

11.4 - 11.5

With the notation $s_{xx} = \sum_{i=1}^m (x_{ri} - \bar{x}_i)^2$, we get

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^m x_{ri}^2}{m s_{xx}} \quad \text{and} \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{s_{xx}}$$

An estimator for σ^2 is given by

$$s^2 = \frac{\sum_{i=1}^m (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{m-2}$$

Also $T_{\beta_0} = \frac{\hat{\beta}_0 - \beta_0}{s \sqrt{\frac{\sum_{i=1}^m x_{ri}^2}{m s_{xx}}}}$ and $T_{\beta_1} = \frac{\hat{\beta}_1 - \beta_1}{s / \sqrt{s_{xx}}}$ are t_{m-2} under the

assumption of normally distributed data.

These can be used for hypothesis testing of the coefficients.

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

The p-value is $2P(T_{\beta_1} \geq |t_{\text{obs}}| \mid H_0)$ and can be compared to the significance level α .

$$\text{For } H_0: \beta_1 = 0 \quad H_1: \beta_1 > 0$$

the p-value is $P(T_{\beta_1} > t_{\text{obs}} \mid H_0)$

Similarly for tests on β_0

$$\text{From } P(-t_{\frac{\alpha}{2}, m-2} \leq \frac{\hat{\beta}_0 - \beta_0}{s \sqrt{\frac{\sum_{i=1}^m x_{ri}^2}{m s_{xx}}}} \leq t_{\frac{\alpha}{2}, m-2}) = 1 - \alpha$$

we get a $100(1-\alpha)\%$ confidence interval for β_0

$$\left(\hat{\beta}_0 - t_{\frac{\alpha}{2}, m-2} s \sqrt{\frac{\sum_{i=1}^m x_{ri}^2}{m s_{xx}}}, \hat{\beta}_0 + t_{\frac{\alpha}{2}, m-2} s \sqrt{\frac{\sum_{i=1}^m x_{ri}^2}{m s_{xx}}} \right)$$

Similarly for β_1

Confidence and Prediction intervals 11.6

A confidence interval for $E[Y|X_i=x_0] = \mu_{Y|X_0}$ can be obtained as follows.

Let $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_0$ be the estimator

$$\hat{Y}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_1 X_0 = \bar{Y} + \hat{\beta}_1 (X_0 - \bar{X}_1)$$

Assuming normally distributed data, \bar{Y} and $\hat{\beta}_1$ are independent.

Hence $\text{Var}(\hat{Y}_0) = \frac{\sigma^2}{m} + (X_0 - \bar{X}_1)^2 \cdot \frac{\sigma^2}{S_{xx}}$

With unknown variance $\frac{\hat{Y}_0 - \mu_{Y|X_0}}{\sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}}} \sim t_{m-2}$

and a $100(1-\alpha)\%$ confidence interval is given by

$$\left(\hat{Y}_0 - t_{\frac{\alpha}{2}, m-2} \sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}}, \hat{Y}_0 + t_{\frac{\alpha}{2}, m-2} \sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}} \right)$$

For a new observation, y_0 , when $X_i = x_0$ we may construct a prediction interval as follows.

$$E[\hat{Y}_0 - y_0] = \beta_0 + \beta_1 x_0 - (\beta_0 + \beta_1 x_0) = 0$$

$$\text{and } \text{Var}(\hat{Y}_0 - y_0) = \text{Var}(\hat{Y}_0) + \text{Var}(y_0) = \frac{\sigma^2}{m} + (X_0 - \bar{X}_1)^2 \cdot \frac{\sigma^2}{S_{xx}} + \sigma^2$$

Hence $\frac{\hat{Y}_0 - y_0}{\sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}}} \sim t_{m-2}$ and a $100(1-\alpha)\%$ prediction interval

$$\left(\hat{Y}_0 - t_{\frac{\alpha}{2}, m-2} \sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}}, \hat{Y}_0 + t_{\frac{\alpha}{2}, m-2} \sqrt{\frac{1}{m} + \frac{(X_0 - \bar{X}_1)^2}{S_{xx}}} \right)$$

Notice that the 1. normal equation gives

$$\sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i) = 0 \quad \text{The sum of residuals are zero}$$

Further we get from the normal equations in general that

$$X'(\underline{y} - \hat{\underline{y}}) = 0 \Rightarrow L'X'(\underline{y} - \hat{\underline{y}}) = 0 \Leftrightarrow \hat{\underline{y}}'(\underline{y} - \hat{\underline{y}}) = 0$$

$$\Leftrightarrow \sum_{i=1}^n \hat{y}_i (\underline{y}_i - \hat{\underline{y}}_i) = 0$$

Decomposing the variation

$$\underline{y}_i - \bar{\underline{y}} = \underline{y}_i - \hat{\underline{y}}_i + \hat{\underline{y}}_i - \bar{\underline{y}}$$

such that $\sum_{i=1}^n (\underline{y}_i - \bar{\underline{y}})^2 = \sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)^2 + 2 \sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)(\hat{\underline{y}}_i - \bar{\underline{y}}) + \sum_{i=1}^n (\hat{\underline{y}}_i - \bar{\underline{y}})^2$

$$\sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)(\hat{\underline{y}}_i - \bar{\underline{y}}) = \sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)\hat{y}_i - \sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)\bar{y}$$

$$= 0 \qquad \qquad 0$$

such that $\sum_{i=1}^n (\underline{y}_i - \bar{\underline{y}})^2 = \sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)^2 + \sum_{i=1}^n (\hat{\underline{y}}_i - \bar{\underline{y}})^2$

$$SS_T = SSE + SSR$$

Total sum of squares = error sum of squares + sum of squares for regression.

Theorem 12.1

In the linear regression model

$$\underline{y} = X\underline{\beta} + \underline{\varepsilon}$$

an unbiased estimator for σ^2 is given by

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)^2}{n - (k+1)} \quad \text{The estimate is } \underline{SSE} = \frac{\sum_{i=1}^n (\underline{y}_i - \hat{\underline{y}}_i)^2}{n - (k+1)}$$

The sum of squares, SST , SSE and SSR are normally gathered in an analysis of variance (ANOVA) table.

Source	SS	d.f.	F
Regression	$SSR = \sum_{i=1}^m (\hat{y}_i - \bar{y})^2$	k	$\frac{SSR}{k} / \frac{SSE}{m-k-1}$
Error	$SSE = \sum_{i=1}^m (y_i - \hat{y}_i)^2$	$m - k - 1$	
Total	$SST = \sum_{i=1}^n (y_i - \bar{y})^2$	$m - 1$	

Let $\epsilon_i \sim N(0, \sigma^2)$ and independent, $i = 1, 2, \dots, m$

The F statistics is used to test the hypothesis

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0 \quad H_1: \text{at least one } \beta_i \text{ is different from 0}$$

i.e. of the regression is significant

$$\text{Reject } H_0 \text{ if } F = \frac{\frac{SSR}{k}}{\frac{SSE}{m-k-1}} \geq f_{\alpha, k, m-k-1}$$

eventually if $P(F \geq F_{\alpha, k, m-k-1}) \leq \alpha$.

12.5 Inference in multiple linear regression

Test for the significant impact on the response given that the other variables are in the model

$$H_0: \beta_j = 0 \quad H_1: \beta_j \neq 0$$

$$\text{More general } H_0: \beta_j = \beta_{j0} \quad H_1: \beta_j \neq \beta_{j0}$$