

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today



Last time

- ▶ Linear least square problems
- ▶ Geometric perspective on the normal equations

Today

- ▶ Orthogonalization with Gram-Schmidt
- ▶ QR decomposition

Announcements

- ▶ Homework 3 has been posted, due Mon, Oct 24 before class

Recap: Least-squares problems

Choosing the least square error, this results in

$$\min_{\mathbf{x}} \|\underline{A\mathbf{x}} - \underline{\mathbf{b}}\|^2,$$

where $\mathbf{x} = (x_1, \dots, x_n)^T$, $\mathbf{b} = (b_1, \dots, b_m)^T$, and $a_{ij} = a_j(t_i)$.

In the following, we study the **overdetermined case**, i.e., $m \geq n$ and $\text{rank}(A) = \underline{n}$

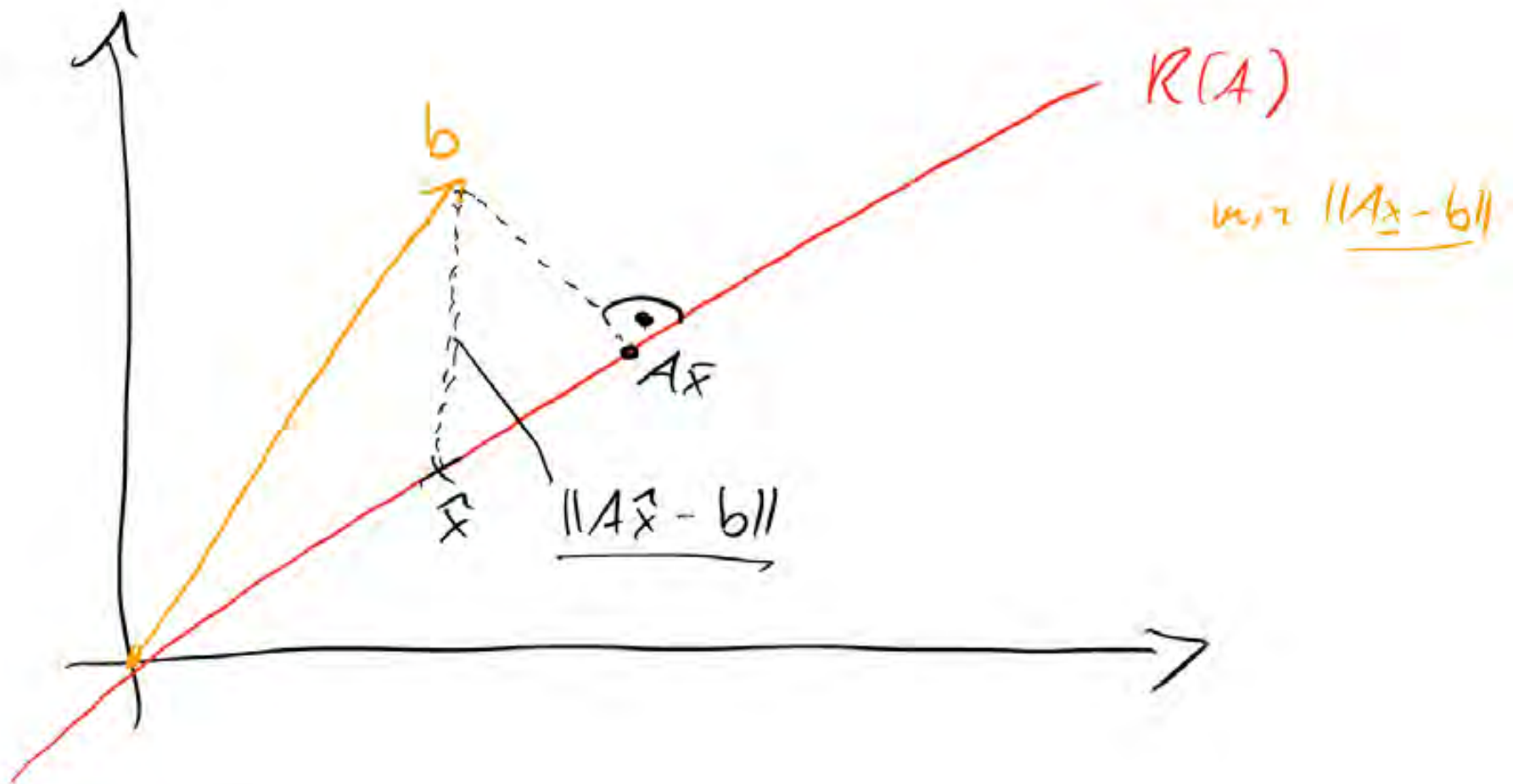
Solving the **normal equations**

$$A^T A \bar{\mathbf{x}} = A^T \mathbf{b}$$

requires:

- ▶ computing $A^T A$ (which is $O(mn^2)$)
- ▶ condition number of $A^T A$ is square of condition number of A ; (problematic for the Choleski factorization)

Recap: Least-squares problems



Recap: Linear least-squares problems

Now for the least-squares problem $\|\mathbf{Ax} - \mathbf{b}\|_2$. The relative condition number κ in the Euclidean norm is bounded by

- ▶ With respect to perturbations in \mathbf{b} :

$$\kappa \leq \frac{\kappa_2(A)}{\cos(\theta)}$$

- ▶ With respect to perturbations in \mathbf{A} :

$$\kappa \leq \underbrace{\kappa_2(A)} + \underbrace{\kappa_2(A)^2}_{\tan(\theta)}$$

Small residual problems, small angle θ $\cos(\theta) \approx 1$, $\tan(\theta) \approx 0$: behavior similar to linear system.

Large residual problems, large angle θ $\cos(\theta) \ll 1$, $\tan(\theta) \approx 1$: behavior very different from linear system because $\kappa_2(A)^2$ shows up

The QR decomposition

Recall that projecting \mathbf{b} onto the column span (range) of \mathbf{A} was the key step \rightsquigarrow let's try to find a numerical method that computes an orthonormal basis $\mathbf{q}_1, \dots, \mathbf{q}_n$ of the rank- n column span of \mathbf{A}

$$\mathbf{A} = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad m \geq n$$

$$\Downarrow$$

$$\underbrace{\begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix}}_{\mathbf{A}} = \underbrace{\begin{bmatrix} | & & | \\ \mathbf{q}_1 & \dots & \mathbf{q}_n \\ | & & | \end{bmatrix}}_{\mathbf{Q}} \underbrace{\begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & & \vdots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_{\mathbf{R}}$$

with an invertible matrix \mathbf{R} so that

$$\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_k) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k), \quad k = 1, \dots, n$$

$$\underbrace{\begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix}}_A = \underbrace{\begin{bmatrix} | & & | \\ \mathbf{q}_1 & \dots & \mathbf{q}_n \\ | & & | \end{bmatrix}}_Q \underbrace{\begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & & \dots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_R$$

\Downarrow leads to system of equations \Downarrow

$$\mathbf{a}_1 = r_{11} \mathbf{q}_1$$

$$\mathbf{a}_2 = r_{12} \mathbf{q}_1 + r_{22} \mathbf{q}_2$$

$$\mathbf{a}_3 = r_{13} \mathbf{q}_1 + r_{23} \mathbf{q}_2 + r_{33} \mathbf{q}_3$$

$$\vdots$$

$$\mathbf{a}_n = r_{1n} \mathbf{q}_1 + r_{2n} \mathbf{q}_2 + \dots + r_{nn} \mathbf{q}_n$$

This motivates a process for computing the basis $\mathbf{q}_1, \dots, \mathbf{q}_n$

- ▶ At step j , we have $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ that span $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_{j-1})$
- ▶ We want to find \mathbf{q}_j orthonormal to $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ so that $\mathbf{q}_1, \dots, \mathbf{q}_j$ spans $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j)$
- ▶ Thus, set

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \dots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

and normalize

$$\mathbf{q}_j = \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|_2}$$

Notice that at step j , the quantities $\mathbf{q}_1^T \mathbf{a}_j, \mathbf{q}_2^T \mathbf{a}_j, \dots, \mathbf{q}_{j-1}^T \mathbf{a}_j$ are the values $r_{j,1}, \dots, r_{j,j-1}$ and r_{jj} is responsible for the normalization and set to

$$r_{jj} = \|\mathbf{a}_j - \sum_{i=1}^{j-1} r_{ij} \mathbf{q}_i\|_2$$

This process is the *classical Gram-Schmidt* procedure to compute the QR factorization;
however, this process is numerically unstable!

Instead of directly computing

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \cdots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

based on \mathbf{a}_j , the *modified* Gram-Schmidt procedure computes \mathbf{v}_j iteratively

$$\begin{aligned} \mathbf{v}_j^{(1)} &= \mathbf{a}_j, \\ \mathbf{v}_j^{(2)} &= \mathbf{v}_j^{(1)} - \mathbf{q}_1 \mathbf{q}_1^T \mathbf{v}_j^{(1)}, && \text{"subtract from } \mathbf{v}_j^{(1)} \text{ what is already in } \mathbf{q}_1\text{"} \\ \mathbf{v}_j^{(3)} &= \mathbf{v}_j^{(2)} - \mathbf{q}_2 \mathbf{q}_2^T \mathbf{v}_j^{(2)}, && \text{"subtract from } \mathbf{v}_j^{(2)} \text{ what is already in } \mathbf{q}_2\text{"} \\ &\vdots \\ \mathbf{v}_j &= \mathbf{v}_j^{(j)} = \mathbf{v}_j^{(j-1)} - \mathbf{q}_{j-1} \mathbf{q}_{j-1}^T \mathbf{v}_j^{(j-1)} \end{aligned}$$

Computing a QR factorization with the modified Gram-Schmidt procedure is stabler than with the classical Gram-Schmidt procedure. However, even the modified Gram-Schmidt procedure can lead to vectors $\mathbf{q}_1, \dots, \mathbf{q}_n$ that are far from orthogonal if the condition number of \mathbf{A} is large (see, Golub et al., Matrix Computations, Section 5.2.9)

Let's recall what the Gram-Schmidt procedure is doing: It is applying a succession of triangular matrices R_k on the right of A so that the resulting matrix

$$\underbrace{A R_1 R_2 \dots R_n}_{R^{-1}} = Q$$

has orthonormal columns and R is upper-triangular.

Instead, we could try to find orthonormal matrices ($\mathbf{X}^T \mathbf{X} = \mathbf{X} \mathbf{X}^T = I$) so that

$$\underbrace{Q_n \dots Q_2 Q_1}_Q A = R$$

is upper-triangular. The product $Q_n \dots Q_2 Q_1 = Q^T$ is orthonormal too and thus $A = QR$ a QR factorization of A .

The Householder method judiciously finds the matrices Q_1, Q_2, \dots, Q_n via so-called Householder reflectors \rightsquigarrow **board**. The Householder method is backward stable.

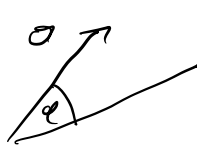
Transform

$$A \rightsquigarrow Q_1 A \rightsquigarrow Q_2 Q_1 A \rightsquigarrow \dots$$

$$k_2(Q) = 1, \quad |\det(Q)| = 1$$

In \mathbb{R}^2 ; two orthogonal transformations

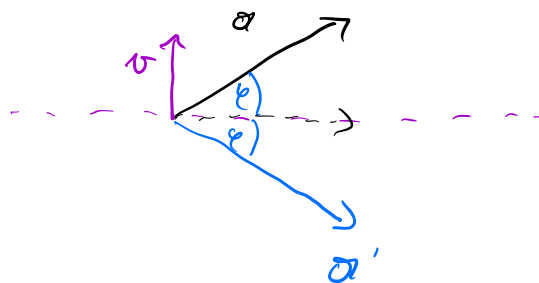
rotation: $\det = 1$


$$\begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

\Rightarrow Givens rotations

\Rightarrow fast books

reflection $\det = -1$



projection

$$a \mapsto \left(I - \frac{vv^T}{v^T v}\right) a$$

reflection

$$a \mapsto \left(I - 2 \frac{vv^T}{v^T v}\right) a$$

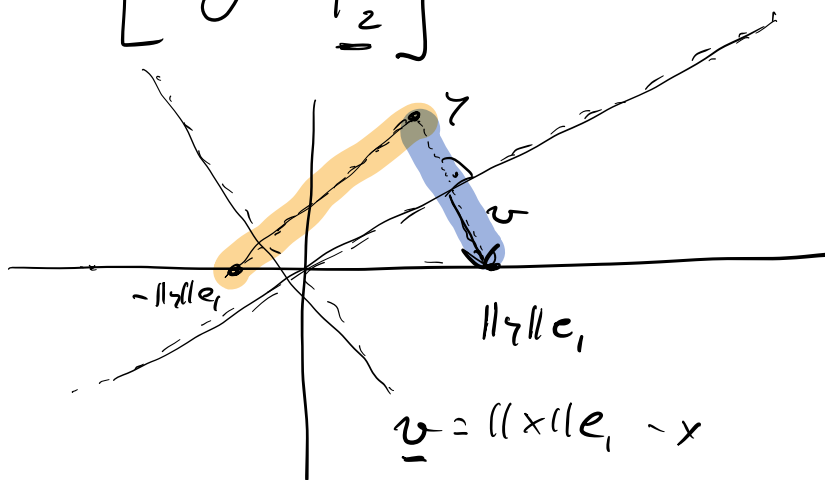
\Rightarrow Householder
reflections

$$\begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{bmatrix} \xrightarrow{P=1, Q_1} \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} \xrightarrow{P=2, Q_2} \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \\ 0 & 0 & x \end{bmatrix}$$

$$\gamma = \begin{bmatrix} x \\ x \\ x \\ x \end{bmatrix} \xrightarrow{F} \begin{bmatrix} x \\ 0 \\ 0 \\ 0 \end{bmatrix} = x e_1$$

find

$$Q_2 = \begin{bmatrix} I_{k-1} & 0 \\ 0 & F_2 \end{bmatrix}$$



$$F_2 = (I - 2 \frac{v v^T}{v^T v})$$

Costs of Householder reflection

$$\underline{2mn^2 - \frac{2}{3}n^3} \text{ FLOPs}$$

to get R

Additionally get Q $O(mn)$

$$\Rightarrow O(mn^2)$$

—

$$A = QR = \begin{matrix} & n \\ n & \end{matrix} \left[\quad \right] \begin{matrix} \diagup \\ \text{wavy line} \\ \diagdown \end{matrix}$$

$$\underline{R_x} = \underline{Q^T b}$$

The QR factorization

All these three algorithms (classical Gram-Schmidt, modified Gram-Schmidt, Householder triangularization) have roughly the FLOPs of $2mn^2$ for an $m \times n$ matrix

Why would we ever want to use (modified) Gram-Schmidt instead of Householder triangularization?

The QR factorization

All these three algorithms (classical Gram-Schmidt, modified Gram-Schmidt, Householder triangularization) have roughly the FLOPs of $2mn^2$ for an $m \times n$ matrix

Why would we ever want to use (modified) Gram-Schmidt instead of Householder triangularization? Gram-Schmidt can be easier to parallelize, for example (Recall that best algorithm depends also on what hardware we want to implement it on.)

Every matrix $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ has a QR factorization. It is unique if we require the diagonal elements of R to be positive.

If $m > n$ and $\mathbf{Q} \in \mathbb{R}^{m \times n}$, then we speak of a reduced QR factorization. Otherwise, we have $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and we speak of a full QR factorization.

```
1: >> A = randn(10, 10); [Q, R] = qr(A);
2: >> size(Q)
3: ans =
4:      10      10
5: >> size(R)
6: ans =
7:      10      10
```

```
1: >> A = randn(10, 4); [Q, R] = qr(A)
2: >> size(Q)
3: ans =
4:      10      10
5: >> size(R)
6: ans =
7:      10       4
8: >>
9: >> [Q, R] = qr(A, 0); % reduced QR
10: >> size(Q)
11: ans =
12:      10       4
13: >> size(R)
14: ans =
15:       4       4
```

Back to our least-squares problem

One would like to avoid the multiplication $A^T A$ and use a suitable factorization of A that avoids solving the normal equation directly:

$$\underline{A} = \underline{Q}\underline{R} = [\underline{Q}_1, \underline{Q}_2] \begin{bmatrix} \underline{R}_1 \\ \underline{0} \end{bmatrix} = \underline{Q}_1 \underline{R}_1,$$

where $Q \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($QQ^T = I$), and $R \in \mathbb{R}^{m \times n}$ consists of an upper triangular matrix and a block of zeros.

How can the QR factorization be used to solve the least-squares problem?

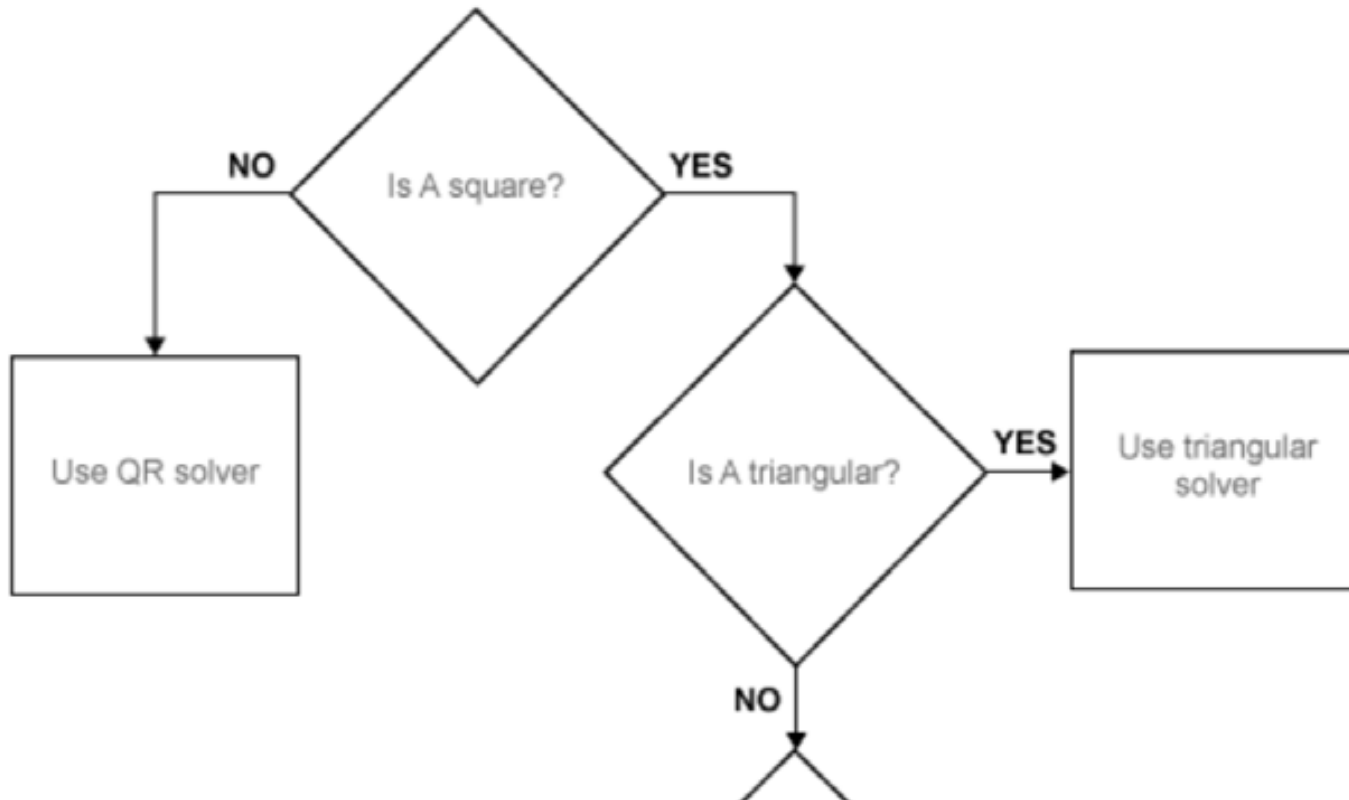
$$\begin{aligned} \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 &= \min_{\mathbf{x}} \|\underline{Q}^T(A\mathbf{x} - \mathbf{b})\|^2 &&= \min_{\mathbf{x}} \left\| \begin{bmatrix} \mathbf{b}_1 - R_1\mathbf{x} \\ \mathbf{b}_2 \end{bmatrix} \right\|^2, \\ &= \min_{\mathbf{x}} \underbrace{\|\mathbf{b}_1 - R_1\mathbf{x}\|^2}_{\text{residual}} + \underbrace{\|\mathbf{b}_2\|^2}_{\text{residual}} \end{aligned}$$

where $Q^T \underline{\mathbf{b}} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}$.

Thus, the least squares solution is $\mathbf{x} = R^{-1}\mathbf{b}_1$ and the residual is $\|\mathbf{b}_2\|$.

Stability of solving least-squares problem with Householder triangularization

Solving a least-squares problem with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$ and $\text{rank}(\mathbf{A}) = n$ via QR factorization computed with Householder triangularization is backward stable.



Eigen decomposition

Eigen decomposition

- ▶ For a square matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, there exists at least one λ such that

$$\mathbf{Ax} = \lambda \mathbf{x} \implies (\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = 0$$

- ▶ Putting the eigenvectors \mathbf{x}_j as columns in a matrix \mathbf{X} , and the eigenvalues λ_j on the diagonal of a diagonal matrix $\mathbf{\Lambda}$, we get

$$\mathbf{AX} = \mathbf{X}\mathbf{\Lambda}$$

- ▶ A matrix is non-defective or diagonalizable if there exist n linearly independent eigenvectors, which means that \mathbf{X} is invertible

$$\mathbf{X}^{-1}\mathbf{AX} = \mathbf{\Lambda}$$

$$\mathbf{A} = \mathbf{X}\mathbf{\Lambda}\mathbf{X}^{-1}$$

- ▶ The transformation from \mathbf{A} to $\mathbf{\Lambda} = \mathbf{X}^{-1}\mathbf{AX}$ is called a similarity transformation and it preserves the eigenvalues.

- ▶ A matrix is unitarily diagonalizable if there exist n linearly independent orthogonal eigenvectors, i.e., if the matrix \mathbf{X} can be chosen to be unitary (orthonormal), $\mathbf{X} = \mathbf{U}$, where $\mathbf{U}^{-1} = \mathbf{U}^H$

$$\underline{\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H}$$

Note that unitary matrices generalize orthogonal matrices to the complex domain, so we use adjoints (conjugate transpose) instead of transpose throughout

- ▶ Theorem: A matrix is unitarily diagonalizable iff it is normal, i.e., it commutes with its adjoint:

$$\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$$

- ▶ Theorem: Hermitian (symmetric) matrices, $\mathbf{A}^H = \mathbf{A}$, are unitarily diagonalizable and have *real* eigenvalues.

- ▶ The usual eigenvectors are more precisely called *right* eigenvectors. There are also *left* eigenvectors corresponding to a given eigenvalue λ

$$\underbrace{\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H} \implies \underbrace{\mathbf{A}^H \mathbf{y} = \bar{\lambda} \mathbf{y}},$$
$$\mathbf{Y}^H \mathbf{A} = \mathbf{\Lambda} \mathbf{Y}^H$$

with conjugate $\bar{\lambda}$ of λ

- ▶ For a matrix that is diagonalizable, observe that

$$\mathbf{Y}^H = \mathbf{X}^{-1}$$

and so the left eigenvectors provide no new information

- ▶ For unitarily diagonalizable matrices, $\mathbf{Y} = (\mathbf{X}^{-1})^H = (\mathbf{X}^H)^H = \mathbf{X} = \mathbf{U}$, so that the left and right eigenvectors coincide.

Numerically finding eigenvalues

For a matrix $A \in \mathbb{C}^{n \times n}$ (potentially real), we want to find $\lambda \in \mathbb{C}$ and $\mathbf{x} \neq 0$ such that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

Most relevant problems:

- ▶ A symmetric (and large)
- ▶ A spd (and large)
- ▶ A stochastic matrix, i.e., all entries $0 \leq a_{ij} \leq 1$ are probabilities, and thus $\sum_j a_{ij} = 1$.

How hard are they to find numerically?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm \rightsquigarrow this is fundamentally different from what we have seen previously when solving systems of *linear* equations! These algorithms (LU, QR) give the *exact* solution in *exact* arithmetic in finite number of steps. We *cannot* expect something similar for computing eigenvalues!

Condition of finding eigenvalues of a matrix

The absolute condition number of determining a simple eigenvalue λ_0 of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with respect to the $\|\cdot\|_2$ is

$$\underline{\kappa_{\text{abs}}} = \frac{1}{|\cos(\angle(\mathbf{x}, \mathbf{y}))|}, \quad \cos(\angle(\mathbf{x}, \mathbf{y})) = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

and the relative condition number is

$$\underline{\kappa_{\text{rel}}} = \frac{\|\mathbf{A}\|}{|\underline{\lambda_0} \cos(\angle(\mathbf{x}, \mathbf{y}))|},$$

where \mathbf{x} is an eigenvector of \mathbf{A} for the eigenvalue λ_0 ($\mathbf{A}\mathbf{x} = \lambda_0\mathbf{x}$) and \mathbf{y} an adjoint eigenvector ($\mathbf{A}^H\mathbf{y} = \bar{\lambda}_0\mathbf{y}$).

Sketch of proof \rightsquigarrow board

(see also Deufhard, Theorem 5.2)

$$Ax = \lambda_0 x, \quad \lambda_0 \text{ is simple}$$

$A \in \mathbb{C}^{n \times n}$, n eigenvalues, $d \leq n$ are distinct

$$\det(A - \lambda I) = \prod_{i=1}^n (\lambda_i - \lambda) = \prod_{i=1}^d (\lambda'_i - \lambda)^{\mu(\lambda'_i)}$$

$\mu(\lambda'_i)$... algebraic multiplicity

Geometric multiplicity: $\#$ linearly indep eigenvectors
 $g(\lambda_i)$ associated with eigenvalue λ_i

$$1 \leq g(\lambda_i) \leq \mu(\lambda_i) \leq n$$

Simple means

$$\mu(\lambda_0) = 1$$

]

We want to show

$$\lambda: A \mapsto \lambda(A)$$

exists, cont. diff. in neighborhood of A

$$\lambda(A) = \lambda_0$$

$$\lambda'(A)C = \frac{\langle Cx_0, \bar{\lambda}_0 \rangle}{\langle x_0, y_0 \rangle}$$

x_0 is eigenvector of A with λ_0

y_0 is eigenvector of A^H with $\bar{\lambda}_0$

-

$$A + \epsilon C$$

$\epsilon \dots$ small

$$C \in \mathbb{C}^{n \times n}$$

Implicit function theorem

$$F(x, \lambda) = Ax - \lambda x$$

$$\frac{\partial}{\partial \lambda} F(x, \lambda) = -x$$

$$\frac{\partial}{\partial x} F(x, \lambda) = A - \lambda I$$

$$J_F(x, \lambda) = [x \mid A - \lambda I] \in \mathbb{C}^{n \times (n+1)}$$

"geo. multi \leq alg. multi" λ_0 simple

$$\ker(A - \lambda_0 I) = \underline{\text{span}\{x_0\}}$$

$$\Rightarrow \text{rank}(A - \lambda_0 I) = \underline{n-1}$$

$$\boxed{x_0 \notin \text{colspan } A - \lambda_0 I}$$

$$J_F = [x_0 \mid A - \lambda_0 I]$$

$$\text{rank}(J_F) = \underline{n}$$

\Rightarrow full rank J_F

\Rightarrow apply implicit function theorem

$$A + tC$$

$$\lambda: A + tC \mapsto \lambda(A + tC)$$

$$\lambda(0) = \lambda_0$$

λ is C^1 over small $(-\varepsilon, \varepsilon)$

$$\underline{\lambda'(A)} : \underline{C} \mapsto \frac{\langle Cx_0, y_0 \rangle}{\langle x_0, y_0 \rangle}$$

$$\| \lambda'(A) \| = \sup_{C \neq 0} \frac{|\langle Cx_0, y_0 \rangle / \langle x_0, y_0 \rangle|}{\|C\|}$$

$$= \dots = \frac{1}{|\cos(\angle(x_0, y_0))|} = h_{\text{obs}}$$

rel.

$$h_{\text{rel}} = \frac{\|A\|}{\|y\|} \frac{1}{|\cos(\angle(x_0, y_0))|}$$

Interpretation

Perturbations of order δ in entries of matrix \mathbf{A} induce changes of the order

$$\delta\lambda = \delta / \cos(\angle(\mathbf{x}_0, \mathbf{y}_0))$$

In particular, for normal matrices* ($\mathbf{A}\mathbf{A}^H = \mathbf{A}^H\mathbf{A}$), we have $\mathbf{x}_0 = \mathbf{y}_0$ and thus $\angle(\mathbf{x}_0, \mathbf{y}_0) = 0$ and thus $\cos(\angle(\mathbf{x}_0, \mathbf{y}_0)) = 1$, which means $\kappa_{\text{abs}} = 1$, which can be considered well conditioned

Finding non-simple eigenvalues can have very high absolute condition number (but can still be done numerically). For a detailed treatment have a look at textbook by Golub et al. on Matrix Computations.

*Equivalent: Have orthonormal eigenbasis of \mathbb{C} ; diagonalizable by unitary matrix.

Bounding error in eigenvalue computation

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a Hermitian matrix and let $(\hat{\lambda}, \hat{\mathbf{x}})$ be a computed approximation of an eigenvalue/eigenvector pair (λ, \mathbf{x}) of \mathbf{A} . Defining the residual

$$\underline{\hat{\mathbf{r}}} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \quad \hat{\mathbf{x}} \neq \mathbf{0},$$

it then follows that

$$\min_{\lambda_i \in \sigma(\mathbf{A})} |\hat{\lambda} - \lambda_i| \leq \frac{\|\hat{\mathbf{r}}\|_2}{\|\hat{\mathbf{x}}\|_2},$$

where $\sigma(\mathbf{A}) = \{\lambda | \lambda \text{ is an eigenvalue of } \mathbf{A}\}$ is the spectrum of \mathbf{A} .

Proof \rightsquigarrow board

What is special about this bound?

Bounding error in eigenvalue computation

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a Hermitian matrix and let $(\hat{\lambda}, \hat{\mathbf{x}})$ be a computed approximation of an eigenvalue/eigenvector pair (λ, \mathbf{x}) of \mathbf{A} . Defining the residual

$$\hat{\mathbf{r}} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \quad \hat{\mathbf{x}} \neq \mathbf{0},$$

it then follows that

$$\min_{\lambda_i \in \sigma(\mathbf{A})} |\hat{\lambda} - \lambda_i| \leq \frac{\|\hat{\mathbf{r}}\|_2}{\|\hat{\mathbf{x}}\|_2},$$

where $\sigma(\mathbf{A}) = \{\lambda | \lambda \text{ is an eigenvalue of } \mathbf{A}\}$ is the spectrum of \mathbf{A} .

Proof \rightsquigarrow board

What is special about this bound?

- ▶ This is an *a posteriori* bound that bounds the error *after* we have computed the result
- ▶ We will see many more residual-based *a posteriori* bounds (broadly speaking: the residual is something we can compute, and if the problem is “well-behaved” then the norm of the residual is a reasonable bound of the norm of the error.)

exists orthogonal eigenbasis $\{v_i\}$

$$\hat{x} = \sum_{i=1}^n \alpha_i v_i$$

$$\alpha_i = v_i^H \hat{x}$$

$$\hat{r} = A \hat{x} - \hat{\lambda} \hat{x} = A \left(\sum_{i=1}^n \alpha_i v_i \right) - \hat{\lambda} \sum_{i=1}^n \alpha_i v_i$$

$$= \sum \alpha_i \lambda_i v_i - \hat{\lambda} \sum \alpha_i v_i$$

$$= \sum \alpha_i (\lambda_i - \hat{\lambda}) v_i$$

$$\frac{\|\hat{r}\|^2}{\|\hat{x}\|^2} = \frac{\sum |\alpha_i|^2 (\lambda_i - \hat{\lambda})^2 \|v_i\|^2}{\sum |\alpha_i|^2 \|v_i\|^2}$$

$$\|v_i\|^2 = 1 \quad \Rightarrow \quad \left(\sum_{i=1}^n \frac{|\alpha_i|^2}{\sum |\alpha_j|^2} \right) (\lambda_i - \hat{\lambda})^2$$

$$= \sum_{i=1}^n \beta_i (\lambda_i - \hat{\lambda})^2$$

$$\beta_i \geq 0, \quad \sum \beta_i = 1$$

$$\begin{aligned}
 \frac{\|r\|^2}{\|x\|^2} &= \sum \beta_i (\lambda_i - \bar{\lambda})^2 \geq \sum \beta_i \min_j (\lambda_j - \bar{\lambda})^2 \\
 &= \min_j (\lambda_j - \bar{\lambda})^2
 \end{aligned}$$

Condition of computing eigenvectors

- ▶ The condition of computing eigenvector \mathbf{x}_i for an eigenvalue λ_i depends on the separation between the eigenvalues

$$\kappa = \frac{1}{\min_{i \neq j} |\lambda_i - \lambda_j|}$$

(Quarteroni et al., Section 5)

- ▶ Computing \mathbf{x}_i can be ill-conditioned if some eigenvalue λ_j is “very close” to the eigenvalue λ_i
- ▶ This indicates that multiple eigenvalues require care. Even for Hermitian matrices eigenvectors can be hard to compute