

where one first computes $y_n = b_n/U_{nn}$, then substitutes *back* into the previous equation to solve for y_{n-1} , and repeats through y_1 .

Here is a Csharp implementation of the in-place¹ back-substitution:

```
static void backsub(matrix U, vector c){
    for(int i=c.size-1; i>=0; i--){
        double sum=0;
        for(int k=i+1; k<c.size; k++) sum+=U[i,k]*c[k];
        c[i]=(c[i]-sum)/U[i,i]; } }
```

For a lower triangular system $\mathbf{L}\mathbf{y} = \mathbf{c}$ the equivalent procedure is called *forward-substitution*,

$$y_i = \frac{1}{L_{ii}} \left(c_i - \sum_{k=1}^{i-1} L_{ik}y_k \right), \quad i = 1, 2, \dots, n. \quad (5)$$

1.3 Reduction to triangular form

Popular algorithms for reducing a square system of linear equations to a triangular form are *LU-decomposition* and *QR-decomposition*.

1.3.1 QR-decomposition

QR-decomposition is a factorization of a matrix into a product of an orthogonal matrix \mathbf{Q} , such that $\mathbf{Q}^T\mathbf{Q} = \mathbf{1}$, where T denotes transposition, and a right triangular matrix \mathbf{R} , such that

$$\mathbf{A} = \mathbf{Q}\mathbf{R}. \quad (6)$$

QR-decomposition can be used to convert (by multiplying with \mathbf{Q}^T from the left) a linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ into the triangular form,

$$\mathbf{R}\mathbf{x} = \mathbf{Q}^T\mathbf{b}, \quad (7)$$

which can be solved directly by back-substitution.

QR-decomposition can also be performed on non-square matrices with few long columns. Generally speaking a rectangular $n \times m$ matrix \mathbf{A} can be represented as a product, $\mathbf{A} = \mathbf{Q}\mathbf{R}$, of an orthogonal $n \times m$ matrix \mathbf{Q} , $\mathbf{Q}^T\mathbf{Q} = \mathbf{1}$, and a right-triangular $m \times m$ matrix \mathbf{R} .

QR-decomposition of a matrix can be computed using several methods, such as Gram-Schmidt orthogonalization, Householder transformation [2], or Givens rotation [1].

Gram-Schmidt orthogonalization *Gram-Schmidt orthogonalization* is an algorithm for orthogonalization of a set of vectors in a given inner product space. It takes a linearly independent set of vectors $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_m\}$ and generates an orthogonal set $\mathbf{Q} = \{\mathbf{q}_1, \dots, \mathbf{q}_m\}$ which spans the same subspace as \mathbf{A} . The algorithm is given as

¹here *in-place* means the right-hand side \mathbf{c} is replaced by the solution \mathbf{y} .

```

for  $i = 1$  to  $m$  :
     $\mathbf{q}_i \leftarrow \mathbf{a}_i / \|\mathbf{a}_i\|$ 
    for  $j = i + 1$  to  $m$  :  $\mathbf{a}_j \leftarrow \mathbf{a}_j - \langle \mathbf{q}_i | \mathbf{a}_j \rangle \mathbf{q}_i$ 

```

where $\langle \mathbf{a} | \mathbf{b} \rangle$ is the inner product of two vectors, and $\|\mathbf{a}\| \doteq \sqrt{\langle \mathbf{a} | \mathbf{a} \rangle}$ is the vector's norm. This variant of the algorithm, where all remaining vectors \mathbf{a}_j are made orthogonal to \mathbf{q}_i as soon as the latter is calculated, is considered to be numerically stable and is referred to as *stabilized* or *modified*.

Stabilized Gram-Schmidt orthogonalization can be used to compute QR-decomposition of a matrix A by orthogonalization of its column-vectors \mathbf{a}_i with the inner product

$$\langle \mathbf{a} | \mathbf{b} \rangle = \mathbf{a}^\top \mathbf{b} \equiv \sum_{k=1}^n (\mathbf{a})_k (\mathbf{b})_k , \quad (8)$$

where n is the length of column-vectors \mathbf{a} and \mathbf{b} , and $(\mathbf{a})_k$ is the k th element of the column-vector,

```

for  $i = 1$  to  $m$  :
     $R_{ii} = \sqrt{\mathbf{a}_i^\top \mathbf{a}_i}$  ;  $\mathbf{q}_i = \mathbf{a}_i / R_{ii}$ 
    for  $j = i + 1$  to  $m$  :
         $R_{ij} = \mathbf{q}_i^\top \mathbf{a}_j$  ;  $\mathbf{a}_j = \mathbf{a}_j - \mathbf{q}_i R_{ij}$  .

```

After orthogonalization the matrices $Q = \{\mathbf{q}_1 \dots \mathbf{q}_m\}$ and R are the sought orthogonal and right-triangular factors of matrix A . A Csharp implementation might look like this,

```

matrix Q=A.copy(), R=new matrix(m,m);
for (int i=0;i<m;i++){
    R[i,i]=Q[i].norm(); /* Q[i] points to the i-th column */
    Q[i]/=R[i,i];
    for (int j=i+1;j<m;j++){
        R[i,j]=Q[i].dot(Q[j]);
        Q[j]-=Q[i]*R[i,j]; } }

```

The factorization is unique under requirement that the diagonal elements of R are positive. For a $n \times m$ matrix the complexity of the algorithm is $O(m^2n)$.

Gram-Schmidt decomposition with column pivoting Pivoted decomposition differs from the ordinary Gram-Schmidt in that at each iteration it takes the largest of the remaining columns and thus introduces the permutation matrix P ,

$$AP = QR , \quad (9)$$

that is (generally) chosen so that the diagonal elements of the R -matrix are decreasing,

$$|R_{11}| \geq |R_{22}| \geq \dots \geq |R_{mm}| . \quad (10)$$

Pivoted QR-decomposition can be used when matrix A is rank deficient or its rank is in doubt. With exact arithmetics if $\text{rank}(A) = k$ then the sub-matrix of R

with rows and columns from $k + 1$ to m would be zero. Numerical determination of rank requires a criterion for deciding when a small diagonal element of R should be treated as zero – a practical choice that depends on both the matrix and the application.

Householder transformation A square matrix H of the form

$$H = 1 - \frac{2}{\mathbf{u}^T \mathbf{u}} \mathbf{u} \mathbf{u}^T \quad (11)$$

is called *Householder matrix*, where the vector \mathbf{u} is called a *Householder vector*. Householder matrices are symmetric and orthogonal,

$$H^T = H, \quad H^T H = 1. \quad (12)$$

The transformation induced by the Householder matrix on a given vector \mathbf{a} ,

$$\mathbf{a} \rightarrow H\mathbf{a}, \quad (13)$$

is called a *Householder transformation* or *Householder reflection*. The transformation changes the sign of the affected vector's component in the \mathbf{u} direction, or, in other words, makes a reflection of the vector about the hyper-plane perpendicular to \mathbf{u} , hence the name.

Householder transformation can be used to zero selected components of a given vector \mathbf{a} . For example, one can zero all components but the first one, such that

$$H\mathbf{a} = \gamma \mathbf{e}_1, \quad (14)$$

where γ is a number and \mathbf{e}_1 is the unit vector in the first direction. The factor γ can be easily calculated,

$$\|\mathbf{a}\|^2 \doteq \mathbf{a}^T \mathbf{a} = \mathbf{a}^T H^T H \mathbf{a} = (\gamma \mathbf{e}_1)^T (\gamma \mathbf{e}_1) = \gamma^2, \quad (15)$$

$$\Rightarrow \gamma = \pm \|\mathbf{a}\|. \quad (16)$$

To find the Householder vector, we notice that

$$\mathbf{a} = H^T H \mathbf{a} = H^T \gamma \mathbf{e}_1 = \gamma \mathbf{e}_1 - \frac{2(\mathbf{u})_1}{\mathbf{u}^T \mathbf{u}} \mathbf{u}, \quad (17)$$

$$\Rightarrow \frac{2(\mathbf{u})_1}{\mathbf{u}^T \mathbf{u}} \mathbf{u} = \gamma \mathbf{e}_1 - \mathbf{a}, \quad (18)$$

where $(\mathbf{u})_1$ is the first component of the vector \mathbf{u} . One usually chooses $(\mathbf{u})_1 = 1$ (for the sake of the possibility to store the other components of the Householder vector in the zeroed elements of the vector \mathbf{a}) and stores the factor

$$\frac{2}{\mathbf{u}^T \mathbf{u}} \equiv \tau \quad (19)$$

separately. With this convention one readily finds τ from the first component of equation (18),

$$\tau = \gamma - (\mathbf{a})_1. \quad (20)$$

where $(\mathbf{a})_1$ is the first element of the vector \mathbf{a} . For the sake of numerical stability the sign of γ has to be chosen opposite to the sign of $(\mathbf{a})_1$,

$$\gamma = -\text{sign}((\mathbf{a})_1) \|\mathbf{a}\|. \quad (21)$$

Finally, the Householder reflection, which zeroes all component of a vector \mathbf{a} but the first, is given as

$$\mathbf{H} = 1 - \tau \mathbf{u} \mathbf{u}^T, \quad \tau = -\text{sign}((\mathbf{a})_1) \|\mathbf{a}\| - (\mathbf{a})_1, \quad (\mathbf{u})_1 = 1, \quad (\mathbf{u})_{i>1} = -\frac{1}{\tau} (\mathbf{a})_i. \quad (22)$$

Now, a QR-decomposition of an $n \times n$ matrix \mathbf{A} by Householder transformations can be performed in the following way:

1. Build the size- n Householder vector \mathbf{u}_1 which zeroes the sub-diagonal elements of the first column of matrix \mathbf{A} , such that

$$\mathbf{H}_1 \mathbf{A} = \left[\begin{array}{c|ccc} \star & \star & \dots & \star \\ \hline 0 & & & \\ \vdots & & & \\ 0 & & \mathbf{A}_1 & \end{array} \right], \quad (23)$$

where $\mathbf{H}_1 = 1 - \tau_1 \mathbf{u}_1 \mathbf{u}_1^T$ and where \star denotes (generally) non-zero matrix elements. In practice one does not build the matrix \mathbf{H}_1 explicitly, but rather calculates the matrix $\mathbf{H}_1 \mathbf{A}$ in-place, consecutively applying the Householder reflection to columns the matrix \mathbf{A} , thus avoiding computationally expensive matrix-matrix operations. The zeroed sub-diagonal elements of the first column of the matrix \mathbf{A} can be used to store the elements of the Householder vector \mathbf{u}_1 while the factor τ_1 has to be stored separately in a special array. This is the storage scheme used by LAPACK and GSL.

2. Similarly, build the size- $(n-1)$ Householder vector \mathbf{u}_2 which zeroes the sub-diagonal elements of the first column of matrix \mathbf{A}_1 from eq. (23). With the transformation matrix \mathbf{H}_2 defined as

$$\mathbf{H}_2 = \left[\begin{array}{c|ccc} 1 & 0 & \dots & 0 \\ \hline 0 & & & \\ \vdots & & & \\ 0 & & 1 - \tau_2 \mathbf{u}_2 \mathbf{u}_2^T & \end{array} \right]. \quad (24)$$

the two transformations together zero the sub-diagonal elements of the two

first columns of matrix A,

$$\mathbf{H}_2\mathbf{H}_1\mathbf{A} = \left[\begin{array}{cc|ccc} \star & \star & \star & \cdots & \star \\ 0 & \star & \star & \cdots & \star \\ \hline 0 & 0 & & & \\ \vdots & \vdots & & & \\ 0 & 0 & & & \mathbf{A}_3 \end{array} \right], \quad (25)$$

3. Repeating the process zero the sub-diagonal elements of the remaining columns. For column k the corresponding Householder matrix is

$$\mathbf{H}_k = \left[\begin{array}{c|c} \mathbf{I}_{k-1} & 0 \\ \hline 0 & 1 - \tau_k \mathbf{u}_k \mathbf{u}_k^\top \end{array} \right], \quad (26)$$

where \mathbf{I}_{k-1} is an identity matrix of size $k - 1$, \mathbf{u}_k is the size- $(n-k+1)$ Householder vector that zeroes the sub-diagonal elements of matrix \mathbf{A}_{k-1} from the previous step. The corresponding transformation step is

$$\mathbf{H}_k \dots \mathbf{H}_2\mathbf{H}_1\mathbf{A} = \left[\begin{array}{c|c} \mathbf{R}_k & \star \\ \hline 0 & \mathbf{A}_k \end{array} \right], \quad (27)$$

where \mathbf{R}_k is a size- k right-triangular matrix.

After $n - 1$ steps the matrix A will be transformed into a right triangular matrix,

$$\mathbf{H}_{n-1} \dots \mathbf{H}_2\mathbf{H}_1\mathbf{A} = \mathbf{R}. \quad (28)$$

4. Finally, introducing an orthogonal matrix $\mathbf{Q} = \mathbf{H}_1^\top \mathbf{H}_2^\top \dots \mathbf{H}_{n-1}^\top$ and multiplying eq. (28) by \mathbf{Q} from the left, we get the sought QR-decomposition,

$$\mathbf{A} = \mathbf{Q}\mathbf{R}. \quad (29)$$

In practice one does not build explicitly the \mathbf{Q} matrix but rather applies the successive Householder reflections stored during the decomposition.

Givens rotations A Givens rotation is a transformation in the form

$$\mathbf{A} \rightarrow \mathbf{G}(p, q, \theta)\mathbf{A}, \quad (30)$$

where \mathbf{A} is the object to be transformed—matrix of vector—and $\mathbf{G}(p, q, \theta)$ is the Givens rotation matrix (also known as Jacobi rotation matrix): an orthogonal

matrix in the form

$$G(p, q, \theta) = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & \cos \theta & \cdots & \sin \theta & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & -\sin \theta & \cdots & \cos \theta & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \begin{array}{l} \leftarrow \text{row } p \\ \\ \leftarrow \text{row } q \end{array} . \quad (31)$$

When a Givens rotation matrix $G(p, q, \theta)$ multiplies a vector \mathbf{x} , only elements x_p and x_q are affected. Considering only these two affected elements, the Givens rotation is given explicitly as

$$\begin{bmatrix} x'_p \\ x'_q \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_p \\ x_q \end{bmatrix} = \begin{bmatrix} x_p \cos \theta + x_q \sin \theta \\ -x_p \sin \theta + x_q \cos \theta \end{bmatrix} . \quad (32)$$

Apparently the rotation can zero the element x'_q , if the angle θ is chosen as

$$\tan \theta = \frac{x_q}{x_p} \Rightarrow \theta = \text{atan2}(x_q, x_p) . \quad (33)$$

A sequence of Givens rotations,

$$G = \prod_{n \geq q > p=1}^m G(p, q, \theta_{qp}) , \quad (34)$$

(where $n \times m$ is the dimension of the matrix A) can zero all elements of a matrix below the main diagonal if the angles θ_{qp} are chosen to zero the elements with indices q, p of the partially transformed matrix just before applying the matrix $G(p, q, \theta_{qp})$. The resulting matrix is obviously the R-matrix of the sought QR-decomposition of the matrix A where $G = Q^T$.

In practice one does not explicitly builds the G matrix but rather stores the θ angles in the places of the corresponding zeroed elements of the original matrix:

```

for (int p=0;p<A.size2;p++){
  for (int q=p+1;q<A.size1;q++){
    double theta=Atan2(A[q,p],A[p,p]);
    double c=Cos(theta),s=Sin(theta);
    for (int k=q;k<A.size2;k++){
      double xp=A[p,k], xq=A[q,k];
      A[p,k]= xp*c+xq*s;
      A[q,k]=-xp*s+xq*c; }
    A[q,p]=theta; } }

```

When solving the linear system $\mathbf{Ax} = \mathbf{b}$ one transforms it into the equivalent triangular system $\mathbf{Rx} = \mathbf{Gb}$ where one calculates \mathbf{Gb} by successively applying the individual Givens rotations with the stored θ -angles:

```

for(int p=0;p<G.size2;p++){
  for(int q=p+1;q<G.size1;q++){
    double theta = G[q,p];
    double c=Cos(theta),s=Sin(theta);
    double bp=b[p], bq=b[q];
    b[p]=+bp*c+bq*s;
    b[q]=-bp*s+bq*c; } }

```

The triangular system $\mathbf{R}\mathbf{x} = \mathbf{G}\mathbf{b}$ is then solved by the ordinary back-substitution. If one needs to build the Q-matrix explicitly, one uses

$$Q_{ij} = \mathbf{e}_i^T \mathbf{Q} \mathbf{e}_j = \mathbf{e}_j^T \mathbf{Q}^T \mathbf{e}_i, \quad (35)$$

where \mathbf{e}_i is the unit vector in the direction i and where again one can use the successive rotations to calculate $\mathbf{Q}^T \mathbf{e}_i$,

Since each Givens rotation only affects two rows of the matrix it is possible to apply a set of rotations in parallel. Givens rotations are also more efficient on sparse matrices.

1.3.2 LU-decomposition

LU-decomposition is a factorization of a square matrix A into a product of a lower triangular matrix L and an upper triangular matrix U,

$$\mathbf{A} = \mathbf{L}\mathbf{U}. \quad (36)$$

The linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ after LU-decomposition of the matrix A becomes $\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{b}$ and can be solved by first solving $\mathbf{L}\mathbf{y} = \mathbf{b}$ for \mathbf{y} and then $\mathbf{U}\mathbf{x} = \mathbf{y}$ for \mathbf{x} with two runs of forward and backward substitutions.

If A is an $n \times n$ matrix, the condition (36) is a set of n^2 equations,

$$\sum_{k=1}^n L_{ik} U_{kj} = A_{ij} \quad |_{i,j=1 \dots n}, \quad (37)$$

for $n^2 + n$ unknown elements of the triangular matrices L and U. The decomposition is thus not unique.

Usually the decomposition is made unique by providing extra n conditions e.g. by the requirement that the elements of the main diagonal of the matrix L are equal one,

$$L_{ii} = 1, \quad i = 1 \dots n. \quad (38)$$

The system (37) with the extra conditions (38) can then be easily solved row after row using the *Doolittle's algorithm*,

```

for i = 1 ... n :
  Lii = 1
  for j = i ... n : Uij = Aij - ∑k<i LikUkj
  for j = i + 1 ... n : Lji =  $\frac{1}{U_{ii}}$  (Aji - ∑k<j LjkUki)

```


In a slightly different *Crout's algorithm* it is the matrix U that has unit diagonal elements,

```

for  $i = 1 \dots n$  :
   $U_{ii} = 1$ 
  for  $j = i \dots n$  :  $L_{ji} = A_{ji} - \sum_{k < i} L_{jk} U_{ki}$ 
  for  $j = i + 1 \dots n$  :  $U_{ij} = \frac{1}{L_{ii}} \left( A_{ji} - \sum_{k < j} L_{jk} U_{ki} \right)$ 

```

Without a proper ordering (permutations) in the matrix, the factorization may fail. For example, it is easy to verify that $A_{11} = L_{11}U_{11}$. If $A_{11} = 0$, then at least one of L_{11} and U_{11} has to be zero, which implies either L or U is singular, which is impossible if A is non-singular. This is however only a procedural problem. It can be removed by simply reordering the rows of A so that the first element of the permuted matrix is nonzero (or, even better, the largest in absolute value among all elements of the column below the diagonal). The same problem in subsequent factorization steps can be removed in a similar way. Such algorithm is referred to as *partial pivoting*. It requires an extra integer array to keep track of row permutations.

1.3.3 Cholesky decomposition

The Cholesky decomposition of a Hermitian positive-definite matrix A is a decomposition in the form

$$A = LL^\dagger, \quad (39)$$

where L is a lower triangular matrix with real and positive diagonal elements, and L^\dagger is the conjugate transpose of L .

For real symmetric positive-definite matrices the decomposition reads

$$A = LL^T, \quad (40)$$

where L is real.

The decomposition can be calculated using the following in-place algorithm,

$$L_{jj} = \sqrt{A_{jj} - \sum_{k=1}^{j-1} L_{jk}^2}, \quad L_{ij} = \frac{1}{L_{jj}} \left(A_{ij} - \sum_{k=1}^{j-1} L_{ik} L_{jk} \right) \Big|_{i>j}. \quad (41)$$

The expression under the square root is always positive if A is real and positive-definite.

When applicable, the Cholesky decomposition is about twice as efficient as LU-decomposition for solving systems of linear equations.

1.4 Determinant of a matrix

LU- and QR-decompositions allow $O(n^3)$ calculation of the determinant of a square matrix. Indeed, for the LU-decomposition,

$$\det A = \det LU = \det L \det U = \det U = \prod_{i=1}^n U_{ii}. \quad (42)$$

For the Gram-Schmidt QR-decomposition

$$\det A = \det QR = \det Q \det R. \quad (43)$$

Since Q is an orthogonal matrix $(\det Q)^2 = 1$,

$$|\det A| = |\det R| = \left| \prod_{i=1}^n R_{ii} \right|. \quad (44)$$

With Gram-Schmidt method one arbitrarily assigns positive sign to diagonal elements of the R -matrix thus removing from the R -matrix the memory of the original sign of the determinant.

However with Givens rotation method the determinant of the individual rotation matrix—and thus the determinant of the total rotation matrix—is equal one, therefore for a square matrix A the QR-decomposition $A = GR$ via Givens rotations allows calculation of the determinant with the correct sign,

$$\det A = \det R \equiv \prod_{i=1}^n R_{ii} \quad (45)$$

1.5 Matrix inverse

The inverse A^{-1} of a square $n \times n$ matrix A can be calculated by solving n linear equations

$$A\mathbf{x}_i = \mathbf{e}_i \Big|_{i=1, \dots, n}, \quad (46)$$

where \mathbf{e}_i is the unit-vector in the i -direction: a column where all elements are equal zero except for the element number i which is equal one. Thus the set of columns $\{\mathbf{e}_i\}_{i=1, \dots, n}$ form the identity matrix. The matrix made of columns \mathbf{x}_i is apparently the inverse of A .

References

- [1] Wallace Givens. Computation of plane unitary rotations transforming a general matrix to triangular form. *J. SIAM*, 6(1):26–50, 1958.
- [2] A.S. Householder. Unitary triangularization of a nonsymmetric matrix. *Journal of the ACM*, 5(4):339–342, 1958.