

A Quick Return to Vectors and Matrices

In the simple linear regression setting:

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Vectors: $\underset{\substack{\sim \\ \text{length } n}}{a} = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$

$$\underset{\sim}{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$$

$$\underset{\sim}{a}' \underset{\sim}{b} = (a_1 \dots a_n) \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = \sum_{i=1}^n a_i b_i$$

Matrices: $C_{n \times 2} = \begin{pmatrix} \underset{\sim}{C_1} & \underset{\sim}{C_2} \\ \vdots & \vdots \\ C_{n1} & C_{n2} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \\ \vdots & \vdots \\ C_{n1} & C_{n2} \end{pmatrix}$

$$\underset{2 \times n}{C} \underset{n \times 1}{b} = \begin{pmatrix} C_{11} \dots C_{1n} \\ C_{21} \dots C_{2n} \\ \vdots \\ C_{n1} \dots C_{n2} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} \underset{\sim}{C_1}' b \\ \underset{\sim}{C_2}' b \\ \vdots \\ C_{n1}' b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n C_{1i} b_i \\ \sum_{i=1}^n C_{2i} b_i \\ \vdots \\ \sum_{i=1}^n C_{ni} b_i \end{pmatrix}$$

How Vectors and Matrices Make Multiple Linear Regression Easier...

$$\underset{\sim}{b}_{p \times 1} = (\underset{\sim}{X}'\underset{\sim}{X})^{-1} \underset{\sim}{X}'\underset{\sim}{y}$$

$$\underset{\sim}{X}_{n \times p} = \begin{pmatrix} \underset{\sim}{x}_0 & \dots & \underset{\sim}{x}_{p-1} \end{pmatrix}$$

$x_{i0} = 1$ for all i

$$\bar{x} = \underbrace{(1 \dots 1)}_{\underset{\sim}{1}_n'} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \times \frac{1}{n}$$

$$\bar{x} = (\underset{\sim}{1}_n' \underset{\sim}{1}_n)^{-1} \underset{\sim}{1}_n' \underset{\sim}{x}$$

$$\underset{\sim}{X}'\underset{\sim}{X} = \begin{pmatrix} \underset{\sim}{x}_0' \underset{\sim}{x}_0 & \underset{\sim}{x}_0' \underset{\sim}{x}_1 & \dots & \underset{\sim}{x}_0' \underset{\sim}{x}_{p-1} \\ \underset{\sim}{x}_1' \underset{\sim}{x}_0 & \underset{\sim}{x}_1' \underset{\sim}{x}_1 & \dots & \underset{\sim}{x}_1' \underset{\sim}{x}_{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ \underset{\sim}{x}_{p-1}' \underset{\sim}{x}_0 & \dots & \underset{\sim}{x}_{p-1}' \underset{\sim}{x}_{p-1} \end{pmatrix}$$

Simple Linear Regression in Matrix and Vector Notation

$$X = \begin{pmatrix} 1_n & \underline{\underline{x}} \end{pmatrix}_{n \times 2} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad b = (X'X)^{-1} X'y$$

$$X'X = \begin{pmatrix} 1_n' 1_n & 1_n' \underline{\underline{x}} \\ 1_n' \underline{\underline{x}} & \underline{\underline{x}}' \underline{\underline{x}} \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}$$

$$(X'X)^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)(\sum x_i)} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix}$$

Simple Linear Regression in Matrix and Vector Notation

$$b = (x'x)^{-1} x'y$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)(\sum x_i)} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)(\sum x_i)}$$

$$\begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)(\sum x_i)}$$

$$\begin{pmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum y_i x_i) \\ (-\sum x_i)(\sum y_i) + n \sum x_i y_i \end{pmatrix}$$

$$\hat{b} = \frac{1}{n \sum x_i^2 - (\sum x_i)(\sum x_i)} \begin{pmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum y_i x_i) \\ (-\sum x_i)(\sum y_i) + n \sum x_i y_i \end{pmatrix}$$

Easy to see the relationship to our previous definition
 if $\bar{x} = 0 \Rightarrow \sum x_i = 0$

$$\hat{b} = \frac{1}{n \sum x_i^2} \begin{pmatrix} \sum x_i^2 \sum y_i \\ n \sum x_i y_i \end{pmatrix} = \begin{pmatrix} \bar{y} \\ \frac{\sum x_i y_i}{\sum x_i^2} \end{pmatrix}$$

This looks like our earlier definition!

Multiple Linear Regression Model with Normal Errors

$$\underset{\substack{n \times 1 \\ \text{response} \\ \text{vector}}}{\mathbf{y}} = \underset{n \times p \text{ covariate matrix}}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{p \times 1}{\boldsymbol{\varepsilon}}$$

first column is a vector of 1's
equivalent to $y_i = \beta_0 + \sum_{k=1}^{p-1} \beta_k X_{ik} + \varepsilon_i$

} Includes $p-1$ covariates

* $E\{\underset{\sim}{\boldsymbol{\varepsilon}}\} = \underset{\sim}{0}$ equivalent to $E\{\varepsilon_i\} = 0$

* $\sigma^2\{\underset{\sim}{\boldsymbol{\varepsilon}}\} = \sigma^2 \mathbf{I}_n$

* $\varepsilon_i \sim N(0, \sigma^2)$

equivalent to $\sigma^2\{\varepsilon_i\} = \sigma^2, \sigma\{\varepsilon_i, \varepsilon_j\} = 0, i \neq j$

$$\Rightarrow E\{\underset{\sim}{\mathbf{y}}\} = \mathbf{X}\boldsymbol{\beta}, \quad \sigma^2\{\underset{\sim}{\mathbf{y}}\} = \sigma^2 \mathbf{I}_n$$

Coefficient of Multiple Determination
 $R^2 = SSR / SSTO = 1 - SSE / SSTO$

We have estimators:

* $\underset{\sim}{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\underset{\sim}{\mathbf{y}}$ * $\underset{\sim}{\hat{\mathbf{y}}} = \mathbf{X}\underset{\sim}{\mathbf{b}}$

* What proportion of the variability in $\underset{\sim}{\mathbf{y}}$ is explained by the covariate

* R^2 is nondecreasing as you add more $\frac{\underset{\sim}{\mathbf{e}}'\underset{\sim}{\mathbf{e}}}{n-p}$
* $\underset{\sim}{\mathbf{e}} = \underset{\sim}{\mathbf{y}} - \underset{\sim}{\hat{\mathbf{y}}}$ * $S^2 = \frac{\underset{\sim}{\mathbf{e}}'\underset{\sim}{\mathbf{e}}}{n-p}$

ANOVA sums-of-squares

$SSE = \underset{\sim}{\mathbf{e}}'\underset{\sim}{\mathbf{e}}$

has $n-p$ degrees of freedom

$SSTO = \sum_{i=1}^n (y_i - \bar{y})^2 = (\underset{\sim}{\mathbf{y}} - \frac{1}{n} \underset{\sim}{\mathbf{1}} \underset{\sim}{\mathbf{1}}' \underset{\sim}{\mathbf{y}})' (\underset{\sim}{\mathbf{y}} - \frac{1}{n} \underset{\sim}{\mathbf{1}} \underset{\sim}{\mathbf{1}}' \underset{\sim}{\mathbf{y}})$

$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (\underset{\sim}{\hat{\mathbf{y}}} - \frac{1}{n} \underset{\sim}{\mathbf{1}} \underset{\sim}{\mathbf{1}}' \underset{\sim}{\hat{\mathbf{y}}})' (\underset{\sim}{\hat{\mathbf{y}}} - \frac{1}{n} \underset{\sim}{\mathbf{1}} \underset{\sim}{\mathbf{1}}' \underset{\sim}{\hat{\mathbf{y}}})$

$MSE = \underset{\sim}{\mathbf{e}}'\underset{\sim}{\mathbf{e}} / (n-p)$ $MSTO = SSTO / n-1$ $MSR = SSR / p-1$

Coefficient of Multiple Determination

$$R^2 = SSR / SSTO = 1 - SSE / SSTO$$

* What proportion of the variability in y is explained by the covariates

* R^2 is nondecreasing as you add more covariates

⇒ If you decide to pick your covariates to maximize R^2 , you will always decide to use all of them!

Brings us to...

Adjusted Coefficient of Multiple Determination

$$R_a^2 = 1 - \frac{MSE}{MSTO} = 1 - \frac{SSE / n - p}{SSTO / n - 1}$$

* R_a^2 can decrease as more covariates are added

⇒ If you were to pick your "best" model using R_a^2 , it won't necessarily include all of your covariates

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

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

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

$$= \sigma^2 (X'X)^{-1}$$

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