

Some algebra underlying the linear regression model

Loïc Yengo

CNSG, University of Queensland

June 29th, 2016

Contents

Introduction

Basic operations with matrices

Linear regression in matrix terms

Efficient matrix inversion

Linear Regression

Definition

For an individual i , we model a quantitative trait y_i as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \quad (i = 1, \dots, n)$$

β_0, \dots, β_p are the regression coefficients and x_{i1}, \dots, x_{ip} will be called covariates.

Objective

- ▶ Estimate the regression coefficients
- ▶ Test assumptions about these coefficients
- ▶ Make predictions for new data

Linear Regression

Definition

For an individual i , we model a quantitative trait y_i as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \quad (i = 1, \dots, n)$$

β_0, \dots, β_p are the regression coefficients and x_{i1}, \dots, x_{ip} will be called covariates.

Objective

- ▶ Estimate the regression coefficients
- ▶ Test assumptions about these coefficients
- ▶ Make predictions for new data

(Part I)
Red pill or blue pill?



What is a matrix?

A n -by- p **matrix** \mathbf{X} is a two-dimensional table with n rows and p columns that contains (real) numbers. n and p are called matrix dimensions.

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & x_{ij} & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

The x_{ij} 's are called matrix entries.

An example

$$\mathbf{X} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \text{ is a 2-by-2 matrix.}$$

A **vector** is a special case of a matrix that has only one column.

What is a matrix?

A n -by- p **matrix** \mathbf{X} is a two-dimensional table with n rows and p columns that contains (real) numbers. n and p are called matrix dimensions.

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & x_{ij} & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

The x_{ij} 's are called matrix entries.

An example

$$\mathbf{X} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \text{ is a 2-by-2 matrix.}$$

A **vector** is a special case of a matrix that has only one column.

We assume that \mathbf{X} is a n -by- p matrix.

0-Transposition

The transpose of matrix \mathbf{X} (denoted ${}^t\mathbf{X}$, \mathbf{X}^T or \mathbf{X}') is a p -by- n matrix defined as

$$\mathbf{X}' = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & x_{ji} & \vdots \\ x_{p1} & \dots & x_{pn} \end{bmatrix}$$

1-Multiplying by a number a

$a\mathbf{X}$ is also n -by- p matrix which entries have all been multiplied by a

$$a\mathbf{X} = \begin{bmatrix} a \times x_{11} & \dots & a \times x_{1p} \\ \vdots & a \times x_{ij} & \vdots \\ a \times x_{n1} & \dots & a \times x_{np} \end{bmatrix}$$

2-Matrix addition

We can only add matrices that have the same dimensions. If \mathbf{Z} is also a n -by- p matrix then

$$\begin{aligned}\mathbf{X} + \mathbf{Z} &= \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & x_{ij} & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix} + \begin{bmatrix} z_{11} & \dots & z_{1p} \\ \vdots & z_{ij} & \vdots \\ z_{n1} & \dots & z_{np} \end{bmatrix} \\ &= \begin{bmatrix} x_{11} + z_{11} & \dots & x_{1p} + z_{1p} \\ \vdots & x_{ij} + z_{ij} & \vdots \\ x_{n1} + z_{n1} & \dots & x_{np} + z_{np} \end{bmatrix}\end{aligned}$$

This rule can be generalized to calculate the sum of any number of matrices.

3-Matrix multiplication

We can only multiply a $n \times p$ -matrix with a $p \times m$ -matrix. If \mathbf{Z} is a p -by- m matrix then \mathbf{XZ} is a n -by- m matrix defined as

$$(\mathbf{XZ})_{i,k} = x_{i,1} \times z_{1,k} + x_{i,2} \times z_{2,k} + \dots + x_{i,p} \times z_{p,k}$$

3-Matrix multiplication

We can only multiply a $n \times p$ -matrix with a $p \times m$ -matrix. If \mathbf{Z} is a p -by- m matrix then \mathbf{XZ} is a n -by- m matrix defined as

$$(\mathbf{XZ})_{i,k} = x_{i,1} \times z_{1,k} + x_{i,2} \times z_{2,k} + \dots + x_{i,p} \times z_{p,k}$$

$$\begin{bmatrix} 1 & -1 & 4 \\ 2 & 5 & 7 \end{bmatrix} \begin{bmatrix} 2 & 1 & 0 \\ 4 & 4 & 3 \\ -2 & 7 & 1 \end{bmatrix} = \begin{bmatrix} -10 & 25 & 1 \\ 10 & 71 & 22 \end{bmatrix}$$

$$1 \cdot 2 + (-1) \cdot 4 + 4 \cdot (-2) = -10$$

$$2 \cdot 1 + 5 \cdot 4 + 7 \cdot 7 = 71$$

Identity matrix and inverse

We call a square matrix, a matrix with as many rows as columns. One special square matrix is the **Identity matrix** defined in dimension n as

$$I_n = n \text{ rows} \left\{ \begin{array}{c} \left[\begin{array}{ccc} 1 & & \\ & 1 & 0 \\ & & \ddots \\ & 0 & 1 \\ & & & 1 \end{array} \right] \\ \underbrace{\hspace{10em}}_{n \text{ columns}} \end{array} \right.$$

A $n \times n$ matrix \mathbf{A} is said **invertible**, if there exists a unique $n \times n$ matrix \mathbf{B} such as $\mathbf{AB} = \mathbf{BA} = I_n$. If such a matrix exist then \mathbf{B} is called **inverse** of \mathbf{A} and is noted \mathbf{A}^{-1} .

How can I tell if a matrix is invertible?

Tough question in general, but we can admit this somewhat satisfactory answer...**A** is invertible if and only if its determinant does equal 0. OK then, but what is the determinant?

Example

$$\mathbf{X} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \text{ then } |\mathbf{X}| = 1 \times 4 - 2 \times 3 = -2.$$

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \text{ then } |\mathbf{X}| = 1 \times 4 - 2 \times 2 = 0$$

The determinant measures how the columns or the rows of **X** are linearly related.

How can I tell if a matrix is invertible?

Tough question in general, but we can admit this somewhat satisfactory answer...**A** is invertible if and only if its determinant does equal 0. OK then, but what is the determinant?

Example

$$\mathbf{X} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \text{ then } |\mathbf{X}| = 1 \times 4 - 2 \times 3 = -2.$$

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \text{ then } |\mathbf{X}| = 1 \times 4 - 2 \times 2 = 0$$

The determinant measures how the columns or the rows of **X** are linearly related.

Matrix rank

The rank of a $n \times p$ matrix \mathbf{A} is the number of linearly independent rows or columns.

$$\text{rank}(\mathbf{A}) \leq \min(n, p).$$

Matrix rank

The rank of a $n \times p$ matrix \mathbf{A} is the number of linearly independent rows or columns.

$$\text{rank}(\mathbf{A}) \leq \min(n, p).$$

When the $n = p$ (square matrices), then \mathbf{A} is invertible if and only if $\text{rank}(\mathbf{A}) = n$ (full rank matrix).

(Part II)
Linear regression in matrix terms

Reformulation

We first defined the linear regression model using the following equation

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \quad (i = 1, \dots, n)$$

this looks pretty much like a matrix product right?

Reformulation

We first defined the linear regression model using the following equation

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \quad (i = 1, \dots, n)$$

this looks pretty much like a matrix product right?

$$y_i = \underbrace{[1, x_{i1}, x_{i2}, \dots, x_{ip}]}_{1 \times (p+1)} \times \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}}_{(p+1) \times 1} + \varepsilon_i \quad (i = 1, \dots, n)$$

Therefore, if we note

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

and

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & x_{ij} & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

The linear regression model can be reformulated as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Estimation

I assume that we know the distribution of the ε_j 's:
 $\varepsilon_j \sim \mathcal{N}(0, \sigma^2)$. Under this assumption we can estimate β using maximum likelihood. The log likelihood is written as

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{X}; \beta) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}'\mathbf{y} - 2\beta'\mathbf{X}'\mathbf{y} + \beta'\mathbf{X}'\mathbf{X}\beta) \\ &= C - \frac{1}{2\sigma^2} [\beta'\mathbf{X}'\mathbf{X}\beta - 2\beta'\mathbf{X}'\mathbf{y}]\end{aligned}$$

To find the maximum with respect to β we have to solve the following equation:

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}; \beta)}{\partial \beta} = -\frac{1}{2\sigma^2} [2\mathbf{X}'\mathbf{X}\beta - 2\mathbf{X}'\mathbf{y}] = \mathbf{0}$$

If the $p \times p$ -matrix $\mathbf{X}'\mathbf{X}$ is invertible, then the maximum likelihood estimate (MLE) of β is given by

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

Otherwise the MLE does not exist. Why would this happen?

1. Colinearity: some variables are too correlated with one another.
2. High dimensionality: the number of covariates (p) exceeds the sample size (n).

(Part III)
Matrix decomposition and inverse

The case of square matrices

Once upon a time, a famous mathematician named Karl Weierstrass proposed the **Spectral Decomposition Theorem** which says in substance that **a symmetric matrix X filled with real numbers can be decomposed as $A = PDP^{-1}$** , where P verifies that $P'P = I$ (i.e. $P^{-1} = P'$) and D is a diagonal matrix.

The case of square matrices

Once upon a time, a famous mathematician named Karl Weierstrass proposed the **Spectral Decomposition Theorem** which says in substance that a symmetric matrix \mathbf{X} filled with real numbers can be decomposed as $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$.

where \mathbf{P} verifies that $\mathbf{P}'\mathbf{P} = \mathbf{I}$ (i.e. $\mathbf{P}^{-1} = \mathbf{P}'$) and \mathbf{D} is a diagonal matrix.

Few more definition/operations

A symmetric matrix \mathbf{A} is a square matrix such as $\mathbf{A}' = \mathbf{A}$.

E.g:

$$\mathbf{X} = \begin{bmatrix} 1 & -3 \\ -3 & 4 \end{bmatrix} \text{ then } \mathbf{X}' = \mathbf{X}$$

One trick we did not mention earlier is that $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ and another a bit similar is $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$.

The case of square matrices

Once upon a time, a famous mathematician named Karl Weierstrass proposed the **Spectral Decomposition Theorem** which says in substance that a symmetric matrix \mathbf{X} filled with real numbers can be decomposed as $\mathbf{A} = \mathbf{PDP}^{-1}$.

where \mathbf{P} verifies that $\mathbf{P}'\mathbf{P} = \mathbf{I}$ (i.e. $\mathbf{P}^{-1} = \mathbf{P}'$) and \mathbf{D} is a diagonal matrix.

Few more definition/operations

A symmetric matrix \mathbf{A} is a square matrix such as $\mathbf{A}' = \mathbf{A}$.

E.g:

$$\mathbf{X} = \begin{bmatrix} 1 & -3 \\ -3 & 4 \end{bmatrix} \text{ then } \mathbf{X}' = \mathbf{X}$$

One trick we did not mention earlier is that $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ and another a bit similar is $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$.

The case of square matrices

Once upon a time, a famous mathematician named Karl Weierstrass proposed the **Spectral Decomposition Theorem** which says in substance that **a symmetric matrix X filled with real numbers can be factorized as $A = PDP^{-1}$.**

where P verifies that $P'P = I$ (i.e. $P^{-1} = P'$) and D is a diagonal matrix.

We can remark that $X'X$ is symmetric because $(X'X)' = X'(X')' = X'X$ and therefore can be decomposed as $X'X = PDP^{-1}$.

The columns of P are called **eigenvectors** of X and the values of the diagonal of D are called **eigenvalues** of X .

Cholesky decomposition

We now assume that matrix \mathbf{A} is symmetric and has a determinant strictly positive. Therefore, we can find a triangular matrix \mathbf{L} such as $\mathbf{A} = \mathbf{L}\mathbf{L}'$. This is also known as the LU decomposition for Lower-Upper.

$$\mathbf{L} = \begin{bmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & \vdots & \ddots & \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{bmatrix}$$

Cholesky decomposition

We now assume that matrix \mathbf{A} is symmetric and has a determinant strictly positive. Therefore, we can find a triangular matrix \mathbf{L} such as $\mathbf{A} = \mathbf{L}\mathbf{L}'$. This is also known as the LU decomposition for Lower-Upper.

$$\mathbf{L} = \begin{bmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & \vdots & \ddots & \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{bmatrix}$$

This decomposition is also interesting because calculating the inverse of a relatively easy.

Singular Vector Decomposition

More generally, any $n \times p$ -matrix \mathbf{X} can be decomposed as $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}'$ where \mathbf{S} is $n \times p$ diagonal matrix and \mathbf{U} and \mathbf{V} are respectively $n \times n$ and $p \times p$ such as $\mathbf{U}'\mathbf{U} = \mathbf{I}_n$ and $\mathbf{V}'\mathbf{V} = \mathbf{I}_n$

$$\mathbf{S} = \begin{bmatrix} s_{11} & & & & \\ 0 & s_{22} & & & \\ \vdots & \vdots & \ddots & & \\ 0 & 0 & \cdots & s_{np} & \\ 0 & 0 & \cdots & 0 & \\ \vdots & \vdots & \vdots & 0 & \\ 0 & 0 & \cdots & 0 & \end{bmatrix}$$

a rectangular diagonal matrix ($n < p$)

QR decomposition

Another useful decomposition of a $n \times p$ matrix is $\mathbf{A} = \mathbf{QR}$ where \mathbf{Q} is an orthonormal matrix, i.e. $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_n$ and \mathbf{R} is an upper triangular matrix.

QR decomposition

Another useful decomposition of a $n \times p$ matrix is $\mathbf{A} = \mathbf{QR}$ where \mathbf{Q} is an orthonormal matrix, i.e. $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_n$ and \mathbf{R} is an upper triangular matrix.

This decomposition has a lower complexity than the SVD or the Cholesky decomposition.

How can these matrices decompositions be obtained?

There are known and efficient algorithms implemented in most statistical suites.

1. in R: **svd()**, **chol()** or **solve()**
2. in C++ (Eigen library): **SelfAdjointEigenSolver()**, **llt()** or **inverse()**

How can these matrices decompositions be obtained?

There are known and efficient algorithms implemented in most statistical suites.

1. in R: **svd()**, **chol()** or **solve()**
2. in C++ (Eigen library): **SelfAdjointEigenSolver()**, **llt()** or **inverse()**

...but you should know these involve **intensive** calculations. For a $n \times p$

1. SVD complexity : $\mathcal{O}(np^2 - p^3/3)$ operations
2. Cholesky Decomposition complexity: $\mathcal{O}(np^2 + p^3)$ operations.

Why do we care if this is already implemented somewhere?

Well...

1. Existing packages can be bugged, so if you know how it works you can fix it!

Why do we care if this is already implemented somewhere?

Well...

1. Existing packages can be bugged, so if you know how it works you can fix it!
2. Improve existing packages or implement faster algorithm for your specific problem

Why do we care if this is already implemented somewhere?

Well...

1. Existing packages can be bugged, so if you know how it works you can fix it!
2. Improve existing packages or implement faster algorithm for your specific problem
3. Understand when an approximation is needed

Some ongoing research in this area

The main driving force in this area is the improvement of computational algorithms to meet the challenges of using large dataset. There are researches on

1. Improving performances for sparse matrices (filled with a lot of zeros) decomposition.
2. Low rank approximations for SVD