

Regression Analysis Tutoring1

Seung Bong Jung

Seoul National University

November 1, 2021

Linear Regression Model

Linear Regression Model

The model: $Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$, where Y is a real-valued response, (x_1, \dots, x_p) is a set of p predictors assumed to be nonrandom, the error ϵ is a random variable with mean $E(\epsilon) = 0$ and finite variance $\sigma^2 = \text{Var}(\epsilon)$

- Assume that we observe the responses Y_i at the preselected predictor values x_{i1}, \dots, x_{ip} , respectively, for $1 \leq i \leq n$, such that

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \text{ where } \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- Then the parameter is $\theta = (\beta, \sigma^2)$ and the likelihood is given by

$$L(\theta) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2\right)$$

where $\beta = (\beta_0, \beta_1, \dots, \beta_p)^t$

Linear Regression Model

- Thus, if we are to get MLE of β , we have to minimize $\|\mathbf{Y} - \mathbf{X}\beta\|^2$ with respect to β , where $\mathbf{Y} = (Y_1, \dots, Y_n)^t$, $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^t$, and $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p)$. We call such minimizer as least squared estimator denoted by $\hat{\beta}$.

Simple Linear Regression

Consider the case when $p = 1$.

- $(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$.
- $\hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 \mathbf{x}$ is the projection of $\mathbf{Y} = (Y_1, \dots, Y_n)^t$ onto the linear subspace in \mathbb{R}^n spanned by $\mathbf{1} = (1, \dots, 1)^t$ and $\mathbf{x} = (x_1, \dots, x_n)^t$.

$$\hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 \mathbf{x} = \Pi(\mathbf{Y} | \mathcal{C}_{1, \mathbf{x}})$$

, where $\mathcal{C}_{1, \mathbf{x}} = \{\beta_0 \mathbf{1} + \beta_1 \mathbf{x} : \beta_0, \beta_1 \in \mathbb{R}\}$

- Orthogonal decomposition of the column space:

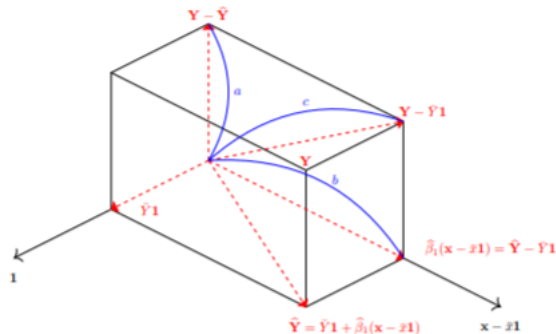
$\mathcal{C}_{1, \mathbf{x}} = \mathcal{C}_{1, \mathbf{x} - \bar{x}\mathbf{1}} = \mathcal{C}_1 \oplus \mathcal{C}_{\mathbf{x} - \bar{x}\mathbf{1}}$ so that

$$\hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 \mathbf{x} = (\hat{\beta}_0 + \hat{\beta}_1 \bar{x}) \mathbf{1} + \hat{\beta}_1 (\mathbf{x} - \bar{x} \mathbf{1})$$

$$\Pi(\mathbf{Y} | \mathcal{C}_{1, \mathbf{x}}) = \Pi(\mathbf{Y} | \mathcal{C}_1) + \Pi(\mathbf{Y} | \mathcal{C}_{\mathbf{x} - \bar{x}\mathbf{1}})$$

$$= \bar{Y} \mathbf{1} + \frac{S_{xy}}{S_{xx}} (\mathbf{x} - \bar{x} \mathbf{1})$$

Decomposition of Sums of Squares



- $$c^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 (\text{SST}), a^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 (\text{SSE})$$

$$b^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 (\text{SSR})$$

Hypothesis Testing for the Slope

- Due to the projection interpretation, for the residual vector $\mathbf{e} = (e_1, \dots, e_n)^t$,

$$\sum_{i=1}^n e_i = \sum_{i=1}^n e_i x_i = \sum_{i=1}^n e_i \hat{Y}_i = 0$$

- Pythagorean theorem: $c^2 = a^2 + b^2$
- Suppose we want to test

$$H_0 : \beta_1 = 0 \text{ versus } H_1 : \beta_1 \neq 0$$

- t-statistic:

$$T \equiv \frac{\frac{\hat{\beta}_1 - 0}{\sqrt{\sigma^2 / S_{xx}}}}{\sqrt{\frac{SSE}{(n-2)\sigma^2}}} \sim t(n-2)$$

- Reject H_0 if $|T| > t_{\alpha/2}(n-2)$.

Hypothesis Testing for the Slope

- F-statistic:

$$F \equiv \frac{\frac{SSR}{1\sigma^2}}{\frac{SSE}{(n-2)\sigma^2}} \sim F(1, n-2)$$

- Reject H_0 if $F > F_\alpha(1, n-2)$.

- Regression coefficients:

$$\beta_1 : \hat{\beta}_1 \pm t_{\alpha/2}(n-2) \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}; \beta_0 : \hat{\beta}_0 \pm t_{\alpha/2}(n-2) \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$$

- Mean Response:

$$\hat{\mu}_x = \hat{\beta}_0 + \hat{\beta}_1 x; \quad \mu_x : \hat{\mu}_x \pm t_{\alpha/2}(n-2) \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}} \right)}$$

- Out-of-sample prediction of a response :

$$\hat{Y}_x = \hat{\beta}_0 + \hat{\beta}_1 x; \quad Y_x : \hat{Y}_x \pm t_{\alpha/2}(n-2) \sqrt{\hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}} \right)}$$

Multiple Linear Regression Models

Consider the case when $p \geq 2$.

- The model: $Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$, $\epsilon \sim N(0, \sigma^2)$
- $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2 = \arg \min_{\beta \in \mathbb{R}^{p+1}} \|\mathbf{Y} - \mathbf{X}\beta\|^2$
where $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^t$, $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p)$ and $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^t$.
- Denote the column space of \mathbf{X} by $\mathcal{C}_{\mathbf{X}}$ (the linear space in \mathbb{R}^n spanned by the columns of the matrix \mathbf{X}).
- Then $\mathbf{X}\hat{\beta} = \Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}})$
- Least squared estimator: $\hat{\beta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}$
- Note $\hat{\beta}$ is also BLUE in multivariate case.

Least Squares Estimation

- Orthogonal decomposition of the column space: Write $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1)$. Then,

$$\mathcal{C}_{\mathbf{X}} = \mathcal{C}_{\mathbf{1}, \mathbf{x}_1 - \Pi_1 \mathbf{x}_1} = \mathcal{C}_{\mathbf{1}} \oplus \mathcal{C}_{\mathbf{x}_1 - \Pi_1 \mathbf{x}_1}$$

where $\Pi_1 = \mathbf{1}(\mathbf{1}^t \mathbf{1})^{-1} \mathbf{1}^t$.

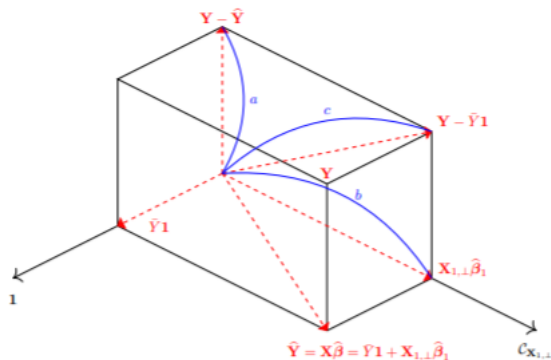
- Writing $(\hat{\beta}_0, \hat{\beta}_1^t)^t = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)^t = \hat{\beta}^t$ and $\mathbf{X}_{1,\perp} = \mathbf{X}_1 - \Pi_1 \mathbf{X}_1$,

$$\begin{aligned}\mathbf{X}\hat{\beta} &= \hat{\beta}_0 \mathbf{1} + \mathbf{X}_1 \hat{\beta}_1 \\ &= (\hat{\beta}_0 + n^{-1} \mathbf{1}^t \mathbf{X}_1 \hat{\beta}_1) \mathbf{1} + \mathbf{X}_{1,\perp} \hat{\beta}_1\end{aligned}$$

- This gives

$$\hat{\beta}_0 = \bar{Y} - n^{-1} \mathbf{1}^t \mathbf{X}_1 \hat{\beta}_1, \quad \hat{\beta}_1 = (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1} \mathbf{X}_{1,\perp}^t \mathbf{Y}$$

Decomposition of Sum of Squares



- $c^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 (\text{SST}), a^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 (\text{SSE})$
 $b^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 (\text{SSR})$

Testing for Significance of Regression

- Due to the projection interpretation, the residual vector \mathbf{e} is orthogonal to $\mathbf{1}$, all columns of \mathbf{X}_1 and $\hat{\mathbf{Y}}$, for all $1 \leq j \leq p$,

$$\sum_{i=1}^n e_i = \sum_{i=1}^n e_i x_{ij} = \sum_{i=1}^n e_i \hat{Y}_i = 0$$

- Pythagorean theorem: $c^2 = a^2 + b^2$ (SST=SSE+SSR)
- Want to test $H_0 : \beta_1 = \mathbf{0}$ versus $H_1 : \beta_1 \neq \mathbf{0}$
- Distribution of $\hat{\beta}_1$: $\hat{\beta}_1 \sim N_p(\beta_1, \sigma^2(\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1})$, so that $(\hat{\beta}_1 - \beta_1)^t (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp}) (\hat{\beta}_1 - \beta_1) \sim \sigma^2 \chi^2(p)$
- $\text{SSE} \sim \sigma^2 \chi^2(n - p - 1)$ and independent with $\hat{\beta}_1$.
- Under H_0 , we have $\hat{\beta}_1^t (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp}) \hat{\beta}_1 \sim \sigma^2 \chi^2(p)$

Testing for Significance of Regression

- F-statistic: $F \equiv \frac{\frac{SSR}{p\sigma^2}}{\frac{SSE}{(n-p-1)\sigma^2}} \sim F(p, n-p-1).$
- Reject H_0 if $F > F_\alpha(p, n-p-1)$

Testing for Individual β_j

- Want to test $H_0 : \beta_j = 0$ versus $H_1 : \beta_j \neq 0$, $1 \leq j \leq p$
- Distribution of $\hat{\beta}_j$: $\hat{\beta}_j \sim N(\beta_j, \sigma^2 \mathbf{1}_j^t (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1} \mathbf{1}_j)$, where $\mathbf{1}_j$ is the unit vector of dimension p with its j th entry being equal to one.
- It follows that $\mathbf{1}_j^t (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1} \mathbf{1}_j = \mathbf{1}_{j+1}^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{1}_{j+1}$
- Thus, writing $C_{jj} = \mathbf{1}_j^t (\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1} \mathbf{1}_j$, we have

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 C_{jj}}} \sim N(0, 1)$$

- Letting, $D_{kk} = \mathbf{1}_k^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{1}_k$, we get $C_{jj} = D_{j+1,j+1}$, $1 \leq j \leq p$ and $D_{11} = \frac{1}{n} + (\frac{\mathbf{1}^t \mathbf{X}_1}{j})(\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1}(\frac{\mathbf{X}_1^t \mathbf{1}}{n})$

Testing for Individual β_j

- Under H_0 ,

$$T \equiv \frac{\frac{\hat{\beta}_j}{\sqrt{\sigma^2 C_{jj}}}}{\sqrt{\frac{\text{SSE}}{(n-p-1)\sigma^2}}} \sim t(n-p-1)$$

- Reject H_0 if $|T| > t_{\alpha/2}(n-p-1)$
- For $j = 0$, under H_0 ,

$$T \equiv \frac{\frac{\hat{\beta}_0}{\sqrt{\sigma^2 D_{11}}}}{\sqrt{\frac{\text{SSE}}{(n-p-1)\sigma^2}}} \sim t(n-p-1)$$

- Similarly, reject H_0 if $|T| > t_{\alpha/2}(n-p-1)$

- Regression coefficients:

$$\beta_j : \hat{\beta}_j \pm t_{\alpha/2}(n - p - 1)\sqrt{\hat{\sigma}^2 C_{jj}}, \quad 1 \leq j \leq p$$

$$\beta_0 : \hat{\beta}_0 \pm t_{\alpha/2}(n - p - 1)\sqrt{\hat{\sigma}^2 D_{11}}$$

- Mean response: Writing

$$C_{\mathbf{z}} = \frac{1}{n} + (\mathbf{z}^t - \frac{\mathbf{1}^t \mathbf{X}_1}{n})(\mathbf{X}_{1,\perp}^t \mathbf{X}_{1,\perp})^{-1}(\mathbf{z}^t - \frac{\mathbf{X}_1^t \mathbf{1}}{n})$$

we get

$$\mu_{\mathbf{z}} : \quad \hat{\mu}_{\mathbf{z}} \pm t_{\alpha/2}(n - p - 1)\sqrt{\hat{\sigma}^2 C_{\mathbf{z}}}$$

- Out-of-sample prediction of a response:

$$Y_{\mathbf{z}} : \quad \hat{Y}_{\mathbf{z}} \pm t_{\alpha/2}(n - p - 1)\sqrt{\hat{\sigma}^2(1 + C_{\mathbf{z}})}$$

Estimator of Sub-vector of Regression Coefficients

Write $\beta^t = (\beta_1^t, \beta_2^t)$ and $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$, so that $\mathbf{X}\beta = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2$.

- Orthogonal decomposition of the column space:

$$\mathcal{C}_{\mathbf{X}} = \mathcal{C}_{\mathbf{X}_1, \mathbf{X}_2 - \Pi_1 \mathbf{X}_2} = \mathcal{C}_{\mathbf{X}_1} \oplus \mathcal{C}_{\mathbf{X}_2 - \Pi_1 \mathbf{X}_2}$$

where $\Pi_1 = \mathbf{X}_1(\mathbf{X}_1^t \mathbf{X}_1)^{-1} \mathbf{X}_1^t$.

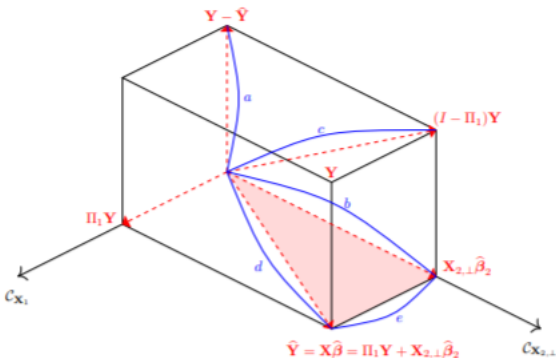
- Writing $\mathbf{X}_{2,\perp} = \mathbf{X}_2 - \Pi_1 \mathbf{X}_2$, we get

$$\begin{aligned}\mathbf{X}\hat{\beta} &= \mathbf{X}_1[\hat{\beta}_1 + (\mathbf{X}_1^t \mathbf{X}_1)^{-1} \mathbf{X}_1^t \mathbf{X}_2 \hat{\beta}_2] + \mathbf{X}_{2,\perp} \hat{\beta}_2 \\ \Pi(\mathbf{Y} | \mathcal{C}_{\mathbf{X}}) &= \Pi_1 \mathbf{Y} + \mathbf{X}_{2,\perp} (\mathbf{X}_{2,\perp}^t \mathbf{X}_{2,\perp})^{-1} \mathbf{X}_{2,\perp}^t \mathbf{Y}\end{aligned}$$

- This gives

$$\hat{\beta}_1 = (\mathbf{X}_1^t \mathbf{X}_1)^{-1} \mathbf{X}_1^t (\mathbf{Y} - \mathbf{X}_2 \hat{\beta}_2), \hat{\beta}_2 = (\mathbf{X}_{2,\perp}^t \mathbf{X}_{2,\perp})^{-1} \mathbf{X}_{2,\perp}^t \mathbf{Y}.$$

Decomposition of Sums of Squares



- Pythagorean theorem: $c^2 = a^2 + b^2$ and $d^2 = b^2 + e^2$

Testing for subsets of Regression Coefficients

- Extra sum of squares due to β_2 :

$$R(\beta_2|\beta_1) = R(\beta_1, \beta_2) - R(\beta_1)$$

where $R(\beta_1)$ and $R(\beta_1, \beta_2)$, respectively, denote the squared norms of $\Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}_1})$ and $\Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}_1, \mathbf{X}_2}) = \Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}})$.

- $R(\beta_2|\beta_1)$ tends to get larger if $\beta_2 \neq 0$ is true.



$$\begin{aligned} R(\beta_2|\beta_1) &= d^2 - e^2 = c^2 - a^2 = b^2 \\ &= \text{SSE}(\beta_1) - \text{SSE}(\beta_1, \beta_2) = b^2 = \|\mathbf{X}_{2,\perp} \hat{\beta}_2\|^2 \end{aligned}$$

where $\text{SSE}(\beta_1) = \|\mathbf{Y} - \Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}_1})\|^2$,
 $\text{SSE}(\beta_1, \beta_2) = \|\mathbf{Y} - \Pi(\mathbf{Y}|\mathcal{C}_{\mathbf{X}})\|^2$

Testing for subsets of Regression Coefficients

Set $\beta_2 \in \mathbb{R}^q$ and assume that $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$. We want test

$$H_0 : \beta_2 = \mathbf{0} \text{ versus } H_1 : \beta_2 \neq \mathbf{0}$$

- Distribution of $\hat{\beta}_2$: $(\mathbf{X}_{2,\perp}^t \mathbf{X}_{2,\perp})^{\frac{1}{2}}(\hat{\beta}_2 - \beta_2) \sim N_q(\mathbf{0}_q, \sigma^2 \mathbf{I}_q)$, so that $(\hat{\beta}_2 - \beta_2)^t (\mathbf{X}_{2,\perp}^t \mathbf{X}_{2,\perp}) (\hat{\beta}_2 - \beta_2) \sim \sigma^2 \chi^2(q)$
- Distribution of SSE: $\text{SSE} \sim \sigma^2 \chi^2(n - p - 1)$
- SSE and $(\hat{\beta}_1, \hat{\beta}_2)$ are independent.
- Under H_0 ,

$$\hat{\beta}_2^t (\mathbf{X}_{2,\perp}^t \mathbf{X}_{2,\perp}) \hat{\beta}_2 = R(\beta_2 | \beta_1) \sim \sigma^2 \chi^2(q)$$

Testing for subsets of Regression Coefficients

- F-statistic:

$$F \equiv \frac{\frac{R(\beta_2|\beta_1)}{q\sigma^2}}{\frac{\text{SSE}}{(n-p-1)\sigma^2}} \sim F(q, n-p-1)$$

- Reject H_0 if $F > F_\alpha(q, n-p-1)$

General Linear Hypothesis Tests

Let \mathbf{A} be a $(p+1) \times q$ matrix of full column rank and \mathbf{c} be q -dimensional vector ($q < p+1$). We want to test

$$H_0 : \mathbf{A}^t \boldsymbol{\beta} = \mathbf{c} \text{ versus } H_1 : \mathbf{A}^t \boldsymbol{\beta} \neq \mathbf{c}$$

- $\hat{\boldsymbol{\beta}}_r = \arg \min_{\mathbf{A}^t \boldsymbol{\beta} = \mathbf{c}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$

$$\begin{aligned} \hat{\boldsymbol{\beta}}_r &= (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y} \\ &\quad - (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{A} [\mathbf{A}^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{A}]^{-1} (\mathbf{A}^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y} - \mathbf{c}) \end{aligned}$$

- Let $R(\boldsymbol{\beta}) = \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2$ and $R(\boldsymbol{\beta}_r) = \|\mathbf{X}\hat{\boldsymbol{\beta}}_r\|^2$. Then $R(\boldsymbol{\beta}) - R(\boldsymbol{\beta}_r) \sim \sigma^2 \chi^2(q)$.
- F-statistic:

$$F \equiv \frac{\frac{R(\boldsymbol{\beta}) - R(\boldsymbol{\beta}_r)}{q\sigma^2}}{\frac{\text{SSE}}{(n-p-1)\sigma^2}} \sim F(q, n-p-1)$$

- Reject H_0 if $F > F_\alpha(q, n-p-1)$.