# Numerical Analysis Course Notes

Isaac Holt

March 6, 2023

# Contents

# 1 Chebyshev Polynomials

**Theorem 1.0.1.** let $w_n(x) = (x - x_0)\ldots(x - x_n)$ with distinct nodes $x_0, \ldots, x_n$, $x_j \subset [-1, 1]$. Then the maximum of $|w_n(x)|$ on $[-1, 1]$ attains its smallest value $(2^{-n})$ iff $x_j$ are the zeros of $T_{n+1}(x)$.

*Proof.* ($\Longleftarrow$): By construction $2^{-n}T_{n+1}(x)$ is a monic polynomial (highest power of $x$ is 1) with $n+1$ roots in $[-1, 1]$. Suppose $S_{n+1}(x) = (x - z_0)\ldots(x - z_n)$ is another monic polynomial such that $\max|S_{n+1}(x)| < 2^{-n} = \max|2^{-n}T_{n+1}(x)|$. Let $q_n(x) := 2^{-n}T_{n+1}(x) - S_{n+1}(x)$. Then $q_n(x) \in P_n$ since the coefficient of $x^{n+1}$ in $T_{n+1}(x)$ and $S_{n+1}(x)$ are both 1 and so cancel out.

Then $q_n(y_j) = 2^{-n}T_{n+1}(y_j) - S_{n+1}(y_j)$ ($y_j$ are the extrema of $T_n(x)$). $|S_{n+1}(y_j)| < 1$ by hypothesis. Therefore $q_n(y_j) > 0$ if $j$ is odd and $< 0$ otherwise.

Since we have $n + 2$ of $y_j$, $q_n$ has at least $n + 1$ zeros. But since $q_n \in P_n$, we must have $q_n(x) = 0$. Therefore $S_{n+1}(x) = 2^{-n}T_{n+1}(x)$. $\square$

**Remark.** To use this in $[a, b]$ instead of $[-1, 1]$, one simply maps $x_j \to a + (x_j + 1)\frac{b-a}{2}$.

**Remark.** Putting the above into Cauchy's error formula, we have

$$\sup|f(x) - p(x)| \le 2^{-n}\left(\frac{b - a}{2}\right)^{n+1}\frac{1}{(n + 1)!}$$

**Remark.** We have by the above theorem, $\max|w_n(x)| = \max|(x - x_0)\ldots(x - x_n)| \ge 2^{-n}$ for any choice of $x_1, \ldots, x_n$, $x_j \subset [-1, 1]$. So $2^{-n}$ is a lower bound for $|w_n(x)|$.

The upper bound is given by $\max|w_n(x)| \le \epsilon|b - a|^n$.

# 2 Root Finding

## 2.1 Bracketing: Bisection

Given $f \in C^0([a, b])$ with $f(a)f(b) < 0$, repeat:

- let $(a_0, b_0) = (a, b)$

- let $m_n = \frac{1}{2}(a_n + b_n)$

- if $f(m_n)f(a_n) \geq 0$, set $(a_{n+1}, b_{n+1}) = (m_n, b_n)$

- otherwise, set $(a_{n+1}, b_{n+1}) = (a_n, m_n)$

$b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$. By the Intermediate Value Theorem, if $f(m_n) \neq 0$, for some $p \in (a_n, b_n)$, $f(p) = 0$.
$|p - m_n| \leq 2^{-(n+1)}(b - a)$

**Remark.** Each time, the width of the interval halves. In principle, we could get an approximation to any desired accuracy, but there are some caveats (e.g. with floating points).

## 2.2 Bracketing: False Position

Suppose we have $|f(b)| \ll |f(a)|$, then we would expect $p$ to be closer to $b$ than to $a$. Instead of $m_n = \frac{1}{2}(a_n + b_n)$, set

$$m_n = b_n - f(b_n)\frac{b_n - a_n}{f(b_n) - f(a_n)}$$

i.e. $m_n$ is the x-intercept of the line from $(a_n, f(a_n))$ to $(b_n, f(b_n))$. This should sometimes give much faster approximation than bisection, but not always.

## 2.3 Aside: Continuity and Convergence

**Definition 2.3.1.** $f : I \to \mathbb{R}$ is continuous at $x \in I$ if for every $\epsilon > 0$, for some $\delta(x, \epsilon)$, $|y - x| < \delta \Rightarrow |f(x) - f(y)| < \epsilon$ for every $y \in B(x)$ ($B(x)$ is an open interval containing $x$).

**Remark.** In general, $\delta$ depends on $\epsilon$ and $x$. When $\delta$ is independent of $x$, $f$ is uniformly continuous.

**Definition 2.3.2.** $f : I \to \mathbb{R}$ is Lipschitz continuous in $I$ if for some $L > 0$, $|f(y) - f(x)| \leq L|y - x|$ for every $x \in I, y \in I$. In this case, $\delta = \epsilon/L$.

**Remark.** $L$ (like $\delta$ above) is not unique. The smallest such $L$ is called the Lipschitz constant of $f$ in $I$.

**Lemma 2.3.3.**

1. If $f$ is differentiable and $I$ is compact, $f$ is Lipschitz in $I$.

2. If $f$ is Lipschitz, $f$ is continuous.

*Proof.*

$$f(y) - f(x) = \int_x^y f'(x)ds$$

$$|f(y) - f(x)| = |\int_x^y f'(x)ds| \leq \int_x^y |f'(x)|ds$$

$$\leq \max_{s \in I} |f'(s)| \int_x^y ds = \max_{s \in I} |f'(s)||y - x|$$

We can take $L = \max_{s \in I} |f'(s)|$ $\qquad \square$

**Remark.** The converses of 1. and 2. are false.

**Remark.**
- When $f$ is continuous in $I$, we write $f \in C^0(I)$.

- When $f$ is differentiable in $I$, we write $f \in C^1(I)$.

- When $f$ is Lipschitz in $I$, we write $f \in C^{0,1}(I)$.

- We can then write $C^1(I) \subsetneq C^{0,1}(I) \subsetneq C^0(I)$.

**Definition 2.3.4.** A sequence $(x_n)$ in $\mathbb{R}^d$ converges to $x$ if for every $\epsilon > 0$, for some $N(\epsilon)$, for every $n \geq N(\epsilon)$, $|x_n - x| < \epsilon$.
This relies on us knowing $x$ in the first place.

**Definition 2.3.5.** A sequence $(x_n)$ is a Cauchy sequence if for every $\epsilon > 0$, for some $N(\epsilon)$, for every $m \geq N, n \geq N$, $|x_n - x_m| < \epsilon$.

**Theorem 2.3.6.** Let $(x_n)$ be a Cauchy sequence in $\mathbb{R}^d$. Then $(x_n)$ converges.

This is useful as it allows us to prove convergence without knowing $x$.

## 2.4 Fixed Point Iterations

We seek $x$ such that $f(x) = 0$ for a function $f$. We rewrite this as

$$x = g(x)$$

We then seek to solve this equation by iterations:

1. pick some $x_0$

2. set $x_{n+1} = g(x_n)$

**Theorem 2.4.1.** (1d local convergence theorem): Let $g \in C'([a, b])$ have a fixed point $x_\star \in [a, b]$ $(g(x_\star) = x_\star)$ with $|g'(x_\star) < 1$. Then for $x_0$ sufficiently close to $x_\star$, the iteeration $x_{n+1} = g(x_n)$ converges to $x_\star$.

*Proof.* Let $g'(x_\star) = L \in (0, 1)$ $(g'(x_\star) < 0$ is analogous). Since $g'$ is continuous at $x_\star$, for every $L' \in (L, 1)$, for some $\delta(L') > 0$, $g'(x) \leq L' < 1$ for every $x \in (x_\star - \delta, x_\star + \delta) = B_\delta$, therefore for every $x \in B_\delta$, $y \in B_\delta$, $|g(x) - g(y)| \leq \sup_{s \in B_\delta} |g'(s)||x - y| = L'|x - y|$ with $L' < 1$.

Let $x_\star \in B_\delta$, then $|g(x) - x_\star| = |g(x) - g(x_\star)| \leq L'|x - x_\star|$ since $x_\star = g(x_\star)$. So $x - x_\star \delta$ as $x \in B_\delta$, so $|g(x) - x_\star| \leq L'\delta \leq \delta$, therefore $g(B_\delta) \subseteq B_\delta$. $\qquad \square$

**Remark.** We do not need to know $x_\star$ to apply the 1d local convergence theorem, we just need to know that $|g'(x)| < 1$ for every $x \in I$ for some interval $I$.

## 2.5 Order of convergence

Order of convergence is a rough measure of how quickly $x_n \to x$. We mainly look at sequences arising from iterations with a nice RHS (so not bisection).

**Definition 2.5.1.** Let $x_n \to x_*$ and assume that $x_n \neq x_*$ for every $n \geq 0$. $x_n \to x_*$ with order at least $\alpha > 1$ if

$$\lim_{n \to \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|^\alpha} = \lambda < todo$$

and with order $\alpha = 1$ if also $\lambda < 1$.

**Example 2.5.2.** $x_n = n^{-\beta}$, $\beta > 0$

$$\frac{|x_{n+1} - x_*|}{|x_n - x_*|} = (\frac{n}{n+1})^\beta \to 1$$

**Example 2.5.3.** $x_n = e^{-n}$

$$\frac{|x_{n+1} - x_*|}{|x_n - x_*|} = 1/e < 1$$

**Example 2.5.4.** $x_n = \frac{1}{n!}$

$$\frac{|x_{n+1} - x_*|}{|x_n - x_*|} = 1/(n+1) \to 0$$

The order of convergence of $x_n \to x_*$ is

$$\alpha = \sup\{\beta : \lim_{n \to \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|^\beta} < todo\}$$

for $\alpha > 1$.

For $\alpha = 1$, we also require that the limit $< 1$.

The convergence is linear if $\alpha = 1$, superlinear if $\alpha > 1$ and sublinear otherwise.

**Remark.** Order of convergence need not be an integer.

**Remark.** We need to know $x_*$ in order to determine order of convergence.

**Remark.** Our definition is not comprehensive for general sequences.

Applying this to iterations:

$x_{n+1} - x_* = g(x_n) - g(x_*) = (x_n - x_*)g'(c_n)$ for some $c \in \text{conv}\{x_n, x_*\}$ by the Mean Value Theorem.

Therefore

$$|x_{n+1} - x_*||x_n - x_*| = |g'(c_n)| \to g'(x_*)$$

We conclude that for $g \in C^2(I)$, the iteration $x_{n+1} = g(x_n)$ converges linearly if $g'(x_*) \neq 0$ and $|g'(x_*)| < 1$, and superlinearly otherwise.

**Proposition 2.5.5.** Let $g \in C^{N+1}(D)$ for some $D \subseteq \mathbb{R}$ and let $g(x_*) = x_*$, with $x_*$ in the interior of $D$.

Then the iteration $x_{n+1} = g(x_n)$ converges to $x_*$ for $x_0$ sufficiently close to $x_*$ with order $N+1$ iff $g'(x_*) = g''(x_*) = \cdots = g^{(N)}(x_*) = 0$ and $g^{(N+1)}(x_*) \neq 0$.

*Proof.* $x_{n+1} - x_* = g(x_n) - g(x_*) = g(x_*) + (x_n - x_*)g'(x_*) + \cdots + \frac{(x_n - x_*)^N}{N!}g^{(N)}(x_*) + \frac{(x_n - x_*)^{N+1}}{(N+1)!}g^{(N+1)}(c_n) - g(x_*) = \frac{(x_n - x_*)^{N+1}}{(N+1)!}g^{(N+1)}(c_n)$. Thus

$$|x_{n+1} - x_*||x_n - x_*|^{N+1} = \frac{|g^{(N+1)}(c_n)|}{(N+1)!} \to \frac{|g^{(N+1)}(x_*)|}{(N+1)!} < todo$$

$\square$

## 2.6 Higher order iterative methods

We want to rearrange $f(x) = 0$ to get faster convergence.

$x_{n+1} = g(x_n) = x_n + \phi(x_n)f(x_n)$ for some $\phi$.

Using the above proposition, we need $g'(x_*) = 0$.

$g'(x_*) = 1 + \phi'(x_*)f(x_*) + \phi(x_*)f'(x_*) = 0$

So if $f'(x_*) \neq 0$, we take $\phi(x) = -\frac{1}{f'(x)}$.

This is the Newton-Raphson method:

$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$.

Multiplication of a computer is parallelisable while division is not. So we can use Newton-Raphson to divide numbers with multiplication.

To compute $x_* = 1/b$ so $f(x_*) = \frac{1}{x_*} - b = 0$, so using Newton-Raphson, $x_{n+1} = x_n - \frac{x_n^{-1} - b}{-x_n^{-2}} = x_n(2 - bx_n)$ which involes only multiplication and subtraction. When taking out the exponent, $x_0 \in \left[\frac{1}{2}, 1\right)$, and this converges quickly.

**Remark.** This only works with floating points, not integers. Floating point division is 5 times faster than integer division.

**Remark.** It can be difficult in practice to determine the interval/domain of convergence for Newton-Raphson. We should always perform a "sanity check" when using it.

**Example 2.6.1.** One advantage of iterative methods is that they also work (in principle) in higher dimensions.

Suppose $f : \mathbb{R}^2 \to \mathbb{R}^2$ has a root at $\underline{p} = (p_1, p_2)$. Using Taylor expansion at the current point $\underline{x_n}$, derive Newton-Raphson (2D):

$$\underline{x_{n+1}} = \underline{x_n} - (Df)^{-1}\underline{f}(\underline{x_n})$$

We write $\underline{x_*} = \underline{p}$ such that $f_1(p_1, p_2) = f_2(p_1, p_2) = 0$. Taylor-expanding at $\underline{x}$:

$$0 = f_1(p_1, p_2) = f_1(x_1, x_2) + (p_1 - x_1)\frac{\partial f_1}{\partial x_1}(x_1, x_2) + (p_2 - x_2)\frac{\partial f_1}{\partial x_2}(x_1, x_2) + O(|\underline{p} - \underline{x}|^2)$$

$$0 = f_2(p_1, p_2) = f_2(x_1, x_2) + (p_1 - x_1)\frac{\partial f_2}{\partial x_1}(x_1, x_2) + (p_2 - x_2)\frac{\partial f_2}{\partial x_2}(x_1, x_2) + O(|\underline{p} - \underline{x}|^2)$$

In matrix form:

$$(0, 0) = (f_1(\underline{x}), f_2(\underline{x})) = (Df)(x_1, x_2) \cdot (x_1 - p_1, x_2 - p_2) + O(|\underline{p} - \underline{x}|^2)$$

Assuming that $Df$ is invertible (equivalently, $f'(\underline{x}) \neq 0$), we can multiply the equation by $(Df)^{-1}$ to get

$$(p_1, p_2) = (x_1, x_2) - (((Df)^{-1})(x_1, x_2)) \cdot (f_1(\underline{x}), f_2(\underline{x})) + O(|\underline{p} - \underline{x}|^2)$$

So

$$p = \underline{x} - (((Df)^{-1})(x_1, x_2))\underline{f(x)} + O(|\underline{p} - \underline{x}|^2)$$

We can use this to construct our iteration by replacing $\underline{x}$ with $\underline{x_n}$ and $\underline{p}$ with $\underline{x_{n+1}}$, and removing the $O(|\underline{p} - \underline{x}|^2)$.

## 2.7  Secant method

One disadvantage of Newton-Raphson is that we need the derivative, $f'$. If $f$ is complicated or is itself computed numerically, we need to approximate $f'$. An alternative method is the secant method.

The secant method approximates $f'$ with:

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

so the iteration becomes

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n)$$

**Remark.** This is a (scalar) two-step method: $x_{n+1} = g(x_n, x_{n-1})$, where $x_n$ and $x_{n-1}$ are needed to calculate $x_{n+1}$.

**Theorem 2.7.1.** Let $f \in C^2$ with $f(x_*) = 0$ and $f'(x_*) \neq 0$. Then the secant method is convergent with order $\alpha = \frac{1+\sqrt{5}}{2}$ for every $x_0 \neq x_1$ sufficiently close to $x_*$.

*Proof.* TODO: See video on Panopto $\qquad\qquad\square$

**Remark.**

1. When implementing the secant method, one must be careful with floating point effects:

$$\frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})}$$

   becomes very inaccurate as $x_n - x_{n-1} \to 0$.

2. $e_n := x_n - x_*$ alternates in sign if $g'(x_k) < 0$ and has the same sign if $g'(x_*) > 0$. From the proof of the method,

$$e_{n+1} = e_n e_{n-1} \frac{f''(\theta_n)}{f'(\phi_n)}$$

   where $\theta_n, \phi_n \in \text{conv}\{x_{n-1}, x_n, x_{n+1}\}$.

   For $e_0 e_1 < 0$ and $n$ sufficiently large, the error $e_n$ follows the pattern $+, +, -$ or $-, -, +$.

# 3 Numerical differentiation

## 3.1 Forward/backward difference

$f'(x) = \lim_{h \to 0} \frac{f(x+h)-f(x)}{h} \approx \frac{f(x+h)-f(x)}{h}$ for some small $h \neq 0$.

More rigorously, we can use Taylor series:

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

where $\xi \in \text{conv}\{x, x+h\}$. So

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

The error $\frac{h}{2}f''(\xi)$ is of order $O(h)$.

For $h > 0$ this is called forward difference.

For $h < 0$, this is called backward difference.

## 3.2 Centred difference

$$f(x \pm h) = f(x) \pm hf'(x) + \frac{h^2}{2}f''(x) \pm \frac{h^3}{6}f^{(3)}(\xi_\pm)$$

Then

$$f'(x) = \frac{f(x+h)-f(x-h)}{2h} + \frac{h^2}{12}\left(f^{(3)}(\xi_+) + f^{(3)}(\xi_-)\right)$$

The error $\frac{h^2}{12}\left(f^{(3)}(\xi_+) + f^{(3)}(\xi_-)\right)$ is of order $O(h^2)$.

**Remark.** It is important to **not take $h$ too small** when computing numerically. There is always a tradeoff between formal analytical accuracy and floating point errors. Formal analytical accuracy is better for smaller $|h|$, floating point errors are worse for smaller $|h|$.

**Remark.** Sometimes, we can only take $h < 0$ or $h > 0$ (not both), e.g. when solving differential equations numerically. If we have an ODE

$$\frac{du}{dt} = F(u)$$

Given $u(0)$ we want to solve for $u(t), t \geq 0$. We approximate $u(t)$ by a $u(t_n)$ with $t_n = n\delta t, n \in \mathbb{N}$, with $\delta t > 0$ small. Then

$$\frac{du}{dt} \approx \frac{u(t_n + \delta t) - u(t_n)}{\delta t} \approx F(u(t_n))$$

## 3.3 Richardson extrapolation

Let $f'(x) - \frac{f(x+h)-f(x)}{h} = f'(x) - R_h^{(1)}(x) = c_1(x)h + c_2(x)h^2 + c_3(x)h^3 + \cdots$ but suppose we cannot compute the $c_k$.

$R_{h/2}^{(1)}(x) = f'(x) - c_1\frac{h}{2} - c_2\frac{h^2}{4} - c_3\frac{h^3}{6} - \cdots$. We can use this to eleminate $c_1$:

$$2R_{h/2}^{(1)}(x) - R_h^{(1)}(x) = f'(x) - c_2'h^2 - c_3'h^3 - \cdots$$

So the error is of order $O(h^2)$. Then

$$f'(x) = R_h^{(2)}(x) + O(h^2)$$

where $R_h^{(2)} = 2R_{h/2}^{(1)}(x) - R_h^{(1)}(x)$.

Now, $R_{h/2}^{(2)}(x) = f'(x) - c_2' \frac{h^2}{4} - c_3' \frac{h^3}{8}$ and we use this to eliminate $c_2'$:

$$4R_{h/2}^{(2)}(x) - R_h^{(2)}(x) = 3f'(x) + O(h^3)$$

So we set

$$R_h^{(3)} := \frac{2^2 R_{h/2}^{(2)}(x) - R_h^{(2)}(x)}{2^2 - 1} = f'(x) + O(h^3)$$

**Remark.**

1. We only need to specify the powers of $h$, not the coefficients $c_k$ as long as they are **non-zero**. So when using centred difference to calculate $R_h^{(1)}$ then $R_h^{(2)}$, $R_h^{(3)}$ will be different.

2. Richardson extrapolation works with many other approximation methods involving a small parameter (not just differentiation).

3. There is no standard notation for this method.

4. Some series expansions have irregular/non-integer powers, e.g. the Airy function $Ai(x)$ which is a solution of $\frac{d^2 f}{dx^2} = xf(x)$.

# 4 Linear systems

Given $A \in M_n(\mathbb{R})$ and $\underline{x} \in \mathbb{R}^n$, we want to solve for $\underline{x}$:

$$A\underline{x} = \underline{b}$$

We want to minimise the error when computing this with floating points, and minimise the amount of computation for: large $n$, for different $\underline{b}$ (with the same $A$) and when $A$ has particular forms.

**Definition 4.0.1.** The **transpose** of a square matrix $A$, $A^T$ is defined as

$$\left(A^T\right)_{i,j} = A_{j,i}$$

**Definition 4.0.2.** $A$ is **symmetric** if $A = A^T$.

**Definition 4.0.3.** $A$ is **skew-symmetric** if $A = -A^T$.

**Definition 4.0.4.** $A$ is **non-singular** if for every $\underline{b}$, for some $\underline{x}$, $A\underline{x} = \underline{b}$.

**Definition 4.0.5.** $A$ is **positive definite** if $(A\underline{x}) \cdot \underline{x} > 0 \quad \underline{x} \neq \underline{0}$.

**Definition 4.0.6.** $A$ is **positive semi-definite** if $(A\underline{x}) \cdot \underline{x} \geq 0 \quad \forall x$.

**Lemma 4.0.7.** If $A$ is positive definite, then $A$ is non-singular.

*Proof.* Omitted. $\qquad\square$

**Lemma 4.0.8.** The converse of the above lemma is false.

*Proof.* $A$ is non-singular $\Rightarrow -A$ is non-singular, but $A$ is positive-definite $\Rightarrow -A$ is negative definite. $\qquad\square$

**Definition 4.0.9.** A matrix $A$ is **lower triangular** if $A_{i,j} = 0 \quad \forall j > i$.

**Definition 4.0.10.** $A$ is **upper triangular** if $A_{i,j} = 0 \quad \forall j < i$.

**Remark.** We often write $U$ for an upper triangular matrix and $L$ for a lower triangular matrix.

**Definition 4.0.11.** To solve $Ux = b$, we can use **backward substitution**. Starting from the last equation,

$$x_n = \frac{1}{U_{n,n}} b_n$$

$$x_{n-1} = \frac{1}{U_{n-1,n-1}}(b_{n-1} - U_{n-1,n}x_n)$$

$$\vdots$$

$$x_j = \frac{1}{U_j}\left(b_j - \sum_{i=j+1}^{n} U_{j,i}x_i\right)$$

**Definition 4.0.12.** Similarly, we can solve $L\underline{x} = \underline{b}$ for a lower triangular matrix $L$ with forward substitution. Starting from the first equation,

$$x_1 = \frac{b_1}{L_{1,1}}$$

$$\vdots$$

$$x_j = \frac{1}{L_{j,j}}\left(b_j - \sum_{i=1}^{j-1} L_{j,i}x_k\right)$$

**Remark.** If $U_{j,j} = 0$ or $L_{j,j} = 0$ for some $j$, then this method doesn't work. This is expected, because in this case $\det U = 0$ (or $\det L = 0$).

**Definition 4.0.13.** If $A$ is neither upper nor lower triangular, then we can transform $A$ into a (usually) upper triangular matrix, by **Gaussian elimination**.

The following operations leave the solution $\underline{x}$ unchanged:

1. Swapping two rows

2. Adding a scalar multiple of a row to another row.

## 4.1 Computational complexity

**Definition 4.1.1.** For simplicity, we assume that each elementary floating point operation takes 1 unit of time, called 1 **flop**.

**Remark.** In practice, binary64 multiplication takes roughly 3 times as long as addition or multiplication, and binary64 division takes roughly 10 times as long as addition or subtraction.

**Definition 4.1.2.** Let $f(n) > 0$ and $g(n) > 0$ for large $n$. We write

$$f(n) \sim o(g(n)) \quad \text{if } \lim_{n \to \infty} \frac{f(n)}{g(n)} = 0$$

**Definition 4.1.3.** Let $f(n) > 0$ and $g(n) > 0$ for large $n$. We write

$$f(n) \sim O(g(n)) \quad \text{if } \limsup_{n \to \infty} \frac{f(n)}{g(n)} < \infty$$

Equivalently,

$$f(n) \sim O(g(n)) \quad \text{if } \exists C, \exists N, \forall n \le N, f(n) \le Cg(n)$$

**Remark.** From these definitions, we have

$$f(n) \sim o(g(n)) \implies f(n) \sim O(g(n))$$

**Example 4.1.4.**

- $100n^3 + 10^6 n^2 \sim O(n^3)$

- $n! + 10^{100} n^{100} \sim O(n!)$

- $n! \sim O(n^{n+\frac{1}{2}} e^{-n})$

**Remark.** In this module, we will be calculate the computational complexity to **leading order**, e.g. we distinguish $3n^2 + 5n$ form $30n^2$ but not from $3n^2 - 2n$.

**Proposition 4.1.5.** Backwards substitution on $U$ where $U$ is an $n \times n$ matrix is an $O(n^2)$ operation.

*Proof.*

- Computing $x_n = b_n / U_{n,n}$ takes 1 flop.

- Computing $x_{n-1} = \frac{(b_{n-1} - U_{n-1,n} x_n)}{U_{n-1,n-1}}$ takes 3 flops.

- ⋮

- Computing $x_1 = \frac{1}{U_1}\left(b_1 - \sum_{i=j+2}^{n} U_{1,i} x_i\right)$ takes $2n - 1$ flops.

So in total, backward substitution takes $1 + 3 + \cdots + 2n - 1 = n^2$ flops. □

**Proposition 4.1.6.** The number of flops needed for solving $Ax = b$ where $A$ is an $n \times n$ matrix is
$$\frac{2}{3}n^3 + \frac{3}{2}n^2 - \frac{7}{6}n$$

*Proof.* To zero the first column for rows $i \in \{2, \ldots, n\}$:

- $\alpha_{i,1} = -A_{i,1}/A_{1,1}$ (1 flop).

- $A_{i,j} \to A_{i,j} + \alpha_{i,1} A_{1,j}$ for $j \in \{2, \ldots, n\}$ and $A_{i,1} = 0$ ($2(n-1)$ flops).

- $b_i \to b_i + \alpha i, 1 b_1$ (2 flops).

This is $2n + 1$ flops in total for each row, and this is done for $n - 1$ rows, so in total there are $(2n + 1)(n - 1)$ flops. For row 2, there are $(2n - 1)(n - 2)$ flops.

In general, to zero column $k$ for rows $i \in \{k + 1, \ldots, n\}$:

- $\alpha_{i,k} = -A_{i,k}/A_{k,k}$ (1 flop).

- $A_{i,j} \to A_{i,j} + \alpha_{i,k} A_{k,j}$ for $j \in \{k + 1, \ldots, n\}$ ($2(n - k)$ flops).

- $b_i \to b_i + \alpha_{i,k} b_k$ (2 flops).

This is $2n - 2k + 3$ flops in total for each row, so for the $n - k$ rows, in total there are $(2n - 2k + 3)(n - k)$ flops.

So the total number of flops for Gaussian elimination is

$$\sum_{k=1}^{n-1}(2n - 2k + 3)(n - k) = \sum_{k=1}^{n-1} 2n^2 - 2kn + 3n - 2nk + 2k^2 - 3k$$

$$= (n - 1)(2n^2 + 3n) - (4n + 3)\sum_{k=1}^{n-1} k + 2\sum_{k=1}^{n-1} k^2$$

$$= (n - 1)(2n^2 + 3n) - (4n + 3)\frac{n(n - 1)}{2}$$

$$+ \frac{1}{3}(n - 1)n(2n - 1)$$

$$= \frac{2}{3}n^3 + \frac{1}{2}n^2 - \frac{7}{6}n$$

Adding the $n^2$ flops from back substitution, in total to solve $Ax = b$, the number of flops needed is
$$\frac{2}{3}n^3 + \frac{3}{2}n^2 - \frac{7}{6}n$$

□

**Remark.** Gaussian elimination with (row) pivotting can transfomr every non-singular matrix into an upper-triangular matrix.

13

## 4.2   Pivoting and roundoffs

**Definition 4.2.1.** To reduce round-off errors, we define **row/partial pivoting**:

- When zeroing column $k$, look for $A_{j,k}$ with $j \in \{k, \ldots, n\}$ with the largest absolute value, and swap row $j$ with row $k$.

**Remark.** Multiplying a pivot row by a large constant does not improve reduction of round-off errors. The **ratio** of the pivot element to the other elements in the row is more important.

**Remark.** For better results, we can also perform column/full pivoting (so we swap rows and columns), but this is much more complicated.

## 4.3   LU decomposition

**Lemma 4.3.1.** Steps of the Gaussian elimination process can be written as matrix multiplication. During Gaussian elimination, let $A^{(k)}$ be the matrix during Gaussian elimination which has zeros below the diagonal in the first $k$ columns.

$$A^{(s)} = F^{(s)} A^{(s-1)}$$

where

$$F_{i,j}^{(s)} = \delta_{i,j} - f_i^{(s)} e_j^{(s)}$$
$$f_i^{(s)} = (0, \ldots, 0, A_{s+1,s}^{s-1}/A_{s,s}^{(s-1)}, \ldots, A_{n,s}^{(s-1)}/A_{s,s}^{(s-1)})$$
$$e_j^{(s)} = \delta_{s,j}$$

*Proof.* To zero column $s$ below the leading diagonal, we do

$$A_{i,j}^{(s)} = A_{i,j}^{(s-1)} - \frac{A_{i,s}^{(s-1)}}{A_{s,s}^{(s-1)}} A_{s,j}^{(s-1)}$$

for $i \in \{s+1, \ldots, n\}$ and $j \in \{s, \ldots, n\}$.

$$A_{i,j}^{(s)} = \sum_k \left( \delta_{i,k} A_{k,j}^{(s-1)} - \frac{A_{i,s}^{(s-1)}}{A_{s,s}^{(s-1)}} \delta_{s,k} A_{k,j}^{(s-1)} \right)$$
$$= \sum_k \left( \delta_{i,k} - \frac{A_{i,s}^{(s-1)}}{A_{s,s}^{(s-1)}} \cdot \delta_{s,k} \right) A_{k,j}^{(s-1)}$$

$\square$

**Definition 4.3.2.** For vectors $u$ and $v$, the **tensor product** of $u$ and $v$ is a matrix defined by

$$(u \otimes v)_{i,j} = u_i v_j$$

**Definition 4.3.3.** A **Frobenius matrix** of index $s$ is a unit lower triangular matrix whose only (possibly) non-zero elements are subdiagonal elements in column $s$ (other than the leading diagonal, which is all 1s).

**Lemma 4.3.4.** Let $F^{(s)} = I - f^{(s)} \otimes e^{(s)}$. Then its inverse is

$$G^{(s)} = I + f^{(s)} \otimes e^{(s)}$$

*Proof.*

$$\sum_k F^{(s)}_{i,k} G^{(s)}_{k,j} = \sum_k \left(\delta_{i,k} - f^{(s)}_i e^{(s)}_k\right)\left(\delta_{k,j} + f^{(s)}_k e^{(s)}_j\right)$$

$$= \sum_k \left(\delta_{i,k} - f^{(s)}_i \delta_{s,k}\right)\left(\delta_{k,j} + f^{(s)}_k \delta_{s,j}\right)$$

$$= \sum_k \left(\delta_{i,j} + f^{(s)}_i \delta_{s,j} - f^{(s)}_i \delta_{s,j} - f^{(s)}_i f^{(s)}_k \delta_{s,k}\delta_{s,j}\right)$$

$$= \delta_{i,j} - f^{(s)}_i f^{(s)}_s \delta_{s,j}$$

But $f^{(s)}_s = 0$ which completes the proof. $\square$

**Lemma 4.3.5.** Let $G^{(s)} = I + f^{(s)} \otimes e^{(s)}$ as above. Then for every $s \in \{1, \ldots, n-1\}$,

$$G^{(1)}G^{(2)}\cdots G^{(s)} = I + \sum_{r=1}^{s} f^{(r)} \otimes e^{(r)}$$

*Proof.* Use induction on $s$.

- For $s = 1$, $G^{(1)} = I + f^{(1)} \otimes e^{(1)}$ trivially.

- Assume the statement is true for some $s$.

- For $s+1$,

$$\left(G^{(1)}\cdots G^{(s)}G^{(s+1)}\right)_{i,j} = \sum_k \left(G^{(1)}\cdots G^{(s)}\right)_{i,k} G^{(s+1)}_{k,j}$$

$$= \left(\sum_k \left(\delta_{i,k}\sum_{r=1}^{s} f^{(r)}_i \delta_{r,k}\right)\right)\left(\delta_{k,j} + f^{(s+1)}_k \delta_{s+1,j}\right)$$

$$= \delta_{i,j} + \sum_{r=1}^{s} f^{(r)}_i \delta_{r,j} + f^{(s+1)}_i \delta_{s+1,j}$$

$$+ \sum_k \left(f^{(s+1)}_k \delta_{s+1,j}\sum_{r=1}^{s} f^{(r)}_i \delta_{r,k}\right)$$

$$= \delta_{i,j} + \sum_{r=1}^{s+1} f^{(r)}_i \delta_{r,j}$$

which completes the induction.

$\square$

**Theorem 4.3.6.** Suppose that $A$ can be reduced to an upper-triangular matrix $U$ by Gaussian elimination without pivoting. During Gaussian elimination, let $A^{(k)}$ be the matrix during Gaussian elimination which has zeros below the diagonal in the first $k$ columns. Then there is a **unit** lower triangular matrix $L$ such that

$$A = LU$$

Also, the subdiagonal elements of $L$ are the coefficients used in the reduction of $A$:

$$\forall i > j, L_{i,j} = A^{(j-1)}{}_{i,j}/A^{(j-1)}{}_{j,j}$$

*Proof.* We write the Gaussian elimination process as

$$U = A^{(n-1)} = F^{(n-1)}A^{(n-2)}$$
$$= F^{(n-1)}F^{(n-2)}A^{(n-3)}$$
$$= F^{(n-1)}F^{(n-2)}\cdots F^{(1)}A^{(0)}$$

Now left-multiplying by $L = G^{(1)}\cdots G^{(n-1)}$:

$$LU = G^{(1)}\cdots G^{(n-1)}F^{(n-1)}F^{(n-2)}\cdots F^{(1)}A^{(0)}$$

But $G^{(k)}F^{(k)} = I$ by Lemma 4.3.4 so this simplifies to $LU = A$. $\qquad\square$

**Theorem 4.3.7.** Every non-singular square matrix $A$ can be written as

$$PA = LU$$

where $L$ and $U$ are lower and upper triangular matrices and $P$ is a permutation matrix (exactly one non-zero element in each row and column).

*Proof.* Too difficult.batyter $\qquad\square$

**Remark.** Once we have $L$ and $U$, we can solve $Ax = LUx = b$ with

1. Solve $Ly = b$ for $y$ using forward substitution.

2. Solve $Ux = y$ for $x$ using backward substitution.

**Example 4.3.8.** $A = [[1,2,3,6],[2,8,6,5],[-4,-8,0,0],[0,12,9,-6]]$. Then $A^{(1)} = [[1,2,3,6],[0,4,0,-7$
$A^{(2)} = [[1,2,3,6],[0,4,0,-7],[0,0,12,24],[0,0,9,15]], A^{(3)} = [[1,2,3,6],[0,40,0,-7],[0,0,12,24],[0,0,0,-$
Then $L = [[1,0,0,0],[L_{2,1},1,0,0],[L_{3,1},L_{3,2},1,0],[L_{4,1},L_{4,2},L_{4,3},1]] = [[1,0,0,0],[2,1,0,0],[-4,0,1,0]$
Then if $b = [1,2,4,4]$, solving $Ly = b$ with forward substitution gives $y = [1,0,8,-2]$.
Then solving $Ux = y$ for $x$ gives $x = [-\frac{10}{3}, \frac{7}{6}, -\frac{2}{3}, \frac{2}{3}]$.

**Remark.** Comparing LU decomposition with Gaussian elimination and back substitution, LU decomposition involves less work when solving for different values of $b$ (same values of $A$). Gaussian elimination and back substitution is $O(n^3)$ for each value of $b$. LU decomposition is $O(n^3)$ for the first value of $b$, but then solving for different values of $b$ afterwards is $O(n^2)$.

**Remark.** LU decomposition is less prone to round off errors than computing $A^{-1}$ and is less computationally expensive for some values of $A$:

- When $A$ is tridiagonal, $L$ and $U$ are as well, but $A^{-1}$ generally isn't.

- When $A$ is sparse (most elements are 0), $L$ and $U$ are also sparse, but $A^{-1}$ generally isn't.

**Proposition 4.3.9.** If $L$ is an $n \times n$ lower triangular, invertible matrix, then $L^{-1}$ is also lower triangular.

*Proof.* Let $M = L^{-1}$, then $LM = I$. Since $L^{-1}$ exists, $L_{i,i} \neq 0$ and $M_{i,i} \neq 0$ $\forall 1 \leq i \leq n$.
We have $I_{i,j} = \delta_{i,j}$ for every $i, j$. Then $0 = I_{1,j} = \delta_{1,j} = L_{1,1}M_{1,j} \Rightarrow M_{1,j} = 0$ $\forall j > 1$ (since $L_{1,1} \neq 0$).

Now for row 2, $1 = \delta_{2,2} = L_{2,1}M_{1,2} + L_{2,2}M_{2,2} + L_{2,3}M_{3,2} + \cdots$ but $L_{2,j} = 0 \forall j > 2$ since $L$ is lower triangular so $1 = L_{2,1}M_{1,2} + L_{2,2}M_{2,2}$. Similarly, $0 = \delta_{2,3} = L_{2,1}M_{1,3} + L_{2,2}M_{2,3} + L_{2,3}M_{3,3} + \cdots = L_{2,1}M_{1,3} + L_{2,2}M_{2,3} = L_{2,2}M_{2,3} \Rightarrow M_{2,3} = 0$. Generally,

$$\delta_{2,j} = \sum_{k=1}^{n} L_{2,k}M_{k,j} = L_{2,1}M_{1,j} + L_{2,2}M_{2,j} + L_{2,3}M_{3,j} + \cdots = L_{2,1}M_{1,j} + L_{2,2}M_{2,j}$$

But for $j > 2$, $M_{1,2}, \ldots, M_{1,n} = 0$. Since $L_{2,2} \neq 0$, $M_{2,j} = 0 \; \forall j > 2$.

This process continues for the rest of the rows by induction: assume that $M_{i,j} = 0 \; \forall i \in \{1, \ldots, s-1\}, j > i$. Now we show that this holds for $i \in \{1, \ldots, s\}$.

$$\delta_{s,j} = \sum_{k=1}^{n} L_{s,k}M_{k,j}$$
$$= \sum_{k=1}^{s} L_{s,k}M_{k,j}$$

For $j > s$,

$$\delta_{s,j} = \sum_{k=1}^{s} L_{s,k}M_{k,j}$$
$$= L_{s,1}M_{1,j} + \cdots + L_{s,s-1}M_{s-1,j} + L_{s,s}M_{s,j}$$

But by assumption, $L_{s,1}M_{1,j} + \cdots + L_{s,s-1}M_{s-1,j} = 0$, so $\delta_{s,j} = L_{s,s}M_{s,j} = 0 \Rightarrow M_{s,j} = 0$ as $L_{s,s} \neq 0$. $\qquad \square$

## 4.4   Normed vector spaces

**Definition 4.4.1.** Given a vector space $V$, a **norm** $||.|| : V \to \mathbb{R}_+$ satisfies the following properties:

1. $||x + y|| \leq ||x|| + ||y|| \quad \forall x, y \in V$.

2. $||c \cdot x|| = |c| \, ||x|| \quad \forall x \in V, c \in \mathbb{C}$.

3. $||x|| = 0 \iff x = 0$ (if this is not satisfied, the norm is called a **seminorm**).

**Example 4.4.2.** Some examples of norms:

1. The $L_2$ norm,
$$||x||_2 := \left( \sum_{k=1}^{n} |x_k|^2 \right)^{1/2}$$

2. The $L_1$ norm,
$$||x||_1 := \sum_{k=1}^{n} |x_k|$$

3. $||x||_\infty := \max_{k \in \{1,\ldots,n\}} |x_k|$

4. For $p \in [1, \infty)$,
$$||x||_p := \left( \sum_{k=1}^{n} |x_k|^p \right)^{1/p}$$

**Definition 4.4.3.** Let $\tilde{V}$ be the set of $n \times n$ matrices over $\mathbb{R}$, then

$$||A||_{\text{row}} := \max_i \sum_{j=1}^{n} |A_{i,j}|$$

$$||A||_{\text{col}} := \max_j \sum_{i=1}^{n} |A_{i,j}|$$

$$||A||_{\text{Frob}} := \left( \sum_{i,j} |A_{i,j}|^2 \right)^{1/2}$$

are the **row-sum norm**, the **column-sum norm**, and the **Frobenius matrix norm**, respectively.

**Definition 4.4.4.** A matrix norm $||.||$ on $\tilde{V}$ is **submultiplicative** if

$$\forall A, B \in \tilde{V}, \quad ||AB|| \leq ||A|| \cdot ||B||$$

**Definition 4.4.5.** Given a vector norm $||.||_{\text{vec}} : V \to \mathbb{R}_+$, the **matrix norm** induced by $||.||_{\text{vec}}$ is defined as

$$||A||_{\text{mat}} = \sup_{x \neq 0} \frac{||Ax||_{\text{vec}}}{||x||_{\text{vec}}} = \max_{||x||_{\text{vec}}=1} ||Ax||_{\text{vec}}$$

**Remark.** Often, we omit vec or mat when the meaning of the norm is clear from context.

**Example 4.4.6.** By linearity of vector norms,

$$||Ax||_{\star} \leq ||A||_{\star} ||x||_{\star}$$

where $|||A||_{\star}$ is the matrix norm induced by the vector norm $||.||_{\star}$.

**Example 4.4.7.** Let $V = \mathbb{R}^2$ and $||.|| = ||.||_2$. Compute $||A||_2$ where

$$A = \begin{pmatrix} 3 & 1 \\ -5 & -1 \end{pmatrix}$$

By definition,
$$\begin{aligned}
||A||_2 &= \max_{||x||_2=1} ||Ax||_2 \\
&= \max_{\theta \in [0,2\pi]} ||A \cdot (\cos(\theta), \sin(\theta))||_2 \\
&= \max_{\theta \in [0,2\pi]} ||(3\cos(\theta) + \sin(\theta), -5\cos(\theta) - \sin(\theta))||_2
\end{aligned}$$

Let

$$\begin{aligned}
f(\theta) &= (3\cos(\theta) + \sin(\theta))^2 + (5\cos(\theta)) + \sin(\theta))^2 \\
&= 34\cos(\theta)^2 + 2\sin(\theta)^2 + 16\cos(\theta)\sin(\theta) \\
&= 16\cos(2\theta) + 8\sin(2\theta) + 18
\end{aligned}$$

so $f'(\theta) = -32\sin(2\theta) + 16\cos(2\theta)$, then $f'(\theta) = 0 \Rightarrow \tan(2\theta) = 1/2$ which gives

$$\max_{\theta} f(\theta) = 18 + \frac{40}{\sqrt{5}}$$

which gives

$$||A||_2 = \sqrt{18 + \frac{40}{\sqrt{5}}}$$

18

**Definition 4.4.8.** A matrix $A$ is **normal** if $AA^T = A^TA$.

**Theorem 4.4.9.** If a matrix $A$ is normal, then $A$ is diagonalisable.

*Proof.* Too long. $\qquad\square$

**Theorem 4.4.10.** If $A$ is symmetric, then its eigenvalues are real, and

$$||A||_2 = \max\{|\lambda| : \lambda \text{ is an eigenvalue of } A\}$$

For non-symmetric $A$,

$$||A||_2 = \max\{|\lambda| : \lambda \text{ is an eigenvalue of } AA^T\}$$

*Proof.* Too long. $\qquad\square$

**Proposition 4.4.11.** The vector 1-norm induces the matrix column-sum norm.

*Proof.* $||x||_1 = \sum_i |x_i|$ and $||A||_{\text{col}} = \max_j \sum_i |A_{i,j}|$.

$$
\begin{aligned}
||Ax||_1 &= \sum_i \left| \sum_j A_{i,j} \right| \\
&\leq \sum_i \sum_j |A_{i,j} x_j| \\
&= \sum_i \sum_j |A_{i,j}||x_j| = \sum_j \left( \sum_i |A_{i,j}| \right) |x_j| \\
&\leq \left( \max_j \sum_i |A_{i,j}| \right) \sum_j |x_j| = ||A||_{\text{col}} ||x||_1
\end{aligned}
$$

Hence,
$$||A||_1 \leq ||A||_{\text{col}}$$

and also
$$
\begin{aligned}
||A||_1 &= \max_{||x||_1=1} ||Ax||_1 \\
&\geq \max_k ||Ae^{(k)}||
\end{aligned}
$$

where $e_j^{(k)} = \delta_{j,k}$. So

$$
\begin{aligned}
||A||_1 &\geq \max_k \sum_i \left| \sum_j A_{i,j} \delta_{j,k} \right| \\
&= \max_k \sum_i |A_{i,k}| = ||A||_{\text{col}}
\end{aligned}
$$

Hence $||A||_1 = ||A||_{\text{col}}$. $\qquad\square$

**Example 4.4.12.** Let

$$A = \begin{bmatrix} 3 & 1 & -4 \\ 1 & -5 & 9 \\ -2 & 6 & 5 \end{bmatrix}$$

so $||A||_{\text{col}} = \max\{3+1+2, 1+5+6, 4+9+5\} = 18$. The max is attained at column 3. Let $\underline{x} = (0,0,1)$ then $A\underline{x} = (-4,9,5)$ and $||A\underline{x}||_1 = 4+9+5 = 18$.

**Proposition 4.4.13.** The vector $\infty$-norm induces the matrix row-sum norm.

*Proof.* We want to show that

$$\max_{||x||_\infty=1} ||Ax||_\infty = \max_i \sum_j |A_{i,j}|$$

First, we show $||A||_\infty \leq ||A||_{\text{row}}$ for every $A$ and $x$:

$$||Ax||_\infty = \max_i \left|\sum_j A_{i,j}x_j\right|$$

$$\leq \max_i \sum_j |A_{i,j}||x_j|$$

$$\leq \max_i \left\{ \max_k |x_k| \sum_j |A_{i,j}| \right\}$$

$$= \max_k |x_k| \max_i \sum_j |A_{i,j}|$$

$$= ||x||_\infty ||A||_{\text{row}}$$

Assuming that $x \neq 0$, we divide by $||x||_\infty$ and take the maximum:

$$||A||_\infty = \max_{||x||_\infty} \frac{||Ax||_\infty}{||x||_\infty} \leq \max_{||x||_\infty \neq 0} ||A||_{\text{row}} = ||A||_{\text{row}}$$

To show equality, for every $A$, we need to find an $x$ that gives equality in the inequalities above. By linearity, we can take $||x||_\infty = 1$. We want to find an $x$ such that

$$\max_k \left|\sum_j A_{k,j}x_j\right| = \max_i \sum_j |A_{i,j}|$$

Let $i^*$ be an $i$ that realises the maximum. Then let

$$x_j = \text{sgn}(A_{i^*,j})$$

Let $k = i^*$, then

$$\sum_j A_{i^*,j}x_j = \sum_j A_{i^*,j} \cdot \text{sgn}(A_{i^*,j}) = \sum_j |A_{i^*,j}| = \max_i \sum_j |A_{i,j}|$$

$\square$

**Remark.**

$$||A|| = \sup_{y \neq 0} \frac{||Ay||}{||x||} \quad \forall x \neq 0$$

Therefore

$$||Ax|| \leq ||A|| \cdot ||x|| \quad \forall x$$

**Example 4.4.14.** (Problems class) Show the equivalence of the induced norm definitions.

$$\sup_{x \neq 0} \frac{||Ax||}{||x||} = \sup_{x \neq 0} \frac{||Ax||/||x||}{||x||/||x||}$$

$$= \sup_{x \neq 0} \frac{||A(x/||x||)||}{||(x/||x||)||}$$

$$= \sup_{x \neq 0} ||A(x/||x||)||$$

$$= \max_{||y||=1} ||Ay||$$

**Example 4.4.15.** (Problems class) Show that $||AB|| \leq ||A|| \cdot ||B||$ for all matrices $A, B$ and every induced norm $|| \cdot ||$.

$$\sup_{x \neq 0, Bx \neq 0} \frac{||ABx||}{||x||} = \sup_{x \neq 0, Bx \neq 0} \frac{||ABx||}{||Bx||} \frac{||Bx||}{||x||}$$

$$\leq \sup_{Bx \neq 0} \frac{||ABx||}{||Bx||} \cdot \sup_{x \neq 0} \frac{||Bx||}{||x||}$$

$$= ||A|| \cdot ||B||$$

When $Bx = 0$,

$$0 = ||ABx|| = ||A|| \cdot ||Bx|| = 0$$

so we can remove the condition $Bx \neq 0$ on the sups in the equation above.

## 4.5  Errors and condition numbers

In the case where we want to solve $Ax = b$ for $x$, but with some floating po0int error in $b$, we have $b + \delta b$ instead of $b$, where $\delta b$ is unknown but bounded (small). So we have

$$A(x + \delta x) = b + \delta b$$

It is not always the case that $\delta x$ is also small. We want to find bounds for $\delta x$ in terms of $\delta b$. We have that $A\delta x = \delta b$. Assuming $A^{-1}$ exists, $\delta x = A^{-1}\delta b$. Using any norm $|| \cdot ||_*$,

$$||\delta x||_* = ||A^{-1}\delta b||_* \leq ||A^{-1}||_*||\delta b||_*$$

Assuming $x \neq 0$, divide by $||x||_*$:

$$\frac{||\delta x||_*}{||x||_*} \leq ||A^{-1}|| \frac{||\delta b||_*}{||x||_*}$$

Now since $Ax = b$,

$$||b||_* = ||Ax||_* \leq ||A||_*||x||_* \implies \frac{1}{||x||_*} \leq ||A||_* \frac{1}{||b||_*}$$

Therefore,

$$\frac{||\delta x||_*}{||x||_*} \leq ||A||_*||A^{-1}||_* \frac{||\delta b||_*}{||b||_*}$$

This means the relative error in $x$ is less than a constant multiplied by the relative error in $b$.

**Definition 4.5.1.** For any norm $|| \cdot ||_*$, we define

$$\kappa_*(A) := ||A||_*||A^{-1}||_*$$

which is called the **\*-condition** number of the matrix $A$.

**Example 4.5.2.** (Problems class) Prove that $\kappa(A) \geq 1$ for every non-singular matrix $A$ and for every induced matrix norm.

$$1 = ||I|| = ||AA^{-1}|| \leq ||A|| \cdot ||A^{-1}||$$

by Example 4.4.15.

**Example 4.5.3.** (Problems class) Prove that

$$||x||_p = \left( \sum_i |x_i|^p \right)^{1/p}$$

is a norm for every $p \in [1, \infty)$.

$$
\begin{aligned}
(||x + y||_p)^p &= \sum_i |x_i + y_i|^p \\
&= \sum_i |x_i + y_i| \cdot |x_i + y_i|^{p-1} \\
&\leq \sum_i (|x_i| + |y_i|)|x_i + y_i|^{p-1} \\
&= \sum_i |x_i| \cdot |x_i + y_i|^{p-1} + \cdots \\
&\leq \left( \sum_i |x_i|^r \right)^{1/r} \left( \sum_j |x_j + y_j|^{q(p-1)} \right)^{1/q} + \cdots
\end{aligned}
$$

by the Holder inequality (TODO: define this inequality). Let $r = p, q = 1/(1 - 1/p) = p/(p-1)$. Then

$$
\begin{aligned}
(||x + y||_p)^p &\leq \left( \sum_i |x_i|^p \right)^{1/p} \left( \sum_j |x_j + y_j|^{(p-1)p/(p-1)} \right)^{(p-1)/p} + \cdots \\
&= ||x||_p(||x + y||)_p^{p-1} + ||y||_p(||x + y||)_p^{p-1}
\end{aligned}
$$

## 4.6  $L^2$ norm, eigenvalues and diagonalisability

**Theorem 4.6.1.** Let $\{\lambda_j\}$ be the eigenvalues of a matrix $A$. Then

$$\max\{|\lambda_j|\} < 1 \iff A^k \to 0 \quad \text{as } k \to \infty$$

**Example 4.6.2.** Let

$$A = \begin{bmatrix} 0.99 & 10^6 \\ 0 & 0.99 \end{bmatrix} \implies A^{128} = \begin{bmatrix} 0.2765\ldots & 3.751 \cdot 10^6 \\ 0 & 0.2765\ldots \end{bmatrix}$$

This does not seem to agree with Theorem 4.6.1, however the theorem is correct, as products and powers of matrices are more complicated than those of complex numbers. For example, for every $z \in \mathbb{C}$, $|z|^k = |z^k|$ but generally $||A^k|| \neq ||A||^k$ for a matrix $A$ and a norm $||\cdot||$. $||A||^k$ could grow before it then tends to zero. This occurs especially when $A$ is not diagonalisable ($A$ is defective, meaning that its eigenvectors do not span $\mathbb{R}^n$).

In fact, Theorem 4.6.1 is easy to prove for diagonalisable $A$, but difficult for defective $A$.

**Definition 4.6.3.** The **spectral radius** of a square matrix $A$ is defined as

$$\rho(A) := \max\{|\lambda_j|\}$$

where the $\lambda_j$ are the eigenvalues of $A$.

**Proposition 4.6.4.** If $A$ is symmetric (or hermitian), then

- All its eigenvalues are real, so $Bv_i = \lambda_i v_i$ where $\lambda_i \in \mathbb{R}$.

- Its eigenvectors $v_i$ form a basis of $\mathbb{R}^n$: every vector $x \in \mathbb{R}^n$ can be written as a linear combination of the eigenvectors of $A$, i.e. $x = \sum_i c_i v_i$ where $v_i$ are the eigenvectors of $A$ and $c_i$ are real numbers.

- Its eigenvectors can be chosen to be orthonormal, i.e. $v_i \cdot v_j = \delta_{ij}$.

- Parseval's identity holds:

$$(||x||_2)^2 = \left( \sum_i c_i v_i, \sum_j c_j v_j \right) = \sum_{i,j} c_i c_j (v_i, v_j) = \sum_i c_i^2$$

**Definition 4.6.5.** A square matrix $A$ is **normal** if it commutes with its transpose:

$$AA^T = A^T A$$

($AA^* = A^* A$ for complex matrices).

**Definition 4.6.6.** A matrix $U$ is unitary if

$$U^* = U^{-1}$$

**Definition 4.6.7.** A matrix $A$ is unitarily diagonalisable if $A = UDU^*$ for some unitary matrix $U$ and diagonal matrix $D$.

**Theorem 4.6.8.** (**The spectral theorem**) A matrix $A$ is unitarily diagonalisable if and only if it is normal.

**Example 4.6.9.** Let $a \neq 1$ and

$$A = \begin{bmatrix} 1 & a \\ 0 & 1 \end{bmatrix}$$

is not normal, as

$$AA^T = \begin{bmatrix} 1 + a^2 & a \\ a & 1 \end{bmatrix} \neq A^T A = \begin{bmatrix} 1 & a \\ a & 1 + a^2 \end{bmatrix}$$

$\lambda = 1$ is the only eigenvalue, with eigenvector $(1, 0)$. $A$ is not unitarily diagonalisable.

**Theorem 4.6.10.** The induced 2-norm of a matrix $A$ is

$$\sqrt{\rho(A^T A)} = \max\{\sqrt{|\lambda|} : \lambda \text{ is an eigenvalue of } A^T A\}$$

*Proof.* For every matrix $A$, $A^T A$ is symmetric, so by Proposition 4.6.4 its eigenvalues are real and its eigenvectors, $\{v_i\}$ for $i \in \{1, \ldots, n\}$, are orthonormal and span $\mathbb{R}^n$. So for every $x \in \mathbb{R}^n$,

$$x = \sum_i c_i v_i$$

for some $c_i \in \mathbb{R}$, $A^T A v_i = c_i v_i$ and $(v_i, v_j) = \delta_{i,j}$.

$$(||Ax||_2)^2 = (Ax, Ax) = (A^T Ax, x) = \left( A^T A \sum_i c_i v_i, \sum_j c_j v_j \right)$$

$$= \left( \sum_i \lambda_i c_i v_i, \sum_k c_j v_j \right)$$

$$= \sum_{i,j} \lambda_i c_i c_j (v_i, v_j)$$

$$= \sum_i \lambda_i c_i^2$$

23

We claim that $\lambda_j \geq 0 \; \forall j$. Indeed, if $\lambda_s < 0$, let $x = v_s$, then

$$(||Av_s||_2)^2 = \lambda_s < 0$$

Write $0 \leq \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$ (counting multiplicities). Then

$$(||Ax||_2)^2 = \sum_i \lambda_i c_i^2 \leq \lambda_n \sum_i c_i^2 = \lambda_n (||x||_2)^2$$

which gives

$$||A||_2 = \sup_{x \neq 0} \frac{||Ax||_2}{||x||_2} \leq \sqrt{\lambda_n}$$

Taking $x = v_n$ gives $(||Ax_n||_2)^2 = \lambda_n(||x_n||_2)^2$, which gives equality in the above inequality.
$\square$

**Remark.** $\rho(A^T A) = \rho(A A^T)$.

**Example 4.6.11.** Let

$$A = \begin{bmatrix} 1 & a \\ 0 & 1 \end{bmatrix} \implies A^T A = \begin{bmatrix} 1 & 0 \\ 0 & 1 + a^2 \end{bmatrix}$$

The eigenvalues of $A^T A$ are given by $0 = (1 - \lambda)(1 + a^2 - \lambda) - a^2$, which gives $\lambda_1 \approx 2$ and $\lambda_2 \approx 2a^2$ for $|a| >> 1$.

## 4.7 Iterative methods for linear systems

This section is concerned with solving the equation $Ax = b$ with $A$ a matrix and $x$ and $b$ vectors, for very large sizes of $A$.

**Definition 4.7.1.** We define **Richardson's method** as the following: for some $\omega > 0$, let

$$x^{(k+1)} = x^{(k)} + \omega(b - Ax^{(k)})$$

Clearly $x^{(k)} = x^{(k+1)} = x$, the exact solution, satisfies this. The **residual** is defined as $r^{(k)} := x^{(k)} - x$. Then

$$\begin{aligned} r^{(k+1)} &= x^{(k+1)} - x \\ &= x^{(k)} - x + \omega(b - Ax^{(k)}) \\ &= r^{(k)} + \omega(b - Ax^{(k)} - Ax) \\ &= (I - \omega A)r^{(k)} \end{aligned}$$

So $r^{(k)} = (I - \omega A)^k r^{(0)}$. Consider the case where $A$ has real eigenvalues which satisfy

$$0 < \lambda_1 \leq \lambda_2 \leq \cdot \leq \lambda_m$$

So the eigenvalues of $I - \omega A$ are also real and satisfy

$$1 - \omega\lambda_m \leq \cdots \leq 1 - \omega\lambda_1 < 1$$

By the spectral radius theorem (TODO: add this to notes), for convergence, we need

$$-1 < 1 - \omega\lambda_m \iff \omega < 2/\lambda_m = 2/\rho(A)$$

So the iteration converges iff the above equation holds.

**Remark.** If $\rho(A)$ is large, then $\omega$ must be small. This is a problem, because if $\omega$ is small, then the iteration takes a long time to converge.

**Definition 4.7.2.** $A$ is called **strictly upper triangular** if $A$ is upper triangular and $A_{i,i}00$ for every $i$. $A$ is called **unit upper triangular triangular** if $A$ is upper triangular and $A_{i,i} = 1$ for every $i$.

**Definition 4.7.3.** We define **Jacobi's method** as follows: assuming $A_{i,i} \neq 0$ for every $i$, we rewrite $Ax = b$ as

$$Dx = (E + F)x + b$$

where $D$ is diagonal, $E$ is strictly lower triangular, $F$ is strictly upper triangular, and we use the iteration

$$Dx^{(k+1)} = (E + F)x^{(k)} + b \iff x^{(k+1)} = D^{-1}((E + F)x^{(k)} + b) =: M_J^{-1}(N_J x^{(k)} + b)$$

Consider the residual $r^{(k)} := x^{(k)} - x$:

$$\begin{aligned}
r^{(k+1)} &= x^{(k+1)} - x \\
&= M^{-1}Nx^{(k)} + M^{-1}b - x \\
&= M^{-1}Nr^{(k)} + M^{-1}Nx + M^{-1}b - x \\
&= M^{-1}Nr^{(k)}
\end{aligned}$$

Arguing as before, this iteration converges iff

$$(M_J^{-1}N_J)^k \to 0 \quad \text{as } k \to \infty$$

The implementation of this method is as follows:

$$x_i^{(k+1)} = \left( \sum_{j \neq i} A_{i,j}x_j^{(k)} + b_i \right) / A_{i,i}$$

**Remark.** Computing eigenvalues is even harder than solving $Ax = b$.

**Definition 4.7.4.** A square matrix $B$ is called **strictly diagonally dominant** if

$$\forall i, \quad |B_{i,i}| < \sum_{j \neq i} |B_{i,j}|$$

$B$ is called **diagonally dominant** if

$$\forall i, \quad |B_{i,i}| \leq \sum_{j \neq i} |B_{i,j}|$$

**Proposition 4.7.5.** The Jacobi method converges if $A$ is strictly diagonal dominant.

*Proof.*

$$\begin{aligned}
(M^{-1}N)_{i,j} &= \sum_k D_{i,k}^{-1}(A - D)_{i,j} \\
&= D_{i,i}^{-1}(A - D)_{i,j} \\
&= \begin{cases} A_{i,i}^{-1}A_{i,j} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

Then consider

$$||M^{-1}N||_\infty = \max_i |M^{-1}N_{i,j}| = \max_i \sum_{j \neq i} |A_{i,j}|/|A_{i,i}| < 1$$

if $A$ is strictly diagonal dominant, henc $(M^{-1}N)^k \to 0$ as $k \to \infty$. $\qquad\square$

**Definition 4.7.6.** The **Gauss-Seidel method** is defined as the iteration

$$x_i^{(k+1)} = \left( b_i - \sum_{j<i} A_{i,j} x_j^{(k+1)} - \sum_{j>i} A_{i,j} x_j^{(k)} \right)$$

$$= x_i^{(k)} + \left( b_i - \sum_{j<i} A_{i,j} x_j^{(k+1)} - \sum_{j\geq i} A_{i,j} x_j^{(k)} \right)$$

which only requires one vector for $x$: when computing $x_i^{(k+1)}$, we use the updated values of $x_j$ for $j < i$ and the old values of $x_j$ for $j \geq i$.

Splitting $A$ into $A = D - E - F$ where $D$ is diagonal, $E$ is strictly lower triangular, and $F$ is strictly upper triangular, we can write

$$(D - E)x^{(k+1)} = Fx^{(k)} + b \iff Mx^{(k+1)} = Nx^{(k)} + b$$

where $M = D - E$ and $N = F$. Consider $r^{(k)} := x^{(k)} - x$, then

$$r^{(k+1)} = M^{-1}(Nx^{(k)} + b) - x = M^{-1}(Nr^{(k)} + Nx + b) - x = M^{-1}Nr^{(k)}$$

**Proposition 4.7.7.** If $A$ is strictly diagonally dominant, then the Gauss-Seidel method converges.

**Remark.** To accelerate convergence, we can sometimes use a parameter $w > 1$ and use

$$x_i^{(k)} + w \left( b_i - \sum_{j<i} A_{i,j} x_j^{(k+1)} - \sum_{j\geq i} A_{i,j} x_j^{(k)} \right)$$

**Example 4.7.8.** Let $W$ be any non-singular matrix and $|| \cdot ||$ be any vector norm. Prove that

$$||x||_W := ||Wx||$$

is a vector norm.

Verify the properties:

- $||cx||_W = ||W(cx)|| = ||c(Wx)|| = |c|||Wx|| = |c|||x||_W$.

- $||x + y||_W = ||W(x + y)|| = ||Wx + Wy|| \leq ||Wx|| + ||Wy|| = ||x||_W + ||y||_W$.

- If $x \neq 0$, then $Wx \neq 0$ as $W$ is non-singular. Hence $||x||_W = 0 \implies x = 0$.

**Example 4.7.9.** Prove that for symmetric $A$, $\rho(A) = ||A||_2$.

For any matrix $B$, $||B||_2 = \rho(BB^T)^1/2$. For symmetric $A = A^T$, we have a full set of eigenvectors with eigenvalues $\lambda_1 \leq \cdots \leq \lambda_n$. Now if $Av = \lambda v$ then $AA^Tv = A^2v = A(\lambda v) = \lambda^2 v$. So if $\lambda^2$ is an eigenvalue of $AA^T$ then $|\lambda|$ is an eigenvalue of $A$ with a largest modulus.

# 5   $L^2$ Approximations

## 5.1   $L^2$ approximations of functions

**Definition 5.1.1.** Let $f$ be a function and $w$ be a weight function, with $w(x) > 0$ except at a finite number of points. The **weighted $L^2$ norm** of $f$ is defined as

$$||f||_{L^2_w(a,b)} := \left( \int_a^b |f(x)|^2 w(x) dx \right)$$

**Remark.** The weighted $L^2$ norm is a seminorm for every function $f$ but is a norm for every continuous $f$.

**Example 5.1.2.** To solve the problem, we can write $p(x) = p_0 \phi(x) + \cdots + p_N \phi_N(x)$ for some undetermined coefficients $p_0, \ldots, p_N$. Then

$$E(p) = \int_a^b |p_0 \phi_0(x) + \cdots + p_N \phi_N(x) - f(x)|^2 w(x) dx$$

$$= \int_a^b \left| \sum_{i=0}^N p_i \phi_i(x) - f(x) \right|^2 w(x) dx$$

$$= \int_a^b \left( \sum_{j,k} p_j p_k \phi_j(x) \phi_k(x) - 2 f(x) \sum_k p_k \phi_k(x) + f(x)^2 \right) w(x) dx$$

$$= ||f||^2 - 2 \sum_k p_k \int_a^b \phi_k(x) f(x) w(x) dx + \sum_{j,k} p_j p_k \int_a^b \phi_j(x) \phi_k(x) w(x) dx$$

$$=: ||f||^2 - 2 \sum_k p_k q_k + \sum_{j,k} p_j p_k A_{j,k}$$

To find the $p$ which minimises this, we vary the coefficients $p_j$ until we find a stationary point. For $i \in \{0, \ldots, N\}$,

$$\frac{\partial E}{\partial p_i} = -2 q_i + 2 \sum_k A_{i,k} p_k + \sum_j p_j A_{j,i} \iff \sum_j A_{i,j} p_j = q_i$$

If we can solve this linear system, we have our stationary point. We show that $A$ is strictly positive definite, so there is a unique $p$ which minimises the error. For every vector $(z_0, \ldots, z_N)$,

$$\sum_{j,k} A_{j,k} z_j z_k = \int_a^b \left( \sum_j z_j \phi_j(x) \right) \left( \sum_k z_k \phi_k(x) \right) w(x) dx$$

$$= \int_a^b \left( \sum_j z_j \phi_j(x) \right)^2 w(x) dx$$

$$> 0 \text{ iff the } \phi_j \text{ are linearly independent}$$

so assuming that $\phi_j$ are linearly independent, $A$ is strictly positive definite. So $E(p)$ has a unique minimum and no maximum.

**Example 5.1.3.** Let $(a, b) = (-1, 1)$ and $w(x) = (1 - x^2)^{-1/2}$. Find a quadratic approximation to $f(x) = \arcsin(x)$. We use the above example to compute

$$q_0 = \int_{-1}^1 x^0 \arcsin(x)(1 - x^2)^{-1/2} dx = 0$$

since the integrand is an odd function. Similarly,

$$q_2 = \int_{-1}^{1} x^2 \arcsin(x)(1 - x^2)^{-1/2} dx = 0$$

Now

$$q_1 = \int_{-1}^{1} x^1 \arcsin(x)(1 - x^2)^{-1/2} dx = \int_{0}^{\pi} \cos(\theta)(\pi/2 - \theta) d\theta = 2$$

So

$$A_{1,1} = \int_{-1}^{1} x^0 x^0 (1 - x^2)^{-1/2} dx = \pi$$

$$A_{1,2} = A_{2,1} = \int_{-1}^{1} x^0 x^1 (1 - x^2)^{-1/2} dx = 0$$

$$A_{2,3} = A_{3,2} = \int_{-1}^{1} x^1 x^2 (1 - x^2)^{-1/2} dx = 0$$

$$A_{3,1} = A_{2,2} = A_{1,3} = \int_{-1}^{1} x^2 x^0 (1 - x^2)^{-1/2} dx = \frac{\pi}{2}$$

$$A_{3,3} = \int_{-1}^{1} x^2 x^2 (1 - x^2)^{-1/2} dx = \frac{3\pi}{8}$$

So

$$A = \pi \begin{pmatrix} 1 & 0 & 1/2 \\ 0 & 1/2 & 0 \\ 1/2 & 0 & 3/8 \end{pmatrix}, \quad q = (0, 2, 0)$$

So $p = \frac{1}{\pi}(0, 4, 0)$ so $p(x) = \frac{4}{\pi}x$.

**Definition 5.1.4.** A map $(\cdot, \cdot) : V \times V \to \mathbb{C}$ is called an **inner product** if $\forall u, v \in V$ and $\alpha, \beta \in \mathbb{C}$,

1. $(\alpha u + \beta u', v) = \alpha(u, v) + \beta(u', v)$.

2. $(u, v) = \overline{(v, u)}$.

3. $(u, u) \geq 0$ and $(u, u) = 0$ iff $u = 0$.

**Remark.** From properties 1. and 2. of an inner product, $(u, \alpha u) = \overline{\alpha}(u, v)$.

**Theorem 5.1.5.** Given an inner product $(\cdot, \cdot)$ on $V$, the map $|| \cdot || : V \to \mathbb{R}_+$ defined as

$$||u|| := (u, u)^{1/2}$$

is a norm.

*Proof.*

- Linearity: use properties 1. and 2.

- Positivity for $u \neq 0$ follows from property 3.

- The triangle inequality follows from Cauchy-Schwarz:

$$0 \leq ||u + \alpha v||^2 = (u + \alpha v, u + \alpha v)$$
$$= ||u||^2 + |\alpha|^2 ||v||^2 + \alpha(v, u) + \overline{\alpha}(u, v)$$

28

Then set $\alpha = -(u,v)/||v||^2$ so $\overline{\alpha} = -(v,u)/||v||^2$ by property 2. Then

$$0 \le ||u||^2 + |(u,v)|^2/||v||^2 - (u,v)(v,u)/||v||^2 - (v,u)(u,v)/||v||^2$$
$$= ||u||^2 - |(u,v)|^2/||v||^2$$
$$\iff |(u,v)|^2 \le ||u||^2||v||^2$$

This proves the Cauchy-Schwarz inequality. Now

$$||u+v||^2 = (u+v, u+v) = ||u||^2 + ||v||^2 + (u,v) + (v,u)$$
$$= ||u||^2 + ||v||^2 + 2|(u,v)|$$
$$\le ||u||^2 + ||v||^2 + 2||u||||v||$$
$$= (||u|| + ||v||)^2$$

so

$\square$

**Example 5.1.6.** Let $V = \mathbb{R}^n$ and let $u = (u_1, \ldots, u_n) = u_1\underline{e_1} + \cdots + u_n\underline{e_n}$. An inner product is defined as

$$(u,v) = \sum_{i=1}^{n} u_i v_i$$

Then a norm is

$$||u|| = \sqrt{\sum_{i=1}^{n} u_i^2}$$

**Example 5.1.7.** Let $V = C^0([a,b])$. An inner product is defined as

$$(u,v)_{L^2w(a,b)} := \int_a^b u(x)v(x)w(x)dx$$

Then a norm is

$$(u,u)_{L^2w(a,b)} = ||u||^2_{L^2w(a,b)}$$

i.e. the $L^2$ norm as before.

**Theorem 5.1.8.** Let $V$ be an inner product space and $X$ be a linear subspace of $V$. Let $f \in V$ be given and suppose $\tilde{p} \in X$ minimises $||f - p||$:

$$E(\tilde{p}) = ||f - \tilde{p}||^2 \le ||f - p||^2 \quad \forall p \in X$$

Then $\tilde{p}$ satisfies

$$(f - \tilde{p}, p) = 0 \quad \forall p \in X$$

*Proof.* (Assume $V$ is a real inner product space). Assume that $(f - \tilde{p}, \hat{p}) = a \ne 0$ for some $\hat{p} \in X$. Let

$$q = \hat{p} + \frac{a}{||\hat{p}||^2}\hat{p} \in X$$

Then

$$||f - q||^2 = ||f - \hat{p} - \frac{a}{||\hat{p}||^2}\hat{p}||^2$$
$$= ||f - \hat{p}||^2 - \frac{2a}{||\hat{p}||^2}(f - \tilde{p}, \hat{p}) + \frac{a^2}{||\hat{p}||^4}||\hat{p}||^2$$
$$= ||f - \tilde{p}||^2 - \frac{a^2}{||\hat{p}||^2}$$
$$< ||f - \tilde{p}||^2$$

But this is a contradiction as then $||f - p||$ is not minimised by $\tilde{p}$. So $a = 0$. $\square$

**Definition 5.1.9.** (**Gram-Schmidt**) Given a basis $\{\phi_k\}$ which is not orthogonal, the algorithm is as follows:

- $\hat{\phi}_0 = \phi_0$.

-
$$\hat{\phi}_k = \phi_k - \sum_{j=0}^{k-1} \frac{(\hat{\phi}_j, \phi_k)}{||\hat{\phi}_j||^2} \hat{\phi}_j$$

**Proposition 5.1.10.** Let $V = C^0(a, b; \mathbb{R})$ with the inner product

$$(u, v)_{L^2 w(a,b)} := \int_a^b u(x)v(x)w(x)dx$$

and $X = P_n$ (set of polynomials of degree $n$). let $P_n = \text{span}\{\phi_0, \dots, \phi_n\}$ with $\phi_k \in P_k$. The following properties hold:

1. $\{\phi_k\}$ is unique up to normalisation: if $\{\hat{\phi}_k\}$ is another orthogonal set with respect to the same inner product, then $\hat{\phi}_k = c_k \phi_k$ for some $c_k \in \mathbb{R}$.

2. $\phi_k$ has exactly $k$ real roots in $(a, b)$.

3. $\{\phi_k\}$ satisfy the recurrence relation $\phi_{-1}(x) := 0, \phi_0(x) = 1$,

$$\phi_{k+1}(x) = \frac{1}{||\phi_k||} x\phi_k(x) - \frac{(x\phi_k, \phi_k)}{||\phi_k||^3} \phi_k(x) - \frac{||\phi_k||}{||\phi_{k-1}||} \phi_{k-1}(x)$$

*Proof.*

1. Let $\phi_k(x) = a_k x^k + \cdots + a_0$, $\tilde{\phi}_k(x) = \tilde{a}_k x^k + \cdots + \tilde{a}_0$. Then $\tilde{\phi}_k(x) - \frac{\tilde{a}_k}{a_k} \phi_k(x) =: b_{k-1} x^{k-1} + \cdots + b_0 =: \delta\phi(x) \in P_{n-1}$. Now,

$$(\delta\phi, \delta\phi) = \left( \tilde{\phi}_k - \frac{\tilde{a}_k}{a_k} \phi_k, \delta\phi \right) = (\tilde{\phi}_k, \delta\phi) - \frac{\tilde{a}_k}{a_k}(\phi_k, \delta\phi) = 0$$

   since $(\tilde{\phi}_k, \phi_m) = 0$ for $k \neq m$ since they are orthogonal. Hence $\delta\phi = 0$.

2. We have $\forall p \in P_{k-1}, (\phi_k, p) = 0$. $\phi_k$ can have at most $k$ sign changes in $(a, b)$. We will show that it has exactly $k$ sign changes. Assume that $\phi_k$ has $m < k$ sign changes at $\{x_j\}_{j=1}^m$. Define
$$s(x) := (x - x_1) \cdots (x - x_m) \in P_m$$

   Then
$$(\phi_k, s) = \int_a^b \phi_k(x)s(x)w(x)dx = 0$$

   since $s \in P_m$ with $m < k$. So $m = k$ and each sign change corresponds to a simple root.

3. Omitted.

$\square$

**Remark.** Properties 1 and 2 imply that the roots of $\phi_k$ are unique. The recurrence relation defined above should be used instead of Gram-Schmidt where possible.

## 5.2   Convergence of $L^2$ approximations

**Theorem 5.2.1. (Weierstrass)** For every $f \in C^0(a, b; \mathbb{R})$ and every $\epsilon > 0$, for some $N \in \mathbb{N}$ and $p \in P_N$,

$$\max_{x \in [a,b]} |f(x) - p(x)| < \epsilon$$

*Proof.* Omitted. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

# 6 Numerical Integration

We want to numerically approximate the integral

$$I[f] = \int_a^b f(x)w(x)dx$$

for $w(x)$ a weight function.

**Definition 6.0.1.** For a function $f$, a **quadrature formlula** is defined as

$$Q_n(f) = (b-a) \sum_{k=0}^n \hat{\sigma}_k f(x_k)$$

where $\hat{\sigma}_k$ are are the **coefficients** and $x_k \in (a,b)$ are the **nodes**.

**Example 6.0.2.** The trapezium rule is

$$Q_1(f) = (b-a)\frac{f(a) + f(b)}{2}$$

Here $x_0 = a$, $x_1 = b$, $\hat{\sigma}_0 = \hat{\sigma}_1 = \frac{1}{2}$.

**Definition 6.0.3.** A quadrature formula $Q_n$ has **degree of exactness** $r$ if

- $Q_n(x^m) = I(x^m)$ for $m \in \{0, \ldots, r\}$ and

- $Q_n(x^{r+1}) \neq I(x^{r+1})$.

**Proposition 6.0.4.**

1. $Q_n$ is linear: $Q_n(f + g) = Q_n(f) + Q_n(g)$.

2. If $Q_n$ has degree of exactness $r$, then $Q_n(p) = I(p)$ for every $p \in P_r$.

**Definition 6.0.5. (Interpolatory quadrature)** Let $w(x) = 1$ and let $\{x_k\}_{k=0}^n$ be such that $x_k < x_{k+1}$. We interpolate $f$ by $p \in P_n$ at the nodes: $p(x_k) = f(x_k)$. Then we compute

$$I_n(f) = \int_a^b p(x)dx$$

Now using the Lagrange interpolation formula,

$$p(x) = \sum_{k=0}^n f(x_k)L_k(x)$$

where

$$L_k(x) = \prod_{j=0, j\neq k}^n \frac{x - x_j}{x_k - x_j}$$

Therefore

$$I_n(f) = \int_a^b \sum_k f(x_k)L_k(x)dx = \sum_k f(x_k) \int_a^b L_k(x)dx$$

Let $t = \frac{x-a}{b-a} \in (0,1)$, $t_k = \frac{x_k-a}{b-a} \in (0,1)$. Then

$$\int_a^b L_k(x) = \int_a^b \prod_{j\neq k} \frac{x - x_j}{x_k - x_j}dx = (b-a)\int_0^1 \prod_{j\neq k} \frac{t - t_j}{t_k - t_j}dt =: \sigma_k$$

**Definition 6.0.6.** If the nodes $x_k$ are equidistant, $x_k = x_{k-1} - x_{k+1} - x_k$, the resulting quadrature is called a **Newton-Cotes formula**. If $x_0 = a$ and $x_n = b$, it is called a **closed Newton-Cotes formula**.

**Example 6.0.7.** For $n = 0$, $L_0(x) = 1$ and $\sigma_0 = 1$. So $I_0(f) = (b-a)F(\xi)$ for any $\xi \in [a, b]$. If $\xi = \frac{b+a}{2}$, this is called the **rectangle method**.

**Example 6.0.8.** For $n = 1$, for the closed Newton-Cotes formula, $t_0 = 0$ and $t_1 = 1$. $L_0(t) = \frac{t-t_1}{t_0-t_1} = 1 - t$ and $L_1(t) = \frac{t-t_0}{t_1-t_0} = t$, so $\sigma_0 = \sigma_1 = \frac{1}{2}$. Hence

$$I_1(f) = (b-a)\frac{f(a) + f(b)}{2}$$

**Example 6.0.9.** For $n = 2$, for the closed Newton-Cotes formula, $t_0 = 0$, $t_1 = \frac{1}{2}$ and $t_2 = 1$. $L_0(t) = \frac{(t-t_1)(t-t_2)}{(t_0-t_1)(t_0-t_2)} = \frac{1}{2}t^2 - t$ and $L_1(t) = \frac{(t-t_0)(t-t_2)}{(t_1-t_0)(t_1-t_2)} = -t^2 + t$ and $L_2(t) = \frac{(t-t_0)(t-t_1)}{(t_2-t_0)(t_2-t_1)} = \frac{1}{2}t^2 - \frac{1}{2}t$. So $\sigma_0 = \sigma_2 = \frac{1}{6}$ and $\sigma_1 = \frac{2}{3}$. Hence

$$I_2(f) = (b-a)\frac{f(a) + 4f(\frac{a+b}{2}) + f(b)}{6}$$

This is called **Simpson's rule**.

**Theorem 6.0.10.** Let $f \in C^{n+1}([a, b])$ and let $p \in P_n$ interpolate $f$ at $\{x_k\}_{k=0}^n \subset [a, b]$. Then for every $x \in [a, b]$,

$$f(x) - p(x) = \frac{(x - x_0) \cdots (x - x_n)}{(n+1)!} f^{(n+1)}(\xi)$$

where $\xi \in (a, b)$ and we denote $\omega_{n+1}(x) = (x - x_0) \cdots (x - x_n)$.

**Theorem 6.0.11.** Let $f \in C^{n+1}([a, b])$. If $I_n$ is an interpolatory quadrature on $[a, b]$, then

$$|I(f) - I_n(f)| \le \frac{1}{(n+1)!} \max_{\xi \in (a,b)} |f^{(n+1)}(\xi)| \int_a^b |(x - x_0) \cdots (x - x_n)| dx$$

**Example 6.0.12.** (Problems class) Suppose that the interpolatory quadrature formula $I_1$ has the nodes $1/4$ and $2/3$ on the integration domain $[0, 1]$. Compute the coefficients of $I_1$.
$I_1(f) = \sum_{k=0}^n \sigma_k f(x_k) = \sigma_0 f(1/4) + \sigma_1 f(2/3)$. We compute

$$\sigma_0 = (b-a) \int_a^b \prod_{j \ne 0} \frac{x - x_j}{x_0 - x_j} dx = \int_0^1 \frac{x - x_1}{x_0 - x_1} dx$$

$$= \int_0^1 \frac{x - 2/3}{1/4 - 2/3} dx = \frac{2}{5}$$

and

$$\sigma_1 = \int_0^1 \frac{x - x_0}{x_1 - x_0} dx = \frac{3}{5}$$

**Example 6.0.13.** (Problems class) Find the coefficient $\sigma_0$ in the open Newton-Cotes formula

$$I_0(f) = \sigma_0 f\left(\frac{a+b}{2}\right)$$

and prove that

$$|I(f) - I_0(f)| \leq \frac{(b-a)^3}{24} \max_{\xi \in (a,b)} |f^{(2)}(\xi)|$$

We use Taylor's theorem to write

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

where $\xi(x) \in \text{conv}\{x, x_0\} \subset [0, 1]$. Then, with $p_0 = f(x_0)$ a constant polynomial,

$$\begin{aligned}
I(f) - I_0(f) &= I(f) - I(p_0) = I(f - p_0) \\
&= \int_a^b (f(x) - p_0(x))dx \\
&= \int_a^b ((x - x_0)f'(x_0) + \frac{(x - x_0)^2}{2!} f''(\xi(x)))dx \\
&= f'(x_0) \int_a^b (x - x_0)dx + \frac{1}{2!} \int_a^b f''(\xi(x))(x - x_0)^2 dx
\end{aligned}$$

So

$$\begin{aligned}
|I(f) - I_0(f)| &= \frac{1}{2} \left| \int_a^b f''(\xi(x))(x - x_0)^2 dx \right| \\
&\leq \frac{1}{2} \int_a^b |f''(\xi(x))|(x - x_0)^2 dx \\
&\leq \frac{1}{2} \max_{\xi \in [0,1]} |f''(\xi)| \int_a^b (x - x_0)^2 dx \\
&= \frac{(b-a)^3}{24} \max_{\xi \in (a,b)} |f^{(2)}(\xi)|
\end{aligned}$$

**Example 6.0.14.** (Problems class) Show that the Fourier-Chebyshev expansion

$$\tilde{p}(x) = \sum_{k=0}^{\infty} \frac{(f, T_k)}{||T_k||^2} T_k(x)$$

where $T_K$ is the Chebyshev polynomial and the inner product is in $(-1, 1)$ with weight $w(x) = (1 - x^2)^{-1/2}$, converges for every $f \in C^2([-1, 1])$.

We will show that the series

$$p_N(x) = \sum_{k=0}^{N} \frac{(f, T_k)}{||T_k||^2} T_k(x)$$

converges absolutely.

$$||p_N||^2 = (p_N, p_N) = \sum_{j,k} c_j c_k (T_k, T_j) = c_0^2 \pi + \sum_{k=1}^{N} c_k^2 \pi/2$$

where $c_j = (f, T_j)/||T_j||^2$. Now

$$
\begin{aligned}
(f, T_k) &= \int_{-1}^{1} f(x) T_k(x)(1 - x^2)^{-1/2} dx \\
&= \int_{0}^{\pi} f(\cos(\theta)) \cos(k\theta) d\theta \\
&= \frac{1}{k} \left[ f(\cos(\theta)) \sin(k\theta) \right]_0^\pi - \frac{1}{k} \int_0^\pi \sin(k\theta) \frac{d}{d\theta} f(\cos(\theta)) d\theta \\
&= \frac{1}{k^2} \left[ \cos(k\theta) \frac{d}{d\theta} f(\cos(\theta)) \right] + \frac{1}{k^2} \int_0^\pi \cos(k\theta) \frac{d^2}{d\theta^2} f(\cos(\theta)) d\theta \\
&= -\frac{1}{k^2} \left[ \cos(k\theta) f'(\cos(\theta)) \sin(\theta) \right]_0^\pi - \frac{1}{k^2} \int_0^\pi \cos(k\theta) \frac{d}{d\theta} (\sin(\theta) f'(\cos(\theta))) d\theta \\
&= -\frac{1}{k^2} \int_0^\pi \cos(k\theta)(\cos(\theta) f'(\cos(\theta) - \sin(\theta)^2 f''(\cos(\theta))) d\theta
\end{aligned}
$$