

STABILIZED ORTHOGONAL FACTORIZATIONS

1. Conditioning of linear algebra problems

Recall that the relative condition number of a mathematical problem $f: X \rightarrow Y$ characterizes the amplification by f of perturbations in its argument

$$\kappa = \lim_{\varepsilon \rightarrow 0} \sup_{\|\delta x\| \leq \varepsilon} \left(\frac{\|f(x + \delta x) - f(x)\|}{\|f(x)\|} / \frac{\|\delta x\|}{\|x\|} \right).$$

Linear combination. The basic operation of linear combination \mathbf{Ax} , $\mathbf{A} \in \mathbb{C}^{m \times n}$, seen as the problem $\mathbb{C}^n \xrightarrow{f} \mathbb{C}^m$ has the condition number

$$\kappa = \sup_{\delta x} \left(\frac{\|\mathbf{A} \delta x\|}{\|\mathbf{Ax}\|} / \frac{\|\delta x\|}{\|x\|} \right) = \sup_{\delta x} \left(\frac{\|\mathbf{A} \delta x\|}{\|\delta x\|} \right) \frac{\|x\|}{\|\mathbf{Ax}\|} = \|\mathbf{A}\| \frac{\|x\|}{\|\mathbf{Ax}\|}.$$

The matrix norm property $\|\mathbf{Ay}\| \leq \|\mathbf{A}\| \|\mathbf{y}\|$ can be used to obtain

$$\|x\| = \|\mathbf{I}_n x\| = \|\mathbf{A}^+ \mathbf{Ax}\| \leq \|\mathbf{A}^+\| \|\mathbf{Ax}\| \Rightarrow \frac{\|x\|}{\|\mathbf{Ax}\|} \leq \|\mathbf{A}^+\|$$

leading to

$$\kappa \leq \|\mathbf{A}\| \|\mathbf{A}^+\| = \kappa(\mathbf{A}),$$

where $\kappa(\mathbf{A})$ is the condition number of the matrix \mathbf{A} . If \mathbf{A} is of full rank with $m > n$, the 2-norm condition number is given by the ratio of largest to smallest singular values.

$$\|\mathbf{A}\| = \sigma_1, \|\mathbf{A}^+\| = 1/\sigma_n \Rightarrow \kappa(\mathbf{A}) = \sigma_1/\sigma_n \geq 1.$$

By convention, if \mathbf{A} is singular, the condition number $\kappa(\mathbf{A}) = \infty$.

Coordinate transformation. For $\mathbf{A} \in \mathbb{C}^{m \times m}$ of full rank, the coordinates of vector $\mathbf{b} \in \mathbb{C}^m$ expressed in the \mathbf{I} basis can be transformed its coordinates $\mathbf{x} \in \mathbb{C}^m$ in the \mathbf{A} basis by solving the linear system $\mathbf{Ax} = \mathbf{Ib}$, with the solution $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ (so written formally, even though the inverse is almost never explicitly computed). This is simply another linear combination of the columns of \mathbf{A}^{-1} , hence the problem $f: \mathbb{C}^m \rightarrow \mathbb{C}^m, f(\mathbf{b}) = \mathbf{A}^{-1}\mathbf{b}$ again has a condition number bounded by the condition number of the matrix \mathbf{A} .

$$\kappa \leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| = \kappa(\mathbf{A}) = \kappa(\mathbf{A}^{-1}).$$

Operator perturbation. Instead of changing the input data as above, the linear mapping represented by the matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ might itself be perturbed. Two mathematical problems may now be formulated:

1. For fixed $\mathbf{b} \in \mathbb{C}^m, f: \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^n, f(\mathbf{A}, \mathbf{b}) = \mathbf{A}^+ \mathbf{b} = \mathbf{x}$. Perturbation of the input \mathbf{A} induces perturbation of \mathbf{x} in order for \mathbf{b} to be kept fixed

$$(\mathbf{A} + \delta \mathbf{A})(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b}.$$

Using $\mathbf{Ax} = \mathbf{b}$, and assuming that $\delta \mathbf{A} \delta \mathbf{x}$ is negligible gives

$$\mathbf{A} \delta \mathbf{x} + \delta \mathbf{A} \mathbf{x} = \mathbf{0} \Rightarrow \delta \mathbf{x} = -\mathbf{A}^+ \delta \mathbf{A} \mathbf{x},$$

hence the relative condition number is

$$\kappa = \frac{\|\mathbf{A}^+ \delta \mathbf{A} \mathbf{x}\|}{\|x\|} \cdot \frac{\|x\|}{\|\delta \mathbf{A}\|} \leq \frac{\|\mathbf{A}^+\| \|\delta \mathbf{A} \mathbf{x}\|}{\|x\|} \cdot \frac{\|x\|}{\|\delta \mathbf{A}\|} \leq \frac{\|\mathbf{A}^+\| \|\delta \mathbf{A}\| \|x\|}{\|x\|} \cdot \frac{\|x\|}{\|\delta \mathbf{A}\|} = \kappa(\mathbf{A}).$$

For all above linear algebra problems the condition number is bounded by the associated matrix condition number. Unitary matrices $\mathbf{Q} \in \mathbb{C}^{m \times m}$ have $\kappa(\mathbf{Q}) = 1$, and define an orthonormal basis for \mathbb{C}^m . A rank-deficient matrix $\mathbf{Z} \in \mathbb{C}^{m \times m}$ has $\kappa(\mathbf{Z}) = \infty$, and corresponds to a linearly dependent vector set $\{z_1, \dots, z_m\}$. The behavior of many numerical approximation procedures based upon linear combinations is determined by condition number of the basis set.

- *Monomial basis with uniform sampling.* Sampling the monomial basis on interval $[a, b]$ at $t_i = ih + a, i = 0, m, h = (b-a)/(m-1)$ leads to the Vandermonde matrix

$$\mathbf{V} = [\mathbf{1} \ t \ \dots \ t^m],$$

an extremely ill-conditioned matrix (Fig.). This can readily be surmised from the example $a=0, b=1$, in which case for large m the last columns of V become ever more colinear to the same e_m vector. Series expansions based on the monomials such as the Taylor series

$$f(t) = f(0) + f'(0)t + \dots + \frac{f^{(n)}(0)}{n!}t^n + \dots$$

are highly sensitive to perturbations, small changes in $f(t)$ lead to large changes in the coordinates $\{f(0), f'(0), \dots\}$.

```

∴ function Vandermonde(a,b,m)
    t=LinRange(a,b,m); v=ones(m,1); V=copy(v)
    for j=2:m
        v = v .* t; V=[V v]
    end
    return V
end;

```

∴

- *Monomial basis with Chebyshev sampling.* Changing the sampling so that points are clustered towards the interval endpoints reduces the condition number at fixed number of sampling points m , but the same limiting behavior for large m is obtained.

```

∴ function VandermondeC(m)
    δ=π/(2*m); θ=LinRange(δ,π-δ,m)
    t=cos.(θ)
    v=ones(m,1); V=copy(v)
    for j=2:m
        v = v .* t; V=[V v]
    end
    return V
end;

```

∴

- *Triangular basis with uniform sampling.* LU-factorization of the monomial basis leads to a different family of polynomials, known as a triangular basis

$$\{1, t-x_1, (t-x_1) \cdot (t-x_2), \dots, (t-x_1) \cdot \dots \cdot (t-x_{m-1})\},$$

where $\{x_1, \dots, x_m\}$ are known as the nodes of the system. These can be chosen to uniformly sample an interval. As to be expected, applying a sequence of non-unitary linear transformations onto an ill-conditioned basis yields even worse conditioning.

```

∴ function Triangular(a,b,m)
    x=LinRange(a,b,m); T=ones(m,1); Tj=copy(T); t=copy(x)
    for j=2:m
        Tj = Tj .* (t .- x[j-1]); T=[T Tj]
    end
    return T
end;

```

∴

- *Triangular basis with Chebyshev sampling.* Adopting Chebyshev sampling ameliorates the conditioning, but only marginally.

```

∴ function TriangularC(m)
    δ=π/(2*m); Θ=LinRange(δ,π-δ,m)
    x=cos.(Θ); T=ones(m,1); Tj=copy(T); t=copy(x)
    for j=2:m
        Tj = Tj .* (t .- x[j-1]); T=[T Tj]
    end
    return T
end;

```

```
∴
```

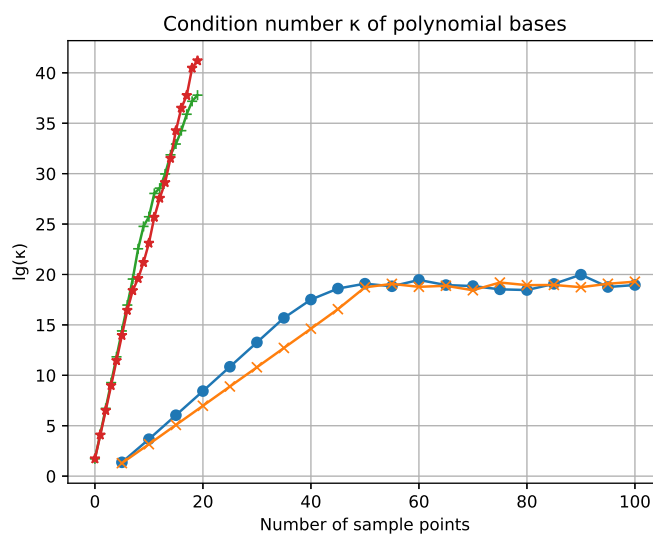


Figure 1. Monomial basis with: (o) uniform sampling, (x) Chebyshev sampling. Triangular basis with: (+) uniform sampling, (*) Chebyshev sampling.

```

∴ mr=5:5:100; κVDMU=log10.(cond.(Vandermonde.(-1,1,mr)));
∴ κVDMC=log10.(cond.(VandermondeC.(mr)));
∴ κTU=log10.(cond.(Triangular.(-1,1,mr)));
∴ κTC=log10.(cond.(TriangularC.(mr)));
∴

```

```

∴ x=collect(mr); clf();
∴ plot(x,κVDMU,"o-",x,κVDMC,"x-",κTU,"+-",κTC,"*-");
∴ grid("on"); title("Condition number κ of polynomial bases");
∴ xlabel("Number of sample points"); ylabel("lg(κ)");
∴ pre=homedir()*"/courses/MATH661/images/";
∴ savefig(pre*"PolyBasesCondNr.eps");
∴

```

2. Orthogonal factorization through Householder reflectors

The Gram-Schmidt procedure constructs an orthogonal factorization by linear combinations of the column vectors of $A \in \mathbb{C}^{m \times n}$, $m \geq n$, $\text{rank}(A) = n$

$$AR_1R_2 \dots R_n = Q \Rightarrow A = QR, R = R_n^{-1} \dots R_1^{-1}.$$

In exact arithmetic $C(Q) = C(A)$ by construction, and $\kappa(Q) = 1$, but the sequence of multiplications with R_1, \dots, R_n might amplify perturbations in A (due for example to floating point representation errors or inherent uncertainty in knowledge of A). The problem $f: \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^{m \times n} \times \mathbb{C}^{n \times n}$, $A \xrightarrow{f} Q, R$ has condition number

$$\kappa = \frac{\|\delta Q\|}{\|Q\|} \cdot \frac{\|A\|}{\|\delta A\|} + \frac{\|\delta R\|}{\|R\|} \cdot \frac{\|A\|}{\|\delta A\|},$$

and numerical experimentation (Fig. 2) readily exhibits large condition numbers.

An alternative approach is to obtain an orthogonal factorization through unitary transformations

$$Q_n \dots Q_1 A = R \Rightarrow A = QR, Q = Q_1^* \dots Q_n^*.$$

Unitary transformations do not modify vector 2-norms or relative orientations

$$\|Qx\|^2 = x^* Q^* Q x = \|x\|^2, (Qy)^*(Qx) = y^* x,$$

and are hence said to be isometric. In Euclidean space reflections and rotations are isometric.

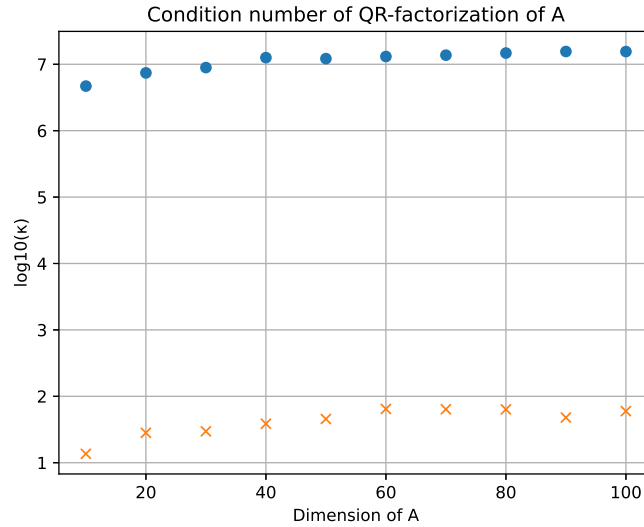


Figure 2. QR-conditioning: (o) modified Gram-Schmidt, (x) Householder.

Construction of an isometric reflection transformation suitable for a QR factorization is represented in Fig. 3. Let vector $x \in \mathbb{C}^{m+1-k}$ represent the portion of the k^{th} column from the diagonal downwards in stage k of reduction of $A \in \mathbb{C}^{m \times n}$ to upper triangular form

$$Q_{k-1} \dots Q_1 A = \begin{bmatrix} R & C \\ \mathbf{0} & B \end{bmatrix}, B = [x \ b_2 \ \dots \ b_{n-k}].$$

The next stage of in reduction to upper triangular form is the isometric transformation of x into $\pm \|x\| e_1$, with $e_1 \in \mathbb{C}^{m+1-k}$ the unit vector along the first direction. With $v = \pm \|x\| e_1 - x$, $q = v / \|v\|$, the projection of x onto the span of v , $C(v)$ is

$$y = P_v x = qq^* x,$$

and its complementary projector onto $N(v^*)$ is

$$z = P_{\perp v} = (I - qq^*)x.$$

The reflector transforming x into $\pm\|x\|e_1$ is obtained by doubling the above displacements, and is known as a Householder reflector

$$H = I - 2qq^*.$$

Of the two possibilities $\pm\|x\|e_1$, the choice

$$v = -\text{sign}(x_1)\|x\|e_1 - x,$$

avoids loss of floating accuracy $x \cong \|x\|e_1$. For $x \in \mathbb{C}^{m+1-k}$, $\text{sign}(x_1) = \exp(\arg(x_1))$.

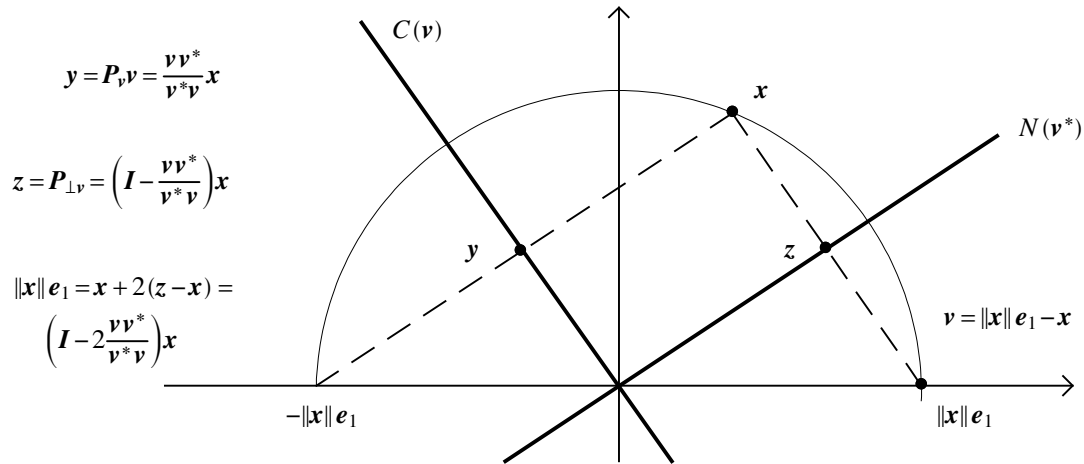


Figure 3. Geometry of Householder reflector

The resulting Householder QR -factorization is given

Input: $A \in \mathbb{C}^{m \times n}$

$Q = \mathbf{0}_{m,n}$

for $k = 1:n$

$x = A[k:m, k]$

$v = \text{sign}(x_1)\|x\| + x$

$q = v/\|v\|$; $Q[k:m, k] = q$

for $j = k:n$

$A[k:m, j] = A[k:m, j] - 2q(q^*A[k:m, j])$

```

function HouseholderQR(A)
    m, n = size(A)
    Q = zeros(m, n); R = copy(A)
    for k = 1:n
        x = R[k:m, k]
        e1 = zeros(size(x)); e1[1] = 1
        v = sign(x[1]) * norm(x) * e1 + x
        q = v / norm(v); Q[k:m, k] = q
        for j = k:n
            aj = R[k:m, j]; c = 2 * q' * aj
            R[k:m, j] = aj - c * q
        end
    end
    return Q, R
end;

```

Note that the above implementation does not return the Q matrix, but rather the Q_1, \dots, Q_n reflectors from which Q can be reconstructed if needed. Usually though, the Q matrix itself is not required, but rather the product Qu which can readily be evaluated as $Q_n \dots Q_1 u$. The Householder reflector algorithm is typically the default procedure in QR -factorizations implemented in software systems, and as seen in (Fig. 2), leads to much better conditioning.

3. Orthogonal factorization through Given rotators

An alternative approach to orthogonal factorization utilizes isometric rotation transformations of the form

$$R(i, k, \theta) = I + (\cos \theta - 1)(e_i e_i^* + e_k e_k^*) - \sin \theta (e_i e_k^* - e_k e_i^*),$$

with the rotation angle θ chosen to nullify the subdiagonal element (i, k) , $i > k$

$$(\mathbf{R}(i, k, \theta) \mathbf{A})_{ik} = a_{kk} \sin \theta + a_{ik} \cos \theta = 0 \Rightarrow \theta_{ik} = \arctan\left(-\frac{a_{ik}}{a_{kk}}\right).$$

Composition of repeated rotations $\mathbf{Q}_{ik} = \mathbf{R}(i, k, \theta_{ik})$ can be organized to lead to an upper triangular matrix

$$\mathbf{Q}_{mn} \dots \mathbf{Q}_{32} \mathbf{Q}_{m1} \dots \mathbf{Q}_{31} \mathbf{Q}_{21} \mathbf{A} = \mathbf{R}.$$

Whereas Householder reflectors work on entire vectors, Givens rotators nullify individual subdiagonal elements. For full matrices Householder reflectors typically require fewer floating point operations, but the special structure of a sparse matrix is better exploited by use of Givens rotators.

Input: $\mathbf{A} \in \mathbb{C}^{m \times n}$

$\mathbf{Q} = \mathbf{0}_{m,n}$

for $k = 1:n$

 for $i = k + 1:m$

$\theta = \arctan(-a_{ik}/a_{kk})$

$c = \cos(\theta)$; $s = \sin(\theta)$

 for $j = k:n$

$u = a_{kj}$; $v = a_{ij}$

$a_{kj} = c u - s v$

$a_{ij} = s u + c v$

```

∴ function GivensQR(A)
    m,n=size(A)
    Q=zeros(m,n); R=copy(A)
    for k=1:n
        for i=k+1:m
            θ = atan(-R[i,k],R[k,k]); Q[i,k]=
            c = cos(θ); s = sin(θ)
            for j=k:n
                u = R[k,j]; v = R[i,j]
                R[k,j]=c*u-s*v
                R[i,j]=s*u+c*v
            end
        end
    end
    return Q,R
end;

```

∴

As in the Householder implementation the above implementation returns data to reconstruct \mathbf{Q} if needed.