r=0 \rightarrow$ there is no linear relationship between x and y . But doesn't imply that there is no relationship between these variables.

Correlation does not imply Causation

F -test under ANOVA for simple linear regression model.

$$H_0: \beta_1 = 0 \text{ vs. } H_a: \beta_1 \neq 0$$

$$\text{Test Statistic: } F = \frac{MSR}{MSE}$$

Decision rule:

Reject H_0 if $F > F_{\alpha/2}$, where $F_{\alpha/2}$ is based on $(1, n-2)$ degrees of freedom.

Error term: ϵ_i

The error terms are independent random

variables that are normally distributed with mean 0 and variance σ^2 .

Underlying Assumptions of the Model

1. Linearity: y and x are linearly related.

2. Normality: Residual errors are normally distributed.

3. Constant variance: Error variance is

constant across all values of the covariate x .

4. Independence: Errors associated with the observations are mutually independent.

Heteroscedasticity

• When the requirement of a constant variance is violated we have a condition of heteroscedasticity.

Examining the residual plots: Deviation of normality and heteroscedasticity.

The residual plot for us to focus on

• Variance is not constant.

• Fins out: the variance is increasing as a function of the covariate x .

• Fins in: the variance is a decreasing function of x .

Deviation of uncorrelatedness:

• Examining the normal probability plot of the residuals.

Cook's distance

• Measure of influence of a data point.

$D_i > 2$: High influence.

Remedial measures (Purifying factors).

Influential on dep. var. on indep. var.

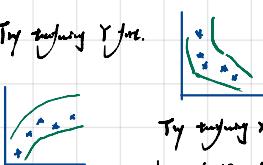
Square Root \sqrt{Y} \sqrt{X}

Log $\log(Y)$ $\log(X)$

Inverse $1/Y$ $1/X$



Solution of regression to remedy non-linearity



Try fitting X and w.r.t. but choose the variable with wider range first.

Try fitting X and w.r.t. first.

$$S \cdot (\hat{\beta}_1) = \frac{s}{\sqrt{n} s_{xx}}$$

$$\hat{\beta}_1 \pm t_{\alpha/2} \cdot S \cdot (\hat{\beta}_1)$$

Residual: $e_i = y_i - \hat{y}_i, i=1,2,\dots,n$

Error Sum of Squares (SSE):

$$\sum_{i=1}^n e_i^2$$

Estimation of σ^2 (Mean Squared Error)

$$\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{n-2}$$

Computational formula for $\hat{\sigma}^2$

$$\hat{\sigma}^2 = \frac{1}{n-2} \left[\sum_{i=1}^n y_i^2 - \hat{\beta}_0 \cdot \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i \right]$$

Inference about model parameters.

Inference about β_1

To test $H_0: \beta_1 = b$

$$\text{Test statistic: } t = \frac{\hat{\beta}_1 - b}{S \cdot (\hat{\beta}_1)}$$

where $S \cdot (\hat{\beta}_1)$ is the standard error of the estimate $\hat{\beta}_1$ and is given by:

$$S \cdot (\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\hat{\sigma}^2}{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}}$$

Decision Rule:

$H_0: \beta_1 = b$ vs. $H_a: \beta_1 > b$

Reject H_0 if $t > t_{\alpha}$

$H_0: \beta_1 = b$ vs. $H_a: \beta_1 < b$

Reject H_0 if $t < -t_{\alpha}$

$H_0: \beta_1 = b$ vs. $H_a: \beta_1 \neq b$

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

- t and $t_{\alpha/2}$ values are based on $(n-2)$ df.
- $H_0: \beta_1 = 0 \rightarrow$ No linear relationship between the dependent variable (y) and covariate (x).
- $H_a: \beta_1 \neq 0 \rightarrow$ Suggest the existence of a linear relationship.

Inference about β_0 : Confidence Interval:

Confidence interval for β_0 with confidence coefficient $1-\alpha$

Inference about β_0 :

- May not be meaningful unless data is available at or near $x=0$.

Confidence interval for β_0

$$\hat{\beta}_0 \pm t_{\alpha/2} \cdot S \cdot (\hat{\beta}_0)$$

where $S \cdot (\hat{\beta}_0)$ is the standard error of $\hat{\beta}_0$ given by:

$$\begin{aligned} S \cdot (\hat{\beta}_0) &= \frac{\hat{\sigma}^2 \sum_{i=1}^n x_i^2}{n \cdot \sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\hat{\sigma}^2 \sum_{i=1}^n x_i^2}{n \left[\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n} \right]} \end{aligned}$$

Prediction of an individual value of y at a given value of $x=x_p$.

Interval estimation of $E(Y)$ at $x=x_p$.

Punkt Estimate of $E(Y)$: $\hat{y}_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$

Confidence Interval: $\hat{y}_p \pm t_{\alpha/2} \cdot S \cdot (\hat{y}_p)$

where $S \cdot (\hat{y}_p)$ is the standard error of the estimate \hat{y}_p given by:

$$S \cdot (\hat{y}_p) = \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Simple: Construct a 95% confidence interval for $E(Y)$ at $x=5$. This is an estimate of the average annual salary of all people with 5 years of experience.

$$\hat{y}_p = -0.675 + 14.229 \cdot 5 = 70.47$$

$$S \cdot (\hat{y}_p) = 19.61 \cdot \sqrt{\frac{1}{12} + \frac{(5-5.5)^2}{59}} = 5.80$$

$$\Rightarrow 70.47 \pm 2.228 \cdot 5.80 \rightarrow (57.55, 83.39)$$

Prediction interval of Y at $x=x_p$

Best predictor of Y : $\hat{y}_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$

Prediction Interval: $\hat{y}_p \pm t_{\alpha/2} \cdot S \cdot (\hat{y}_p)$

Standard Error of \hat{y}_p :

$$S \cdot (\hat{y}_p) = \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$\hat{S} \cdot (\hat{y}_p) = \hat{\sigma} + S \cdot (\hat{y}_p)$$

Simple: Construct a 95% prediction interval of Y at $x=5$. The prediction interval provides a salary estimate of an individual with 5 years of experience.

$$\hat{y}_p = -0.675 + 14.229 \cdot 5 = 70.47$$

$$S \cdot (\hat{y}_p) = 19.61 \cdot \sqrt{\frac{1}{12} + \frac{(5-5.5)^2}{59}} = 20.45$$

$$\Rightarrow 70.47 \pm 2.228 \cdot 20.45 \rightarrow (24.91, 116.03)$$

Note: The prediction interval for Y at $x=5$ is significantly wider than the confidence interval for $E(Y)$ at $x=5$.

Calculator: Inverse T: invT(0.975, 6) $\rightarrow 2.4469$

Inverse F: invF(0.95, 1, 10) $\rightarrow 4.8646$

Inference about the dependent variable:

Two types of inference:

- Inference about the mean of y (expected value of y or $E(y)$) at a given value of $x = x_m$.