

ABE 6933 SML, Fall 2020

A Matrix Algebra Approach to Linear Regression

Instructor: Dr. Nikolay Bliznyuk

29 September, 2020

Slides under "[4].lectures\probability.and.stats.slides"

Matrix as a Rectangular Array

A matrix with r rows and c columns is a rectangular array. It will be represented either in full form

$$\mathbf{A}_{r \times c} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1j} & \cdots & a_{1c} \\ a_{21} & a_{22} & \cdots & a_{2j} & \cdots & a_{2c} \\ \vdots & \vdots & & \vdots & & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{ij} & \cdots & a_{ic} \\ \vdots & \vdots & & \vdots & & \vdots \\ a_{r1} & a_{r2} & \cdots & a_{rj} & \cdots & a_{rc} \end{bmatrix},$$

or in abbreviated form

$$\mathbf{A} = [a_{ij}] \quad i = 1, \dots, r; \quad j = 1, \dots, c;$$

or simply by a boldface symbol, such as \mathbf{A} .

Transpose of a Matrix: an Illustration

The transpose of a matrix \mathbf{A} is another matrix, denoted by changing corresponding columns and rows of the matrix \mathbf{A} . For example, if

$$\mathbf{A}_{3 \times 2} = \begin{bmatrix} 2 & 5 \\ 7 & 10 \\ 3 & 4 \end{bmatrix},$$

then the transpose of \mathbf{A} is

$$\mathbf{A}'_{2 \times 3} = \begin{bmatrix} 2 & 7 & 3 \\ 5 & 10 & 4 \end{bmatrix}.$$

Transpose of a Matrix

$$\mathbf{A}_{r \times c} = \begin{bmatrix} a_{11} & \cdots & a_{1c} \\ \vdots & & \vdots \\ a_{r1} & \cdots & a_{rc} \end{bmatrix} = [a_{ij}], \quad i = 1, \dots, r; \quad j = 1, \dots, c.$$

$$\mathbf{A}'_{c \times r} = \begin{bmatrix} a_{11} & \cdots & a_{r1} \\ \vdots & & \vdots \\ a_{1c} & \cdots & a_{rc} \end{bmatrix} = [a_{ji}], \quad j = 1, \dots, c; \quad i = 1, \dots, r.$$

In R, the call to compute \mathbf{A}' is `t(A)`.

Addition and Subtraction of Matrices

In general, if

$$\mathbf{A}_{r \times c} = [a_{ij}] \quad \text{and} \quad \mathbf{B}_{r \times c} = [b_{ij}], \quad i = 1, \dots, r; \quad j = 1, \dots, c,$$

then

$$\mathbf{A}_{r \times c} + \mathbf{B}_{r \times c} = [a_{ij} + b_{ij}] \quad \text{and} \quad \mathbf{A}_{r \times c} - \mathbf{B}_{r \times c} = [a_{ij} - b_{ij}].$$

In \mathbb{R} , standard $+/ -$ operations apply, e.g., $(\mathbf{A} - \mathbf{B})$ so long as the dimensions are conformable.

Matrix Multiplication - Example

$$\begin{aligned}\mathbf{AB} &= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix} \\ &= \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{bmatrix}.\end{aligned}$$

Matrix Multiplication in General

In general, if \mathbf{A} has dimension $r \times c$ and \mathbf{B} has dimension $c \times s$, the product \mathbf{AB} is a matrix of dimension $r \times s$ whose element in the i th row and j th column is

$$\sum_{k=1}^c a_{ik} b_{kj},$$

so that

$$\mathbf{AB}_{r \times s} = \left[\sum_{k=1}^c a_{ik} b_{kj} \right] \quad i = 1, \dots, r; \quad j = 1, \dots, s$$

In R, $\mathbf{AB} = \mathbf{A} \%*\% \mathbf{B}$.

Matrix Multiplication: Alternative Views

$$\begin{aligned}\mathbf{AB} &= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix} \\ &= \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{bmatrix}.\end{aligned}$$

Let $\mathbf{C} = \mathbf{AB}$. Recall that $\mathbf{C}_{ij} = \sum_{k=1}^c a_{ik}b_{kj}$.

Column form representation: We can represent columns of \mathbf{C} as linear combinations of the columns of \mathbf{A} : $\mathbf{C}_j = \mathbf{A} \cdot \mathbf{B}_j$, where \mathbf{B}_j is the j th column of \mathbf{B} so that $\mathbf{C} = [\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_s]$.

We can represent \mathbf{C} in the *outer product form*, i.e., $\mathbf{C} = \sum_{k=1}^c \mathbf{A}[, k] \cdot \mathbf{B}[k,]$ using R notation, where $\mathbf{A}[, k]$ is the k th column of \mathbf{A} and $\mathbf{B}[k,]$ is the k th row of \mathbf{B} .

The Identity Matrix

The identity matrix, denoted by \mathbf{I} , is a diagonal matrix whose elements on the main diagonal are all 1s. Premultiplying or postmultiplying any $r \times c$ matrix \mathbf{A} by the identity matrix (of conformable dimensions) leaves \mathbf{A} unchanged. For example,

$$\mathbf{IA} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

Similarly,

$$\mathbf{AI} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

In general, for any $r \times r$ matrix \mathbf{A} we have $\mathbf{AI} = \mathbf{IA} = \mathbf{A}$.
In R, the $r \times r$ identity matrix can be created as `diag(r)`.

Linear Dependence/Multicollinearity

Let \mathbf{A} be a matrix with columns $\mathbf{A}_1, \dots, \mathbf{A}_c$.

If one can find scalars $\lambda_1, \dots, \lambda_c$, not all zero, such that

$$\mathbf{A} \cdot \boldsymbol{\lambda} = \lambda_1 \mathbf{A}_1 + \lambda_2 \mathbf{A}_2 + \dots + \lambda_c \mathbf{A}_c = \mathbf{0},$$

where $\mathbf{0}$ denotes the zero column vector, the column vectors are *linearly dependent*.

If the only set of scalars for which the equality holds is $\lambda_1 = 0, \dots, \lambda_c = 0$, the columns are *linearly independent*.

Let \mathbf{X} be the design matrix for a multiple linear regression problem; i.e., columns of \mathbf{X} are predictors/features.

Collinearity/multicollinearity occurs when the columns matrix \mathbf{X} are linearly dependent (loosely, contain redundant information).

Linear Dependence: an Illustration

Consider the following matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 5 & 1 \\ 2 & 2 & 10 & 6 \\ 3 & 4 & 15 & 1 \end{bmatrix}.$$

If $\lambda_1 = 5, \lambda_2 = 0, \lambda_3 = -1, \lambda_4 = 0$, then

$$5 \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 0 \begin{bmatrix} 2 \\ 2 \\ 4 \end{bmatrix} - 1 \begin{bmatrix} 5 \\ 10 \\ 15 \end{bmatrix} + 0 \begin{bmatrix} 1 \\ 6 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

In general, linear dependence is not restricted to the situations where one column is a multiple of another column.

Rank of a Matrix

The rank of a matrix is defined to be the maximum number of linearly independent columns in the matrix.

We know that the rank of \mathbf{A} in our earlier example cannot be 4, since the four columns are linearly dependent.

We can, however, find three columns (1, 2, and 4) which are linearly independent. There are no scalars $\lambda_1, \lambda_2, \lambda_4$ such that $\lambda_1 \mathbf{A}_1 + \lambda_2 \mathbf{A}_2 + \lambda_4 \mathbf{A}_4 = \mathbf{0}$ other than $\lambda_1 = \lambda_2 = \lambda_4 = 0$. Thus, the rank of \mathbf{A} in our example is 3.

The rank of a matrix is unique and can equivalently be defined as the maximum number of linearly independent rows. It follows that the rank of an $r \times c$ matrix cannot exceed $\min(r, c)$, the minimum of the two values r and c .

Inverse of a Matrix

For a square matrix \mathbf{A} of full rank, the inverse of \mathbf{A} is a matrix \mathbf{A}^{-1} such that $\mathbf{A}^{-1}\mathbf{A} = \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$.

For example, the inverse of the matrix

$$\mathbf{A}_{2 \times 2} = \begin{bmatrix} 2 & 4 \\ 3 & 1 \end{bmatrix}$$

is

$$\mathbf{A}_{2 \times 2}^{-1} = \begin{bmatrix} -.1 & .4 \\ .3 & -.2 \end{bmatrix}$$

since

$$\mathbf{A}^{-1}\mathbf{A} = \begin{bmatrix} -.1 & .4 \\ .3 & -.2 \end{bmatrix} \begin{bmatrix} 2 & 4 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{I}$$

or

$$\mathbf{A}\mathbf{A}^{-1} = \begin{bmatrix} 2 & 4 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} -.1 & .4 \\ .3 & -.2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{I}.$$

Example of an Inverse of a 3x3 Matrix

If

$$\mathbf{B}_{3 \times 3} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & k \end{bmatrix},$$

then

$$\mathbf{B}^{-1} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & k \end{bmatrix}^{-1} = \begin{bmatrix} A & B & C \\ D & E & F \\ G & H & K \end{bmatrix},$$

where

$$\begin{aligned} A &= (ek - fh)/Z & B &= -(bk - ch)/Z & C &= (bf - ce)/Z \\ D &= -(dk - fg)/Z & E &= (ak - cg)/Z & F &= -(af - cd)/Z \\ G &= (dh - eg)/Z & H &= -(ah - bg)/Z & K &= (ae - bd)/Z \end{aligned}$$

and

$$Z = a(ek - fh) - b(dk - fg) + c(dh - eg).$$

Z is called the determinant of the matrix \mathbf{B} .

Solving Systems of Linear Equations

A solution to the system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ is a vector \mathbf{x}^* for which the identity $\mathbf{A}\mathbf{x}^* = \mathbf{b}$ is satisfied.

One is generally interested in solving the systems of equations where \mathbf{A} is a square $r \times r$ matrix. Such equations have a unique solution for a general right-hand side \mathbf{b} if and only if \mathbf{A} is of full rank (i.e., invertible), in which case the solution $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$.

Although $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$ is the "mathematical/theoretical" solution, in practice finding \mathbf{A}^{-1} to solve the linear system is generally a bad idea (from the standpoint of numerical accuracy).

In R, the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ may be solved as `solve(A,b)`. In the background, this linear system is solved using efficient matrix factorizations of \mathbf{A} , e.g., the "LU" factorization.

In the rare cases where \mathbf{A}^{-1} itself is needed (e.g., when we need a covariance matrix of our least squares estimator $\hat{\beta}$), it can be found as `solve(A)`.

Common Matrix Algebra Identities

$$\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$$

$$(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$$

$$(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$$

$$\mathbf{C}(\mathbf{A} + \mathbf{B}) = \mathbf{CA} + \mathbf{CB}$$

$$\lambda(\mathbf{A} + \mathbf{B}) = \lambda\mathbf{A} + \lambda\mathbf{B}$$

$$(\mathbf{A}')' = \mathbf{A}$$

$$(\mathbf{A} + \mathbf{B})' = \mathbf{A}' + \mathbf{B}'$$

$$(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$$

$$(\mathbf{ABC})' = \mathbf{C}'\mathbf{B}'\mathbf{A}'$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$$

$$(\mathbf{ABC})^{-1} = \mathbf{C}^{-1}\mathbf{B}^{-1}\mathbf{A}^{-1}$$

$$(\mathbf{A}^{-1})^{-1} = \mathbf{A}$$

$$(\mathbf{A}')^{-1} = (\mathbf{A}^{-1})'$$

Expectation of a Random Vector or a Matrix

In general, for a random vector \mathbf{Y} the expectation is

$$\mathbf{E}\{\mathbf{Y}\} = [E\{Y_i\}] \quad i = 1, \dots, n.$$

For a random matrix \mathbf{Y} with dimension $n \times p$, the expectation is

$$\mathbf{E}\{\mathbf{Y}\}_{n \times p} = [E\{Y_{ij}\}] \quad i = 1, \dots, n; \quad j = 1, \dots, p.$$

Expectation of a Random Vector: SLR/MLR Example

$$\begin{array}{c} \mathbf{Y} \\ n \times 1 \end{array} = \begin{array}{c} \mathbf{E}\{\mathbf{Y}\} \\ n \times 1 \end{array} + \begin{array}{c} \boldsymbol{\varepsilon} \\ n \times 1 \end{array}$$
$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} E\{Y_1\} \\ E\{Y_2\} \\ \vdots \\ E\{Y_n\} \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} = \begin{bmatrix} E\{Y_1\} + \varepsilon_1 \\ E\{Y_2\} + \varepsilon_2 \\ \vdots \\ E\{Y_n\} + \varepsilon_n \end{bmatrix}$$

Example: Covariance Matrix I

Let \mathbf{Y} be a random vector consisting of three rvs Y_1, Y_2, Y_3 .

The variances of the three rvs, $\sigma^2\{Y_i\}$, and the covariances between any two of the rvs, $\sigma\{Y_i, Y_j\}$, are assembled in the *covariance matrix* of \mathbf{Y} , denoted by $\sigma^2\{\mathbf{Y}\}$ as follows:

$$\sigma^2\{\mathbf{Y}\} = \begin{bmatrix} \sigma^2\{Y_1\} & \sigma\{Y_1, Y_2\} & \sigma\{Y_1, Y_3\} \\ \sigma\{Y_2, Y_1\} & \sigma^2\{Y_2\} & \sigma\{Y_2, Y_3\} \\ \sigma\{Y_3, Y_1\} & \sigma\{Y_3, Y_2\} & \sigma^2\{Y_3\} \end{bmatrix}.$$

Notice that $\sigma\{Y_2, Y_1\} = \sigma\{Y_1, Y_2\}$, since $\sigma\{Y_i, Y_j\} = \sigma\{Y_j, Y_i\}$ for all $i \neq j$, $\sigma^2\{\mathbf{Y}\}$ is a symmetric matrix.

In this course, the terms “covariance matrix” and “variance-covariance matrix” are used interchangeably.

Example: Covariance Matrix II

It follows readily that:

$$\sigma^2\{\mathbf{Y}\} = \mathbf{E} \{ (\mathbf{Y} - \mathbf{E}\{\mathbf{Y}\}) \cdot (\mathbf{Y} - \mathbf{E}\{\mathbf{Y}\})' \}.$$

For our illustration, we have

$$\sigma^2\{\mathbf{Y}\} = \mathbf{E} \left\{ \begin{bmatrix} Y_1 - E\{Y_1\} \\ Y_2 - E\{Y_2\} \\ Y_3 - E\{Y_3\} \end{bmatrix} \cdot [Y_1 - E\{Y_1\}, Y_2 - E\{Y_2\}, Y_3 - E\{Y_3\}] \right\}.$$

If we define $\mathbf{Z} = (\mathbf{Y} - \mathbf{E}\{\mathbf{Y}\})$, then $\sigma^2\{\mathbf{Y}\} = \mathbf{E}(\mathbf{Z} \cdot \mathbf{Z}')$ and $[\sigma^2\{\mathbf{Y}\}]_{ij} = \text{Cov}(Y_i, Y_j) = \text{Cov}(Y_j, Y_i) = [\sigma^2\{\mathbf{Y}\}]_{ji}$.

$$\text{Cov}(Y_i, Y_j) = \mathbf{E} \left((Y_i - E(Y_i)) \cdot (Y_j - E(Y_j)) \right)$$

Covariance Matrix - General Case

The covariance matrix for a general $n \times 1$ random vector \mathbf{Y} is

$$\sigma^2\{\mathbf{Y}\} = \begin{bmatrix} \sigma^2\{Y_1\} & \sigma\{Y_1, Y_2\} & \cdots & \sigma\{Y_1, Y_n\} \\ \sigma\{Y_2, Y_1\} & \sigma^2\{Y_2\} & \cdots & \sigma\{Y_2, Y_n\} \\ \vdots & \vdots & & \vdots \\ \sigma\{Y_n, Y_1\} & \sigma\{Y_n, Y_2\} & \cdots & \sigma^2\{Y_n\} \end{bmatrix}.$$

Notice again that $\sigma^2\{\mathbf{Y}\}$ is a symmetric matrix, i.e.,
 $[\sigma^2\{\mathbf{Y}\}]_{ij} = [\sigma^2\{\mathbf{Y}\}]_{ji}.$

For notational transparency, the covariance matrix of \mathbf{Y} is denoted here as $\sigma^2\{\mathbf{Y}\}$. A more common notation is $Var(\mathbf{Y}) = \Sigma_{\mathbf{Y}}$ whenever there are multiple random vectors under consideration, or $Var(\mathbf{Y}) = \Sigma$ if there is no ambiguity.

Expectation and Covariance for a Linear Transformation

Frequently, we shall encounter a random vector \mathbf{W} which is obtained by premultiplying the random vector \mathbf{Y} by a constant matrix \mathbf{A} (a matrix whose elements are fixed):

$$\mathbf{W} = \mathbf{A}\mathbf{Y}.$$

Here, \mathbf{W} is called a linear transformation of \mathbf{Y} .

Some basic results for this case are

$$\mathbf{E}\{\mathbf{A}\} = \mathbf{A}$$

$$\mathbf{E}\{\mathbf{W}\} = \mathbf{E}\{\mathbf{A}\mathbf{Y}\} = \mathbf{A}\mathbf{E}\{\mathbf{Y}\}$$

$$\sigma^2\{\mathbf{W}\} = \sigma^2\{\mathbf{A}\mathbf{Y}\} = \mathbf{A}\sigma^2\{\mathbf{Y}\}\mathbf{A}',$$

where $\sigma^2\{\mathbf{Y}\}$ is the variance-covariance matrix of \mathbf{Y} .

Expectation for a Linear Transformation

$$\mathbf{E}\{\mathbf{W}\}_{2 \times 1} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} E\{Y_1\} \\ E\{Y_2\} \end{bmatrix} = \begin{bmatrix} E\{Y_1\} - E\{Y_2\} \\ E\{Y_1\} + E\{Y_2\} \end{bmatrix}$$

Covariance Matrix for a Linear Transformation

$$\begin{aligned}\sigma^2\{\mathbf{W}\}_{2 \times 2} &= \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma^2\{Y_1\} & \sigma\{Y_1, Y_2\} \\ \sigma\{Y_2, Y_1\} & \sigma^2\{Y_2\} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} \sigma^2\{Y_1\} + \sigma^2\{Y_2\} - 2\sigma\{Y_1, Y_2\} & \sigma^2\{Y_1\} - \sigma^2\{Y_2\} \\ \sigma^2\{Y_1\} - \sigma^2\{Y_2\} & \sigma^2\{Y_1\} + \sigma^2\{Y_2\} + 2\sigma\{Y_1, Y_2\} \end{bmatrix}\end{aligned}$$

Simple Linear Regression Model Using Equations

Our model for the individual observations Y_i is

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, n.$$

This implies

$$Y_1 = \beta_0 + \beta_1 X_1 + \varepsilon_1,$$

$$Y_2 = \beta_0 + \beta_1 X_2 + \varepsilon_2,$$

$$\vdots$$

$$Y_n = \beta_0 + \beta_1 X_n + \varepsilon_n.$$

Simple Linear Regression using Matrix Algebra

$$\begin{aligned} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} &= \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \\ &= \begin{bmatrix} \beta_0 + \beta_1 X_1 \\ \beta_0 + \beta_1 X_2 \\ \vdots \\ \beta_0 + \beta_1 X_n \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 X_1 + \varepsilon_1 \\ \beta_0 + \beta_1 X_2 + \varepsilon_2 \\ \vdots \\ \beta_0 + \beta_1 X_n + \varepsilon_n \end{bmatrix} \end{aligned}$$

Multiple Linear Regression Model Using Matrix Algebra I

$$\begin{aligned}\mathbf{Y}_{n \times 1} &= \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} & \mathbf{X}_{n \times p} &= \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,p-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,p-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,p-1} \end{bmatrix} \\ \boldsymbol{\beta}_{p \times 1} &= \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} & \boldsymbol{\varepsilon}_{n \times 1} &= \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}\end{aligned}$$

Multiple Linear Regression Model Using Matrix Algebra II

In matrix terms, a multiple linear regression (MLR) model is

$$\underset{n \times 1}{\mathbf{Y}} = \underset{n \times p}{\mathbf{X}} \cdot \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}, \quad \text{where}$$

\mathbf{Y} is a vector of responses,

$\boldsymbol{\beta}$ is a vector of parameters/coefficients,

\mathbf{X} is a matrix of constants (the design matrix), and

$\boldsymbol{\varepsilon}$ is a vector of uncorrelated errors with expectation $\mathbf{E}\{\boldsymbol{\varepsilon}\} = \mathbf{0}$ and covariance matrix

$$\underset{n \times n}{\boldsymbol{\sigma}^2\{\boldsymbol{\varepsilon}\}} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}.$$

Least Squares (LS) Estimation for the MLR

$$Q(\mathbf{b}) = \sum_{i=1}^n (Y_i - b_0 - b_1 X_{i1} - \cdots - b_{p-1} X_{i,p-1})^2.$$

The least squares (LS) solution/estimators are those values of b_0, b_1, \dots, b_{p-1} that minimize the SSE, here denoted by Q . Define

$$\underset{p \times 1}{\mathbf{b}} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{p-1} \end{bmatrix}.$$

In matrix notation,

$$Q(\mathbf{b}) = (\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b}) = \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2.$$

Expanding, we obtain

$$Q(\mathbf{b}) = \mathbf{Y}'\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{Y} - \mathbf{Y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{Y}'\mathbf{Y} - 2\mathbf{Y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}.$$

Normal Equations and the LS Estimator

To find the minimizer of $Q(\mathbf{b})$, differentiate $Q(\mathbf{b})$ wrt \mathbf{b} :

$$\frac{\partial Q(\mathbf{b})}{\partial \mathbf{b}} = -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\mathbf{b},$$

set the derivative to $\mathbf{0}$ and solve.

Notice the solution must satisfy $\mathbf{X}' \cdot (\mathbf{Y} - \mathbf{X}\mathbf{b}) = \mathbf{0}$.

The least squares *normal equations* are $\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y}$.

To solve, premultiply both sides by $(\mathbf{X}'\mathbf{X})^{-1}$ (assume this exists):

$$(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}.$$

Since $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} = \mathbf{I}$ and $\mathbf{I}\mathbf{b} = \mathbf{b}$, we then find the solution

$$\underset{m \times 1}{\mathbf{b}^*} = \underset{m \times m}{(\mathbf{X}'\mathbf{X})^{-1}} \underset{m \times 1}{\mathbf{X}'\mathbf{Y}}, \quad \text{the LS estimator that minimizes } Q(\mathbf{b}).$$

Statistical Model for MLR

How to solve this problem “statistically”? Assume

$$Y_i = \underbrace{\mathbf{x}_i' \boldsymbol{\beta}}_{\sum_{j=1}^p x_{ij} \beta_j} + \varepsilon_i,$$

where ε_i 's are independent $\text{Normal}(0, \sigma^2)$ rvs. In vector-matrix form,

$$\underset{n \times 1}{\mathbf{Y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}.$$

Here, $Y_i \sim \text{Normal}(\mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$; the Y_i 's are independent but not identically distributed.

Since we know the joint pdf of the Y_i 's, we can estimate $\boldsymbol{\beta}$ and σ^2 using the MLE.

Statistical Estimation by MLE I

Step 1: write down the likelihood:

$$\begin{aligned} L(y_1, \dots, y_n \mid \beta, \sigma^2) &= \prod_{i=1}^n f_i(y_i \mid \beta, \sigma^2) \\ &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y_i - \mathbf{x}_i'\beta)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^n (y_i - \mathbf{x}_i'\beta)^2}_{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)}\right). \end{aligned}$$

Step 2: obtain the log-likelihood:

$$\begin{aligned} \ell(\beta, \sigma^2) &= \ln(L(y_1, \dots, y_n \mid \beta, \sigma^2)) \\ &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta). \end{aligned}$$

Statistical Estimation by MLE II

Step 3. Find the gradient of $\ell(\beta, \sigma^2)$ with respect to β and σ^2 , set the gradient to 0, solve for β and σ^2 .

$$\begin{aligned}\frac{\partial \ell(\beta, \sigma^2)}{\partial \beta} &= -\frac{1}{2\sigma^2} 2\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) \\ \frac{\partial \ell(\beta, \sigma^2)}{\partial \sigma^2} &= -\frac{n}{2} \frac{1}{\sigma^2} - \frac{1}{2} \frac{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)}{(\sigma^2)^2}\end{aligned}$$

The solution is $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ and $\hat{\sigma}^2 = \frac{1}{n}(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})$.

Step 4. Make sure we found the maximizers of $\ell(\beta, \sigma^2)$ by showing that all eigenvalues of the matrix of second derivatives of $\ell(\beta, \sigma^2)$ —known as the Hessian—are negative. (Equivalently, $(-1) \cdot \text{Hessian}$ is positive definite.)

The Vector of Fitted Values and the Hat Matrix

Notice that the expressions for the LS and ML estimators coincide. Let's use $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$ for notational transparency.

Let's express the vector of fitted values $\hat{\mathbf{Y}}$ using the formula for $\hat{\beta}$:

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = \mathbf{H}\mathbf{Y},$$

where

$$\mathbf{H}_{n \times n} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'.$$

Notice that \mathbf{H} is symmetric ($\mathbf{H} = \mathbf{H}'$) and idempotent, i.e.,

$$\mathbf{H} = \mathbf{H} \cdot \mathbf{H}.$$

The Hat Matrix and the Vector of Residuals

We can express the vector of residuals \mathbf{e} as

$$\underset{n \times 1}{\mathbf{e}} = \mathbf{Y} - \hat{\mathbf{Y}} = (\underset{n \times 1}{\mathbf{I}} - \overset{n \times n}{\mathbf{H}}) \cdot \underset{n \times 1}{\mathbf{Y}},$$

where \mathbf{H} is the hat matrix. The matrix $(\mathbf{I} - \mathbf{H})$, like the matrix \mathbf{H} , is symmetric and idempotent.

The variance-covariance matrix of the vector of residuals \mathbf{e} is

$$\sigma^2\{\mathbf{e}\} = \sigma^2(\mathbf{I} - \mathbf{H})$$

and is estimated by

$$\mathbf{s}^2\{\mathbf{e}\} = \hat{\sigma}^2(\mathbf{I} - \mathbf{H}),$$

where $\hat{\sigma}^2 = SSE/(n - p) = \mathbf{e}'\mathbf{e}/(n - p)$ is referred to as the MSE in ANOVA tables.

Distributional Results when $\varepsilon \sim MVN(\mathbf{0}, \sigma^2 \mathbf{I})$

The ML estimator is $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = \mathbf{A}\mathbf{Y}$. Hence

$$\begin{aligned}\mathbf{E}(\hat{\beta}) &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \cdot \mathbf{E}(\mathbf{Y}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\beta = \beta, \\ \sigma^2\{\hat{\beta}\} &= \mathbf{A} \cdot \sigma^2\{\mathbf{Y}\} \cdot \mathbf{A}' = \dots = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}.\end{aligned}$$

Additionally,

$$\begin{aligned}\hat{\beta} &\sim \text{Multivariate.Normal}(\beta, \sigma^2\{\hat{\beta}\}), \quad \text{and} \\ \mathbf{a}'\hat{\beta} &\sim \text{Normal}(\mathbf{a}'\beta, \mathbf{a}'\sigma^2\{\hat{\beta}\}\mathbf{a}),\end{aligned}$$

where \mathbf{a} is a column vector of constants.

Lastly,

$$\frac{\mathbf{a}'\hat{\beta} - \mathbf{a}'\beta}{\hat{\sigma}\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}} \sim t(n-p).$$

Estimation/Prediction of the Mean Response

For given values of X_1, \dots, X_{p-1} , denoted by $x_{h,1}, \dots, x_{h,p-1}$, the mean response is denoted by $\mathbf{E}\{Y_h\}$. We define the vector \mathbf{x}_h

$$\mathbf{x}_h = \begin{bmatrix} 1 \\ x_{h,1} \\ \vdots \\ x_{h,p-1} \end{bmatrix},$$

so that the mean response to be estimated is

$$\mathbf{E}\{Y_h\} = \mathbf{x}_h' \boldsymbol{\beta}.$$

The estimated mean response corresponding to \mathbf{x}_h is

$$\hat{Y}_h = \mathbf{x}_h' \hat{\boldsymbol{\beta}}.$$

Estimation/Prediction of the Mean Response

This estimator $\hat{Y}_h = \mathbf{x}'_h \hat{\boldsymbol{\beta}}$ is unbiased:

$$\mathbf{E} \left\{ \hat{Y}_h \right\} = \mathbf{x}'_h \boldsymbol{\beta} = \mathbf{E} \{ Y_h \}$$

and its variance is

$$\sigma^2 \left\{ \hat{Y}_h \right\} = \sigma^2 \mathbf{x}'_h (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_h.$$

This variance can be expressed as a function of the variance-covariance matrix of the estimated regression coefficients

$$\sigma^2 \left\{ \hat{Y}_h \right\} = \mathbf{x}'_h \sigma^2 \{ \hat{\boldsymbol{\beta}} \} \mathbf{x}_h.$$

Confidence Interval for the Mean of the Response Y_h

Notice that the variance $\sigma^2 \left\{ \hat{Y}_h \right\}$ is a function of the covariance matrix $\sigma^2 \{ \hat{\beta} \} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$.

The estimated variance $s^2 \left\{ \hat{Y}_h \right\}$ is given by

$$s^2 \left\{ \hat{Y}_h \right\} = MSE \cdot \left(\mathbf{x}_h' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_h \right) = \mathbf{x}_h' \mathbf{s}^2 \{ \mathbf{b} \} \mathbf{x}_h.$$

The $(1 - \alpha)$ confidence limits for $E \{ Y_h \}$ are

$$\hat{Y}_h \pm t(1 - \alpha/2; n - p) s \left\{ \hat{Y}_h \right\}.$$

Here, $MSE = \frac{SSE}{(n-p)} = \hat{\sigma}^2$ is the square of the "residual standard error" (RSE) reported by R in the summary of an `lm` object.

Prediction Interval for the Response Y_h

The $(1 - \alpha)$ prediction limits for a new observation Y_h corresponding to \mathbf{x}_h , the specified values of the covariates, are

$$\hat{Y}_h \pm t(1 - \alpha/2; n - p) \cdot s\{\text{pred}\},$$

where

$$s^2\{\text{pred}\} = MSE + s^2 \left\{ \hat{Y}_h \right\} = MSE \cdot \left(1 + \mathbf{x}_h' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_h \right)$$

and $s^2 \left\{ \hat{Y}_h \right\}$ is given above.

In R, point-level predictions \hat{Y}_h , confidence intervals and prediction intervals can be obtained using the function `predict.lm`.