
Comp Prelim Study Material

Trevor Squires
trevorsquires9@gmail.com

Last updated January 8, 2023

Contents

| | | |
|-----------|---|-----------|
| I | List of Topics | 1 |
| 1 | Math 8600 | 2 |
| 2 | Math 8610 | 3 |
| II | 8600 Notes | 4 |
| 1 | Scientific Computing Overview | 6 |
| 1 | Error | 6 |
| 2 | Solving Problems and Algorithm Properties | 6 |
| 2 | Roundoff Errors | 8 |
| 1 | Representation | 8 |
| 2 | Consequences of Finite Number System | 8 |
| 3 | Roots to Nonlinear Equations | 10 |
| 1 | Bisection Method | 11 |
| 2 | Fixed Point Iteration | 11 |
| 3 | Newton's Method | 12 |
| 4 | Secant Method | 12 |
| 4 | Direct Methods of Solving Linear Systems | 14 |
| 1 | Gaussian Elimination | 14 |
| 2 | Forward/Backward Substitution | 14 |
| 3 | LU Factorization | 14 |
| 4 | Pivoting Strategies | 14 |
| 5 | Condition Number | 15 |
| 5 | Least Squares Problem | 16 |
| 1 | An Analytical Solution | 16 |
| 2 | Naive Numerical Solution | 16 |
| 3 | Other Numerical Methods | 17 |
| 6 | Polynomial interpolation | 18 |
| 1 | Monomial interpolation | 19 |
| 2 | Lagrange Interpolation | 19 |
| 3 | Newton Polynomial Interpretation | 20 |
| 4 | Error in Polynomial Interpolation | 21 |
| 7 | Piecewise Polynomial Interpolation | 23 |
| 1 | Broken line and piecewise Hermite interpolation | 23 |
| 2 | Piecewise cubic interpolation | 23 |
| 3 | Cubic Spline Interpolation | 24 |

| | | |
|------------|---|-----------|
| 8 | Numerical Differentiation | 26 |
| 1 | Taylor Series Approximations | 26 |
| 2 | Richardson Extrapolation | 27 |
| 3 | Lagrange Polynomial Approximations | 27 |
| 9 | Numerical Integration | 28 |
| 1 | Basic Quadrature Rules | 28 |
| 2 | Quadrature Error | 29 |
| 3 | Composite Numerical Integration | 29 |
| 4 | Gaussian Quadratures | 30 |
| 5 | Adaptive Quadrature | 31 |
| 6 | Romberg Integration | 31 |
| 10 | Differential Equations | 33 |
| 1 | Euler Methods | 33 |
| 1.1 | Preliminaries | 33 |
| 1.2 | Method Errors | 34 |
| 1.3 | Stability | 35 |
| 2 | Runge-Kutta Methods | 35 |
| III | 8610 Notes | 37 |
| 11 | Conditioning and Stability | 39 |
| 1 | Conditioning of a Problem | 39 |
| 2 | Conditioning of a System of Equations | 40 |
| 3 | Conditioning of Eigenvalues of Matrices | 40 |
| 4 | Conditioning of Roots of a Polynomial | 41 |
| 5 | Algorithm Stability | 41 |
| 6 | Stability of Linear Solvers | 42 |
| 7 | Conditioning of GE/LU factorization | 43 |
| 12 | QR Factorization | 45 |
| 1 | Properties of QR | 45 |
| 2 | QR Factorization via Gram-Schmidt | 46 |
| 3 | QR Factorization via Orthogonal Transformations | 47 |
| 4 | QR Factorization via Given's Rotation | 49 |
| 5 | Given's Rotation Computational Cost | 49 |
| 13 | Singular Value Decomposition | 51 |
| 1 | SVD Review | 51 |
| 1.1 | Low Rank Approximations | 51 |
| 2 | Computing an SVD | 52 |
| 2.1 | Naive Idea | 52 |
| 2.2 | Golub-Kahan Bidiagonalization | 53 |
| IV | Exercises | 54 |
| A | 8600 Exercises | 55 |
| 1 | Scientific Computing Fundamentals | 55 |
| 1.1 | Numerical Algorithms | 55 |
| 1.2 | Roundoff Errors | 56 |
| 1.3 | Nonlinear Equations of One Variable | 57 |
| 2 | Numerical Systems Analysis | 58 |
| 2.1 | Direct Methods for Linear Systems | 58 |
| 2.2 | Linear Least Squares Problems | 60 |
| 3 | Numerical Approximation | 61 |

| | | |
|----------|---|-----------|
| 3.1 | Polynomial Approximation | 61 |
| 3.2 | Piecewise Polynomial Interpolation | 63 |
| 3.3 | Numerical Differentiation | 64 |
| 3.4 | Numerical Integration | 64 |
| B | 8610 Exercises | 68 |
| 1 | Numerical Linear Algebra Fundamentals | 68 |
| 2 | Conditioning and Stability | 71 |
| 3 | QR and Linear Least Squares | 76 |

Part I
List of Topics

Below are a list of topics for Math 8600 and 8610 taken directly from the school's webpage.

1 Math 8600

1. Scientific computing
 - Floating point number system
 - Floating point arithmetic
 - Sensitivity and conditioning
2. Systems of Linear Equations
 - Back solving and forward solving
 - Gauss transformation
 - LU decomposition
 - Cholesky decomposition
 - Band matrix
 - Vector norm, matrix norm, and condition number
 - Sensitivity of a solution of a linear system
3. Linear Least Squares
 - Existence of a solution of a linear least squares problem
 - Normal equations
 - QR decomposition
 - Singular value decomposition
4. Nonlinear equations
 - Rate of convergence
 - Bisection method
 - Regular falsi method
 - Fixed point iteration
 - Newton's method
 - Secant method
5. Interpolation
 - Polynomial interpolation
 - method of undetermined coefficients
 - Lagrange interpolation
 - Neville's algorithm
 - error analysis
 - Piecewise polynomial interpolation
 - hermite cubic interpolation
 - cubic spline interpolation
6. Numerical Integration and Differentiation
 - Newton-Cotes quadrature
 - Gaussian quadrature
 - Composite and adaptive quadrature
 - Richardson's extrapolation

- Romberg integration

7. Initial Value Problems for Ordinary Differential Equations

- Introduction
- One step method
 - Euler method
 - Taylor method
 - Runge-Kutta method
 - order of accuracy and error analysis
- Multi-step method
 - Adams methods
 - predictor-corrector method
- Stability
- Stiff equation

2 Math 8610

1. Conditioning and Stability

- Condition and condition number
- Forward and backward stability
- Growth factor and stability of LU and other similar factorizations

2. QR factorization and linear least squares

- QR factorization by modified Gram-Schmidt
- QR factorization by Householder reflectors and Givens rotations
- Linear least squares problem and solution algorithms

3. Singular value decompositions (SVD)

- Definition of SVD, and its important relations and properties
- Golub-Khan bidiagonalization and the equivalence of SVD on two symmetric eigenvalue problems
- Applications such as low-rank approximation

4. Eigenvalue problems and algorithms

- Diagonalization, complex and real Schur form
- Reduction to upper Hessenberg/tridiagonal form
- Shifted QR iteration and its important relations and properties
- Simultaneous iteration and Arnoldi/Lanczos method for computing several eigenvalues

5. Iterative methods for large sparse linear systems

- Conjugate gradient (CG) method
- Generalized minimal residual (GMRES) method
- Preconditioned linear systems

6. Iteration complexity for all non-iterative algorithms

Part II
8600 Notes

This part contains a condensed set of notes from 8600. The notes are primarily motivated by Ascher and Grief the textbook for 8600. These are not intended to be a replacement to the comp courses, but rather as a supplement. Motivation, big ideas, and clarity will be emphasized, but will be fairly light on details. The fundamentals will be covered, but one is advised to attempt the problems in later sections for full exposure. The listed topics in the previous section will be the backbone of these notes.

Chapter 1

Scientific Computing Overview

We define **scientific computing** as the development and studying of numerical algorithms for solving mathematical problems. In a standard setting, one might model a naturally occurring problem using a mathematical model. This model, usually continuous, is sometimes difficult (and often outright impossible) to describe finitely. A natural step from here is to approximate the continuous (infinite) model with a finite dimensional one - that can be solved on the computer. **Numerical Analysis** is the study of of such approximations and resulting algorithms.

1 Error

Definition. Error is unavoidable in scientific computing. We quantify it with absolute and relative error. Let u be the true value of some quantity and v be the approximation. Then we define relative error as

$$\text{Err}_r(u, v) := \frac{|u - v|}{|u|}$$

and absolute error as

$$\text{Err}_A(u, v) := |u - v|$$

In addition to quantifying error, it may be helpful to describe them. We can identify three main sources of error

1. Modeling - A mathematical formulation rarely exactly describes a real life phenomenon. In practice, instruments do not record exact measurements and human error is omnipresent. These issues give us modeling errors.
2. Approximation - Although many problems can be described using an infinite process, this is not always easy to work with. An replacement of a infinite process with a finite one introduces approximation errors. Even further, we can say that
 - Discretization errors arise from discretizations of a continuous process
 - Convergence errors arise from the termination or truncation of an infinite process
3. Roundoff - Because we desire to use computers to solve our models, we must work in finite precision. This finite precision leads to roundoff errors.

For this note, we will primarily be concerned with roundoff errors.

2 Solving Problems and Algorithm Properties

Once a mathematical model is formulated, an algorithm is sought after to solve the model. We may describe an algorithm by a few characteristics

1. Accuracy - the ability to provide an accurate solution upon termination

2. Efficiency - the effort required to provide an accurate solution. This is usually measured in flops of number of function evaluations
3. Robustness - the reliability of an algorithm. A good algorithm will work in most cases and be able to describe the situations in which it doesn't. It will be reliable and stable.

On the other hand, there are issues that are problem dependent and not algorithm dependent. The most important of these is conditioning.

Definition. Problem conditioning refers to the propensity of change of a solution given a change in input data.

We usually describe conditioning using one of two straight-forward terms

1. Ill-conditioned - a small perturbation in the data would produce a large difference in the result
2. Well-conditioned - the solution is resistant to small changes in the input data

Loosely speaking, we say that the condition number is $\frac{\text{relative output}}{\text{relative input}}$. A well conditioned problem has a condition number close to one. This reasoning has an additional meaning for functions. Note that

$$\text{cond number} = \frac{\frac{f(x) - f(\hat{x})}{f(x)}}{\frac{x - \hat{x}}{x}} = f'(\eta) \frac{x}{f(x)} \leq \max f'(x) \frac{x}{f(x)}$$

Chapter 2

Roundoff Errors

In this section, we discuss the aforementioned roundoff errors. While the details are of importance and can be followed closely for an in-depth analysis, the main take away from this section is when roundoff errors become an issue.

1 Representation

A real number $x \in \mathbb{R}$ can be written as

$$x = \pm(1.d_1d_2d_3\dots d_t) \cdot \beta^e$$

where $d_i \in \{0, 1\}$, e is the integer exponent, and β is the base of representation. Letting $\beta = 2$, this is referred to as the standard binary format. Digits after the decimal point are called the mantissa. That is, to refer to a number on a computer (in binary), we need to know

- the sign
- the mantissa
- the exponent

Note here that the first digit is always a 1 for normalization. Given such a representation, we say that this represents

$$(x) = \pm\left(\frac{d_1}{\beta^0} + \dots + \frac{d_t}{\beta^{t-1}}\right) \cdot \beta^e$$

However, it is far more common to work in the IEEE standard than anything else. For this reason, future discussion will focus on the IEEE standard (with $\beta = 2$).

Definition. The IEEE standard provides specific usage of the 64 bits of precision provided by most modern computers. It breaks them down in the following way

1. 1 bit for the sign of the number
2. 52 bits for the mantissa
3. 11 bits for the exponent

2 Consequences of Finite Number System

There are a number of issues that stem from this representation. A few are listed below

- There are largest and smallest numbers
- There are a finite number of representable quantities
- The absolute difference between two closest is not the same for all numbers
- Every number has a potential error of $\varepsilon_{\text{mach}} = 2^{-52} \approx 10^{-16}$. This error can accumulate if not accounted for and ruin algorithm accuracy.

Possible error sources

Some operations are more prone to significant errors than others (all operations have errors, but not all errors are meaningful). We list a few of these below

1. Adding large and small numbers. In general, the smaller a number is, the more likely it is affected by a large relative error. Adding a large and smaller number, i.e. $1 + 10^{-16}$, is likely to produce large relative error.
2. Product/division of numbers close to 0. If $y \ll 1$ then xy and $\frac{x}{y}$ can have large relative and absolute error.
3. Subtraction of similarly sized numbers. If $x \approx y$, then $x - y \approx 0$, but due to cancellation error, can be as large as $2\varepsilon_{\text{mach}}$. The relative error is even worse.

Chapter 3

Roots to Nonlinear Equations

In this section we discuss methods of finding roots to nonlinear equations. Per the results of the Abel Ruffini Theorem, there exists no algebraic solution to a general polynomial of degree 5 or higher. As such, we should not expect our solutions to be finite processes, but rather iterative ones. With this in mind, it makes sense to first consider the standard nuances of iterative algorithms: stopping criterion and desired properties.

When solving for roots of nonlinear equations, there are a few stopping criterion that may become useful. Let $\varepsilon > 0$ be some specialized tolerance and $\{x_k\}_{k=1}^N$ be the iterates of a nonlinear solver used to solve $f(x) = 0$.

1. Absolute tolerance

$$|x_n - x_{n-1}| < \varepsilon$$

2. Relative tolerance

$$|x_n - x_{n-1}| < \varepsilon |x_n|$$

3. Functional tolerance

$$|f(x_n)| < \varepsilon$$

It is important to note that there is no dominance among these 3 criteria. We simply list them all because some methods of convergence analysis may be more suitable for finding bounds of specific stopping criterion. Furthermore, we would like our algorithm to (in addition to the previous discussion) satisfy some of the following

- Small requirements on the smoothness of f
- Little dependence on function evaluations
- Generalizes easy
- Robust

Last but not least, we more accurately describe what it means for an algorithm to be slow (or fast).

Definition. There are 3 basic types of convergence for iterative methods

1. Linear Convergence

$$|x_{n+1} - x^*| < C |x_n - x^*|$$

for $C < 1$.

2. Superlinear Convergence

$$|x_{n+1} - x^*| < \rho_n |x_n - x^*|$$

where $\rho_n \rightarrow 0, n \rightarrow \infty$.

3. Quadratic Convergence

$$|x_{n+1} - x^*| < M |x_n - x^*|^2$$

With these in mind, we look at a few methods in particular

1 Bisection Method

The motivation behind the bisection method is quite simple. For a continuous function f , if $f(a) < 0$ and $f(b) > 0$, then by the IVT, there must exist a c between these points such that $f(c) = 0$. By querying f at the right places, we can reduce the search space. That is, if a and b are the points from above and $f(c) < 0$, then our search interval is reduced from $[a, b]$ to $[c, b]$. Otherwise, it becomes $[a, c]$. In order to minimize the expected interval change, we should choose c such that $c = \frac{a+b}{2}$. Then with each iteration (one function evaluation), our search interval decreases by a factor of 2.

That is, with each iteration, the absolute error decreases by a factor of 2 since

$$|x_{n+1} - x^*| < 0.5 |x_n - x^*| < (0.5)^n |x_0 - x^*|$$

That is, the Bisection Method is a linearly convergent algorithm. Thus, to achieve an $\varepsilon > 0$ tolerance, we need N to satisfy

$$N > \log\left(\frac{x_0 - x^*}{\varepsilon}\right) > \log\left(\frac{b_0 - a_0}{\varepsilon}\right)$$

We summarize the Bisection Method with the following pros and cons

Pros

- Robust (Globally convergent)
- Only uses function evaluations
- Only needs a fixed number of iterations

Cons

- Not very efficient
- Does not generalize

2 Fixed Point Iteration

Fixed Point Iteration, or FPI, is a method of solving $g(x) = x$ for some function g . You may wonder why this discussion belongs in the root solving section of these notes. Hopefully it is not too hard to see that for any equation $f(x) = 0$, we can find a (not necessarily unique) equation $g(x) = x$ such that x^* solves the former if and only if x^* also solves the latter. FPI considers the sequence $\{x_n\}$ where $x_{n+1} = g(x_n)$. If this sequence converges, then we have found a point x_n such that $g(x_n) = x_n$. The conditions on which this sequence converge are demonstrated in the following theorem.

Theorem 3.1. *If $g \in C([a, b])$ and $a \leq g(x) \leq b$ for all $x \in [a, b]$, then there is a fixed point x^* of g in $[a, b]$. Additionally, if $|g'(x)| < \rho$ for some $\rho < 1$ and any $x \in [a, b]$, then x^* is unique.*

Proof. Hint: Consider $\phi(x) = g(x) - x$ and use IVT. □

For convergence, note that

$$\begin{aligned} |x_{k+1} - x^*| &= |g(x_k) - g(x^*)| \\ &\leq |g'(c_k)| |x_k - x^*| \\ &\leq |\rho| |x_k - x^*| \end{aligned}$$

That is, FPI is a linearly convergent method with $C = \rho$ from the theorem.

Pros

- Easy to generalize
- Only uses function evaluations

Cons

- Not very efficient
- May be difficult to find a good function g

3 Newton's Method

Newton's method is derived from the Taylor expansion of f . Instead of looking for roots of f why not find a root of a linear approximation of f which is a considerably easier task? The linear approximation is available via the Taylor expansion

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + f''(\eta_k)(x - x_k)^2$$

for some η_k . Letting $x = x^*$ and ignoring the error term, this becomes

$$0 \approx f(x_k) + f'(x_k)(x^* - x_k)$$

or that

$$x^* \approx x_k - \frac{f(x_k)}{f'(x_k)}$$

Of course, this is only an approximation so we instead adopt the scheme

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

We can prove that under suitable conditions, this is a quadratically convergent algorithm.

Theorem 3.2. *Let $f(x^*) = 0, f'(x^*) \neq 0, f \in C^2$. Then for x_0 chosen close enough to x^* , Newton's method converges at least quadratically.*

Proof. Consider the Taylor expansion of $f(x^*)$ at $f(x_k)$

$$0 = f(x_k) + f'(x_k)(x^* - x_k) + \frac{1}{2}f''(c_k)(x^* - x_k)^2$$

Rearranging we obtain

$$-\frac{f(x_k)}{f'(x_k)} = x^* - x_k + \frac{f''(c_k)}{2f'(x_k)}(x^* - x_k)^2$$

or that

$$\frac{|e_{k+1}|}{|e_k|^2} = \frac{x_{k+1} - x^*}{(x^* - x_k)^2} = \frac{f''(c_k)}{2f'(x_k)} = M$$

which is our definition of quadratically convergent. □

Pros

- Easy to generalize
- Extremely fast

Cons

- Requires gradient
- Only locally convergent

4 Secant Method

One of the major cons of Newton's method is that it requires us to compute gradient values at each step. The secant method tackles this problem directly by using the secant line to approximate the gradient. The update for the secant method is

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}$$

Unfortunately, this approximation loses the quadratic convergence. However, the secant method still achieves a superlinear convergence.

Pros

- Easy to generalize
- Super linear
- No gradient computations

Cons

- Slower than Newton's
- Only locally convergent

Chapter 4

Direct Methods of Solving Linear Systems

Systems of linear equations arise in every mathematical field on a consistent basis. This section overviews a few basic techniques of solving such systems directly.

1 Gaussian Elimination

The most naive approach to solving a system $Ax = b$ is applying a series of row reductions Q such that $QAx = Ix = Qb$. In this way, the solution can be immediately read off as $x = Qb$. One may notice that these row operations Q represent the operation of an inverse for A . Indeed, this method simply computes $x = A^{-1}b$. The computational cost for such an approach is $\mathcal{O}(n^3)$ for the construction of Q and then an addition $\mathcal{O}(n^2)$ for the matrix-vector multiplication Qb . We will see that this is far from the best approach.

2 Forward/Backward Substitution

An equivalently naive approach is instead of reducing A to the identity matrix, why not simply reduce it to an upper (or lower) triangular matrix and then solve the corresponding system? That is, row operations R are performed to transform A into an upper triangular matrix T . Once in upper triangular form, the system $Tx = Rb$ can be solved in $\mathcal{O}(n^2)$. The construction, however, again takes $\mathcal{O}(n^3)$.

3 LU Factorization

A more sustainable method would be to decompose A into upper and lower triangular matrices U and L . In this way, the system $Ax = b$ becomes $LUx = b$. Accordingly, one could solve $Ly = b$ and then $Ux = y$, i.e. one forward substitution and one backward substitution. As it turns out, this takes the same amount of work as the reduction to an upper triangular matrix **and** is reusable for future $Ax = b$ solves. Furthermore, LU factorization takes less time than finding an inverse, and is more numerically stable.

To compute an LU factorization, simply reduce A to an upper triangular matrix U . For the lower triangular matrix L , the entry L_{ij} $i > j$ is simply the multiplication factor used in the Gaussian elimination process. That is, if $A_{ii} \cdot l_{ij} + A_{ij} = 0$, then $L_{ij} = -l_{ij}$. The diagonal elements are 1 and the remaining elements are computed during the elimination scheme!

4 Pivoting Strategies

Gaussian Elimination is not always stable. For instance, if $A_{ii} \approx \varepsilon_{\text{mach}}$, then a resulting decomposition could be incredibly prone to errors. To avoid this, we may want to swap rows before doing the elimination

at every step. That is, interchange rows such that the diagonal element is as large as it can be. This will reduce the instability previously mentioned. However, it no longer is true that $LU = A$. To keep up with the row changes, we must introduce a permutation matrix P such that $LU = PA$. We then simply solve $LU = Pb$ instead. This decomposition is called LU factorization with partial pivoting.

Additional Comments

As the name suggests, partial pivoting is not the only pivoting strategy, nor does it guarantee stability. However, in practice, it tends to perform perfectly fine and is faster than methods that do guarantee stability (complete pivoting). One should always consider the tradeoffs involved when deciding on algorithm. Furthermore, there exists classes of matrices (diagonally dominant, SPD, etc) that are innately resistant to numerical error in decomposition techniques and thus only require basic LU factorization.

5 Condition Number

So far, we have looked at direct methods for solving $Ax = b$. But in the presence of roundoff errors, we should not expect even a direct method to produce a perfectly accurate solution. We would like some cheap method to estimate the relative error $\frac{\|x - \hat{x}\|}{\|x\|}$. One such quantity is the residual

$$\hat{r} = b - A\hat{x}$$

It can be shown that even with a small residual, the error can still be incredibly large. That being said, how can we know how good our solution ever is? Well note that

$$\hat{r} = b - A\hat{x} = Ax - A\hat{x} = A(x - \hat{x})$$

so

$$x - \hat{x} = A^{-1}\hat{r}$$

and by Cauchy Schwarz, we have

$$\|x - \hat{x}\| \leq \|A^{-1}\| \|\hat{r}\|$$

Furthermore, since $\|b\| \leq \|A\| \|x\|$, it follows that

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq \kappa(A) \frac{\|\hat{r}\|}{\|b\|}$$

where $\kappa(A) = \|A\| \|A^{-1}\|$. In words, the relative error in the solution is bounded by the condition number of the matrix A times the relative error in the residual! Thus, for an ill-conditioned matrix, even a good residual will not guarantee a small relative error.

Another approach to error analysis is **backward** error analysis. The computed solution \hat{x} of $Ax = b$ can be viewed as the exact solution to a slightly perturbed problem $(A + \delta A)x = b + \delta b$. In this way, we see that

$$\hat{r} = b - A\hat{x} = (\delta A)\hat{x} - \delta b$$

Plugging this into the previous equation (and assuming that the perturbation is small, i.e. $\|\delta A\| < 1/\|A^{-1}\|$) we see that

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A)(\|\delta A\|/\|A\|)} \left(\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right)$$

In summary, a stable algorithm is responsible for producing a small residual. This will yield an acceptably small error in the solution if the problem is well conditioned, i.e. has a small condition number.

Chapter 5

Least Squares Problem

Last section, algorithms for solving $Ax = b$ were introduced. Now consider the following scenario: we seek a solution to $Ax = b$, but in this case, the number of rows of A are greater than the number of columns, i.e. $A \in \mathbb{R}^{m \times n}$ and $m > n$. Here, the system is overdetermined. If $b \notin \text{col}(A)$, then there exists no feasible solution.

In many disciplines, a solution to $Ax \approx b$ would still be desirable if $Ax = b$ cannot be solved. There are many ways to describe a "good" approximation, but the way we will proceed with is rather than solving $Ax = b$, we may instead solve $x = \text{argmin}_x \|Ax - b\|_2$. Noting that multiplication by a constant factor and squaring does not change the minimizer (only the minimum value), this is equivalent to

$$x^* = \text{argmin}_x \frac{1}{2} \|Ax - b\|_2^2$$

which is what we will call the least squares problem.

1 An Analytical Solution

Let $\phi(x) = \frac{1}{2} \|Ax - b\|_2^2$. Then, our least squares problem is of the form $\min \phi(x)$, an optimization formulation. From optimization, we know that since ϕ is a convex quadratic function, necessary and sufficient conditions for x^* to be a minimizer are $\nabla \phi(x^*) = 0$ and $\nabla^2 \phi(x^*) \succ 0$. It is easy to compute $\nabla \phi(x) = A^T Ax - A^T b$ and $\nabla^2 \phi(x) = A^T A$. We may assume A to be full rank (if not, then the rows of A contain redundant information) and so $x^T A^T Ax = \|Ax\|^2 > 0$ for $x \neq 0$ and so $A^T A$ is positive definite. Thus, the Hessian of ϕ is positive definite everywhere. Consequently, our sufficient and necessary condition for x^* to solve the least squares problem is

$$A^T Ax = A^T b$$

As it turns out, solving $Ax \approx b$ is equivalent to solving another linear system $A^T Ax = A^T b$. This small derivation is the backbone of many areas of data science and machine learning and also allows us to transfer any knowledge of solving linear systems to solving least squares problems.

2 Naive Numerical Solution

Hopefully by now, we know better than to simply compute $x = (A^T A)^{-1} A^T b$ and call it a day. Yet, even more stable methods such as LU factorization with partial pivoting discussed previously may have trouble solving this system. To investigate, we must look at the conditioning of this problem.

Recall that the conditioning of solving a linear system is strongly correlated with the condition number of the matrix. We need to compute $\kappa_2(A^T A)$ to fully understand the difficulties. Let $A = U \Sigma V^T$ be a singular value decomposition of A . Then

$$\kappa_2(A^T A) = \kappa_2(V \Sigma^T \Sigma V^T) = \frac{\sigma_1^2}{\sigma_2^2} = \kappa_2(A)^2$$

We see here that the conditioning of solving this linear system is the squared condition number of solving a linear system of just A . In many scenarios, A is a matrix of data and cannot be assumed to be well-conditioned. If A is poorly conditioned, i.e. $\kappa_2(A) \approx \sqrt{\varepsilon_{\text{mach}}}$, then any solution computed via $x = (A^T A)^{-1} A^T b$ will have no meaningful digits. We may rely on this approach for small systems with small condition numbers, but for larger matrices, we must find another approach.

3 Other Numerical Methods

The two suggestions here rely on QR and SVD decomposition methods. Properties and computation of these factorizations will be covered later, but we only rely on the basics here. For QR factorization, the solution via normal equations reduces to be

$$x = (A^T A)^{-1} A^T b = (R^T Q^T Q R)^{-1} R^T Q^T b = (R^T R)^{-1} R^T Q^T b = R^{-1} Q^T b$$

Now suppose $R = U \Sigma V^T$ is an SVD of R . Then $A = QR = (QU) \Sigma V^T$ is a valid SVD for A . Thus, the singular values of R are the same as those of A and therefore $\kappa_2(R) = \kappa_2(A)$. Indeed, solving the least squares problem in this manner has condition number $\kappa_2(A)$ instead of $\kappa_2(A)^2$.

Alternatively, let $A = U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T$ be a full SVD of A . Then by properties of orthogonal matrices

$$\begin{aligned} \|Ax - b\| &= \left\| U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T x \right\| \\ &= \left\| \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T x - U^T b \right\| \end{aligned}$$

Letting $U^T b = \begin{bmatrix} y \\ z \end{bmatrix}$, it follows that

$$\|Ax - b\| = \|y - \Sigma V^T x\| + \|z\|$$

Thus, in order to solve the least squares problem, it suffices to minimize $\|y - \Sigma V^T x\|$ instead. This yields system

$$\Sigma V^T x = y$$

to solve. Again, note that $I \Sigma V^T$ is a valid SVD for the matrix in question so it has the same singular values as A . Thus, $\kappa_2(A) = \kappa_2(\Sigma V^T)$. This is yet another more stable approach to solving the least squares problem. Note that these decompositions of A will take time to compute and are thus more costly operations wise but yield more stable numerical solutions.

Chapter 6

Polynomial interpolation

In this section we build up the foundation for the next few sections. Many topics in numerical analysis such as ODE's and integration rely on polynomial interpolation to work their magic. Here, we describe some basic techniques and a few results that will come in handy for future sections.

Function approximation can roughly be broken down into two categories: data fitting and approximating functions. They are different, but the distinction is subtle. **Data Fitting** is the process of finding a function that "fits" some data points. We use the term fit here loosely because there are different ways a function can fit a dataset. One such example is interpolation, i.e. the function passes through each point exactly. However, a function may fit in the least squares sense where the function is simpler, but may not interpolate exactly as is the case with linear regressions. **Approximating functions** are exactly that - functions that approximate other functions. It should be noted that the difference between these two is that the latter is identical to those of data fitting once we specify the data points.

For approximating functions, we generally assume a linear form

$$v(x) = \sum_{j=0}^n c_j \phi_j(x)$$

where $\{c_j\}_{j=0}^n$ are unknown coefficients and $\{\phi_j\}_{j=0}^n$ basis functions for the space we wish to approximate in. Note that $v(x)$ is linear wrt the basis functions and not x itself. Furthermore, we assume that the basis functions are linearly independent.

By default, we assume that the number of data points and the number of basis functions used are equal. If there are fewer basis functions than data, then the resulting system to solve for the coefficients c_j would be overdetermined, and the problem would be reduced to a least squares one previously covered.

Suppose we have data points $\{x_j, y_j\}_{j=0}^n$ and basis functions $\{\phi_j\}_{j=0}^n$. Then the coefficients can be found by solving the linear system

$$\begin{bmatrix} \phi_0(x_0) & \dots & \phi_n(x_0) \\ \vdots & \ddots & \vdots \\ \phi_0(x_n) & \dots & \phi_n(x_n) \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix}$$

Again, we have reduced our current problem (polynomial interpolation) to one that we have already solved (linear systems). Thus, the same principles from previous chapters, such as conditioning and stability, are present here.

There is one additional point to be made about interpolation. Although we make the argument above that polynomial interpolation is immediately understood via solutions of linear systems, there is also the point of application. A polynomial interpolant isn't constructed to sit and be observed. It is generally evaluated. That is, there are two steps to polynomial interpolation: construction and evaluation. Most of construction is inherited from solutions of linear systems. The next few sections will analyze different choice of basis functions for construction with evaluation being considered on a "as necessary" basis.

1 Monomial interpolation

Perhaps the simplest basis for the $n + 1$ degree polynomial space is the monomial basis, $\phi_j := x^j$. An interpolating polynomial p would then have to satisfy $p(x_i) = \sum_{j=0}^n c_j x^j$ for some c_j . Assuming (for now) that $x_i \neq x_j$ for $i \neq j$, the corresponding linear system for coefficients c_j is

$$\begin{bmatrix} 1 & \dots & x_0^n \\ \vdots & \ddots & \vdots \\ 1 & \dots & x_n^n \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix}$$

The coefficient matrix X in question here is known as Vandermonde matrix. From linear algebra, we know that

$$\det X = \prod_{i=0}^{n-1} \prod_{j=i+1}^n x_j - x_i$$

that is, under our distinct data point assumption, the determinant is non-zero and the corresponding interpolant is unique as described in the following theorem.

Theorem 6.1. *For any real data points $\{x_j, y_j\}_{j=0}^n$ with distinct abscissae x_i there exists a unique polynomial $p(x)$ of degree at most n which satisfies the interpolation conditions*

$$p(x_i) = y_i, i = 0, 1, \dots, n$$

We can summarize the monomial interpolation technique in just a few points

1. The coefficients computed may completely change if we only slightly modify the interpolation problem (more on this later)
2. The data matrix X is often ill-conditioned as n grows large or as the data points themselves spread out.
3. The construction stage requires $\mathcal{O}(n^3)$ steps, but the evaluation can be done as quickly as $\mathcal{O}(n)$ with only roughly $2n$ flops per point.

It is important to point out that the latter two disadvantages are not prevalent in small datasets. Monomial interpolation is a perfectly acceptable method for such cases.

2 Lagrange Interpolation

Monomial interpolation is quite straight-forward. It would be the first approach of any naive attempt. In fact, one of its major upsides is that it is easy to understand. Lagrange interpolation, in contrast, is not so. Nonetheless, let us proceed as intuitively as possible.

The main computational drawback of monomial interpolation is solving the linear system. What if we found a polynomial basis such that $c_j = y_j$? Then, such a representation would be easy to manipulate and significantly reduce the construction costs. These polynomial bases are the Lagrange polynomials $L_j(x)$ which satisfy

$$L_j(x_i) = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

Thus, by letting

$$p(x) = \sum_{j=0}^n y_j L_j(x)$$

we have formed our polynomial interpolant. Indeed, it satisfies the interpolation condition because

$$p(x_i) = \sum_{j=0}^n y_j L_j(x_i) = y_i$$

for any i .

With the construction stage a mere formality, what left is there to do? This is where evaluation becomes important. Let us look at the Lagrange polynomials.

$$L_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^n \frac{(x - x_i)}{(x_j - x_i)}$$

Indeed, the n roots of the polynomial must be x_i for $i \neq j$. To ensure $L_j(x_j) = 1$, we must also divide out by $x_j - x_i$ for every $i \neq j$. Thus, giving us the representation above. With all of this in hand, we are ready to do evaluation.

Notice that we can construct the denominators of the $n + 1$ Lagrange polynomials without the use of an evaluation point x . Let us compute

$$\rho_j = \prod_{i \neq j} (x_j - x_i), w_j = \frac{1}{\rho_j}$$

This requires roughly n^2 flops. We call the w_j **barycentric weights**. For evaluation, we may define the function

$$\psi(x) = \prod_{i=0}^n x - x_i$$

to obtain the interpolant

$$p(x) = \psi(x) \sum_{j=0}^n \frac{w_j y_j}{x - x_j}$$

For any given argument x , the above takes roughly $5n$ flops.

We can simplify this slightly further. Note that the function $f(x) \equiv 1$, it must be that y_j for all j . Since f is a degree 0 polynomial, it must be the degree n interpolating polynomial for any n . Thus,

$$1 = \psi(x) \sum_{j=0}^n \frac{w_j \cdot 1}{x - x_j}$$

That is, ψ can be computed using quantities that are already used. This brings us to the final representation used for evaluation

$$p(x) = \frac{\sum_{j=0}^n \frac{w_j y_j}{x - x_j}}{\sum_{j=0}^n \frac{w_j}{x - x_j}}$$

for any x .

3 Newton Polynomial Interpretation

The previous two basis functions fail to be flexible with respect to a growing dataset and also do not make it very easy to compute error in the interpolant. The Newton polynomial basis approach, however, does. We can view the Newton polynomial basis as a compromise of monomial and Lagrange: set

$$\phi_j(x) = \prod_{i=0}^{j-1} x - x_i$$

for $j = 0, 1, \dots, n$. Here we see that by construction of the basis functions, an interpolant constructed in this way is adaptive. That is, to compute the k th coefficient, only the first k data points are required. This allows us the flexibility of not having all data at once as is often the case with lab experiments.

Furthermore, the coefficients are themselves solutions of a lower triangular system. One could theoretically form such a system and solve for the coefficients, or use the following: let

$$f[x_i] = f(x_i)$$

$$f[x_i, \dots, x_j] = \frac{f[x_{i+1}, \dots, x_j] - f[x_i, \dots, x_{j-1}]}{x_j - x_i}$$

and then set $f[x_0, \dots, x_j] = c_j$. The c_j 's are known as divided differences. Note that the textbook these notes are transcribed from (and the corresponding lecturer notes) show how one can form a divided difference table recursively, but is no more than fancy bookkeeping and will not be shown here.

4 Error in Polynomial Interpolation

Let us briefly discuss possible error in the approximations. We define an error function of an interpolant $p_n(x)$ as

$$e_n(x) = f(x) - p_n(x)$$

Newton's approach allows us to cleverly compute such an error. The error at a new point x is simply the difference between the polynomial interpolant already computed evaluated at x and the polynomial interpolant that includes x . Mathematically, that is

$$f(x) = p_{n+1}(x) = p_n(x) + f[x_0, \dots, x_n, x]\phi_n(x)$$

and so

$$e_n(x) = f(x) - p_n(x) = f[x_0, \dots, x_n, x]\phi_n(x)$$

Furthermore, because the interpolating polynomial is unique, this is precisely the error no matter which basis is used. Unfortunately, this depends on the data and the evaluation point x . We continue our search for a more general error bound.

Assuming f is smooth enough, we may replace the divided differences by their corresponding derivatives to yield

$$f(x) - p_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!}\phi_n(x)$$

for some ξ . The only remaining unknowns are $\phi_n(x)$ and $f^{(n+1)}(\xi)$. By upper bounding these, we obtain the following result.

Theorem 6.2. *If p_n interpolates f at the $n+1$ points x_0, \dots, x_n and f has $n+1$ bounded derivatives on an interval $[a, b]$ containing these points, then for each $x \in [a, b]$ there exists a point $\xi = \xi(x) \in [a, b]$ such that*

$$f(x) - p_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!}\phi_n(x)$$

with error bound

$$\max_{a \leq x \leq b} |f(x) - p_n(x)| \leq \frac{1}{(n+1)!} \max_{a \leq t \leq b} |f^{(n+1)}(t)| \max_{a \leq s \leq b} \prod_{i=0}^n |s - x_i| \quad (6.1)$$

This error bound is so important that it is one of the few equations in this monograph that will earn an equation number.

Minimizing Error

Suppose that we are tasked with approximating a function, but are allowed to choose our interpolating points ourselves. Notice that from the error bound computed previously, this is our only hope at minimizing the error anyways. It suffices to choose points such that the product of distances between points is minimized. Such a choice is provided by the **Chebyshev points**. These points are defined on the interval $[-1, 1]$ by

$$x_i = \cos\left(\frac{2i+1}{2(n+1)}\pi\right)$$

for $i = 0, \dots, n$. For a general interval $[a, b]$, apply the transformation

$$x_i \leftarrow a + \frac{b-a}{2}(x_i + 1)$$

For points in the $[-1, 1]$ interval, the interpolant using the Chebyshev points satisfies¹

$$\max_{-1 \leq x \leq 1} |f(x) - p_n(x)| \leq \frac{1}{2^n(n+1)!} \max_{-1 \leq t \leq 1} |f^{(n+1)}(t)|$$

¹see appendix for Chebyshev polynomial discussion

Chapter 7

Piecewise Polynomial Interpolation

The contrast between the error of an ordinary polynomial interpolant and that of one with Chebyshev points should indicate that polynomial interpolation can be quite bad if we are not allowed to choose our points cleverly. In such a case, equation (6.1) shows us that by decreasing the size of the interval, we can also reduce the error¹. Piecewise polynomial interpolation accomplishes this by reducing the overarching problem into many smaller ones. This decomposition will also allow our interpolant to be globally flexible. That is, if one were to change a data realization y_i , only the subinterval containing the corresponding x_i will be affected. This is not the case for polynomial interpolation covered in the previous section.

That is, we may divide an interval $[a, b]$ into a smaller number of subintervals by the partition

$$a = t_0 < t_1 < \dots < t_r = b$$

and use a low degree polynomial interpolant for each subinterval. Call these interpolants $s_i(x)$. These are then patched together to form a continuous global interpolating curve which satisfies

$$v(x) = s_i(x)$$

for any $t_i \leq x \leq t_{i+1}$, $i = 0, \dots, r - 1$.

1 Broken line and piecewise Hermite interpolation

The simplest piecewise interpolation is piecewise linear, "broken line" interpolation. In this approach, the s_i 's are simply linear functions. We can even compute the global error on the interval $[a, b]$ using equation (6.1). Let $v(x)$ represent the global interpolant. For any argument $t_i \leq x \leq t_{i+1}$, the error is obtained by $f(x) - v(x) = f(x) - s_i(x)$. For a linear interpolant, this error from (6.1) is

$$f(x) - s_i(x) = \frac{f''(\xi)}{2!} (x - t_i)(x - t_{i+1})$$

The maximum value of the RHS is achieved at $x = \frac{t_i + t_{i+1}}{2}$. Letting $h = \max_{1 \leq i \leq r} t_{i+1} - t_i$, it follows that

$$|f(x) - v(x)| \leq \frac{h^2}{8} \max_{a \leq \xi \leq b} |f''(\xi)|$$

Unfortunately, only enforcing continuity forces us to give up the hopes of differentiability. In some cases, this is not so desirable. Let us look at techniques that ensure differentiability.

2 Piecewise cubic interpolation

The reason we cannot guarantee differentiability is that a linear function only has two degrees of freedom. With n interpolants, that is only $2n$ degrees of freedom. The interpolating conditions are $s_i(x_i) = f(x_i)$

¹Note that a simple rescaling of the x-axis will not work. Convince yourself why

and $s_i(x_{i+1}) = f(x_{i+1})$ for each i which is a total of $2n$ conditions. Thus, there is no room for differentiability. The natural way to alleviate this is to move to a higher degree polynomial. By far the most common choice outside of linear is cubic. By writing $s_i(x) = a_i + b_i(x - t_i) + c_i(x - t_i)^2 + d_i(x - t_i)^3$, there are a total of $4n$ degrees of freedom to enforce constraints for. The next two sections will cover ways to choose these conditions.

Piecewise cubic Hermite interpolation

Recall that the continuity conditions take up $2n$ of the $4n$ degrees of freedom available. If the values $f'(t_i)$ are provided, we can further ensure that the derivative of our interpolant $v'(x)$ is also continuous in an analogous way. This provides $2n$ more conditions which totals to our degrees of freedom. These conditions written explicitly are of course

$$s'_i(t_i) = f'(t_i) \text{ and } s'_i(t_{i+1}) = f'(t_{i+1})$$

This interpolation technique is known as Hermite cubic interpolation. One very handy property of such an interpolant is that each piece can be computed independently. That is, the approximation is completely local. Furthermore, one should expect that using more points does indeed provide us with a better error estimate. Indeed, it would do the reader well to prove the following theorem.

Theorem 7.1. *Let v interpolate f at the $n + 1$ points $x_0 < \dots < x_n$ and define $h = \max_{1 \leq i \leq n} x_i - x_{i-1}$ and assume f has as many bounded derivatives as necessary for the bounds below on an interval $[a, b]$ containing these points. Then, using a local constant, linear, or Hermite cubic interpolation, for each $x \in [a, b]$ the interpolation error is bounded by*

$$\begin{aligned} |f(x) - v(x)| &\leq \frac{h}{2} \max_{a \leq \xi \leq b} |f'(\xi)| && \text{piecewise constant} \\ |f(x) - v(x)| &\leq \frac{h^2}{8} \max_{a \leq \xi \leq b} |f''(\xi)| && \text{piecewise linear} \\ |f(x) - v(x)| &\leq \frac{h^4}{384} \max_{a \leq \xi \leq b} |f''''(\xi)| && \text{piecewise cubic Hermite} \end{aligned}$$

3 Cubic Spline Interpolation

Of course, the main disadvantage with Hermite cubics is the need for derivative values. As with everything in numerical analysis, it is best to have a methodology for when these are not available. Enter cubic splines. Our setting is again the same: approximate each interval with an interpolating cubic. However, with cubic splines, the remaining conditions are used to ensure that $v \in C^3$. That is, the conditions to be satisfied by the cubic spline are

$$\begin{aligned} s_i(x_i) &= f(x_i) && i = 0, \dots, n-1 \\ s_i(x_{i+1}) &= f(x_{i+1}) && i = 0, \dots, n-1 \\ s'_i(x_{i+1}) &= s'_{i+1}(x_{i+1}) && i = 0, \dots, n-2 \\ s''_i(x_{i+1}) &= s''_{i+1}(x_{i+1}) && i = 0, \dots, n-2 \end{aligned}$$

Note that the latter two equations only provide $2n - 2$ conditions because they only apply to the interior intersections. Thus there are even more variants that handle these remaining conditions differently.

1. The free boundary approach, giving a **natural spline**:

$$v''(x_0) = v''(x_n) = 0$$

Although popular, there is no reason to believe that this additional constraint improves the interpolation. It is widely used for its simplicity.

2. If f' is available on interval ends (such as with some boundary value problems) then the clamped boundary may be considered, specified by

$$v'(x_0) = f'(x_0), v'(x_n) = f'(x_n)$$

The interpolant here is known as the **complete spline**.

3. The third alternative is called **not-a-knot**. This approach ensures continuity of the third derivative the spline interpolant at the the nearest interior break points x_1 and x_{n-1} . That is,

$$s_0'''(x_1) = s_1'''(x_1) \text{ and } s_{n-2}'''(x_{n-1}) = s_{n-1}'''(x_{n-1})$$

Constructing the cubic spline

The trade-off of ensuring differentiability is that now each interpolant depends on those immediately next to it. The effect diminishes significantly the further out we go, but one serious downside is that the construction is not so simple anymore. Unlike before with Hermite cubics, the coefficients in the splines must be computed simultaneously. The derivation for the following algorithm is not given, but one is encouraged to try and explain it on your own. To obtain coefficients $\{a_i, b_i, c_i, d_i\}_{i=0}^{n-1}$, simply

1. Set $a_i = y_i = f(x_i)$.
2. Construct and solve a tridiagonal system of equations for the unknowns c_0, \dots, c_n using the two boundary conditions of your choosing and the equations

$$f[x_{i-1}, x_i] - h_{i-1}(2c_{i-1} + c_i) + 2h_{i-1}c_{i-1} + h_{i-1}(c_i - c_{i-1}) = f[x_i, x_{i+1}] - \frac{h_i}{3}(2c_i + c_{i+1})$$

3. Set the coefficients d_i for $i = 0, \dots, n - 1$ by

$$d_i = \frac{c_{i+1} - c_i}{3h_i}$$

4. Set the coefficients b_i for $i = 0, \dots, n - 1$ by

$$b_i = f[x_i, x_{i+1}] - \frac{h_i}{3}(2c_i + c_{i+1})$$

The corresponding interpolant for $x_i \leq x \leq x_{i+1}$ is

$$s_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3$$

Chapter 8

Numerical Differentiation

In my opinion, this is perhaps the easiest and most straight-forward section in this document. Here, we discuss the topic of numerical differentiation. Although a very routine and often easy task for early level college students, differentiation is not so trivial numerically. As a powerful tool in ODE's (a later chapter), numerical differentiation is a topic that is essential for any scientific computing text. However, most of the ingredients have already been discussed. For this reason, this section is rather short, but do not mistake this for irrelevance.

1 Taylor Series Approximations

Let f be a sufficiently smooth function with x_0 existing within its domain. We seek a cheap, accurate approximation to $f'(x_0)$, preferably through functional evaluations. As with many things in computational courses, a solid first approach is to consider a Taylor series expansion. Indeed, note that

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(\xi)$$

for some $\xi \in (x_0, x_0 + h)$. Solving for $f'(x_0)$ gives

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} + \frac{h}{2}f''(\xi) \approx \frac{f(x_0 + h) - f(x_0)}{h}$$

which will refer to as a one-sided, two point formula. We say one-sided since the function evaluations are only points larger (or smaller) than x_0 and two point because there are two function evaluations involved. This particular approximation is commonly known as the forward difference formula for $f'(x_0)$ and has truncation error $\mathcal{O}(h)$.

There is nothing particular about the Taylor series expansion chosen here, barring the fact that $f'(x_0)$ lies readily available. Indeed, we may arrive at a different approximation in a similar way using a two point, centered formula. Let

$$\begin{aligned} f(x_0 + h) &= f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(x_0) + \frac{h^3}{6}f'''(\xi_1) \\ f(x_0 - h) &= f(x_0) - hf'(x_0) + \frac{h^2}{2}f''(x_0) - \frac{h^3}{6}f'''(\xi_2) \end{aligned}$$

be two Taylor expansions. Then subtracting the second from the first and solving for $f'(x_0)$, we obtain

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} - \frac{h^2}{6}f'''(\xi)$$

Here, we used the intermediate value theorem in (??) to arrive at our error term. Note that we only used two evaluations of f , but still obtained a second order accurate approximation. A downside to this is that the method is two-sided and the points are more spread apart than those in the forward difference formula.

As long as our tolerance for long, redundant Taylor series expansions does not falter, we may continue to compute both higher order approximations as well as approximations of higher orders. However, we will not spend too much time here, only noting the possible approach.

2 Richardson Extrapolation

Rather than deal with complicated Taylor series expansions, which will unfortunately come back to us, let us take a look at another approach for generating higher order approximations. Suppose we have two order q approximation formulas for $f'(x_0)$, $g_1(x_0, h)$ and $g_2(x_0, h)$ with errors $e_1 h^q$ and $e_2 h^q$. Then for constants c_1, c_2 such that $c_1 e_1 + c_2 e_2 = 0$, the approximation

$$f'(x_0) = \frac{c_1 g_1(x_0, h) + c_2 g_2(x_0, h)}{2}$$

has error of order $\mathcal{O}(h^{q+1})$. This technique is known as Richardson extrapolation.

If the Taylor series approach was not compelling enough to convince you that the generation of highly accurate approximation formulas is quite simple, then Richardson's extrapolation should leave you satisfied. As with anything else in this text though, it does not come without a cost. Such formulas proposed here require "sufficiently smooth" functions. For high order formulas, we required many nicely bounded derivatives of f , which may not always be possible. Furthermore, approximations through Richardson extrapolation is the formulas are not compact, and often use points from a wider area than necessary.

3 Lagrange Polynomial Approximations

It would be a shame if we did all this polynomial interpolation work previously and never put it to use. Fear not, Lagrange interpolation provides us with remarkably easy to produce differentiation formulas - so much so that you may wonder why we even bothered with the first two sections. Recall that the degree n interpolating polynomial through point x_0, \dots, x_n is given by

$$p_n(x) = \sum_{i=0}^n f(x_i) L_i(x)$$

with error term $e_n(x) = \prod_{i=0}^n (x - x_i) f[x_0, \dots, x_n]$. Thus, by differentiating, we get

$$f'(x) = \sum_{i=0}^n f(x_i) L'_i(x)$$

with error term $e'_n(x)$. Evaluating at x_0 gives the almost-too-clean-to-be-comp expression

$$f'(x_0) = \sum_{i=0}^n f(x_i) L'_i(x_0)$$

with error $e'_n(x_0)$. As you may have guessed it, neither the error term nor the derivatives of the Lagrange polynomials are particularly easy to evaluate in the sense that I certainly don't want to include their expressions. However, rest assured that it can be done. Last but not least, under an equidistant point assumption, the error is $\mathcal{O}(h^n)$.

Chapter 9

Numerical Integration

Just as the previous chapter, numerical integration is also crucial in the development of numerical ODE's discussed next. Furthermore, much of this chapter is motivated by simple polynomial interpolation covered previously. However, unlike the previous section, analysis in numerical integration is not as easily derived. This comes from the simple fact that it is in general much easier to differentiate symbolically than it is to integrate. Consequently, motivation will be key here while direct computation will be left as exercises.

1 Basic Quadrature Rules

Continuing with the assumptions of smoothness asserted previously, our problem is now to find a numerical solution to $\int_a^b f(x)dx$. Similar to before, note that $\int_a^b p_n(x)dx$ approximates $\int_a^b f(x)dx$. Using our Lagrange interpolation form, we have that

$$\int_a^b f(x)dx \approx \int_a^b p_n(x)dx = \int_a^b \sum_{i=0}^n f(x_i)L_i(x)dx = \sum_{i=0}^n f(x_i) \int_a^b L_i(x)dx$$

A numerical integral of this form is called a quadrature with $a_i = \int_a^b L_i(x)dx$ and x_i called the quadrature weights and nodes, respectively. Here, we may choose our x_i freely and use our Lagrange interpolation polynomials to compute the quadrature weights. It is important to note here that once an a and b are fixed, the quadrature weights are computed once, and only once. That is, they do not depend on the integrand itself. Let us show a few examples.

Example 9.1 (Trapezoidal Rule). Set $n = 1$ and interpolate at the ends $x_0 = a$ and $x_1 = b$. Then

$$L_0(x) = \frac{x-b}{a-b}, L_1(x) = \frac{x-a}{b-a}$$

with weights

$$\int_a^b \frac{x-b}{a-b} dx = \frac{b-a}{2}$$
$$\int_a^b \frac{x-a}{b-a} dx = \frac{b-a}{2}$$

The resulting interpolant is known as the trapezoidal rule and has the following form

$$I_f \approx Q_f = \frac{b-a}{2} [f(a) + f(b)]$$

Example 9.2 (Simpson Rule). Letting $n = 2$, $x_0 = a$, $x_1 = \frac{a+b}{2}$, $x_2 = b$ gives us the Simpson rule given by

$$I_f \approx Q_f = \frac{b-a}{6} \left[f(a) + 4f\left(\frac{b+a}{2}\right) + f(b) \right]$$

A quadrature rule based on polynomial interpolation at equidistant abscissae are referred to as Newton-Cotes forms. These two examples have the special property that $f(a)$ and $f(b)$ are explicitly used in the computation. Quadrature rules satisfying this are called closed formulas. An example of an open formula is the midpoint rule given by

$$I_f \approx Q(f) = (b-a)f\left(\frac{a+b}{2}\right)$$

2 Quadrature Error

The error discussion of quadrature rules is where this section differs most drastically from the previous. Recall that the error induced by quadrature rules based on polynomial interpolation takes the form

$$E(f) = \int_a^b f(x)dx - \int_a^b p_n(x)dx = \int_a^b f(x) - p_n(x)dx = \int_a^b f[x_0, x_1, \dots, x_n, x] \prod_{i=0}^n (x - x_i)dx$$

In the past, there were simple cases where $E(f)$ could be computed at least up to a constant. Unfortunately, the best we can muster here is on a case-by-case basis. The interested reader should attempt to replicate the results below (see exercises for example). Nonetheless, we present the error bounds without proof

$$\begin{aligned} E_{\text{trap}}(f) &= \frac{f''(\eta)}{12}(b-a)^3 \\ E_{\text{simp}}(f) &= \frac{f'''(\zeta)}{90}\left(\frac{b-a}{2}\right)^5 \\ E_{\text{mid}}(f) &= \frac{f''(\eta)}{24}(b-a)^3 \end{aligned}$$

In addition to error, another quantity used to measure quadrature rules is the precision, or degree of accuracy. If a quadrature formula satisfies $E(f) = 0$ for all polynomials of degree p or less, then we say that the quadrature $Q(f)$ has precision p . For example, trapezoidal and midpoint have precision 1, and the Simpson rule has precision 3¹. If the success of the midpoint rule catches your attention, then you have a healthy amount of skepticism. Indeed, the midpoint rule only uses one function compared to the trapezoidal's two, but maintains the same precision and only slightly higher error (a constant of 2). The trapezoidal rule makes its presence known in composite numerical integration.

3 Composite Numerical Integration

In order to decrease the error of a quadrature, we have two realistic options. One is simply use a quadrature formula with a higher ordered error term. This approach tends to suffer greatly because in order to achieve such, more evaluation points are needed. Should this not bother you, it is suggested that you go and read the chapter on polynomial interpolation again. The only remaining option is to decrease the interval size. That is, apply the same technique we did with polynomial interpolation. However, here, there is no need to enforce smoothness conditions and as such the composite numerical integration formulas are a seamless transition from their single interval counterparts.

Thus, in its simplest form, divide the interval $[a, b]$ into r equal subintervals of length $h = \frac{b-a}{r}$ each. Then by additivity of integration,

$$\int_a^b f(x)dx = \sum_{i=1}^r \int_{a+(i-1)h}^{a+ih} f(x)dx = \sum_{i=1}^r \int_{t_{i-1}}^{t_i} f(x)dx$$

where $t_i = a + ih$. Consequently, the error is simply the error committed by each individual subinterval. Although what follows may have a different form than before, make no mistake, it is simply the quadrature

¹In fact, it is possible to derive all these rules independently by simply enforcing these degrees of accuracy along with the abscissae

applied to each individual subinterval. We proceed nevertheless. The composite trapezoid rule yields

$$\begin{aligned}\int_a^b f(x)dx &\approx \frac{h}{2} \sum_{i=1}^r [f(t_{i-1}) + f(t_i)] \\ &= \frac{h}{2} \left[f(a) + f(b) + 2 \sum_{i=1}^{r-1} f(t_i) \right]\end{aligned}$$

with error

$$E(f) = \sum_{i=1}^r \left(-\frac{f''(\eta_i)}{12} h^3 \right) = -\frac{f''(\eta)}{12} (b-a)h^2$$

Before deriving the composite Simpson and midpoint rules, note the beauty of closed form quadrature rules here. The function evaluations on the beginning and end of each subinterval are repeated. That is, although the trapezoidal rule requires 2 function evaluations per interval, the composite version only requires $r+1$ evaluations for r subintervals. An open point formula, such as midpoint, will not be able to reuse evaluations and is less efficient in the generalization. Indeed, we continue with both the composite Simpson and midpoint rules, respectively

$$\begin{aligned}\int_a^b f(x)dx &\approx \frac{h}{3} \left[f(a) + f(b) + 2 \sum_{k=1}^{r/2-1} f(t_{2k}) + 4 \sum_{k=1}^{r/2} f(t_{2k-1}) \right] \\ \int_a^b f(x)dx &\approx h \sum_{i=1}^r f(a + (i-1/2)h)\end{aligned}$$

You are free to compute their errors.

4 Gaussian Quadratures

Up until this point we have been considering Newton-Cotes quadrature formulas that are based on equidistant evaluation points. But we have (hopefully) learned that such a choice of abscissae can perform quite poorly. Indeed, recall that the error of a quadrature based on polynomial interpolation is

$$E(f) = \int_a^b f[x_0, x_1, \dots, x_n, x] \prod_{i=0}^n (x - x_i) dx$$

Suppose that f is a polynomial of degree $m \leq n$. Then

$$f[x_0, x_1, \dots, x_n, x] = \frac{f^{(n+1)}(\xi)}{(n+1)!} = 0$$

That is, for any $n+1$, we will generate a formula of precision n . We would optimistically hope that by choosing the $n+1$ parameters carefully, we can obtain $2n+1$ precision. To attempt this, recall that the Legendre polynomials $\{\phi_i(x)\}_{i=0}^{n+1}$ form an orthogonal basis for polynomials of degree at most n in the sense that

$$\int_a^b \phi_i(x)\phi_j(x)dx = 0, i \neq j$$

and

$$\int_a^b g(x)\phi_{n+1}(x)dx = 0$$

for any polynomial g of degree less than or equal to n . By setting x_i to be the i^{th} root of the Legendre polynomial $\phi_{n+1}(x)$, we have that

$$E(f) = \int_a^b K f[x_0, x_1, \dots, x_n, x] \phi_{n+1}(x) dx$$

for an appropriate constant K . Since for a function f of degree m , $f[x_0, x_1, \dots, x_n, x]$ is of degree $m - n - 1$, it follows that

$$E(f) = \int_a^b K f[x_0, x_1, \dots, x_n, x] \phi_{n+1}(x) dx = 0$$

for any polynomial f that satisfies $m - n - 1 \leq n \Leftrightarrow m \leq 2n + 1$. That is, this quadrature rule has precision $2n + 1$, the most we can logically expect². This is known as the Gaussian quadrature since the roots of the Legendre polynomials are called Gauss points. Note that such a procedure requires us to compute both the weights and the nodes, but, just as with the weights, the nodes only need to be computed for a single interval $[a, b]$ and degree n - not for a particular integrand.

Furthermore, it is often the case that we wish to approximate $\int_a^b f(x)w(x)dx$ for some weight function w . Rather than directly apply the above techniques, suppose we seek a solution of the form $\sum_{i=0}^n a_i f(x_i)$, i.e. the evaluations are only on f . We may follow the same procedure as the Gauss quadrature above, except the nodes must be the roots of an orthogonal polynomial basis with respect to $w(x)$. That is, if $\{\alpha_i\}_{i=0}^{n+1}$ is a basis for polynomials of degree n or less and

$$\int_a^b \alpha_i(x)\alpha_j(x)w(x)dx = 0, i \neq j$$

then letting the nodes of our quadrature formula be the roots of α_{n+1} will produce a rule with precision $2n + 1$ for approximating $\int_a^b f(x)w(x)dx$. We may use Gram-Schmidt to compute such an orthogonal basis.

5 Adaptive Quadrature

We can certainly conjure up examples where a function is wild over a subinterval $[c, d]$, but moderately mild everywhere else in $[a, b]$. In such a case, we would be required to use many subintervals in composite quadrature techniques to ensure a small error over the wild interval. But this would be significant overkill for the rest of the interval $[a, b]$. A natural solution to this is to simply only use as many subintervals as necessary to achieve a particular error goal. In order to achieve such an optimistic task, we need to know when a particular quadrature estimate is good or bad. We can do this by approximating its error.

Consider a rule with error written as $E(f) = E(f; h) = Kh^{q-1} + \mathcal{O}(h^q)$ for some constant K . Let us compute two approximations R_1 and R_2 using the same quadrature rule with h and $h/2$. The first error $I_f - R_1$ is approximately Kh^q while the second $I_f - R_2$ is $K(h/2)^q = Kh^q/2^q$. Then we may compute the errors in terms of R_1 and R_2 as

$$\begin{aligned} I_f - R_1 &\approx \frac{2^q}{2^q - 1}(R_2 - R_1) \\ I_f - R_2 &\approx \frac{1}{2^q - 1}(R_2 - R_1) \end{aligned}$$

Since R_2 is the better approximation, if $I_f - R_2 < \varepsilon$ for some specified tolerance ε , then our approximation is good. If not, then we can simply split our bad interval in 2 and recursively perform the same quadrature rule over the subintervals until they satisfy the tolerance. This allows us to use many subintervals where necessary, and few when not. There are a few subtleties here that a robust program will address³, but hopefully the idea is clear enough to understand.

6 Romberg Integration

Lastly, we present a procedure known as Romberg integration. It can be thought of as the integration version of Richardson extrapolation both in results and derivation. It can be shown (and perhaps one should) that the error for the trapezoidal rule⁴ can be written as

$$E(f; h) = K_1 h^2 + \dots + K_s h^{2s} + \mathcal{O}(h^{2s+1})$$

²It is important to point out that these nodes are chosen to maximize the precision of a rule and not necessarily the error, should it have any.

³such as the error being only an approximation, making sure to reuse function evaluations, etc

⁴Consider tracking our optimism with the trapezoidal rule throughout this chapter. It is an excellent example of the difference between mathematically feasible and computationally feasible approaches.

| $\mathcal{O}(h^2)$ | $\mathcal{O}(h^4)$ | $\mathcal{O}(h^6)$ | \dots | $\mathcal{O}(h^{2s})$ |
|--------------------|--------------------|--------------------|----------|-----------------------|
| $R_{1,1}$ | | | | |
| $R_{2,1}$ | $R_{2,2}$ | | | |
| $R_{3,1}$ | $R_{3,2}$ | $R_{3,3}$ | | |
| \vdots | \vdots | \vdots | \ddots | |
| $R_{s,1}$ | $R_{s,2}$ | $R_{s,3}$ | \dots | $R_{s,s}$ |

Figure 9.1: Romberg Table of Integration

for some constants K_i . As with Richardson extrapolation, we can also extrapolate the trapezoidal rule for higher order approximations. Consider approximations $R_{1,1}$ and $R_{2,1}$ based on $h_1 = h = b - a$ and $h_2 = \frac{h}{2}$. The errors are

$$\begin{aligned} E_{1,1} &= K_1 h^2 + K_2 h^4 + \dots \\ E_{2,1} &= K_1 (h/2)^2 + K_2 (h/2)^4 + \dots \end{aligned}$$

We can see that $E_{1,1} - 4E_{2,1} = \mathcal{O}(h^4)$. Coming back to our quadratures, it follows that

$$E_{1,1} - 4E_{2,1} = (4R_{2,1} - R_{1,1}) - 3I_f$$

Therefore,

$$I_f = \frac{4R_{2,1} - R_{1,1}}{3} + \mathcal{O}(h^4)$$

That is, $R_{2,2} = R_{2,1} + \frac{R_{2,1} - R_{1,1}}{3}$ is a $\mathcal{O}(h^4)$ accurate approximation to I_f . Not unlike Richardson extrapolation, we can continue this process to theoretically achieve an arbitrary small error. In fact, such a process is done cleanly using a Romberg table shown in Table 9.1. Note that the table can be generated each row at a time. That is, we may generate approximations in an adaptive manner. Much of the same problems persist, however. For very small h , roundoff errors start to dominate. Also, for rough functions, the lower order terms may not necessarily dominate as proposed since the constants themselves depend on high order derivatives of f . Thus, Romberg integration is most reliably applied to smooth functions with a need for high degree accuracy.

Chapter 10

Differential Equations

The materials presented in this chapter have an acquired taste. On one hand, all of the build in the previous chapters takes form in differential equations making it a particularly strong section to serve as a general review. On the other hand, much of the analysis becomes less elegant and more "left to the reader". Furthermore, there is quite a bit of content in just simple ordinary differential equations - so much so that in this chapter we often present a logical derivation or example of an idea and leave natural extensions of it to the problem sets to follow. Nonetheless, differential equations are without a doubt one of the most important tools an applied mathematician can have.

1 Euler Methods

1.1 Preliminaries

In this chapter we seek a solution to $y'(t) = f(t, y)$ for $a \leq t \leq b$. Here, we are looking for a particular function y , or rather $y(t)$ for all $t \in [a, b]$. We typically refer to the independent variable as t since a majority of differential equations are time based, but this need not be the case. Let us look at two examples.

Example 10.1. Consider the function $f(t, y) = t - y$ defined over $t \geq 0$. This gives the ODE

$$y'(t) = -y(t) + t$$

Like many numerical solutions, it is much easier to verify a solution than to find it. Indeed, please verify that this has solution $y(t) = t - 1 + \alpha e^{-t}$ for any scalar α . In order for the solution to be unique, we would require one additional piece of information to fix α . One such condition is $y(0) = c$ which enforces $\alpha = c + 1$. This type of specification is typically referred to as an initial value, or a trajectory.

For most of this chapter, we will only consider scalar ODE's simply because of their simplicity. However, it is entirely possible, and likely, that this is not the case in practice. For most systems of ODE's, they can be handled in a very similar manner as scalar ones.

Example 10.2. Consider a tiny ball of mass 1 attached at the end of a rigid, massless rod of length $r = 1$. At the other end the rod's position is fixed. Denote θ the angle between the pendulum and the negative vertical axis. A simplified model of the motion of the pendulum is governed by

$$\theta'' = -g \sin(\theta)$$

where $g = 9.81$. We can write this ODE as a first order system. Let

$$y_1(t) = \theta(t), y_2(t) = \theta'(t)$$

Then $y_1' = y_2$ and $y_2' = -g \sin(y_1)$. This defines

$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, f(t, y) = \begin{bmatrix} y_2 \\ -g \sin(y_1) \end{bmatrix}, c = \begin{bmatrix} \theta(0) \\ \theta'(0) \end{bmatrix}$$

We will use the simplest numerical method for approximately solving initial value ODE's, Euler's method, to introduce necessary terminology and core fundamentals. In later sections, where applicable, we will covered more advanced machinery as we build our theory. Consider searching for an approximate solution via equidistant abscissae defined as $t_0 = a, t_i = a + ih$ with $h = \frac{b-a}{N}$. Recall our forward difference formula for approximating a derivative

$$y'(t_i) = \frac{y(t_{i+1}) - y(t_i)}{h} - \frac{h}{2}y''(\xi)$$

From the differential equation, it follows that

$$y(t_{i+1}) = y(t_i) + hf(t_i, y(t_i)) + \frac{h^2}{2}y''(\xi_1)$$

Let us denote y_i an approximate solution to $y(t_i)$. It then makes sense to approximate y_{i+1} via

$$y_{i+1} = y_i + hf(t_i, y_i)$$

which is the celebrated forward Euler method. This simple method allows us to step through time and compute an approximation based on the previous one. Note that the initial value is necessary to do so. But who's to say that we can't replicate such a method with the backward difference formula instead? Indeed, using

$$y'(t_{i+1}) = \frac{y(t_{i+1}) - y(t_i)}{h} - \frac{h}{2}y''(\xi_2)$$

we arrive at

$$y_{i+1} = y_i + hf(t_{i+1}, y_{i+1})$$

At first glance, one may not have much an issue with this method. However, it belongs to the class of implicit equations. That is, solving for y_{i+1} , the unknown value, is not so easy anymore. The ease of computing the forward Euler is due to it being an explicit formula, while backward Euler is an implicit method. For now, it suffices to say that the implicit method has added an unnecessary headache. We will further discuss implicit methods in sections to follow.

Example 10.3. This example will need to be showcased somewhere so lets take a look at it now. Consider the differential equation $y' = \lambda y, y(0) = 1$ for some scalar λ . It can be verified that the true solution is $y = e^{\lambda t}$. Applying forward Euler we obtain

$$\begin{aligned} y_{i+1} &= y_i + hf(t_i, y_i) \\ &= y_i + h\lambda y_i \\ &= (1 + h\lambda)^{i+1}y(0) \\ &= (1 + h\lambda)^{i+1} \end{aligned}$$

We will show later that as $h \rightarrow 0$, this method converges to the true solution under a few assumptions.

1.2 Method Errors

There are two crucial types of errors for an approximation in this chapter. The first is the local truncation error d_i which is the amount by which the exact solution fails to satisfy the difference equation. The order of accuracy q is the smallest positive integer such that the local truncation error is $\mathcal{O}(h^q)$. The second is the global error e_i , the amount an approximation differs from the true solution at a given point, defined by

$$e_i = y(t_i) - y_i$$

For example, by construction of the forward Euler method, we have

$$d_i = \frac{y(t_{i+1}) - y(t_i)}{h} - f(t_i, y(t_i)) = \frac{h}{2}y''(\xi)$$

Thus, it is first order accurate with $q = 1$.

We say that a method converges if the maximum global error tends to 0 as h tends to 0. Consider the forward Euler global error. By subtracting

$$\begin{aligned}d_i &= \frac{y(t_{i+1}) - y(t_i)}{h} - f(t_i, y(t_i)) \\0 &= \frac{y_{i+1} - y_i}{h} - f(t_i, y_i)\end{aligned}$$

we obtain

$$d_i = \frac{e_{i+1} - e_i}{h} - [f(t_i, y(t_i)) - f(t_i, y_i)]$$

If $y''(t)$ is bounded by some M over $[a, b]$ and $f(t, y)$ is L -Lipschitz in y , then

$$\begin{aligned}e_{i+1} &= e_i + h[f(t_i, y(t_i)) - f(t_i, y_i)] + hd_i \\&\leq e_i + hLe_i + \frac{Mh}{2} \\&\leq \dots \leq (1 + hL)^{i+1}e_0 + \frac{Mh}{2} \sum_j^i (1 - hL)^j \\&= \frac{Mh}{2} \sum_j^i (1 - hL)^j\end{aligned}$$

which tends to 0 as $h \rightarrow 0$. Thus, under some moderate conditions, the forward Euler method is convergent.

1.3 Stability

We must make one last remark regarding stability before moving to more advanced methods. Consider the test equation ODE introduced previously $y' = \lambda y$. When $\lambda < 0$ the solution $e^{-|\lambda|t}$ decays as $t \rightarrow \infty$. Thus, it is reasonable to expect for any approximate solution y_{i+1} , we should have

$$|y_{i+1}| \leq |y_i|$$

Recall that for forward Euler applied to this ODE, we have

$$y_{i+1} = (1 + h\lambda)y_i$$

Thus, in order for our seemingly harmless condition to hold, it must be that

$$h \leq \frac{2}{\lambda}$$

which is an awfully strict condition! Any ODE that requires more stringent conditions for stability than for accuracy is called strict.

One should know by now that if the entire story of forward vs backward Euler methods were told in the previous two sections, then we would have never bothered to introduce backward Euler in the first place. Applying backward Euler to the test equation, this imposed constraint simplifies to

$$\frac{1}{1 - h\lambda} \leq 1$$

which holds for any $h > 0, \lambda < 0$. Indeed, the implicit nature of the backward Euler allows for a much more flexible selection of h . Although the test equation may appear a bit arbitrary, this analysis is quite fundamental for determining the stability of a method. We will explore more on stability later.

2 Runge-Kutta Methods

With both Euler methods being only first order accurate, there is a definite need for higher order methods. One such family is the Runge-Kutta (RK) methods. In short, consider integrating the differential equation to obtain

$$y(t_{i+1}) = y(t_i) + \int_{t_i}^{t_{i+1}} f(t, y(t)) dt$$

RK methods are built around the idea of approximating the integral using some methods derived in previous the previous chapter. For example, the implicit trapezoidal method gives

$$y_{i+1} = y_i + \frac{h}{2}(f(t_i, y_i) + f(t_{i+1}, y_{i+1}))$$

Unfortunately, as an implicit method, this has some serious drawbacks (most notably being slow). We can "remedy" these issues by introducing the explicit trapezoidal method. It is obtained by, you guessed it, approximating y_{i+1} on the right hand side of the implicit trapezoid method by an explicit method, forward Euler. Thus

$$y_{i+1} = y_i + \frac{h}{2}(f(t_i, y_i) + f(t_{i+1}, y_i + hf(t_i, y_i)))$$

gives us an explicit (order 2) method. It should come as no surprise that the explicit trapezoid rule is only conditionally stable, however¹. We will quickly present the remaining RK methods based on our favorite quadrature rules the midpoint and Simpson rules. For midpoint we have implicit midpoint given by

$$y_{i+1} = y_i + hf(t_{i+1/2}, y_{i+1/2})$$

where $t_{i+1/2} = 0.5(t_i + t_{i+1})$ and $y_{i+1/2} = 0.5(y_i + y_{i+1})$, and its explicit counterpart

$$y_{i+1} = y_i + hf(t_{i+1/2}, y_i + \frac{h}{2}f(t_i, y_i))$$

The RK method based on the Simpson rule is the most popular of the RK methods and even named RK4. It is typically written as the explicit procedure

$$\begin{aligned} Y_1 &= y_i \\ Y_2 &= y_i + \frac{h}{2}f(t_i, Y_1) \\ Y_3 &= y_i + \frac{h}{2}f(t_{i+1/2}, Y_2) \\ Y_4 &= y_i + \frac{h}{2}f(t_{i+1/2}, Y_3) \\ y_{i+1} &= y_i + \frac{h}{6}(f(t_i, Y_1) + 2f(t_{i+1/2}, Y_2) + 2f(t_{i+1/2}, Y_3) + f(t_{i+1}, Y_4)) \end{aligned}$$

¹There is no free lunch.

Part III
8610 Notes

This part is a culmination of notes from Xue's 8610 class and draw much comparison to Trefethen's Numerical Linear Algebra textbook. Similar to the 8600 part, the primary emphasis will be on motivation and intuition with detailed analysis left to homeworks and practice questions. The first two chapters were scribed while enrolled in the class and thus contain more detailed explanations/examples. However, do not draw the wrong conclusion from the brevity of the final three chapters. They are fundamental to numerical linear algebra and make up the bulk of the course.

Chapter 11

Conditioning and Stability

In this chapter, we turn to a systematic discussion of two fundamental issues of numerical analysis. **Conditioning** is the perturbation behavior of a mathematical problem. It is, loosely speaking, the sensitivity of the solution (output) to a mathematical problem with respect to the problem data (input). It is an intrinsic property of a problem. **Stability** pertains to the perturbation behavior of an algorithm used to solve that problem on a computer. Likewise, stability is the ability of an algorithm to produce a "reasonable" computed solution to a math problem under a small perturbation of input data. It is an intrinsic property of an algorithm.

1 Conditioning of a Problem

Consider the mapping $f : X \rightarrow Y$ where X, Y are normed vector spaces and $X = \{\text{all valid input data}\}, Y = \{\text{all possible solutions}\}$. Assume f is continuous and typically differentiable.

Given a problem data point x , let $y = f(x)$ be the corresponding solution. Now let the problem data change from x to $x + \Delta x$ and the solution becomes $f(x + \Delta x)$. The absolute condition number of the problem at x is

$$\kappa_a = \lim_{\delta \rightarrow 0} \sup_{\|\Delta x\| \leq \delta} \frac{\|f(x + \Delta x) - f(x)\|}{\|\Delta x\|} = \sup_{\Delta x} \frac{\|f(x + \Delta x) - f(x)\|}{\|\Delta x\|}$$

Likewise, the relative condition number is defined as

$$\kappa_r = \lim_{\delta \rightarrow 0} \sup_{\|\Delta x\| \leq \delta} \frac{\|x\| \|f(x + \Delta x) - f(x)\|}{\|f(x)\| \|\Delta x\|} = \sup_{\Delta x} \frac{\|x\| \|f(x + \Delta x) - f(x)\|}{\|f(x)\| \|\Delta x\|}$$

From now on, assume that both x and $f(x)$ are vectors and $\frac{\partial f_i}{\partial x_j}$ exist for all i, j . Then

$$f(x + \Delta x) - f(x) = J_f(x)\Delta x + \mathcal{O}(\Delta x)^2 \approx J_f(x)\Delta x$$

In this instance, the absolute condition number is

$$\kappa_a = \sup_{\Delta x} \frac{\|J_f(x)\Delta x + \mathcal{O}(\Delta x)^2\|}{\|\Delta x\|} = \sup_{\Delta x} \frac{\|J_f(x)\Delta x\|}{\|\Delta x\|} = \|J_f(x)\|$$

Similarly, one can show that the corresponding relative condition number is

$$\kappa_r = \sup_{\Delta x} \frac{\|J_f(x)\| \|x\|}{\|f(x)\|}$$

Examples

Example 11.1. Let $f(x) = \alpha x$.

$$\kappa_r = \frac{|\alpha| \|x\|}{|\alpha x|} = 1$$

So this function is well conditioned by any standard.

Example 11.2. Let $f(x) = \sqrt{x}, n \in \mathbb{N}^+, x > 0$.

$$\kappa_r = \frac{\left| \frac{1}{n} x^{1/n-1} \right| |x|}{|x^{1/n}|} = \frac{1}{n}, \text{ but } \kappa_a = \frac{1}{n} \left| x^{1/n-1} \right|$$

In this example, we see that the relative condition number is quite reasonable, but the absolute conditioning of the function is quite poor when $x \ll 1$.

Example 11.3. Let $f(x) = \tan(x)$ and $x = 10^{20}$. Then

$$\kappa_r = \frac{\sec^2(x) |x|}{|\tan(x)|} = \frac{1 + \tan^2(x)}{|\tan(x)|} |x| \geq 2 |x| = 2 \cdot 10^{20}$$

which is a terrible relative error! Conversely, $f(x) = \tan^{-1}(x)$ is perfectly conditioned.

Example 11.4. Let $f(x) = x_1 - x_2$. Then

$$\kappa_r = \frac{\|J_f(x)\| \|x\|}{\|f(x)\|} = \frac{2 \cdot \max(x_1, x_2)}{|x_1 - x_2|}$$

under the inf-norm since $J_f(x) = [1, -1]$. We see here that if $x_1 \approx x_2$ and x_1, x_2 are not close to 0, then the relative conditioning of the subtraction operation is very ill-conditioned.

Example 11.5. Let $f(x) = Ax$.

$$\kappa_r = \|A\| \frac{\|x\|}{\|Ax\|}$$

If we assume A to be square and nonsingular, then

$$\kappa_r = \|A\| \frac{\|A^{-1}Ax\|}{\|Ax\|} \leq \|A\| \|A^{-1}\|$$

Similarly, if we were to let $f(x) = A^{-1}x, \kappa_r \leq \|A^{-1}\| \|A\|$. That is to say, the sensitivity of solving $Ay = b$ for slightly perturbed y or b is bounded by the same condition number.

2 Conditioning of a System of Equations

A natural question extending from the previous section is how sensitive is $f(x) = Ax$ to changes in the coefficient matrix A ? Let y be the solution to $Ay = b$ for some fixed $b \in \mathbb{R}^n$ and $y + \Delta y$ be the solution to a slightly perturbed problem $(A + \Delta A)(y + \Delta y) = b$. From the latter, we have that

$$b = Ay + \Delta Ay + A\Delta y + \Delta A\Delta y \approx Ay + \Delta Ay + A\Delta y$$

since we can drop the double infinitesimal term in the limit. Now substituting $Ay = b$ we have that $\Delta Ay = -A\Delta y$ or that $\Delta y = -A^{-1}\Delta Ay$. Taking norms and applying Cauchy-Schwarz provides us with $\Delta y \leq \|A^{-1}\| \|\Delta A\| \|y\|$. Thus,

$$\kappa_r = \sup_{\Delta A} \frac{\|\Delta y\| / \|y\|}{\|\Delta A\| / \|A\|} \leq \|A^{-1}\| \|A\|$$

3 Conditioning of Eigenvalues of Matrices

Consider the matrix $A = \begin{bmatrix} 1 & \frac{1}{\varepsilon} \\ 0 & 1 \end{bmatrix}$. It is not difficult to show that $\lambda_1 = \lambda_2 = 1$. However, through

a small perturbation we may represent the above as $\hat{A} = \begin{bmatrix} 1 & \frac{1}{\varepsilon} \\ \varepsilon & 1 \end{bmatrix}$. In this case, we see that $\lambda_1 = 0, \lambda_2 = 2$, a significant difference. In general, for non-symmetric matrices, certain eigenvalues could be very sensitive. Let λ be a simple eigenvalue of A and v and w be the corresponding right and left eigenvectors. That is, $Av = \lambda v$ and $w^H A = \lambda w^H$. Furthermore, set E to be a small perturbation of A such that $(A + E)(v + \Delta v) = (\lambda + \Delta\lambda)(v + \Delta v)$. Then, $|\Delta\lambda| = \frac{\|E\|_2}{\cos \angle(v, w)}$. In the symmetric case, $v = w$ in direction and therefore $\kappa_a = 1$.

4 Conditioning of Roots of a Polynomial

Recall from before that to interpolate a set of points using monomial basis is terrible as it requires the solve of a dense $n \times n$ matrix that has high condition number. We will further reinforce the idea that the monomial basis should very rarely be used.

Consider the solving $x^2 - 2x + 1 = 0$. Clearly $x_1 = x_2 = 0$. Now, for the slightly perturbed problem $x^2 - 2x + 1 - \delta^2$, we have $x_1 = 1 - \delta, x_2 = 1 + \delta$ in exact arithmetic. However, if $\delta < \sqrt{\varepsilon}$, then the roots become $x_1 = x_2 = 1$ as the original problem is unchanged in double precision ($\delta^2 < \varepsilon_{mach}$). Here, a relative perturbation in one coefficient of magnitude $\mathcal{O}(\delta^2)$ produces a perturbation in the roots of magnitude $\mathcal{O}(\delta)$. So condition number of the roots are $\lim_{\delta \rightarrow 0} \frac{\mathcal{O}(\delta)}{\mathcal{O}(\delta^2)} = \infty$.

In general, if a polynomial in the monomial basis has a repeated root of multiplicity m , then a perturbation in absolute value in the polynomial coefficients of $\mathcal{O}(\delta^m)$ is enough to warrant an error of $\mathcal{O}(\delta)$ in the roots.

Example 11.6. Let's take a look at a more general example. Consider the Wilkinson Polynomial $p(x) = (x - 1)(x - 2) \dots (x - 23)(x - 24) = x^{24} + \dots + a_1x + a_0$. How sensitive are the roots to perturbations in coefficients? Consider p as a function of the coefficients **and** x . It follows by Taylor expansion that

$$\begin{aligned} 0 &= p(x_j + \Delta x_j; a_0, \dots, a_i + \Delta a_i, \dots, a_{23}) - p(x_j; a_0, \dots, a_{23}) \\ &= p(x_j; a_0, \dots, a_{23}) + \frac{\partial P}{\partial x_j} \Big|_{(x_j; a_0, \dots, a_{23})} \Delta x_j + \frac{\partial P}{\partial a_i} \Big|_{(x_j; a_0, \dots, a_{23})} \Delta a_i - p(x_j; a_0, \dots, a_{23}) \end{aligned}$$

Therefore,

$$\Delta x_j = - \frac{\frac{\partial P}{\partial a_i} \Delta a_i}{\frac{\partial P}{\partial x_j}} = - \frac{x_j^i \Delta a_i}{p'(x_j)}$$

And the relative condition is consequently

$$\kappa_r = \lim_{\Delta a_i \rightarrow 0} \frac{|\Delta x_j / x_j|}{|\Delta a_i / a_i|} = \frac{x_j^{i-1} a_i}{p'(x_j)}$$

For instance, if $i = j$, then $\kappa_r = 3.54 \cdot 10^{15}$. That is, a perturbation of machine precision on this coefficient makes it practically impossible to find the corresponding root!

5 Algorithm Stability

When solving a problem numerically, the conditioning of the problem is only half the battle. We also need to ensure that our algorithm is stable. Let $y = f(x)$ be the true solution and $\hat{y} = \hat{f}(x)$ be the computed solution by some numerical algorithm.

We naturally would like $\frac{\|f(x) - \hat{f}(x)\|}{\|f(x)\|}$ to be small. An algorithm is called accurate if this can be achieved for any valid input. However, if the problem is ill-conditioned, such an expectation is a bit ambitious. In this case, a reasonable expectation is that the algorithm is stable. That is, for each input x , there exists some Δx such that

$$\left\| \frac{\hat{f}(x) - f(x + \Delta x)}{f(x + \Delta x)} \right\| = \mathcal{O}(\varepsilon_{mach}) \text{ and } \frac{\|\Delta x\|}{\|x\|} = \mathcal{O}(\varepsilon_{mach})$$

The difference between the computed solution and the true solution is referred to as the forward error. Another stronger type of stability is called backward stability. In backward stability, we require that the computed solution is **exactly** the solution to a slightly perturbed problem. That is, $\hat{f}(x) = f(x + \Delta x)$ for some Δx satisfying $\frac{\|\Delta x\|}{\|x\|} = \mathcal{O}(\varepsilon_{mach})$.

Example 11.7. Let's look at an example of an algorithm that is not backward stable: Gaussian Elimination/LU factorization without pivoting. Let $A = \begin{bmatrix} \frac{\varepsilon}{2} & -1 \\ 1 & 1 \end{bmatrix}$. Its LU factorization in exact arithmetic is

$$A = \begin{bmatrix} 1 & 0 \\ \frac{2}{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} \frac{\varepsilon}{2} & -1 \\ 0 & 1 + \frac{2}{\varepsilon} \end{bmatrix}$$

However, when we perform the same task numerically, we obtain

$$\hat{A} = \begin{bmatrix} 1 & 0 \\ \frac{2}{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} \frac{\varepsilon}{2} & -1 \\ 0 & \frac{2}{\varepsilon} \end{bmatrix} = \begin{bmatrix} \frac{\varepsilon}{2} & -1 \\ 1 & 0 \end{bmatrix}$$

Therefore, $\Delta A = A - \hat{A} \implies \frac{\|\Delta A\|_\infty}{\|A\|_\infty} = \frac{1}{2} \gg \mathcal{O}(\varepsilon)$, i.e. the algorithm is not backward stable.

6 Stability of Linear Solvers

In this section, we will consider the stability of a few algorithms to solve systems of linear equations ($Ax = b$). We will take a look at

1. Gaussian Elimination/LU factorization with partial pivoting
2. Cramer's Rule
3. Compute A^{-1} then $x = A^{-1}b$
4. QR factorization: $x = A^{-1}b = R^{-1}Q^{-1}b = R^{-1}(Q^T b)$

Let \hat{x} be the computed solution generated by these algorithms. We know that the forward error $\frac{\|\hat{x} - x\|}{\|x\|}$ could be large if A is ill-conditioned. We will instead explore the backward error. We define the backward error as the solution to $\min_{\Delta A} \frac{\|\Delta A\|}{\|A\|}$ s.t. $(A + \Delta A)\hat{x} = b$. Our question now is how can we solve such a problem?

Theorem 11.1. *Let $r = b - A\hat{x}$ be the residual of the computed solution. Then*

$$\min_{\Delta A} \frac{\|\Delta A\|_2}{\|A\|_2} = \frac{\|r\|_2}{\|A\|_2 \|\hat{x}\|_2}$$

To complete the proof, we will need to make use of the following lemma.

Lemma 11.1. *Let $u, v \in \mathbb{R}^n$. Then $\|uv^T\|_2 = \|u\|_2 \|v\|_2$.*

Proof. Let $u, v \in \mathbb{R}^n$. To prove the lemma, first set \tilde{u}, \tilde{v} to be unit vectors in the direction of u, v respectively, i.e. $\tilde{u}\|u\|_2 = u, \tilde{v}\|v\|_2 = v$. Set U, V to be orthonormal basis extensions of u, v . Denote E to be the zero $n \times n$ matrix with a 1 in the first entry. Then $\tilde{u}\tilde{v}^T = \tilde{A} = UEV^T$. Scaling up to A , we see that $\|u\|_2 \|v\|_2$ is the only nontrivial singular value of A . Thus, it must be the 2-norm. \square

Proof of Theorem 11.1. Let $r = b - A\hat{x}$. From $(A + \Delta A)\hat{x} = b$, we get simplify to $r = \Delta A\hat{x}$. Therefore, $\|r\|_2 = \|b - A\hat{x}\|_2 \leq \|\Delta A\|_2 \|\hat{x}\|_2$. It directly follows that

$$\frac{\|r\|_2}{\|A\|_2 \|\hat{x}\|_2} \leq \frac{\|\Delta A\|_2}{\|A\|_2}$$

To achieve equality, consider the 1 rank matrix $\Delta A = \frac{r\hat{x}^T}{\hat{x}^T \hat{x}}$. By the lemma,

$$\|\Delta A\|_2 = \frac{\|r\|_2 \|\hat{x}\|_2}{\|\hat{x}\|_2^2} = \frac{\|r\|_2}{\|\hat{x}\|_2}$$

And so the relative backward error is

$$\frac{\|\Delta A\|_2}{\|A\|_2} = \frac{\|r\|_2}{\|A\|_2 \|\hat{x}\|_2}$$

which completes our proof. \square

Example 11.8. Recall the previously mentioned algorithms. Let A be the Hilbert matrix of order 10. Set $x = [1, \dots, 1]^T, b = Ax$. We now solve for \hat{x} using the above algorithms and compute $\min_{\Delta A} \frac{\|\Delta A\|_2}{\|A\|_2}$. The results are summarized in the table below

| Algorithm | Relative Backward Error | Relative Forward Error |
|------------------|-------------------------|------------------------|
| GEPP/LUPP | 10^{-16} | 3.3×10^{-4} |
| Cramer's Rule | 10^{-6} | 5.7×10^{-4} |
| $A^{-1}b$ | 10^{-5} | 10^{-3} |
| QR Factorization | 10^{-16} | 1.2×10^{-3} |

Clearly, Cramer's Rule and solving via inverses are not backward stable. However, these two methods appear to perform decently well in terms of forward error. We say that all of these algorithms are forward stable. A loose definition of forward stability is an algorithm is forward stable if it produces a forward error similar to that of the forward error produced by a backward stable algorithm!

7 Conditioning of GE/LU factorization

Continuing with our stability exploration of linear solvers, let A be a nonsingular square matrix of order n and assume that no zero pivot arises during factorization in exact arithmetic such that $A = LU$. Then for sufficiently small $\varepsilon_{\text{mach}}$, the factorization can also be completed successfully in floating point arithmetic. Furthermore, let \hat{L} and \hat{U} be the computed factors of LU decomposition. Then, it can be shown that

$$\hat{L}\hat{U} = \hat{A} = A + \Delta A$$

where ΔA satisfies $\frac{\|\Delta A\|}{\|\hat{L}\| \|\hat{U}\|} = \mathcal{O}(\varepsilon)$. Similarly, let $|A|$ be the matrix obtained by component-wise absolute value operation. It can be shown that

$$|\Delta A| \leq \frac{n\varepsilon_{\text{mach}}}{1 - n\varepsilon_{\text{mach}}} \cdot \left| \hat{L} \right| \left| \hat{U} \right|$$

holds component-wise. Furthermore, suppose these factors were used in the forward/backward substitutions of solving $Ax = b$. Then \hat{x} satisfies

$$(A + \Delta A)\hat{x} = b \text{ where } |\Delta A| \leq \frac{3n\varepsilon_{\text{mach}}}{1 - 3n\varepsilon_{\text{mach}}} \left| \hat{L} \right| \left| \hat{U} \right|$$

Thus, for backwards stability of both factorizing A and solving $Ax = b$, we need

$$\frac{\|\Delta A\|}{\|A\|} = \mathcal{O}(\varepsilon) \text{ or } |\Delta A| \leq \mathcal{O}(\varepsilon) |A|$$

Therefore, the stability is dependent on whether or not we have $\left| \left| \hat{L} \right| \right| \left| \left| \hat{U} \right| \right| \leq C_n \|A\|$ or $\left| \hat{L} \right| \left| \hat{U} \right| \leq C_n |A|$. Note that this is equivalent to C_n 'not being too large' because combining the above yields

$$\frac{|\Delta A|}{|A|} \leq \frac{n\varepsilon_{\text{mach}}}{1 - n\varepsilon_{\text{mach}}} C_n$$

To explore how large C_n would be, let $A^{(k)}$ be the intermediate matrix during factorization. Define

$$\rho_n := \frac{\max |a_{ij}^{(k)}|}{\max |a_{ij}|}$$

Example 11.9. Let's look at an example for computing ρ_n . Consider the factorization

$$A = \begin{bmatrix} \varepsilon & -1 \\ 1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} \varepsilon & -1 \\ \frac{1}{\varepsilon} & \frac{1}{\varepsilon} + 1 \end{bmatrix}$$

after one pivot. Thus,

$$L = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix}, U = \begin{bmatrix} \varepsilon & -1 \\ 0 & \frac{1}{\varepsilon} + 1 \end{bmatrix}$$

It follows that $\rho_n \leq 1 + \frac{1}{\varepsilon} = \mathcal{O}(\frac{1}{\varepsilon})$ which is very bad.

In general, we have the following result (that we present without proof) for LU factorization without pivoting.

$$\| |L| |U| \|_{\infty} \leq [1 + 2(n^2 - n)\rho_n] \|A\|_{\infty}$$

If ρ_n is small, $|L| |U|$ will also be small. For sufficiently small \hat{L}, \hat{U} satisfy a similar relation. Therefore, the magnitude of ρ_n determines the backwards stability of GE/LU. To achieve backwards stability, however, we must use pivoting to control ρ_n . At the very least, we need to ensure ρ_n depends on n only. To accomplish this, we use partial pivoting.

Let $PA = LU$ be the exact LU factorization of A with partial pivoting. Then, the computed factors satisfy $\hat{L}\hat{U} = \hat{P}A + \Delta A$ where $\frac{\|\Delta A\|}{\|A\|} = \mathcal{O}(\rho_n \varepsilon_{\text{mach}})$ and $\rho_n \leq 2^{n-1}$. In general, if GEPP is used to solve $Ax = b$, then the computed solution satisfies

$$(A + \Delta A)\hat{x} = b \text{ where } \|\Delta A\|_{\infty} \leq \frac{3n^2 \rho_n \varepsilon_{\text{mach}}}{1 - 3n \varepsilon_{\text{mach}}} \|A\|_{\infty}$$

While this may look promising at first, note that $\rho_n \leq 2^{n-1}$ is in many ways unacceptable. As n grows large, even partial pivoting cannot guarantee backward stability. However, in practice, this upper bound is far from attained. In fact, only for very constructed examples does the growth factor ever reach unacceptable levels. With this in mind, we instead say that GEPP is backwards stable for a fixed n .

The above paragraph is not to say that all matrices suffer this unfortunate fate. We conclude this section by reiterating the idea that matrices with special structure (SPD, row/diagonal dominance) do not require pivoting at all to be backwards stable.

Chapter 12

QR Factorization

We begin the next chapter with a look into one of (if not the) most important concepts in all of numerical linear algebra: QR factorization.

1 Properties of QR

Given a matrix $A \in \mathbb{R}^{m \times n}$ ($m \geq n$), the *reduced* QR factorization produces $Q_1 \in \mathbb{R}^{m \times n}$ (same size as A) and $R_1 \in \mathbb{R}^{n \times n}$ (upper triangular) such that $A = Q_1 R_1$ where Q_1 has orthonormal columns (i.e. $Q_1^T Q_1 = I_n$).

Furthermore, we can also compute a full QR factorization of the form $A = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$ where

$$\begin{aligned} [Q_1 \ Q_2]^T [Q_1 \ Q_2] &= I_m \\ Q_1^T Q_1 &= I_n \\ Q_2^T Q_2 &= I_{m-n} \\ Q_1^T Q_2 &= O_{n \times (m-n)} \end{aligned}$$

Let $A = QR$ be a reduce QR factorization of A . Multiplying each on the right by e_k , we obtain the k th columns. Thus, we have

$$a_k = Ae_k = QRe_k = \sum_{i=1}^k q_i r_{ik}$$

That is, the columns of A can be written as a linear combination of those in Q . Does the reverse direction hold? Theorem 12.1 sheds light on this topic.

Theorem 12.1. *Let $A = QR$ be a reduced QR factorization of A . If A has full column rank, then $r_{jj} \neq 0$ for all j and consequently R is nonsingular. Therefore, $AR^{-1} = Q$, i.e. the columns of Q can be written as linear combinations of those in A .*

Proof. Assume for the sake of contradiction that $r_{kk} = 0$ for some $1 \leq k \leq n$. Then R is singular. Thus the column space of QR cannot be n , contradicting the fact that A has linearly independent columns. \square

In general, if A is of rank k and the first $l \leq k$ columns of A are linearly independent, but the $(l+1)$ st column is a linear combinations of the first l columns of A , then the top left submatrix of R is nonsingular, but $r_{l+1,l+1} = 0$. Thus, q_{l+1} cannot be written as a linear combination of A_1, \dots, A_{l+1} .

Proof. For a proof, consider the proof of Theorem 12.1 restricted to only the first l columns. \square

2 QR Factorization via Gram-Schmidt

We will now consider methods in which we can compute the QR factorization. Directly enforcing the requirements of $A = QR$, we can derive $A_1 = r_{11}Q_1 \implies q_1 = A_1/r_{11}$. Continuing in this manner, we can write the k step as $Q_k = \frac{A_k - \sum_{i=1}^{k-1} Q_i r_{ik}}{r_{kk}}$. This is in fact just the Gram-Schmidt algorithm applied to the columns of A . The steps of this Gram-Schmidt Process are summarized in the Algorithm 12.1. A

Algorithm 12.1 Classical Gram-Schmidt

```

for  $k = 1, \dots, n$  do
  for  $i < k$  do
     $r_{ik} = Q_i^T A_k$ 
  end for
   $\tilde{A}_k = A_k - \sum_{i=1}^{k-1} Q_i r_{ik}$ 
   $|r_{kk}| = \|\tilde{A}_k\|_2$ 
   $Q_k = \tilde{A}_k / r_{kk}$ 
end for

```

few remarks about the algorithm,

1. Notice that we set the modulus of $|r_{kk}|$, but not the vector itself. That is, the QR factorization is unique up to scaling by a complex number with modulus 1.
2. Recall that if A does not have full column rank, then we will have some $r_{jj} = 0$. To continue in this case, we may choose any unit vector $u \notin \text{span}\{Q_1, \dots, Q_{j-1}\}$ and orthogonalize it against $\{Q_1, \dots, Q_{j-1}\}$.

This process is referred to as the Classical Gram-Schmidt algorithm. It is wildly unstable in practice. Because of this, we from now on use Modified Gram-Schmidt detailed in Algorithm 12.2. For Modified

Algorithm 12.2 Modified Gram-Schmidt

```

for  $k = 1, \dots, n$  do
   $Q_k = A_k$ 
end for
for  $k = 1, \dots, n$  do
   $r_{kk} = \|Q_k\|_2$ 
   $Q_k = Q_k / r_{kk}$ 
  for  $j = k + 1, \dots, n$  do
     $r_{kj} = Q_k^T Q_j$ 
     $Q_j = Q_j - r_{kj} Q_k$ 
  end for
end for

```

GS, we have the following result

$$\|Q^T Q - I_n\| = \mathcal{O}(\kappa_2(A)\varepsilon_{\text{mach}})$$

From this result, we see that we may not be able to guarantee that the columns of our output are orthogonal, the most important property of GS. It is for this reason that common practice is to simply perform GS twice for reliable orthogonality.

GS Computational Cost

We will now turn our attention to the computation cost of QR factorization via Gram-Schmidt. Recall that at step k of Gram-Schmidt, we compute the following

$$\begin{aligned}\tilde{A}_k &= A_k - \sum_{i=1}^{k-1} Q_i(Q_i^T A_k) \\ A_k &= \tilde{A}_k / \|\tilde{A}_k\|\end{aligned}$$

The computational cost of such is

$$[m + (m - 1) + m + m](k - 1) + 3m$$

Thus, for all k , we have

$$\begin{aligned}\text{cost} &= \sum_{k=1}^n [4m - 1](k - 1) + 3m \\ &= (4m - 1)\frac{(n - 1)(n)}{2} + 3mn \\ &\approx 2mn^2\end{aligned}$$

And again, for increased reliability, we perform this operation twice.

3 QR Factorization via Orthogonal Transformations

The recommended approach for QR factorization is by orthogonal transformations, i.e. Householder transformations and Givens rotations. We will begin with Householder. While our task may be motivated by QR, let us first consider a slightly simpler problem: given a vector x , project x onto the space spanned by e_1 . Define $v = x - \|x\|_2 e_1$ and note that

$$x - 2 \text{proj}_v(x) = x - v = \|x\|_2 e_1$$

where $\text{proj}_v(x) = \frac{v^T x}{v^T v} v$. Substituting this into the previous gives

$$x - 2v \frac{v^T x}{v^T v} = (I - 2 \frac{vv^T}{v^T v})x = \|x\|_2 e_1$$

The matrix $H(x) = (I - \frac{vv^T}{v^T v})x$ is called the Householder reflector. $H(x)$ has many useful properties such as

1. Symmetry. $H^T = H, H^* = H$
2. Involutory. $H^2 = I - 4 \frac{vv^T}{v^T v} + 4 \frac{vv^T}{v^T v} = I$
3. Orthogonal. $H^T H = H^2 = I$

Property (3) is perhaps the most important of the Householder reflector. Orthogonality plays a major role in both QR factorization and SVD decomposition. A few nice consequences of orthogonality of a matrix Q are

- Multiplication by Q does not affect norm. $\|Qv\|^2 = (Qv)^T Qv = v^T Q^T Qv = v^T v = \|v\|^2$
- Eigenvalues have modulus 1. For an eigenvalue/eigenvector pair λ, v , $\|v\| = \|Qv\| = \|\lambda v\| = \lambda \|v\| \implies |\lambda| = 1$
- Perfect conditioning. $\kappa_2(Q) = \frac{\text{largest eigenvalue}}{\text{smallest eigenvalue}} = 1$

With the help of Householder transformations, we can compute QR factorizations. Recall that $H(x)x = \|x\|_2 e_1$. Let $A \in \mathbb{R}^{m \times n}$ be a fully dense matrix with $m > n$. Then applying $H(A_1)A$ eliminates the lower triangular portion of the first column, i.e.

$$A = \begin{bmatrix} x & x & x \\ \vdots & \vdots & \vdots \\ x & x & x \end{bmatrix} \xrightarrow{Q_1 = H(A_1)A} \begin{bmatrix} x & x & x \\ 0 & x & x \\ \vdots & \vdots & \vdots \\ 0 & x & x \end{bmatrix}$$

We then apply Householder again to the first column of the first principle minor of the transformed matrix. We now have $Q_2 = \begin{bmatrix} I_1 & O_{1 \times n} \\ O_{n \times 1} & H([O_{n \times 1} I_n]) \end{bmatrix}$ and

$$A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ \vdots & \vdots & \vdots \\ 0 & x & x \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ \vdots & \vdots & \vdots \\ 0 & 0 & x \end{bmatrix}$$

We may continue in this fashion until we have transformed A into an upper triangular matrix R . Since $AQ_1Q_2 \dots Q_n = R$, R is the finished product of our transformations of A and Q is the orthogonal matrix $Q_nQ_{n-1} \dots Q_1$.

In addition to the derivation of Householder reflectors, we must also consider the practical side of using Householder transformations to compute QR factors. Below we list some of the more computational ideas associated with the algorithm.

- When applying the Householder reflectors, the Householder reflector should never be explicitly formed, e.g.

$$HU = (I - 2\frac{vv^T}{v^Tv})U = U - v(\frac{2}{v^Tv}(v^TU))$$

- Note that $Q_nQ_{n-1} \dots Q_1A = \begin{bmatrix} R \\ 0 \end{bmatrix} \implies A = Q_1 \dots Q_n \begin{bmatrix} R \\ 0 \end{bmatrix}$. Denoting Q to be the product of the individual Q_i factors, we see that $A = Q \begin{bmatrix} R \\ 0 \end{bmatrix} = [Q_L Q_R] \begin{bmatrix} R \\ 0 \end{bmatrix} = Q_L R$. Thus, for a reduced QR factorization, we only need to compute Q_L ! Observe that

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix} = Q \begin{bmatrix} I \\ 0 \end{bmatrix}$$

Therefore, to explicitly compute Q_L , simply apply Q_n, \dots, Q_1 to the left of I_n .

- In many problems, such as linear least squares, explicit computation of Q_L is not necessary! Consider solving $\min \|b - Ax\|_2$ by QR factorization where A is full rank. We have previous seen that either $x = (A^T A)^{-1} A^T b$ or $x = R^{-1} Q^T b$ would suffice. However, the later in Householder reduces to $x = R^{-1} ([I_n 0_{n \times (m-n)}] Q_n \dots Q_1 b)$. Thus, there is no need to explicitly form Q_L either.

Householder Computational Cost

As always, it is essential to know the computational cost of an algorithm. Algorithm 12.3 showcases some pseudocode of the process. We can enumerate the flops necessary for each steps as in the table below

| Step | Flops Required |
|---|--------------------|
| $\ x_k\ _2$ | $2(m - k + 1)$ |
| $\ x_k\ _2$ | 1 |
| $\ x_k\ _2 - x_k$ | $(m - k + 1)$ |
| $v_k / \ v_k\ _2$ | $3(m - k + 1)$ |
| $2v_k^T A(k : m, j)$ | $2(m - k + 1) + 1$ |
| $A(k : m, j) - v_k(2v_k^T A(k : m, j))$ | $2(m - k + 1) + 1$ |

where the latter two steps are done $n - k$ times. In total, this is roughly

$$\begin{aligned} \sum_{k=1}^n \{6(m - k + 1) + 1 + [4(m - k + 1) + 1](n - k)\} &\approx \sum_{k=1}^n 4(m - k + 1) + 1)(n - k) \\ &= 4(m - n + 1) \sum_{k=1}^n n - k + 4 \sum_{k=1}^n (n - k)^2 \\ &= 4(m - n + 1) \frac{(n - 1)n}{2} + 4 \frac{(n - 1)n(2n - 1)}{6} \\ &\approx 2mn^2 - \frac{2}{3}n^3 \end{aligned}$$

A few comments must be made here

- This is only for the R factor. We still need to do additional work to get our Q factor.
- In certain problems, like Linear Least Squares mentioned above, we can avoid doing this work.
- Note that when $m = n$, the cost of solving Linear Least Squares using QR factorization via Householder transformations costs $\frac{4}{3}n^3$ as opposed to $\frac{2}{3}n^3$ for GEPP.
- Modified GSQR is even worse since it requires $2mn^2$ flops.

Algorithm 12.3 Householder Transformation Phase 1

```

v=zeros(m,n)
for k=1,...,n do
  x_k=A(k:m,k);
  v_k=||x_k||_2 e_1 - x_k;
  v_k=v_k/||v_k||_2
  A(k:m,k)=||x_k||_2 e_1
  A(k:m,k+1:n)=A(k:m,k+1:n)-v_k 2(v_k^T(A(k:m,k+1:n)))
end for
  
```

Lastly, for QR factorization via Householder reflections, we have Theorem 12.2 regarding the stability of Householder presented without proof.

Theorem 12.2. Let \hat{Q}, \hat{R} be the computed factors of A by Householder such that $\hat{Q}\hat{R} = \hat{A} = A + \Delta A$. Then, $\frac{\|\Delta A\|}{\|A\|} = \mathcal{O}(\varepsilon_{mach})$ and $\|\hat{Q}^T \hat{Q} - I_n\| = \mathcal{O}(\varepsilon_{mach})$.

4 QR Factorization via Given's Rotation

Although Householder transformations are the defacto method of performing QR factorization, there is occasionally a more efficient method, Given's rotation. For a given vector $\begin{bmatrix} a \\ b \end{bmatrix} \in \mathbb{R}^2$, define $G = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}$, $c = \frac{a}{\sqrt{a^2+b^2}}$, $s = \frac{b}{\sqrt{a^2+b^2}}$. As a result, $G \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ca + sb \\ -sa + cb \end{bmatrix} = \begin{bmatrix} \sqrt{a^2+b^2} \\ 0 \end{bmatrix}$. If $a = b = 0$, $G = I_2$. It is easy to see that $G^T G = I$. Thus, this is an orthogonal transformation.

5 Given's Rotation Computational Cost

We can use this transformation to compute a QR factorization. We simply need to apply the Given's Rotation to all nonzero entries of the lower triangular part of A .

| Step | Flops Required |
|---|-------------------------|
| Construction of Given's Transformation | 6 |
| $G \begin{bmatrix} a \\ b \end{bmatrix}$ multiplication | 6 |
| Apply G to a row | $6 + 6(n - 1)$ |
| Apply G to entire matrix | $(6 + 6(n - 1))(m - 1)$ |

Thus, summing over all columns, we get a total cost of

$$\begin{aligned}
 6 \sum_{k=1}^{n-2} (n-k)(m-k-1) &= 6 \sum_{k=0}^{n-2} (n-k)[(n-k)(m-n-1)] \\
 &= 6 \sum_{k=0}^{n-2} (n-k)^2 + 6(m-n-1) \sum_{k=0}^{n-2} n-k \\
 &= 6 \left(\frac{n(n+1)(2n+1)}{6} - 1 \right) + 6(m-n-1) \frac{(n+2)(n-1)}{2} \\
 &\approx 2n^3 + 3(m-n)n^2 \\
 &= 3mn^2 - n^3
 \end{aligned}$$

flops. Compare this with Householder transformations. It is roughly 50% worse! However, consider the case when A is upper Hessenberg. It can be shown that only $3n^2$ flops are needed. Similarly, if A is tridiagonal, we only need roughly $14n$ flops for (phase 1) QR factorization. That is, Given's rotation actually outperforms Householder transformations matrices with special structures.

Chapter 13

Singular Value Decomposition

Like QR factorizations, singular value decompositions play an important role in numerical linear algebra. A very natural and reliable way of solving virtually any numerical linear algebra problem is simply by asking: what if we take the SVD?

1 SVD Review

Assume that $A \in \mathbb{R}^{m \times n}$ with $m \geq n$. Then a singular value decomposition (SVD) of A is $A = U\Sigma V^T$ where $U \in \mathbb{R}^{m \times m}$, $\Sigma \in \mathbb{R}^{m \times n}$, $V \in \mathbb{R}^{n \times n}$, U and V are orthogonal matrices, $\sigma_i \geq 0$. We can rewrite this as $A = \begin{bmatrix} U_L & U_R \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T = U_L \Sigma V^T$ to obtain a reduced SVD.

Theorem 13.1. *Every matrix $A \in \mathbb{R}^{m \times n}$ has a full SVD and the singular values are uniquely determined by A itself.*

Proof. If A is the zero matrix, then the result is trivially satisfied with identities. If $A \neq 0$, then $\|A\| > 0$. Let $v_1 \in \mathbb{R}^n, \|v_1\|_2 = 1$ be such that $\|Av_1\| = \frac{\|Av_1\|}{\|v_1\|} = \|A\|$. Define $\sigma_1 := \|A\| > 0$ and $u_1 = \frac{1}{\sigma_1} Av_1$. It follows that $\|u_1\| = 1$. Now suppose that $V_2 \in \mathbb{R}^{n \times (n-1)}$, $U \in \mathbb{R}^{m \times (m-1)}$ be such that $\begin{bmatrix} v_1 & V_2 \end{bmatrix}, \begin{bmatrix} u_1 & U_1 \end{bmatrix}$ are orthogonal matrices. Consider $A_1 = \begin{bmatrix} u_1 & U_2 \end{bmatrix}^T A \begin{bmatrix} v_1 & V_2 \end{bmatrix} = \begin{bmatrix} u_1^T Av_1 & u_1^T Av_2 \\ V_2^T Av_1 & V_2^T Av_2 \end{bmatrix} = \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \end{bmatrix}$. Inducting on B gives the full SVD. \square

The SVD has a few nice properties.

1. $\text{rank}(A)$ = number of non-zero singular values.
2. Assume $\text{rank}(A) = r \leq n$. Then $\text{range}(A) = \text{span}(u_1, \dots, u_r)$ and $\text{null}(A) = \text{span}(v_{r+1}, \dots, v_n)$.
3. $\|A\| = \sigma_1, \|A\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$.
4. $\lambda_i(A^T A) = \sigma_i^2(A)$.
5. If $A = A^T$, then $|\lambda_i| = \sigma_i$.
6. $|\det A| = |\det U \Sigma V^T| = |\det U| |\det \Sigma| |\det V^T| = \prod \sigma_i$.

1.1 Low Rank Approximations

For a rank r matrix A , we have that

$$A = \begin{bmatrix} u_1 & \dots & u_r \end{bmatrix} \Sigma \begin{bmatrix} v_r^T \\ \vdots \\ v_n^T \end{bmatrix}$$

Assume with loss of generality that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ and define $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ where $k \leq r$. Then we the following theorem characterizes low rank approximations of A .

Theorem 13.2. A_k is the best rank k approximation to A in the 2-norm, i.e.

$$\|A - A_k\|_2 = \inf_{B \in \mathbb{R}^{m \times n}} \|A - B\|_2 = \sigma_{k+1}$$

Proof. Suppose for the sake of contradiction that there exists some B of rank k such that $\|A - B\|_2 < \|A - A_k\|_2 = \sigma_{k+1}$. Then B has a $n - k$ dimensional null space W and so for any $w \in W$, it follows that

$$\|Aw\| = \|Aw - Bw\| = \|(A - B)w\| < \|w\| \sigma_{k+1}$$

Now consider the space V' spanned by the first $k + 1$ right singular vectors of A . For any $v \in V'$,

$$\|Av\| \geq \sigma_{k+1} \|v\|$$

Since the sum of these dimensions is greater than n , there must be some nonzero element in both, a contradiction. Thus, A_k is the minimizer that we seek. \square

2 Computing an SVD

The last topic for this short chapter is an incomplete one. We have previously discussed many ways to utilize an SVD, but no methods of computing such a form. The reason for this is simple: it is difficult to do well.

2.1 Naive Idea

Indeed, one could note that for a matrix A with SVD $A = U\Sigma V^T$,

$$A^T A = V\Sigma^2 V'$$

That is, the singular values of A are the square roots of the eigenvalues of $A^T A$. Then this problem is reduced to an eigenvalue decomposition (which will be covered in depth). Unfortunately, such a method fails to be stable. A backward stable algorithm for computing singular values would obtain $\hat{\sigma}_k$ satisfying

$$\hat{\sigma}_k = \sigma_k(A + \delta A), \frac{\|\delta A\|}{\|A\|} = \mathcal{O}(\varepsilon_{\text{mach}})$$

which, combined with previous perturbation analysis $|\sigma_k(A + \delta A) - \sigma_k(A)| \leq \|\delta A\|_2$ gives

$$|\hat{\sigma}_k - \sigma_k| = \mathcal{O}(\varepsilon_{\text{mach}} \|A\|)$$

Now finding $\lambda_k(A^T A)$ in the same fashion gives

$$|\hat{\lambda}_k - \lambda_k| = \mathcal{O}(\varepsilon_{\text{mach}} \|\mathbb{A}^T A\|) = \mathcal{O}(\varepsilon_{\text{mach}} \|A\|^2)$$

Taking square roots gives

$$|\hat{\sigma}_k - \sigma_k| = \mathcal{O}(\varepsilon_{\text{mach}} \|A\|^2 / \sigma_k)$$

which is noticeably worse by a factor of $\|A\|/\sigma_k$. That is, algorithms aside, computing eigenvalues of $A^T A$ to give singular values of A is a bad idea.

However, all hope is not lost. Consider the matrix $H = \begin{bmatrix} 0 & H^T \\ H & 0 \end{bmatrix}$. Since $A = U\Sigma V^T$ implies $AV = U\Sigma$ and $A^T U = V\Sigma$, we have

$$\begin{bmatrix} 0 & H^T \\ H & 0 \end{bmatrix} \begin{bmatrix} V & V \\ U & -U \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}$$

which is an eigenvalue decomposition of H . Thus, the singular values of A are the absolute value of the eigenvalues of H . Please convince yourself that this method is more stable. Unfortunately, the matrix H is a square matrix of dimension $m+n$. When $m \gg n$, this is not very troubling, but still a consideration. The standard SVD algorithms are based around this approach, but never form such matrix explicitly.

2.2 Golub-Kahan Bidiagonalization

The algorithm by Golub-Kahan arrives at an SVD through two steps aptly named Phase 1 and Phase 2. Phase 1 attempts to reduce the matrix into a bidiagonal form through orthogonal transformations while Phase 2 takes the bidiagonal matrix into a diagonal one through an iterative process (also via orthogonal transformations). The Phase 2 operation is very finicky and is not covered.

The Phase 1 procedure, however, is not difficult to understand at all. We seek to transform a $m \times n$ matrix A into bidiagonal form. Since an SVD is of the form $A = U\Sigma V^T$, we have the freedom of letting the orthogonal operations perform on the left and right of A be different. That is to say, we can continuously apply distinct Householder transformations to the left and right of A to zero out columns and rows. Applied to a 4×3 matrix, this looks like

$$A = \begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{bmatrix} \xrightarrow{U_1^T} \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} \xrightarrow{V_1} \begin{bmatrix} x & x & 0 \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} \xrightarrow{U_2^T} \begin{bmatrix} x & x & 0 \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \end{bmatrix} \xrightarrow{U_3^T} \begin{bmatrix} x & x & 0 \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{bmatrix}$$

At the end, we will have applied $n = 3$ reflectors on the left and $n - 2 = 1$ reflectors on the right. The total cost is effectively twice that of QR factorization since the process is just two QR schemes applied to A and A^T . Thus, the total cost is roughly $4mn^2 - \frac{4}{3}n^3$.

Part IV
Exercises

Appendix A

8600 Exercises

1 Scientific Computing Fundamentals

1.1 Numerical Algorithms

1. **Ascher and Greif problem 1.4.** Assess the conditioning of the problem of evaluating

$$g(x) = \tanh(cx) = \frac{\exp(cx) - \exp(-cx)}{\exp(cx) + \exp(-cx)}$$

near $x = 0$ as the positive parameter c grows.

Solution. Approach 1 of 2. Let $|x| \ll 1$ and $\bar{x} = 0$ be a small perturbation of x . Clearly, $g(\bar{x}) = 0$. Consider $g(x) - g(\bar{x})$. It follows that

$$g(x) - g(\bar{x}) = g(x) \approx cx - \frac{(cx)^3}{3}$$

via the Taylor expansion of $g(x)$ about $x = 0$. Since $|x| \ll 1$,

$$cx - \frac{(cx)^3}{3} \approx cx = c(x - \bar{x})$$

That is, the absolute condition number is linearly dependent on c . Thus, as c grows large, this evaluation becomes ill-conditioned.

Approach 2 of 2. We can directly evaluate $|g'(0)| = \frac{|4c|}{4} = c > 0$ at $x = 0$ (which matches our analysis from before). The relative condition number is

$$\lim_{x \rightarrow 0} \frac{xg'(x)}{g(x)} = c \cdot \lim_{x \rightarrow 0} \frac{x}{g(x)} = c \frac{1}{g'(0)} = 1$$

which completes our analysis.

2. **Ascher and Greif problem 1.5.** Consider evaluating the integral

$$y_n = \int_0^1 \frac{x^n}{x + 10} dx$$

- (a) Derive a formula for approximately computing these integrals based on evaluating y_{n-1} given y_n .

Solution. Solving the previous equation for y_{n-1} gives

$$y_{n-1} = \frac{1}{10} \left(\frac{1}{n} - y_n \right)$$

- (b) Show that for any given value $\varepsilon > 0$ and positive integer n_0 , there exists an integer $n_1 \geq n_0$ such that for taking $y_{n_1} = 0$ as a starting value will produce integral evaluations y_n with absolute error smaller than ε for all $0 < n \leq n_0$.

Solution. Fix $\varepsilon > 0$ and let $n_0 \in \mathbb{N}$. Denote the absolute error $|y_n - \hat{y}_n|$ as ξ_n . We can compute that

$$\xi_{n-1} = -\frac{1}{10}\xi_n$$

That is,

$$\xi_{n_0} = \xi_{n_1} \left(\frac{1}{10}\right)^{n_1 - n_0}$$

and we can fix n_1 accordingly.

- (c) Argue that your algorithm is stable.

Solution. The algorithm is clearly stable as the absolute error decreases in each iteration.

1.2 Roundoff Errors

1. **Ascher and Greif problem 2.10.** The function $f_1(x, \delta) = \cos(x + \delta) - \cos(x)$ can be transformed into another form $f_2(x, \delta)$ using the trigonometric formula

$$\cos(\phi) - \cos(\psi) = -2 \sin\left(\frac{\phi + \psi}{2}\right) \sin\left(\frac{\phi - \psi}{2}\right)$$

- (a) Show that, analytically, $f_1(x, \delta)/\delta$ and $f_2(x, \delta)/\delta$ are effective approximations for $-\sin(x)$ for δ sufficiently small.

Solution. Using the Taylor expansion of $\cos(x + \delta)$ at x , we see that

$$f_1(x, \delta)/\delta = (-\sin(x)) + \mathcal{O}(\delta)$$

So for small δ , $f_1(x, \delta)/\delta \approx -\sin(x)$

- (b) Derive $f_2(x, \delta)$.

Solution. Using the formula given, we compute

$$f_2(x, \delta) = -2 \sin\left(\frac{2x + \delta}{2}\right) \sin\left(\frac{\delta}{2}\right)$$

- (c) When computing both f_1 and f_2 , one notices they are not as similar as "analytically equal" would imply. Explain the difference.

Solution. There is an associated error in the evaluation of f_1 that is not evident in f_2 . The first function must compute the difference of two values close in modulus, $\cos(x + \delta)$ and $\cos(x)$. Since subtraction is an ill-conditioned problem, this introduces large relative error. However, f_2 circumvents this issue by computing $x + \delta - x = \delta$ in exact arithmetic.

2. **Ascher and Greif problem 2.13.** Consider the linear system

$$\begin{bmatrix} a & b \\ b & a \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

with $a, b > 0$

- (a) If $a \approx b$, what is the numerical difficulty in solving this system?

Solution. It is not difficult to see that $x = \frac{a}{a^2 - b^2}$ and $y = -\frac{b}{a^2 - b^2}$. When $a \approx b$, these computations produce large relative (and absolute) errors since $a^2 \approx b^2$.

- (b) Suggest a numerically stable formula for computer $z = x + y$ given a and b .

Solution. The trap solution is to simply compute x and y and add them together. However, as discussed above, both x and y will have large errors associated with them and so will z . A better solution is to note that in exact arithmetic

$$z = \frac{a}{a^2 - b^2} - \frac{b}{a^2 - b^2} = \frac{1}{a + b}$$

which is a much less error-prone operation. Big idea: when dealing with roundoff errors, try to handle the troubled operations in exact arithmetic as much as possible.

1.3 Nonlinear Equations of One Variable

1. **Ascher and Greif problem 3.1.** Apply the bisection routine to find the root of the function

$$f(x) = \sqrt{x} - 1.1$$

starting from interval $[0, 1]$ with tolerance equal to $1e - 8$.

- (a) How many iterations are required?

Solution. At least $n = 27$ are required to converge within the tolerance. This is obtained from $n = \lceil \log(\frac{b-a}{2 \cdot tol}) \rceil$.

- (b) What is the resulting absolute error? Could this absolute error be predicted by our convergence analysis?

Solution. The absolute error is $8.94 \cdot 10^{-10}$. We could have clearly predicted that it would be below $1e - 8$, otherwise we would not have stopped.

2. **Ascher and Greif problem 3.3.** Consider the fixed point iteration $x_{k+1} = g(x_k)$ and let all the assumptions of the fixed point theorem hold. Use a Taylor's series expansion to show that the order of convergence depends on how many of the derivatives of g vanish at $x = x^*$. Use your result to state how fast (at least) a fixed point iteration is expected to converge if $g'(x^*) = \dots = g^{(r)}(x^*) = 0$ where the integer $r \geq 1$ is given.

Solution. Consider the absolute error of the fixed point iteration above $e_k = |x_k - x^*|$. It follows that

$$e_{k+1} = |x_{k+1} - x^*| = |g(x_k) - g(x^*)|$$

because x^* is our fixed point. Expanding $g(x_k)$ using a Taylor series gives us

$$e_{k+1} = \left| -g(x^*) + g(x^*) + g'(x^*)(x_k - x^*) + \frac{g''(x^*)}{2!}(x_k - x^*)^2 + \dots + \frac{g^{(r+1)}(\eta)}{(r+1)!}(x_k - x^*)^r \right|$$

Then, if $g'(x^*) = \dots = g^{(r)}(x^*) = 0$, we obtain

$$e_{k+1} = \frac{g^{(r+1)}(\eta)e_k^{r+1}}{(r+1)!}$$

In general, if the first r derivatives of g at x^* are 0, then the rate of convergence is $r + 1$. This is precisely why Newton's Method has quadratic convergence.

3. **Ascher and Greif problem 3.4.** Consider the function $g(x) = x^2 + \frac{3}{16}$.

- (a) This function has two fixed points. What are they?

Solution. Solving $x = x^2 + \frac{3}{16}$ yields $x = \{1/4, 3/4\}$.

- (b) Consider the fixed point iteration $x_{k+1} = g(x_k)$ for this g . For which of the points you have found in (a) can you be sure that the iterations will converge to that fixed point? Briefly justify your answer. You may assume that the initial guess is sufficiently close to the fixed point.

Solution. Evaluating $g'(x) = 2x$ at x_1, x_2 gives us

$$\begin{aligned} g'(x_1) &= 0.5 \\ g'(x_2) &= 1.5 \end{aligned}$$

Since $g'(x_1) < 1$ and $g([a, b]) \subset [a, b]$, we know that if g is a contraction, then the fixed point iteration will converge to x_1 for sufficiently close x_0 . However, since $g'(x_2) > 1$, fixed point iteration may not necessarily converge.

- (c) For the point or points you found in (b), roughly how many iterations will be required to reduce the convergence error by a factor of 10^7 ?

Solution. Since $g'(x_1) = 0.5$, the convergence rate will be exactly the same as bisection. That is, it will take at least 4 iterations to reduce it by a factor of 10. This could be manually done by computing $\frac{-1}{\log_{10} 0.5}$

4. **Ascher and Greif problem 3.7.** Consider Steffensen's method

$$x_{k+1} = x_k - \frac{f(x_k)}{g(x_k)}$$

where

$$g(x) = \frac{f(x+f(x)) - f(x)}{f(x)}$$

(a) Show that in general, the method converges quadratically to a root of $f(x)$.

Proof. Let $F(x) = x - \frac{f(x)^2}{f(x+f(x))-f(x)}$. Then Steffensen's method is simply the fixed point iteration with $g = F$. From question 3.3, we know that if $F'(x^*) = 0$, then Steffensen's method will converge quadratically. Thus, it suffices to show that $F'(x^*) = 0$.

Via Taylor expansion, we can write $f(x+f(x)) = f(x) + f'(x)f(x) + \frac{f''(\eta)}{2}f(x)^2$. Now consider

$$F(x) = x - \frac{f(x)}{f'(x) + \frac{f''(\eta)}{2}f(x)}$$

Through elementary operations, we can reformulate this to be

$$\begin{aligned} F(x) - F(x^*) &= x - x^* - \frac{f(x) - f(x^*)}{f'(x) + \frac{f''(\eta)}{2}f(x)} \\ \frac{F(x) - F(x^*)}{x - x^*} &= 1 - \frac{f(x) - f(x^*)}{x - x^*} \cdot \frac{1}{f'(x) + \frac{f''(\eta)}{2}f(x)} \end{aligned}$$

Now, letting $x \rightarrow x^*$ gives

$$F'(x^*) = 1 - \frac{f'(x^*)}{f'(x^*)} = 0$$

which was our intent to show. □

(b) Compare the method's efficiency to the efficiency of the secant method.

Solution. First and foremost, the secant method has superlinear convergence whereas Steffensen's is quadratic. Furthermore, the secant method requires computation of $f(x_k)$ and $f(x_{k-1})$ (and thus storage of previous values). Steffensen's requires both $f(x_k)$ and $f(x_k + f(x_k))$ so roughly twice as many function evaluations. Steffensen's method is effectively replacing the secant approximation of the gradient with a more accurate one as we approach the root at the cost of doubling the function evaluations.

2 Numerical Systems Analysis

2.1 Direct Methods for Linear Systems

1. **Question from Dr. Xue.** Let $A = \begin{bmatrix} 2 & 1 & 1 & 4 \\ -3 & -1 & 3 & 2 \\ -5 & -1 & 2 & 5 \\ 4 & 2 & 3 & 1 \end{bmatrix}$ and $b = \begin{bmatrix} 4 \\ 3 \\ 8 \\ 1 \end{bmatrix}$.

(a) Solve the linear system $Ax = b$ by Gaussian elimination without pivoting.

Solution.

Begin with

$$\begin{aligned} \left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ -3 & -1 & 3 & 2 & 3 \\ -5 & -1 & 2 & 5 & 8 \\ 4 & 2 & 3 & 1 & 1 \end{array} \right] &\mapsto \left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ 0 & 0.5 & 4.5 & 8 & 9 \\ 0 & 1.5 & 4.5 & 15 & 18 \\ 0 & 0 & 1 & -7 & -7 \end{array} \right] &\mapsto \left[\begin{array}{cccc|c} 2 & 0 & -8 & -12 & -14 \\ 0 & 0.5 & 4.5 & 8 & 9 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & -7 & -7 \end{array} \right] \\ &\mapsto \left[\begin{array}{cccc|c} 2 & 0 & 0 & -4 & -6 \\ 0 & 0.5 & 0 & -3.5 & -4.5 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] &\mapsto \left[\begin{array}{cccc|c} 2 & 0 & 0 & 0 & -2 \\ 0 & 0.5 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] \end{aligned}$$

And therefore $x = [-1, 2, 0, 1]^T$.

- (b) Find the LU factorization without pivoting, then solve $Ax = b$ by two triangular solves.

Solution. Now using LU factorization,

$$\left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ -3 & 1 & 3 & 2 & 3 \\ -5 & -1 & 2 & 5 & 8 \\ 4 & 2 & 3 & 1 & 1 \end{array} \right] \mapsto \left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ 0 & 0.5 & 4.5 & 8 & 9 \\ 0 & 1.5 & 4.5 & 1 & 18 \\ 0 & 0 & 1 & -7 & -7 \end{array} \right] \mapsto \left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ 0 & 0.5 & 4.5 & 8 & 9 \\ 0 & 0 & -9 & -9 & -9 \\ 0 & 0 & 1 & -7 & -7 \end{array} \right] \mapsto \left[\begin{array}{cccc|c} 2 & 1 & 1 & 4 & 4 \\ 0 & 0.5 & 4.5 & 8 & 9 \\ 0 & 0 & -9 & -9 & -9 \\ 0 & 0 & 0 & -8 & -8 \end{array} \right] = U$$

And

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1.5 & 1 & 0 & 0 \\ -2.5 & 3 & 1 & 0 \\ 2 & 0 & -1/9 & 1 \end{bmatrix}$$

Back solving for $Ly = b$ gives $y = [4, 9, -9, -8]^T$ and consequently $Ux = y$ yields $x = [-1, 2, 0, 1]^T$ which is consistent with part (a).

2. Question from Dr. Xue.

- (a) Show that the diagonal elements of SPD matrix must be positive

Solution. Let A be a SPD $n \times n$ matrix and $x = e_k$ where $1 \leq k \leq n$. Then the argument

$$a_{kk} = e_k^T A e_k = x^T A x \geq 0$$

shows that $a_{kk} \geq 0$ for all k . That is, the diagonal elements of symmetric positive definite matrix must be positive.

- (b) Why are the elements in the L factor uniformly bounded without using pivoting?

Solution. The calculation of l_{jj} for the Cholesky factorization is as follows

$$l_{jj} = \left(a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 \right)^{1/2}$$

We can rearrange this to be

$$\sum_{k=1}^j l_{jk}^2 = a_{jj}$$

for all j, k . Therefore,

$$l_{jk} \leq \max_j \sqrt{a_{jj}}$$

So the coefficients of L are uniformly bounded.

- (c) Evaluate the arithmetic cost of Cholesky factorization of a full SPD matrix.

Solution. The arithmetic cost of Cholesky factorization can be broken into two parts: the diagonal and the off-diagonal elements. Each diagonal clearly takes $2(j-1) + 1$ operations, while off diagonal elements require $[2(j-1) + 1](n-j)$ operations for a row of L . This sum can be simplified as follows

$$\begin{aligned} \sum_{j=1}^n 2(j-1) + 1 + [2(j-1) + 1](n-j) &= \sum_{j=1}^n (2j-1)(1+n-j) \\ &= 2 \sum_{j=1}^n j(n+1-j) - \sum_{j=1}^n (n+1-j) \\ &= \frac{(2n+1)(n)(n+1)}{6} + \frac{n(n+1)}{2} - \frac{n(n+1)}{2} \\ &= \frac{(2n+1)(n)(n+1)}{6} \end{aligned}$$

After some tedious algebra (indeed, I skipped a few steps here for brevity).

3. **Ascher and Greif Problem 5.22.** Given that a and b are two real positive numbers, the eigenvalues of the symmetric tridiagonal matrix are $\lambda_j = a + 2b \cos(\frac{\pi j}{n+1})$, $j = 1, \dots, n$.

(a) Find $\|A\|_\infty$

Solution. Since the matrix A is tridiagonal, every row has the same sum except rows 1 and n . The sum on rows 1 and n are $a + b$ and all other rows have sum $a + 2b$. Because $a, b > 0$, we conclude that $\|A\|_\infty = a + 2b$.

(b) Show that if A is strictly diagonally dominant, then it is symmetric positive definite.

Solution. We know that for each $j = 1, \dots, n$, $\lambda_j = a + 2b \cos(\frac{\pi j}{n+1})$. Under the assumption that A is strictly diagonally dominant, we have $a > 2b$. Since $|\cos(x)| \leq 1$, it follows that $\lambda_j = a + 2b \cos(\frac{\pi j}{n+1}) > 0$. Since A is also clearly symmetric, A is symmetric positive definite.

(c) Suppose $a > 0$ and $b > 0$ are such that A is symmetric positive definite. Find the condition number $\kappa_2(A)$.

Solution. For a symmetric positive definite matrix, we know that

$$\kappa_2(A) = \frac{\max \lambda_i}{\min \lambda_i} = \frac{a + 2b \cos(\frac{\pi}{n+1})}{a + 2b \cos(\frac{n\pi}{n+1})}$$

Letting $n \rightarrow \infty$, we see that the numerator evaluates to $a + 2b$ and the denominator to $a - 2b$ in the limiting sense. Thus, under the assumption of symmetric positive definite, for sufficiently large n ,

$$\kappa_2(A) \approx \frac{a + 2b}{a - 2b}$$

2.2 Linear Least Squares Problems

1. **Ascher and Greif Problem 6.5.**

(a) Why can't one directly extend the LU decomposition to a long and skinny matrix in order to solve the least squares problem?

Solution. We can certainly extend LU factorization to non-square matrices. That is, for an $m \times n$ matrix A , we can write it of the form $LU = A$ where L is an $m \times n$ unit lower triangular matrix and U is an $n \times n$ upper triangular matrix. However, this decomposition does us no good. Substituting in place of A for the normal equations gives us

$$U^T L^T L U x = U^T L^T b$$

In this manner, we cannot utilize the triangularity of our matrices to arrive at a solution any faster. However, it is perfectly acceptable to use LU factorization to solve $Bx = y$ where

$B = A^T A$ and $y = A^T b$. In fact, that is precisely what the algorithm does, only using the Cholesky factorization which is simply a special case of LU factorization on symmetric positive definite matrices.

(b) When writing

$$x = (A^T A)^{-1} A^T b = (R^T Q^T Q R)^{-1} R^T Q^T b = (R^T R)^{-1} R^T Q^T b = R^{-1} Q^T b$$

we have somehow moved from the conditioning $\kappa(A)^2$ to $\kappa(A)$. Where does such a change take place?

Solution. The improvement is based on the fact that R has the same conditioning as A . That is, in the final step, solving $Rx = Q^T b$ is the same conditioning as $Ax = b$. We will show that the singular values of R and A are the same. We can prove this straight from the pseudo-inverse form of the condition number. First note that

$$A^\dagger = (A^T A)^{-1} A^T = (R^T R)^{-1} (R^T Q^T) = R^{-1} Q$$

Then by definition

$$\kappa(A) = \|A\| \|A^{-1}\| = \|QR\| \|R^{-1}Q\|$$

Since Q is orthogonal and multiplication by it does not change the norm,

$$\kappa(A) = \|QR\| \|R^{-1}Q\| = \|R\| \|R^{-1}\| = \kappa(R)$$

Additionally, we could show that A and R have the same singular values by showing they have a near-equivalent SVD decomposition. The last step in the equation shown above achieves this system by reducing $R^{-1}R$ in exact arithmetic.

2. **Ascher and Greif Problem 6.6(b).** Let Q be $m \times n$ with orthonormal columns such that $A = QR$. Show that the diagonal elements of R all satisfy $r_{ii} \neq 0$ for $i = 1, \dots, n$ if and only if A has full column rank for this economy size decomposition.

Solution. Since Q is orthogonal, it follows that $Q^T A = R$. We know that Q has full rank so the singularity of R depends directly on the column rank of A . We conclude by noting that R is singular if and only if there exists $r_{ii} = 0$ for some i .

3 Numerical Approximation

3.1 Polynomial Approximation

1. **Ascher and Greif Problem 10.4(a).** Given $n+1$ data pairs $\{x_i, y_i\}_{i=0}^n$, define for $j = 0, 1, \dots, n$ the functions $\rho_j = \prod_{i \neq j} x_j - x_i$ and let also $\phi(x) = \prod_{i=0}^n x - x_i$. Show that

$$\rho_j = \phi'(x_j)$$

Solution. The product rule gives us $\phi'(x) = \sum_{j=0}^n \prod_{i=0, i \neq j}^n x - x_i$ with $\phi'(x_j) = \prod_{i=0, i \neq j}^n x_j - x_i = \rho_j$.

2. **Ascher and Greif Problem 10.6.** Given the four data points $(-1, 1), (0, 1), (1, 2), (2, 0)$, determine the interpolating cubic polynomial using the monomial, Lagrange, and Newton basis. Show that each representation gives the same polynomial.

Solution. Using Matlab, it is quite easy to compute the closed form expression of these three representations. In the monomial basis,

$$p_M(x) = -\frac{2}{3}x^3 + \frac{1}{2}x^2 + \frac{7}{6}x + 1$$

For Lagrange we have

$$p_L(x) = -\frac{1}{6}x(x-1)(x-2) + \frac{1}{2}(x+1)(x-1)(x-2) - (x+1)(x)(x-2)$$

And finally, using Newton's basis

$$p_N(x) = -\frac{2}{3}(x+1)(x)(x-1) + \frac{1}{2}(x+1)x + 1$$

We should expect them to be the same, and the fact that two of them share a leading coefficient at a quick glance should reassure this thought.

Again, using Matlab, one could quickly evaluate that (in a slight abuse of notation)

$$\begin{aligned} p_M([-1, 0, 1, 2]) &= [1, 1, 2, 0] \\ p_L([-1, 0, 1, 2]) &= [1, 1, 2, 0] \\ p_N([-1, 0, 1, 2]) &= [1, 1, 2, 0] \end{aligned}$$

And since the degree 3 polynomials all coincide at 4 different points, they must be the same polynomial. If this feels a little off, that's because we constructed the polynomials to agree at these four points, so of course they would be the same.

3. **Ascher and Greif Problem 10.7.** For Newton basis, prove that $c_k = f[x_0, x_1, \dots, x_k]$ satisfies the recursion formula stated in the notes.

Solution. Let p_k, p'_k be polynomials that interpolate f at x_0, \dots, x_k and x_1, \dots, x_{k+1} respectively. We claim that we can write the polynomial interpolating x_0, \dots, x_{k+1} , denoted p_{k+1} , as

$$p_{k+1} = \frac{(x-x_0)p'_k + (x_{k+1}-x)p_k}{x_{k+1}-x_0}$$

Indeed, for $j = 0$,

$$p_{k+1}(x_0) = \frac{(x_0-x_0)p'_{k+1}(x_0) + (x_{k+1}-x_0)f(x_0)}{x_{k+1}-x_0} = f(x_0)$$

and for $1 < j < k$,

$$p_{k+1}(x_j) = \frac{(x_j-x_0)f(x_j) + (x_{k+1}-x_j)f(x_j)}{x_{k+1}-x_0} = f(x_j)$$

and for $j = k+1$

$$p_{k+1}(x_{k+1}) = \frac{(x_{k+1}-x_0)f(x_{k+1}) + (x_{k+1}-x_{k+1})f(x_j)}{x_{k+1}-x_0} = f(x_{k+1})$$

So p_{k+1} is the unique interpolating polynomial of $f(x)$ at x_0, \dots, x_{k+1} , and its leading coefficient is precisely $f[x_0, x_1, \dots, x_{k+1}]$. Note from the definition of p_{k+1} , we see that its leading coefficient is the leading coefficient of $\frac{p'_k - p_k}{x_{k+1} - x_0}$. That is,

$$f[x_0, \dots, x_{k+1}] = \frac{f[x_1, \dots, x_{k+1}] - f[x_0, \dots, x_k]}{x_{k+1} - x_0}$$

4. **Ascher and Greif Problem 10.13.** Let $(\hat{x}_0, \dots, \hat{x}_k)$ be a permutation of the abscissae (x_0, \dots, x_k) . Show that

$$f[\hat{x}_0, \dots, \hat{x}_k] = f[x_0, \dots, x_k]$$

Solution. Let $\hat{x}_0, \dots, \hat{x}_k$ be a permutation of x_0, \dots, x_k . Interpolate at x_0, \dots, x_k and call the interpolating polynomial $p_k(x)$. Note that no matter which permutation of our abscissa points we use, the interpolating polynomial is the same. Thus, we can write it as

$$p_k(x) = \sum_{i=0}^k f[x_0, \dots, x_i] \phi_i(x) = \sum_{i=0}^k f[\hat{x}_0, \dots, \hat{x}_i] \hat{\phi}_i(x)$$

Differentiating each expression k times leaves us with only the leading coefficient. That is,

$$p_n^{(k)}(x) = f[x_0, \dots, x_k] = f[\hat{x}_0, \dots, \hat{x}_k]$$

5. **Question from Dr. Xue.** Show that the Lagrange basis polynomials satisfy $\sum_{i=0}^n L_i(x) \equiv 1$ for all x .

Solution. Let $f(x) = 1$ and $p_n(x)$ be its degree n interpolation. Using the Lagrange basis functions, we get that

$$p_n(x) = \sum_{i=0}^n 1 \cdot L_i(x)$$

However, $f(x)$ is a polynomial of degree $0 < n$ that (clearly) interpolates these points as well, and we know that all interpolating polynomials of degree less than or equal to n are unique. That is,

$$f(x) = \sum_{i=0}^n L_i(x)$$

Another way to see this is noting by the polynomial interpolation error, we have that

$$f(x) - p_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i) = 0$$

So $f(x) \equiv p_n(x)$. Indeed, these two methods are essentially saying the same thing. That is, if f is a degree n polynomial, its degree k interpolation is itself if $k \geq n$.

3.2 Piecewise Polynomial Interpolation

1. **Ascher and Greif Problem 11.3.** Let $f \in C^3[a, b]$ be given at equidistant points $x_i = a + ih$, $i = 0, 1, \dots, n$ where $nh = b - a$. Assume further that $f'(a)$ is given as well.

- (a) Construct an algorithm for a C^1 piecewise quadratic interpolation of the given values. Thus, the interpolating function is written as

$$v(x) = s_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2, x_i \leq x \leq x_{i+1}$$

for $i = 0, 1, \dots, n-1$ and your job is to specify an algorithm for determining the $3n$ coefficients.

Solution. Proceed as follows for each i

- i. Set $a_i = f(x_i)$
 - ii. Set $b_i = b_{i-1} + 2hc_{i-1}$ if $i > 0$ or $b_0 = f'(a)$ otherwise
 - iii. Set $c_i = \frac{f(x_{i+1}) - f(x_i) - hb_{i-1} - 2h^2c_{i-1}}{h^2}$
- (b) How accurate do you expect this approximation to be as a function of h ? Justify.

Solution. The error will be of order 3 with respect to h . In fact, the error is bounded above by the error of a degree 2 interpolating polynomial. That is,

$$\text{Err}(x) = |f(x) - p(x)| \leq \frac{f'''(\xi)}{3!} \max_{x \in [t_{i-1}, t_i]} (x - t_{i-1})(x - t_i)^2$$

The maximization problem on the right is achieved at $x = \frac{1}{3}(2t_{i-1} + t_i)$, which still results in an order of 3 with respect to h when substituting this back in.

2. **Ascher and Greif Problem 11.5.** Derive the matrix problem for cubic spline interpolation with the not-a-knot condition. Show that the matrix is tridiagonal and strictly diagonally dominant.

Solution. Assume for the sake of simplicity that $h_i = h$ for all i ¹. To solve for the cubic spline coefficients, we need to solve the following system for the c_i 's

$$\begin{aligned} hc_{i-1} + 2hc_i + hc_{i+1} &= 3(f[x_i, x_{i+1}] - f[x_{i-1}, x_i]) \\ \frac{c_2 - c_1}{h} &= \frac{c_1 - c_0}{h} \\ \frac{c_n - c_{n-1}}{h} &= \frac{c_{n-1} - c_{n-2}}{h} \end{aligned}$$

¹The general case is just messier

which corresponds to the following tridiagonal, strictly diagonally dominant matrix problem after elimination

$$\begin{bmatrix} 1 & 2 & -1 & \dots & \dots & 0 \\ h & 2h & h & \dots & \dots & \vdots \\ \vdots & h & 2h & h & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & \dots & \dots & h & 2h & h \\ 0 & \dots & \dots & -1 & 2 & 1 \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ \vdots \\ c_{n-1} \\ c_n \end{bmatrix} = \begin{bmatrix} 0 \\ \psi_1 \\ \vdots \\ \vdots \\ \psi_{n-1} \\ 0 \end{bmatrix}$$

where $\psi_i = 3(f[x_i, x_{i+1}] - f[x_{i-1}, x_i])$.

3.3 Numerical Differentiation

1. **Ascher and Greif Problem 14.7.** Show that

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0 - \frac{h}{2})}{h + \frac{h}{2}}$$

decreases linearly and not faster.

Solution. Let

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(x_0) + \mathcal{O}(h^3) \quad f(x_0 - \frac{h}{2}) = f(x_0) - \frac{h}{2}f'(x_0) + \frac{h^2}{8}f''(x_0) + \mathcal{O}(h^3)$$

be the two relevant Taylor expansions. By subtracting the two and solving for $f'(x_0)$, note that the second order error expressions do not cancel. That is,

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - \frac{h}{2})}{\frac{3}{2}h} + \mathcal{O}\left(\frac{h^2}{h}\right)$$

Thus, the order is first order accurate only.

3.4 Numerical Integration

1. **Question from Dr. Xue.** Assume that a quadrature rule has both positive and negative weights, denoted by w_k with $k \in S^+ \subset \{0, 1, \dots, n\}$ and $k \in S^- = \{0, 1, \dots, n\} \setminus S^+$ respectively. Suppose that the largest positive and negative (in absolute value) weights both to infinity as $n \rightarrow \infty$. For any given $\delta > 0$ (small) and $M > 0$ (large), show that there exists n and $f, g \in C[a, b]$ with $\|f - g\|_\infty \leq \delta$ such that $|Q(f) - Q(g)| \geq M$.

Solution. Let $\varepsilon > 0, f(x) \equiv 1, Q$ be a quadrature rule satisfying the problem statement, and $\delta > 0$. Denote k^+ and k^- the indices corresponding to the largest and smallest weights of Q . Define $g(x) \equiv f(x), \forall x \in [a, b] \setminus ((x_{k^+} - \varepsilon, x_{k^+} + \varepsilon) \cap (x_{k^-} - \varepsilon, x_{k^-} + \varepsilon))$. For $x \in (x_{k^+} - \varepsilon, x_{k^+} + \varepsilon)$, let $g(x)$ be the piecewise function passing through the points $(x_{k^+} - \varepsilon, 1), (x_{k^+}, 1 + \delta)$ and $(x_{k^+} + \varepsilon, 1)$ and likewise for the interval $(x_{k^-} - \varepsilon, x_{k^-} + \varepsilon)$ Clearly, $\|f - g\|_\infty = \delta$. But

$$|Q(f) - Q(g)| = \left| \sum_{k=0}^n (f(x_k) - g(x_k))w_k \right| = \delta(w_{k^+} - w_{k^-})$$

as $\varepsilon \rightarrow 0$. Since $w_{k^+} \rightarrow \infty$ and $w_{k^-} \rightarrow -\infty$ as $n \rightarrow \infty$, we can choose n large enough s.t. $Q(f) - Q(g) \geq M$ for any $M \in \mathbb{R}$.

2. **Question from Dr. Xue.** Show that the weights of the Gauss quadrature are all positive.

Solution. Let $f(x) = L_k^2(x)$ where $L_k(x)$ is the k^{th} Lagrange interpolant through the roots of the Legendre polynomials x_0, \dots, x_n . Specifically, $L_k(x_k) = 1, L_k(x_i) = 0, i \neq k$. Since $f(x)$ is of degree $2n < 2n + 1$, the Gauss quadrature is exact. That is,

$$\int_a^b L_k^2(x) dx = \sum_j w_j L_k^2(x_j) = w_k$$

However, since $L_k^2(x) \geq 0, w_k \geq 0$ for any arbitrary k which is what we wanted to show.

3. Question from Dr. Xue.

- (a) Use the 2nd order Taylor expansion at $\frac{a+b}{2}$ to show that the error of the midpoint rule is $\frac{(b-a)^3}{24} f''(\eta)$ for some $\eta \in (a, b)$.

Solution. We begin by expanding $f(x)$ at $x = \frac{a+b}{2}$.

$$f(x) = f\left(\frac{a+b}{2}\right) + f'\left(\frac{a+b}{2}\right)\left(x - \frac{a+b}{2}\right) + \frac{f''(c)}{2}\left(x - \frac{a+b}{2}\right)^2$$

for some $c \in (a, b)$. Then the computation below follows

$$\begin{aligned} \int_a^b f(x) dx - Q(f) &= \int_a^b f\left(\frac{a+b}{2}\right) + f'\left(\frac{a+b}{2}\right)\left(x - \frac{a+b}{2}\right) + \frac{f''(c)}{2}\left(x - \frac{a+b}{2}\right)^2 dx - (b-a)f\left(\frac{a+b}{2}\right) \\ &= \int_a^b f'\left(\frac{a+b}{2}\right)\left(x - \frac{a+b}{2}\right) + \frac{f''(c)}{2}\left(x - \frac{a+b}{2}\right)^2 dx \\ &= \int_a^b \frac{f''(c)}{2}\left(x - \frac{a+b}{2}\right)^2 dx \\ &= \frac{f''(c)(b-a)^3}{24} \end{aligned}$$

which proves the result.

- (b) Let $p_1(x)$ be the linear interpolant of $f(x)$ at $x_0 = a$ and $x_1 = b$. Use the error $f(x) - p_1(x)$ to show that the error of the trapezoidal rule is $-\frac{(b-a)^3}{12} f''(\eta)$.

Solution. Per the problem, let $p_1(x)$ be the linear interpolant of $f(x)$ at $x_0 = a, x_1 = b$. Then we directly have that

$$\begin{aligned} \int_a^b f(x) dx - Q(f) &= \int_a^b f(x) dx - \frac{b-a}{2}[f(a) + f(b)] \\ &= \int_a^b f(x) - p_1(x) dx \\ &= \int_a^b \frac{f''(c)}{2}(x-a)(x-b) dx \\ &= \frac{-(b-a)^3}{12} f''(c) \end{aligned}$$

where we used that $f(x) - p_1(x) = \frac{f''(c)}{2}(x-a)(x-b)$ for some $c \in (a, b)$.

Question from Dr. Xue. Derive Simpson's rule by hand.

Solution. For a Newton Cote's rule with 3 nodes, we know that $Q(f) = \sum_{k=0}^2 f(x_k) \int_a^b L_k(x) dx$ where $x_k = \{a, \frac{a+b}{2}, b\}$. By direct computation, we have

$$\begin{aligned} L_0 &= \frac{(x - \frac{a+b}{2})(x - b)}{(a - \frac{a+b}{2})(a - b)} \\ L_1 &= \frac{(x - a)(x - b)}{(\frac{a+b}{2} - a)(\frac{a+b}{2} - b)} \end{aligned}$$

And then

$$\begin{aligned}w_0 &= \int_a^b L_0(x) dx = \frac{b-a}{6} \\w_1 &= \int_a^b L_1(x) dx = 4 \frac{(b-a)}{6} \\w_2 &= \int_a^b L_2(x) dx = \frac{b-a}{6}\end{aligned}$$

where w_2 follows by symmetry.

4. **Question from Dr. Xue.** Derive the 3-point Gauss-Legendre quadrature for $\int_{-1}^1 f(x) dx$ and the 2-point Gauss-Chebyshev quadrature for $\int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx$.

Solution. For our Gauss-Legendre quadrature, we can use our recurrence relation to obtain the Legendre polynomials. That is, set $\phi_0(x) = 1$, $\phi_1(x) = x$ and recursively define

$$\phi_{j+1} = \frac{2j+1}{j+1} x \cdot \phi_j(x) - \frac{1}{j+1} \cdot \phi_{j-1}(x)$$

This set up gives us the polynomials

$$\begin{aligned}\phi_2(x) &= \frac{3}{2}x^2 - \frac{1}{2} \\ \phi_3(x) &= \frac{5}{2}x^3 - \frac{3}{2}x\end{aligned}$$

The latter has roots $x = \{0, \pm\sqrt{3/5}\}$ which become our nodes x_i . We can then generate the Lagrange polynomials via these roots to obtain

$$\begin{aligned}L_0 &= \frac{(x-0)(x-\sqrt{3/5})}{(-\sqrt{3/5}-0)(-2\sqrt{3/5})} \\ L_1 &= \frac{-(x^2-3/5)5}{3}\end{aligned}$$

Again, because of the symmetry of $[-1, 1]$, we need not compute L_2 . Integrating each of these over $[-1, 1]$ gives us

$$\begin{aligned}w_0 &= \frac{5}{9} \\ w_1 &= \frac{8}{9} \\ w_2 &= w_0\end{aligned}$$

Thus, our quadrature rule is

$$Q(f) = \frac{5}{9}f(-\sqrt{3/5}) + \frac{8}{9}f(0) + \frac{5}{9}f(\sqrt{3/5})$$

Our Gauss-Chebyshev quadrature is not as easy. To find its nodes, we set $w(x) = \frac{1}{\sqrt{1-x^2}}$ and find the orthogonal polynomials with respect to inner product $\langle f, g \rangle = \int_{-1}^1 f(x)g(x)w(x)dx$. Through not difficult, but tedious, Gram-Schmidt computation on $\{1, x, x^2\}$, we arrive at

$$\phi_2(x) = x^2 - 1/2$$

This gives us nodes $x = \{-1/\sqrt{2}, 1/\sqrt{2}\}$ ². Luckily, our weight function is also symmetric so we simply need to solve for w_0 and we are done. A few more computations show us that

$$\int_{-1}^1 \frac{x - 1/\sqrt{2}}{-2/\sqrt{2}} w(x) dx = \frac{\pi}{2}$$

Thus, $w_0 = w_1 = \frac{\pi}{2}$ and our Gauss-Chebyshev quadrature is

$$Q(f) = \frac{\pi}{2}[f(-1/\sqrt{2}) + f(1/\sqrt{2})]$$

5. **Ascher and Greif Problem 15.9.** Derive and use a 2-node weighted Gauss quadrature to integrate $\int_0^1 \frac{e^x}{\sqrt{x}} dx$.

Solution. This problem is not too much different than the previous Gauss-Chebyshev one. For this reason, I will not delineate my steps as well and will just provide computer-aided values. We begin by defining the inner product

$$\langle f, g \rangle = \int_0^1 f(x)g(x)w(x)dx$$

where $w(x) = \frac{1}{\sqrt{x}}$. We can then use Gram-Schmidt again on $\{1, x, x^2\}$ to arrive at the following set of orthogonal set of polynomials w.r.t $\langle \cdot, \cdot \rangle$

$$\begin{aligned}\phi_0(x) &= 1 \\ \phi_1(x) &= x - \frac{1}{3} \\ \phi_2(x) &= x^2 - \frac{6}{7}x + \frac{3}{35}\end{aligned}$$

The roots of $\phi_2(x)$ are computed numerically to be $x = \{0.115587, 0.741557\}$. This leads to the following Lagrange polynomials

$$\begin{aligned}L_0(x) &= \frac{x - x_1}{x_0 - x_1} \\ L_1(x) &= \frac{x - x_0}{x_1 - x_0}\end{aligned}$$

In this case, the weight function is not symmetric about $[0, 1]$ so we must compute each weight individually. Doing so results in

$$\begin{aligned}w_0 &= \int_0^1 L_0(x)dx = 1.3043 \\ w_1 &= \int_0^1 L_1(x)dx = 0.6957\end{aligned}$$

Indeed, computing the quadrature of $\frac{e^x}{\sqrt{x}}$ using the rule defined by x_0, x_1, w_0, w_1 gives

$$I_f = \int_0^1 \frac{e^x}{\sqrt{x}} dx \approx w_0 x_0 + w_1 x_1 = 2.924539758$$

which has 3 correct decimal digits at only 2 function evaluations.

²If one were a bit smarter than me, one could reason that the nodes of this quadrature must be the Chebyshev points for $n = 1$ in advance. Unfortunately, I had to find out the hard way

Appendix B

8610 Exercises

1 Numerical Linear Algebra Fundamentals

1. **Trefethen and Bau Problem 2.1.** Show that if a matrix A is both triangular and unitary, then it is diagonal.

Solution. Suppose A is upper (lower) triangular. Then, A^{-1} is also upper (lower) triangular. This can be shown through elementary row reductions on A . Since A is unitary, we conclude that $A^{-1} = A^T$ is upper (lower) triangular. This is only possible if A is diagonal.

2. **Trefethen and Bau Problem 2.2.** The Pythagorean theorem asserts that for a set of n orthogonal vectors $\{x_i\}$,

$$\left\| \sum_{i=1}^n x_i \right\|^2 = \sum_{i=1}^n \|x_i\|^2$$

- (a) Prove this in the case $n = 2$.

Solution. We can directly compute

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2 + 2\langle y, x \rangle.$$

Since $\langle y, x \rangle = 0$, by orthogonality, the result follows immediately.

- (b) Prove the general case.

Solution. The previous question provides us a base case for an induction argument. Assume the statement holds for $n - 1$. Then it follows that

$$\begin{aligned} \|(x_1 + x_2 + \cdots + x_{n-1}) + x_n\|^2 &= \|x_n\|^2 + \|(x_1 + \cdots + x_{n-1})\|^2 \\ &= \|x_n\|^2 + \|x_1\|^2 + \cdots + \|x_{n-1}\|^2 \end{aligned}$$

by the induction hypothesis.

3. **Trefethen and Bau Problem 2.3.** Let $A \in \mathbb{C}^{m \times m}$ be hermitian.

- (a) Prove that all eigenvalues of A are real.

Solution. Let (λ, x) be an eigenvalue/eigenvector pair. Then

$$\begin{aligned} \bar{\lambda}x^*x &= (Ax)^*x \\ &= x^*A^*x \\ &= x^*Ax \\ &= \lambda x^*x. \end{aligned}$$

Hence, $\lambda = \bar{\lambda}$ and λ must be real.

- (b) Prove that if x and y are eigenvectors corresponding to distinct eigenvalues, then x and y are orthogonal.

Solution. Notice that for eigenvectors x, y corresponding to eigenvalues $\lambda_1 \neq \lambda_2$, we have the following

$$\bar{\lambda}_1 \lambda_2 \langle y, x \rangle = x^* A^* A y = \lambda_2 x^* A y = \lambda_2^2 \langle y, x \rangle$$

from which we conclude that either $\langle y, x \rangle = 0$, or $\lambda_1 = \bar{\lambda}_1 = \lambda_2$. Here we used the fact that each eigenvalue of A is real proved previously. The latter case implies $\lambda_1 = \lambda_2$ which we assumed to be false. Thus, $\langle y, x \rangle = 0$.

4. **Trefethen and Bau Problem 2.4.** What can be said about the eigenvalues of a unitary matrix?

Solution. Let (λ, x) be an eigenvalue/eigenvector pair of a unitary matrix A . Then

$$\begin{aligned} x &= A^* A x \\ \lambda A^* x & \\ \lambda \bar{\lambda} x & \end{aligned}$$

the adjoint operation preserves eigenvectors, but conjugates eigenvalues. Thus, the modulus of eigenvalues of a unitary matrix are either 0 or 1.

5. **Trefethen and Bau Problem 2.5.** Let $S \in \mathbb{C}^{m \times m}$ be skew-hermitian.

- (a) Show that the eigenvalues of S are purely imaginary.

Solution. Follow the previous argument for hermitian matrices.

- (b) Show that $I - S$ is nonsingular.

Solution. Suppose $I - S$ were singular. Then there exists some x such that $Sx = x$. That is, 1 is an eigenvalue of S . We previously showed that all eigenvalues are purely imaginary, thus $I - S$ must be nonsingular.

- (c) Show that the matrix $Q = (I - S)^{-1}(I + S)$ is unitary.

Solution. We simply need to verify the adjoint is the inverse. We compute

$$\begin{aligned} (I - S)^{-1}(I + S)((I - S)^{-1}(I + S))^* &= (I - S)^{-1}(I + S)(I + S)^*(I - S)^{-*} \\ &= (I - S)^{-1}(I + S)(I - S)(I + S)^{-1} \\ &= (I - S)^{-1}(I - S^2)(I + S)^{-1} \\ &= (I - S)^{-1}(I - S)(I + S)(I + S)^{-1} \\ &= I \end{aligned}$$

to conclude the result.

6. **Trefethen and Bau Problem 2.6.** Show that if $A = I + uv^*$ is nonsingular, then its inverse has the form $I + \alpha uv^*$ for vectors $u, v \in \mathbb{C}^m$.

Solution. We can directly compute

$$\begin{aligned} (I + uv^*)(I + \alpha uv^*) &= I + (\alpha + 1)uv^* + \alpha u(v^*u)v^* \\ &= I + (\alpha + 1 + \alpha v^*u)uv^*. \end{aligned}$$

It is clear that we choose α such that

$$\alpha + 1 + \alpha v^*u = 0 \Leftrightarrow \alpha = -1/(1 + v^*u)$$

for $v^*u \neq -1$, then $A^{-1} = I + \alpha uv^*$. If $v^*u = -1$, then note that

$$Au = Iu + uv^*u = u - u = 0.$$

Hence, A is singular and $\text{span}(u) \in \text{null}(A)$. Since uv^* is rank-one, $\dim(\text{null}(A)) = 1$ and so $\text{span}(u) = \text{null}(A)$.

7. **Trefethen and Bau Problem 3.1.** Prove that if W is an arbitrary nonsingular matrix, the function $\|\cdot\|_W := \|Wx\|$ is a vector norm.

Solution. Let $\alpha \in \mathbb{R}$ and $x, y \in \mathbb{C}^m$. We verify the following

$$\begin{aligned}\|\alpha x\|_W &= \|Wx\alpha\| = |\alpha| \|Wx\| = |\alpha| \|x\|_W \\ \|x+y\|_W &= \|W(x+y)\| \leq \|Wx\| + \|Wy\| = \|x\|_W + \|y\|_W \\ \|x\|_W &= \|Wx\| \geq 0\end{aligned}$$

and note that if $\|Wx\| = 0$, then $Wx = 0$ and $x = 0$ by nonsingularity of W .

8. **Trefethen and Bau Problem 3.2.** Let $\|\cdot\|$ be any norm. Show that $\rho(A) \leq \|A\|$ where $\rho(A)$ is the spectral radius of A and denotes the largest eigenvector of A .

Solution. Let (λ, v) be an eigenvalue/eigenvector pair of A . Noting $Ax = \lambda x$, by taking norms, we have that

$$|\lambda| \|x\| = \|Ax\| \leq \|A\| \|x\|.$$

Dividing by $\|x\|$, we conclude that $|\lambda| \leq \|A\|$ for any eigenvalue λ . Thus, $\rho(A) \leq \|A\|$.

9. **Trefethen and Bau Problem 3.4.** Let B be a submatrix of a $m \times n$ matrix A that is a $\mu \times \nu$ matrix obtained by selecting selecting certain rows and columns of A . Show that $\|B\|_p \leq \|A\|_p$ for any $1 \leq p \leq \infty$.

Solution. Let $\{c_i\}$ be the set of indices that correspond to the selected rows of A in B and similarly let $\{d_i\}$ be the indices of the the columns selected. Let U be a $\mu \times m$ matrix such that for every i the c_i^{th} column of U is e_i , the standard basis vector in \mathbb{R}^m . Let V be the matrix representing this procedure for the set $\{d_i\}$. Then $UAV = B$. Both U and V have norm 1 for any p-norm, thus $\|B\|_p = \|UAV\|_p \leq \|A\|_p$

10. **Trefethen and Bau Problem 3.5.** Prove or disprove: For $E = uv^*$, $\|E\|_F = \|u\|_F \|v\|_F$.

Solution. True. It is easy to see that

$$\begin{aligned}\|E\|_F^2 &= \text{tr}((uv^*)^*(uv^*)) \\ &= \text{tr}(v^*u^*uv) \\ &= \text{tr}(\|u\|_2^2 \|v\|_2^2) \\ &= \|u\|_2^2 \|v\|_2^2 \\ &= \|u\|_F^2 \|v\|_F^2.\end{aligned}$$

Taking square root concludes the result.

11. **Trefethen and Bau Problem 4.2.** Suppose A is an $m \times n$ matrix and B is the $n \times m$ matrix obtained by rotating A 90 degrees clockwise on paper. Do A and B have the same singular values?

Solution. Note that B and A^T have the same columns but in a different order. Thus, $BP = A^T$ for some permutation matrix P . Since any permutation matrix is orthogonal, B and A^T have the same singular values and consequently so do B and A .

12. **Trefethen and Bau Problem 4.4.** Two matrices $A, B \in \mathbb{C}^{m \times m}$ are unitarily equivalent if $A = QBQ^*$. Prove or disprove: A and B are unitarily equivalent if and only if they have the same singular values.

Solution. True. Let A, B be unitarily equivalent matrices and $A = U\Sigma V^*$ be an SVD of A . Then $B = Q^*AQ = Q^*U\Sigma V^*Q$ is a valid SVD for B . Thus, A and B have the same singular values. For the converse, assume A and B have the same singular values. Then let $A = U_1\Sigma V_1^*$ and $B = U_2\Sigma V_2^*$ be two valid SVDS for A and B respectively. By rearrangement, it follows that $U_1U_2^*BV_2V_1^* = A$ and thus A and B are unitarily equivalent.

13. **Trefethen and Bau Problem 5.2.** Show that the set of full-rank matrices is a dense subset of $C^{m \times m}$.

Solution. Let $A \in C^{m \times m}$ and let $A = U\Sigma V^*$ be a (reduced) SVD of A . Denote $A_k = U(\Sigma + \frac{1}{k}I)V^*$. Since

$$\lim_{k \rightarrow \infty} \|A - A_k\|_2 = \lim_{k \rightarrow \infty} \left\| U\Sigma V^* - U\left(\Sigma + \frac{1}{k}I\right)V^* \right\|_2 = \lim_{k \rightarrow \infty} \frac{1}{k} = 0,$$

$\{A_k\} \rightarrow A$ as $k \rightarrow \infty$. Thus, the set of full rank matrices is a dense subset of $C^{m \times m}$.

14. **Trefethen and Bau Problem 5.4.** Suppose A has SVD $A = U\Sigma V^*$. Find an eigendecomposition of $\begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix}$.

Solution.

$$\begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix} \begin{bmatrix} V & V \\ U & -U \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}$$

2 Conditioning and Stability

1. Question from Dr. Xue

- (a) Find the absolute and relative condition numbers of $f(x) = e^{-2x}$ and $f(x) = \ln^3(x)$. For what values of x are these functions sensitive to perturbations?

Solution. For $f(x) = e^{-2x}$,

$$\begin{aligned} \kappa_A &= \|J_f(x)\| = |2e^{-2x}| \\ \kappa_r &= \frac{\|J_f(x)\| |x|}{|e^{-2x}|} = 2|x| \end{aligned}$$

That is, the absolute condition number is never sensitive since e^{-2x} is bounded above by 1. The relative condition number is only sensitive when x gets really large. When $f(x) = \ln^3(x)$, we have

$$\begin{aligned} \kappa_A &= \frac{3 \ln^2(x)}{|x|} \\ \kappa_r &= \frac{\kappa_A |x|}{|\ln(x)| \ln^2(x)} = \frac{3}{\ln(x)} \end{aligned}$$

The absolute condition number becomes large as $x \rightarrow 0^+$. The relative conditioning behaves in a similar manner when $x \rightarrow 1$. Both are incredibly sensitive to perturbations around their respective critical values.

- (b) Let $x_1, x_2 \in \mathbb{R}^+$ and $f(x_1, x_2) = x_1^{x_2}$. Find the relative condition number of $f(x)$, and for what range of values of x_1 and x_2 is this problem ill-conditioned?

Solution. Let $f(x_1, x_2) = x_1^{x_2}$. Then

$$\begin{aligned} \kappa_r &= \frac{\|J_f(x)\|_\infty \|x\|_\infty}{\|f(x)\|} = \frac{\max\{x_1^{x_2-1}x_2, x_1^{x_2-1}x_1 \log x_2\} \max\{x_1, x_2\}}{|x_1^{x_2-1}| |x_1|} \\ &= \frac{\max\{x_2, x_1 \log x_2\} \max\{x_1, x_2\}}{|x_1|} \end{aligned}$$

From this construction we see that if $x_1 \log x_2 > x_2$, then our relative condition number is simply $\log x_2 \max\{x_1, x_2\}$. This is only ill-conditioned when both x_1 and x_2 are very (very) large. However, if $x_1 \log x_2 < x_2$, then our relative condition number is $\frac{x_2^2}{x_1}$ which can clearly grow very fast when $x_1 \rightarrow 0$ or $x_2 \rightarrow \infty$.

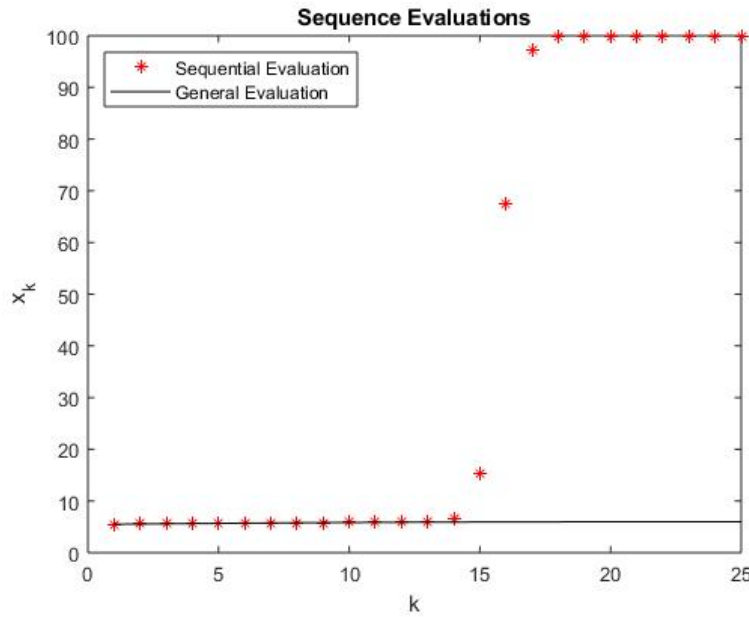


Figure B.1: Sequence x_k

2. **Question from Dr. Xue.** Consider the recurrence $x_k = 111 - \frac{1130 - \frac{3000}{x_{k-1}}}{x_k}$ who general solution is $x_k = \frac{100^{k+1}a + 6^{k+1}b + 5^{k+1}c}{100^k a + 6^k b + 5^k c}$, where a, b and c depend on the initial values. Given $x_0 = \frac{11}{2}$ and $x_1 = \frac{61}{11}$, we have $a = 0, b = c = 1$.

(a) Show that this gives a monotonically increasing sequence to 6.

Solution. Consider a rescaling by a factor of 6^k for a fixed k . Then

$$\frac{6^{k+1} + 5^{k+1}}{6^k + 5^k} = \frac{6 + 5(\frac{5}{6})^k}{1 + (\frac{5}{6})^k}$$

It follows that

$$\lim_{k \rightarrow \infty} \frac{6^{k+1} + 5^{k+1}}{6^k + 5^k} = \lim_{k \rightarrow \infty} \frac{6 + 5(\frac{5}{6})^k}{1 + (\frac{5}{6})^k} = 6$$

For monotonicity, note that as a function of k , $f(k) = \frac{6^{k+1} + 5^{k+1}}{6^k + 5^k}$ has a derivative of $f'(k) = \frac{30^k \log(\frac{6}{5})}{(6^k + 5^k)^2}$ which is positive for all nonnegative values of k .

(b) Implement this recurrence on MATLAB, plot $\{x_k\}$, compare with the exact solution. What is the condition number of the limit of this particular sequence as a function of x_0 and x_1 ?

Solution. The conditioning of the sequence as a function of x_0 and x_1 is infinity. For any perturbation of size $\varepsilon > 0$ on the inputs x_0, x_1 , we see that the limit jumps from 6 to 100 (since a becomes nonzero). That is, $\frac{100-6}{\varepsilon} \rightarrow \infty$ as $\varepsilon \rightarrow 0$. We see this behavior in Figure B.1. After only a few iterations, the sequence shoots to 100. The general solution avoids this problem by setting $a = 0$ in exact arithmetic.

3. **Question from Dr. Xue.** Let $p_{24}(x) = (x - 1)(x - 2) \cdots (x - 24) = a_0 + a_1x + \cdots + a_{24}x^{24}$. Evaluate the relative condition number of the k -th root $x_k = k$ subject to the perturbation of a_k for $k = 16, 17, 18, 19, 20$ and find the root that is most sensitive to the perturbation of the corresponding coefficient. Use MATLAB to compute the roots and compare them to the true roots.

Solution. Recall from class that for p_{24} , we have that

$$\kappa_r = \frac{x_j^{i-1} a_i}{p'(x_j)}$$

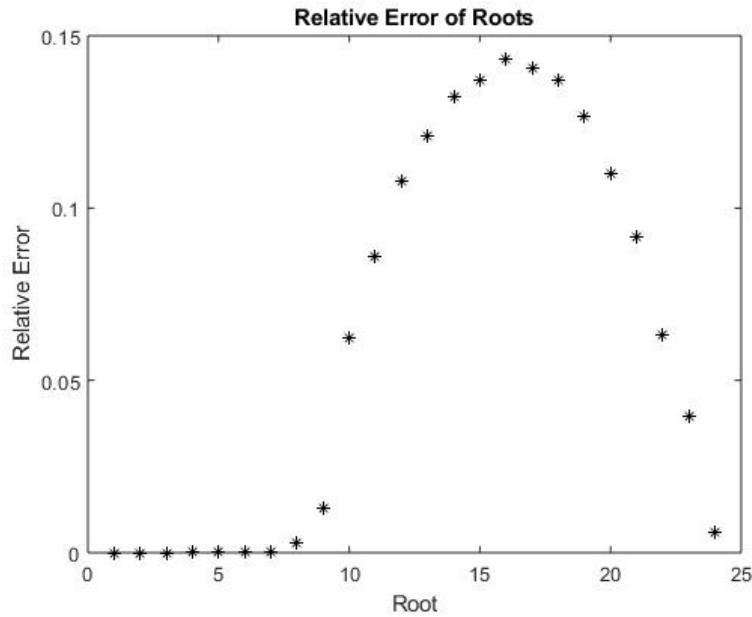


Figure B.2: Root Relative Errors

which simplifies to $\frac{x_k^{k-1} a_k}{p'(x_k)}$ for this problem since $i = j = k$. The computed results are summarized in Table B.1. We can clearly see that $k = 16$ is the most sensitive root to perturbations of the

| Root | Relative Condition Number |
|----------|---------------------------|
| $k = 16$ | $-2.12e16$ |
| $k = 17$ | $7.88e14$ |
| $k = 18$ | $-1.77e15$ |
| $k = 19$ | $-1.59e15$ |
| $k = 20$ | $-2.31e15$ |

Table B.1: Relative Conditioning of Roots

corresponding coefficient. Using MATLAB to compute the roots, we see that numerous roots had a relative error of over 0.1. Because of the large relative condition number at these roots, this is already terrible. Figure B.2 shows the relative error of the computed roots.

4. **Question from Dr. Xue.** Let x_0, \dots, x_n be $n + 1$ equidistant points on $[-1, 1]$ where $x_0 = -1, x_n = 1$. Use MATLAB's `vander` to generate Vandermonde matrices for $n = 9, 19, 29, 39$. Let $x = [1, \dots, 1]^T$ and $b = Ax$. Pretend that we do not know x and use numerical algorithms to solve for x . Let \hat{x} be the computed solution. Compute the relative forward errors and the smallest relative backward errors for GEPP, QR factorization, Cramer's Rule, $A^{-1}b$, and GE without pivoting. Comment on the forward/backward stability of these methods.

Solution. Table B.2 presents the numerical results for $n = 39$. We immediately see that both GEPP and QR factorization appear to be numerically backward stable. The final three algorithms cannot make such claim, but two of them, $A^{-1}b$ and Cramer's Rule appear to at least be forward stable, i.e. they produce forward errors similar to the forward errors of a backwards stable algorithm. However, for GE without pivoting, it is neither forward nor backward stable, for such a large n . Admittedly, it is likely the large n that creates the most problems. The forward error for $n = 9$ of GE without pivoting is on the order of 10^{-13} . Full evaluations can be found in Figure B.3.

5. **Question from Dr. Xue.** Though pivoting is needed for factorizing general matrices, it is not needed for symmetric positive definite and diagonally dominant matrices.

| Algorithm | Forward Error | Backward Error |
|------------------|---------------|----------------|
| GEPP | 4.1723 | $0.2113e - 16$ |
| QR Factorization | 4.5289 | $0.0408e - 15$ |
| $A^{-1}b$ | 4.7941 | 0.0065 |
| Cramer's Rule | 2.3003 | 0.0784 |
| GE without pivot | $1.828e28$ | 0.0134 |

Table B.2: Algorithm Stability

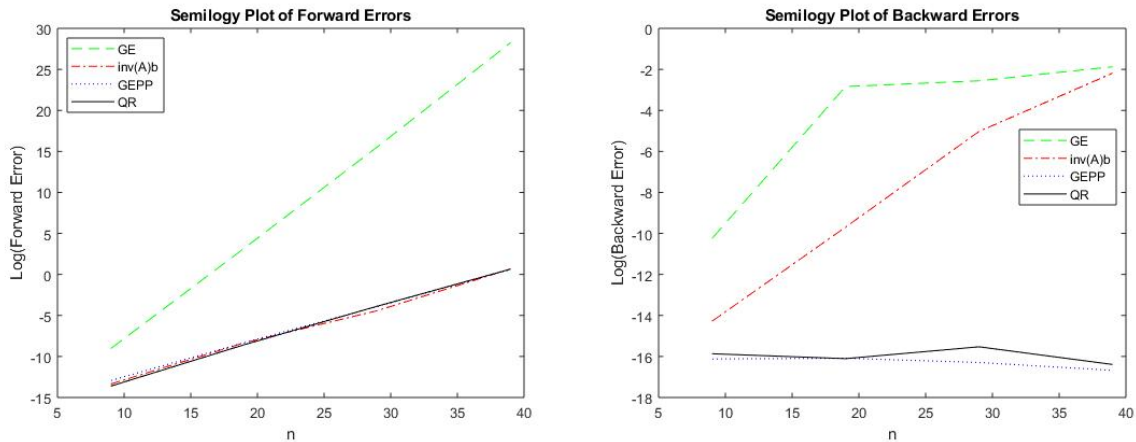


Figure B.3: Stability Graphs

(a) For a symmetric positive definite A , with the one-step Cholesky factorization

$$A = \begin{bmatrix} a_{11} & w^T \\ w & K \end{bmatrix} = \begin{bmatrix} \sqrt{a_{11}} & 0 \\ \frac{w}{\sqrt{a_{11}}} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & K - \frac{ww^T}{a_{11}} \end{bmatrix} \begin{bmatrix} \sqrt{a_{11}} & \frac{w^T}{\sqrt{a_{11}}} \\ 0 & I \end{bmatrix} = R_1^T A_1 R_1$$

show that the submatrix $K - \frac{ww^T}{a_{11}}$ is symmetric positive definite. Consequently the factorization can be completed without break-down. Then, show that $\|R\|_2 = \|A\|_2^{\frac{1}{2}}$, which means the elements in R are uniformly bounded by that of $\|A\|$. Explain why this observation leads to the backward stability of Cholesky factorization.

Solution. Since $\det(R_1^T) = \det(R_1) = \sqrt{a_{11}} > 0$, and all other leading principal minors of R_1, R_1^T are positive (the submatrix is the identity in both cases), these matrices must be invertible and symmetric positive definite. Thus, their inverses must also be SPD. So $A_1 = R_1^{-T} A R_1^{-1}$ is also positive definite since it is the product of 3 SPD matrices. Finally, because $K - \frac{ww^T}{a_{11}}$ is a principal minor of A_1 which itself is SPD, it must also be SPD. Recall from class that the magnitude of ρ_n controls the backward stability of an LU factorization algorithm.

To see the equivalence of norms, let $R = U\Sigma V^T$ be a singular value decomposition of R . We can then compute

$$A = R^T R = V\Sigma^2 V^T$$

which is a singular value decomposition of A . From here we see that the singular values of R are the square root of the corresponding singular values in A . Because both R and A are symmetric, this is also true for the eigenvalues. Thus,

$$\|R\|_2 = \rho(R) = \rho(A)^{\frac{1}{2}} = \|A\|_2^{\frac{1}{2}}$$

Because the elements of R are uniformly bounded by $\|A\|$, the growth factor is incredibly well behaved for Cholesky factorization and it is consequently backwards stable.

- (b) Suppose that $A = \begin{bmatrix} \alpha & w^T \\ v & C \end{bmatrix}$ is column diagonally dominant, with one step LU factorization $A = \begin{bmatrix} 1 & 0 \\ \frac{v}{\alpha} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & C - \frac{1}{\alpha}vw^T \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ 0 & I \end{bmatrix}$. Show that the submatrix $C - \frac{1}{\alpha}vw^T$ is also column diagonally dominant, and no pivoting is needed.

Solution. For brevity, let $D = C - \frac{1}{\alpha}vw^T$ and consequently, $A = \begin{bmatrix} \alpha & w^T \\ v & D + \frac{vw^T}{\alpha} \end{bmatrix}$. Suppose for the sake of contradiction that there exists some column j such that $D_{jj} < \sum_{i=1, i \neq j}^n D_{ij}$, i.e. D is not strictly diagonally dominant. Note that since A is column diagonally dominant, we have $\sum_{i=1}^n v_i < \alpha$. We may rearrange this to obtain $\frac{1}{\alpha} \left[\sum_{i=1, i \neq j}^n v_i w_j + v_j w_j \right] < w_j$. Combining the previous two inequalities gives us

$$D_{jj} + \frac{v_j w_j}{\alpha} < \sum_{i=1, i \neq j}^n \left(D_{ij} + \frac{v_i w_j}{\alpha} \right) + w_j$$

which is precisely the condition that column j of A is not strictly diagonally dominant, a contradiction.

- (c) Show that worst-case growth factor $\rho_n = 2^{n-1}$ for GEPP. However, we construct matrices with random elements, each are independent samples from the normal distribution of means 0 and standard deviation $\frac{1}{\sqrt{n}}$. Let $n = 32, 64, \dots, 2048$, and for each n , repeat the experiment 5000 times. Find the percent of experiments when $\rho_n > \sqrt{n}$. Comment on the chance of having a large ρ_n .

Solution. It is clear from the results in Table B.3 that the likelihood of having a troublesome growth factor is very small. Of the results tallied, they only indicate an instance of $\rho_n > \sqrt{n}$, which itself is not incredibly intimidating. However, it is important to note that many applications in practice are not on matrices sampled from a normal distribution. Often times, systems are banded or at least sparse. If I had written an LU factorization algorithm that took advantage of these special structures, I would test this as well, but I have not. Nonetheless, GEPP is most likely a very backward stable algorithm in practice.

| n | Bad Growth Factor Count |
|------|-------------------------|
| 32 | 15 |
| 64 | 13 |
| 128 | 19 |
| 256 | 16 |
| 512 | 9 |
| 1024 | 8 |
| 2048 | 5 |

Table B.3: Large Growth Factor Frequencies

6. **Question from Dr. Xue.** Consider the eigenvalue problem $Av = \lambda v$. Let $(\hat{\lambda}, \hat{v})$ be a computed eigenpair, which is assumed to be the exact eigenpair of a perturbed matrix $A + \Delta A$. Show that the minimum 2-norm of all ΔA is $\frac{\|A\hat{v} - \hat{\lambda}\hat{v}\|_2}{\|\hat{v}\|_2}$ and find a particular ΔA whose 2-norm is the minimum.

Solution. Consider the perturbed equation $(A + \Delta A)\hat{v} = \hat{\lambda}\hat{v}$. Rearranging gives $\Delta A\hat{v} = \hat{\lambda}\hat{v} - A\hat{v}$.

It follows that $\|A\hat{v} - \hat{\lambda}\hat{v}\|_2 = \|\Delta A\hat{v}\|_2 \leq \|\Delta A\|_2 \|\hat{v}\|_2$. Thus,

$$\frac{\|A\hat{v} - \hat{\lambda}\hat{v}\|_2}{\|\hat{v}\|_2} \leq \|\Delta A\|_2$$

To find a matrix satisfying this inequality, we will need a lemma. Let $u, v \in \mathbb{R}^n$. Then $\|uv^T\|_2 = \|u\|_2 \|v\|_2$.

Proof. Let $u, v \in \mathbb{R}^n$. To prove the lemma, first set \tilde{u}, \tilde{v} to be unit vectors in the directions of u, v respectively, i.e. $\tilde{u}\|u\|_2 = u, \tilde{v}\|v\|_2 = v$. Let U, V to be orthonormal basis extensions of u, v and $A = uv^T$. Denote E to be the zero $n \times n$ matrix with a 1 in the top left entry. Then $\tilde{u}\tilde{v}^T = \tilde{A} = UEV^T$ is a singular value decomposition of \tilde{A} . Scaling up to A , we see that $\|u\|_2 \|v\|_2$ is the only nontrivial singular value of A . Thus, it must be the 2-norm. \square

Now set $r = \hat{\lambda}\hat{v} - A\hat{v}$. Then consider $\Delta A = \frac{r\hat{v}^T}{\hat{v}^T\hat{v}}$. From the lemma we have $\|r\hat{v}^T\|_2 = \|r\|_2 \|\hat{v}\|_2$. Consequently, it follows that

$$\|\Delta A\|_2 = \frac{\|r\|_2 \|\hat{v}\|_2}{\|\hat{v}\|_2^2} = \frac{\|A\hat{v} - \hat{\lambda}\hat{v}\|_2 \|\hat{v}\|_2}{\|\hat{v}\|_2^2} = \frac{\|A\hat{v} - \hat{\lambda}\hat{v}\|_2}{\|\hat{v}\|_2}$$

Thus we have found a matrix satisfying the minimum 2-norm.

3 QR and Linear Least Squares

1. **Trefethen and Bau Problem 6.1.** Prove that $I - 2P$ is unitary for an orthogonal projector P and provide a geometric interpretation.

Solution. Note that

$$(I - 2P)(I - 2P)^* = I - 2P^* - 2P + 4PP^* = I - 4P + 4P = I.$$

The other direction holds similarly. Thus, $I - 2P$ is unitary. Recall that $I - P$ is a projection onto the null space of P . Then $I - 2P = I - P - P$ is simply moving in this direction, twice. This amounts to a reflection across the null space of P . Reflections preserve distance and are therefore unitary.

2. **Trefethen and Bau Problem 6.2** Let E be the $m \times m$ matrix such that $E = (I + F)/2$ where F is the $m \times m$ matrix that flips (x_1, \dots, x_m) to (x_m, \dots, x_1) . Classify E as a projector.

Solution. It is easy to argue $F^2 = I$. Then

$$E^2 = \frac{1}{4}(I + F)(I + F) = \frac{1}{4}(I + 2F + F^2) = \frac{1}{2}(I + F) = E$$

So E is a projector. Both I and F are clearly symmetric, so E is an orthogonal projector.

3. **Trefethen and Bau Problem 6.3.** Given $A \in \mathbb{C}^{m \times n}$ with $m \geq n$, show that A^*A is nonsingular if and only if A has full rank.

Solution. Suppose A is not full rank. Then there exists nonzero x such that $Ax = 0$ and thus $A^*Ax = A^*0 = 0$. Now let A^*A have non-full rank. Then there exists nonzero x such that $A^*Ax = 0$. Suppose that A were full rank. Then so too is A^* . Thus, $A^*Ax = 0$ implies $Ax = 0$. But A was assumed to have full rank so $x \neq 0$. This is a contradiction and consequently A is not full rank.

4. **Trefethen and Bau Problem 6.5.** Let $P \in \mathbb{C}^{m \times m}$ be a nonzero projector. Show that $\|P\|_2 \geq 1$ with equality if and only if P is an orthogonal projector.

Solution. From properties of projectors, we see that $\|P\| = \|P^2\| \leq \|P\|^2$ and thus $\|P\| \geq 1$. The singular values of P are simply the eigenvalues of P^*P . For an orthogonal projector, $P^*P = P^2 = P$. Since $Px = \lambda x \implies Px = \lambda Px$, it is clear that any eigenvalue of $P = P^*P$ must have modulus 1. Thus, $\|P\|_2 = 1$.

5. **Trefethen and Bau Problem 7.2.** Let $A \in \mathbb{R}^{m \times n}$ be a matrix with the property that columns 1, 3, 5, 7, ... are orthogonal to columns 2, 4, 6, ... In a reduced QR factorization, what special structure does R possess?

Solution. The columns of Q satisfy $q_i \in \text{span}(\{a_i\})$ for $i = j \pmod 2, i \leq j$. Consequently, $r_{ij} = \langle q_i, a_j \rangle = 0$ whenever $i \not\equiv j \pmod 2$. To see this, note that $q_1 \in \text{span}(a_1)$ by construction and $q_2 \perp a_2 - r_{12}q_1 = a_2$. Now assume the statement holds for all $t \leq 2i$. Then

$$q_{2i+1} \perp a_{2i+1} - \sum_{j=0}^{i-1} r_{2j+1,2i+1} q_{2j+1} - \sum_{j=0}^{i-1} r_{2j,2i+1} q_{2j} = a_{2i+1} - \sum_{j=0}^{i-1} r_{2j+1,2i+1} q_{2j+1}$$

from the induction hypothesis since the second sum is all zeros. The case for q_{2i+2} can be done similarly. Thus, the statement holds for all $t \leq n$ by induction.

6. **Trefethen and Bau Problem 7.5.** Let $A \in \mathbb{R}^{m \times n}$ and let $A = QR$ be a reduced QR factorization. Suppose R has k nonzero diagonal entries for $0 \leq k \leq n$. What does this imply about the rank of A ?

Solution. For any j such that $r_{jj} = 0$, it must be that $a_j - \sum_{i=1}^{j-1} r_{ij} q_i = 0$, that is, $a_j \in \text{span}(q_1, \dots, q_{j-1}) = \text{span}(a_1, \dots, a_{j-1})$. Thus, $\text{rank}(A) \leq k$. To see that $\text{rank}(A) \geq k$, note that since Q is unitary (invertible), $\text{rank}(A) = \text{rank}(R)$. Since for any j such that $r_{jj} \neq 0$ Re_j are k linearly independent vectors in the column space of R , we have that $\text{rank}(A) = \text{rank}(R) \geq k$. Thus, $\text{rank}(A) = k$.

7. **Trefethen and Bau Problem 10.1.** Determine the eigenvalues, determinant, and singular values of a Householder reflector. For the eigenvalues, give both a geometric and algebraic proof.

Solution. Let us begin with a geometric interpretation of the Householder transformation. We would like to find a transformation that maps x to $\|x_1\| e_1$. Consider the vector $v = \|x_1\| e_1 - x$. The projection of x onto the hyperplane orthogonal to the vector v is $\text{proj}(x) = x - v \left(\frac{v^T x}{v^T v} \right)$. However, if we are to project to $\|x_1\| e_1$, we need to go twice this distance. Thus, our Householder reflection is characterized by $x - 2v \left(\frac{v^T x}{v^T v} \right)$ and we say that our Householder reflect $H = (I - \frac{2vv^T}{v^T v})$. This geometric interpretation immediately motivates many properties such as symmetry and orthogonality. One can also see eigenvalues from this interpretation. Since our projection is a reflection across a hyperplane (of dimension $m - 1$), all but 1 direction is fixed. That is, $\lambda_i = 1$ for $i = 1, \dots, m - 1$ and $\lambda_m = -1$.

For the algebraic solutions, note that if $y \perp v$. Then $Hy = y - 2 \frac{v^T y}{v^T v} v = y - 2 \cdot 0 = y$. Since $\dim \text{span}(y) = 1, \dim \text{span}(y)^\perp = m - 1$. Thus, $\lambda = 1$ is an eigenvalue with multiplicity $m - 1$. Noting that $Hv = -v$ shows that the remaining eigenvalue is -1 . The determinant is simply the product of all of these eigenvalues. Thus, $\det H = 1^{m-1} \cdot -1 = -1$. Since H is symmetric (because vv^T is symmetric), we know that the singular values are the absolute value of the eigenvalues. That is, $\sigma_i = 1$ for all $i = 1, \dots, m$.

8. **Trefethen and Bau Problem 10.4.** Consider the orthogonal matrices

$$F = \begin{bmatrix} -\cos \theta & \sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, J = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

Describe what geometric effects left-multiplication by F and J have on the plane \mathbb{R}^2 .

Solution. For F , since $\text{tr}(F) = 0$ and $\det(F) = -1$, we can conclude that $\lambda_{1,2} = \pm 1$. Thus, F is an orthogonal matrix with eigenvalues ± 1 . Hence, F is a reflection. For J , notice that $\det J = 1$ and consider vectors x and Jx . Let t be the angle between these two vectors. We can compute that

$$\cos t = \frac{\langle x, Jx \rangle}{\|x\|_2 \|Jx\|_2} = \frac{\langle x, Jx \rangle}{\|x\|_2^2}$$

since J is unitary. Evaluating the numerator, we continue with

$$\cos t = \frac{\cos \theta (x_1^2 + x_2^2)}{\|x\|_2^2} = \cos \theta.$$

Hence, $t = \theta$. Thus, J has the effect of a rotation on a vector in \mathbb{R}^2 . Plugging in $\theta = \pi/2$ and $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, we see that this is a clockwise rotation.

9. **Question from Dr. Xue.** Consider the Givens rotation $G = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$. Give a geometric interpretation of the action of G on a vector in \mathbb{R}^2 . Determine both the eigenvalues and the singular values.

Solution. Recall that the matrix $G = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$ rotates a vector in \mathbb{R}^2 clockwise by an angle of θ . If θ is the angle between a point x and the x-axis, then Gx rotates x to the x-axis. We accomplish this by setting θ to be the angle such that $\cos(\theta) = \frac{x_1}{\sqrt{x_1^2 + x_2^2}}$ and $\sin(\theta) = \frac{x_2}{\sqrt{x_1^2 + x_2^2}}$

Algebraically, we can compute the eigenvalues as the roots of the equation

$$\cos^2(\theta) + \lambda^2 - 2\lambda \cos(\theta) + \sin^2(\theta) = 0$$

Letting $x = (a, b)^T$, this becomes

$$1 + \lambda - \frac{2\lambda a}{\sqrt{a^2 + b^2}} = 0$$

after substituting θ . By the quadratic formula, we obtain

$$\lambda_{1,2} = \frac{\frac{2a}{\sqrt{a^2 + b^2}} \pm \sqrt{\frac{2a}{\sqrt{a^2 + b^2}} - 4}}{2} = \frac{2a}{\sqrt{a^2 + b^2}} \pm \frac{2b}{\sqrt{a^2 + b^2}}i = \cos(\theta) \pm i \sin(\theta)$$

as our eigenvalues. The singular values can be found through construction of a singular value decomposition. First note that for any α, β , the vectors $u = (-\alpha, \beta)^T, v = (\beta, \alpha)^T$ are orthogonal since

$$\begin{bmatrix} -\alpha & \beta \end{bmatrix} \begin{bmatrix} \beta \\ \alpha \end{bmatrix} = -\alpha\beta + \beta\alpha = 0$$

Thus,

$$\begin{bmatrix} -\cos(\theta) & \sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} I_2 = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix} = G$$

is a valid SVD for G . From this decomposition, we immediately read off $\sigma_{1,2} = \pm 1$.

10. **Trefethen and Bau Problem 11.1** Suppose $A \in \mathbb{R}^{m \times n}$ matrix has the form

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

where $A_1 \in \mathbb{R}^{n \times n}$ is nonsingular. Show that $\|A^\dagger\|_2 \leq \|A_1^{-1}\|_2$.

Solution. Let $A = U\Sigma V^T$ be a reduced SVD of A . Then

$$A^\dagger = (A^T A)^{-1} A^T = V \Sigma^{-2} V^T V \Sigma U^T = V \Sigma^{-1} U^T$$

is a valid SVD for A^\dagger and hence $\|A^\dagger\|_2 = 1/\sigma(A)$ where $\sigma(X)$ represents the smallest singular value of a matrix X . Similarly, it can be shown that $\|A_1^{-1}\|_2 = 1/\sigma(A_1)$. Thus, it suffices to show that $\sigma(A_1) \leq \sigma(A)$. To see this, compute

$$\sigma(A) = \min_{\|x\|=1} \|Ax\| = \min_{\|x\|=1} \left\| \begin{pmatrix} A_1 x \\ A_2 x \end{pmatrix} \right\| \geq \min_{\|x\|=1} \|A_1 x\| = \sigma(A_1)$$

and we conclude that $\|A^\dagger\|_2 \leq \|A_1^{-1}\|_2$.

11. **Trefethen and Bau Problem 11.3.** Take $m = 50, n = 12$. Using MATLAB's `linspace`, define t to be the m vector corresponding to linearly spaced grid points from 0 to 1. Using MATLAB's `vander` and `fliplr`, define A to be the $m \times n$ matrix associated with the least squares fitting on this grid by a polynomial of degree $n - 1$. Take b to be the function $\cos(4t)$ evaluated on the grid. What do you observe from the 6 methods?

- (a) Normal equations (using MATLAB's `\`)
- (b) QR factorization by MGS
- (c) QR factorization by Householder transformations
- (d) QR factorization computed by MATLAB's `qr`
- (e) $x = A \backslash b$ in MATLAB
- (f) SVD, using MATLAB's `svd`

Solution. The results are summarized in Table B.4. To generate the table, I assumed that the solution $x = A \backslash b$ computed using MATLAB's backslash operator was the true solution. Each vector of coefficients were then compared to the solution generated with the backslash, and then normed. The results are a bit unsurprising. The normal equations and MGS look to be numerically unstable while Householder, MATLAB QR factorization, and MATLAB SVD all seem to be relatively close to the true solution. I am unsure as to how we only managed to get the square root of machine precision, but it is likely due to the fact that $x = A \backslash b$ is, in fact, not the true solution.

| Algorithm | Norm of Errors |
|-------------|----------------|
| Normal Eq | $1.06e - 1$ |
| MGS | $1.526e - 1$ |
| Householder | $2.85e - 8$ |
| MATLAB QR | $3.31e - 8$ |
| MATLAB SVD | $3.20e - 8$ |

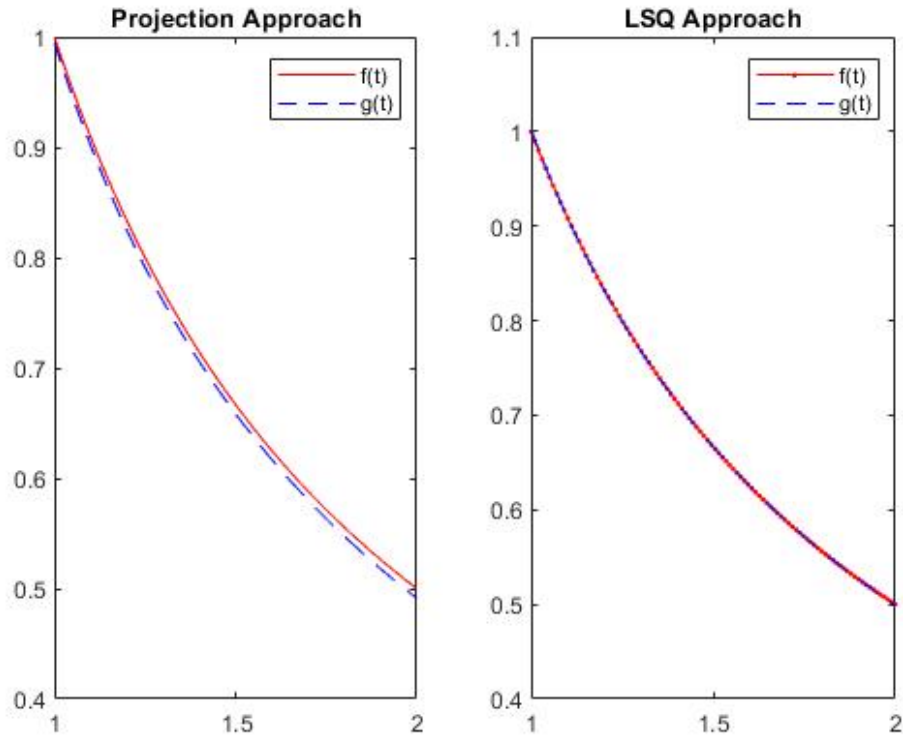
Table B.4: Norm of Coefficient Errors

12. **Trefethen and Bau Problem 11.2.** How closely can $f(x) = \frac{1}{x}$ be measured in the L^2 norm by linear combinations of $e^x, \sin(x)$, and $\Gamma(x)$ over $[1, 2]$?

Solution. We wish to project $f(x)$ to some space P spanned by $e^x, \sin(x)$, and $\Gamma(x)$. For an orthonormal basis $\{\phi_i(x)\}_{i=1}^n$, the solution to $\min_{g \in P} \|f - g\|_*$ is

$$g(x) = \sum_{i=1}^n \langle \phi_i(x), f(x) \rangle \phi_i(x)$$

Hence, we only need to find an orthonormal basis for P with respect to the L^2 norm over $[1, 2]$. By applying Gram-Schmidt to our functions, we can construct an orthonormal basis $\{\phi_i(x)\}_{i=1}^n$ of P . I have left the details in the code, but provided Figure B.4 to show the results. The results look strong to the naked eye, but has an error with norm around 0.05.

Figure B.4: Approximations of $f(x)$

After being reminded that we can tackle this problem via LSQ, I tried a different approach. Recall that we can write our approximation problem as

$$\begin{bmatrix} f_1(x) & f_2(x) & f_3(x) \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} = f(x)$$

where $f_i(x)$ are our basis functions. Naively applying the normal equation technique to this LSQ problem gives us the equation $A^T A x = A^T b$ where $[A^T A]_{ij} = \langle f_i, f_j \rangle = \int_0^1 f_i(x) f_j(x) dx$ and $[A^T b]_i = \langle f_i, f \rangle = \int_0^1 f_i(x) f(x) dx$. Using MATLAB's backslash operator on this linear system gives the coefficients $x = [-0.1078, 0.0092, 1.2872]$. Since $\text{cond}(A^T A) \approx 10^4$, we can expect these coefficients to have around $16 - 4 = 12$ digits of accuracy, which is enough to provide us with the nice approximation shown in the right pane of Figure B.4.