Math 6610: Analysis of Numerical Methods, I The QR decomposition

Department of Mathematics, University of Utah

Fall 2025

Resources: Trefethen and Bau 1997, Lectures 20, 21, 23

Atkinson 1989, Chapter 1

Salgado and Wise 2022, Chapter 3

The main goal of orthogonalization:

Given $\{a_j\}_{j\in[n]}\subset\mathbb{C}^m$ with $n\leqslant m$, compute $\{q_j\}_{j\in[n]}$ such that:

$$\langle \boldsymbol{q}_{j}, \boldsymbol{q}_{k} \rangle = \delta_{j,k},$$

$$\mathrm{span}\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_n\}=\mathrm{span}\{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_n\}$$

The main goal of orthogonalization:

Given $\{a_j\}_{j\in[n]}\subset\mathbb{C}^m$ with $n\leqslant m$, compute $\{q_j\}_{j\in[n]}$ such that:

$$\langle \boldsymbol{q}_j, \boldsymbol{q}_k \rangle = \delta_{j,k},$$
 $\operatorname{span}\{\boldsymbol{a}_1, \dots, \boldsymbol{a}_n\} = \operatorname{span}\{\boldsymbol{q}_1, \dots, \boldsymbol{q}_n\}$

Why?

One reason is that the \mathbb{C}^m -orthogonal projector onto $\mathrm{span}\{a_1,\ldots,a_n\}$ is given by,

$$oldsymbol{P} = oldsymbol{Q} oldsymbol{Q}^*, \hspace{1cm} oldsymbol{Q} = \left(egin{array}{cccc} ig| & ig| &$$

One essentially explicit algorithm to orthogonalize is Gram-Schmidt.

This algorithm is "triangular orthogonalization": I.e., it's an algorithm that orthogonalizes vectors by accessing them in a triangular pattern.

Input: n vectors $\{a_j\}_{j\in[n]}$. (Assume they're linearly independent for now.) Output: n vectors $\{q_j\}_{j\in[n]}$.

One essentially explicit algorithm to orthogonalize is Gram-Schmidt.

This algorithm is "triangular orthogonalization": I.e., it's an algorithm that orthogonalizes vectors by accessing them in a triangular pattern.

```
Input: n vectors \{a_j\}_{j\in[n]}. (Assume they're linearly independent for now.) Output: n vectors \{q_j\}_{j\in[n]}.
```

Basic idea is induction:

- 1. Set $q_1 = \frac{a_1}{\|a_1\|_2}$, and j = 1.
- 2. Since $\{q_1, \ldots, q_i\}$ have been computed:
 - $lackbox{ }$ Define $\widetilde{m{q}}_{j+1}=(m{I}-m{P}_j)m{a}_{j+1}$, where $m{P}_j$ is the orthogonal projector onto $\{m{q}_1,\ldots,m{q}_j\}$.
 - $\blacktriangleright \mathsf{ Set } \boldsymbol{q}_{j+1} = \widetilde{\boldsymbol{q}}_{j+1} / \|\widetilde{\boldsymbol{q}}_{j+1}\|_2.$
- 3. If j = n 1, quit. Otherwise, $j \leftarrow j + 1$ and go back to step 2.

If the input vectors are linearly independent, this procedure cannot fail. (At least, not in exact arithmetic....)

$$\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_n\} \longrightarrow \{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_n\},$$

 $\operatorname{span}\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_j\} = \operatorname{span}\{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_j\} \quad j \in [n]$

$$\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_n\} \longrightarrow \{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_n\}\,,$$

 $\operatorname{span}\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_j\} = \operatorname{span}\{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_j\} \quad j \in [n]$

We can rewrite this to explicitly express the original vectors $m{a}_j$ in terms of the orthogonalized vectors $m{q}_j$.

At each step:

$$oldsymbol{q}_j = rac{1}{r_{j,j}} \left(oldsymbol{I} - oldsymbol{P}_j
ight) oldsymbol{a}_j \quad \Longrightarrow \quad oldsymbol{a}_j = r_{j,j} oldsymbol{q}_j + \sum_{k \in [j-1]} r_{k,j} oldsymbol{q}_k, \quad r_{k,j} \coloneqq oldsymbol{q}_k^* oldsymbol{a}_j = \left\langle oldsymbol{a}_j, oldsymbol{q}_k
ight
angle,$$

where $r_{j,j} = \|(\boldsymbol{I} - \boldsymbol{P}_j)\boldsymbol{a}_j\|_2$.

$$\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_n\} \longrightarrow \{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_n\},$$

 $\operatorname{span}\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_j\} = \operatorname{span}\{\boldsymbol{q}_1,\ldots,\boldsymbol{q}_j\} \quad j \in [n]$

We can rewrite this to explicitly express the original vectors $m{a}_j$ in terms of the orthogonalized vectors $m{q}_j$.

At each step:

$$oldsymbol{q}_j = rac{1}{r_{j,j}} \left(oldsymbol{I} - oldsymbol{P}_j
ight) oldsymbol{a}_j \quad \Longrightarrow \quad oldsymbol{a}_j = r_{j,j} oldsymbol{q}_j + \sum_{k \in [j-1]} r_{k,j} oldsymbol{q}_k, \quad r_{k,j} \coloneqq oldsymbol{q}_k^* oldsymbol{a}_j = \left\langle oldsymbol{a}_j, oldsymbol{q}_k
ight
angle,$$

where $r_{j,j} = \|(\boldsymbol{I} - \boldsymbol{P}_j)\boldsymbol{a}_j\|_2$.

If $A \in \mathbb{C}^{m \times n}$ has a_j as columns, and $Q \in \mathbb{C}^{m \times n}$ has q_j as columns, then this is equivalent to,

$$\boldsymbol{A} = \boldsymbol{Q}\boldsymbol{R}, \tag{R}_{j,k} = r_{j,k}.$$

Columnwise: this expression is a record of how to reconstruct columns of A from the orthonormal columns of Q.

In fact, these computations implies the following result:

Theorem

Let $A \in \mathbb{C}^{m \times n}$ be any matrix. Then there exists a unitary matrix $Q \in \mathbb{C}^{m \times m}$, and an upper-triangular matrix $R \in \mathbb{C}^{m \times n}$ such that

$$A = QR$$
.

If A has full rank, then the diagonal entries of R can be chosen to be positive.

In fact, these computations implies the following result:

Theorem

Let $A \in \mathbb{C}^{m \times n}$ be any matrix. Then there exists a unitary matrix $Q \in \mathbb{C}^{m \times m}$, and an upper-triangular matrix $R \in \mathbb{C}^{m \times n}$ such that

$$A = QR$$
.

If A has full rank, then the diagonal entries of R can be chosen to be positive.

We previously considered $n = \operatorname{rank}(A) \leqslant m$. The remaining m-n columns of Q are an(y) orthonormal completion of $\{q_1, \dots q_n\}$.

Other cases:

- m = rank(A) < n: Only m vectors are linearly independent.
 - $Q \in \mathbb{C}^{m \times m}$, and R is short+fat.
- $m\geqslant n>r={\rm rank}({\boldsymbol A})$: Columns of ${\boldsymbol A}$ are dependent.

The first r columns of Q span the range of A.

 $m{R}$ is upper triangular, but the main diagonal may have zeros, and the last m-r rows of $m{R}$ vanish.

- "Thin" QR: If rank(A) < n, then Q has only r columns.

In summary:

Given $\{a_j\}_{j\in[n]}\subset\mathbb{C}^m$, compute $\{q_j\}_{j\in[n]}$ such that:

$$\left\langle \boldsymbol{q}_{j}, \boldsymbol{q}_{k} \right\rangle = \delta_{j,k}, \hspace{1cm} \mathrm{span} \left\{ \boldsymbol{a}_{1}, \ldots, \boldsymbol{a}_{n} \right\} = \mathrm{span} \left\{ \boldsymbol{q}_{1}, \ldots, \boldsymbol{q}_{n} \right\}$$

Any algorithm to accomplish this (e.g., Gram-Schmidt) implies:

$$m{A} = m{Q}m{R}, \hspace{1cm} m{A} = \left(egin{array}{ccccc} ig| & ig| & & ig| \ m{a}_1 & m{a}_2 & \cdots & m{a}_n \ ig| & ig| & & ig| \end{array}
ight), \hspace{1cm} m{Q} = \left(egin{array}{cccc} ig| & ig| & & ig| \ m{q}_1 & m{q}_2 & \cdots & m{q}_n \ ig| & & ig| & & ig| \end{array}
ight),$$

with R upper triangular.

In summary:

Given $\{a_j\}_{j\in[n]}\subset\mathbb{C}^m$, compute $\{q_i\}_{j\in[n]}$ such that:

$$\langle \boldsymbol{q}_j, \boldsymbol{q}_k \rangle = \delta_{j,k},$$
 $\operatorname{span}\{\boldsymbol{a}_1, \dots, \boldsymbol{a}_n\} = \operatorname{span}\{\boldsymbol{q}_1, \dots, \boldsymbol{q}_n\}$

Any algorithm to accomplish this (e.g., Gram-Schmidt) implies:

$$m{A} = m{Q}m{R}, \hspace{1cm} m{A} = \left(egin{array}{ccccc} ig| & ig| & & ig| & & ig| \ m{a}_1 & m{a}_2 & \cdots & m{a}_n \ ig| & ig| & & ig| \end{array}
ight), \hspace{1cm} m{Q} = \left(egin{array}{cccc} ig| & ig| & & ig| & m{q}_1 \ ig| & & ig| & & ig| \end{array}
ight),$$

with R upper triangular. "Classical" Gram-Schmidt numerically does:

$$oldsymbol{u}_j = oldsymbol{a}_j - oldsymbol{P}_{j-1} oldsymbol{a}_j, \qquad \qquad oldsymbol{q}_j = rac{oldsymbol{u}_j}{\|oldsymbol{u}_j\|_2}, \qquad \qquad \operatorname{range}(oldsymbol{P}_j) = \operatorname{span}\{oldsymbol{q}_1, \ldots, oldsymbol{q}_j\}.$$

It turns out this is unstable 😇

$$oldsymbol{u}_j = oldsymbol{a}_j - oldsymbol{P}_{j-1} oldsymbol{a}_j, \qquad \qquad oldsymbol{q}_j = \sup\{oldsymbol{q}_1, \dots, oldsymbol{q}_j\}.$$

The cause of numerical instability is that, if a_j is nearly parallel to $\mathrm{span}\{q_1,\ldots,q_{j-1}\}$, this projection step can produce numerically incorrect results. (More precisely, the inner products $r_{k,j}$ for k < j don't accurately represent the coordinates of a_j .)

$$oldsymbol{u}_j = oldsymbol{a}_j - oldsymbol{P}_{j-1} oldsymbol{a}_j, \qquad \qquad oldsymbol{q}_j = rac{oldsymbol{u}_j}{\|oldsymbol{u}_j\|_2}, \qquad \qquad \operatorname{range}(oldsymbol{P}_j) = \operatorname{span}\{oldsymbol{q}_1, \ldots, oldsymbol{q}_j\}.$$

The cause of numerical instability is that, if a_j is nearly parallel to $\mathrm{span}\{q_1,\ldots,q_{j-1}\}$, this projection step can produce numerically incorrect results. (More precisely, the inner products $r_{k,j}$ for k < j don't accurately represent the coordinates of a_j .)

This problem can be fixed with a "modified" version of Gram-Schmidt, which computes $r_{k+1,j}$ by first orthogonalizing a_j against q_k :

- 1. Set $u_j = a_j$, $j \in [n]$
- 2. Set $q_1 = \frac{a_1}{\|a_1\|_2}$, and j = 1.
- 3. For $\ell > 1$: Set $r_{1,\ell} = \boldsymbol{q}_1^* \boldsymbol{u}_{\ell}$, and $\boldsymbol{u}_{\ell} = \boldsymbol{u}_{\ell} r_{1,\ell} \boldsymbol{q}_1$.
- 4. Since u_{j+1} is orthogonal to q_k , $k \in [j]$:
 - Set $q_{j+1} = u_{j+1}/\|u_{j+1}\|_2$
 - For $\ell > j+1$: Set $r_{j+1,\ell} = q_{j+1}^* u_{\ell}$.
 - For $\ell > j+1$: Set $\boldsymbol{u}_{\ell} = \boldsymbol{u}_{\ell} r_{j+1,\ell} \boldsymbol{q}_{j+1}$.
- 5. If j = n 1, quit. Otherwise, $j \leftarrow j + 1$ and go back to step 5.

Thus, the projections are computed "one at a time".

We've seen "classical" (unstable) Gram-Schmidt and "modified" Gram-Schmidt.

In terms of stability, modified Gram-Schmidt effectively fixes the problem.

Both of these operations are "triangular orthogonalization", i.e., they perform the operation,

$$A\mapsto AR^{-1}=Q.$$

In terms of (ℓ^2 -type) conditioning, this operation suffers a condition number of $\kappa(R)$.

We've seen "classical" (unstable) Gram-Schmidt and "modified" Gram-Schmidt.

In terms of stability, modified Gram-Schmidt effectively fixes the problem.

Both of these operations are "triangular orthogonalization", i.e., they perform the operation,

$$A\mapsto AR^{-1}=Q.$$

In terms of (ℓ^2 -type) conditioning, this operation suffers a condition number of $\kappa(\mathbf{R})$.

Instead of attempting to compute Q, we could attempt to compute R. This amounts to "orthogonal triangularization". The reason to cnosider this is that the operation,

$$A\mapsto Q^*A=R,$$

is a much more well-conditioned operation.

There are two high-level strategies for orthogonal triangularization:

- Givens rotations: performs 2×2 unitary operations
- Householder reflectors: performs dimension-n unitary operations

Let ${m P}$ be an orthogonal projection matrix. Then $I-2{m P}$ is Hermitian, unitary, and involutory (is its own inverse).

Let P be an orthogonal projection matrix. Then I-2P is Hermitian, unitary, and involutory (is its own inverse).

Thus, application of this matrix, $\boldsymbol{x}\mapsto (\boldsymbol{I}-2\boldsymbol{P})\boldsymbol{x}$, is well-conditioned.

In particular, if $oldsymbol{P}$ is a rank-1 projector, then there is a unit vector $oldsymbol{v}$ such that

$$P = vv^*$$
.

(And in particular, $x \mapsto (I - 2P)x$ does <u>not</u> require (expensive) matrix-vector multiplications.)

Our main use of these reflectors is the following:

Given $\boldsymbol{x} \in \mathbb{C}^m$, we want to achieve:

$$oldsymbol{x} \stackrel{ ext{Householder reflector}}{\longrightarrow} \|oldsymbol{x}\|e^{i heta}oldsymbol{e}_1,$$

for some $\theta \in [0, 2, \pi)$.

Our main use of these reflectors is the following:

Given $x \in \mathbb{C}^m$, we want to achieve:

$$oldsymbol{x} \stackrel{ ext{Householder reflector}}{\longrightarrow} \|oldsymbol{x}\| e^{i heta} oldsymbol{e}_1,$$

for some $\theta \in [0, 2, \pi)$.

This is achieved by the reflector $I - 2vv^*$, with v given by

$$m{v} = rac{m{x} - \|m{x}\|e^{i heta}m{e}_1}{\|m{x} - \|m{x}\|e^{i heta}m{e}_1\|},$$

for arbitrary θ .

Our main use of these reflectors is the following:

Given $x \in \mathbb{C}^m$, we want to achieve:

$$oldsymbol{x} \stackrel{ ext{Householder reflector}}{\longrightarrow} \|oldsymbol{x}\| e^{i heta} oldsymbol{e}_1,$$

for some $\theta \in [0, 2, \pi)$.

This is achieved by the reflector $I - 2vv^*$, with v given by

$$m{v} = rac{m{x} - \|m{x}\|e^{i heta}m{e}_1}{\|m{x} - \|m{x}\|e^{i heta}m{e}_1\|},$$

for arbitrary θ .

For numerical stability, this reflector should make large changes to x, rather than small changes. The largest change is achieved by selecting

$$e^{i\theta} = -\frac{x_1}{|x_1|}.$$

We now have the following procedure:

Given $\boldsymbol{x} \in \mathbb{C}^m$, we compute $\boldsymbol{v} \in \mathbb{C}^m$ such that

$$(\boldsymbol{I} - 2\boldsymbol{P})\boldsymbol{x} = (\boldsymbol{I} - \boldsymbol{v}\boldsymbol{v}^*)\boldsymbol{x} = c\,\boldsymbol{e}_1,$$

for some scalar $c \in \mathbb{C}$.

We now have the following procedure:

Given $\boldsymbol{x} \in \mathbb{C}^m$, we compute $\boldsymbol{v} \in \mathbb{C}^m$ such that

$$(\mathbf{I} - 2\mathbf{P})\mathbf{x} = (\mathbf{I} - \mathbf{v}\mathbf{v}^*)\mathbf{x} = c\,\mathbf{e}_1,$$

for some scalar $c \in \mathbb{C}$.

Put another way: we can, via an efficiently-applicable unitary transform, map $m{x}$ to $m{e}_1.$

The idea now is that one

- first reflects the first column to $oldsymbol{e}_1$
- then reflects rows 2 through n of column 2 to $oldsymbol{e}_2$
- then reflects rows 3 through n of column 3 to $oldsymbol{e}_3$
- <u>-</u> :

We now have the following procedure:

Given $\boldsymbol{x} \in \mathbb{C}^m$, we compute $\boldsymbol{v} \in \mathbb{C}^m$ such that

$$(\boldsymbol{I} - 2\boldsymbol{P})\boldsymbol{x} = (\boldsymbol{I} - \boldsymbol{v}\boldsymbol{v}^*)\boldsymbol{x} = c\,\boldsymbol{e}_1,$$

for some scalar $c \in \mathbb{C}$.

Put another way: we can, via an efficiently-applicable unitary transform, map $m{x}$ to $m{e}_1.$

The idea now is that one

- first reflects the first column to e_1
- then reflects rows 2 through n of column 2 to $oldsymbol{e}_2$
- then reflects rows 3 through n of column 3 to $oldsymbol{e}_3$

- :

We expect Householder reflectors to be stable since we are simply applying unitary (well-conditioned) matrices to A.

(This is in fact what a typical standard implementations of QR decomposition uses.)

In real arithmetic, a Givens rotation is a 2×2 matrix defined by a 3-tuple $(i, j, \theta) \in [m] \times [n] \times [0, 2\pi)$:

$$G_{\{i,j\},\{i,j\}} = R(\theta)$$
 $R(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$

For real matrices, this accomplishes a rotation of θ radians in the two-dimensional (i,j) plane.

In real arithmetic, a Givens rotation is a 2×2 matrix defined by a 3-tuple $(i, j, \theta) \in [m] \times [n] \times [0, 2\pi)$:

$$m{G}_{\{i,j\},\{i,j\}} = m{R}(heta)$$
 $m{R}(heta) = egin{pmatrix} \cos heta & -\sin heta \\ \sin heta & \cos heta \end{pmatrix}$

For real matrices, this accomplishes a rotation of θ radians in the two-dimensional (i,j) plane.

An alternative to Householder reflectors:

- In column 1, use row 1 to eliminate row j via a Givens rotation, for $j=2,\ldots,n$.
- In column 2, use row 2 to eliminate row j via a Givens rotation, for $j=3,\ldots,n$.
- _ :

In real arithmetic, a Givens rotation is a 2×2 matrix defined by a 3-tuple $(i, j, \theta) \in [m] \times [n] \times [0, 2\pi)$:

$$m{G}_{\{i,j\},\{i,j\}} = m{R}(heta)$$
 $m{R}(heta) = egin{pmatrix} \cos heta & -\sin heta \\ \sin heta & \cos heta \end{pmatrix}$

For real matrices, this accomplishes a rotation of θ radians in the two-dimensional (i,j) plane.

An alternative to Householder reflectors:

- In column 1, use row 1 to eliminate row j via a Givens rotation, for $j=2,\ldots,n$.
- In column 2, use row 2 to eliminate row j via a Givens rotation, for $j=3,\ldots,n$.

- :

While both Householder reflectors and Givens rotations are both effective (well-conditioned), Householder reflectors are generally employed for generic dense QR decomposition operations since they require fewer conceptual and computational steps.

(To zero out column 1, we need just one Householder reflector, but n-1 Givens rotations.)

If $A \in \mathbb{C}^{m \times n}$ and $b \in \mathbb{C}^n$, we are interested in computing the least-squares solution to

$$Ax = b$$

This arises in several situations, e.g., data fitting.

If $A \in \mathbb{C}^{m \times n}$ and $b \in \mathbb{C}^n$, we are interested in computing the least-squares solution to

$$Ax = b$$

This arises in several situations, e.g., data fitting.

The following is a result we have essentially already proven:

Theorem

Suppose $A \in \mathbb{C}^{m \times n}$ has full column rank (n). Then, for any $b \in \mathbb{C}^n$, there is a unique solution x that solves

$$\underset{\boldsymbol{x} \in \mathbb{C}^n}{\arg\min} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2.$$

Furthermore, this solution x is the unique solution to $A^*Ax = A^*b$, and the residual r := b - Ax is orthogonal to range(A).

The system $A^*Ax = A^*b$ is called the *normal equations*.

While the normal equations are typically useful for analysis, they are typically not used for computation.

$$A = QR \implies x = R^{-1}Q^*b.$$

In most cases, the QR decomposition is used through the above procedure, largely for stability reasons.

While the normal equations are typically useful for analysis, they are typically not used for computation.

$$A = QR \implies x = R^{-1}Q^*b.$$

In most cases, the QR decomposition is used through the above procedure, largely for stability reasons.

One can still use this same idea if A doesn't have full column rank, but has to be more careful.

E.g., ${m Q}$ should only span the range of ${m A}$.

Finally, we note that QR can fail when A does not have full column rank, rank(A) = r < n.

One standard way to address this is by column pivoting.

If doing triangular orthogonalization (Gram-Schmidt-type): At orthogonalization step j, identify column index k satisfying:

$$\|\boldsymbol{a}_{k} - \boldsymbol{P}_{j-1}\boldsymbol{a}_{k}\|_{2} \geqslant \|\boldsymbol{a}_{\ell} - \boldsymbol{P}_{j-1}\boldsymbol{a}_{\ell}\|_{2},$$
 $\ell \geqslant j$

Then interchange (permute) columns k and j. Repeat at every orthogonalization step. This results in the factorization:

$$AP = QR$$
.

Finally, we note that QR can fail when A does not have full column rank, rank(A) = r < n.

One standard way to address this is by column pivoting.

If doing orthogonal triangularization (Householder/Givens): At step j, with $S = \{j, j+1, \ldots, n\}$, identify column index k satisfying:

$$\|\boldsymbol{A}_{S,k}\|_2 \geqslant \|\boldsymbol{A}_{S,\ell}\|_2, \qquad \qquad \ell \geqslant j.$$

Then interchange (permute) columns k and j. Repeat at every step. This (again) results in the factorization:

$$AP = QR$$
.

References I D07-S16(a)



Atkinson, Kendall (1989). An Introduction to Numerical Analysis. New York: Wiley. ISBN: 978-0-471-62489-9.



Salgado, Abner J. and Steven M. Wise (2022). *Classical Numerical Analysis: A Comprehensive Course*. Cambridge: Cambridge University Press. ISBN: 978-1-108-83770-5. DOI: 10.1017/9781108942607.



Trefethen, Lloyd N. and David Bau (1997). *Numerical Linear Algebra*. SIAM: Society for Industrial and Applied Mathematics. ISBN: 0-89871-361-7.