

Inference about Regression Parameters

Under our multiple linear regression model,

$$\underline{b} = (X'X)^{-1}X'y$$

* are unbiased $E\{\underline{b}\} = \underline{\beta}$ (equivalent to

$$\begin{aligned} * \sigma^2\{\underline{b}\} &= \begin{pmatrix} \sigma^2\{b_0\} & \sigma\{b_0, b_1\} & \dots & \sigma\{b_0, b_{p-1}\} \\ \sigma\{b_0, b_1\} & \sigma^2\{b_1\} & \dots & \sigma\{b_1, b_{p-1}\} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma\{b_0, b_{p-1}\} & \sigma\{b_1, b_{p-1}\} & \dots & \sigma^2\{b_{p-1}\} \end{pmatrix} \\ &= \sigma^2 (X'X)^{-1} \end{aligned}$$

$$\begin{pmatrix} E\{b_0\} \\ E\{b_1\} \\ \vdots \\ E\{b_{p-1}\} \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix}$$

* can get an estimate of $\sigma^2\{\underline{b}\}$, $s^2\{\underline{b}\} = s^2 (X'X)^{-1}$
(MSE)

$$s^2\{\underline{b}\} = s^2(X'X)^{-1} = \begin{pmatrix} s^2\{b_0\} & s\{b_0, b_1\} & \dots & s\{b_0, b_{p-1}\} \\ s\{b_0, b_1\} & s^2\{b_1\} & \dots & s\{b_1, b_{p-1}\} \\ \vdots & \vdots & \ddots & \vdots \\ s\{b_0, b_{p-1}\} & s\{b_1, b_{p-1}\} & \dots & s^2\{b_{p-1}\} \end{pmatrix}$$

Can use $s\{\underline{b}\}$ for:

- * Interval estimation of β_k
- * tests for β_k
- * interval estimation for fitted values

* Interval estimation of β_k

Under our normal multiple linear regression model

$$\frac{b_k - \beta_k}{s\{b_k\}} \sim t_{n-p} \quad \text{for } k=0, 1, \dots, p-1$$

confidence limits for β_k with $1-\alpha$ confidence coefficient are $(b_k + t(\alpha/2, n-p)s\{b_k\}, b_k + t(1-\alpha/2, n-p)s\{b_k\})$

* Interval estimation of β_k

under our normal multiple linear regression model

$$\frac{b_k - \beta_k}{s\{b_k\}} \sim t_{n-p} \text{ for } k=0, 1, \dots, p-1$$

confidence limits for β_k with $1-\alpha$ confidence coefficient are

same as $(b_k + t(\alpha/2, n-p) s\{b_k\}, b_k + t(1-\alpha/2, n-p) s\{b_k\})$

$(b_k - t(1-\alpha/2, n-p) s\{b_k\}, b_k + t(1-\alpha/2, n-p) s\{b_k\})$

because t -distribution is symmetric about 0,

$$t(\alpha/2, n-p) = -t(1-\alpha/2, n-p)$$

* Tests for β_k

To test $H_0: \beta_k = 0$ versus $H_a: \beta_k \neq 0$

we use the test statistic $t^* = b_k / s\{b_k\}$

Our decision rule for a level- α test will be

* If $|t^*| \leq t(1-\alpha/2, n-p)$ then conclude H_0

* If $|t^*| > t(1-\alpha/2, n-p)$ then conclude H_a

We can also use a type of F-test to test H_0 versus H_a
but we're going to wait until chapter 7 for this!

* interval estimation for fitted values

Under our normal multiple linear regression model, given values of X_1, X_2, \dots, X_{p-1} for a new observation denoted by $X_{n1}, X_{n2}, \dots, X_{np-1}$,

$$E\{Y_n\} = \beta_0 + \sum_{k=1}^{p-1} \beta_k X_{nk}$$

↑
average
response
of observation
n

$$= \tilde{X}_n' \tilde{\beta}$$

where

$$\tilde{X}_n = \begin{pmatrix} 1 \\ X_{n1} \\ X_{n2} \\ \vdots \\ X_{np-1} \end{pmatrix}$$

$$\hat{Y}_n = \tilde{X}_n' \tilde{b} = b_0 + \sum_{k=1}^{p-1} b_k X_{nk}$$

↑
estimated
mean response

$$E\{\hat{Y}_n\} = E[Y_n] = \tilde{X}_n' \tilde{\beta}$$

$$\sigma^2\{\hat{Y}_n\} = \sigma^2 \tilde{X}_n' (X'X)^{-1} \tilde{X}_n$$

$$\sigma^2\{\hat{Y}_n\} = \sigma^2 \tilde{X}_n' (X'X)^{-1} \tilde{X}_n$$

$$= \tilde{X}_n' \sigma^2 \tilde{b} \tilde{X}_n$$

The 1- α confidence limits for $E\{Y_n\}$ are $(\hat{Y}_n \pm t_{(\alpha/2, n-p)} s\{\hat{Y}_n\}, \hat{Y}_n \pm t_{(1-\alpha/2, n-p)} s\{\hat{Y}_n\})$

Interval Estimation for future values of the response

$$y_h = E\{y_h\} + \varepsilon_h$$

$$\sigma^2\{\text{pred}\} = \sigma^2 + \sigma^2\{\hat{y}_h\}$$

$$s^2\{\text{pred}\} = s^2 + s^2\{\hat{y}_h\} = s^2 \left(1 + \underline{x}_h' (X'X)^{-1} \underline{x}_h \right)$$

The $1-\alpha$ prediction limits for a new observation

$y_{h(\text{new})}$ corresponding to \underline{x}_h are

$$\left(\hat{y}_h + t(\alpha/2, n-p) s\{\text{pred}\}, \hat{y}_h + t(1-\alpha/2, n-p) s\{\text{pred}\} \right)$$

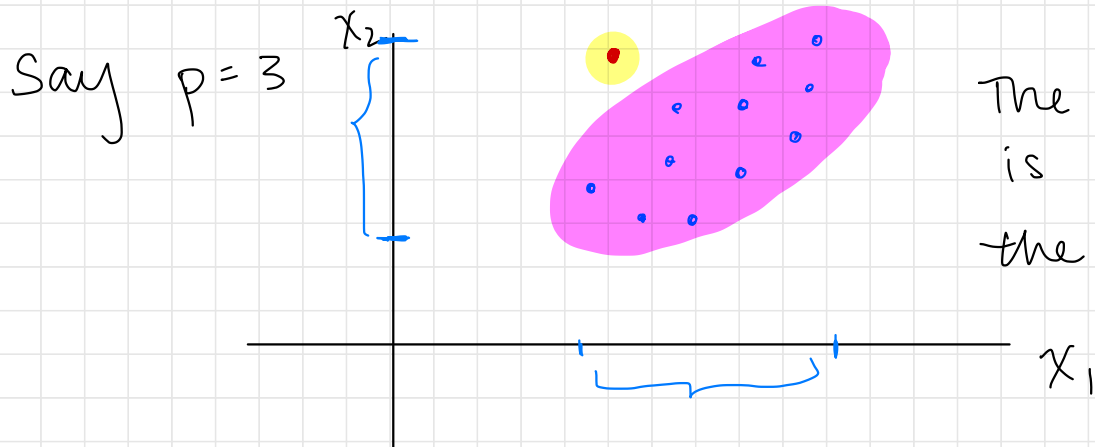
Interval estimation for average of m future values of the response? $\sigma^2\{\text{pred}_{\text{mean}}\} = \frac{\sigma^2}{m} + \sigma^2\{\hat{y}_h\}$

Extrapolation / Defining Scope in Multiple Linear Regression

Say $p=2 \Rightarrow$ simple linear regression

data contains $\tilde{X} = \begin{pmatrix} 11 \\ 21 \\ 18 \\ 32 \\ 13 \end{pmatrix}$

Is $x_n = 0$ within the scope?



The red point is not within the scope of the model

Diagnostics / "Remedial Measures"

Transformations we may need to take and/or additional covariates we may need to include to improve the plausibility of multiple linear regression assumptions

* scatter plots of response and covariates

- Is linearity reasonable?
- Is constant variance reasonable?
- Is our model complete/are the relevant predictors/covariates included?

* scatter plots of residuals from a candidate model and fitted values, observed response values, covariates

Example - Duane Photo Studios Example

y_i : sales in city i

x_{i1} : proportion of population under 17 in city i

x_{i2} : average per capita disposable income
in city i

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$