

The goal is to find values of the β 's that minimize the magnitude of $\vec{\epsilon}$.

$$\|\vec{\epsilon}\|_2 = \sqrt{\vec{\epsilon}^T \vec{\epsilon}} = \sqrt{\sum_{i=1}^n \epsilon_i^2} = \sqrt{S(\beta_0, \beta_1, \dots, \beta_p)}$$

L_2 -norm

Because the square root function increases monotonically then the values of the β 's that minimize $\|\vec{\epsilon}\|_2$ are the same ones that minimize $S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n \epsilon_i^2$. Thus, what we've worked out geometrically, is equivalent to least squares minimization.

The vector in $L(x)$ that is closest to \vec{y} is obtained by making the difference $\vec{y} - \vec{\mu} > \vec{\epsilon}$ perpendicular (orthogonal) to $L(x)$. Thus we require $\vec{y} - \vec{\mu} = \vec{\epsilon}$ to be orthogonal to every vector in $L(x)$. In particular, to $\vec{I}, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_p$. This gives:

$$\begin{aligned} \vec{I}^T (\vec{y} - \vec{\mu}) &= 0 \\ \vec{x}_1^T (\vec{y} - \vec{\mu}) &= 0 \\ \vec{x}_2^T (\vec{y} - \vec{\mu}) &= 0 \\ &\vdots \\ \vec{x}_p^T (\vec{y} - \vec{\mu}) &= 0 \end{aligned}$$

As before, this system of equations can be written as:

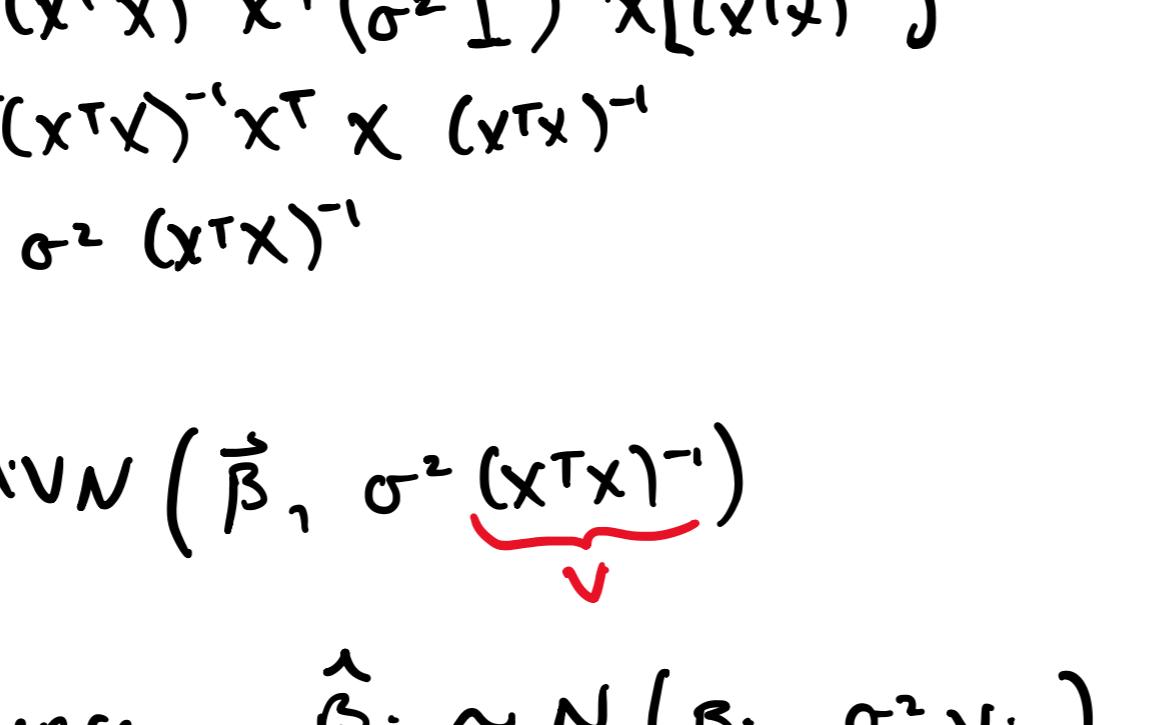
$$X^T (\vec{y} - \vec{\mu}) = \vec{0}$$

As before, solving this equation yields $\hat{\beta} = (X^T X)^{-1} X^T \vec{y}$.

Related Quantities

We call the vector $\hat{\mu} = X\hat{\beta}$ the fitted values.

We call the vector $\vec{\epsilon} = \vec{y} - \hat{\mu} = \vec{y} - X\hat{\beta}$ the residuals.



* The fitted value vector is the orthogonal projection of \vec{y} onto $L(x)$.

The residuals, as in the case of SLR, are used to estimate σ^2 , the error variance:

$$\begin{aligned} \hat{\sigma}^2 &= \frac{\sum_{i=1}^n \epsilon_i^2}{n-p-1} = \frac{\vec{\epsilon}^T \vec{\epsilon}}{n-p-1} \quad \text{"residual degrees of freedom"} \\ &\quad = (\text{sample size}) - (\text{# of estimated parameters}) \\ \hat{\sigma} &= \sqrt{\frac{\sum_{i=1}^n \epsilon_i^2}{n-p-1}} = \sqrt{\frac{\vec{\epsilon}^T \vec{\epsilon}}{n-p-1}} \\ &\quad = n - (p+1) \\ &\quad = n - p - 1 \end{aligned}$$

Inference for $\hat{\beta} = (X^T X)^{-1} X^T \vec{y}$

Since $\vec{y} \sim MVN(\vec{\mu} = X\vec{\beta}, \sigma^2 I_{n \times n})$ and because $\hat{\beta}$ is a matrix multiple of \vec{y} , $\hat{\beta}$ also follows a multivariate normal distribution.

$$\begin{aligned} E[\hat{\beta}] &= E[(X^T X)^{-1} X^T \vec{y}] = (X^T X)^{-1} X^T E[\vec{y}] \\ &= (X^T X)^{-1} X^T X \vec{\beta} \\ &= \vec{\beta} \end{aligned}$$

$$\therefore E \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\begin{aligned} \text{Var}[\hat{\beta}] &= \text{Var}[(X^T X)^{-1} X^T \vec{y}] \\ &= (X^T X)^{-1} X^T \text{Var}[\vec{y}] [X(X^T X)^{-1} X^T]^T \\ &= (X^T X)^{-1} X^T (\sigma^2 I) X [X(X^T X)^{-1}]^T \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

$$\therefore \hat{\beta} \sim MVN(\vec{\beta}, \sigma^2 (X^T X)^{-1})$$

As a consequence, $\hat{\beta}_j \sim N(\beta_j, \sigma^2 v_{jj}) \quad j = 0, 1, 2, \dots, p$

↑
jth diagonal element of V

$$\therefore \frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{v_{jj}}} \sim N(0, 1)$$

$$\Rightarrow \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{v_{jj}}} \sim t_{(n-p-1)} \quad (\text{since } \frac{(n-p-1)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{(n-p-1)})$$

We use this result to calculate test statistics and confidence intervals

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

$$\text{test statistic: } t = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_{jj}}} \quad \text{SE}(\hat{\beta}_j)$$

$$\text{p-value} = 2P(t_{(n-p-1)} \geq |t|)$$

A $(1-\alpha) \times 100\%$ CI for β_j is

$$\hat{\beta}_j \pm t_{(n-p-1), (1-\alpha/2)} \times \hat{\sigma} \sqrt{v_{jj}} \quad \text{SE}(\hat{\beta}_j)$$

$$\text{where } P(W \leq t_{(n-p-1), (1-\alpha/2)}) = 1 - \frac{\alpha}{2} \quad \text{where } W \sim t_{(n-p-1)}$$

Interpreting the β 's

$$E[y] = \mu = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

Notice that $E[y | x_1 = x_2 = x_3 = \dots = x_p = 0] = \beta_0$. Therefore, β_0 is interpreted as the expected response when all of the explanatory variables are equal to zero. And then, we interpret $\hat{\beta}_0$ as the estimated expected response when all explanatory variables are equal to zero.

Next, notice $E[y | x_j = x] = \beta_0 + \beta_1 x_1 + \dots + \beta_j x + \dots + \beta_p x_p$

$$E[y | x_j = x+1] = \beta_0 + \beta_1 x_1 + \dots + \beta_j(x+1) + \dots + \beta_p x_p$$

So $E[y | x_j = x+1] - E[y | x_j = x] = \beta_j$. Therefore β_j is interpreted as the expected change in the response corresponding to a unit increase in x_j , when all other explanatory variables stay the same. And then, we interpret $\hat{\beta}_j$ as the estimate of the expected change in y for a unit increase in x_j , all else held constant.

this is new