# C6.1 Numerical Linear Algebra

## Yuji Nakatsukasa*

Welcome to numerical linear algebra (NLA)! NLA is a beautiful subject that combines mathematical rigor, amazing algorithms, and an extremely rich variety of applications.

What is NLA? In a sentence, it is a subject that deals with the numerical solution (i.e., using a computer) of linear systems $Ax = b$ (given $A \in \mathbb{R}^{n \times n}$ (i.e., a real $n \times n$ matrix) and $b \in \mathbb{R}^n$ (real $n$-vector), find $x \in \mathbb{R}^n$) and eigenvalue problems $Ax = \lambda x$ (given $A \in \mathbb{R}^{n \times n}$, find $\lambda \in \mathbb{C}$ and $x \in \mathbb{C}^n$), for problems that are too large to solve by hand ($n \geq 4$ is already large; we aim for $n$ in the thousands or even millions). This can rightfully sound dull, and some mathematicians (those purely oriented?) tend to get turned off after hearing this — how could such a course be interesting compared with other courses offered by the Oxford Mathematical Institute? I hope and firmly believe that at the end of the course you will all agree that there is more to the subject than you imagined. The rapid rise of data science and machine learning has only meant that the importance of NLA is still growing, with a vast number of problems in these fields requiring NLA techniques and algorithms. It is perhaps worth noting also that these fields have had enormous impact on the direction of NLA, in particular the recent and very active field of randomized algorithm was born in light of needs arising from these extremely active fields.

In fact NLA is a truly exciting field that utilises a huge number of ideas from different branches of mathematics (e.g. matrix analysis, approximation theory, and probability) to solve problems that actually matter in real-world applications. Having said that, the number of prerequisites for taking the course is the bare minimum; essentially a basic understanding of the fundamentals of linear algebra would suffice (and the first lecture will briefly review the basic facts). If you've taken the Part A Numerical Analysis course you will find it helpful, but again, this is not necessary.

The field NLA has been blessed with many excellent books on the subject. These notes will try to be self-contained, but these references will definitely help. There is a lot to learn; literally as much as you want to.

- Trefethen-Bau (97) [33]: Numerical Linear Algebra

    - covers essentials, beautiful exposition

- Golub-Van Loan (12) [14]: Matrix Computations

---

*Last update: December 1, 2021. Please report any corrections or comments on these lecture notes to `nakatsukasa@maths.ox.ac.uk`

- classic, encyclopedic

- Horn and Johnson (12) [20]: Matrix Analysis (& Topics in Matrix Analysis (86) [19])

    - excellent theoretical treatise, little numerical treatment

- J. Demmel (97) [8]: Applied Numerical Linear Algebra

    - impressive content

- N. J. Higham (02) [17]: Accuracy and Stability of Algorithms

    - bible for stability, conditioning

- H. C. Elman, D. J. Silvester, A. J. Wathen (14) [11]: Finite Elements and Fast Iterative Solvers

    - PDE applications of linear systems, Krylov methods and preconditioning

This course covers the fundamentals of NLA. We first discuss the singular value decomposition (SVD), which is a fundamental matrix decomposition whose importance is only growing. We then turn to linear systems and eigenvalue problems. Broadly, we will cover

- Direct methods ($n \lesssim 10{,}000$): Sections 5–10 (except 8)

- Iterative methods ($n \lesssim 1{,}000{,}000$, sometimes larger): Sections 11–13

- Randomized methods ($n \gtrsim 1{,}000{,}000$): Sections 14–16

in this order. Lectures 1–4 cover the fundamentals of matrix theory, in particular the SVD, its properties and applications.

This document consists of 16 sections. Very roughly speaking, one section corresponds to one lecture (though this will not be followed strictly at all).

# Contents

*Notation.* For convenience below we list the notation that we use throughout the course.

- $\lambda(A)$: the set of eigenvalues of $A$. If a natural ordering exists (e.g. $A$ is symmetric so $\lambda$ is real), $\lambda_i(A)$ is the $i$th (largest) eigenvalue.

- $\sigma(A)$: the set of singular values of $A$. $\sigma_i(A)$ always denotes the $i$th largest singular value. We often just write $\sigma_i$.

- $\mathrm{diag}(A)$: the vector of diagonal entries of $A$.

- We use capital letters for matrices, lower-case for vectors and scalars. Unless otherwise specified, $A$ is a given matrix, $b$ is a given vector, and $x$ is an unknown vector.

- $\|\cdot\|$ denotes a norm for a vector or matrix. $\|\cdot\|_2$ denotes the spectral (or 2-) norm, $\|\cdot\|_F$ the Frobenius norm. For vectors, to simplify notation we sometimes use $\|\cdot\|$ for the spectral norm (which for vectors is the familiar Euclidean norm).

- Span(A) denotes the span or range of the column space of $A$. This is the subspace consisting of vectors of the form $Ax$.

- We reserve $Q$ for an orthonormal (or orthogonal) matrix. $L, (U)$ are often lower (upper) triangular.

- $I$ always denotes the identity matrix. $I_n$ is the $n \times n$ identity when the size needs to be specified.

- $A^T$ is the transpose of the matrix; $(A^T)_{ij} = A_{ji}$. $A^*$ is the (complex) conjugate transpose $(A^*)_{ij} = \bar{A}_{ji}$.

- $\succ, \succeq$ denote the positive (semi)definite ordering. That is, $A \succ (\succeq)0$ means $A$ is positive (semi)definite (abbreviated as PD, PSD), i.e., symmetric and with positive (nonnegative) eigenvalues. $A \succ B$ means $A - B \succ 0$.

We sometimes use the following shorthand: alg for algorithm, eigval for eigenvalue, eigvec for eigenvector, singval for singular value, and singvec for singular vector, iff for "if and only if".

# 0   Introduction, why $Ax = b$ and $Ax = \lambda x$?

As already stated, NLA is the study of numerical algorithms for problems involving matrices, and there are only two main problems(!):

1. Linear system
$$Ax = b.$$

   Given a (often square $m = n$ but we will discuss $m > n$ extensively, and $m < n$ briefly at the end) matrix $A \in \mathbb{R}^{m \times n}$ and vector $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ such that $Ax = b$.

2. Eigenvalue problem
$$Ax = \lambda x.$$

   Given a (always![1]) square matrix $A \in \mathbb{R}^{n \times n}$ find $\lambda$: eigenvalues (eigval), and $x \in \mathbb{R}^n$: eigenvectors (eigvec).

We'll see many variants of these problems; one worthy of particular mention is the SVD, which is related to eigenvalue problems but given its ubiquity has a life of its own. (So if there's a third problem we solve in NLA, it would definitely be the SVD.)

It is worth discussing *why* we care about linear systems and eigenvalue problems.

The primary reason is that many (in fact most) problems in scientific computing (and even machine learning) boil down to linear problems:

- Because that's often the only way to deal with the scale of problems we face today! (and in future)

- For linear problems, so much is understood and reliable algorithms are available[2].

---

[1] There are exciting recent developments involving eigenvalue problems for rectangular matrices, but these are outside the scope of this course.

[2] A pertinent quote is Richard Feynman's "Linear systems are important because we can solve them". Because we can solve them, we do all sorts of tricks to reduce difficult problems to linear systems!

A related important question is where and how these problems arise in real-world problems.

Let us mention a specific context that is relevant in data science: optimisation. Suppose one is interested in minimising a high-dimensional real-valued function $f(x) : \mathbb{R}^n \to \mathbb{R}$ where $n \gg 1$.

A successful approach is to try and find critical points, that is, points $x_*$ where $\nabla f(x_*) = 0$. Mathematically, this is a non-linear high-dimensional root-finding problem of finding $x \in \mathbb{R}^n$ such that $\nabla f(x) =: F(x) = 0$ (the vector $0 \in \mathbb{R}^n$) where $F : \mathbb{R}^n \to \mathbb{R}^n$. One of the most commonly employed methods for this task is Newton's method (which some of you have seen in Prelims Constructive Mathematics). This boils down to

- Newton's method for $F(x) = 0$, $F : \mathbb{R}^n \to \mathbb{R}^n$ nonlinear:

    1. Start with initial guess $x^{(0)} \in \mathbb{R}^n$, set $i = 0$
    2. Find Jacobian matrix $J \in \mathbb{R}^{n \times n}$, $J_{ij} = \frac{\partial F_i(x)}{\partial x_j}\big|_{x=x^{(0)}}$
    3. Update $x^{(i+1)} := x^{(i)} - J^{-1}F(x^{(i)})$, $i \leftarrow i + 1$, go to step 2 and repeat

    Note that the main computational task is to find the vector $y = J^{-1}F(x^{(i)})$, which is a linear system $Jy = F(x^{(i)})$ (which we solve for the vector $y$)

What about eigenvalue problems $Ax = \lambda x$? Google's pagerank is a famous application (we will cover this if we have time). Another example the Schrödinger equation of physics and chemistry. Sometimes a nonconvex optimisation problem can be solved by an eigenvalue problem.

Equally important is principal component analysis (PCA), which can be used for data compression. This is more tightly connected to the SVD.

Other sources of linear algebra problems include differential equations, optimisation, regression, data analysis, ...

# 1 Basic LA review

We start with a review of key LA facts that will be used in the course. Some will be trivial to you while others may not. You might also notice that some facts that you have learned in a core LA course will not be used in this course. For example we will never deal with finite fields, and determinants only play a passing role.

## 1.1 Warmup exercise

Let $A \in \mathbb{R}^{n \times n}$ ($n \times n$ square matrix). (or $\mathbb{C}^{n \times n}$; the difference hardly matters in most of this course[3]). Try to think of statements that are equivalent to $A$ being nonsingular. Try to come up with as many conditions as possible, before turning the page.

---

[3]While there are a small number of cases where the distinction between real and complex matrices matters, in the majority of cases it does not, and the argument carries over to complex matrices by replacing $\cdot^T$ with

Here is a list: The following are equivalent.

1. $A$ is nonsingular.

2. $A$ is invertible: $A^{-1}$ exists.

3. The map $A : \mathbb{R}^n \to \mathbb{R}^n$ is a bijection.

4. All $n$ eigenvalues of $A$ are nonzero.

5. All $n$ singular values of $A$ are positive.

6. $\operatorname{rank}(A) = n$.

7. The rows of $A$ are linearly independent.

8. The columns of $A$ are linearly independent.

9. $Ax = b$ has a solution for every $b \in \mathbb{C}^n$.

10. $A$ has no nonzero null vector.

11. $A^T$ has no nonzero null vector.

12. $A^*A$ is positive definite (not just semidefinite).

13. $\det(A) \neq 0$.

14. An $n \times n$ matrix $A^{-1}$ exists such that $A^{-1}A = I_n$. (this, btw, implies (iff) $AA^{-1} = I_n$, a nontrivial fact)

15. ... (what did I miss?)

## 1.2   Structured matrices

We will be discussing lots of structured matrices. For square matrices,

- Symmetric: $A_{ij} = A_{ji}$ (Hermitian: $A_{ij} = \bar{A}_{ji}$)

  - The most important property of symmetric matrices is the symmetric eigenvalue decomposition $A = V \Lambda V^T$; $V$ is orthogonal $V^T V = V V^T = I_n$, and $\Lambda$ is a diagonal matrix of eigenvalues $\Lambda = \operatorname{diag}(\lambda_1, \ldots, \lambda_n)$.
  - symmetric positive (semi)definite $A \succ (\succeq)0$: symmetric and all positive (nonnegative) eigenvalues.

---

.*. Therefore for the most part, we lose no generality in assuming the matrix is real (which slightly simplifies our mindset). Whenever necessary, we will highlight the subtleties that arise resulting from the difference between real and complex. (For the curious, these are the Schur form/decomposition, $LDL^T$ factorisation and eigenvalue decomposition for (real) matrices with complex eigenvalues.)

- Orthogonal: $AA^T = A^T A = I$ (Unitary: $AA^* = A^* A = I$). Note that for square matrices, $A^T A = I$ implies $AA^T = I$.

- Skew-symmetric: $A_{ij} = -A_{ji}$ (skew-Hermitian: $A_{ij} = -\bar{A}_{ji}$).

- Normal: $A^T A = AA^T$. (Here it's better to discuss the complex case $A^* A = AA^*$: this is a necessary and sufficient condition for diagonalisability under a unitary transformation, i.e., $A = U\Lambda U^*$ where $\Lambda$ is diagonal and $U$ is unitary.)

- Tridiagonal: $A_{ij} = 0$ if $|i - j| > 1$.

- Upper triangular: $A_{ij} = 0$ if $i > j$.

- Lower triangular: $A_{ij} = 0$ if $i < j$.

For (possibly nonsquare) matrices $A \in \mathbb{C}^{m \times n}$, (usually $m \geq n$).

- (upper) Hessenberg: $A_{ij} = 0$ if $i > j + 1$. (we will see this structure often.)

- "orthonormal": $A^T A = I_n$, and $A$ is (tall) rectangular. (This isn't an established name—we could call it "matrix with orthonormal columns" every time it appears—but we use these matrices all the time in this course, so we need a consistent shorthand name for it.)

- sparse: most elements are zero. $\mathrm{nnz}(A)$ denotes the number of nonzero elements in $A$. Matrices that are not sparse are called *dense*.

Other structures: Hankel, Toeplitz, circulant, symplectic,... (we won't use these in this course)

## 1.3  Matrix eigenvalues: basics

$$Ax = \lambda x, \quad A \in \mathbb{R}^{n \times n}, (0 \neq) x \in \mathbb{R}^n$$

- Example: $\begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

  This matrix has an *eigenvalue* $\lambda = 4$, with corresponding *eigenvector* $x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ (together they are an *eigenpair*).

  An $n \times n$ matrix always has $n$ eigenvalues (not always $n$ linearly independent eigenvectors); In the example above, $(\lambda, x) = \left( 1, \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \right), \left( 1, \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} \right)$ are also eigenpairs.

- The eigenvalues are the roots of the characteristic polynomial $\det(\lambda I - A) = 0$: $\det(\lambda I - A) = \prod_{i=1}^{n}(\lambda - \lambda_i)$.

- According to Galois theory, eigenvalues cannot be computed exactly for matrices with $n \geq 5$. **But we still want to compute them!** In this course we will (among other things) explain how this is done in practice by the *QR algorithm*, one of the greatest hits of the field.

## 1.4 Computational complexity (operation counts) of matrix algorithms

Since NLA is a field that aspires to develop practical algorithms for solving matrix problems, it is important to be aware of the computational cost (often referred to as complexity) of the algorithms. We will discuss these as the algorithms are developed, but for now let's examine the costs for basic matrix-matrix multiplication. The cost is measured in terms of flops (floating-point operations), which counts the number of additions, subtractions, multiplications, and divisions (all treated equally) performed.

In NLA the constant in front of the leading term in the cost is (clearly) important. It is customary (for good reason) to only track the leading term of the cost. For example, $n^3 + 10n^2$ is abbreviated to $n^3$.

- Multiplying two $n \times n$ matrices $AB$ costs $2n^3$ flops. More generally, if $A$ is $m \times n$ and $B$ is $n \times k$, then computing $AB$ costs $2mnk$ flops.

- Multiplying a vector to an $m \times n$ matrix $A$ costs $2mn$ flops.

# Norms

We will need a tool (or metric) to measure how big a vector or matrix is. Norms give us a means to achieve this. Surely you have already seen some norms (e.g. the vector Euclidean norm). We will discuss a number of norms for vectors and matrices that we will use in the upcoming lectures.

## 1.5 Vector norms

For vectors $x = [x_1, \ldots, x_n]^T \in \mathbb{C}^n$

- $p$-norm $\|x\|_p = (|x_1|^p + |x_2|^p + \cdots + |x_n|^p)^{1/p}$ $\quad (1 \leq p \leq \infty)$

  - Euclidean norm=2-norm $\|x\|_2 = \sqrt{|x_1|^2 + |x_2|^2 + \cdots + |x_n|^2}$
  - 1-norm $\|x\|_1 = |x_1| + |x_2| + \cdots + |x_n|$
  - $\infty$-norm $\|x\|_\infty = \max_i |x_i|$

Of particular importance are the three cases $p = 1, 2, \infty$. In this course, we will see $p = 2$ the most often.

A norm needs to satisfy the following axioms:

- $\|\alpha x\| = |\alpha| \|x\|$ for any $\alpha \in \mathbb{C}$ (homogeneity),

- $\|x\| \geq 0$ and $\|x\| = 0 \Leftrightarrow x = 0$ (nonnegativity),

- $\|x + y\| \leq \|x\| + \|y\|$ (triangle inequality).

The vector $p$-norm satisfies all these, for any $p$.

Here are some useful inequalities for vector norms. A proof is left for your exercise and is highly recommended. (Try to think when each equality is satisfied.) For $x \in \mathbb{C}^n$,

- $\frac{1}{\sqrt{n}} \|x\|_2 \leq \|x\|_\infty \leq \|x\|_2$

- $\frac{1}{\sqrt{n}} \|x\|_1 \leq \|x\|_2 \leq \|x\|_1$

- $\frac{1}{n} \|x\|_1 \leq \|x\|_\infty \leq \|x\|_1$

Note that with the 2-norm, $\|Ux\|_2 = \|x\|_2$ for any unitary $U$ and any $x \in \mathbb{C}^n$. Norms with this property are called **unitarily invariant**.

The 2-norm is also induced by the inner product $\|x\|_2 = \sqrt{x^T x}$. An important property of inner products is the Cauchy-Schwarz inequality $|x^T y| \leq \|x\|_2 \|y\|_2$ (which can be directly proved but is perhaps best to prove in general setting)[4]. When we just say $\|x\|$ for a vector we mean the 2-norm.

## 1.6 Matrix norms

We now turn to norms of matrices. As you will see, many (but not the Frobenius and trace norms) are defined via the vector norms (these are called *induced* norms).

- $p$-norm $\|A\|_p = \max_x \frac{\|Ax\|_p}{\|x\|_p}$

  - 2-norm=spectral norm(=Euclidean norm) $\|A\|_2 = \sigma_{\max}(A)$ (largest singular value; see Section 2)

  - 1-norm $\|A\|_1 = \max_i \sum_{j=1}^{m} |A_{ji}|$

  - $\infty$-norm $\|A\|_\infty = \max_i \sum_{j=1}^{n} |A_{ij}|$

- Frobenius norm $\|A\|_F = \sqrt{\sum_i \sum_j |A_{ij}|^2}$
  (2-norm of vectorisation)

---

[4] Just in case, here's a proof: for any scalar $c$, $\|x - cy\|^2 = \|x\|^2 - 2cx^T y + c^2 \|y\|^2$. This is minimised wrt $c$ at $c = \frac{x^T y}{\|y\|^2}$ with minimiser $\|x\|^2 - \frac{(x^T y)^2}{\|y\|^2}$. Since this must be $\geq 0$, the CS inequality follows.

- trace norm=nuclear norm $\|A\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i(A)$. (this is the maximum trace of $Q^T A$ where $Q$ is orthonormal, hence the name)

Colored in red are **unitarily invariant** norms $\|A\|_* = \|UAV\|_*, \|A\|_F = \|UAV\|_F, \|A\|_2 = \|UAV\|_2$ for any unitary/orthogonal $U, V$.

Norm axioms hold for each of these. Useful inequalities include: For $A \in \mathbb{C}^{m \times n}$, (exercise; it is instructive to study the cases where each of these equalities holds)

- $\frac{1}{\sqrt{n}}\|A\|_\infty \le \|A\|_2 \le \sqrt{m}\|A\|_\infty$

- $\frac{1}{\sqrt{m}}\|A\|_1 \le \|A\|_2 \le \sqrt{n}\|A\|_1$

- $\|A\|_2 \le \|A\|_F \le \sqrt{\min(m,n)}\|A\|_2$

A useful property of $p$-norms is that they are subordinate, i.e., $\|AB\|_p \le \|A\|_p\|B\|_p$ (problem sheet). Note that not all norms satisfy this, e.g. with the max norm $\|A\|_{\max} = \max_{i,j}|A_{ij}|$, with $A = [1,1]$ and $B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ one has $\|AB\|_{\max} = 2$ but $\|A\|_{\max} = \|B\|_{\max} = 1$.

## 1.7 Subspaces and orthonormal matrices

A key notion that we will keep using throughout is a **subspace** $\mathcal{S}$. In this course we will almost exclusively confine ourselves to subspaces of $\mathbb{R}^n$, even though they generalize to more abstract vector spaces. A subspace is the set of vectors that can be written as a linear combination **basis vectors** $v_1, \ldots, v_d$, which are assumed to be linearly independent (otherwise there is a basis with fewer vectors). That is, $x \in \mathcal{S}$ iff $\sum_{i=1}^d c_i v_i$ (where $c_i$ are scalars, $\mathbb{R}$ or $\mathbb{C}$). The integer $d$ is called the *dimension* of the subspace. We also say the subspace is *spanned* by the vectors $v_1, \ldots, v_d$, or that $v_1, \ldots, v_d$ spans the subpsace.

How does one represent a subspace? An obvious answer is to use the basis vectors $v_1, \ldots, v_d$. This sometimes becomes cumbersome, and a common and convenient way to represent the subspace is to use a (tall-skinny) rectangular matrix $V \in \mathbb{R}^{n \times d} = [v_1, v_2, \ldots, v_d]$, as $\mathcal{S} = \text{span}(V)$ (or sometimes just "subspace $V$") which means the subspace of vectors that can be written as $Vc$, where $c$ is a 'coefficient' vector $c \in \mathbb{R}^d$.

It will be (not necessary but) convenient to represent subspaces using an orthonormal matrix $Q \in \mathbb{R}^{n \times d}$. (once we cover the QR factorisation, you'll see that there is no loss of generality in doing so).

An important fact about subspaces of $\mathbb{R}^n$ is the following:

**Lemma 1.1** *Let $V_1 \in \mathbb{R}^{n \times d_1}$ and $V_2 \in \mathbb{R}^{n \times d_2}$ each have linearly independent column vectors. If $d_1 + d_2 > n$, then there is a nonzero intersection between two subspaces $\mathcal{S}_1 = \text{span}(V_1)$ and $\mathcal{S}_2 = \text{span}(V_2)$, that is, there is a nonzero vector $x \in \mathbb{R}^n$ such that $x = V_1 c_1 = V_2 c_2$ for some vectors $c_1, c_2$.*

This is straightforward but important enough to warrant a proof.

**Proof:** Consider the matrix $M := [V_1, V_2]$, which is of size $n \times (d_1 + d_2)$. Since $d_1 + d_2 > n$ by assumption, this matrix has a right null vector[5] $c \neq 0$ such that $Mc = 0$. Splitting $c = \begin{bmatrix} c_1 \\ -c_2 \end{bmatrix}$ we have the required result. $\qquad\square$

Let us conclude this review with a list of useful results that will be helpful. Proofs (or counterexample) should be straightforward.

- $(AB)^T = B^T A^T$

- If $A, B$ invertible, $(AB)^{-1} = B^{-1} A^{-1}$

- If $A, B$ square and $AB = I$, then $BA = I$

- $\begin{bmatrix} I_m & X \\ 0 & I_n \end{bmatrix}^{-1} = \begin{bmatrix} I_m & -X \\ 0 & I_n \end{bmatrix}$

- Neumann series: if $\|X\| < 1$ in any norm,

$$(I - X)^{-1} = I + X + X^2 + X^3 + \cdots$$

- For a square $n \times n$ matrix $A$, the trace is $\text{Trace}(A) = \sum_{i=1}^n A_{i,i}$ (sum of diagonals). For any $X, Y$ such that $XY$ is square, $\text{Trace}(XY) = \text{Trace}(YX)$ (quite useful). For $B \in \mathbb{R}^{m \times n}$, we have $\|B\|_F^2 = \sum_i \sum_j |B_{ij}|^2 = \text{Trace}(B^T B)$.

- Triangular structure (upper or lower) is invariant under addition, multiplication, and inversion. That is, triangular matrices form a ring (in abstract algebra; don't worry if this is foreign to you).

- Symmetry is invariant under addition and inversion, *but not multiplication*; $AB$ is usually not symmetric even if $A, B$ are.

# 2 SVD: the most important matrix decomposition

We now start the discussion of the most important topic of the course: the singular value decomposition (SVD). The SVD exists for any matrix, square or rectangular, real or complex.

We will prove its existence and discuss its properties and applications in particular in low-rank approximation, which can immediately be used for compressing the matrix, and therefore data.

The SVD has many intimate connections to symmetric eigenvalue problems. Let's start with a review.

---

[5]If this argument isn't convincing to you now, probably the easiest way to see this is via the SVD; so stay tuned and we'll resolve this in footnote 9, once you've seen the SVD!

- **Symmetric eigenvalue decomposition**: Any symmetric matrix $A \in \mathbb{R}^{n \times n}$ has the decomposition

$$A = V \Lambda V^T \tag{1}$$

  where $V$ is orthogonal, $V^T V = I_n = V V^T$, and $\Lambda = \text{diag}(\lambda_1, \ldots, \lambda_n)$ is a diagonal matrix of eigenvalues.

$\lambda_i$ are the eigenvalues, and $V$ is the matrix of eigenvectors (its columns are the eigenvectors).

The decomposition (1) makes two remarkable claims: the eigenvectors can be taken to be orthogonal (which is true more generally of *normal* matrices s.t. $A^* A = A A^*$), and the eigenvalues are real.

It is worth reminding you that eigenvectors are not uniquely determined: (assuming they are normalised s.t. the 2-norm is 1) their signs can always be flipped. (And actually more, when there are eigenvalues that are multiple $\lambda_i = \lambda_j$; in this case the eigenvectors span a subspace whose dimension matches the multiplicity. For example, any vector is an eigenvector of the identity matrix $I$).

Now here is the protagonist of this course: Be sure to spend dozens (if not hundreds) of hours thinking about it!

**Theorem 2.1 (Singular Value Decomposition (SVD))** *Any matrix $A \in \mathbb{R}^{m \times n}$ has the decomposition :*

$$A = U \Sigma V^T. \tag{2}$$

*Here $U^T U = V^T V = I_n$ (assuming $m \geq n$ for definiteness), $\Sigma = diag(\sigma_1, \ldots, \sigma_n)$, $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n \geq 0$.*

$\sigma_i$ (always nonnegative) are called the *singular values* of $A$. The *rank* of $A$ is the number of positive singular values. The columns of $U$ are called the *left singular vectors*, and the columns of $V$ are the *right singular vectors*.

Writing $U = [u_1, \ldots, u_n]$ and $V = [v_1, \ldots, v_n]$, we have an important alternative expression for $A$: $A = \sum_{i=1}^n \sigma_i u_i v_i^T$. We will use this expression repeatedly in what follows. Also note that we always order the singular values in nonincreasing order, so $\sigma_1$ is always the largest singular value.

The SVD tells us that any (tall) matrix can be written as orthonormal-diagonal-orthogonal. Roughly, ortho-normal/gonal matrices can be thought of as rotations or reflection, so the SVD says the action of a matrix can be thought of as a rotation/reflection followed by magnification (or shrinkage), followed by another rotation/reflection.

**Proof:**    For SVD[6] ($m \geq n$ and assume full-rank $\sigma_n > 0$ for simplicity): Take Gram matrix $A^T A$ (symmetric) and its eigendecomposition $A^T A = V \Lambda V^T$ with $V$ orthogonal. $\Lambda$ is nonnegative, and $(AV)^T (AV) =: \Sigma^2$ is diagonal, so $AV\Sigma^{-1} =: U$ is orthonormal. Right-multiply by $\Sigma V^T$ to get $A = U\Sigma V^T$. $\qquad \square$

---

[6]I like to think this is the shortest proof out there.

It is also worth mentioning the "full" SVD: $A = [U, U_\perp] \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T$ where $[U, U_\perp] \in \mathbb{R}^{m \times m}$ is square and orthogonal. Essentially this follows from the (thin) SVD (2) by filling in $U$ in (2) with its orthogonal complement $U_\perp$ (whose construction can be done via the Householder QR factorsation of Section 6).

## 2.1 (Some of the many) applications and consequences of the SVD: rank, column/row space, etc

From the SVD one can immediately read off a number of important properties of the matrix. For example:

- The *rank* $r$ of $A \in \mathbb{R}^{m \times n}$, often denoted rank($A$): this is the number of nonzero (positive) singular values $\sigma_i(A)$. rank($A$) is also equal to the number of linearly independent rows or the number of independent columns, as you probably learnt in your first course on linear algebra (exercise).

  – We can always write $A = \sum_{i=1}^{\text{rank}(A)} \sigma_i u_i v_i^T$.

  Important: An $m \times n$ matrix $A$ that is of rank $r$ can be written as an (outer) product of $m \times r$ and $r \times n$ matrices:

$$A_r = \boxed{U_r} \; \boxed{\Sigma_r} \; \boxed{\quad V_r^T \quad}$$

  To see this, note from the SVD that $A = \sum_{i=1}^{n} \sigma_i u_i v_i^T = \sum_{i=1}^{r} \sigma_i u_i v_i^T$ (since $\sigma_{r+1} = 0$), and so $A = U_r \Sigma_r V_r^T$ where $U_r = [u_1, \dots, u_r]$, $V_r = [v_1, \dots, v_r]$, and $\Sigma_r$ is the leading $r \times r$ submatrix of $\Sigma$.

- Column space (Span(A), linear subspace spanned by vectors $Ax$): span of $U = [u_1, \dots, u_r]$, often denoted Span($U$).

- Row space (Span($A^T$)): row span of $v_1^T, \dots, v_r^T$, Span($V$)$^T$.

- Null space Null(A): span of $v_{r+1}, \dots, v_n$; as $Av = 0$ for these vectors. Null(A) is empty (or just the 0 vector) if $m \geq n$ and $r = n$. (When $m < n$ the full SVD is needed to describe Null(A).)

Aside from these and other applications, the SVD is also a versatile theoretical tool. Very often, a good place to start in proving a fact about matrices is to first consider its SVD; you will see this many times in this course. For example, the SVD can give solutions immediately for linear systems and least-squares problems, though there are more efficient ways to solve these problems.

## 2.2 SVD and symmetric eigenvalue decomposition

As mentioned above, the SVD and the eigenvalue decomposition for symmetric matrices are closely connected. Here are some results that highlight the connections between $A = U\Sigma V^T$ and symmetric eigenvalue decomposition. (We assume $m \geq n$ for definiteness)

- $V$ is an eigvector matrix of $A^T A$. (To verify, see proof of SVD)

- $U$ is an eigvector matrix (for nonzero eigvals) of $AA^T$ (be careful with sign flips; see below)

- $\sigma_i = \sqrt{\lambda_i(A^T A)}$ for $i = 1, \ldots, n$

- If $A$ is symmetric, its singular values $\sigma_i(A)$ are the absolute values of its eigenvalues $\lambda_i(A)$, i.e., $\sigma_i(A) = |\lambda_i(A)|$.

  Exercise: What if $A$ is unitary, skew-symmetric, normal matrices, triangular? (problem sheet)

- Jordan-Wieldant matrix $\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}$: This matrix has eigenvalues $\pm\sigma_i(A)$, and $m-n$ copies of eigenvalues at 0. Its eigenvector matrix is $\begin{bmatrix} U & U & U_0 \\ V & -V & 0 \end{bmatrix}$, where $A^T U_0 = 0$ ($U_0$ is empty when $m = n$). This matrix, along with the Gram matrix $A^T A$, is a very useful tool when one tries to extend a result on symmetric eigenvalues to an analogue in terms of the SVD (or vice versa).

## 2.3 Uniqueness etc

We have established the existence of the SVD $A = U\Sigma V^T$. A natural question is: is it unique? In other words, are the factors $U, \Sigma, V$ uniquely determined by $A$?

It is straightforward to see that the singular vectors are not uniquely determined. Most obviously, the singular vectors can be flipped in signs, just like eigenvectors. However, note that the signs of $u_i, v_i$ are not entirely arbitrary: if we replace $u_i$ by $-u_i$, the same needs to be done for $v_i$ in order to satisfy $A = \sum_{i=1}^n \sigma_i u_i v_i^T$. Essentially, once the sign (or rotation) of $u_i$ (or $v_i$) is fixed, $v_i$ (or $u_i$) is determined uniquely.

More generally, in the presence of multiple singular values (i.e., $\sigma_i = \sigma_{i+1}$ for some $i$), there is a higher degree of freedom in the SVD. Again this is very much analogous to eigenvalues (recall the discussion on eigenvectors of the identity). Here think of what the SVD is for an orthogonal matrix: there is an enormous amount of degrees of freedom in the choice of $U$ and $V$. The singular values $\sigma_i$, on the other hand, are always unique, as they are the eigenvalues of the Gram matrix.

## 3 Low-rank approximation via truncated SVD

While the SVD has a huge number of applications, undoubtedly the biggest reason that makes it so important in computational mathematics is its optimality for low-rank approximation.

To discuss this topic we need some preparation. We will make heavy use of the spectral norm of matrices. We start with an important characterisation of $\|A\|_2$ in terms of the singular value(s), which we previously stated but did not prove.

**Proposition 3.1**

$$\|A\|_2 = \max_x \frac{\|Ax\|_2}{\|x\|_2} = \max_{\|x\|_2=1} \|Ax\|_2 = \sigma_1(A).$$

**Proof:** Use the SVD: For any $x$ with unit norm $\|x\|_2 = 1$,

$$
\begin{aligned}
\|Ax\|_2 &= \|U\Sigma V^T x\|_2 \\
&= \|\Sigma V^T x\|_2 \quad \text{by unitary invariance} \\
&= \|\Sigma y\|_2 \quad \text{with } \|y\|_2 = 1 \\
&= \sqrt{\sum_{i=1}^{n} \sigma_i^2 y_i^2} \\
&\leq \sqrt{\sum_{i=1}^{n} \sigma_1^2 y_i^2} = \sigma_1 \|y\|_2^2 = \sigma_1.
\end{aligned}
$$

Finally, note that taking $x = v_1$ (the leading right singular vector), we have $\|Av_1\|_2 = \sigma_1$. $\square$

Similarly, the Frobenius norm can be expressed as $\|A\|_F = \sqrt{\sum_i \sum_j |A_{ij}|^2} = \sqrt{\sum_{i=1}^{n} (\sigma_i(A))^2}$, and the trace norm is (by definition) $\|A\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i(A)$. (exercise) In general, norms that are *unitarily invariant* can be characterized by the singular values [20].

Now to the main problem: Given $A \in \mathbb{R}^{m \times n}$, consider the problem of finding a *rank-r matrix* (remember; these are matrices with $r$ nonzero singular values) $A_r \in \mathbb{R}^{m \times n}$ that best approximates $A$. That is, find the minimiser $A_r$ for

$$\text{argmin}_{\text{rank}(A_r) \leq r} \|A - A_r\|_2. \tag{3}$$

It is definitely worth visualizing the situation:



We immediately see that a low rank approximation (when possible) is beneficial in terms of the storage cost when $r \ll m, n$. Instead of storing $mn$ entries for $A$, we can store entries for $U_r, \Sigma_r, V_r$ ($(m+n+1)r$ entries) to keep the low-rank factorisation without losing information.

Low-rank approximation can also bring computational benefits. For example, in order to compute $Ax$ for a vector $x$, by noting that $A_r x = U_r(\Sigma_r(V_r^T x))$, one needs only $O((m+n)r)$

operations[7] instead of $O(mn)$. The utility and prevalence of low-rank matrices in data science is remarkable.

Here is the solution for (3): Truncated SVD, defined via $A_r = \sum_{i=1}^{r} \sigma_i u_i v_i^T (= U_r \Sigma_r V_r^T)$. This is the matrix obtained by truncating (removing) the trailing terms in the expression $A = \sum_{i=1}^{n} \sigma_i u_i v_i^T$. Pictorially,

$$A = \underbrace{\begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & \cdots & * & * \end{bmatrix}}_{\sigma_1 u_1 v_1} + \underbrace{\begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & \cdots & * & * \end{bmatrix}}_{\sigma_2 u_2 v_2} + \cdots + \underbrace{\begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & \cdots & * & * \end{bmatrix}}_{\sigma_n u_n v_n},$$

$$A_r = \underbrace{\begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & \cdots & * & * \end{bmatrix}}_{\sigma_1 u_1 v_1} + \cdots + \underbrace{\begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & \cdots & * & * \end{bmatrix}}_{\sigma_r u_r v_r}.$$

In particular, we have

**Theorem 3.1** *For any $A \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1(A) \geq \sigma_2(A) \geq \cdots \geq \sigma_n(A) \geq 0$, and any nonnegative integer[8] $r < \min(m, n)$,*

$$\|A - A_r\|_2 = \sigma_{r+1}(A) = \min_{rank(B) \leq r} \|A - B\|_2. \tag{4}$$

Before proving this result, let us make some observations.

- Good approximation $A \approx A_r$ is obtained iff $\sigma_{r+1} \ll \sigma_1$.

- Optimality holds for any unitarily invariant norm: that is, the norms in (3) can be replaced by e.g. the Frobenius norm. This is surprising, as the low-rank approximation problem $\min_{\text{rank}(B) \leq r} \|A - B\|$ does depend on the choice of the norm (and for many problems, including the least-squares problem in Section 6.5, the norm choice has a significant effect on the solution). The proof for this fact is nonexaminable, but if curious see [20] for a complete proof.

- A prominent application of low-rank approximation is PCA (principal component analysis) in statistics and data science.

- Many matrices have explicit or hidden low-rank structure (nonexaminable, but see e.g. [34]).

---

[7]I presume you've seen the big-Oh notation before—$f(n) = O(n^k)$ means there exists a constant $C$ s.t. $f(n)/n^k \leq C$ for all sufficiently large $n$. In NLA it is a convenient way to roughly measure the operation count.

[8]If $r \geq \min(m, n)$ then we can simply take $A_r = A$.

**Proof:** of Theorem 3.1:

1. Since $\text{rank}(B) \leq r$, we can write $B = B_1 B_2^T$ where $B_1, B_2$ have $r$ columns.

2. It follows that $B_2^T$ (and hence $B$) has a null space of dimension at least $n - r$. That is, there exists an orthonormal matrix $W \in \mathbb{C}^{n \times (n-r)}$ s.t. $BW = 0$. Then $\|A - B\|_2 \geq \|(A - B)W\|_2 = \|AW\|_2 = \|U\Sigma(V^TW)\|_2$. (Why does the first inequality hold?)

3. Now since $W$ is $(n - r)$-dimensional, by Lemma 1.1 there is an intersection between $W$ and $[v_1, \ldots, v_{r+1}]$, the $(r + 1)$-dimensional subspace spanned by the leading $r + 1$ right singular vectors.

   We will use this type of argument again, so to be more precise: the matrix $[W, v_1, \ldots, v_{r+1}]$ is "fat" rectangular, so must have a null vector. That is, $[W, v_1, \ldots, v_{r+1}]\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$ has a nonzero solution $x_1, x_2$; then $Wx_1$ is such a vector[9]. We scale it so that it has unit norm $\|Wx_1\|_2 = \|[v_1, \ldots, v_{r+1}]x_2\|_2 = 1$, that is, $\|x_1\|_2 = \|x_2\|_2 = 1$.

4. Note that $U\Sigma V^T W x_1 = U\Sigma V^T[v_1, \ldots, v_{r+1}]x_2$ and $V^T[v_1, \ldots, v_{r+1}] = \begin{bmatrix} I_{r+1} \\ 0 \end{bmatrix}$, so $\|U\Sigma V^T W x_1\|_2 = \|U\begin{bmatrix} \Sigma_{r+1} \\ 0 \end{bmatrix} x_2\|_2$, where $\Sigma_{r+1}$ is the leading (top-left) $(r + 1) \times (r + 1)$ part of $\Sigma$. As $U$ is orthogonal this is equal to $\|\Sigma_{r+1}y_1\|_2$, and $\|\Sigma_{r+1}y_1\|_2 \geq \sigma_{r+1}$ can be verified by direct calculation.

   Finally, for the reverse direction, take $B = A_r$. □

## 3.1 Low-rank approximation: image compression

A classical and visually pleasing illustration of low-rank approximation by the truncated SVD is image compression. (Admittedly this is slightly outdated—it is not the state-of-the-art way of compressing images.)

The idea is that greyscale images can be represented by a matrix where each entry indicates the intensity of a pixel. The matrix can then be compressed by finding a low-rank approximation[10], resulting in image compression. (Of course there are generalisations for color images.)

Below in Figure 1 we take the Oxford logo, represent it as a matrix and find its rank-$r$ approximation, for varying $r$. (The matrix being a mere $500 \times 500$, its SVD is easy to compute in a fraction of a second; see Section 10.2 for how this is done.) We then reconstruct the image from the low-rank approximations to visualise them.

---

[9]Let us now resolve the cliffhanger in footnote 5. The claim is that any "fat" $m \times n$ ($m < n$) matrix $M$ has a right null vector $y \neq 0$ such that $My = 0$. To prove this, use the full SVD $M = U\Sigma V^T$ to see that $Mv_n = 0$.

[10]It is somewhat surprising that images are approximable by low-rank matrices. See https://www.youtube.com/watch?v=9BYsNpTCZGg for a nice explanation.

original      rank 1      rank 5

rank 10      rank 20      rank 50

Figure 1: Image compression by low-rank approximation via the truncated SVD.

We see that as the rank is increased the image becomes finer and finer. At rank 50 it is fair to say the image looks almost identical to the original. The original matrix is $500 \times 500$, so we still achieve a significant amount of data compression in the matrix with $r = 50$.

# 4 Courant-Fischer minmax theorem

Continuing on SVD-related topics, we now discuss a very important and useful result with far-reaching ramifications: the Courant-Fischer (C-F) minimax characterisation.

**Theorem 4.1** *The $i$th largest[11] eigenvalue $\lambda_i$ of a symmetric matrix $A \in \mathbb{R}^{n \times n}$ is (below $x \neq 0$)*

$$\lambda_i(A) = \max_{\dim \mathcal{S}=i} \min_{x \in \mathcal{S}} \frac{x^T A x}{x^T x} \quad \left( = \min_{\dim \mathcal{S}=n-i+1} \max_{x \in \mathcal{S}} \frac{x^T A x}{x^T x} \right) \tag{5}$$

*Analogously, for any rectangular $A \in \mathbb{C}^{m \times n}(m \geq n)$, we have*

$$\sigma_i(A) = \max_{\dim \mathcal{S}=i} \min_{x \in \mathcal{S}} \frac{\|Ax\|_2}{\|x\|_2} \quad \left( = \min_{\dim \mathcal{S}=n-i+1} \max_{x \in \mathcal{S}} \frac{\|Ax\|_2}{\|x\|_2} \right). \tag{6}$$

It would take some time to get a hang of what the statements mean. One helpful way to look at it is perhaps to note that inside the maximum in (6) the expression is $\min_{x \in \mathcal{S}, \|x\|_2 = 1} \|Ax\|_2 = \min_{Q^T Q = I_i, \|y\|_2 = 1} \|AQy\|_2 = \sigma_{\min}(AQ) = \sigma_i(AQ)$, where $\text{span}(Q) = \mathcal{S}$. The C-F theorem says $\sigma_i(A)$ is equal to the maximum possible value of this over all subspaces $\mathcal{S}$ of dimension $i$.

**Proof:** We will prove (6). A proof for (5) is analogous and a recommended exercise.

1. Fix $S$ and let $V_i = [v_i, \ldots, v_n]$. We have $\dim(\mathcal{S}) + \dim(\text{span}(V_i)) = i + (n-i+1) = n+1$, so $\exists$ intersection $w \in S \cap V_i$, $\|w\|_2 = 1$.

2. For this $w$, we have $\|Aw\|_2 = \|\text{diag}(\sigma_i, \ldots, \sigma_n)(V_i^T w)\|_2 \leq \sigma_i$; thus $\sigma_i \geq \min_{x \in \mathcal{S}} \frac{\|Ax\|_2}{\|x\|_2}$.

3. For the reverse inequality, take $S = [v_1, \ldots, v_i]$, for which $w = v_i$.

$\square$

## 4.1 Weyl's inequality

As an example of the many significant ramifications of the C-F theorem, we present *Weyl's theorem*[12] (or Weyl's inequality), an important perturbation result for singular values and eigenvalues of symmetric matrices.

**Theorem 4.2** *Weyl's inequality*

---

[11] exact analogues hold for the $i$th *smallest* eigenvalue and singular values.

[12] Hermann Weyl was one of the prominent mathematicians of the 20th centry.

- *For the singular values of any matrix $A$,*

  - $\sigma_i(A + E) \in \sigma_i(A) + [-\|E\|_2, \|E\|_2]$ *for all $i$.*
  - *Special case:* $\|A\|_2 - \|E\|_2 \le \|A + E\|_2 \le \|A\|_2 + \|E\|_2$

- *For eigenvalues of a symmetric matrix $A$, $\lambda_i(A + E) \in \lambda_i(A) + [-\|E\|_2, \|E\|_2]$ for all $i$.*

(Proof: exercise; almost a direct consequence of C-F.)

The upshot is that singular values and eigenvalues of symmetric matrices are insensitive to perturbation; a property known as being *well conditioned.*

This is important because this means a backward stable algorithm (see Section 7) computes these quantities with essentially full precision.

### 4.1.1 Eigenvalues of nonsymmetric matrices are sensitive to perturbation

It is worth remarking that eigenvalues of nonsymmetric matrices can be far from well conditioned! Consider for example the Jordan block $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. By perturbing this to $\begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix}$ one gets eigenvalues that are perturbed by $\sqrt{\epsilon}$, a magnification factor of $1/\sqrt{\epsilon} \gg 1$. More generally, consider the eigenvalues of a Jordan block and its perturbation

$$
J = \begin{bmatrix} 1 & 1 & & \\ & 1 & \ddots & \\ & & \ddots & 1 \\ & & & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad J + E = \begin{bmatrix} 1 & 1 & & \\ & 1 & \ddots & \\ & & \ddots & 1 \\ \epsilon & & & 1 \end{bmatrix}
$$

$\lambda(J) = 1$ ($n$ copies), but we have $|\lambda(J + E) - 1| \approx \epsilon^{1/n}$. For example when $n = 100$, an $10^{-100}$ perturbation in $J$ would result in a $0.1$ perturbation in all the eigenvalues!

(nonexaminable) This is pretty much the worst-case situation. In the generic case where the matrix is diagonalizable, $A = X \Lambda X^{-1}$, with an $\epsilon$-perturbation the eigenvalues get perturbed by $O(c\epsilon)$, where the constant $c$ depends on the so-called *condition number* $\kappa_2(X)$ of the eigenvector matrix $X$ (see Section 7.2, and [8]).

## 4.2 More applications of C-F

(Somewhat optional) Let's explore more applications of C-F. A lot more can be proved; see [19] for many more results and examples along these lines.

**Example 4.1**  
- $\sigma_i\left(\begin{bmatrix} A_1 \\ A_2 \end{bmatrix}\right) \ge \max(\sigma_i(A_1), \sigma_i(A_2))$

  *Proof (sketch): $LHS = \max_{\dim \mathcal{S} = i} \min_{x \in \mathcal{S}, \|x\|_2 = 1} \left\| \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} x \right\|_2$, and for any $x$, $\left\| \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} x \right\|_2 \ge \max(\|A_1 x\|_2, \|A_2 x\|_2)$.*

- $\sigma_i([A_1 \quad A_2]) \geq \max(\sigma_i(A_1), \sigma_i(A_2))$

  *Proof:* $LHS = \max_{\dim \mathcal{S}=i} \min_{\left[\begin{smallmatrix} x_1 \\ x_2 \end{smallmatrix}\right] \in \mathcal{S}, \left\| \left[\begin{smallmatrix} x_1 \\ x_2 \end{smallmatrix}\right] \right\|_2 = 1} \left\| [A_1 \quad A_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|_2$, *while*

  $\sigma_i(A_1) = \max_{\dim \mathcal{S}=i, range(\mathcal{S}) \in range(\left[\begin{smallmatrix} I_n \\ 0 \end{smallmatrix}\right])} \min_{\left[\begin{smallmatrix} x_1 \\ x_2 \end{smallmatrix}\right] \in \mathcal{S}, \left\| \left[\begin{smallmatrix} x_1 \\ x_2 \end{smallmatrix}\right] \right\|_2 = 1} \left\| [A_1 \quad A_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|_2$. *Since*

  *the latter imposes restrictions on $\mathcal{S}$ to take the maximum over, the former is at least as big.*

## 4.3   (Taking stock) Matrix decompositions you should know

Let us now take stock to review the matrix decompositions that we have covered, along with those that we will discuss next.

- SVD $A = U\Sigma V^T$

- Eigenvalue decomposition $A = X\Lambda X^{-1}$

  - Normal: $X$ unitary $X^*X = I$
  - Symmetric: $X$ unitary and $\Lambda$ real

- Jordan decomposition: $A = XJX^{-1}$, $J = \mathrm{diag}\left( \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix} \right)$

- Schur decomposition $A = QTQ^*$: $T$ upper triangular

- QR: $Q$ orthonormal, $U$ upper triangular

- LU: $L$ lower triangular, $U$ upper triangular

Red: Orthogonal decompositions, stable computation available

The Jordan decomposition is mathematically the ultimate form that any matrix can be reduced to by a similarity transformation, but numerically it is not very useful—one of the problems is that Jordan decompositions are very difficult to compute, as an arbitrarily small perturbation can change the eigenvalues and block sizes by a large amount (recall the discussion in Section 4.1.1).

# 5   Linear systems $Ax = b$

We now (finally) start our discussion on direct algorithms in NLA. The fundamental idea is to *factorisa* a matrix into a product of two (or more) simpler matrices. This foundational idea has been named one of the top 10 algorithms of the 20th centry [9].

The sophistication of the state-of-the-art implementation of direct methods is simply astonishing. For instance, a $100 \times 100$ dense linear system or eigenvalue problem can be solved in less than a milisecond on a standard laptop. Imagine solving it by hand!

We start with solving linear systems, unquestionably the most important problem in NLA (for applications). In some sense we needn't spend too much time here, as you must have seen much of the material (e.g. Gaussian elimination) before. However, the description of the LU factorisation given below will likely be different from the one that you have seen before. We have chosen a nonstandard description as it reveals its connection to low-rank approximation.

Let $A \in \mathbb{R}^{n \times n}$. Suppose we can decompose (or factorise) $A$ into (here and below, $*$ denotes entries that are possibly nonzero).

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} = \begin{bmatrix} * & & & & \\ * & * & & & \\ * & * & * & & \\ * & * & * & * & \\ * & * & * & * & * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \\ & & & & * \end{bmatrix} = LU$$

Here $L$ is lower triangular, and $U$ is upper triangular. How can we find $L, U$?

To get started, consider rewriting $A$ as

$$A = \begin{bmatrix} * \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \end{bmatrix} + \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} * \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \end{bmatrix}}_{L_1 U_1} + \underbrace{\begin{bmatrix} 0 \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} 0 & * & * & * & * \end{bmatrix}}_{L_2 U_2} + \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix} = \cdots$$

Namely, the first step finds $L_1, U_1$ such that:

$$A = \underbrace{\begin{bmatrix} * \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \end{bmatrix}}_{L_1 U_1} + \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}.$$

Specifying the elements, the algorithm can be described as (taking $a = A_{11}$)

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} & A_{15} \\ A_{21} \\ A_{31} \\ A_{41} \\ A_{51} \end{bmatrix} = \begin{bmatrix} L_{11} \\ L_{21} \\ L_{31} \\ L_{41} \\ L_{51} \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} & U_{13} & U_{14} & U_{15} \end{bmatrix} + \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} 1 \\ A_{21}/a \\ A_{31}/a \\ A_{41}/a \\ A_{51}/a \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} & A_{15} \end{bmatrix}}_{=L_1 U_1} + \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}$$

Here we've assumed $a \neq 0$; we'll discuss the case $a = 0$ later. Repeating the process gives

$$A = \begin{bmatrix} * \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \end{bmatrix} + \begin{bmatrix} 0 \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} 0 & * & * & * & * \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} 0 & 0 & * & * & * \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ * \\ * \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & * & * \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ * \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 & * \end{bmatrix}$$

$$= \quad L_1 U_1 \quad + \quad L_2 U_2 \quad + \quad L_3 U_3 \quad + \quad L_4 U_4 \quad + \quad L_5 U_5$$

$$= [L_1, L_2, \ldots, L_5] \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_5 \end{bmatrix} = \begin{bmatrix} * & & & & \\ * & * & & & \\ * & * & * & & \\ * & * & * & * & \\ * & * & * & * & * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \\ & & & & * \end{bmatrix},$$

an LU factorisation as required.

Note the above expression for $A$; clearly the $L_i U_i$ factors are rank-1 matrices; the LU factorisation can be thought of as writing $A$ as a sum of (structured) rank-1 matrices.

## 5.1 Solving $Ax = b$ via LU

Having found an LU factorisation $A = LU$, one can efficiently solve an $n \times n$ linear system $Ax = b$: *First solve $Ly = b$, then $Ux = y$. Then $b = Ly = LUx = Ax$.*

- These are *triangular* linear systems, which are easy to solve and can be done in $O(n^2)$ flops.

- Triangular solve is always backward stable: e.g. $(L + \Delta L)\hat{y} = b$ (see Higham's book)

The computational cost is

- For LU: $\frac{2}{3}n^3$ flops (floating-point operations).

- Triangular solve is $O(n^2)$.

Note that once we have an LU factorisation we can solve another linear system with respect to the same matrix with only $O(n^2)$ additional operations.

## 5.2 Pivoting

Above we've assumed the diagonal element (pivot) $a \neq 0$—when $a = 0$ we are in trouble! In fact not every matrix has an LU factorisation. For example, there is no LU for $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$. We need a remedy. In practice, a remedy is needed whenever $a$ is small.

The idea is to *permute* the rows, so that the largest element of the (first active) column is brought to the pivot. This process is called *pivoting* (sometimes *partial pivoting*, to emphasize the difference from *complete pivoting* wherein both rows and columns are permuted[13]). This

---

[13]While we won't discuss complete pivoting further, you might be interested to know that when LU with complete pivoting is applied to a low-rank matrix (with rapidly decaying singular values), one tends to find a good low-rank approximation, almost as good as truncated SVD.

results in $PA = LU$, where $P$ is a *permutation matrix*: orthogonal matrices with only 1 and 0s (every row/column has exactly one 1); applying $P$ would reorder the rows (with $PA$) or columns ($AP$)).

Thus solving $Ax_i = b_i$ for $i = 1, \ldots, k$ requires $\frac{2}{3}n^3 + O(kn^2)$ operations instead of $O(kn^3)$.

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} & A_{15} \\ A_{21} \\ A_{31} \\ A_{41} \\ A_{51} \end{bmatrix} = \begin{bmatrix} 1 \\ A_{21}/a \\ A_{31}/a \\ A_{41}/a \\ A_{51}/a \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} & A_{15} \\ \\ \\ \\ \end{bmatrix} + \begin{bmatrix} & & & \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix}$$

When $a = 0$, remedy: pivot, permute rows such that the largest entry of first (active) column is at the top. $\Rightarrow PA = LU$, $P$: permutation matrix

- for $Ax = b$, solve $PAx = Pb \Leftrightarrow LUx = Pb$

- cost still $\frac{2}{3}n^3 + O(n^2)$

In fact, one can show that any nonsingular matrix $A$ has a pivoted LU factorisation. (proof: exercise). This means that any linear system that is computationally feasible can be solved by pivoted LU factorisation.

- Even with pivoting, unstable examples exist (Section 7), but almost always stable in practice and used everywhere.

- Stability here means $\hat{L}\hat{U} = PA + \Delta A$ with small $\|\Delta A\|$; see Section 7.

## 5.3  Cholesky factorisation for $A \succ 0$

If $A \succ 0$ (symmetric positive definite[14] (S)PD$\Leftrightarrow \lambda_i(A) > 0$ for all $i$), two simplifications happen in LU:

- We can take $U_i = L_i^T =: R_i$ by symmetry

- No pivot needed as long as $A$ is PD

$$A = \underbrace{\begin{bmatrix} * \\ * \\ * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * & * \end{bmatrix}}_{R_1 R_1^T} + \underbrace{\begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}}_{\text{also PD}} \tag{7}$$

Notes:

- $\frac{1}{3}n^3$ flops, half as many as LU

---

[14]Positive definite matrices are so important a class of matrices; also full of interesting mathematical properties, enough for a nice book to be written about it [4].

- diag($R$) no longer 1's (clearly)

- $A$ can be written as $A = R^T R$ for some $R \in \mathbb{R}^{n \times n}$ iff $A \succeq 0$ ($\lambda_i(A) \geq 0$)

- Indefinite case: when $A = A^*$ but $A$ not PSD (i.e., negative eigenvalues are present), $\exists\, A = LDL^*$ where $D$ diagonal (when $A \in \mathbb{R}^{n \times n}$, $D$ can have $2 \times 2$ diagonal blocks), $L$ has 1's on diagonal. This is often called the "LDLT factorisation".

- It's not easy at this point to see why the second (red) term above is also PD; one way to see this is via the uniqueness of the Cholesky factorisation; see next section.

Therefore, roughly speaking, symmetric linear systems can be solved with half the effort of a general non-symmetric system.

# 6 QR factorisation and least-squares problems

We've seen that the LU factorisation is a key step towards solving $Ax = b$. For an overdetemined problem (least-squares problems, the subject of Section 6.5), we will need the *QR factorisation*. For any $A \in \mathbb{R}^{m \times n}$, there exists a factorisation

$$ A = Q\ R $$

where $Q \in \mathbb{R}^{m \times n}$ is orthonormal $Q^T Q = I_n$, and $R \in \mathbb{R}^{n \times n}$ is upper triangular.

- Many algorithms available: Gram-Schmidt, Householder QR, CholeskyQR, ...

- various applications: least-squares, orthogonalisation, computing SVD, manifold retraction...

- With Householder, pivoting $A = QRP$ not needed for numerical stability.

  - but pivoting gives rank-revealing QR (nonexaminable).

## 6.1 QR via Gram-Schmidt

No doubt you have seen the Gram-Schmidt (G-S) process. What you might not know is that when applied to the columns of a matrix $A$, it gives you a QR factorisation $A = QR$.

Gram-Schmidt: Given $A = [a_1, a_2, \ldots, a_n] \in \mathbb{R}^{m \times n}$ (assume full-rank $\operatorname{rank}(A) = n$), find orthonormal $[q_1, \ldots, q_n]$ s.t. $\operatorname{span}(q_1, \ldots, q_n) = \operatorname{span}(a_1, \ldots, a_n)$

More precisely, the algorithm performs the following: $q_1 = \frac{a_1}{\|a_1\|}$, then $\tilde{q}_2 = a_2 - q_1 q_1^T a_2$, $q_2 = \frac{\tilde{q}_2}{\|\tilde{q}_2\|}$, (orthogonalise and normalise)

repeat for $j = 3, \dots, n$: $\tilde{q}_j = a_j - \sum_{i=1}^{j-1} q_i q_i^T a_j$, $q_j = \frac{\tilde{q}_j}{\|\tilde{q}_j\|}$.

**This gives a QR factorisation!** To see this, let $r_{ij} = q_i^T a_j$ $(i \neq j)$ and $r_{jj} = \|a_j - \sum_{i=1}^{j-1} r_{ij} q_i\|$,

$$q_1 = \frac{a_1}{r_{11}}$$
$$q_2 = \frac{a_2 - r_{12} q_1}{r_{22}}$$
$$q_j = \frac{a_j - \sum_{i=1}^{j-1} r_{ij} q_i}{r_{jj}}$$

which can be written equivalently as

$$a_1 = r_{11} q_1$$
$$a_2 = r_{12} q_1 + r_{22} q_2$$
$$a_j = r_{1j} q_1 + r_{2j} q_2 + \cdots + r_{jj} q_j.$$

This in turn is $\boxed{A} = \boxed{Q}\,\boxed{R}$, where $Q^T Q = I_n$, and $R$ upper triangular.

- But this isn't the recommended way to compute the QR factorisation, as it's numerically unstable; see Section 7.7.1 and [17, Ch. 19,20].

## 6.2   Towards a stable QR factorisation: Householder reflectors

There is a beautiful alternative algorithm for computing the QR factorisation: *Householder QR factorisation*. In order to describe it, let us first introduce Householder reflectors. These are the class of matrices $H$ that are symmetric, orthogonal and can be written as a rank one update of the identity

$$H = I - 2vv^T, \qquad \|v\| = 1$$

- $H$ is orthogonal and symmetric: $H^T H = H^2 = I$. Its eigenvalues are 1 ($n-1$ copies) and $-1$ (1 copy).

- For any given $u, w \in \mathbb{R}^n$ s.t. $\|u\| = \|w\|$ and $u \neq v$, $H = I - 2vv^T$ with $v = \frac{w-u}{\|w-u\|}$ gives $Hu = w$ ($\Leftrightarrow u = Hw$, thus 'reflector')

It follows that by choosing the vector $v$ appropriately, one can perform a variety of operations to a given vector $x$. A primary example is the particular Householder reflector

$$H = I - 2vv^T, \qquad v = \frac{x - \|x\|e}{\|x - \|x\|e\|}, \qquad e = [1, 0, \ldots, 0]^T,$$

which satisfies $Hx = [\|x\|, 0, \ldots, 0]^T$. That is, an arbitrary vector $x$ is mapped by $H$ to a multiple of $e = [1, 0, \ldots, 0]^T$.

(I hope the picture above is helpful—the reflector reflects vectors about the hyperplane described by $v^T x = 0$.)

In summary, we have the useful result

**Lemma 6.1** *For any $x \in \mathbb{R}^n$ not in the form $[\pm\|x\|, 0, \ldots, 0]^T$, there exists a Householder reflector $H = I - 2vv^T$ where $\|v\| = 1$ such that $Hx = [\|x\|, 0, \ldots, 0]^T$.*

## 6.3 Householder QR factorisation

Now we describe how to use the Householder reflectors in order to compute a QR factorisation of a given matrix $A \in \mathbb{R}^{m \times n}$.

The first step to obtain QR is to find $H_1$ s.t. $H_1 a_1 = \begin{bmatrix} \|a_1\| \\ 0 \\ \vdots \\ 0 \end{bmatrix}$,

and repeat to get $H_n \cdots H_2 H_1 A = R$ upper triangular, then $A = (H_1 \cdots H_{n-1} H_n)R = QR$

Here is a pictorial illustration: start with

$$A = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}.$$

We apply a sequence of Householder reflectors

29

$$H_1 A = (I - 2v_1 v_1^T)A = \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \end{bmatrix}, \qquad H_2 H_1 A = (I - 2v_2 v_2^T)H_1 A = \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & * & * \\ & & * & * \end{bmatrix},$$

$$H_3 H_2 H_1 A = \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \end{bmatrix}, \qquad H_n \cdots H_3 H_2 H_1 A = \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \end{bmatrix}.$$

Note the zero pattern $v_k = [\underbrace{0, 0, \ldots, 0}_{k-1 \ 0\text{'s}}, *, *, \ldots, *]^T$. We have

$$H_n \cdots H_2 H_1 A = \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & \end{bmatrix} = \begin{bmatrix} R \\ 0 \end{bmatrix}$$

To obtain a QR factorisation of $A$ we simply invert the Householder reflectors, noting that they are both orthogonal and symmetric; therefore the inverse is itself. This yields

$$\Leftrightarrow A = (H_1^T \cdots H_{n-1}^T H_n^T) \begin{bmatrix} R \\ 0 \end{bmatrix} =: Q_F \begin{bmatrix} R \\ 0 \end{bmatrix};$$ which is a **full** QR, wherein $Q_F$ is square orthogonal.

Moreover, writing $Q_F = [Q \ Q_\perp]$ where $Q \in \mathbb{R}^{m \times n}$ is orthonormal, we also have $A = QR$ ('**thin**' QR or just QR); this is more economical especially when $A$ is tall-skinny. In a majority of cases in computational mathematics, this is the object that we wish to compute).

We note some properties of Householder QR.

- Cost $\frac{4}{3}n^3$ flops with Householder-QR (twice that of LU).

- Unconditionally backward stable: the computed version satisfies $\hat{Q}\hat{R} = A + \Delta A$, $\|\hat{Q}^T \hat{Q} - I\|_2 = \epsilon$ (Section 7).

- The algorithm gives a constructive proof for the existence of a full QR $A = QR$. It also gives, for example, a proof of the existence of the orthogonal complement of the column space of an orthonormal matrix $U$.

- To solve $Ax = b$, solve $Rx = Q^T b$ via triangle solve.
  $\rightarrow$ Excellent method, but twice slower than LU (so it is rarely used)

- The process is aptly called orthogonal triangularisation. (By contrast, Gram-Schmidt and CholeskyQR[15] are triangular orthogonalisation).

---

[15]This algorithm does the following: $A^T A = R^T R$ (Cholesky), then $Q = AR^{-1}$. As stated it's a very fast but unstable algorithm.

## 6.4  Givens rotations

Householder QR is an excellent method for computing the QR factorisation of a general matrix, and it is widely used in practice. However, each Householder reflector acts globally–it affects all the entries of the (active part of) the matrix. For structured matrices—such as sparse matrices—sometimes there is a better tool to reduce the matrix to triangular form (and other forms) by working more locally. Givens rotations give a convenient tool for this. They are matrices of the form

$$G = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \quad c^2 + s^2 = 1.$$

Designed to 'zero' one element at a time. For example to compute the QR factorisation for an upper Hessenberg matrix, one can perform

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix}, \quad G_1 A = \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix}, \quad G_2 G_1 A = \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix},$$

$$G_3 G_2 G_1 A = \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \\ & & & * & * \end{bmatrix}, \quad G_4 G_3 G_2 G_1 A = \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \\ & & & & * \end{bmatrix} =: R.$$

This means $A = G_1^T G_2^T G_3^T G_4^T R$ is the QR factorisation. (note that Givens rotations are orthogonal but not symmetric—so its inverse is $G^T$, not $G$).

- $G$ acts locally on two rows (when left-multiplied; two columns if right-multiplied)

- Non-neighboring rows/cols allowed. For example, a rotation acting on the $i, j$th columns would have $c, s$ values in the $(i, i), (i, j), (j, i), (j, j)$ entries. Visually,

$$G_{i,j} = \begin{bmatrix} 1 & & & & & & & & & & \\ & \ddots & & & & & & & & & \\ & & 1 & & & & & & & & \\ & & & \cos(\theta) & & & & \sin(\theta) & & & \\ & & & & 1 & & & & & & \\ & & & & & \ddots & & & & & \\ & & & & & & 1 & & & & \\ & & & -\sin(\theta) & & & & \cos(\theta) & & & \\ & & & & & & & & 1 & & \\ & & & & & & & & & \ddots & \\ & & & & & & & & & & 1 \end{bmatrix}$$

## 6.5  Least-squares problems via QR

So far we have discussed linear systems wherein the coefficient matrix is square. However, in many situations in data science and beyond, there is a good reason to over-sample to obtain a

robust solution (for instance, in the presence of noise in measurements). For example, when there is massive data and we would like to fit the data with a simple model, we will have many more equations than the degrees of freedom. This leads to the so-called *least-squares problem* which we will discuss here.

Given $A \in \mathbb{R}^{m \times n}, m \geq n$ and $b \in \mathbb{R}^m$, a least-squares problem seeks to find $x \in \mathbb{R}^n$ such that the residual is minimised:

$$\min_x \left\| \begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline x \\ \hline \end{array} - \begin{array}{|c|} \hline b \\ \hline \end{array} \right\|_2 \tag{8}$$

- 'Overdetermined' linear system; attaining equality $Ax = b$ is usually impossible

- Thus the goal is to try minimise the *residual* $\|Ax - b\|$; usually $\|Ax - b\|_2$ but sometimes e.g. $\|Ax - b\|_1$ is of interest. Here we focus on $\|Ax - b\|_2$.

- Throughout we assume full rank condition $\text{rank}(A) = n$; this makes the solution unique and is generically satisfied. If not, the problem will have infinitely many minimisers (and a standard practice is to look for the minimum-norm solution).

Here is how we solve the least-squares problem (8).

**Theorem 6.1** *Let $A \in \mathbb{R}^{m \times n}, m > n$ and $b \in \mathbb{R}^m$, with $\text{rank}(A) = n$. The least-squares problem $\min_x \|Ax - b\|_2$ has solution given by $x = R^{-1}Q^T b$, where $A = QR$ is the (thin) QR factorisation.*

**Proof:** Let $A = [Q \ Q_\perp] \begin{bmatrix} R \\ 0 \end{bmatrix} = Q_F \begin{bmatrix} R \\ 0 \end{bmatrix}$ be 'full' QR factorisation. Then

$$\|Ax - b\|_2 = \|Q_F^T(Ax - b)\|_2 = \left\| \begin{bmatrix} R \\ 0 \end{bmatrix} x - \begin{bmatrix} Q^T b \\ Q_\perp^T b \end{bmatrix} \right\|_2$$

so $x = R^{-1}Q^T b$ is solution. $\square$

This also gives an algorithm (which is essentially the workhorse algorithm used in practice):

1. Compute **thin** QR factorisation $A = QR$ (using Householder QR)

2. Solve linear system $Rx = Q^T b$.


- This process is backward stable. That is, the computed $\hat{x}$ solution for $\min_x \|(A + \Delta A)x + (b + \Delta b)\|_2$ (see Higham's book Ch.20)

- Unlike square system $Ax = b$, one really needs QR: LU won't do the job at all.

One might wonder why we chose the 2-norm in the least-squares formulation (8). Unlike for low-rank approximation (where the truncated SVD is a solution for any unitarily invariant norm[16]) the choice of the norm does matter and affects the properties of the solution $x$ significantly. For example, an increasingly popular choice of norm is the 1-norm, which tends to promote sparsity in the quantity to be minimised. In particular, if we simply replace the tune alarm with the 1 norm, the solution tends to give a residual that is sparse. (nonexaminable)

## 6.6 QR-based algorithm for linear systems

It is straightforward to see that the exact same algorithm can be applied for solving square linear systems. Is this algorithm good? Absolutely! It turns out that it is even better than the LU-based method in that backward stability can be guaranteed (which isn't the case with pivoted LU). However, it is unfortunately twice expensive; which is the reason LU is used in the vast majority of cases for solving linear systems.

Another very stable algorithm is to compute the SVD $A = U\Sigma V^T$ and take $x = V\Sigma^{-1}U^T b$. This is even more expensive than via QR (by $\approx$ x10).

## 6.7 Solution of least-squares via normal equation

There is another way to solve the least-squares problem, by the so-called normal equation. We've seen that
$$\min_x \|Ax - b\|_2, \qquad A \in \mathbb{R}^{m\times n}, m \geq n$$
$x = R^{-1}Q^T b$ is the solution $\Leftrightarrow x$ solution for $n \times n$ **normal equation**

$$(A^T A)x = A^T b$$

- $A^T A \succeq 0$ (always) and $A^T A \succ 0$ if rank$(A) = n$; then PD linear system; use Cholesky to solve.

- This is fast! but NOT backward stable; $\kappa_2(A^T A) = (\kappa_2(A))^2$ where $\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$ **condition number** (next topic)

In fact, more generally, given a linear least-squares approximation problem $\min_{p\in\mathcal{P}} \|f - p\|$ in an inner-product space (of which (8) is a particular example; other examples include polynomial approximation of a function with the inner product e.g. $\langle f, g \rangle = \int_{-1}^{1} f(x)g(x)dx$) the solution is characterized by the property that the residual is orthogonal to the subspace from which a solution is sought, that is, $\langle q, f - p \rangle = 0$ for all $q \in \mathcal{P}$. To see this, consider the problem of approximating $f$ with $p$ in a subspace $\mathcal{P}$. Let $p_* \in \mathcal{P}$ be such that the residual $f - p_+$ is orthogonal to any element in $\mathcal{P}$. Then for any $q \in \mathcal{P}$, we have $\|f - (p_* + q)\|^2 =$

---

[16]Just to be clear, if one uses a norm that is not unitarily invariant (e.g. 1-norm), the truncated SVD may cease to be a solution for the low-rank approximation problem.

$\|f - p_*\|^2 - 2\langle f - p_*, q \rangle + \|q\|^2 = \|f - p_*\|^2 + \|q\|^2 \geq \|f - p_*\|^2$, proving $p_*$ is a minimiser (it is actually unique).

Since we mentioned Cholesky, let us now revisit (7) and show why the second term there must be PSD. A PD matrix has an eigenvalue decomposition $A = VD^2V^T = (VDV^T)^2 = (VDV^T)^T(VDV^T)$. Now let $VDV^T = QR$ be the QR factorisation. Then $(VDV^T)^T(VDV^T) = R^TR$ (this establishes the existence of Cholesky). But now the 0-structure in (7) means the first term must be $rr^T$ where $r^T$ is the first row of $R$, and hence the second term must be $R_2^TR_2$, which is PSD. Here $R^T = [r \ R_2^T]$.

## 6.8 Application of least-squares: regression/function approximation

To illustrate the usefulness of least-squares problems as compared with linear systems here let's consider a function approximation problem.

Given function $f : [-1, 1] \to \mathbb{R}$,

Consider approximating via polynomial $f(x) \approx p(x) = \sum_{i=0} c_i x^i$.

Very common technique: **Regression**: this is a very widely applicable problem in statistics and data science.

1. Sample $f$ at points $\{z_i\}_{i=1}^m$, and

2. Find coefficients $c$ defined by Vandermonde system $Ac \approx f$,

$$\begin{bmatrix} 1 & z_1 & \cdots & z_1^n \\ 1 & z_2 & \cdots & z_2^n \\ \vdots & \vdots & & \vdots \\ 1 & z_m & \cdots & z_m^n \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_n \end{bmatrix} \approx \begin{bmatrix} f(z_1) \\ f(z_2) \\ \vdots \\ f(z_m) \end{bmatrix}.$$

- Numerous applications, e.g. in statistics, numerical analysis, approximation theory, data analysis!

**Illustration** We illustrate this with an example where we approximate the function $f(x) = 1 + \sin(10x)\exp(x)$ (which we suppose we don't know but we can sample it).



$m = n = 11$ (degree 10 polynomial)



$m = 100, n = 11$

We observe that with 11 (equispaced) sample points, the degree-10 polynomial is deviating from the 'true function' quite a bit. With many more sample points the situation significantly improves. This is not a cherry-picked example but a phenomenon that can be mathematically proved; look for "Lebesgue constant" if interested (nonexaminable).

# 7  Numerical stability

An important aspect that is very often overlooked in numerical computing is *numerical stability*. Very roughly, it concerns the quality of a solution obtained by a numerical algorithm, given that computation on computers is done not exactly but with rounding errors. So far we have mentioned stability here and there in passing, but in this section it will be our focus.

Let us first look at an example where roundoff errors play a visible role to affect the computed solution of a linear system.

The situation is complicated. For example, let $A = U\Sigma V^T$, where $U = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$, $\Sigma = \begin{bmatrix} 1 & \\ & 10^{-15} \end{bmatrix}$, $V = I$, and let $b = A\begin{bmatrix} 1 \\ 1 \end{bmatrix}$. That is, we are solving a linear system whose solution is $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$.

If we solve this in MATLAB using $\mathtt{x = A\backslash b}$, the output is $\begin{bmatrix} 1.0000 \\ 0.94206 \end{bmatrix}$. Quite different from the exact solution! Did something go wrong? Did MATLAB or the algorithm fail? The answer is NO, MATLAB and the algorithm (LU) performed just fine. This is a ramification of ill-conditioning, not instability. Make sure that after covering this section, you will be able to explain what happened in the above example.

## 7.1  Floating-point arithmetic

The IEEE (double precision) floating point arithmetic is by far the most commonly used model of computation adopted in computation, and we will assume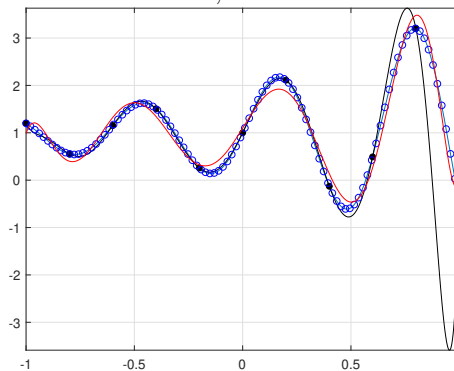 its use here. We will not get into its details as it becomes too computer scientific, rather than mathematical. Just to give a very sketchy introduction

- Computers store number in base 2 with finite/fixed memory (bits)

- Numbers not exactly representable with finite bits in base 2, including irrational numbers, are stored inexactly (rounded), e.g. $1/3 \approx 0.333...$ The unit with which rounding takes place is the *machine precision*, often denoted by $\epsilon$ (or $u$ for unit roundoff). In the most standard setting of IEEE double-precision arithmetic, $u \approx 10^{-16}$.

- Whenever calculations (addition, subtraction, multiplication, division) are performed, the result is rounded to the nearest floating-point number (rounding error); this is where numerical errors really creep in

- Thus the accuracy of the final error is nontrivial; in pathological cases, it is rubbish!

To get an idea of how things can go wrong and how serious the final error can be, here are two examples with MATLAB:

- $((\texttt{sqrt(2)})^2 - 2) * \texttt{1e15} = \texttt{0.4441}$ (should be 0..)

- $\sum_{n=1}^{\infty} \frac{1}{n} \approx 30$ (should be $\infty$..)

For matrices, there are much more nontrivial and surprising phenomena than these. An important (but not main) part of numerical analysis/NLA is to study the effect of rounding errors. This topic can easily span a whole course. By far the best reference on this topic is Higham's book [17].

In this section we denote by $fl(X)$ a computed version of $X$ (fl stands for floating point). For basic operations such as addition and multiplication, one has $fl(x + y) = x + y + c\epsilon$ where $|c| \leq \max(|x|, |y|)$ and $fl(xy) = xy + c\epsilon$ where $|c| \leq \max(|xy|)$.

## 7.2   Conditioning and stability

It is important to solidify the definition of **stability** (which is a property of an algorithm) and **conditioning** (which concerns the sensitivity of a problem and has nothing to do with the algorithm used to solve it).

- Conditioning is the sensitivity of a problem (e.g. of finding $y = f(x)$ given $x$) to perturbation in inputs, i.e., how large $\kappa := \sup_{\delta x} \|f(x + \delta x) - f(x)\| / \|\delta x\|$ is in the limit $\delta x \to 0$. (Very informally, one can think of conditioning as the largest directional derivative).

  (this is the *absolute* condition number; equally important is the *relative* condition number $\kappa_r := \sup_{\delta x} \frac{\|f(x+\delta x)-f(x)\|}{\|f(x)\|} / \frac{\|\delta x\|}{\|x\|}$ )

- (Backward) Stability is a property of an algorithm, which describes if the computed solution $\hat{y}$ is a 'good' solution, in that it is an exact solution of a nearby input, that is, $\hat{y} = f(x + \Delta x)$ for a small $\Delta x$: if $\|\Delta x\|$ can be shown to be small, $\hat{y}$ is a backward stable solution. If an algorithm is guarantee to output a backward stable solution, that algorithm is called backward stable. Throughout this section, $\Delta$ denotes a quantity that is small relative to the quantity to follow: $\|\Delta X\| / \|X\| = O(\epsilon)$, where $\epsilon$ is the machine precision.

To repeat, conditioning is intrinsic in the problem. Stability is a property of an algorithm. Thus we will never say "this problem is backward stable" or "this algorithm is ill-conditioned". We can say "this problem is ill/well-conditioned", or "this algorithm is/isn't

(backward) stable". If a problem is ill-conditioned $\kappa \gg 1$, and the computed solution is no very accurate, then one should blame the problem, not the algorithm. In such cases, a backward stable solution (see below) is usually still considered a good solution.

Notation/convention in this section: $\hat{x}$ denotes a computed approximation to $x$ (e.g. of $x = A^{-1}b$). $\epsilon$ denotes a small term $O(u)$, on the order of unit roundoff/working precision; so we write e.g. $u, 10u, (m+n)u, mnu$ all as $\epsilon$. (In other words, here we assume $m, n \ll u^{-1}$.)

- Consequently (in this lecture/discussion) the norm choice does not matter for the discussion.

## 7.3 Numerical stability; backward stability

Let us dwell more on (backward) stability, because it is really at the heart of the discussion on numerical stability. The word *backward* is key here and it probably differs from the natural notion of stability that you might first think of.

For a computational task $Y = f(X)$ (given input $X$, compute $Y$) and computed approximant $\hat{Y}$,

- Ideally, error $\|Y - \hat{Y}\|/\|Y\| = \epsilon$: but this is seldom true, and often impossible!
  ($u$: unit roundoff, $\approx 10^{-16}$ in standard double precision)

- Good alg has Backward stability $\hat{Y} = f(X + \Delta X)$, $\frac{\|\Delta X\|}{\|X\|} = \epsilon$ "exact solution of a slightly wrong input".

- Justification: The input (matrix) is usually inexact anyway, as storing it on a computer as a floating-point object already incurs rounding errors! Consequently, $f(X + \Delta X)$ is just as good at $f(X)$ at approximating $f(X_*)$ where $\|\Delta X\| = O(\|X - X_*\|)$.

  We shall 'settle with' such solution, though it may not mean $\hat{Y} - Y$ is small.

- Forward stability[17] $\|Y - \hat{Y}\|/\|Y\| = O(\kappa(f)u)$ "error is as small as backward stable alg".

- Another important notion: mixed forward-backward stability: "The computed output is a slightly perturbed solution for a slightly perturbed problem".

---

[17]The definition here follows Higham's book [17]. The phrase is sometimes used to mean small error; However, as hopefully you will be convinced after this section, it is very often impossible to get a solution of full accuracy if the original problem was ill-conditioned. The notion of backward stability that we employ here (and in much of numerical analysis) is therefore much more realistic.

## 7.4  Matrix condition number

The best way to illustrate conditioning is to look at the conditioning of linear systems. In fact it leads to the following definition, which arises so frequently in NLA that it merits its own name: the *condition number* of a matrix.

**Definition 7.1** *The matrix condition number is defined by*

$$\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}(\geq 1).$$

*That is, $\kappa_2(A) = \frac{\sigma_1(A)}{\sigma_n(A)}$ for $A \in \mathbb{R}^{m \times n}, \quad m \geq n$.*

Let's see how this arises:

**Theorem 7.1** *Consider a backward stable solution for $Ax = b$, s.t. $(A + \Delta A)\hat{x} = b$ with $\|\Delta A\| \leq \epsilon\|A\|$ and $\kappa_2(A) \ll \epsilon^{-1}$ (so $\|A^{-1}\Delta A\| \ll 1$). Then we have*

$$\frac{\|\hat{x} - x\|}{\|x\|} \leq \epsilon\kappa_2(A) + O(\epsilon^2).$$

**Proof:**  By Neumann series

$$(A + \Delta A)^{-1} = (A(I + A^{-1}\Delta A))^{-1} = (I - A^{-1}\Delta A + O(\|A^{-1}\Delta A\|^2))A^{-1}$$

So $\hat{x} = (A+\Delta A)^{-1}b = A^{-1}b - A^{-1}\Delta A A^{-1}b + O(\|A^{-1}\Delta A\|^2) = x - A^{-1}\Delta A x + O(\|A^{-1}\Delta A\|^2)$, Hence

$$\|x - \hat{x}\| \lesssim \|A^{-1}\Delta A x\| \leq \|A^{-1}\|\|\Delta A\|\|x\| \leq \epsilon\|A\|\|A^{-1}\|\|x\| = \epsilon\kappa_2(A)\|x\|.$$

$\square$

In other words, even with a backward stable solution, one would only have $O(\kappa_2(A)\epsilon)$ relative accuracy in the solution. If $\kappa_2(A)\epsilon > 1$, the solution may be rubbish! But the NLA view is that's not the fault of the algorithm, the blame is on the problem being so ill-conditioned.

### 7.4.1  Backward stable+well conditioned=accurate solution

We've seen that backward stability does not necessarily imply the solution is accurate. There is a happier side of this argument and useful rule of thumb that can be used to estimate the accuracy of a computed solution using backward stability and conditioning.

Suppose

- $Y = f(X)$ is computed backward stably i.e., $\hat{Y} = f(X + \Delta X)$, $\|\Delta X\| = \epsilon$.

- Conditioning $\|f(X) - f(X + \Delta X)\| \lesssim \kappa\|\Delta X\|$.

Then      (this is the absolute version, relative version possible)

$$\|\hat{Y} - Y\| \lesssim \kappa\epsilon.$$

'proof':
$$\|\hat{Y} - Y\| = \|f(X + \Delta X) - f(X)\| \lesssim \kappa\|\Delta X\|\|f(X)\| = \kappa\epsilon.$$

Here is how to interpret the result: If the problem is well-conditioned $\kappa = O(1)$, this immediately implies good accuracy of the solution! But otherwise the solution might have poor accuracy—but it is still the exact solution of a nearby problem. This is often as good as one can possibly hope for.

The reason this is only a rule of thumb and not exactly rigorous is that conditioning only examines the asymptotic behavior, where the perturbation is infinitesimally small. Nonetheless it often gives an indicative estimate for the error and sometimes we can get rigorous bounds if we know more about the problem. Important examples include the following:

- Well-conditioned linear system $Ax = b$, $\kappa_2(A) \approx 1$.

- Eigenvalues of symmetric matrices (via Weyl's bound $\lambda_i(A+E) \in \lambda_i(A)+[-\|E\|_2, \|E\|_2])$.

- Singular values of any matrix $\sigma_i(A + E) \in \sigma_i(A) + [-\|E\|_2, \|E\|_2]$.

Indeed, these problems are well-conditioned, so can be solved with extremely high accuracy, essentially to working precision $O(u)$ (times a small factor, usually bounded by something like $\sqrt{n}$).

Note: eigvecs/singvecs can be highly ill-conditioned even for the simplest problems. Again, think of the identity. (Those curious are invited to look up the Davis-Kahan $\sin\theta$ theorem [7].)

## 7.5    Stability of triangular systems

### 7.5.1    Backward stability of triangular systems

While we will not be able to discuss in detail which NLA algorithm is backward stable etc, we will pick a few important examples.

One fact is worth knowing (a proof is omitted; see [33] or [17]): triangular linear systems can be solved in a backward stable manner. This fact is important as these arise naturally in the solution of linear systems: $Ax = b$ is solved via $Ly = b$, $Ux = y$ (triangular systems), as we've seen in Section 5.

Let $R$ denote a matrix that is (upper or lower) triangular. The computed solution $\hat{x}$ for a (upper/lower) triangular linear system $Rx = b$ solved via back/forward substitution is backward stable, i.e., it satisfies

$$(R + \Delta R)\hat{x} = b, \qquad \|\Delta R\| = O(\epsilon\|R\|).$$

Proof: Trefethen-Bau or Higham (nonexaminable but interesting).
Notes:

- The backward error can be bounded componentwise.

- Using the previous rule-of-thumb, this means $\|\hat{x} - x\|/\|x\| \lesssim \epsilon \kappa_2(R)$.

  - (unavoidably) poor worst-case (and attainable) bound when ill-conditioned
  - often better with triangular systems; the behavior of triangular linear systems keep surprising experts!

### 7.5.2 (In)stability of $Ax = b$ via LU with pivots

We have discussed how to solve a linear system using the LU factorisation. An obvious question given the context is: is it backward stable? This question has a fascinating answer, and transpires to touch on one of the biggest open problems in the field.

Fact (proof nonexaminable): Computed $\hat{L}\hat{U}$ satisfies $\frac{\|\hat{L}\hat{U} - A\|}{\|L\|\|U\|} = \epsilon$.

(note: not $\frac{\|\hat{L}\hat{U} - A\|}{\|A\|} = \epsilon$)

- By stability of triangular systems $(L + \Delta L)(U + \Delta U)\hat{x} = b$. Now if $\|L\|\|U\| = O(\|A\|)$, then it follows that $\Rightarrow \hat{x}$ backward stable solution (exercise).

**Question**: Does $LU = A + \Delta A$ or $LU = PA + \Delta A$ with $\|\Delta A\| = \epsilon\|A\|$ (i.e., $\|L\|\|U\| = O(\|A\|)$) hold?

Without pivot $(P = I)$: $\|L\|\|U\| \gg \|A\|$ unboundedly (e.g. $\begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}$) unstable.
With pivots:

- Worst-case: $\|L\|\|U\| \gg \|A\|$ grows exponentially with $n$, unstable.

  - growth governed by that of $\|L\|\|U\|/\|A\| \Rightarrow \|U\|/\|A\|$.

- In practice (average case): perfectly stable.

  - Hence this is how $Ax = b$ is solved, despite alternatives with guaranteed stability exist (but slower; e.g. via SVD, or QR (next)).

Resolution/explanation: among biggest open problems in numerical linear algebra!

### 7.5.3 Backward stability of Cholesky for $A \succ 0$

We've seen that symmetric positive definite matrices can be solved with half the effort using the Cholesky factorisation. As a bonus, in this case it turns out that stability can be guaranteed.

The key fact is that the Cholesky factorisation $A = R^T R$ for $A \succ 0$

- succeeds without pivot (the active matrix is always positive definite).

- $R$ never contains entries $> \sqrt{\|A\|}$.

It follows that computing the Cholesky factorisation (assuming $A \succ 0$) is backward stable! Hence positive definite linear system $Ax = b$ can always be solved in a backward stable manner via Cholesky.

## 7.6 Matrix multiplication is not backward stable

Here is perhaps a shock—matrix matrix multiplication, one of the most basic operations, is in general not backward stable.

Let's start with the basics.

- Vector-vector multiplication is backward stable: $fl(y^T x) = (y + \Delta y)(x + \Delta x)$; in fact $fl(y^T x) = (y + \Delta y)x$. (direct proof is possible)

- It immediately follows that matrix-vector multiplication is also backward stable: $fl(Ax) = (A + \Delta A)x$.

- But it is not true to say matrix-matrix multiplication is backward stable; which would require $fl(AB)$ to be equal to $(A + \Delta A)(B + \Delta B)$. This may not be satisfied!

In general, what we can prove is a bound for the forward error $\|fl(AB) - AB\| \leq \epsilon\|A\|\|B\|$, so $\|fl(AB) - AB\|/\|AB\| \leq \epsilon \min(\kappa_2(A), \kappa_2(B))$ (proof: problem sheet).

This is great news when $A$ or $B$ is orthogonal (or more generally square and well-conditioned): say if $A = Q$ is orthogonal, then we have

$$\|fl(QB) - QB\| \leq \epsilon\|B\|,$$

so it follows that $fl(QB) = QB + \epsilon\|B\|$, hence defining $\Delta B = Q^T \epsilon\|B\|$ we have $fl(QB) = Q(B + \Delta B)$, that is, **orthogonal multiplication is backward stable** (this argument proves this for left multiplication; orthogonal right-multiplication is entirely analogous).

One of the reasons backward stability of matrix-matrix multiplication fails to hold is that there are not enough parameters in the backward error to account for the computational error incurred in the matrix multiplication. Each matrix-vector product is backward stable; but we cannot concentrate the backward errors into a single term without potentially increasing the backward error's norm.

On the other hand, we have seen that a linear system can be solved in a backward stable fashion (by QR if necessary; in practice LU with pivoting suffices). This means the inverse can be applied to a vector in a backward stable fashion: The computed $\hat{x}$ satisfies $(A + \Delta A)^{-1}\hat{x} = b$.

One might wonder, can we not solve $n$ linear systems in order to get a backward stable inverse? However, if one solves $n$ linear systems with $b = e_1, e_2, \ldots, e_n$, the solutions will satisfy $(A + \Delta_i A)^{-1}\hat{x}_i = e_i$, where $\Delta_i A$ depends on $i$. There is no uniform $\Delta A$ that is $O(\epsilon)$ such that $(A + \Delta A)^{-1}[\hat{x}_1, \ldots, \hat{x}_n] = I$.

**Orthogonality matters for stability**   A happy and important exception is with orthogonal matrices $Q$ (or more generally with well-conditioned matrices):

$$\frac{\|fl(QA) - QA\|}{\|QA\|} \leq \epsilon, \qquad \frac{\|fl(AQ) - AQ\|}{\|AQ\|} \leq \epsilon.$$

They are also backward stable:

$$fl(QA) = QA + \epsilon \quad \Leftrightarrow \quad fl(QA) = Q(A + \Delta A).$$
$$fl(AQ) = AQ + \epsilon \quad \Leftrightarrow \quad fl(AQ) = (A + \Delta A)Q.$$

Hence algorithms involving ill-conditioned matrices are unstable (e.g. eigenvalue decomposition of non-normal matrices, Jordan form, etc), whereas those based on orthogonal matrices are stable. These include

- Householder QR factorisation, QR-based linear system (next subsection).

- **QR algorithm** for $Ax = \lambda x$.

- **Golub-Kahan** algorithm for $A = U\Sigma V^T$.

- **QZ algorithm** for $Ax = \lambda Bx$.

Section 8 onwards treats our second big topic, eigenvalue problems. This includes discussing algorithms shown above in boldface.

## 7.7   Stability of Householder QR

Householder QR has excellent numerical stability, basically because it's based on orthogonal transformations. With Householder QR, the computed $\hat{Q}, \hat{R}$ satisfy

$$\|\hat{Q}^T\hat{Q} - I\| = O(\epsilon), \quad \|A - \hat{Q}\hat{R}\| = O(\epsilon\|A\|),$$

and (of course) $R$ is upper triangular.

Rough proof: Essentially the key idea is that multiplying by an orthogonal matrix is very stable and Householder QR is based on a sequence of multiplications by orthogonal matrices. To give a bit more detail,

- Each reflector satisfies $fl(H_iA) = H_iA + \epsilon_i\|A\|$.

- Hence $(\hat{R} =)fl(H_n \cdots H_1A) = H_n \cdots H_1A + \epsilon\|A\|$.

- $fl(H_n \cdots H_1) =: \hat{Q}^T = H_n \cdots H_1 + \epsilon$.

- Thus $\hat{Q}\hat{R} = A + \epsilon\|A\|$.

Notes:

- This doesn't mean $\|\hat{Q}-Q\|, \|\hat{R}-R\|$ are small at all! Indeed $Q, R$ are as ill-conditioned as $A$ [17, Ch. 20].

- Solving $Ax = b$ and least-squares problems via the QR factorisation is stable.

### 7.7.1 (In)stability of Gram-Schmidt

(Nonexaminable) A somewhat surprising fact is that the Gram-Schmidt algorithm, when used for computing the QR factorisation, is not backward stable. Namely, orthogonality of the computed $\hat{Q}$ matrix is not guaranteed.

- Gram-Schmidt is subtle:

    - plain (classical) version: $\|\hat{Q}^T \hat{Q} - I\| \le \epsilon(\kappa_2(A))^2$.

    - modified Gram-Schmidt (orthogonalise 'one vector at a time'): $\|\hat{Q}^T \hat{Q} - I\| \le \epsilon\kappa_2(A)$.

    - Gram-Schmidt twice (G-S again on computed $\hat{Q}$) is excellent: $\|\hat{Q}^T \hat{Q} - I\| \le \epsilon$.

# 8 Eigenvalue problems

First of all, recall that $Ax = \lambda x$ has no explicit solution (neither $\lambda$ nor $x$); huge difference from $Ax = b$ for which $x = A^{-1}b$.

From a mathematical viewpoint this marks an enormous point of departure: something that is is explicitly written vs. something that has no closed-from solution.

From a practical viewpoint the gulf is much smaller, because we have an extremely reliable algorithm for eigenvalue problems; so robust that it is essentially bulletproof provably backward stable.

Before we start describing the QR algorithm let's discuss a few interesting properties of eigenvalues.

- Eigenvalues are the roots of characteristic polynomial (Abel's theorem).

- For any polynomial $p$, there exist (infinitely many) matrices whose eigvals are roots of $p$.

- Here is a nice and useful fact[18]: Let $p(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_1 x + a_0$, $a_i \in \mathbb{C}$. Then
$p(\lambda) = 0 \Leftrightarrow \lambda$ eigenvalue of

$$
C = \begin{bmatrix}
-a_{n-1} & -a_{n-2} & \ldots & -a_1 & -a_0 \\
1 & & & & \\
& 1 & & & \\
& & \ddots & & \\
& & & 1 & 0
\end{bmatrix} \in \mathbb{C}^{n \times n}.
$$

---

[18] A great application of this is univariate optimisation: to minimise a polynomial $p(x)$, one can find the critical points by solving for $p'(x) = 0$, via the companion eigenvalues.

- So no finite-step algorithm exists for $Ax = \lambda x$.

Eigenvalue algorithms are necessarily iterative and approximate.

- Same for SVD, as $\sigma_i(A) = \sqrt{\lambda_i(A^T A)}$.

- But this doesn't mean they're inaccurate!

Usual goal: compute the Schur decomposition $A = UTU^*$: $U$ unitary, $T$ upper triangular.

## 8.1 Schur decomposition

**Theorem 8.1** *Let $A \in \mathbb{C}^{n \times n}$ (arbitrary square matrix). Then there exists a unitary matrix $U \in \mathbb{C}^{n \times n}$ s.t.*

$$A = UTU^*, \tag{9}$$

*with $T$ upper triangular. The decomposition (9) is called the Schur decomposition.*

Note that

- $\mathrm{eig}(A) = \mathrm{eig}(T) = \mathrm{diag}(T)$.

- $T$ diagonal iff $A$ is normal, i.e., $A^* A = AA^*$

**Proof:** Let $Av = \lambda_1 v$ and find $U_1 = [v_1, V_\perp]$ unitary. Then $AU_1 = U_1 \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix} \Leftrightarrow$

$U_1^* AU_1 = \begin{bmatrix} * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix}$. Repeat on the lower-right $(n-1) \times (n-1)$ part to get

$U_{n-1}^* U_{n-2}^* \dots U_1^* AU_1 U_2 \dots U_{n-1} = T$. $\qquad\qquad\square$

The reason we usually take the Schur form to be the goal is that its computation can be done in a backward stable manner. Essentially it boils down to the fact that the decomposition is based on orthogonal transformations.

- For normal matrices $A^* A = AA^*$, $T$ must be diagonal and Schur form is automatically diagonalised.

- For nonnormal $A$, if diagonalisation $A = X\Lambda X^{-1}$ really necessary, done via starting with the Schur decomposition and further reducing $T$ by solving Sylvester equations; but this process involves nonorthogonal transformations and is not backward stable (nonexaminable).

- The Schur decomposition is among the few examples in NLA where the difference between $\mathbb{C}$ and $\mathbb{R}$ matters a bit. When working only in $\mathbb{R}$ one cannot always get a triangular $T$; it will have $2 \times 2$ blocks in the diagonal (these blocks have complex eigenvalues). This is essentially because real matrices can have complex eigenvalues. This is still not a major issue as eigenvalues of $2 \times 2$ matrices can be computed easily.

(This marks the end of "the first half" of the course (i.e., up until the basic facts about eigenvalues, but not its computation). This information is relevant to MMSC students.)

## 8.2   The power method for finding the dominant eigenpair $Ax = \lambda x$

We now start describing the algorithms for solving eigenvalue problems. The first algorithm that we introduce is surprisingly simple and is based on the idea of keep multiplying the matrix $A$ to an arbitrary vector.

This algorithm by construction is designed to compute only a single eigenvalue and its corresponding eigenvector, namely the dominant one. It is not able, at least as presented, to compute all the eigenvalues.

Despite its limitation and simplicity it is an extremely important algorithm and the underlying idea is in fact a basis for the ultimate QR algorithm that we use in order to compute all eigenvalues of a matrix. So here is the algorithm, called the *power method*.

---

**Algorithm 8.1** The power method for computing the dominant eigenvalue and eigenvector of $A \in \mathbb{R}^{n \times n}$.

1: Let $x \in \mathbb{R}^n$ :=random initial vector (unless specified)
2: Repeat: $x = Ax$, $x = \frac{x}{\|x\|_2}$.
3: $\hat{\lambda} = x^T A x$ gives an estimate for the eigenvalue.

---

- Convergence analysis: suppose $A$ is diagonalisable (generic assumption). We can write $x_0 = \sum_{i=1}^n c_i v_i$, $Av_i = \lambda_i v_i$ with $|\lambda_1| > |\lambda_2| > \cdots$. Then after $k$ iterations,

$$x = C \sum_{i=1}^n \left( \frac{\lambda_i}{\lambda_1} \right)^k c_i v_i \to C c_1 v_1 \quad \text{as } k \to \infty \text{ for some scalar } C$$

- Converges geometrically $(\lambda, x) \to (\lambda_1, x_1)$ with **linear rate** $\frac{|\lambda_2|}{|\lambda_1|}$

- What does this imply about $A^k = QR$ as $k \to \infty$? First column of $Q \to v_1$ under mild assumptions.

Notes:

- Google pagerank & Markov chain are linked to the power method.

- As we'll see, the power method is basis for refined algorithms (QR algorithm, Krylov methods (Lanczos, Arnoldi,...))

### 8.2.1 Digression (optional): Why compute eigenvalues? Google PageRank

Let us briefly digress and talk about a famous application that at least used to solve an eigenvalue problem in order to achieve a familiar task: Google web search.

Once Google receives a user's inquiry with keywords, it needs to rank the relevant webpages, to output the most important pages first.

Here, the 'importance' of websites is determined by the dominant eigenvector of column-stochastic matrix (i.e., column sums are all 1)

$$A = \alpha P + \frac{(1-\alpha)}{n} \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$$

$P$: scaled adjacency matrix ($P_{ij} = 1$ if $i, j$ connected by an edge, 0 otherwise, and then scaled s.t. column sums are 1), $\alpha \in (0, 1)$

**PageRank**

image from wikipedia

To solve this approximately (note that high accuracy isn't crucial here—getting the ordering wrong isn't the end of the world), Google does (at least in the old days) a few steps of power method: with initial guess $x_0$, $k = 0, 1, \ldots$

1. $x_{k+1} = Ax_k$

2. $x_{k+1} = x_{k+1}/\|x_{k+1}\|_2$,     $k \leftarrow k + 1$, repeat.

- $x_k \rightarrow$ PageRank vector $v_1$ : $Av_1 = \lambda_1 v_1$. As $A$ is a nonnegative matrix $A_{ij} \geq 0$, the dominant eigenvector $v_1$ can be taken to be nonnegative (by the Perron-Frobenius theorem; nonexaminable).

For more on pagerank etc, see Gleich [13] if interested (nonexaminable).

### 8.2.2 Shifted inverse power method

We saw that the convergence of the power method is governed by $|\lambda_1/\lambda_2|$, the ratio between the absolute value of dominant and the next dominant eigenvalue.

If this ratio is close to 1 the power method would converge very slowly. A natural question arises: can one speed up the convergence in such cases?

This leads to the *inverse power method*, or more generally shifted inverse power method, which can compute eigenpairs with much improved speed provided that the parameters (shifts) are chosen appropriately.

**Algorithm 8.2** Shifted inverse power method for computing the eigenvalue and eigenvector of $A \in \mathbb{R}^{n \times n}$ closest to a prescribed value $\mu \in \mathbb{C}$.

1: Let $x \in \mathbb{R}^n :=$ random initial vector, unless specified.
2: Repeat: $x := (A - \mu I)^{-1} x$, $x = x / \|x\|$.
3: $\hat{\lambda} = x^T A x$ gives an estimate for the eigenvalue.

Note that the eigenvalues of the matrix $(A - \mu I)^{-1}$ are $(\lambda(A) - \mu)^{-1}$.

- By the same analysis as above, shifted-inverse power method converges with **improved linear rate** $\frac{|\lambda_{\sigma(1)} - \mu|}{|\lambda_{\sigma(2)} - \mu|}$ to the eigenpair closest to $\mu$. Here $\sigma$ denotes a permutation of $\{1, \ldots, n\}$ such that $|\lambda_{\sigma(1)} - \mu|$ minimises $|\lambda_i - \mu|$ over $i$.

- $\mu$ can change adaptively with the iterations. The choice $\mu := x^T A x$ gives the *Rayleigh quotient iteration*, with **quadratic** convergence $\|Ax^{(k+1)} - \lambda^{(k+1)} x^{(k+1)}\| = O(\|Ax^{(k)} - \lambda^{(k)} x^{(k)}\|^2)$; this is further improved to cubic convergence if $A$ is symmetric.

It is worth emphasising that the improved convergence comes at a cost: one step of shifted inverse power method requires a solution of a linear system, which is clearly more expensive than power method.

# 9 The QR algorithm

We'll now describe an algorithm called the QR algorithm that is used universally for solving eigenvalue problems of moderate size, e.g. by MATLAB's `eig`. Given $A \in \mathbb{R}^{n \times n}$, the algorithm

- Finds all eigenvalues (approximately but reliably) in $O(n^3)$ flops,

- Is backward stable.

Sister problem: Given $A \in \mathbb{R}^{m \times n}$ compute SVD $A = U\Sigma V^T$

- 'ok' algorithm: $\text{eig}(A^T A)$ to find $V$, then normalise $AV$

- there's a better (but still closely related) algorithm: Golub-Kahan bidiagonalisation, discussed later in Section 10.2.

## 9.1 QR algorithm for $Ax = \lambda x$

As the name suggests, the QR algorithm is based on the QR factorisation of a matrix. Another key fact is that the eigenvalues of a product of two matrices remain the same when the order of the product is reversed ($\text{eig}(AB) = \text{eig}(BA)$, problem sheet). The QR algorithm essentially is a simple combination of these two ideas, and the vanilla version is deceptively

simple: basically take the QR factorisation, swap the order, take the QR and repeat the process. Namely,

---

**Algorithm 9.1** The QR algorithm for finding the Schur decomposition of a square matrix $A$.

---

1: Set $A_1 = A$.
2: Repeat: $A_1 = Q_1 R_1, \quad A_2 = R_1 Q_1, \quad A_2 = Q_2 R_2, \quad A_3 = R_2 Q_2, \quad \dots$

---

Notes:

- $A_k$ are all similar: $A_{k+1} = Q_k^T A_k Q_k$

- We shall 'show' that $A \to$ triangular **triangular** (diagonal if $A$ normal), under weak assumptions.

- Basically: $QR(\text{factorise}) \to RQ(\text{swap}) \to QR \to RQ \to \cdots$

- Fundamental work by Francis (61,62) and Kublanovskaya (63)

- Truly Magical algorithm!

  - The algorithm is backward stable, as based on orthogonal transforms: essentially, $RQ = fl(Q^T(QR)Q) = Q^T(QR + \Delta(QR))Q$.
  - always converges (with shifts) in practice, but a global proof is currently unavailable(!)
  - uses 'shifted inverse power method' (rational functions) without inversions

Again, the QR algorithm performs: $A_k = Q_k R_k$, $A_{k+1} = R_k Q_k$, repeat.
Let's look at its properties.

**Theorem 9.1** *For $k \geq 1$,*

$$A_{k+1} = (Q^{(k)})^T A Q^{(k)}, \qquad A^k = (Q_1 \cdots Q_k)(R_k \cdots R_1) =: Q^{(k)} R^{(k)}.$$

Proof : recall $A_{k+1} = Q_k^T A_k Q_k$, repeat.

Proof by induction: $k = 1$ trivial.
Suppose $A^{k-1} = Q^{(k-1)} R^{(k-1)}$. We have

$$A_k = (Q^{(k-1)})^T A Q^{(k-1)} = Q_k R_k.$$

Then $A Q^{(k-1)} = Q^{(k-1)} Q_k R_k$, and so

$$A^k = A Q^{(k-1)} R^{(k-1)} = Q^{(k-1)} Q_k R_k R^{(k-1)} = Q^{(k)} R^{(k)} \square$$

48

**QR algorithm and power method**  We can deduce that the QR algorithm is closely related to the power method. Let's try explain the connection. By Theorem 9.1, $Q^{(k)}R^{(k)}$ is the QR factorisation of $A^k$: as we saw in the analysis of the power method, the columns of $A^k$ are 'dominated by the leading eigenvector' $x_1$, where $Ax_1 = \lambda_1 x_1$.

In particular, consider $A^k[1, 0, \ldots, 0]^= A^k e_n$:

- $A^k e_n = R^{(k)}(1,1)Q^{(k)}(:,1)$, which is parallel to the first column of $Q^{(k)}$.

- By the analysis of the power method, this implies $Q^{(k)}(:,1) \to x_1$

- Hence by $\boxed{A_{k+1} = (Q^{(k)})^A Q^{(k)}}$ , $A_k(:,1) \to [\lambda_1, 0, \ldots, 0]^T$.

This tells us why the QR algorithm would eventually compute the eigenvalues—at least the dominant ones. One can even go further and argue that once the dominant eigenvalue has converged, we can expect the next dominant eigenvalue to start converging too—as due to the nature of the QR algorithm based on orthogonal transformations, we are then working in the orthogonal complement; and so on and so forth until the matrix $A$ becomes upper triangular (and in the normal case this becomes diagonal), completing the solution of the eigenvalue problem. But there is much better news.

**QR algorithm and inverse power method**  We have seen that the QR algorithm is related to the power method. We have also seen that the power method can be improved by a shift-and-invert technique. A natural question arises: can we introduce a similar technique in the QR algorithm? This question has an elegant and remarkable answer: not only is this possible but it is possible without ever inverting a matrix or solving a linear system (isn't that incredible!?).

Let's try and explain this. We start with the same expression for the QR iterates:

$$\boxed{A^k = (Q_1 \cdots Q_k)(R_k \cdots R_1) =: Q^{(k)}R^{(k)}} , \qquad \boxed{A_{k+1} = (Q^{(k)})^T A Q^{(k)}.}$$

Now take inverse: $A^{-k} = (R^{(k)})^{-1}(Q^{(k)})^T$,

Conjugate transpose: $(A^{-k})^T = Q^{(k)}(R^{(k)})^{-*}$

$\Rightarrow$ QR factorisation of matrix $(A^{-k})^T$ with eigvals $r(\lambda_i) = \boxed{\lambda_i^{-k}}$

$\Rightarrow$ Connection also with (unshifted) inverse power method. Note that no matrix inverse is performed in the algorithm.

- This means the final column of $Q^{(k)}$ converges to minimum left eigenvector $x_n$ with rate $\frac{|\lambda_n|}{|\lambda_{n-1}|}$, hence $A_k(n,:) \to [0, \ldots, 0, \lambda_n]$.

- (Very) fast convergence if $|\lambda_n| \ll |\lambda_{n-1}|$.

- Can we achieve this situation? **Yes by shifts**.

**QR algorithm with shifts and shifted inverse power method**  We are now ready to reveal the connection between the shift-and-invert power method and the QR algorithm with shifts.

First, here is the QR algorithm with shifts.

---

**Algorithm 9.2** The QR algorithm with shifts for finding the Schur decomposition of a square matrix $A$.

---

1: Set $A_1 = A$, $k = 1$.
2: $A_k - s_k I = Q_k R_k$ (QR factorisation)
3: $A_{k+1} = R_k Q_k + s_k I$, $\quad k \leftarrow k + 1$, repeat.

---

1. $A_k - s_k I = Q_k R_k$ (QR factorisation)

2. $A_{k+1} = R_k Q_k + s_k I$, $\quad k \leftarrow k + 1$, repeat.

Roughly, if $s_k \approx \lambda_n$, then $A_{k+1} \approx \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ & & & & \lambda_n \end{bmatrix}$ by the argument just made.

**Theorem 9.2**

$$\prod_{i=1}^{k}(A - s_i I) = Q^{(k)} R^{(k)} \left(= (Q_1 \cdots Q_k)(R_k \cdots R_1)\right).$$

**Proof:**  Suppose true for $k-1$. Then the QR alg computes $(Q^{(k-1)})^T (A - s_k I) Q^{(k-1)} = Q_k R_k$, so $(A - s_k I) Q^{(k-1)} = Q^{(k-1)} Q_k R_k$, hence

$$\prod_{i=1}^{k}(A - s_i I) = (A - s_k I) Q^{(k-1)} R^{(k-1)} = Q^{(k-1)} Q_k R_k R^{(k-1)} = Q^{(k)} R^{(k)}.$$

Inverse conjugate transpose: $\prod_{i=1}^{k}(A - s_i I)^{-*} = Q^{(k)}(R^{(k)})^{-*}$. $\qquad\qquad\square$

- This means the algorithm (implicitly) finds the QR factorisation of a matrix with eigvals $r(\lambda_j) = \boxed{\prod_{i=1}^{k} \frac{1}{\lambda_j - s_i}}$ .

- This reveals the intimate connection between shifted QR and shifted inverse power method, hence  rational approximation .

- Ideally, we would like to choose $s_k \approx \lambda_n$ to accelerate convergence. This is done by choosing $s_k$ to be the bottom-right corner entry[19] of $A_k$; which is sensible given that with the QR algorithm, it tends to converge to an eigenvalue rapidly (with a few steps of the unshifted QR, it tends to converge to the smallest eigenvalue).

---

[19]In production code a so-called Wilkinson shift is often used, which computes the eigenvalue of the $2 \times 2$ bottom-right submatrix of $A_k$, and takes $s_k$ to be the one closer to the bottom-right entry of $A_k$.

## 9.2 QR algorithm preprocessing: reduction to Hessenberg form

We've seen the QR iterations drives colored entries to 0 (esp. red ones)

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}$$

- Hence $A_{n,n} \to \lambda_n$, so choosing $s_k = A_{n,n}$ is sensible.

- This reduces #QR iterations to $O(n)$ (empirical but reliable estimate).

- But each iteration of the QR algorithm is $O(n^3)$ for QR, overall $O(n^4)$.

- We next discuss a preprocessing technique to reduce the overall cost to $O(n^3)$.

The idea is to initially apply a series of deterministic orthogonal transformations that reduces the matrix to a form that is closer to upper triangular. This is done before starting the QR iterates, hence called a preprocessing step.

More precisely, to improve the cost of QR factorisation, we first reduce the matrix via orthogonal Householder transformations as follows:

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}, \quad H_1 A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix}, \quad H_1 = I - 2v_1 v_1^T, \ v_1 = \begin{bmatrix} 0 \\ * \\ * \\ * \\ * \end{bmatrix}$$

Then $H_1 A H_1 = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix}$. Repeat with $H_2 = I - 2v_2 v_2^T, v_2 = [0,0,*,*,*]^T$, ...:

$$H_2 H_1 A H_1 H_2 = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & * & * & * \end{bmatrix}, \qquad H_3 H_2 H_1 A H_1 H_2 H_3 = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix},$$

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{H_1} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix} \xrightarrow{H_2} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & * & * & * \end{bmatrix} \xrightarrow{H_3} \cdots \xrightarrow{H_{n-2}} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix}.$$

We have thus transformed the matrix $A$ into upper Hessenberg form.

- This helps QR iterations preserve structure: if $A_1 = QR$ Hessenberg, then so is $A_2 = RQ$

- Using Givens rotations, each QR iter is $O(n^2)$ (not $O(n^3)$).

The remaining task (done by shifted QR): drive subdiagonal $*$ to 0. $A$ is thus reduced to a triangular form, revealing the Schur decomposition of $A$ (if one traces back all the orthogonal transformations employed).

Once bottom-right $|*| < \epsilon$,

$$
\begin{bmatrix}
* & * & * & * & * \\
* & * & * & * & * \\
  & * & * & * & * \\
  &   & * & * & * \\
  &   &   & * & *
\end{bmatrix}
\approx
\begin{bmatrix}
* & * & * & * & * \\
* & * & * & * & * \\
  & * & * & * & * \\
  &   & * & * & * \\
  &   &   &   & *
\end{bmatrix}
$$

and continue with shifted QR on $(n-1) \times (n-1)$ block, repeat.

- Empirically, the shifted QR algorithm needs $2 - -4$ iterations per eigenvalue. Overall, the cost is $O(n^3), \approx 25n^3$ flops.

### 9.2.1   The (shifted) QR algorithm in action

Let's see how the QR algorithm works with a small example. The plots show the convergence of the subdiagonal entries $|A_{i+1,i}|$ (note that their convergence signifies the convergence of the QR algorithm as we initially reduce the matrix to Hessenberg form).

No shift (plain QR)    QR with shifts

In light of the connection to rational functions as discussed above, here we plot the underlying functions (red dots: eigvals). The idea here is that we want the functions to take large values at the target eigenvalue (at the current iterate) in order to accelerate convergence.

### 9.2.2 (Optional) QR algorithm: other improvement techniques

We have completed the description of the main ingredients of the QR algorithm. Nonetheless, there are a few more bells and whistles that go into to a production code. We will not get into too much detail but here is a list of the key players.

- Double-shift strategy for $A \in \mathbb{R}^{n \times n}$

    - $(A - sI)(A - \bar{s}I) = QR$ using only real arithmetic

- Aggressive early deflation                                    [Braman-Byers-Mathias 2002 [5]]

    - Examine lower-right (say $100 \times 100$) block instead of $(n, n-1)$ element
    - dramatic speedup ($\approx \times 10$)

- Balancing $A \leftarrow DAD^{-1}$, $D$: diagonal

    - Aims at reducing $\|DAD^{-1}\|$: often yields better-conditioned eigenvalues.

53

Finally, let us mention a cautionary tale:

- For nonsymmetric $A$, global convergence is NOT established(!)

  – Of course it always converges in practice.. proving convergence is another big open problem in numerical linear algebra.

  – For symmetric $A$, global convergence analysis is available, see [28, Ch. 8].

# 10 QR algorithm continued

## 10.1 QR algorithm for symmetric $A$

so far we have not assumed anything about the matrix besides that it is square. This is for good reason—because the QR algorithm is applicable to solving any eigenvalue problem. However, in many situations the eigenvalue problem is *structured*. In particular the case where the matrix is symmetric arises very frequently in practice and it comes with significant simplification of the QR algorithm, so it deserves a special mention.

- Most importantly, symmetry immediately implies that the initial reduction to Hessenberg form reduces $A$ to tridiagonal:

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \overset{Q_1}{\rightarrow} \begin{bmatrix} * & * & & & \\ * & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \\ & * & * & * & * \end{bmatrix} \overset{Q_2}{\rightarrow} \begin{bmatrix} * & * & & & \\ * & * & * & & \\ & * & * & * & * \\ & & * & * & * \\ & & * & * & * \end{bmatrix} \overset{Q_3}{\rightarrow} \begin{bmatrix} * & * & & & \\ * & * & * & & \\ & * & * & * & \\ & & * & * & * \\ & & & * & * \end{bmatrix}$$

- QR steps for tridiagonal: requires $O(n)$ flops instead of $O(n^2)$ per step.

- Powerful alternatives are available for tridiagonal eigenvalue problem (divide-conquer [Gu-Eisenstat 95], HODLR [Kressner-Susnjara 19],...)

- Cost: $\frac{4}{3}n^3$ flops for eigvals, $\approx 10n^3$ for eigvecs (store Givens rotations to compute eigvecs).

- Another approach (nonexaminable): spectral divide-and-conquer (Nakatsukasa-Freund, Higham); which is all about using a carefully chosen <span style="color:red">rational approximation</span> to find orthogonal matrices $V_i$ such that

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \overset{V_1}{\rightarrow} \begin{bmatrix} * & * & * & & \\ * & * & * & & \\ * & * & * & & \\ & & & * & * \\ & & & * & * \end{bmatrix} \overset{V_2}{\rightarrow} \begin{bmatrix} * & & & & \\ & * & * & & \\ & * & * & & \\ & & & * & \\ & & & & * \end{bmatrix} \overset{V_3}{\rightarrow} \begin{bmatrix} * & & & & \\ & * & & & \\ & & * & & \\ & & & * & \\ & & & & * \end{bmatrix} = \Lambda.$$

## 10.2 Computing the SVD: Golub-Kahan's bidiagonalisation algorithm

The key ideas of the QR algorithm turns out to be applicable for the computation of the SVD. Clearly this is a major step—given the importance of the SVD. In particular the SVD algorithm has strong connections to the symmetric QR algorithm. This is perhaps not surprising given the strong connection between the SVD and symmetric eigenvalue problems as we discussed previously.

A noteworthy difference is that in the SVD $A = U\Sigma V^T$ the two matrices $U$ and $V$ are allowed to be different. The algorithm respects this and instead of initially reducing the matrix to tridiagonal form, it reduces it to a so-called bidiagonal form.

Here is how it works. Apply Householder reflectors from left and right (different ones) to **bidiagonalise**

$$A \to B = H_{L,n} \cdots H_{L,1} A H_{R,1} H_{R,2} \cdots H_{R,n-2}$$

$$A \overset{H_{L,1}}{\to}
\begin{bmatrix} \star & \star & \star & \star \\ & \star & \star & \star \\ & \star & \star & \star \\ & \star & \star & \star \\ & \star & \star & \star \end{bmatrix}
\overset{H_{R,1}}{\to}
\begin{bmatrix} \star & \star & & \\ & \star & \star & \star \\ & \star & \star & \star \\ & \star & \star & \star \\ & \star & \star & \star \end{bmatrix}
\overset{H_{L,2}}{\to}
\begin{bmatrix} \star & \star & & \\ & \star & \star & \star \\ & & \star & \star \\ & & \star & \star \\ & & \star & \star \end{bmatrix}
\overset{H_{R,2}}{\to}
\begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & & \star & \star \\ & & \star & \star \\ & & \star & \star \end{bmatrix}
\overset{H_{L,3}}{\to}
\begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & & \star & \star \\ & & & \star \\ & & & \star \end{bmatrix}
\overset{H_{L,4}}{\to}
\begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & & \star & \star \\ & & & \star \\ & & & \end{bmatrix}
= B,$$

- Since the transformations are all orthogonal multiplications, singular values are preserved $\sigma_i(A) = \sigma_i(B)$.

- Once bidiagonalised, one can complete the SVD as follows:

    - Mathematically, do QR alg on $B^T B$ (symmetric tridiagonal)
    - More elegant: divide-and-conquer [Gu-Eisenstat 1995] or dqds algorithm [Fernando-Parlett 1994] (nonexaminable)

- Cost: $\approx 4mn^2$ flops for singular values $\Sigma$, $\approx 20mn^2$ flops to also compute the singular vectors $U, V$.

## 10.3 (Optional but important) QZ algorithm for generalised eigenvalue problems

An increasingly important class of eigenvalue problems is the so-called generalised eigenvalue problems involving two matrices. You have probably not seen them before, but the number of applications that boil down to a generalised eigenvalue problem has been increasing rapidly.

A generalised eigenvalue problem is of the form

$$Ax = \lambda Bx, \qquad A, B \in \mathbb{C}^{n \times n}$$

- The matrices $A, B$ are given. The goal is to find the eigenvalues $\lambda$ and their corresponding eigenvectors $x$.

- There are usually (incl. when $B$ is nonsingular) $n$ eigenvalues, which are the roots of $\det(A - \lambda B)$.

- When $B$ is invertible, one can reduce the problem to $B^{-1}Ax = \lambda x$.

- Important case: $A, B$ symmetric, $B$ positive definite: in this case $\lambda$ are all real.

QZ algorithm: look for unitary $Q, Z$ s.t. $QAZ, QBZ$ both upper triangular.

- Then $\text{diag}(QAZ)/\text{diag}(QBZ)$ are the eigenvalues.

- Algorithm: first reduce $A, B$ to Hessenberg-triangular form.

- Then implicitly do QR to $B^{-1}A$ (without inverting $B$).

- Cost: $\approx 50n^3$.

- See [14] for details.

## 10.4   (Optional) Tractable eigenvalue problems

Beyond generalised eigenvalue problems there are more exotic generalisations of eigenvalue problems and reliable algorithms have been proposed for solving them.

Thanks to the remarkable developments in NLA, the following problems are 'tractable' in that reliable algorithms exist for solving them, at least when the matrix size $n$ is modest (say in the thousands).

- Standard eigenvalue problems $\boxed{Ax = \lambda x}$

    - symmetric ($4/3n^3$ flops for eigvals, $+9n^3$ for eigvecs)
    - nonsymmetric ($10n^3$ flops for eigvals, $+15n^3$ for eigvecs)

- SVD $\boxed{A = U\Sigma V^T}$ for $A \in \mathbb{R}^{m \times n}$: ($\frac{8}{3}mn^2$ flops for singvals, $+20mn^2$ for singvecs)

- Generalised eigenvalue problems $\boxed{Ax = \lambda Bx}$, $A, B \in \mathbb{C}^{n \times n}$

- Polynomial eigenvalue problems, e.g. (degree $k = 2$) $P(\lambda)x = \boxed{(\lambda^2 A + \lambda B + C)x = 0}$, $A, B, C \in \mathbb{C}^{n \times n}{:}\approx 20(nk)^3$

- Nonlinear problems, e.g. $N(\lambda)x = (A\exp(\lambda) + B)x = 0$

    - often solved via approximating by polynomial $N(\lambda) \approx P(\lambda)$
    - more difficult: $A(x)x = \lambda x$: eigenvector nonlinearity

Further speedup is often possible when structure present (e.g. sparse, low-rank)

# 11 Iterative methods: introduction

This section marks a point of departure from previous sections. So far we've been discussing direct methods. Namely, we've covered the LU-based linear system solution and the QR algorithm for eigenvalue problems[20].

Direct methods are

- Incredibly reliable, and backward stable

- Works like magic if $n \lesssim 10000$

- But not if $n$ is larger!

A 'big' matrix problem is one for which direct methods aren't feasible. Historically, as computers become increasingly faster, roughly

- 1950: $n \geq 20$

- 1965: $n \geq 200$

- 1980: $n \geq 2000$

- 1995: $n \geq 20000$

- 2010: $n \geq 100000$

- 2020: $n \geq 500000$ ($n \geq 50000$ on a standard desktop)

was considered 'too large for direct methods'. While it's clearly good news that our ability to solve problems with direct methods has been improving, the scale of problems we face in data science has been growing at the same (or even faster) pace! For such problems, we need to turn to alternative algorithms: we'll cover **iterative** and **randomized** methods. We first discuss iterative methods, with a focus on *Krylov subspace methods.*

**Direct vs. iterative methods**   Broadly speaking, the idea of iterative methods is to:

- Gradually refine the solution iteratively.

- Each iteration should be (a lot) cheaper than direct methods, usually $O(n^2)$ or less.

- Iterative methods can be (but not always) much faster than direct methods.

- Tends to be (slightly) less robust, nontrivial/problem-dependent analysis.

- Often, after $O(n^3)$ work it still gets the exact solution (ignoring roundoff errors). But one would hope to get a (acceptably) good solution long before that!
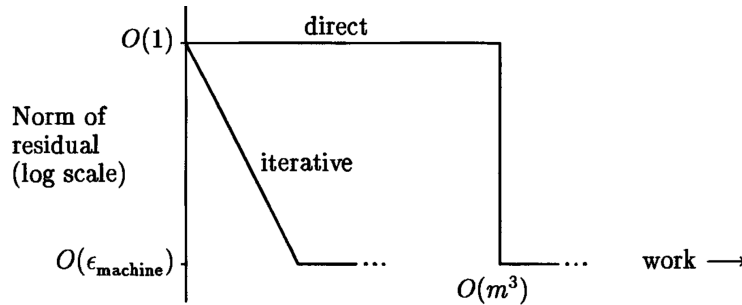
Figure 2: Rough comparison of direct vs. iterative methods, image from [Trefethen-Bau [33]]

Each iteration of most iterative methods is based on multiplying $A$ to a vector; clearly cheaper than an $O(n^3)$ direct algorithm. We'll focus on **Krylov subspace methods**. (Other iterative methods we won't get into details include the Gauss-Seidel, SOR and Chebyshev semi-iterative methods.)

## 11.1 Polynomial approximation: basic idea of Krylov

The big idea behind Krylov subspace methods is to approximate the solution in terms of a polynomial of the matrix times a vector. Namely, in Krylov subspace methods, we look for an (approximate) solution $\hat{x}$ (for $Ax = b$ or $Ax = \lambda x$) of the form (after $k$th iteration)

$$\hat{x} = p_{k-1}(A)v ,$$

where $p_{k-1}$ is a polynomial of degree (at most) $k-1$, that is, $p_{k-1}(z) = \sum_{i=0}^{k-1} c_i z^i$ for some coefficients $c_i \in \mathbb{C}$, and $v \in \mathbb{R}^n$ is the *initial vector* (usually $v = b$ for linear systems, for eigenproblems $v$ is usually a random vector). That is, $p_{k-1}(A) = \sum_{i=0}^{k-1} c_i A^i$.

Natural questions:

- Why would this be a good idea?

    - Clearly, 'easy' to compute.
    - One example: recall power method $\hat{x} = A^{k-1}v = p_{k-1}(A)v$
      Krylov finds a "better/optimal" polynomial $p_{k-1}(A)$. When the eigenvalues of $A$ are well-behaved, we'll see that $O(1)$ iterations suffices for convergence.
    - We'll see more cases where Krylov is powerful.

- How to turn this idea into an algorithm?

    - Find an orthonormal basis: Arnoldi (next), Lanczos.

---

[20]Note that the QR algorithm is iterative so it is not exactly correct to call it direct; however the nature of its consistency and robustness together with the fact that the cost is more less predictable has rightfully earned it the classification as a direct method.

## 11.2  Orthonormal basis for $\mathcal{K}_k(A, b)$

Goal: Find approximate solution $\hat{x} = p_{k-1}(A)b$, i.e. in Krylov subspace

$$\mathcal{K}_k(A, b) := \text{span}([b, Ab, A^2b, \ldots, A^{k-1}b])$$

You would want to convince yourself that any vector in the Krylov subspace can be written as a polynomial of $A$ times the vector $b$.

   An important and non-trivial step towards finding a good solution is to form an orthonormal basis for the Krylov subspace.

   First step: form an orthonormal basis $Q$, s.t. solution $x \in \mathcal{K}_k(A, b)$ can be written as $x = Qy$

- Naive idea: Form matrix $[b, Ab, A^2b, \ldots, A^{k-1}b]$, then compute its QR factorisation.

    - $[b, Ab, A^2b, \ldots, A^{k-1}b]$ is usually terribly conditioned! Dominated by leading eigvec
    - $Q$ is therefore extremely ill-conditioned, and hence inaccurately computed

- Much better solution: Arnoldi iteration

    - Multiply $A$ once at a time to the latest orthonormal vector $q_i$
    - Then orthogonalise $Aq_i$ against previous $q_j$'s ($j = 1, \ldots, i - 1$) (as in Gram-Schmidt).

## 11.3  Arnoldi iteration

Here is a pseudocode of the Arnold iteration. Essentially, what it does is multiply the matrix $A$, orthogonalise against the previous vectors, and repeat.

---

**Algorithm 11.1** The Arnoldi iteration for finding an orthonormal basis for Krylov subspace $\mathcal{K}_k(A, b)$.

---

1: Set $q_1 = b/\|b\|_2$
2: For $k = 1, 2, \ldots,$
3:   set $v = Aq_k$
4:     for $j = 1, 2, \ldots, k$
5:       $h_{jk} = q_j^T v$, $v = v - h_{jk}q_j$ % orthogonalise against $q_j$ via (modified) Gram-Schmidt

6:     end for
7:   $h_{k+1,k} = \|v\|_2$, $q_{k+1} = v/h_{k+1,k}$
8: End for

---

- After $k$ steps, $AQ_k = Q_{k+1}\tilde{H}_k = Q_k H_k + q_{k+1}[0, \ldots, 0, h_{k+1,k}]$, with $Q_k = [q_1, q_2, \ldots, q_k], Q_{k+1} = [Q_k, q_{k+1}]$, $\operatorname{span}(Q_k) = \operatorname{span}([b, Ab, \ldots, A^{k-1}b])$

$$\boxed{A}\,\boxed{Q_k} = \boxed{Q_{k+1}}\,\boxed{\tilde{H}_k}, \quad \tilde{H}_k = \underbrace{\begin{bmatrix} h_{1,1} & h_{1,2} & \ldots & & h_{1,k} \\ h_{2,1} & h_{2,2} & \ldots & & h_{2,k} \\ & \ddots & & & \vdots \\ & & h_{k,k-1} & & h_{k,k} \\ & & & & h_{k+1,k} \end{bmatrix}}_{\mathbb{R}^{(k+1)\times k} \text{ upper Hessenberg}}, \quad Q_{k+1}^T Q_{k+1} = I_{k+1}$$

- Cost is $k$ $A$-multiplications $+ O(k^2)$ inner products ($O(nk^2)$ flops)

## 11.4 Lanczos iteration

When $A$ is symmetric, Arnoldi simplifies to

$$AQ_k = Q_k T_k + q_{k+1}[0, \ldots, 0, h_{k+1,k}],$$

where $T_k$ is symmetric tridiagonal (proof: just note $H_k = Q_k^T A Q_k$ in Arnoldi)

$$\boxed{A}\,\boxed{Q_k} = \boxed{Q_{k+1}}\,\boxed{\tilde{T}_k}, \quad \tilde{T}_k = \underbrace{\begin{bmatrix} t_{1,1} & t_{1,2} & & \\ t_{2,1} & t_{2,2} & \ddots & \\ & \ddots & \ddots & t_{k-1,k} \\ & & t_{k,k-1} & t_{k,k} \\ & & & t_{k+1,k} \end{bmatrix}}_{\mathbb{R}^{(k+1)\times k} \text{ symmetric tridiagonal}}, \quad Q_{k+1}^T Q_{k+1} = I_{k+1}$$

- The vectors $q_k$ form a 3-term recurrence $t_{k+1,k}q_{k+1} = (A - t_{k,k})q_k - t_{k-1,k}q_{k-1}$. Orthogonalisation is necessary only against last two vectors $q_k, q_{k-1}$

- Significant speedup over Arnoldi; cost is $k$ $A$-multiplication plus $O(k)$ inner products ($O(nk)$).

- In floating-point arithmetic, sometimes the computed $Q_k$ loses orthogonality and re-orthogonalisation may be necessary (nonexaminable, see e.g. Demmel [8])

## 11.5 The Lanczos algorithm for symmetric eigenproblem

We are now ready to describe one of the most successful algorithms for large-scale symmetric eigenvalue problems: the Lanczos algorithm. In simple words, it finds an eigenvalue and eigenvector in a Krylov subspace by a projection method, called the Rayleigh-Ritz process.

---

**Algorithm 11.2 Rayleigh-Ritz**: given symmetric $A$ and orthonormal $Q$, find approximate eigenpairs

1: Compute $Q^T A Q$.
2: Eigenvalue decomposition $Q^T A Q = V \hat{\Lambda} V^T$.
3: Approximate eigenvalues $\operatorname{diag}(\hat{\Lambda})$ (Ritz values) and eigenvectors $QV$ (Ritz vectors).

---

This is a **projection** method (similar alg is available for SVD).

Now we can describe the Lanczos algorithm as follows:
Lanczos algorithm=Lanczos iteration+Rayleigh-Ritz

---

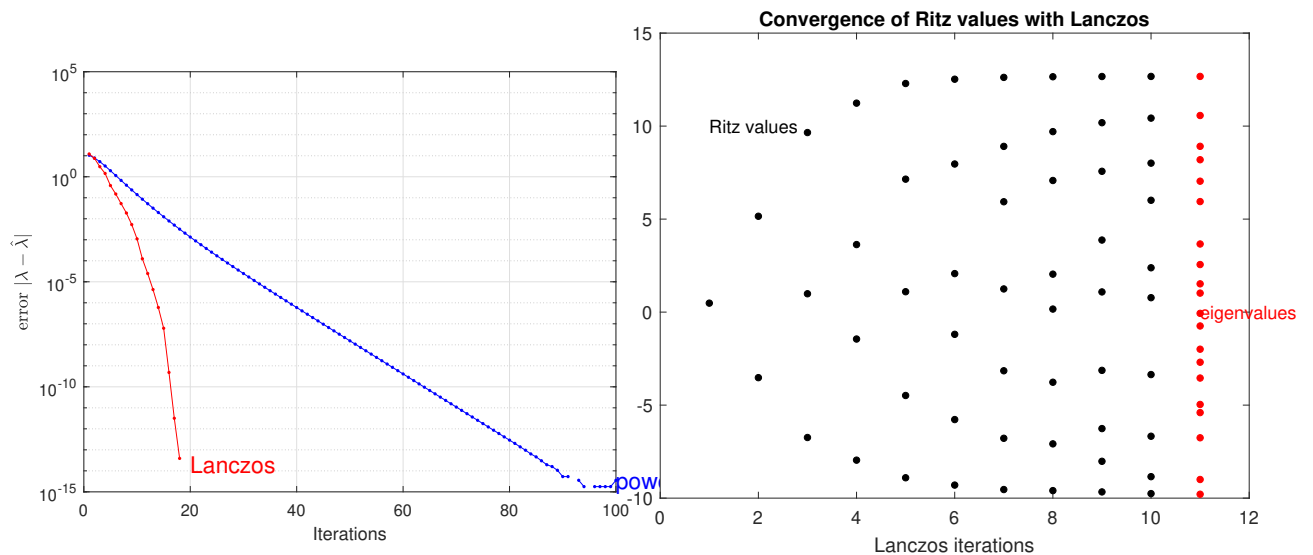**Algorithm 11.3 Lanczos** algorithm: $A \in \mathbb{R}^{n \times n}$ find (extremal) eigenpairs.

---
1: Perform Lanczos iterations to obtain $AQ_k = Q_k T_k + q_{k+1}[0, \ldots, 0, h_{k+1,k}]$.
2: Compute the eigenvalue decomposition of $T_k = V_k \hat{\Lambda} V_k^T$ (Rayleigh-Ritz with subspace $Q_k$).
3: diag($\hat{\Lambda}$) are the approximate eigenvalues (Ritz values), and the columns of $Q_k V_k$ are the approximate eigenvectors (Ritz vectors).

---

- In this case $Q = Q_k$, so simply $Q_k^T A Q_k = T_k$ (tridiagonal eigenproblem)

- Very good convergence is usually observed to extremal eigenpairs. To see this:

    - Recall from Courant-Fischer $\lambda_{\max}(A) = \max_x \frac{x^T A x}{x^T x}$

    - Hence $\lambda_{\max}(A) \geq \underbrace{\max_{x \in \mathcal{K}_k(A,b)} \frac{x^T A x}{x^T x}}_{\text{Lanczos output}} \geq \underbrace{\frac{v^T A v}{v^T v}, \quad v = A^{k-1}b}_{\text{power method}}$

    - Same for $\lambda_{\min}$, similar for e.g. $\lambda_2$

This is admittedly a very crude estimate for the convergence of Lanczos. A thorough analysis turns out to be very complicated; if interested see for example [26].

**Experiments with Lanczos** Symmetric $A \in \mathbb{R}^{n \times n}, n = 100$, Lanczos/power method with random initial vector $b$



Convergence to dominant eigenvalue

Convergence of Ritz values (approximate eigenvalues)

The same principles of projecting the matrix onto a Krylov subspace applies to nonsymmetric eigenvalue problems. Essentially it boils down to finding the eigenvalues of the upper Hessenberg matrix $H$ arising in the Arnoldi iteration, rather than the tridiagonal matrix as in Lanczos.

# 12    Arnoldi and GMRES for $Ax = b$

This is an exciting section as we will describe the GMRES algorithm. This algorithm has been so successful that in the 90s the paper that introduced it was the most cited paper in all of applied mathematics.

Just as the Lanczos method for eigenvalue problems is a simple projection with the matrix onto the Krylov subspace, GMRES attempts to find an approximate solution in the Krylov subspace that is in some sense the best possible.

Idea (very simple!): minimise residual in Krylov subspace:          [Saad-Schulz 86 [29]]

$$x = \text{argmin}_{x \in \mathcal{K}_k(A,b)} \|Ax - b\|_2.$$

In order to solve this, The algorithm cleverly takes advantage of the structure afforded by Arnoldi.

Algorithm: Given $AQ_k = Q_{k+1}\tilde{H}_k$ and writing $x = Q_k y$, rewrite as

$$\min_y \|AQ_k y - b\|_2 = \min_y \|Q_{k+1}\tilde{H}_k y - b\|_2$$

$$= \min_y \left\| \begin{bmatrix} \tilde{H}_k \\ 0 \end{bmatrix} y - \begin{bmatrix} Q_{k+1}^T \\ Q_{k+1,\perp}^T \end{bmatrix} b \right\|_2$$

$$= \min_y \left\| \begin{bmatrix} \tilde{H}_k \\ 0 \end{bmatrix} y - \|b\|_2 e_1 \right\|_2, \quad e_1 = [1, 0, \dots, 0]^T \in \mathbb{R}^n$$

( where $[Q_{k+1}, Q_{k+1,\perp}]$ is orthogonal; we're using the same trick as in least-squares)

- Minimised when $\|\tilde{H}_k y - Q_{k+1}^T b\|_2$ is minimised; this is a Hessenberg least-squares problem.

- Solve via the QR-based approach described in Section 6.5. Here it's even simpler as the matrix is Hessenberg: $k$ Givens rotations yields the QR factorisation, then a triangular solve will complete the solution, so $O(k^2)$ work in addition to Arnoldi; recall Section 6.4.

## 12.1    GMRES convergence: polynomial approximation

We now study the convergence of GMRES. As with any Krylov subspace method, the analysis is based on the theory of polynomial approximation.

**Theorem 12.1 (GMRES convergence)** *Assume that $A$ is diagonalisable, $A = X\Lambda X^{-1}$. Then the kth GMRES iterate $x_k$ satisfies*

$$\|Ax_k - b\|_2 \leq \kappa_2(X) \min_{p \in \mathcal{P}_k, p(0)=1} \max_{z \in \lambda(A)} |p(z)| \|b\|_2.$$

62

**Proof:** Recall that $x_k \in \mathcal{K}_k(A, b) \Rightarrow x_k = p_{k-1}(A)b$, where $p_{k-1}$ is a polynomial of degree at most $k - 1$. Hence GMRES solution is

$$\min_{x_k \in \mathcal{K}_k(A,b)} \|Ax_k - b\|_2 = \min_{p_{k-1} \in \mathcal{P}_{k-1}} \|Ap_{k-1}(A)b - b\|_2$$

$$= \min_{\tilde{p} \in \mathcal{P}_k, \tilde{p}(0)=0} \|(\tilde{p}(A) - I)b\|_2$$

$$= \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)b\|_2$$

If $A$ is diagonalizable, $A = X\Lambda X^{-1}$,

$$\|p(A)\|_2 = \|Xp(\Lambda)X^{-1}\|_2 \leq \|X\|_2 \|X^{-1}\|_2 \|p(\Lambda)\|_2$$
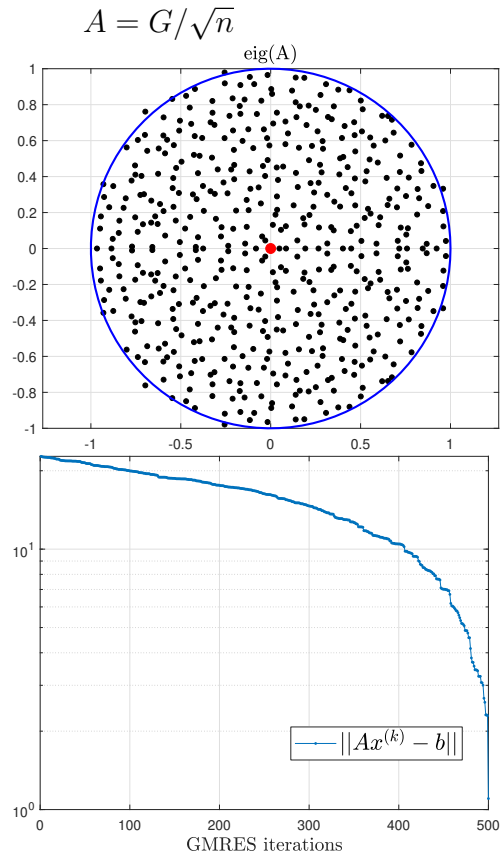$$= \kappa_2(X) \max_{z \in \lambda(A)} |p(z)|$$

(recall that $\lambda(A)$ is the set of eigenvalues of $A$)  $\square$

Noting that $\|p(A)b\|_2 \leq \|p(A)\|_2 \|b\|_2$, here is the interpretation of this analysis: We would like to find a polynomial s.t. $p(0) = 1$ and $|p|$ is small at the eigenvalues of $A$, that is, $|p(\lambda_i)|$ is small for all $i$.

The question becomes, what type of eigenvalues are 'nice' for this to hold?

**GMRES example**  $G$: Gaussian random matrix ($G_{ij} \sim N(0, 1)$, i.i.d.) $G/\sqrt{n}$, $n = 1000$: eigvals in unit disk.

$A = 2I + G/\sqrt{n}$,
$p(z) = 2^{-k}(z-2)^k$

$A = G/\sqrt{n}$



63

**Initial vector.** Sometimes a good initial guess $x_0$ for $x_*$ is available. In this case we take the initial residual $r_0 = Ax_0 - b$, and work in the affine space $x_k = x_0 + \mathcal{K}_k(A, r_0)$. All the analysis above can be modified readily to allow for this situation with essentially the same conclusion (clearly if one has a good $x_0$ by all means that should be used).

## 12.2 When does GMRES converge fast?

Recall GMRES solution satisfies (assuming $A$ is diagonalisable+nonsingular)

$$\min_{x \in \mathcal{K}_k(A,b)} \|Ax - b\|_2 = \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)b\|_2 \leq \kappa_2(X) \max_{z \in \lambda(A)} |p(z)| \|b\|_2.$$

$\max_{z \in \lambda(A)} |p(z)|$ is small when

- $\lambda(A)$ are clustered away from 0
    - a good $p$ can be found quite easily
    - e.g. example above


- When $\lambda(A)$ takes $k (\ll n)$ distinct values
    - Then convergence in $k$ GMRES iterations (why?)

## 12.3 Preconditioning for GMRES

We've seen that GMRES is great if the eigenvalues are clustered away from 0. If this is not true, GMRES can really require close to (or equal to!) the full $n$ iterations for convergence. This is undesirable—at this point we've spent more computation than a direct method! We need a workaround.

The idea of preconditioning is to instead of solving

$$Ax = b,$$

find a "preconditioner" $M \in \mathbb{R}^{n \times n}$ and solve

$$MAx = Mb$$

Of course, $M$ cannot be an arbitrary matrix. It has to be chosen very carefully. Desiderata of $M$ are

- $M$ is simple enough s.t. applying $M$ to vector is easy (note that each GMRES iteration requires $MA$-multiplication), and one of

    1. $MA$ has clustered eigenvalues away from 0.
    2. $MA$ has a small number of distinct nonzero eigenvalues.
    3. $MA$ is well-conditioned $\kappa_2(MA) = O(1)$; then solve the normal equation $(MA)^T MAx = (MA)^T Mb$.

**Preconditioners: examples**

- ILU (Incomplete LU) preconditioner: $A \approx LU, M = (LU)^{-1} = U^{-1}L^{-1}$, $L, U$ 'as sparse as $A$' $\Rightarrow MA \approx I$ (hopefully; 'cluster away from 0').

- For $\tilde{A} = \begin{bmatrix} A & B \\ C & 0 \end{bmatrix}$, set $M = \begin{bmatrix} A^{-1} & \\ & (CA^{-1}B)^{-1} \end{bmatrix}$. Then if $M$ nonsingular, $M\tilde{A}$ has eigvals$\in \{1, \frac{1}{2}(1 \pm \sqrt{5})\} \Rightarrow$ 3-step convergence.  $\qquad$ [Murphy-Golub-Wathen 2000 [23]]

- Multigrid-based, operator preconditioning, ...

- A "perfect" preconditioner is $M = A^{-1}$; as then preconditioned GMRES will converge in one step. Obviously this $M$ isn't easy to apply (if it was then we're done!). Preconditioning, therefore, can be regarded as an act of efficiently approximating the inverse.


$\qquad$ Finding effective preconditioners is a never-ending research topic.
Prof. Andy Wathen is our Oxford expert!

## 12.4 Restarted GMRES

Another practical GMRES technique is restarting. For $k$ iterations, GMRES costs $k$ matrix multiplications$+O(nk^2)$ for orthogonalisation $\rightarrow$ Arnoldi eventually becomes expensive.

$\qquad$ Practical solution: restart by solving 'iterative refinement':

1. Stop GMRES after $k_{\max}$ (prescribed) steps to get approx. solution $\hat{x}_1$.

2. Solve $A\tilde{x} = b - A\hat{x}_1$ via GMRES. (This is a linear system with a different right-hand side).

3. Obtain solution $\hat{x}_1 + \tilde{x}$.

Sometimes multiple restarts are needed.

## 12.5 Arnoldi for nonsymmetric eigenvalue problems

Arnoldi for eigenvalue problems: Arnoldi iteration+Rayleigh-Ritz (just like Lanczos alg)

1. Compute $Q^T A Q$.

2. Eigenvalue decomposition $Q^T A Q = X\hat{\Lambda}X^{-1}$.

3. Approximate eigenvalues diag($\hat{\Lambda}$). (Ritz values) and eigenvectors $QX$ (Ritz vectors).

As in Lanczos, $Q = Q_k = \mathcal{K}_k(A, b)$, so simply $Q_k^T A Q_k = H_k$ (Hessenberg eigenproblem, note that this is ideal for the QR algorithm as the preprocessing step can be skipped).

Which eigenvalues are found by Arnoldi? We give a rather qualitative answer:

- First note that Krylov subspace is invariant under shift: $\mathcal{K}_k(A, b) = \mathcal{K}_k(A - sI, b)$.

- Thus any eigenvector that power method applied to $A - sI$ converges to should be contained in $\mathcal{K}_k(A, b)$.

- To find other (e.g. interior) eigvals, one can use shift-invert Arnoldi: $Q = \mathcal{K}_k((A - sI)^{-1}, b)$.

# 13 Lanczos and Conjugate Gradient method for $Ax = b$, $A \succ 0$

Here we introduce the conjugate gradient (CG) method. CG is not only important in practice but has historical significance in that it was the very first Krylov algorithm to be introduced (and initially made a big hype, but then took a while to be recognized as a competitive method).

First recall that when $A$ is symmetric, Lanczos gives $Q_k, T_k$ such that $AQ_k = Q_k T_k + q_{k+1}[0, \ldots, 0, 1]$, $T_k$: tridiagonal.

The idea of CG is as follows: when $A \succ 0$ PD, solve $Q_k^T (AQ_k y - b) = T_k y - Q_k^T b = 0$, and $x = Q_k y$.

This is known as "Galerkin orthogonality": it imposes that the residual $Ax - b$ is orthogonal to $Q_k$.

- $T_k y = Q_k^T b$ is a tridiagonal linear system, so it requires only $O(k)$ operations to solve.

- Three-term recurrence reduces cost to $O(k)$ $A$-multiplications, making the orthogonalisation cost almost negligible.

- The CG algorithm minimises $A$-norm of error in the Krylov subspace $x_k = \operatorname{argmin}_{x \in Q_k} \|x - x_*\|_A$ ($Ax_* = b$): writing $x_k = Q_k y$, we have

$$(x_k - x_*)^T A(x_k - x_*) = (Q_k y - x_*)^T A(Q_k y - x_*)$$
$$= y^T (Q_k^T A Q_k) y - 2b^T Q_k y + b^T x_*,$$

minimiser is $y = (Q_k^T A Q_k)^{-1} Q_k^T b$, so $Q_k^T (AQ_k y - b) = 0$.

  - Note $\|x\|_A = \sqrt{x^T A x}$ defines a norm (exercise)
  - More generally, for inner-product norm $\|z\|_M = \sqrt{\langle z, z \rangle_M}$, $\min_{x = Qy} \|x_* - x\|_M$ attained when $\langle q_i, x_* - x \rangle_M = 0$, $\forall q_i$ (cf. Part A Numerical Analysis).

## 13.1 CG algorithm for $Ax = b$, $A \succ 0$

We've described the CG algorithm conceptually. To derive the practical algorithm some clever manipulations are necessary. We won't go over them in detail but here is the outcome:

Set $x_0 = 0$, $r_0 = -b$, $p_0 = r_0$ and do for $k = 1, 2, 3, \ldots$

$$\alpha_k = \langle r_k, r_k \rangle / \langle p_k, Ap_k \rangle$$
$$x_{k+1} = x_k + \alpha_k p_k$$
$$r_{k+1} = r_k - \alpha_k Ap_k$$
$$\beta_k = \langle r_{k+1}, r_{k+1} \rangle / \langle r_k, r_k \rangle$$
$$p_{k+1} = r_{k+1} + \beta_k p_k$$

where $r_k = b - Ax_k$ (residual) and $p_k$ (search direction). $x_k$ is the CG solution after $k$ iterations.

One can show among others (exercise/sheet)

- $\mathcal{K}_k(A, b) = \text{span}(r_0, r_1, \ldots, r_{k-1}) = \text{span}(x_1, x_2, \ldots, x_k)$ (also equal to $\text{span}(p_0, p_1, \ldots, p_{k-1})$)

- $r_j^T r_k = 0$, $j = 0, 1, 2, \ldots, k-1$

Thus $x_k$ is $k$th CG solution, satisfying Galerkin orthogonality $Q_k^T(Ax_k - b) = 0$: residual is orthogonal to the (Krylov) subspace.

## 13.2 CG convergence

Let's examine the convergence of the CG iterates.

**Theorem 13.1** *Let $A \succ 0$ be an $n \times n$ positive definite matrix and $b \in \mathbb{R}^n$. Let $e_k := x_* - x_k$ be the error after the $k$th CG iteration ($x_*$ is the exact solution $Ax_* = b$). Then*

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq 2 \left( \frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k.$$

**Proof:** We have $e_0 = x_*$ ($x_0 = 0$), and

$$\frac{\|e_k\|_A}{\|e_0\|_A} = \min_{x \in \mathcal{K}_k(A,b)} \|x_k - x_*\|_A / \|x_*\|_A$$
$$= \min_{p_{k-1} \in \mathcal{P}_{k-1}} \|p_{k-1}(A)b - A^{-1}b\|_A / \|e_0\|_A$$
$$= \min_{p_{k-1} \in \mathcal{P}_{k-1}} \|(p_{k-1}(A)A - I)e_0\|_A / \|e_0\|_A$$
$$= \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)e_0\|_A / \|e_0\|_A$$
$$= \min_{p \in \mathcal{P}_k, p(0)=1} \left\| V \begin{bmatrix} p(\lambda_1) & & \\ & \ddots & \\ & & p(\lambda_n) \end{bmatrix} V^T e_0 \right\|_A / \|e_0\|_A.$$

67

Now $(\text{blue})^2 = \sum_i \lambda_i p(\lambda_i)^2 (V^T e_0)_i^2 \leq \max_j p(\lambda_j)^2 \sum_i \lambda_i (V^T e_0)_i^2 = \max_j p(\lambda_j)^2 \|e_0\|_A^2$.

We've shown

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq \min_{p \in \mathcal{P}_k, p(0)=1} \max_j |p(\lambda_j)| \leq \min_{p \in \mathcal{P}_k, p(0)=1} \max_{x \in [\lambda_{\min}(A), \lambda_{\max}(A)]} |p(x)|$$

To complete the proof, in the next subsection we will show that

$$\min_{p \in \mathcal{P}_k, p(0)=1} \max_{x \in [\lambda_{\min}(A), \lambda_{\max}(A)]} |p(x)| \leq 2 \left( \frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k. \tag{10}$$
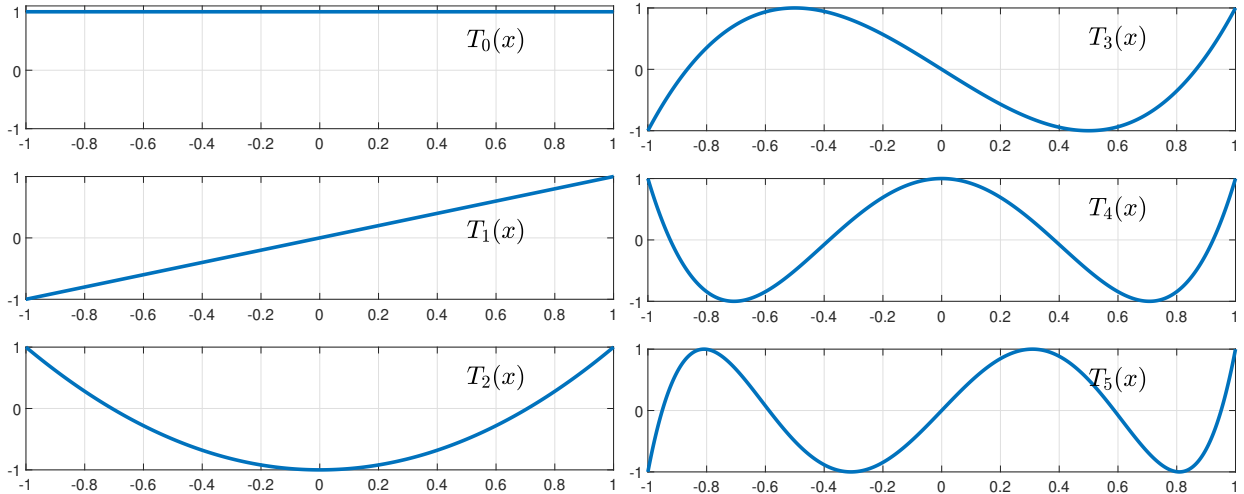
$\square$

- Note that $\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} (=: \frac{b}{a})$.

- The above bound is obtained using Chebyshev polynomials on $[\lambda_{\min}(A), \lambda_{\max}(A)]$. This is a class of polynomials that arise in a varieties of contexts in computational maths. Let's next look at their properties.

### 13.2.1 Chebyshev polynomials

For $z = \exp(i\theta)$, $x = \frac{1}{2}(z + z^{-1}) = \cos\theta \in [-1, 1]$, $\theta = \text{acos}(x)$, $T_k(x) = \frac{1}{2}(z^k + z^{-k}) = \cos(k\theta)$. $T_k(x)$ is a polynomial in $x$:

$$\frac{1}{2}(z + z^{-1})(z^k + z^{-k}) = \frac{1}{2}(z^{k+1} + z^{-(k+1)}) + \frac{1}{2}(z^{k-1} + z^{-(k-1)}) \Leftrightarrow \underbrace{2x T_k(x) = T_{k+1}(x) + T_{k-1}(x)}_{\text{3-term recurrence}}$$



These polynomials grow very fast outside the interval (here the 'standard' $[-1, 1]$). For example, plots on $[-2, 1]$ look like

Here's a nice plot of several Chebyshev polynomials:



### 13.2.2  Properties of Chebyshev polynomials
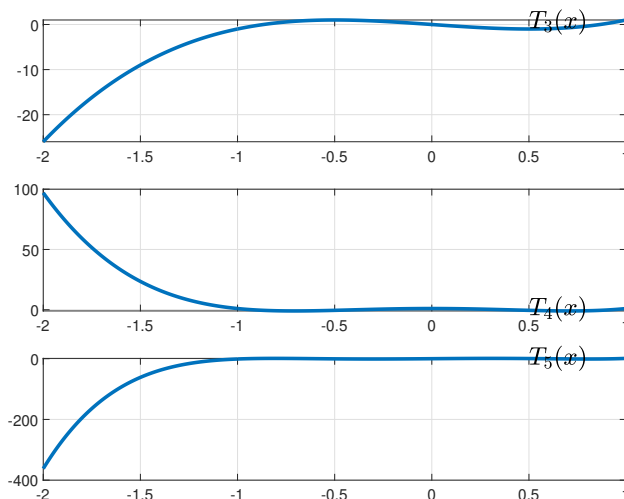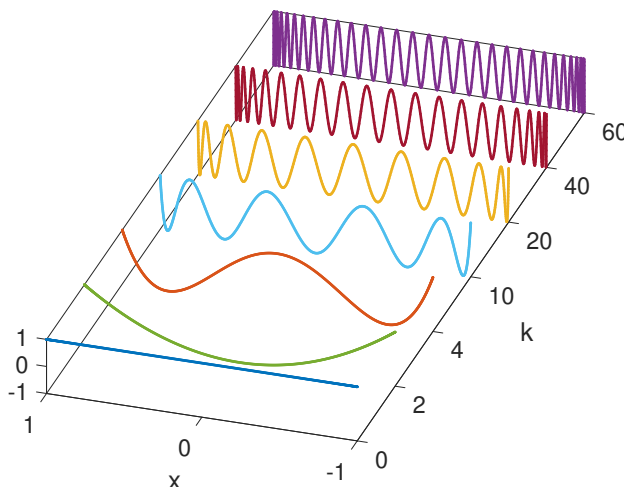
For $z = \exp(i\theta)$, $x = \frac{1}{2}(z+z^{-1}) = \cos\theta \in [-1,1]$, $\theta = \mathrm{acos}(x)$, $T_k(x) = \frac{1}{2}(z^k + z^{-k}) = \cos(k\theta)$.

- Inside $[-1,1]$, $|T_k(x)| \le 1$

- Outside $[-1,1]$, $|T_k(x)| \gg 1$ grows rapidly with $|x|$ and $k$(fastest growth among $\mathcal{P}_k$)

Shift+scale s.t. $p(x) = c_k T_k(\frac{2x-b-a}{b-a})$ where $c_k = 1/T_k(\frac{-(b+a)}{b-a})$ so $p(0) = 1$. Then

- $|p(x)| \le 1/|T_k(\frac{b+a}{b-a})|$ on $x \in [a,b]$

69

- $T_k(z) = \frac{1}{2}(z^k + z^{-k})$ with $\frac{1}{2}(z + z^{-1}) = \frac{b+a}{b-a} \Rightarrow z = \frac{\sqrt{b/a}+1}{\sqrt{b/a}-1} = \frac{\sqrt{\kappa_2(A)}+1}{\sqrt{\kappa_2(A)}-1}$, so

$$|p(x)| \leq 1/T_k(\frac{b+a}{b-a}) \leq 2\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k.$$

This establishes (10), completing the proof of Theorem 13.1. $\qquad\square$

For much more about $T_k$, see C6.3 Approximation of Functions

## 13.3   MINRES: symmetric (indefinite) version of GMRES

When the matrix is symmetric but not positive definite, GMRES can still be simplified although some care is needed as the Matrix $A$ ceases to define a norm unlike for positive definite matrices.

Symmetric analogue of GMRES: MINRES (minimum-residual method) for $A = A^T$ (but not necessarily $A \succ 0$)

$$x = \mathrm{argmin}_{x \in \mathcal{K}_k(A,b)} \|Ax - b\|_2.$$

Algorithm: Given $AQ_k = Q_{k+1}\tilde{T}_k$ and writing $x = Q_k y$, rewrite as

$$\min_y \|AQ_k y - b\|_2 = \min_y \|Q_{k+1}\tilde{T}_k y - b\|_2$$

$$= \min_y \left\| \begin{bmatrix} \tilde{T}_k \\ 0 \end{bmatrix} y - \begin{bmatrix} Q_k^T \\ Q_{k,\perp}^T \end{bmatrix} b \right\|_2$$

$$= \min_y \left\| \begin{bmatrix} \tilde{T}_k \\ 0 \end{bmatrix} y - \|b\|_2 e_1 \right\|_2, \quad e_1 = [1,0,\ldots,0]^T \in \mathbb{R}^n$$

( where $[Q_k, Q_{k,\perp}]$ orthogonal; same trick as in least-squares)

- Minimised when $\|\tilde{T}_k y - \tilde{Q}_k^T b\| \to \min$; tridiagonal least-squares problem

- Solve via QR ($k$ Givens rotations)+ tridiagonal solve, $O(k)$ in addition to Lanczos

### 13.3.1   MINRES convergence

As in GMRES, we can examine the MINRES residual in terms of minimising the values of a polynomial at the eigenvalues of $A$.

$$\min_{x \in \mathcal{K}_k(A,b)} \|Ax - b\|_2 = \min_{p_{k-1} \in \mathcal{P}_{k-1}} \|Ap_{k-1}(A)b - b\|_2 = \min_{\tilde{p} \in \mathcal{P}_k, \tilde{p}(0)=0} \|(\tilde{p}(A) - I)b\|_2$$

$$= \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)b\|_2.$$

Since $A = A^T$, $A$ is diagonalisable $A = Q\Lambda Q^T$ with $Q$ orthogonal, so

$$\|p(A)\|_2 = \|Qp(\Lambda)Q^T\|_2 \leq \|Q\|_2 \|Q^T\|_2 \|p(\Lambda)\|_2$$

$$= \max_{z \in \lambda(A)} |p(z)|.$$

70

Interpretation: (again) find polynomial s.t. $p(0) = 1$ and $|p(\lambda_i)|$ small
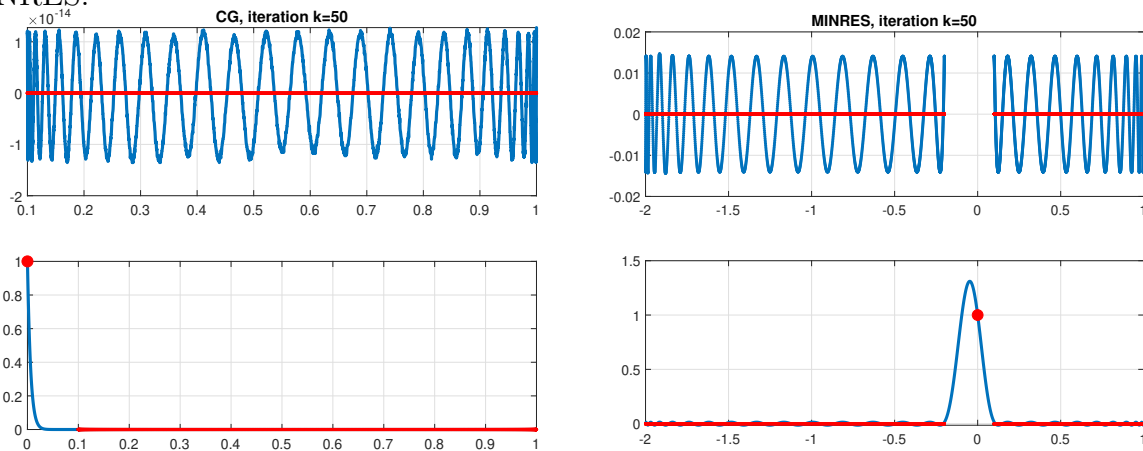
$$\frac{\|Ax - b\|_2}{\|b\|_2} \leq \min_{p \in \mathcal{P}_k, p(0)=1} \max |p(\lambda_i)|.$$

One can prove (nonexaminable)

$$\min_{p \in \mathcal{P}_k, p(0)=1} \max |p(\lambda_i)| \leq 2 \left( \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^{k/2}.$$

- obtained by Chebyshev+Möbius change of variables [Greenbaum's book 97]

- minimisation needed on positive **and** negative sides, hence slower convergence when $A$ is indefinite

**CG and MINRES, optimal polynomials** Here are some optimal polynomials for CG and MINRES. Note that $\kappa_2(A) = 10$ in both cases; observe how much faster CG is than MINRES.



## 13.4 Preconditioned CG/MINRES

The preceding analysis suggests that the linear system

$$Ax = b, \quad A = A^T (\succ 0)$$

may not converge rapidly with CG or MINRES if the eigenvalues are not distributed in a favorable manner (i.e., clustered away from 0).

In this case, as with GMRES, a workaround is to find a good preconditioner $M$ such that "$M^T M \approx A^{-1}$" and solve

$$M^T A M y = M^T b, \quad My = x$$

As before, desiderata of $M$:

- $M^T A M$ is easy to multiply to a vector.

- $M^T A M$ has clustered eigenvalues away from 0.

Note that reducing $\kappa_2(M^T A M)$ directly implies rapid convergence.

- It is possible to implement preconditioned CG with just $M^T M$ (no need to find $M$).

# 14 Randomized algorithms in NLA

In this final part of this lecture series we will talk about randomised algorithms. This takes us to the forefront of research in the field: A major idea in numerical linear algebra since 2005 or so has been the use of randomisation where in a matrix sketch is used in order to extract information about the matrix, for example in order to construct a low-rank approximation.

So far, all algorithms have been deterministic (always same output).

- Direct methods (LU for $Ax = b$, QR alg for $Ax = \lambda x$ or $A = U\Sigma V^T$) are

  - Incredibly reliable, backward stable.
  - Works like magic if $n \lesssim 10000$.
  - But not beyond; cubic complexity $O(n^3)$ or $O(mn^2)$.

- Iterative methods (GMRES, CG, Arnoldi, Lanczos) are

  - Very fast when it works (nice spectrum etc).
  - Otherwise, not so much; need for preconditioning.

- Randomized algorithms

  - Output differs at every run.
  - Ideally succeed with enormous probability, e.g. $1 - \exp(-cn)$.
  - Often by far the fastest&only feasible approach.
  - Not for all problems—active field of research.

Why do we need randomisation? Randomised algorithms are somehow attached to a negative connotation that the algorithm is not reliable and always produces different results. Well, this is a valid concern. We hope to remove such concerns in what follows. The reason randomisation is necessary is the sheer size of datasets that we face today, with the rise of data science. As you will see, one core idea of randomisation is that by allowing for a very small error (which can easily be in the order of machine precision) instead of getting an exact solution, one can often dramatically speed up the algorithm.

We'll cover two NLA topics where randomisation has been very successful: **low-rank approximation (randomized SVD)**, and overdetermined **least-squares problems**

**Gaussian matrices** In randomized algorithms we typically introduce a random matrix. For the analysis, Gaussian matrices $G \in \mathbb{R}^{m \times n}$ are the most convenient (not always for computation). These are matrices whose entries are drawn independently from the standard normal (Gaussian) distribution $G_{ij} \sim N(0, 1)$. We cannot do justice to random matrix theory (there is a course in Hilary term!), but here is a summary of the properties of Gaussian matrices that we'll be using.

- A useful fact about Gaussian random matrices $G$ is that its distribution is invariant under orthogonal transformations. That is, if $G$ is Gaussian, so is $QG$ and $GQ$, where $Q$ is any orthogonal matrix independent of $G$. To see this (nonexaminable): note that a sum of Gaussian (scalar) random variables is Gaussian, and by independence the variance is simply the sum of the variances. Now let $g_i$ denote the $i$th column of $G$. Then $\mathbb{E}[(Qg_i)^T(Qg_i)] = \mathbb{E}[g_i^T g_i] = I$, so each $Qg_i$ is multivariate Gaussian with the same distribution as $g_i$. Independence of $Qg_i, Qg_j$ is immediate.

- Another very useful fact is that when the matrix is rectangular $m > n$ (or $m < n$), the singular values are known to lie in the interval $[\sqrt{m} - \sqrt{n}, \sqrt{m} + \sqrt{n}]$. This is a consequence of the Marchenko-Pastur rule, which we discuss a bit more later.

## 14.1 Randomized SVD by Halko-Martinsson-Tropp

We start with what has been arguably the most successful usage of randomisation in NLA, low-rank approximation. Probably the best reference is the paper by Halko, Martinsson and Tropp [16]. This paper has been enormously successful, with over 3000 Google Scholar citations as of 2021. See also the recent survey by the same authors [22].

The algorithm itself is astonishingly simple:

---

**Algorithm 14.1 Randomised SVD** (HMT): given $A \in \mathbb{R}^{m \times n}$ and rank $r$, find a rank-$r$ approximation $\hat{A} \approx A$.

---

1: Form a random matrix $X \in \mathbb{R}^{n \times r}$, usually $r \ll n$.
2: Compute $AX$.
3: Compute the QR factorisation $AX = QR$.

4: $\boxed{\phantom{xx}A\phantom{xx}} \approx \boxed{Q}\, \boxed{\phantom{x}Q^T A\phantom{x}}\,(= (QU_0)\Sigma_0 V_0^T)$ is a rank-$r$ approximation.

---

Here, $X$ is a random matrix taking independent and identically distributed (iid) entries. A convenient choice (for the theory, not necessarily for computation) is a Gaussian matrix, with iid entries $X_{ij} \sim N(0, 1)$.

Here are some properties of the HMT algorithm:

- $O(mnr)$ cost for dense $A$.

- Near-optimal approximation guarantee: for any $\hat{r} < r$,

$$\mathbb{E}\|A - \hat{A}\|_F \leq \left(1 + \frac{r}{r - \hat{r} - 1}\right)\|A - A_{\hat{r}}\|_F,$$

where $A_{\hat{r}}$ is the rank $\hat{r}$-truncated SVD (expectation w.r.t. random matrix $X$).

This is a remarkable result; make sure to pause and think about what it says! The approximant $\hat{A}$ has error $\|A - \hat{A}\|_F$ that is within a factor $\left(1 + \frac{r}{r-\hat{r}-1}\right)$ of the optimal truncated SVD, for a slightly lower rank $\hat{r} < r$ (say, $\hat{r} = 0.9r$).

Goal: understand this, or at least why $\mathbb{E}\|A - \hat{A}\| = O(1)\|A - A_{\hat{r}}\|$.

## 14.2 Pseudoinverse and projectors

To understand why the HMT algorithm works, we need to introduce two notions: the pseudo inverse and (orthogonal and oblique) projectors.

Given $M \in \mathbb{R}^{m \times n}$ with economical SVD $M = U_r \Sigma_r V_r^T$ ($U_r \in \mathbb{R}^{m \times r}, \Sigma_r \in \mathbb{R}^{r \times r}, V_r \in \mathbb{R}^{n \times r}$ where $r = \mathrm{rank}(M)$ so that $\Sigma_r \succ 0$), the **pseudoinverse** $M^\dagger$ is

$$M^\dagger = V_r \Sigma_r^{-1} U_r^T \in \mathbb{R}^{n \times m}.$$

- $M^\dagger$ satisfies $MM^\dagger M = M$, $M^\dagger M M^\dagger = M^\dagger$, $MM^\dagger = (MM^\dagger)^T$, $M^\dagger M = (M^\dagger M)^T$ (these are often taken to be the definition of the pseudoinverse—the above definition is much simpler IMO).

- $M^\dagger = M^{-1}$ if $M$ nonsingular.

- $M^\dagger M = I_n (MM^\dagger = I_m)$ if $m \geq n (m \geq n)$ and $M$ is full rank.

A square matrix $P \in \mathbb{R}^{n \times n}$ is called a **projector** if $P^2 = P$.

- $P$ is always diagonalisable and all eigenvalues are 1 or 0. (think why this is?)

- $\|P\|_2 \geq 1$ and $\|P\|_2 = 1$ iff $P = P^T$; in this case $P$ is called an orthogonal projector, and $P$ can be written as $P = QQ^T$ where $Q$ is orthonormal.

- One can easily show that $I - P$ is another projector as $(I - P)^2 = I - P$, and unless $P = 0$ or $P = I$, we have $\|I - P\|_2 = \|P\|_2$:
Schur form $QPQ^* = \begin{bmatrix} I & B \\ 0 & 0 \end{bmatrix}$, $Q(I - P)Q^* = \begin{bmatrix} 0 & -B \\ 0 & I \end{bmatrix}$; See Szyld 2006 [32] for many more about projections.

## 14.3 HMT approximant: analysis (down from 70 pages!)

We are now in a position to explain why the HMT algorithm works. The original HMT paper is over 70 pages with long analysis. Here we attempt to condense the arguments to the essence (and with a different proof). Recall that our low-rank approximation is $\hat{A} = QQ^T A$, where $AX = QR$. Goal: $\|A - \hat{A}\| = \|(I_m - QQ^T)A\| = O(\|A - A_{\hat{r}}\|)$.

1. $QQ^T AX = AX$ ($QQ^T$ is **orthogonal projector** onto $\text{span}(AX)$). Hence $(I_m - QQ^T)AX = 0$, so $A - \hat{A} = (I_m - QQ^T)A(I_n - XM^T)$ for any $M \in \mathbb{R}^{n \times r}$.

   The idea then is to choose $M$ cleverly such that the expression $(I_m - QQ^T)A(I_n - XM^T)$ can be shown to have small norm.

2. Set $M^T = (V^T X)^\dagger V^T$ where $V = [v_1, \ldots, v_{\hat{r}}] \in \mathbb{R}^{n \times \hat{r}}$ is the top right singular vectors of $A$ ($\hat{r} \leq r$).

   Recall that $\hat{r}$ is any integer bounded by $r$.

3. $VV^T(I - XM^T) = VV^T(I - X(V^TX)^\dagger V^T) = 0$ if $V^T X$ full row-rank (this is a generic assumption), so $A - \hat{A} = (I_m - QQ^T)A(I - VV^T)(I_n - XM^T)$.

4. Taking norms yields $\|A - \hat{A}\|_2 = \|(I_m - QQ^T)A(I - VV^T)(I_n - XM^T)\|_2 = \|(I_m - QQ^T)U_2\Sigma_2 V_2^T(I_n - XM^T)\|_2$ where $[V, V_2]$ is orthogonal (and $A = [U, U_2]\begin{bmatrix}\Sigma & \\ & \Sigma_2\end{bmatrix}[V, V_2]^T$ is the SVD), so

$$\|A - \hat{A}\|_2 \leq \|\Sigma_2\|_2\|(I_n - XM^T)\|_2 = \underbrace{\|\Sigma_2\|_2}_{\text{optimal rank-}\hat{r}} \boxed{\|XM^T\|_2}$$

It remains to prove $\|XM^T\|_2 = O(1)$. To see why this should hold with high probability, we need a result from random matrix theory.
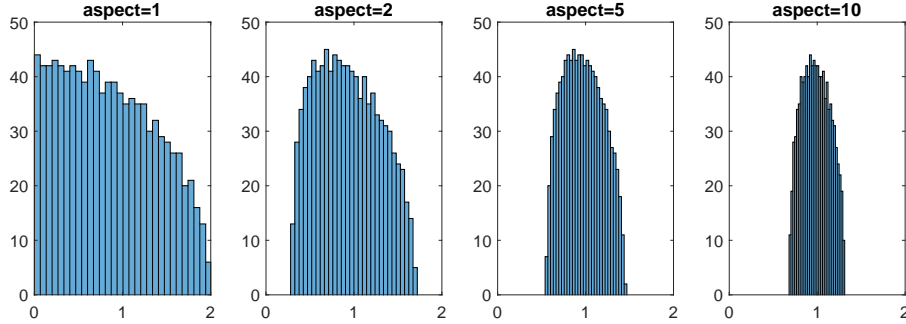
## 14.4 Tool from RMT: Rectangular random matrices are well conditioned

A final piece of information required to complete the puzzle is the *Marchenko-Pastur law*, a classical result in random matrix theory, which we will not be able to prove here (and hence is clearly nonexaminable). We refer those interested to the part C course on Random Matrix Theory. However, understanding the statement and the ability to use this fact is indeed examinable.

The key message is easy to state: a **rectangular random matrix is well conditioned** with extremely high probability. This fact is enormously important and useful in a variety of contexts in computational mathematics.

Here is a more precise statement:

**Theorem 14.1 (Marchenko-Pastur)** *The singular values of random matrix $X \in \mathbb{R}^{m \times n}$ $(m \geq n)$ with iid $X_{ij}$ (mean 0, variance 1) follow* Marchenko-Pastur (M-P) *distribution (proof nonexaminable), with density* $\sim \frac{1}{x}\sqrt{((\sqrt{m} + \sqrt{n}) - x)(x - (\sqrt{m} - \sqrt{n}))}$, *and support* $[\sqrt{m} - \sqrt{n}, \sqrt{m} + \sqrt{n}]$.



Histogram of singular values of random Gaussian matrices with varying aspect ratio $m/n$.

$\sigma_{\max}(X) \approx \sqrt{m} + \sqrt{n}$, $\sigma_{\min}(X) \approx \sqrt{m} - \sqrt{n}$, hence $\kappa_2(X) \approx \frac{1 + \sqrt{m/n}}{1 - \sqrt{m/n}} = O(1)$.

Proof: omitted. (strictly speaking, the theorem concerns the limit $m, n \to \infty$; but the result holds in the nonasymptotic limit with enormous probability [6]).

As stated above, this is a key fact in many breakthroughs in computational maths! Examples include

- Randomized SVD, Blendenpik (randomized least-squares)

- (nonexaminable:) Compressed sensing (RIP) [Donoho 06, Candes-Tao 06], Matrix concentration inequalities [Tropp 11], Function approx. by least-squares [Cohen-Davenport-Leviatan 13]

- (nonexaminable:) You might have heard of the Johnson-Lindenstrauss (JL) Lemma. A host (thousands) of papers have appeared that use JL to prove interesting results. It turns out that many of these can be equally well be proven (and IMO more easily/naturally understood) using MP.

$\|XM^T\|_2 = O(1)$ Let's get back to the HMT analysis. Recall that we've shown for $M^T = (V^T X)^\dagger V^T$ where $X \in \mathbb{R}^{n \times r}$ is random, that

$$\|A - \hat{A}\|_2 \leq \|\Sigma_2\|_2 \|(I_n - XM^T)\|_2 = \underbrace{\|\Sigma_2\|_2}_{\text{optimal rank-}\hat{r}} \|XM^T\|_2.$$

Now $\|XM^T\|_2 = \|X(V^T X)^\dagger V^T\|_2 = \|X(V^T X)^\dagger\|_2 \leq \|X\|_2 \|(V^T X)^\dagger\|_2$.

Now let's analyse the (standard) case where $X$ is random Gaussian $X_{ij} \sim \mathcal{N}(0, 1)$. Then

- $V^T X$ is another Gaussian matrix (an important fact about Gaussian matrices is that orthogonal×Gaussian=Gaussian (in distribution); this is nonexaminable but a nice exercise), hence $\|(V^T X)^\dagger\| = 1/\sigma_{\min}(V^T X) \lesssim 1/(\sqrt{r} - \sqrt{\hat{r}})$ by M-P.

- $\|X\|_2 \lesssim \sqrt{n} + \sqrt{r}$ by M-P.[21]

Together we get $\|XM^T\|_2 \lesssim \frac{\sqrt{n}+\sqrt{r}}{\sqrt{r}-\sqrt{\hat{r}}} = "O(1)"$.

Remark:

- When $X$ is a non-Gaussian random matrix, the performance is similar, but is harder to analyze. A popular choice is the so-called SRFT matrices, which use the FFT (fast Fourier transform) and can be applied to $A$ with $O(mn \log m)$ cost rather than $O(mnr)$.

## 14.5 Precise analysis for HMT (nonexaminable)

A slightly more elaborate analysis again using random matrix theory will give us a very sharp bound on the expected value of the error $E_{\mathrm{HMT}} =: A - \hat{A}$. (this again is non-examinable).

**Theorem 14.2 (Reproduces HMT 2011 Thm.10.5)** *If $X$ is Gaussian, for any $\hat{r} < r$,*
$\mathbb{E}\|E_{\mathrm{HMT}}\|_F \le \sqrt{\mathbb{E}\|E_{\mathrm{HMT}}\|_F^2} = \sqrt{1 + \frac{r}{r-\hat{r}-1}}\|A - A_{\hat{r}}\|_F.$

**Proof:**  First ineq: Cauchy-Schwarz. $\|E_{\mathrm{HMT}}\|_F^2$ is

$$\|A(I - VV^T)(I - \mathcal{P}_{X,V})\|_F^2 = \|A(I - VV^T)\|_F^2 + \|A(I - VV^T)\mathcal{P}_{X,V}\|_F^2$$
$$= \|\Sigma_2\|_F^2 + \|\Sigma_2 \mathcal{P}_{X,V}\|_F^2 = \|\Sigma_2\|_F^2 + \|\Sigma_2(V_\perp^T X)(V^T X)^\dagger V^T\|_F^2.$$

Now if $X$ is Gaussian then $V_\perp^T X \in \mathbb{R}^{(n-\hat{r}) \times r}$ and $V^T X \in \mathbb{R}^{\hat{r} \times r}$ are independent Gaussian. Hence by [HMT Prop. 10.1] $\mathbb{E}\|\Sigma_2(V_\perp^T X)(V^T X)^\dagger\|_F^2 = \frac{r}{r-\hat{r}-1}\|\Sigma_2\|_F^2$, so

$$\mathbb{E}\|E_{\mathrm{HMT}}\|_F^2 = \left(1 + \frac{r}{r - \hat{r} - 1}\right)\|\Sigma_2\|_F^2.$$

$\square$

Note how remarkable the theorem is—the 'lazily' computed approximant is nearly optimal up to a factor $\sqrt{1 + \frac{r}{r-\hat{r}-1}}$ for a near rank $\hat{r}$ (one can take, e.g. $\hat{r} = 0.9r$).

## 14.6 Generalised Nyström

We wish to briefly mention an algorithm for low-rank approximation that is even faster than HMT, especially when $r \gg 1$.

Let $X \in \mathbb{R}^{n \times r}$ be a random matrix as before; and set another random matrix $Y \in \mathbb{R}^{n \times (r+\ell)}$, and                    [Nakatsukasa arXiv 2020 [24]]

$$\hat{A} = \boxed{(AX(Y^T AX)^\dagger Y^T)A} = \mathcal{P}_{AX,Y}A.$$

---

[21]This and the next line has been corrected from $\|X\|_2 \lesssim \sqrt{m} + \sqrt{r}$ and $\|XM^T\|_2 \lesssim \frac{\sqrt{m}+\sqrt{r}}{\sqrt{r}-\sqrt{\hat{r}}} = "O(1)"$ on 24 March 2022.

Then $\hat{A}$ is another rank-$r$ approximation to $A$, and $A - \hat{A} = (I - \mathcal{P}_{AX,Y})A = (I - \mathcal{P}_{AX,Y})A(I - XM^T)$; choose $M$ s.t. $XM^T = X(V^TX)^\dagger V^T = \mathcal{P}_{X,V}$. Then $\mathcal{P}_{AX,Y}, \mathcal{P}_{X,V}$ are (nonorthogonal) projections, and
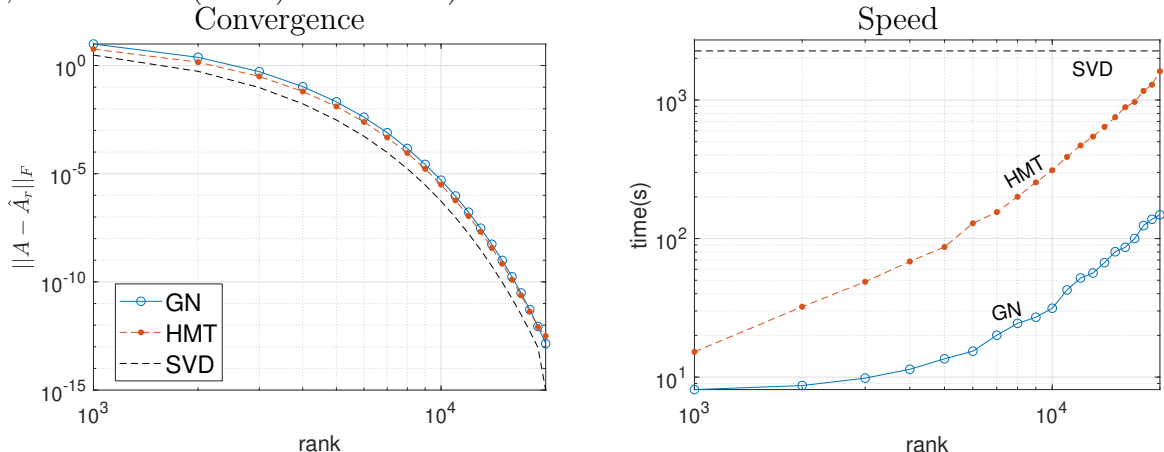
$$\begin{aligned}
\|A - \hat{A}\| &= \|(I - \mathcal{P}_{AX,Y})A(I - \mathcal{P}_{X,V})\| \\
&\leq \|(I - \mathcal{P}_{AX,Y})A(I - VV^T)(I - \mathcal{P}_{X,V})\| \\
&\leq \|A(I - VV^T)(I - \mathcal{P}_{X,V})\| + \|\mathcal{P}_{AX,Y}A(I - VV^T)(I - \mathcal{P}_{X,V})\|.
\end{aligned}$$

- Note that the $\|A(I - VV^T)(I - \mathcal{P}_{X,V})\|$ term is the exact same as in the HMT error.

- Extra term $\|\mathcal{P}_{AX,Y}\|_2 = O(1)$ as before if $c > 1$ in $Y \in \mathbb{R}^{m \times cr}$ (again, by Marchenko-Pastur).

- Overall, about $(1 + \|\mathcal{P}_{AX,Y}\|_2) \approx (1 + \frac{\sqrt{n} + \sqrt{r+\ell}}{\sqrt{r+\ell} - \sqrt{r}})$ times bigger expected error than HMT, **still near-optimal** and **much faster** $O(mn \log n + r^3)$.

Here's some experiments with HMT and GN (generalized Nyström ).

Let $A = U\Sigma V^T$ be a dense $30,000 \times 30,000$ matrix with geometrically decaying singular values $\sigma_i$, and choose $U, V$ to be random (orthogonal factors of a random Gaussian matrix).

We use HMT and GN to find low-rank approximations to $A$ varying the rank $r$. We also compare with MATLAB's SVD, which gives the optimal truncated SVD (up to numerical errors, which are $O(10^{-15})$ so invisible).



We see that randomised algorithms can outperform the standard SVD significantly.

## 14.7 MATLAB code

Implementing the HMT and GN algorithms (at least with Gaussian matrices, which don't always give optimal speed performance) is very easy!

Here's some sample MATLAB code. We can set up a matrix with exponentially decaying singular values (the matrix is constructed by computing $A = U\Sigma V^T$, where $U, V$ are orthonormal and $\Sigma$ is a geometric series from $10^{-100}$ to $1$).

```
n = 1000; % size
A = gallery('randsvd',n,1e100); % geometrically decaying singvals
r = 200;  % rank
```

Then to do HMT as follows:

```
X = randn(n,r);
AX = A*X;
[Q,R] = qr(AX,0); % QR fact.
At = Q*(Q'*A);
```

which with high probability gives me an excellent approximation

```
norm(At-A,'fro')/norm(A,'fro')
ans = 1.2832e-15
```

And for Generalized Nyström :

```
X = randn(n,r); Y = randn(n,1.5*r);
AX = A*X;   YA = Y'*A;   YAX = YA*X;
[Q,R] = qr(YAX,0);  % stable pseudo-inverse via the QR factorisation
At = (AX/R)*(Q'*YA);
```

```
norm(At-A,'fro')/norm(A,'fro')
ans = 2.8138e-15
```

Both algorithms give an excellent low-rank approximation to $A$.

# 15  Randomized least-squares: Blendenpik

Our final topic is a randomised algorithm for the least-squares problems that are highly overdetermined.

[Avron-Maymounkov-Toledo 2010 [1]]

$$\min_x \|Ax - b\|_2, \qquad \boxed{A} \in \mathbb{R}^{m \times n}, \ m \gg n$$

- Traditional method: normal eqn $x = (A^T A)^{-1} A^T b$ or $A = QR, x = R^{-1}(Q^T b)$, both require $O(mn^2)$ cost.

- Randomized: generate random $G \in \mathbb{R}^{4n \times m}$, and $\boxed{\ \ \ G\ \ \ } \boxed{A} = \boxed{\hat{Q}}\ \boxed{\hat{R}}$

(QR factorisation), then solve $\min_y \|(A\hat{R}^{-1})y - b\|_2$'s normal eqn via Krylov.

- $O(mn \log m + n^3)$ cost using fast FFT-type transforms[22] for $G$.
- Crucially, $A\hat{R}^{-1}$ is well-conditioned. Why? Marchenko-Pastur (next)

## 15.1 Explaining Blendenpik via Marchenko-Pastur

Let us prove that $\kappa_2(A\hat{R}^{-1}) = O(1)$ with high probability. A key result, once again, is M-P.

Claim: $A\hat{R}^{-1}$ is well-conditioned with $\boxed{G}\,\boxed{A} = \boxed{\hat{Q}}\,\boxed{\hat{R}}$ (QR factorisation)

Let's prove this for $G \in \mathbb{R}^{4n \times m}$ Gaussian:

Proof: Let $A = QR$. Then $GA = (GQ)R =: \tilde{G}R$

- $\boxed{\tilde{G}}$ is $4n \times n$ rectangular Gaussian, hence well-conditioned.

- So by M-P, $\kappa_2(\tilde{R}^{-1}) = O(1)$ where $\tilde{G} = \tilde{Q}\tilde{R}$ is the QR factorisation.

- Thus $\tilde{G}R = (\tilde{Q}\tilde{R})R = \tilde{Q}(\tilde{R}R) = \tilde{Q}\hat{R}$, so $\hat{R}^{-1} = R^{-1}\tilde{R}^{-1}$.

- Hence $A\hat{R}^{-1} = Q\tilde{R}^{-1}$, so $\kappa_2(A\hat{R}^{-1}) = \kappa_2(\tilde{R}^{-1}) = O(1)$.

## 15.2 Blendenpik: solving $\min_x \|Ax - b\|_2$ using $\hat{R}$

We have $\kappa_2(A\hat{R}^{-1}) =: \kappa_2(B) = O(1)$; defining $\hat{R}x = y$, $\min_x \|Ax - b\|_2 = \min_y \|(A\hat{R}^{-1})y - b\|_2 = \min_y \|By - b\|_2$.

- $B$ is well-conditioned$\Rightarrow$in normal equation

$$B^T B y = B^T b \tag{11}$$

$B^T B$ is also well-conditioned $\kappa_2(B^T B) = O(1)$; so positive definite and well-conditioned.

- Thus we can solve (11) via CG (or LSQR [27], a more stable variant in this context; nonexaminable)

  - exponential convergence, $O(1)$ iterations! (or $O(\log \frac{1}{\epsilon})$ iterations for $\epsilon$ accuracy)
  - each iteration requires $w \leftarrow Bw$ and $w \leftarrow B^T w$, consisting of $w \leftarrow \hat{R}^{-1}w$ ($n \times n$ triangular solve) and $w \leftarrow Aw$ ($m \times n$ matrix-vector multiplication); $O(mn)$ cost overall

---

[22]The FFT (fast Fourier transform) is one of the important topics that we can't treat properly—for now just think of it as a matrix-vector multiplication that can be performed in $O(n \log n)$ flops rather than $O(n^2)$.)

## 15.3  Blendenpik experiments

Let's illustrate our findings. Since Blendenpik finds a preconditioner such that $AR^{-1}$ is well-conditioned (regardless of $\kappa_2(A)$), we expect the convergence of CG to be independent of $\kappa_2(A)$. This is indeed what we see here:
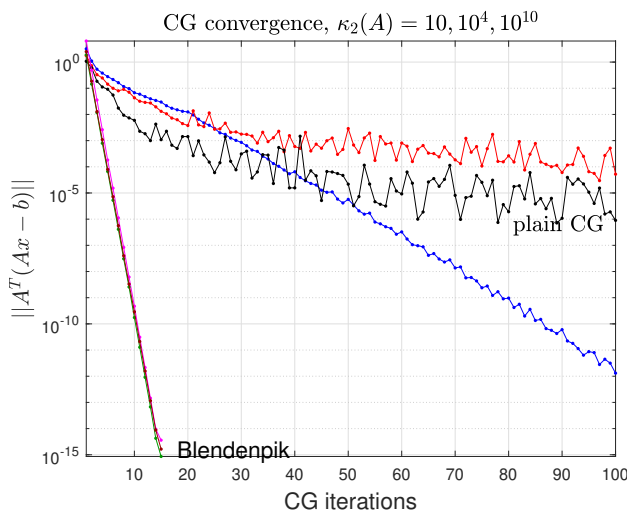


Figure 3: CG for $A^T A x = A^T b$ vs. Blendenpik $(AR^{-1})^T(AR^{-1})x = (AR^{-1})^T b$, $m = 10000, n = 100$

In practice, Blendenpik gets $\approx \times 5$ speedup over classical (Householder-QR based) method when $m \gg n$.

## 15.4  Sketch and solve for $\min_x \|Ax - b\|_2$

In closing, let us descibe a related, simpler algorithm. To solve a least-squares problem

$$\underset{x}{\text{minimize}} \, \|Ax - b\|_2, \tag{12}$$

where $A \in \mathbb{C}^{n \times r}$, $n \gg r$, one can *sketch and solve* the problem [35]: draw a random matrix $G \in \mathbb{C}^{\tilde{r} \times n}$ where $\tilde{r} = O(r) \ll n$, and solve the sketched problem

$$\underset{x}{\text{minimize}} \, \|G(Ax - b)\|_2. \tag{13}$$

In some cases, the solution of this problem is already a good enough approximation to the original problem.

We've taken $G$ to be Gaussian above. For a variety of choices of $G$, with high probability the solutions for (13) and (12) can be shown to be similar in that they have a comparable residual $\|Ax - b\|_2$.

Let's understand why the solution for (13) should be good, again using Marchenko-Pastur. Let $[A\, b] = QR \in \mathbb{C}^{m \times (n+1)}$ be a thin QR factorization, and suppose that the sketch $G \in \mathbb{C}^{\tilde{r} \times n}$

is Gaussian. Then $GQ$ is rectangular Gaussian (again using the rotational invariance of Gaussian random matrices), so well-conditioned. Suppose without loss of generality that $G$ is scaled so that $1 - \delta \leq \sigma_i(GQ) \leq 1 + \delta$ for some $\delta < 1$ (note that here $\delta$ isn't an $O(u)$ quantity, $\delta = 0.5$, say, is a typical number), and therefore (since $Qv = [A\ b]\tilde{v}$ for some $\tilde{v}$) for any $\tilde{v} \in \mathbb{C}^{n+1}$ we have $(1 - \delta)\|[A\ b]\tilde{v}\|_2 \leq \|G[A\ b]\tilde{v}\|_2 \leq (1 + \delta)\|G[A\ b]\tilde{v}\|_2$. Taking $\tilde{v} = \begin{bmatrix} x \\ -1 \end{bmatrix}$ it follows that for any vector $x \in \mathbb{C}^n$ we have

$$(1 - \delta)\|Ax - b\|_2 \leq \|G(Ax - b)\|_2 \leq (1 + \delta)\|Ax - b\|_2.$$

Consequently, the minimizer $x_s$ of $\|G(Ax - b)\|_2$ for (13) also minimizes $\|Ax - b\|_2$ for (12) up to a factor $\frac{1+\delta}{1-\delta}$. If the residual of the original problem can be made small, say $\|Ax - b\|_2/\|b\|_2 = 10^{-10}$, then $\|Ax_G - b\|_2/\|b\|_2 \leq \frac{1+\delta}{1-\delta} \times 10^{-10}$, which with a modest and typical value $\delta = \frac{1}{2}$ gives $3 \times 10^{-10}$, giving an excellent least-squares fit. If $A$ is well-conditioned, this also impliees the solution $x$ is close to the exact solution.

## 15.5  Randomized algorithm for $Ax = b$, $Ax = \lambda x$?

We have seen that randomization can be very powerful for low-rank approximation and least-squares problems. What about the two core problems in NLA, linear systems $Ax = b$ and eigenvalue problems $Ax = \lambda x$? A recent preprint [25] proposes the use of sketching (as described above) for GMRES (for $Ax = b$) and Rayleigh-Ritz (for eigenproblems) combined with a generation of a basis for a Krylov subspace that is not orthonormal. This is certainly not the end of the story.

Randomized algorithms promise to be employed in more and more applications and problems. We will almost surely see more breakthroughs in randomized NLA. Who be the inventors? Would the young and fresh brains like to take a shot?

# 16  Conclusion and discussion

We have tried to give a good overview of the field of numerical linear algebrain this course. However a number of topics have been omitted completely due to lack of time. Here is an incomplete list of other topics.

## 16.1  Important (N)LA topics not treated

These are clearly nonexaminable but definitely worth knowing if you want to get seriously into the field. These will be discuss in lecture if and only if time permits, and in any case only superficially.

- tensors                                                                    [Kolda-Bader 2009 [21]]

- FFT (values↔coefficients map for polynomials)        [e.g. Golub and Van Loan 2012 [14]]

- sparse direct solvers                                                [Duff, Erisman, Reid 2017 [10]]

- multigrid  [e.g. Elman-Silvester-Wathen 2014 [11]]

- fast (Strassen-type) matrix multiplication etc  [Strassen 1969 [31]+many follow-ups]

- functions of matrices  [Higham 2008 [18]]

- generalised, polynomial/nonlinear eigenvalue problems  [Guttel-Tisseur 2017 [15]]

- perturbation theory (Davis-Kahan [7] etc)  [Stewart-Sun 1990 [30]]

- compressed sensing (this deals with $Ax = b$ where $A$ is very 'fat')  [Foucart-Rauhut 2013 [12]]

- model order reduction  [Benner-Gugercin-Willcox 2015 [3]]

- communication-avoiding algorithms  [e.g. Ballard-Demmel-Holtz-Schwartz 2011 [2]]

## 16.2   Course summary

This information is provided for MSc students who take two separate exams based on the 1st and 2nd halfs.

1st half

- SVD and its properties (Courant-Fischer etc), applications (low-rank)

- Direct methods (LU) for linear systems and least-squares problems (QR)

- Stability of algorithms


2nd half

- Direct method (QR algorithm) for eigenvalue problems, SVD

- Krylov subspace methods for linear systems (GMRES, CG) and eigenvalue problems (Arnoldi, Lanczos)

- Randomized algorithms for SVD and least-squares


## 16.3   Related courses you can take

Courses with significant intersection with NLA include

- C6.3 Approximation of Functions: Chebyshev polynomials/approximation theory

- C7.7 Random Matrix Theory: for theoretical underpinnings of Randomized NLA

- C6.4 Finite Element Method for PDEs: NLA arising in solutions of PDEs

- C6.2 Continuous Optimisation: NLA in optimisation problems

and many more: differential equations, data science, optimisation, machine learning,... NLA is everywhere in computational mathematics.

Thank you for your interest in NLA!

# References

[1] H. Avron, P. Maymounkov, and S. Toledo. Blendenpik: Supercharging LAPACK's least-squares solver. *SIAM J. Sci. Comp.*, 32(3):1217–1236, 2010.

[2] G. Ballard, J. Demmel, O. Holtz, and O. Schwartz. Minimizing communication in numerical linear algebra. *SIAM J. Matrix Anal. Appl.*, 32(3):866–901, 2011.

[3] P. Benner, S. Gugercin, and K. Willcox. A survey of projection-based model reduction methods for parametric dynamical systems. *SIAM Rev.*, 57(4):483–531, 2015.

[4] R. Bhatia. *Positive Definite Matrices*. Princeton University Press, 2009.

[5] K. Braman, R. Byers, and R. Mathias. The multishift QR algorithm. Part II: Aggressive early deflation. *SIAM J. Matrix Anal. Appl.*, 23:948–973, 2002.

[6] K. R. Davidson and S. J. Szarek. Local operator theory, random matrices and banach spaces. *Handbook of the geometry of Banach spaces*, 1(317-366):131, 2001.

[7] C. Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.*, 7(1):1–46, 1970.

[8] J. Demmel. *Applied Numerical Linear Algebra*. SIAM, Philadelphia, USA, 1997.

[9] J. Dongarra and F. Sullivan. Guest editors' introduction: The top 10 algorithms. *IEEE Computer Architecture Letters*, 2(01):22–23, 2000.

[10] I. S. Duff, A. M. Erisman, and J. K. Reid. *Direct Methods for Sparse Matrices*. Oxford University Press, 2017.

[11] H. C. Elman, D. J. Silvester, and A. J. Wathen. *Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics*. Oxford University Press, USA, 2014.

[12] S. Foucart and H. Rauhut. *A Mathematical Introduction to Compressive Sensing*. Springer, 2013.

[13] D. F. Gleich. Pagerank beyond the web. *SIAM Rev.*, 57(3):321–363, 2015.

[14] G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, 4th edition, 2012.

[15] S. Güttel and F. Tisseur. The nonlinear eigenvalue problem. *Acta Numer.*, 26:1–94, 2017.

[16] N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.

[17] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, PA, USA, second edition, 2002.

[18] N. J. Higham. *Functions of Matrices: Theory and Computation*. SIAM, Philadelphia, PA, USA, 2008.

[19] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1991.

[20] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, second edition, 2012.

[21] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Rev.*, 51(3):455–500, 2009.

[22] P.-G. Martinsson and J. A. Tropp. Randomized numerical linear algebra: Foundations and algorithms. *Acta Numer.*, pages 403—572, 2020.

[23] M. F. Murphy, G. H. Golub, and A. J. Wathen. A note on preconditioning for indefinite linear systems. *SIAM J. Sci. Comp.*, 21(6):1969–1972, 2000.

[24] Y. Nakatsukasa. Fast and stable randomized low-rank matrix approximation. arXiv:2009.11392.

[25] Y. Nakatsukasa and J. A. Tropp. Fast & accurate randomized algorithms for linear systems and eigenvalue problems. arXiv 2111.00113.

[26] C. C. Paige. Error analysis of the Lanczos algorithm for tridiagonalizing a symmetric matrix. *IMA J. Appl. Math.*, 18(3):341–349, 1976.

[27] C. C. Paige and M. A. Saunders. LSQR: An algorithm for sparse linear equations and sparse least squares. *ACM Trans. Math. Soft.*, 8(1):43–71, 1982.

[28] B. N. Parlett. *The Symmetric Eigenvalue Problem*. SIAM, Philadelphia, 1998.

[29] Y. Saad and M. H. Schultz. GMRES - A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM J. Sci. Stat. Comp.*, 7(3):856–869, 1986.

[30] G. W. Stewart and J.-G. Sun. *Matrix Perturbation Theory (Computer Science and Scientific Computing)*. Academic Press, 1990.

[31] V. Strassen. Gaussian elimination is not optimal. *Numer. Math.*, 13:354–356, 1969.

[32] D. B. Szyld. The many proofs of an identity on the norm of oblique projections. *Numerical Algorithms*, 42(3-4):309–323, 2006.

[33] L. N. Trefethen and D. Bau. *Numerical Linear Algebra*. SIAM, Philadelphia, 1997.

[34] M. Udell and A. Townsend. Why are big data matrices approximately low rank? *SIAM Journal on Mathematics of Data Science*, 1(1):144–160, 2019.

[35] D. P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends®  in Theoretical Computer Science*, 10(1–2):1–157, 2014.