

Matrix decompositions for regression analysis

Douglas Bates

2010-09-07 Tue

Contents

1	Matrix decompositions	1
1.1	Orthogonal matrices	1
1.1.1	Preserving lengths	1
1.1.2	The determinant of an orthogonal matrix	1
1.2	The QR decomposition	1
2	Least squares estimation	2
2.1	Spherical normal distributions and least squares estimators	2
2.2	Comparison to the usual text-book formulas	3
2.3	R functions related to the QR decomposition	3
3	Related matrix decompositions	4
3.1	The Cholesky decomposition	4
3.2	Evaluation of the Cholesky decomposition	4
3.3	The singular value decomposition	5

1 Matrix decompositions

1.1 Orthogonal matrices

An *orthogonal* $n \times n$ matrix, \mathbf{Q} has the property that its transpose is its inverse. That is,

$$\mathbf{Q}'\mathbf{Q} = \mathbf{Q}\mathbf{Q}' = \mathbf{I}_n$$

That is, the columns of \mathbf{Q} must be orthogonal to each other and all have unit length, and the same for the rows.

Two consequences of this property are that the transformation from \mathbb{R}^n to \mathbb{R}^n determined by \mathbf{Q} or \mathbf{Q}' preserves lengths and that the determinant of \mathbf{Q} is ± 1 .

1.1.1 Preserving lengths

For any $\mathbf{x} \in \mathbb{R}^n$

$$\|\mathbf{Q}\mathbf{x}\|^2 = (\mathbf{Q}\mathbf{x})'\mathbf{Q}\mathbf{x} = \mathbf{x}'\mathbf{Q}'\mathbf{Q}\mathbf{x} = \mathbf{x}'\mathbf{x} = \|\mathbf{x}\|^2$$

Thus the transformation determined by \mathbf{Q} or by \mathbf{Q}' must be a *rigid* transformation, composed of reflections or rotations.

1.1.2 The determinant of an orthogonal matrix

The determinants of $n \times n$ matrices \mathbf{A} and \mathbf{B} , written $|\mathbf{A}|$ and $|\mathbf{B}|$, are scalars with the property that

$$|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$$

Furthermore, the determinant of a *diagonal* matrix or a *triangular* matrix is simply the product of its diagonal elements. Also, $|\mathbf{A}| = |\mathbf{A}'|$.

From these properties we can derive

$$1 = |\mathbf{I}_n| = |\mathbf{Q}\mathbf{Q}'| = |\mathbf{Q}||\mathbf{Q}'| = |\mathbf{Q}|^2 \quad \Rightarrow \quad |\mathbf{Q}| = \pm 1.$$

As described on the Wikipedia page, the determinant of \mathbf{Q} is the volume of the parallelepiped spanned by the columns (or by the rows) of \mathbf{Q} . We know that the columns of \mathbf{Q} are *orthonormal* hence they span a unit volume. The sign indicates whether the transformation preserves orientation. In two dimensions a rotation preserves orientation and a reflection reverses the orientation.

1.2 The QR decomposition

An $n \times p$ matrix \mathbf{X} has a QR decomposition consisting of an orthogonal $n \times n$ matrix \mathbf{Q} and an $n \times p$ matrix \mathbf{R} that is zero below the main diagonal. In the cases we consider in regression analysis the matrix \mathbf{X} is the model matrix where n is the number of observations, p is the number of coefficients (or *parameters*) in the model and $n \geq p$. Our basic model is that the observed responses \mathbf{y} are the realization of a vector-valued random variable \mathcal{Y} with distribution

$$\mathcal{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}_0, \sigma^2\mathbf{I}_n)$$

for some unknown value $\boldsymbol{\beta}_0$ of the coefficient vector $\boldsymbol{\beta}$.

The QR decomposition of the model matrix \mathbf{X} is written

$$\mathbf{X} = \mathbf{Q} \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix} = [\mathbf{Q}_1 \quad \mathbf{Q}_2] \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_1 \mathbf{R} \quad (1)$$

where \mathbf{R} is a $p \times p$ upper triangular matrix, \mathbf{Q}_1 is the first p columns of \mathbf{Q} and \mathbf{Q}_2 is the last $n - p$ columns of \mathbf{Q} .

That fact that matrices \mathbf{Q} and \mathbf{R} must exist is proved by construction. The matrix \mathbf{Q} is the product of p *Householder reflections* (see the Wikipedia page for the QR decomposition). The process is similar to the Gram-Schmidt orthogonalization process, but more flexible and numerically stable. If the diagonal elements of \mathbf{R} are all non-zero (in practice this means that none of them is very small in absolute value) then \mathbf{X} has *full column rank* and the columns of \mathbf{Q}_1 form an *orthonormal basis* of the *column span* of \mathbf{X} .

The implementation of the QR decomposition in R guarantees that any near-zero elements on the diagonal of \mathbf{R} are rearranged by column permutation to occur in the trailing columns. That is, if the rank of \mathbf{X} is $k < p$ then the first k columns of \mathbf{Q} form an orthogonal basis for the column span of $\mathbf{X}\mathbf{P}$ where \mathbf{P} is a $p \times p$ permutation matrix.

Our text often mentions rank-deficient cases where $\text{rank}(\mathbf{X}) = k < p$. In practice these occur rarely because the process of building the model matrix in R involves a considerable amount of analysis of the model formula to remove the most common cases of rank deficiency. Nevertheless, rank deficiency can occur and is detected and handled in the `lm` function in R.

2 Least squares estimation

The reason that we are interested in a QR decomposition of \mathbf{X} is because the probability model is linked to geometric properties of the response \mathbf{y} and the column span of \mathbf{X} .

2.1 Spherical normal distributions and least squares estimators

A distribution of the form

$$\mathcal{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}_0, \sigma^2 \mathbf{I}_n) \quad (2)$$

is called a *spherical normal* distribution because the contours of constant probability density are spheres centered at $\mathbf{X}\boldsymbol{\beta}_0$. In other words, the probability density function of \mathcal{Y} , evaluated at \mathbf{y} is determined by the distance, $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0\|^2$ with larger distances corresponding to lower densities.

The *maximum likelihood estimate*, $\hat{\boldsymbol{\beta}}$ of the coefficient vector, $\boldsymbol{\beta}$, is

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \quad (3)$$

These estimates are also called the *least squares estimates* of $\boldsymbol{\beta}$.

Because multiplication by an orthogonal matrix like \mathbf{Q}' preserves lengths we can write

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \arg \min_{\boldsymbol{\beta}} \|\mathbf{Q}'\mathbf{y} - \mathbf{Q}'\mathbf{X}\boldsymbol{\beta}\|^2 = \arg \min_{\boldsymbol{\beta}} \|\mathbf{c}_1 - \mathbf{R}\boldsymbol{\beta}\|^2 + \|\mathbf{c}_2\|^2 \quad (4)$$

where $\mathbf{c}_1 = \mathbf{Q}'_1\mathbf{y}$ is the first p elements of $\mathbf{Q}\mathbf{y}$ and $\mathbf{c}_2 = \mathbf{Q}'_2\mathbf{y}$ is the last $n - p$ elements. If $\text{rank}(\mathbf{X}) = p$ then $\text{rank}(\mathbf{R}) = p$ and \mathbf{R}^{-1} exists and we can write $\hat{\boldsymbol{\beta}} = \mathbf{R}^{-1}\mathbf{c}_1$ (although you don't actually calculate \mathbf{R}^{-1} to solve the triangular linear system $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{c}_1$ for $\hat{\boldsymbol{\beta}}$).

In a model fit by the `lm` or `aov` functions in R there is a component `effects` which is $\mathbf{Q}'\mathbf{y}$. The component `qr` is a condensed form of the QR decomposition of the model matrix \mathbf{X} . The matrix \mathbf{R} is embedded in there but the matrix \mathbf{Q} is a virtual matrix represented as a product of Householder reflections and not usually evaluated explicitly.

2.2 Comparison to the usual text-book formulas

Most text books state that the least squares estimates are

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (5)$$

giving the impression that $\hat{\boldsymbol{\beta}}$ is calculated this way. It isn't.

If you substitute $\mathbf{X} = \mathbf{Q}_1\mathbf{R}$ in 5 you get

$$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{R}'\mathbf{R})^{-1}\mathbf{R}'\mathbf{Q}_1'\mathbf{y} = \mathbf{R}^{-1}(\mathbf{R}')^{-1}\mathbf{R}'\mathbf{Q}_1'\mathbf{y} = \mathbf{R}^{-1}\mathbf{Q}_1'\mathbf{y},$$

our previous result.

Whenever you see $\mathbf{X}'\mathbf{X}$ in a formula you should mentally replace it by $\mathbf{R}'\mathbf{R}$ and similarly replace $(\mathbf{X}'\mathbf{X})^{-1}$ by $\mathbf{R}^{-1}(\mathbf{R}')^{-1}$ then see if you can simplify the result.

For example, the variance of the least squares estimator $\hat{\boldsymbol{\beta}}$ is

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2\mathbf{R}^{-1}(\mathbf{R}')^{-1}$$

The R function `chol2inv` calculates $\mathbf{R}^{-1}(\mathbf{R}')^{-1}$ directly from \mathbf{R} (not a big deal in most cases but when p is very large it should be faster and more accurate than evaluating \mathbf{R}^{-1} explicitly). The determinant of $\mathbf{X}'\mathbf{X}$ is

$$|\mathbf{X}'\mathbf{X}| = |\mathbf{R}'\mathbf{R}| = |\mathbf{R}|^2 = \left(\prod_{i=1}^p r_{i,i}\right)^2$$

The fitted values $\hat{\mathbf{y}}$ are $\mathbf{Q}_1\mathbf{Q}_1'\mathbf{y}$ and thus the *hat matrix* (which puts a “hat” on \mathbf{y} by transforming it to $\hat{\mathbf{y}}$) is the $n \times n$ matrix $\mathbf{Q}_1\mathbf{Q}_1'$. Often we are interested in the diagonal elements of the hat matrix, which are the sums of the squares of rows of \mathbf{Q}_1 . (In practice you don't want to calculate the entire $n \times n$ hat matrix just to get the diagonal elements when n could be very large.)

The residuals, $\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}}$, are calculated as $\hat{\mathbf{e}} = \mathbf{Q}_2\mathbf{Q}_2'\mathbf{y}$.

The matrices $\mathbf{Q}_1\mathbf{Q}_1'$ and $\mathbf{Q}_2\mathbf{Q}_2'$ are *projection matrices*, which means that they are symmetric and *idempotent*. (A square matrix \mathbf{A} is idempotent if $\mathbf{A}\mathbf{A} = \mathbf{A}$.) When $\text{rank}(\mathbf{X}) = p$, the hat matrix $\mathbf{Q}_1\mathbf{Q}_1'$ projects any vector in \mathbb{R}^n onto the column span of \mathbf{X} . The other projection, $\mathbf{Q}_2\mathbf{Q}_2'$, is onto the subspace orthogonal to the column span of \mathbf{X} (see the figure on the front cover of the text).

2.3 R functions related to the QR decomposition

As mentioned above, every time you fit a linear model with `lm` or `aov` or `lm.fit`, the returned object contains a `qr` component. This is a condensed form of the decomposition, only slightly larger than \mathbf{X} itself. Its class is “qr”.

There are several extractor functions for a “qr” object: `qr.R()`, `qr.Q()` and `qr.X()`, which regenerates the original matrix. By default `qr.Q()` returns the matrix called \mathbf{Q}_1 above with p

columns but you can specify the number of columns desired. Typical alternative choices are n or $\text{rank}(\mathbf{X})$.

The `$rank` component of a "qr" object is the computed rank of \mathbf{X} (and, hence, of \mathbf{R}). The `$pivot` component is the permutation applied to the columns. It will be `1:p` when $\text{rank}(\mathbf{X}) = p$ but when $\text{rank}(\mathbf{X}) < p$ it may be other than the identity permutation.

Several functions are applied to a "qr" object and a vector or matrix. These include `qr.coef()`, `qr.qy()`, `qr.qty()`, `qr.resid()` and `qr.fitted()`. The `qr.qy()` and `qr.qty()` functions multiply an n -vector or an $n \times m$ matrix by \mathbf{Q} or \mathbf{Q}' without ever forming \mathbf{Q} . Similarly, `qr.fitted()` creates $\mathbf{Q}_1\mathbf{Q}'_1\mathbf{x}$ and `qr.resid()` creates $\mathbf{Q}_2\mathbf{Q}'_2\mathbf{x}$ without forming \mathbf{Q} .

The `is.qr()` function tests an object to determine if it is of class "qr".

3 Related matrix decompositions

3.1 The Cholesky decomposition

The Cholesky decomposition of a positive definite symmetric matrix, which means a $p \times p$ symmetric matrix \mathbf{A} such that $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$ for all non-zero $\mathbf{x} \in \mathbb{R}^p$ is of the form

$$\mathbf{A} = \mathbf{R}'\mathbf{R} = \mathbf{L}\mathbf{L}'$$

where \mathbf{R} is an upper triangular $p \times p$ matrix and $\mathbf{L} = \mathbf{R}'$ is lower triangular. The two forms are the same decomposition: it is just a matter of whether you want \mathbf{R} , the factor on the right, or \mathbf{L} , the factor on the left. Generally statisticians write the decomposition as $\mathbf{R}'\mathbf{R}$.

The decomposition is only determined up to changes in sign of the rows of \mathbf{R} (or, equivalently, the columns of \mathbf{L}). For definiteness we require positive diagonal elements in \mathbf{R} .

When $\text{rank}(\mathbf{X}) = p$ the Cholesky decomposition \mathbf{R} of $\mathbf{X}'\mathbf{X}$ is the equal to the matrix \mathbf{R} from the QR decomposition up to changes in sign of rows. The matrix $\mathbf{X}'\mathbf{X}$ matrix is obviously symmetric and it is positive definite because

$$\mathbf{x}'(\mathbf{X}'\mathbf{X})\mathbf{x} = \mathbf{x}'(\mathbf{R}'\mathbf{R})\mathbf{x} = \|\mathbf{R}\mathbf{x}\|^2 \geq 0$$

with equality only when $\mathbf{R}\mathbf{x} = \mathbf{0}$, which, when $\text{rank}(\mathbf{R}) = p$, implies that $\mathbf{x} = \mathbf{0}$.

3.2 Evaluation of the Cholesky decomposition

The R function `chol()` evaluates the Cholesky decomposition. As mentioned above `chol2inv()` creates $(\mathbf{X}'\mathbf{X})^{-1}$ directly from the Cholesky decomposition of $\mathbf{X}'\mathbf{X}$.

Generally the QR decomposition is preferred to the Cholesky decomposition for least squares problems because there is a certain loss of precision when forming $\mathbf{X}'\mathbf{X}$. However, when n is very large you may want to build up $\mathbf{X}'\mathbf{X}$ using blocks of rows. Also, if \mathbf{X} is *sparse* it is an advantage to use sparse matrix techniques to evaluate and store the Cholesky decomposition.

The `Matrix` package for R provides even more capabilities related to the Cholesky decomposition, especially for sparse matrices.

For everything we will do in Statistics 849 the QR decomposition should be the method of choice.

3.3 The singular value decomposition

Another decomposition related to orthogonal matrices is the singular value decomposition (or SVD) in which the matrix \mathbf{X} is reduced to a diagonal form

$$\mathbf{X} = \mathbf{U}_1 \mathbf{D} \mathbf{V}' = \mathbf{U} \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix} \mathbf{V}'$$

where \mathbf{U} is an $n \times n$ orthogonal matrix, \mathbf{D} is a $p \times p$ diagonal matrix with non-negative diagonal elements (which are called the *singular values* of \mathbf{X}) and \mathbf{V} is a $p \times p$ orthogonal matrix. As for \mathbf{Q} and \mathbf{Q}_1 , \mathbf{U}_1 consists of the first p columns of \mathbf{U} . For definiteness we order the diagonal elements of \mathbf{D} , which must be non-negative, in decreasing order.

The singular value decomposition of \mathbf{X} is related to the eigendecomposition or spectral decomposition of $\mathbf{X}'\mathbf{X}$ because

$$\mathbf{X}'\mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}'_1 \mathbf{U}_1 \mathbf{D} \mathbf{V}' = \mathbf{V} \mathbf{D}^2 \mathbf{V}'$$

implying that the eigenvalues of $\mathbf{X}'\mathbf{X}$ are the squares of the singular values of \mathbf{X} and the right singular vectors, which are the columns of \mathbf{V} are also the eigenvectors of $\mathbf{X}'\mathbf{X}$

Calculation of the SVD is an iterative (as opposed to a direct) computation and potentially more computing intensive than the QR decomposition, although modern methods for evaluating the SVD are very good indeed.

Symbolically we can write the least squares solution in the full-rank case as

$$\hat{\boldsymbol{\beta}} = \mathbf{V} \mathbf{D}^{-1} \mathbf{U}'_1 \mathbf{y}$$

where \mathbf{D}^{-1} is a diagonal matrix whose diagonal elements are the inverses of the diagonal elements of \mathbf{D} .

The pseudo-inverse or *generalized inverse* of \mathbf{X} , written \mathbf{X}^- is calculated from the pseudo-inverse of the diagonal matrix, \mathbf{D} . In theory the diagonal elements of \mathbf{D}^- are $1/d_{i,i}$ when $d_{i,i} \neq 0$ and 0 when $d_{i,i} = 0$. However, we can't count on $d_{i,i}$ being 0 even when, in theory, it should be. You need to decide when the singular values are close to zero, which is actually a very difficult problem. At best we can use some heuristics, based on the ratio of $d_{i,i}/d_{1,1}$ to decide when a diagonal element is "effectively zero".

The use of the pseudo-inverse seems to be a convenient way to handle rank-deficient \mathbf{X} matrices but, as mentioned above, the best way to handle rank-deficient \mathbf{X} matrices is not to produce them in the first place.