

Deck 1: A Primer on Numerical Linear Algebra

Math 7870: Topics in Randomized Numerical Linear Algebra

Spring 2026

Akil Narayan

Notation

Some consistent notation we'll use:

- Non-negative integers: $n, m, p, r \in \mathbb{N}_0$
- Index sets: $[n] = \{1, 2, \dots, n\}$
- Vectors: $\mathbf{v} \in \mathbb{C}^n$ (often we'll specialize to \mathbb{R}^n for notational simplicity)
- Matrices: $\mathbf{A} \in \mathbb{C}^{m \times n}$
- Slicing vectors: Given, $S \in [n]$, then $\mathbf{v}_S \in \mathbb{C}^{|S|}$ is the S -sliced entries from \mathbf{v}
- Slicing matrices: Given, $S \in [n]$, then $\mathbf{A}_{*S} \in \mathbb{C}^{m \times |S|}$ is the S -sliced columns from \mathbf{A}
- Row slices: \mathbf{A}_{T*} for $T \subset [m]$. The matrix \mathbf{A}_{TS} also makes sense.
- Matrix (conjugate) transpose/adjoint, determinant, trace: $\mathbf{A}^* \in \mathbb{C}^{n \times m}$, $\det(\mathbf{A})$, $\text{tr}(\mathbf{A})$.
- Standard inner product on vectors: $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{w}^* \mathbf{v} \in \mathbb{C}$.
- Vectors are orthogonal (in ℓ^2) if $\langle \mathbf{v}, \mathbf{w} \rangle = 0$.

Some basic properties:

- 1-element slices produces entries: $S = \{i\}$ means $\mathbf{v}_S = v_i$. $T = \{j\}$ means $\mathbf{A}_{TS} = A_{ji}$.
- $|\det \mathbf{A}| = 1$ iff \mathbf{A} has orthonormal columns.
- Norms $\|\cdot\|$ on vectors or matrices are non-negative, order-1 positively homogeneous, convex functions with trivial zero level set. (Cf. the standard ℓ^p norm.)

Matrix classifications and structure

- Square/rectangular/wide/tall
- (skew-)Hermitian
- Unitary/orthogonal
- Normal
- Diagonalizable
- Sparse
- (Orthogonal) projectors
- Circulant

Matrix structure is very useful

E.g.: A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is circulant if it satisfies,

$$\left(\begin{array}{ccccc} a_1 & a_2 & a_3 & a_4 & a_5 \\ a_5 & a_1 & a_2 & a_3 & a_4 \\ a_4 & a_5 & a_1 & a_2 & a_3 \\ a_3 & a_4 & a_5 & a_1 & a_2 \\ a_2 & a_3 & a_4 & a_5 & a_1 \end{array} \right) \quad \mathbf{A} = \left(\begin{array}{ccccc} a_1 & a_n & a_{n-1} & \cdots & a_2 \\ a_2 & a_1 & a_n & \cdots & a_3 \\ \ddots & \ddots & \ddots & \ddots & \ddots \\ a_n & a_{n-1} & a_{n-2} & \cdots & a_1 \end{array} \right), \quad \text{i.e.: } \begin{cases} \mathbf{A}_{*,i+1} = \mathbf{P} \mathbf{A}_{*,i}, & i \in [n-1] \\ \mathbf{A}_{*,1} = \mathbf{P} \mathbf{A}_{*,n}. \end{cases}$$

where \mathbf{P} is the down-shift permutation matrix:

$$\mathbf{P} = (\mathbf{e}_2 \ \mathbf{e}_2 \ \mathbf{e}_3 \ \dots \ \mathbf{e}_n \ \mathbf{e}_1) = \begin{pmatrix} & & & & 1 \\ & & & & \\ 1 & & & & \\ & & & & 1 \\ & & & \ddots & \\ & & & & 1 \end{pmatrix}, \quad \mathbf{P} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_{n-1} \\ v_n \end{pmatrix} = \begin{pmatrix} v_n \\ v_1 \\ \vdots \\ v_{n-2} \\ v_{n-1} \end{pmatrix}$$

The punchline: circulant matrices are diagonalizable by the discrete Fourier transform!
(Why is this useful...?)

Why focus on linear algebra?

The canonical reasons:

- Solving linear systems
- Vector manipulation (e.g., orthogonalizing, identifying subspaces)
- Constructing operators: projectors, discretization of differential operators

Some less obvious reasons:

- Compression, dimension reduction
- High-order optimization schemes
- Low-rank approximations
- Training for AI/LLM/neural architectures

There are “classical” applications where linear algebra is known to be useful (physics-based modeling, classical statistics, computer vision/animation).

But it’s quite hard to undersell the importance of (numerical) linear algebra in modern computational applications (data science/analysis, machine learning, language models).

Computational building blocks

Linear algebra (LA), roughly, studies the mathematics of objects in linear (sub)spaces, which include operators.
(i.e., vectors, matrices, and linear-type operations involving these)

Numerical linear algebra (NLA), roughly, studies computational methods/algorithms that one actually uses to accomplish some linear algebraic operations.
(e.g., computing the determinant of a matrix)

The core, classical components of NLA are matrix decompositions.

Matrix decompositions/factorizations

(LU factorization)

$$A = \left(\begin{array}{c} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{array} \right) \left(\begin{array}{c} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{array} \right)$$

(QR factorization)

$$A = \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right) \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right)$$

(Eigenvalue decomposition)

$$A \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right) = \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right) \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right)$$

(Singular value decomposition)

$$A = \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right) \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right) \left(\begin{array}{c} \text{blue} \\ \text{orange} \\ \text{green} \\ \text{grey} \\ \text{red} \end{array} \right)$$

(P)LU

Given $\mathbf{A} \in \mathbb{C}^{m \times n}$, its LU decomposition is,

$$\mathbf{A} = \mathbf{L}\mathbf{U}, \quad \mathbf{L} \in \mathbb{C}^{m \times p}, \quad \mathbf{U} \in \mathbb{C}^{p \times n}, \quad p = \min\{m, n\},$$

where \mathbf{L} is *lower triangular*, and \mathbf{U} is *upper triangular*.

The LU decomposition need not exist for an arbitrary matrix \mathbf{A} .

(It exists iff $\det \mathbf{A}_{[q],[q]} \neq 0$ for all $q \in [p]$.)

A more general LU-type decomposition is the *pivoted* LU decomposition. This is one of the decompositions:

$$\mathbf{A}\mathbf{P}_2 = \mathbf{L}\mathbf{U}, \quad \mathbf{P}_2 \in \mathbb{C}^{n \times n},$$

$$\mathbf{P}_1\mathbf{A} = \mathbf{L}\mathbf{U}, \quad \mathbf{P}_1 \in \mathbb{C}^{m \times m},$$

$$\mathbf{P}_1\mathbf{A}\mathbf{P}_2 = \mathbf{L}\mathbf{U}, \quad \mathbf{P}_1 \in \mathbb{C}^{m \times m}, \quad \mathbf{P}_2 \in \mathbb{C}^{n \times n},$$

where both \mathbf{P}_1 and \mathbf{P}_2 are row and column *permutation* matrices, respectively.

(A permutation matrix is a unitary matrix, where each row/column is a cardinal unit vector.)

Pivoted LU decompositions exist for any matrix.

(P)LU – in practice and uses

The computation of the LU decomposition is “easy”: it’s Gaussian elimination.

(This is actually how the direct computation is done on a computer.)

The permutations are chosen iteratively so that elements in pivot locations have as large a magnitude as possible.

The most popular pivoting strategy is row(-only) pivoting.

Why compute an LU decomposition?

- This is how invertible linear systems are solved: $\mathbf{Ax} = \mathbf{b} \longrightarrow \mathbf{x} = \mathbf{U}^{-1} \mathbf{L}^{-1} \mathbf{P}^* \mathbf{b}$.
- This is how “simultaneous” systems are solved (and how matrix inverses are computed).
- This is how determinants are computed: $\det \mathbf{A} = (\det \mathbf{P}_1)(\det \mathbf{L})(\det \mathbf{U})$.
- This is how (quasi-optimal) low-rank approximations are built. (“skeletonization”, “empirical interpolation”)
- If \mathbf{A} is Hermitian and positive-definite, then its LU decomposition ($\mathbf{U} = \mathbf{L}^*$) is called the *Cholesky decomposition*, which is quite useful when working with these classes of matrices.

QR

Given $\mathbf{A} \in \mathbb{C}^{m \times n}$, its QR decomposition is,

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \quad \mathbf{Q} \in \mathbb{C}^{m \times m}, \quad \mathbf{R} \in \mathbb{C}^{m \times n},$$

where \mathbf{Q} is unitary ($\mathbf{Q}^*\mathbf{Q} = \mathbf{I}_m$), and \mathbf{R} is upper triangular.

There is also an “economical”/“thin” QR decomposition, mainly useful when $m > n$, which truncates columns of \mathbf{Q} corresponding to all-zero rows of \mathbf{R} .

There's also a (column-)pivoted version of QR:

$$\mathbf{AP} = \mathbf{Q}\mathbf{R}, \quad \mathbf{P} \in \mathbb{C}^{m \times m},$$

where \mathbf{P} is chosen to ensure that the diagonal elements of \mathbf{R} are non-decreasing.

QR – in practice and uses

Explicit computation of the QR decomposition is “just” orthogonalizing the column vectors of \mathbf{A} . Gram-Schmidt orthogonalization can do this – but this is a numerically unstable procedure.

“Modified” Gram-Schmidt fixes the instability, but is not really used in practice: there are procedures that employ a sequence of *unitary* transforms to compute the QR decomposition:

- Householder reflectors
- Givens rotations

This is how any modern implementation of the QR decomposition works in practice.

Why compute a QR decomposition?

- To orthogonalize vectors and/or compute orthogonal projection matrices
- This is how one could solve invertible linear systems: $\mathbf{A}\mathbf{x} = \mathbf{b} \longrightarrow \mathbf{x} = \mathbf{R}^{-1}\mathbf{Q}^*\mathbf{b}$.
- This is how one solves (linear) least squares problems: $\mathbf{A}\mathbf{x} = \mathbf{b} \longrightarrow \mathbf{x} = \mathbf{R}^{-1}(\mathbf{Q}_{*,[n]})^*\mathbf{b}$.
- This is a core ingredient in computing eigenvalues (!).

Eigenvalue decompositions

Given $\mathbf{A} \in \mathbb{C}^{n \times n}$, its eigenvalue decomposition is given by,

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^{-1}, \quad \mathbf{V}, \Lambda \in \mathbb{C}^{n \times n},$$

where Λ is a diagonal matrix. The diagonal elements of Λ are the *eigenvalues* of \mathbf{A} , and the columns of \mathbf{V} are the corresponding *eigenvectors* of \mathbf{A} .

Square matrices having an eigenvalue decomposition are *diagonalizable*. Not all square matrices are diagonalizable. (But “most” are.)

Diagonalizable matrices are “just” diagonal matrices, when the mapping \mathbf{A} is represented in the right coordinate system.

There is a special class of matrices for which even more is true: \mathbf{A} is normal if $\mathbf{A}\mathbf{A}^* = \mathbf{A}^*\mathbf{A}$. A matrix is normal iff it's *unitarily diagonalizable*, i.e., is diagonalizable with \mathbf{V} a unitary matrix.

Normal matrices are diagonal matrices, when the input/output coordinates are simply rotated/reflected.

(For example, Hermitian and skew-Hermitian matrices are normal.)

Eigenvalue decompositions – in practice and uses

Computing eigenvalues/eigenvectors is a big (+ difficult) business.

When \mathbf{A} is normal, the business is not too bad. (Roughly speaking, the spectrum can be computed from \mathbf{A} through a well-conditioned operation.)

The so-called “non-symmetric” eigenvalue problem is really hard – computing eigenvalues can be “arbitrarily difficult”.

Here's a vague sense of how QR is used to compute eigenvalues:

- Suppose we knew \mathbf{V} .
- Compute $\mathbf{V} = \mathbf{Q}\mathbf{R}$.
- How could we compute the spectrum of $\mathbf{Q}^* \mathbf{A} \mathbf{Q}$?

Why compute eigenvalues?

- Transforming a complicated matrix to a diagonal makes it easy to understand what the matrix is “doing”.
- When \mathbf{A} is Hermitian and positive-definite, the spectrum tells us a *lot* about how to perform compression and low-rank approximation.
- We can compute singular values....

An interlude: Hermitian positive (semi)-definite matrices

Hermitian matrices $\mathbf{A} \in \mathbb{C}^{n \times n}$ that are *positive semi-definite*, i.e., $\mathbf{x}^* \mathbf{A} \mathbf{x} \geq 0$ for all \mathbf{x} , are so important that they deserve their own discussion.

The standard abbreviation for these matrices is “SPD”. An SPD matrix \mathbf{A} has n real eigenvalues:

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \geq 0.$$

These eigenvalues are informative about \mathbf{A} : $\|\mathbf{A}\|_2 = \lambda_1$, and \mathbf{A} is singular iff $\lambda_n = 0$.

SPD matrices are ubiquitous, with the simplest examples being covariance matrices, (graph) Laplacian matrices, kernel matrices, Gram matrices, ...

The cone of Hermitian matrices has a useful partial ordering, the *Loewner order*:

$$\mathbf{A} \leq \mathbf{B} \iff \mathbf{B} - \mathbf{A} \text{ is SPD}$$

This ordering plays a nice role in matrix functional analysis: A sensible way to define a matrix function is through its spectral decomposition:

$$f : \mathbb{R} \rightarrow \mathbb{R}, \quad \mathbf{A} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^*, \quad f(\mathbf{A}) := \mathbf{U} f(\boldsymbol{\Lambda}) \mathbf{U}^* \quad (f(\boldsymbol{\Lambda}) \text{ diagonal, componentwise evaluation}).$$

An interesting question: are there *operator monotone* functions on SPD matrices?
I.e., if $\mathbf{A}, \mathbf{B} \succeq 0$ and $\mathbf{A} \leq \mathbf{B}$, is $f(\mathbf{A}) \leq f(\mathbf{B})$?

SVD

Given $\mathbf{A} \in \mathbb{C}^{m \times n}$, its singular value decomposition (SVD) is,

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*, \quad \mathbf{U} \in \mathbb{C}^{m \times m}, \quad \Sigma \in \mathbb{R}^{m \times n}, \quad \mathbf{V} \in \mathbb{C}^{n \times n},$$

where both \mathbf{U} and \mathbf{V} are unitary, and Σ is diagonal with its non-negative entries arranged in non-increasing order.

The diagonal elements of Σ are the *singular values*, $\sigma_1, \sigma_2, \dots, \sigma_p$, with $p = \min\{m, n\}$.

The columns of \mathbf{U} and \mathbf{V} are the left- and right-singular vectors of \mathbf{A} , respectively.

Sometimes, the SVD is *truncated* to remove unnecessary columns/rows from \mathbf{U}/\mathbf{V}^* , e.g., those corresponding to zero singular values.

Generally, singular values and eigenvalues are unrelated.

E.g., a matrix can have all zero eigenvalues but some non-zero singular values.

Similarly, singular vectors and eigenvectors are generally unrelated.

The SVD and the eigenvalue decomposition of \mathbf{A} coincide iff \mathbf{A} is Hermitian and positive semi-definite.

SVD – in practice and uses

The SVD is computed through an eigenvalue decomposition of the (normal!) matrix $\mathbf{A}^* \mathbf{A}$:

$$\mathbf{A} \mathbf{A}^* = \mathbf{U} \Lambda \mathbf{U}^*, \quad \Lambda = \Sigma^2, \quad \mathbf{V}^* = \Sigma^{-1} \mathbf{U}^* \mathbf{A}$$

For simplicity, we've assumed above that Σ is invertible. (Σ is invertible iff \mathbf{A} has full rank.)

Therefore: computing eigenvalues of Hermitian matrices is sufficient to allow us to compute the SVD.

Why compute the SVD?

- The SVD tells you “almost everything” you need to know about a matrix: its (co)range/(co)kernel, its rank, its norm for any unitarily invariant norm.
- One can solve linear systems, solve least-squares problems, compute determinants, determine matrix rank, etc.
- Truncated SVDs are *optimal* low-rank approximations of \mathbf{A} in any unitarily invariant norm, and norms of sequence of truncated singular value are corresponding low-rank approximation errors.

Software and implementation

The core routines for computing essentially everything we've discussed is standardized through existing, nearly bulletproof software:

- (BLAS) Basic linear algebra subroutines: low-level addition/subtraction of vectors, scalar multiplication, matrix-vector and matrix-matrix multiplication.
- (LAPACK) The linear algebra package: eigenvalues, linear solvers, matrix factorization (LU, QR, Cholesky, EVD, SVD)

These are *modular* routines: LAPACK uses BLAS for its building block.

E.g., GPU-enabled implementations of dense linear algebra factorizations can, “in principle”, be implemented by recoding just the BLAS portion without changing LAPACK. (Cf. cuBLAS)

This is all great: what's the problem?

All of the above is *dense* linear algebra: matrices are stored and manipulated as explicit arrays, where each array entry is explicitly stored and arithmetically exercised.

Modern practice and bottlenecks

The problem: We always want our hammer to smack bigger, tougher nails.

Using these established routines:

- For $n \times n$ matrices, matmat multiplications have $\mathcal{O}(n^3)$ complexity, matvecs are $\mathcal{O}(n^2)$.
- For $\mathbf{A} \in \mathbb{C}^{m \times n}$, $m \geq n$, computing its LU, QR, or SVD factorizations requires $\mathcal{O}(mn^2)$ effort, with $\mathcal{O}(mn)$ memory.
- For $\mathbf{A} \in \mathbb{C}^{n \times n}$, solving a corresponding linear system requires $\mathcal{O}(n^3)$ effort, and $\mathcal{O}(n^2)$ memory.

These are (effectively) optimal complexities for these operations.

For context:

- Datasets often feature $m, n \sim 10^6$.
- ML architectures routinely have 10^6 parameters. (We need to compute matvecs with these parameters, store/use gradients, or maybe even Hessians?)
- In physics-based models, we'd *love* to routinely solve linear systems with $n \gg 10^{10}$.
- Even if matrices/vectors are *sparse* and you can get away with storing them with reasonable memory requirements, many intermediate NLA quantities are *dense*.
- In general, matrices/vectors do not have special structure/sparsity that we can exploit.

This is the problem: the scale of modern problems defy naive usage of dense NLA procedures.

The potential, pitfalls, and promise of *randomized* NLA

PSA: There are no free lunches here.

If we are willing to give up nothing, we cannot really gain much more.

But there are some things we might be willing to give up:

- Maybe we only need an answer to low precision
- Maybe we only need some part of an answer
- Maybe we're willing to cope with a possibility, hopefully minute, of algorithmic failure

It is here where *randomized* methods hold promise:

- I can get “close” to the answer by randomly approximating the problem
- I can get “part” of the answer by randomly compressing the problem
- Randomness entails some possibility of failure

The foundation of randomized NLA methods is laid in probability and statistics....