# Numerical Analysis Lecture Notes

Rebekah Dix

October 17, 2018

# Contents

# 1 Results from Real Analysis

**Theorem 1.** (The Mean Value Theorem) Suppose $f$ is a real-valued function, defined and continuous on the closed interval $[a, b] \in \mathbb{R}$ and $f$ differentiable on the open interval $(a, b)$. Then there exists a number $\xi \in (a, b)$ such that

$$f(b) - f(a) = f'(\xi)(b - a) \tag{1}$$

**Theorem 2.** (Taylor's Theorem) Suppose that $n$ is a nonnegative integer, and $f$ is a real-valued function, defined and continuous on the closed interval $[a, b]$ of $\mathbb{R}$, such that the derivatives of $f$ of order up to and including $n$ are defined and continuous on the closed interval $[a, b]$. Suppose further that $f^{(n)}$ is differentiable on the open interval $(a, b)$. Then, for each value of $x \in [a, b]$, there exists a number $\xi = \xi(x)$ in the open interval $(a, b)$ such that

$$f(x) = f(a) + (x - a)f'(a) + \cdots + \frac{(x - a)^n}{n!}f^{(n)}(a) + \frac{(x - a)^{n+1}}{(n + 1)!}f^{(n+1)}(\xi) \tag{2}$$

# 2 Solution of equations by iteration

## 2.1 Simple Iteration

**Theorem 3.** (Existence of Root) Let $f$ be a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line. Assume further, that $f(a)f(b) \leq 0$; then, there exists $\xi$ in $[a, b]$ such that $f(\xi) = 0$.

*Proof.* The condition $f(a)f(b) \leq 0$ implies that $f(a)$ and $f(b)$ have opposite signs, or one of them is 0. If either $f(a)$ or $f(b)$ is 0, then we've found a root. Suppose that both endpoints are non-zero (in which case they have opposite signs). In this case, 0 must belong to the open interval whose endpoints are $f(a)$ and $f(b)$. The intermediate value theorem gives the existence of a root in the open interval $(a, b)$. Thus, in both cases, a zero is guaranteed. $\square$

- The converse of Theorem 3 is clearly false.

**Theorem 4.** (Brouwer's Fixed Point Theorem) Suppose that $g$ is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and let $g(x) \in [a, b]$ for all $x \in [a, b]$. Then, there exists $\xi \in [a, b]$ such that $\xi = g(\xi)$. $\xi$ is called a fixed point of the function $g$.

*Proof.* Define a function $f(x) = x - g(x)$. If we find a root $\xi$ of $f$, then $\xi$ is a fixed point of $g$. Then,

$$f(a)f(b) = (a - g(a))(b - g(b)) \leq 0 \tag{3}$$

By assumption, $a \leq g(a), g(b) \leq b$. Therefore, the first term is negative and the second term is positive. Therefore, $f(a)f(b) \leq 0$. By Theorem 3, there exists a $\xi \in [a, b]$ such that $f(\xi) = 0$. Then, for this $\xi$, $g(\xi) = \xi$. $\square$

**Definition 1.** (Simple Iteration) Suppose that $g$ is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and let $g(x) \in [a, b]$ for all $x \in [a, b]$. Given that $x_0 \in [a, b]$, the recursion defined by

$$x_{k+1} = g(x_k) \tag{4}$$

is called simple iteration; the numbers $x_k$, $k \geq 0$, are referred to as iterates.

- If this sequence converges, the limit must be a fixed of $g$, since $g$ is continuous on a closed interval. Note that

$$\xi = \lim_{k \to \infty} x_{k+1} = \lim_{k \to \infty} g(x_k) = g\left(\lim_{k \to \infty} x_k\right) = g(\xi) \tag{5}$$

**Definition 2.** (Contraction) Let $g$ be a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line. Then, $g$ is said to be a contraction on $[a, b]$ if there exists a constant $L$ such that $0 < L < 1$ and

$$|g(x) - g(y)| \leq L|x - y| \quad \forall x, y \in [a, b] \tag{6}$$

**Theorem 5.** (Contraction Mapping Theorem) Suppose that $g$ is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and let $g(x) \in [a, b]$ for all $x \in [a, b]$. Suppose $g$ is a contraction on $[a, b]$. Then, $g$ has a unique fixed point $\xi$ in the interval $[a, b]$. Moreover, the sequence $(x_k)$ defined by simple iteration converges to $\xi$ as $k \to \infty$ for any starting value $x_0$ in $[a, b]$.

Let $\epsilon > 0$ be a certain tolerance, and let $k_0(\epsilon)$ denote the smallest positive integer such that $x_k$ is no more than $\epsilon$ away from the fixed point $\xi$ (i.e. $|x_k - \xi| \leq \epsilon$) for all $k \geq k_0(\epsilon)$. Then,

$$k_0(\epsilon) \leq \left\lfloor \frac{\ln|x_1 - x_0| - \ln(\epsilon(1 - L))}{\ln(1/L)} \right\rfloor + 1 \tag{7}$$

*Proof.* Let $E_k = |x_k - \xi|$ be the error at $k$. Then

$$|x_{k+1} - \xi| = |g(x_k) - g(\xi)|$$
$$< L|x_k - \xi|$$

Therefore

$$E_k \leq L^k E_0 \tag{8}$$

Since $L < 1$, $L^k \to 0$ as $k \to \infty$. $\qquad \square$

**Definition 3.** (Stable, Unstable Fixed Point) Suppose that $g$ is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and let $g(x) \in [a, b]$ for all $x \in [a, b]$, and let $\xi$ denote a fixed point of $g$. $\xi$ is a stable fixed point of $g$ if the sequence $(x_k)$ defined by the iteration $x_{k+1} = g(x_k)$, $k \geq 0$, converges to $\xi$ whenever the starting value $x_0$ is sufficiently close to $\xi$. Conversely, if no sequence $(x_k)$ defined by this

iteration converges to $\xi$ for any starting value $x_0$ close to $\xi$, except for $x_0 = \xi$, then we say that $\xi$ is an unstable fixed point of $g$.

- With this definition, a fixed point may be neither stable nor unstable.

**Definition 4.** (Rate of Convergence) Suppose $\xi = \lim_{k \to \infty} x_k$. Define $E_k = |x_k - \xi|$.

- The sequence $(x_k)$ converges to $\xi$ linearly if there exists a number $\mu \in (0, 1)$ such that

$$\lim_{k \to \infty} \frac{E_{k+1}}{E_k} = \mu \tag{9}$$

- The sequence $(x_k)$ converges to $\xi$ superlineraly if $\mu = 0$. That is, the sequence of $\mu_k$ generated at each step $\to 0$ as $k \to \infty$.

- The sequence $(x_k)$ converges to $\xi$ with order $q$ if there exists a $\mu > 0$ such that

$$\lim_{k \to \infty} \frac{E_{k+1}}{E_k^q} = \mu \tag{10}$$

In particular, if $q = 2$, then the sequence converges quadratically.

## 2.2 Newton's Method

**Definition 5.** (Newton's Method) Newton's method for the solution of $f(x) = 0$ is defined by

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \tag{11}$$

Geometrically, $(x_{n+1}, 0)$ is the intersection of the $x$-axis and the tangent of the graph of $f$ at $(x_n, f(x_n))$.
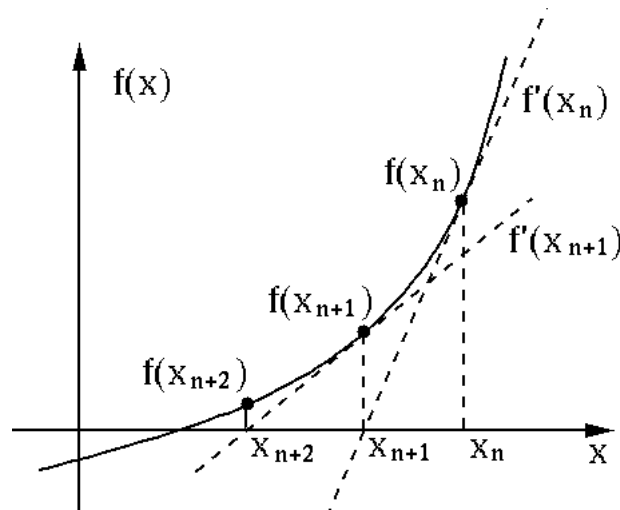


Figure 1: Geometric Interpretation of Newton's Method in $\mathbb{R}$

5

Intuitively, the fixed points of this iteration $g$ will be stable.
We can show that $|g'(\xi)| < 1$.

$$g'(x) = 1 - \frac{f' \cdot f' - f \cdot f''}{(f')^2}$$

$$= 1 - \left(1 - \frac{f(x) \cdot f''(x)}{(f'(x))^2}\right)$$

$$= \frac{f(x) \cdot f''(x)}{(f'(x))^2}$$

Therefore

$$|g'(\xi)| = \left|\frac{f(\xi) \cdot f''(\xi)}{(f'(\xi))^2}\right| = 0 < 1 \tag{12}$$

**Theorem 6.** (Convergence of Newton's Method) Suppose that $f$ is a continuous real-valued function with continuous second derivative $f''$ defined on the closed interval $I_\delta = [\xi - \delta, \xi + \delta]$, $\delta > 0$, such that $f(\xi) = 0$ and $f''(\xi) \neq 0$. Additionally suppose that there exists a positive constant $A$ such that

$$\frac{|f''(x)|}{|f'(y)|} \leq A \quad \forall x, y, \in I_\delta \tag{13}$$

If initially

$$|\xi - x_0| \leq h = \min(\delta, \frac{1}{A}) \tag{14}$$

then the sequence $(x_k)$ defined by Newton's method converges quadratically to $\xi$.

*Proof.* We first compute the Taylor expansion of $f(\xi)$, expanding about the point $x_k \in I_\delta$, where $|\xi - x_k| \leq h = \min(\delta, \frac{1}{A})$. Thus

$$f(\xi) = f(x_k) + (\xi - x_k)f'(x_k) + \frac{(\xi - x_k)^2}{2}f''(\eta_k) \tag{15}$$

where $\eta_k$ is between $\xi$ and $x_k$. Recall that $f(\xi) = 0$. We can use this fact and the definition of Newton's iteration to rearrange the above expansion as

$$\xi - x_{k+1} = -\frac{(\xi - x_k)^2 f''(\eta_k)}{2f'(x_k)} \tag{16}$$

A small modification to this equation allows us to derive a relationship between adjacent errors

$$E_{k+1} = \frac{f''(\eta_k)}{2f'(x_k)}E_k^2 \tag{17}$$

Recall by assumption we have that $|\xi - x_k| \leq h = \min(\delta, \frac{1}{A})$ and $\frac{|f''(x)|}{|f'(y)|} \leq A \quad \forall x, y, \in I_\delta$.

Therefore,

$$|E_{k+1}| = \frac{1}{2}\left|\frac{f''(\eta_k)}{f'(x_k)}\right||E_k|^2 \le \frac{1}{2}|E_k| \tag{18}$$

We are given that $|\xi - x_0| \le h = \min(\delta, \frac{1}{A})$, so that induction gives that

$$|E_k| = |\xi - x_k| \le \frac{1}{2^k}h \tag{19}$$

Therefore $(x_k)$ converges to $\xi$ as $k \to \infty$.

To show convergence is quadratic, notice that

$$\lim_{k\to\infty}\frac{|E_{k+1}|}{|E_k|} = \lim_{k\to\infty}\frac{1}{2}\frac{|f''(\eta_k)|}{|f'(x_k)|}$$
$$= \frac{1}{2}\frac{|f''(\xi)|}{|f'(\xi)|} = \mu \le \frac{A}{2}.$$

This shows that convergence is quadratic. $\qquad\qquad\qquad\qquad\qquad\qquad\square$

## 2.3 Secant Method

Observe that Newton's method requires us to know the first derivative $f'$ of $f$. In applications, we might not know $f'$ or it could be expensive to calculate. This motivates approximating the $f'(x_k)$ in Newton's method with

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}} \tag{20}$$

**Definition 6.** (Secant Method) The secant method is defined by

$$x_{k+1} = x_k - f(x_k)\frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})} \tag{21}$$
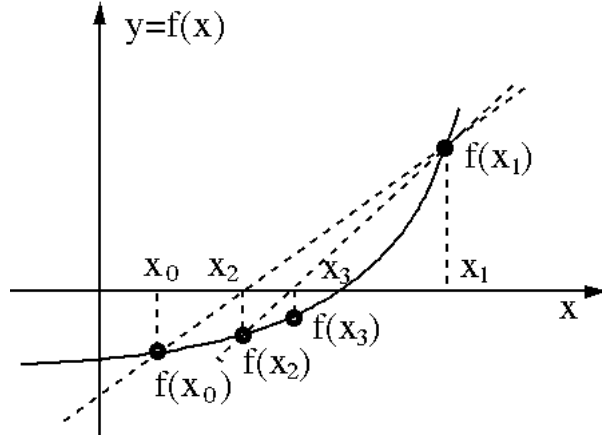
Figure 2: Geometric Interpretation of Secant Method in $\mathbb{R}$

**Theorem 7.** (Convergence of Secant Method) Suppose that $f$ is a real-valued function, defined and continuously differentiable on an interval $I = [\xi - h, \xi + h]$, $h > 0$, with center point $\xi$. Suppose further that $f(\xi) = 0$, $f'(\xi) \neq 0$. Then, the sequence $(x_k)$ defined by the secant method converges at least linearly to $\xi$ provided that $x_0$ and $x_1$ are sufficiently close to $\xi$.

*Proof.* Without loss of generality, assume that $\alpha = f'(\xi) > 0$ in a small neighborhood of $\xi$. We'll choose this neighborhood such that

$$0 < \frac{3}{4}\alpha < f'(x) < \frac{5}{4}\alpha \tag{22}$$

for all $x$ in the interval.

Recall that the secant method is defined by

$$x_{k+1} = x_k - f(x_k)\frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})} \tag{23}$$

We can repeatedly use the mean value theorem to approximate each of these terms. First observe that

$$\frac{f(x_k) - f(\xi)}{x_k - \xi} = f'(\eta_k) \tag{24}$$

for some $\eta_k$ between $x_k$ and $\xi$. Since $f(\xi) = 0$, this equation implies that

$$f'(x_k) = f'(\eta_k)(x_k - \xi) \tag{25}$$

Next observe that

$$\frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}} = f'(\theta_k) \tag{26}$$

for some $\theta_k$ between $x_k$ and $x_{k-1}$. Therefore, we can put these pieces together to observe

8

that,

$$x_{k+1} = x_k - \frac{f'(\eta_k)(x_k - \xi)}{f'(\theta_k)} \tag{27}$$

To show convergence, we can compare successive error terms.

$$
\begin{aligned}
E_{k+1} &= x_{k+1} - \xi \\
&= E_k - \frac{f'(\eta_k)}{f'(\theta_k)} E_k \\
&= \left(1 - \frac{f'(\eta_k)}{f'(\theta_k)}\right) E_k
\end{aligned}
$$

Therefore

$$
\begin{aligned}
\frac{E_{k+1}}{E_k} &= 1 - \frac{f'(\eta_k)}{f'(\theta_k)} \\
&< 1 - \frac{5\alpha/4}{3\alpha/4} \\
&= \frac{2}{3} \\
&< 1
\end{aligned}
$$

Therefore the secant method converges at least linearly. $\qquad\square$

# 3 Solution of systems of linear equations

## 3.1 LU Decomposition

## 3.2 Least Squares

Given a system of equations $Ax = b$, the least squares problem is

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 \tag{28}$$

We can expand the objective function out as

$$
\begin{aligned}
\|Ax - b\|_2^2 &= (Ax - b)^T (Ax - b) \\
&= x^T A^T A x - 2b^T A x + b^T b
\end{aligned}
$$

To find the $x$ that minimizes this expression we find the $x$ that satisfies $\nabla_x F = 0$. That is

$$\nabla_x F = 0 = 2A^T A x - 2A^T b \tag{29}$$

Therefore the minimizer is $x = (A^T A)^{-1} A^T b$. $(A^T A)^{-1} A^T$ is called the pseudo-inverse of $A$. If $A$ is square and invertible, then the pseudo-inverse equals $A^{-1}$.

## 3.3 Gram-Schimdt Orthogonalization

Algorithm: Denote the columns of $A$ by $a_i$.

1. $q_1 = a_1$. Then normalized by $q_1 = \frac{q_1}{\|q_1\|}$.

2. $q_2 = a_2 - \langle q_1, a_2 \rangle q_1$. Then normalize by $q_2 = \frac{q_2}{\|q_2\|}$. It's simple to verify that $q_2 \perp q_1$.

3. For an arbitrary $k$, $q_k = a_k - \langle a_k, q_1 \rangle q_1 - \langle a_k, q_2 \rangle q_2 - \ldots - \langle a_k, q_{k-1} \rangle q_{k-1}$. Then normalize by $q_k = \frac{q_k}{\|q_k\|}$.

We can observe the following properties:

1. $\|q_i\| = 1$ (this follows directly)

2. $q_i \perp q_j$ for all $i \neq j$

3. $q_k \in span(a_1, \ldots, a_k)$ and $a_k \in span(q_1, \ldots, q_k)$ so that $span(a_1, \ldots, a_k) = span(q_1, \ldots, q_k)$.

[[Write proof for 2]].

## 3.4 QR Factorization

**Definition 7.** (Unitary Matrix) A matrix $Q = [q_1 \ldots q_n] \in \mathbb{R}^{m \times n}$ is unitary if and only if $\langle q_i, q_j \rangle = \delta_{ij}$.

Observations about this definition:

1. $Q^T Q = I$

2. If $Q$ is square, then $Q^T = Q^{-1}$.

### 3.4.1 Application to Least Squares

Suppose that we can write $A = QR$, where $A \in \mathbb{R}^{m \times n}$, $Q \in \mathbb{R}^{m \times n}$ and unitary, and $R \in \mathbb{R}^{n \times n}$ and upper triangular. Then the least squares solution to $Ax = b$ is given by

$$
\begin{aligned}
x &= (A^T A)^{-1} A^T b \\
&= (R^T Q^T Q R)^{-1} R^T Q^T b \\
&= (R^T R)^{-1} R^T Q^T b \\
\implies (R^T R)x &= R^T Q^T b \\
Rx &= Q^T b \qquad \text{(assume $R$ is invertible (i.e. no zeros on the diagonal))}
\end{aligned}
$$

We can then solve for $x$ using back substitution, which is $\mathcal{O}(n^2)$.

## 3.5 Norms and Condition Numbers

**Definition 8.** (Norm) Suppose that $\mathcal{V}$ is a linear space over the field $\mathbb{R}$. The *nonnegative* real-valued function $\|\cdot\|$ is a norm on $\mathcal{V}$ if the following axioms are satisfied: Fix $v \in \mathcal{V}$

1. Positivity: $\|v\| = 0$ if and only if $v = 0$

2. Scale Preservation: $\|\alpha v\| = |\alpha| \|v\|$ for all $\alpha \in \mathbb{R}$

3. Triangle Inequality: $\|v + w\| \leq \|v\| + \|w\|$.

**Example 1.** (Examples of Norms)

1. 1-norm:
$$\|v\|_1 = \sum_{i=1}^{n} |v_i| = |v_1| + \cdots + |v_n| \tag{30}$$

2. 2-norm:
$$\|v\|_2 = \left( \sum_{i=1}^{n} v_i^2 \right)^{\frac{1}{2}} = \sqrt{v_1^2 + \cdots + v_n^2} = \sqrt{v^T v} \tag{31}$$

3. $\infty$-norm
$$\|x\|_\infty = \max_{i=1,\ldots,n} |v_i| \tag{32}$$

4. $p$-norm
$$\|v\|_p = \left( \sum_{i=1}^{n} |v_i|^p \right)^{\frac{1}{p}} \tag{33}$$

For the $p$-norm, proving the triangle inequality follows from the Minkowski's inequality.

**Definition 9.** (Operator Norm) Let $A$ be an $m \times n$ matrix. That is, $A$ is a linear transformation form $\mathbb{R}^n$ to $\mathbb{R}^m$. Then the operator norm (or subordinate matrix norm) of $A$ is
$$\|A\|_{p,q} = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|_q}{\|x\|_p}. \tag{34}$$

Observations about this definition:

1. It's easy to check that this definition of the operator norm satisfies the properties of a norm given in Definition 8. For the triangle inequality, observe that

$$\|(A+B)x\|_p \leq \|Ax\|_p + \|Bx\|_p \qquad \text{(from Minkowski's inequality)}$$
$$\implies \frac{\|(A+B)x\|_p}{\|x\|_p} \leq \frac{\|Ax\|_p}{\|x\|_p} + \frac{\|Bx\|_p}{\|x\|_p}$$

Taking the supremum of both sides over $x$ shows that $\|A + B\|_p \leq \|A\|_p + \|B\|_p$.

11

2. The definition immediately implies that for an arbitrary $x \in \mathbb{R}^n$, $x \neq 0$,

$$\|Ax\|_q \leq \|A\|_{p,q}\|x\|_p \tag{35}$$

We can generalize this inequality to claim that

$$\|AB\| \leq \|A\|\|B\| \tag{36}$$

for conformable matrices $A, B$. Indeed, fix $0 \neq x \in R^n$. Then

$$\|ABx\| \leq \|A\|\|Bx\| \leq \|A\|\|B\|\|x\| \tag{37}$$

We can divide all inequalities by $\|x\|$ to see that for all $x \neq 0$,

$$\frac{\|ABx\|}{\|x\|} \leq \|A\|\|B\| \tag{38}$$

Taking the supremum over $x$ on the left hand side shows that $\|AB\| \leq \|A\|\|B\|$.

**Theorem 8.** (The 1-norm of a matrix is the largest absolute-value column sum) Let $A \in \mathbb{R}^{m \times n}$ and denote the columns of $A$ by $a_j$, $j = 1, \ldots, n$. Then $\|A\|_1 = \max_{j=1,\ldots,n} \sum_{i=1}^m |a_{ij}| = \max_{j=1,\ldots,n}\|a_j\|$.

*Proof.* Fix $x \in \mathbb{R}^n$. Let $C = \max_{j=1,\ldots,n} \sum_{i=1}^m |a_{ij}|$. First consider the product $A \cdot x$. The $i$th element is $\sum_{j=1}^n a_{ij}x_j$. Then

$$\|Ax\|_1 = \sum_{i=1}^m |(Ax)_i| = \sum_{i=1}^m |\sum_{j=1}^n a_{ij}x_j|$$

$$\leq \sum_{i=1}^m \sum_{j=1}^n |a_{ij}||x_j| \qquad \text{(triangle inequality)}$$

$$= \sum_{j=1}^n |x_j| \left( \sum_{i=1}^m |a_{ij}| \right) \qquad \text{(interchange order of summation, assumed finite)}$$

$$\leq C\|x\|_1$$

Therefore $\frac{\|Ax\|_1}{\|x\|_1} \leq C$ for all $x$. Next, we find an $x$ such we achieve equality with $C$. Call index $J$ the index such that $\|a_J\|_1 = C = \max_{j=1,\ldots,n} \sum_{i=1}^m |a_{ij}|$. Then let $e_J$ be the $n$-vector of zeros with a 1 in the $J$th entry. Clearly $\|e_J\|_1 = 1$. But then

$$\|Ae_J\|_1 = \|a_J\|_1 = C \tag{39}$$

In sum, we first showed that for all $x \in \mathbb{R}^n$

$$\frac{\|Ax\|_1}{\|x\|_1} \leq C \tag{40}$$

We then found an $x \in \mathbb{R}^n$ such that $\frac{\|Ax\|_1}{\|x\|_1} = C$. Therefore

$$\|A\|_1 = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|_1}{\|x\|_1} = C = \max_{j=1,\ldots,n} \sum_{i=1}^{m} |a_{ij}| = \max_{j=1,\ldots,n} \|a_j\| \tag{41}$$

$\square$

**Theorem 9.** (The $\infty$-norm of a matrix is the largest absolute-value row sum) Let $A \in \mathbb{R}^{m \times n}$ and denote the rows of $A$ by $b_i$, $i = 1,\ldots,m$. Then $\|A\|_\infty = \max_{i=1,\ldots,m} \sum_{j=1}^{n} |a_{ij}| = \max_{i=1,\ldots,m} \|b_i\|$.

*Proof.* Fix $x \in \mathbb{R}^n$. Let $C = \max_{i=1,\ldots,m} \sum_{j=1}^{n} |a_{ij}|$.

$$\|Ax\|_\infty = \max_{i=1,\ldots,m} |\sum_{j=1}^{n} a_{ij}x_j|$$

$$\leq \max_{i=1,\ldots,m} \sum_{j=1}^{n} |a_{ij}||x_j| \qquad \text{(by the triangle inequality)}$$

$$\leq \max_{i=1,\ldots,m} \sum_{j=1}^{n} |a_{ij}|\|x\|_\infty \qquad \text{(since } |x_j| \leq \|x\|_\infty \text{ for all } j\text{)}$$

$$= C\|x\|_\infty$$

Next, we find an $x$ such we achieve equality with $C$. Call $I$ the index for which $\|b_I\|_\infty = C$. Define

$$x_j = \begin{cases} 1 & a_{Ij} > 0 \\ -1 & a_{Ij} < 0 \end{cases} \tag{42}$$

Observe that $\|x\|_\infty = 1$. Then

$$|A \cdot x|_I = |b_I^T \cdot x|$$

$$= |\sum_{j=1}^{m} a_{Ij}x_j|$$

$$= |\sum_{j=1}^{m} a_{Ij}|$$

$$= C$$

We then found an $x \in \mathbb{R}^n$ such that $\frac{\|Ax\|_\infty}{\|x\|_\infty} = C$. Therefore

$$\|A\|_\infty = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_\infty} = C = \max_{i=1,\ldots,m} \sum_{j=1}^{n} |a_{ij}| = \max_{i=1,\ldots,m} \|b_i\| \tag{43}$$

$\square$

**Theorem 10.** (The 2-norm of a symmetric positive definite matrix is the maximum absolute value of its eigenvalues) Let $A$ be a positive definite $n \times n$ matrix. Then

$$\|A\|_2 = \max_{i=1,\ldots,n} |\lambda_i| \tag{44}$$

*Proof.* Since $A$ is positive definite, $A$ has $n$ distinct eigenvalues, which implies that it has $n$ linearly independent eigenvectors. Therefore, for an arbitrary $x \in \mathbb{R}^n$, we can write $x$ as a linearly combination of the eigenvectors $x_1, \ldots, x_n$. Then

$$x = c_1 x_1 + \cdots + c_n x_n$$
$$Ax = c_1 A x_1 + \cdots + c_n A x_n$$
$$= c_1 \lambda_1 x_2 + \cdots + c_n \lambda_n x_n$$

We can normalize the eigenvectors of $A$ so that $x_i^T x_i = 1$. Then $\|Ax\|_2 = \sqrt{\sum_{i=1}^n c_i^2 \lambda_i^2}$ and $\|x\|_2 = \sqrt{\sum_{i=1}^n c_i^2}$. Therefore

$$\frac{\|Ax\|_2}{\|x\|_2} = \sqrt{\frac{\sum_{i=1}^n c_i^2 \lambda_i^2}{\sum_{i=1}^n c_i^2}} \le \max_i |\lambda_i| = |\lambda_I| \tag{45}$$

Now we'll find an $x$ such that we actually achieve equality. Call $I$ the index of the maximum absolute value of an eigenvalue. Then, consider the eigenvector associated with this eigenvalue, called $x_i$. Then

$$\frac{\|A x_I\|_2}{\|x_I\|_2} = \frac{|\lambda_I| \|x_I\|}{\|x_I\|} = |\lambda_I| \tag{46}$$

This shows that $\|A\|_2 = \max_i |\lambda_i|$. $\qquad \square$

**Theorem 11.** (The 2-norm of a matrix $A_{m \times n}$ equals its largest singular value) Let $A$ be an $m \times n$ matrix and denote the eigenvalues of the matrix $B = A^T A$ by $\lambda_i$, $i = 1, \ldots, n$. Then

$$\|A\|_2 = \max_i \sqrt{\lambda_i} \tag{47}$$

The square roots of the (nonnegative) eigenvalues of $A^T A$ are referred to as the singular values of $A$.

### 3.5.1 Conditioning

Conditioning helps us quantify the sensitivity of the output to perturbations of the input. In what follows, let $f$ be a mapping from a subset $D$ of a normed linear space $\mathcal{V}$ to another normed linear space $\mathcal{W}$.

**Definition 10.** (Absolute Condition Number)

$$Cond(f) = \sup_{x,y \in D, x \neq y} \frac{\|f(x) - f(y)\|}{\|x - y\|} \tag{48}$$

**Definition 11.** (Absolute Local Condition Number)

$$Cond_x(f) = \sup_{x+\delta x \in D, \delta x \neq 0} \frac{\|f(x + \delta x) - f(x)\|}{\|\delta x\|} \tag{49}$$

The previous two definitions depend on the magnitudes of $f(x)$ and $x$. In applications, it's often better to rescale as follows

**Definition 12.** (Relative Local Condition Number)

$$cond_x(f) = \sup_{x+\delta x \in D, \delta x \neq 0} \frac{\|f(x + \delta x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|} \tag{50}$$

**In these definitions, if $f$ is differentiable then we can replace the differences with the appropriate derivatives.**

**Example 2.** (Example of condition numbers) Let $D$ be a subinterval of $[0, \infty)$ and $f(x) = \sqrt{x}$. Then $f'(x) = \frac{1}{2\sqrt{x}}$.

1. If $D = [1, 2]$, then $Cond(f) = \frac{1}{2}$.

2. If $D = [0, 1]$, then $Cond(f) = \infty$.

3. If $D = (0, \infty)$, then the the absolute local condition number of $f$ at $x \in D$ is

$$Cond_x(f) = \frac{1}{2\sqrt{x}} \tag{51}$$

Thus as $x \to$, $Cond_x(f) \to \infty$, and as $x \to \infty$, $Cond_x(f) \to 0$.

4. If $D = (0, \infty)$, then the relative local condition number of $f$ is $cond_x(f) = 1/2$ for all $x \in D$.

**Definition 13.** (Condition Number of a Nonsingular Matrix) The condition number of a nonsingular matrix $A$ is defined by

$$\kappa(A) = \|A\| \|A^{-1}\| \tag{52}$$

If $\kappa(A) \gg 1$, the matrix is said to be ill-conditioned.

Observations about this definition:

1. $\kappa(A) = \kappa(A^{-1})$

2. For all $A$, $\kappa(A) \geq 1$. This follows because

$$1 = \|I\| = \|AA^{-1}\| \leq \|A\|\|A^{-1}\| \tag{53}$$

3. The condition number of a matrix is unaffected by scaling all its elements by multiplying by a nonzero constant.

4. There is a condition number for each norm, and the size of the condition number is strongly dependent on the choice of norm.

# 4 Special Matrices

## 4.1 Symmetric Positive Definite Matrices

**Definition 14.** (Symmetric, Positive Definite, spd) The real matrix $A$ is said to be symmetric if $A = A^T$. A square $n \times n$ matrix is called positive definite if

$$x^T A x > 0 \tag{54}$$

for all $x \in \mathbb{R}^n$, $x \neq 0$.

**Theorem 12.** (Properties of spd matrices) Let $A$ be an $n \times n$ real, spd matrix. Then

1. $a_{ii} > 0$ for all $i = 1, \ldots, n$ (the diagonal elements of $A$ are positive).

2. $Ax_i = \lambda_i x_i \implies \lambda_i \in \mathbb{R}_{>0}, x \in \mathbb{R}^n \setminus \{0\}$ (the eigenvalues of $A$ are real and positive, and the eigenvectors of $A$ belong to $\mathbb{R}^n \setminus \{0\}$).

3. $x_i \perp x_j$ if $\lambda_i \neq \lambda_j$ (the eigenvectors of distinct eigenvalues of $A$ are orthogonal)

4. $\det(A) > 0$ (the determinant of $A$ is positive)

5. Every submatrix $B$ of $A$ obtained by deleting any set of rows and the corresponding set of columns from $A$ is symmetric and positive definite (in particular, every principal submatrix is positive definite).

*Proof.* We prove each claim in the theorem as follows

1. Let $e_i$ be the $i$th canonical basis vector in $\mathbb{R}^n$. Then

$$a_{ii} = e_i^T A e_i > 0 \tag{55}$$

since $A$ is pd. A few observations: this only relies on $A$ being pd. $e_i^T A$ picks out the $i$th row of $A$. $Ae_i$ picks out the $i$th column of $A$.

2. We'll first show that the eigenvalues of $A$ are real. Suppose $\lambda$, $x$ are an eigenvalue/vector pair of $A$. Thus $Ax = \lambda x$. We can conjugate this equation to find that $\bar{A}\bar{x} = A\bar{x} = \bar{\lambda}\bar{x}$ (thus complex eigenvalues of real valued matrices come in conjugate pairs). Then

$$x^T A \bar{x} = \bar{\lambda} x^T \bar{x}$$
$$x^T A^T \bar{x} = (Ax)^T \bar{x} = \lambda x^T \bar{x}$$

Since $A = A^T$, we know that $\lambda x^T \bar{x} = \bar{\lambda} x^T \bar{x}$. As long as $x \neq 0$, then $x^T \bar{x} \neq 0$. Therefore $\bar{\lambda} = \lambda$, which shows $\lambda \in \mathbb{R}$.

The fact that the eigenvector associated with $\lambda$ has real elements follows by noting that all elements of the singular matrix $A - \lambda I$ are real numbers. Therefore, the columns of $A - \lambda I$ are linearly dependent in $\mathbb{R}^n$. Hence there exists an $x \in \mathbb{R}^n$ such that $(A - \lambda I)x = 0$.

This proof only requires that $A$ is symmetric – therefore any real, symmetric matrix has real eigenvalues/vectors.

Next we'll show the eigenvalues of $A$ are positive. Suppose $\lambda$, $x$ are an eigenvalue/vector pair of $A$. Then

$$0 < x^T A x = \lambda x^T x \tag{56}$$

Since $x \neq 0$ and $x^T x$ is positive (it's actually the squared 2-norm of $x$), then $\lambda > 0$. Note that this part of the proves requires $A$ be pd.

3. Let $\lambda_i, \lambda_j$ be distinct eigenvalues of $A$, and $x_i, x_j$ the corresponding eigenvectors. Then

$$x_i^T A x_j = \lambda_j x_i^T x_j$$
$$x_i^T A^T x_j = (Ax_i)^T x_j = \lambda_i x_i^T x_j$$

Since $A$ is symmetric, these two string of equalities must be equal. We can subtract them to find that

$$(\lambda_i - \lambda_j) x_i^T x_j = 0 \tag{57}$$

Since we assumed $\lambda_i \neq \lambda_j$, then it must be that $x_i^T x_j$. Therefore $x_i \perp x_j$. This proof again only relies on the symmetry of $A$. which is the product of the diagonal elements of the matrix (the eigenvalues).

4. This follows from the fact that the determinant of $A$ is equal to the product of its eigenvalues. Or, we can write $A$ in terms of its eigenvalue decomposition. Thus

$$A = X \Lambda X^{-1} \tag{58}$$

Therefore

$$\det(A) = \det(X) \det(\Lambda) \det(X)^{-1} = \det(\Lambda) \tag{59}$$

Or, we can write $A$ in terms of its eigenvalue decomposition. Thus

$$A = X\Lambda X^{-1} \tag{60}$$

Therefore

$$\det(A) = \det(X)\det(\Lambda)\det(X)^{-1} = \det(\Lambda) \tag{61}$$

5. Let $I \subset \{1, 2, \ldots, n\}$ be a subset of indices and let $B = A_{II}$. $A$ is symmetric, so that $A_{II} = A_{II}^T$. Therefore $B$ is symmetric. Let $x \in \mathbb{R}^n$ and define a vector $y$ that is 0 for the indices not included in $I$ and follows the value of $x$ for the indices included in $I$. Therefore, $x^T B x = y^T A y > 0$ since $A$ is pd.

$\square$

## 4.2 Cholesky Factorization

**Theorem 13.** If $A$ is spd, then there exists a lower diagonal matrix $L$ such that $A = LL^T$. This is called the Cholesky decomposition.

---
**Algorithm 1** Cholesky Factorization

---
**Require:** $A \in \mathbb{R}^{n \times n}$, SPD
     $L_1 \leftarrow \sqrt{a_{11}}$
     **for** $k \leftarrow 2, 3, \ldots, n$ **do**
        Solve $L_{k-1} l_k = a_k$ for $l_k$
        $l_{kk} \leftarrow \sqrt{a_{kk} - l_k^T l_k}$
        $L_k \leftarrow \begin{pmatrix} L_{k-1} & 0 \\ l_k^T & l_{kk} \end{pmatrix}$
     **end for**

---

Notation:

- $L_{k-1}$: the first $k - 1 \times k - 1$ upper left corner of $L$

- $a_k$: the first $k - 1$ entries in column $k$ of $A$

- $l_k$: the first $k - 1$ entries in column $k$ of $L^T$ [[?]]

- $a_{kk}, l_{kk}$: the $kk$ entries of $A$ and $L$, respectively

## 4.3 Banded Matrices and Differential Equations

Consider the two-point boundary value problem

$$u'' + 2u' = -1, \quad u(x = 0) = 0, u(x = 1) = 0 \tag{62}$$

where $x \in [0, 1]$.

Define a sequence of grid points $\{x_i\}_{i=0}^{N+1}$. We can approximate the derivative of $u$ at each point on the grid as follows

$$u'(x_j) = \lim_{\delta \to 0} \frac{u(x_j + \delta) - u(x_j - \delta)}{2\delta}$$
$$\approx \frac{u(x_{j-1}) - u(x_{j+1})}{2\Delta x}$$

where we use the centered difference quotient of order 2. Similarly, we can approximate the second derivative of $u$ each each point in the domain as

$$u''(x_j) = \lim_{\delta \to 0} \frac{u'(x_j + \delta) - u'(x_j - \delta)}{2\delta}$$
$$\approx \frac{u'(x_{j+1}) - u'(x_{j-1})}{2\Delta x}$$
$$= \frac{\frac{u(x_{j+2}) - u(x_j)}{2\Delta x} - \frac{u(x_j) - u(x_{j-2})}{2\Delta x}}{2\Delta x}$$
$$= \frac{u(x_{j+2}) - 2u(x_j) + u(x_{j-2})}{4\Delta x}$$

Let's instead use the grid points adjacent to $x_j$:

$$u''(x_j) \approx \frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1})}{\Delta x} \tag{63}$$

Then, going back to the initial differential equation, for $x_j$, we have

$$\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1})}{\Delta x} + 2\frac{u(x_{j-1}) - u(x_{j+1})}{2\Delta x} = -1 \tag{64}$$

Let

$$U = \begin{bmatrix} u(x_1) \\ u(x_2) \\ \vdots \\ u(x_N) \end{bmatrix} \tag{65}$$

and let $u_i = u(x_i)$. In this notation, the differential equation at $x_j$ can be written as

$$\frac{u_{j+1} - 2u_j + u_{j-1}}{\Delta x^2} + \frac{u_{j-1} - u_{j+1}}{\Delta x} = -1 \tag{66}$$

We can put these equations together into a matrix. Each row will only have 3 non-zero

entries at $j - 1$, $j$, and $j + 1$. Thus the $j$th row is

$$\begin{pmatrix} 0 & 0 & \cdots & \frac{1}{\Delta x^2} + \frac{1}{\Delta x} & \frac{-2}{\Delta x^2} & \frac{1}{\Delta x^2} - \frac{1}{\Delta x} & 0 & 0 & \cdots \end{pmatrix} \tag{67}$$

Thus stacking these rows together will give a tridiagonal matrix. Call this matrix $A$. Then we have that

$$AU = -1 \tag{68}$$

# 5 Simultaneous nonlinear equations

## 5.1 Analysis Preliminaries

**Definition 15.** (Cauchy Sequence) A sequence $(x^{(k)}) \subset \mathbb{R}^n$ is called a Cauchy sequence in $\mathbb{R}^n$ if for any $\epsilon > 0$ there exists a positive integer $k_0 = k_0(\epsilon)$ such that

$$\|x^{(k)} - x^{(m)}\|_\infty < \epsilon \quad \forall k, m \geq k_0(\epsilon) \tag{69}$$

**Remark 1.** $\mathbb{R}^n$ is **complete** in the sense that every Cauchy sequence $(x^{(k)})$ converges to some $\tilde{\zeta} \in \mathbb{R}^n$.

**Definition 16.** (Continuous function) Let $D \subset \mathbb{R}^n$ be nonempty and $f : D \to \mathbb{R}^n$. Given $\xi \in D$, $f$ is continuous at $\xi$ if for every $\epsilon > 0$, there exists a $\delta = \delta(\epsilon) > 0$ such that for every $x \in B(\xi; \delta) \cap D$

$$\|f(x) - f(\xi)\|_\infty < \epsilon \tag{70}$$

**Lemma 1.** Let $D \subset \mathbb{R}^n$ be nonempty and $f : D \to \mathbb{R}^n$ be defined and continuous on $D$. If $(x^{(k)}) \subset D$ converges in $\mathbb{R}^n$ to $\xi \in D$, then $f(x^{(k)})$ also converges to $f(\xi)$.

We want to find a vector $x \in \mathbb{R}^n$ such that $f(x) = 0$.

**Example 3.** Consider the linear system

$$Ax = b \tag{71}$$

Then $A : \mathbb{R}^n \to \mathbb{R}^m$. Let $f(x) = Ax - b$.

**Example 4.** Let

$$f = \begin{bmatrix} x_1^2 + x_2^2 - 1 \\ 5x_1^2 + 21x_2^2 - 9 \end{bmatrix} \tag{72}$$

Note that $x_1^2 + x_2^2 = 1$ is the 0 level set of $f$, and is the unit circle. $5x_1^2 + 21x_2^2 = 9$ is the 0 level set of $f$ and is an ellipse.

This function has four zeros

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \pm\sqrt{3}/2 \\ \pm\frac{1}{2} \end{bmatrix} \tag{73}$$

## 5.2 Simultaneous iteration

**Example 5.**

**Definition 17.** (Lipschitz condition, constant, and contraction) Let $D$ be a closed subset of $\mathbb{R}^n$ and $g : D \to D$. If there exists a positive constant $L$ such that

$$\|g(x) - g(y)\|_\infty \leq L\|x - y\|_\infty \tag{74}$$

for all $x, y \in D$, then $g$ satisfies the Lipschitz condition on $D$ in the $\infty$-norm. $L$ is called the Lipschitz constant. If $L \in (0, 1)$, then $g$ is called a contraction on $D$ in the $\infty$-norm.

Observations about this definition:

- Any function $g$ that satisfies the Lipschitz condition on $D$ is continuous on $D$ (to see this, set $\delta = \frac{\epsilon}{L}$).

- If $g$ satisfies the Lipschitz condition on $D$ in the $\infty$-norm, then it also does in the $p$-norm for $p \in [1, \infty)$ and vice-versa. However the size of $L$ depends on the choice of norm.

**Theorem 14.** (Contraction Mapping Theorem in $\mathbb{R}^n$) Suppose $D$ is a closed subset of $\mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}^n$ is defined on $D$, and $g(D) \subset D$. Suppose further that $g$ is a contraction on $D$ in the $\infty$-norm. Then,

1. $g$ has a unique fixed point $\boldsymbol{\xi} \in D$

2. The sequence $(\boldsymbol{x}^{(k)})$ defined by $\boldsymbol{x}^{(k+1)} = g(\boldsymbol{x}^k)$ converges to $\boldsymbol{\xi}$ for any starting value $x^{(0)} \in D$.

*Proof.* The proof has three parts:

1. First prove uniqueness, assuming existence of a fixed point.

2. Prove the iteration generates a Cauchy sequence (then convergence to some $\boldsymbol{\xi}$ follows from the completeness of the space).

3. Show $\boldsymbol{\xi}$ is indeed the fixed point.

Uniqueness: Suppose $\boldsymbol{\xi}, \boldsymbol{\eta}$ are both fixed points of $g$ in $D$. Then,

$$\begin{aligned}
\|\boldsymbol{\xi} - \boldsymbol{\eta}\|_\infty &= \|g(\boldsymbol{\xi}) - g(\boldsymbol{\eta})\| && (\boldsymbol{\xi}, \boldsymbol{\eta} \text{ are fixed points}) \\
&\leq L\|\boldsymbol{\xi} - \boldsymbol{\eta}\|_\infty && (g \text{ is a contraction on } D)
\end{aligned}$$

We can rearrange this to see that $(1 - L)\|\boldsymbol{\xi} - \boldsymbol{\eta}\|_\infty \leq 0$. By assumption, $L \in (0, 1)$, and the norm of a quantity is always weakly positive. Therefore, $\|\boldsymbol{\xi} - \boldsymbol{\eta}\|_\infty = 0$ which implies $\boldsymbol{\xi} = \boldsymbol{\eta}$.

Convergence: Assuming $g$ has a fixed point $\boldsymbol{\xi} \in D$, the sequence $\boldsymbol{x}^{(k+1)} = g(\boldsymbol{x}^k)$ will converge to $\boldsymbol{\xi}$ for any $\boldsymbol{x}^{(0)} \in D$. This follows because

$$\|\boldsymbol{x}^{(k)} - \boldsymbol{\xi}\|_\infty \leq L^k \frac{1}{1-L} \|\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(0)}\|_\infty \tag{75}$$

Since $L \in (0,1)$, $\lim_{k\to\infty} L^k = 0$, and therefore

$$\lim_{k\to\infty} \|\boldsymbol{x}^{(k)} - \boldsymbol{\xi}\|_\infty = 0 \tag{76}$$

Existence: First observe that if $\boldsymbol{x}^{(0)}$ belongs to $D$, then $\boldsymbol{x}^{(k+1)} = g(\boldsymbol{x}^k) \in D$ for all $k \geq 1$ since $g(D) \subset D$ (this is important since the proof relies on $g$ being a contraction on $D$). Next, consider the distance between two adjacent terms on the sequence $\boldsymbol{x}^{(k+1)} = g(\boldsymbol{x}^k)$

$$\begin{aligned}
\|\boldsymbol{x}^{(k)} - \boldsymbol{x}^{(k-1)}\|_\infty &= \|g(\boldsymbol{x}^{(k-1)}) - g(\boldsymbol{x}^{(k-2)})\|_\infty && \text{(definition of } g\text{)} \\
&\leq L\|\boldsymbol{x}^{(k-1)} - \boldsymbol{x}^{(k-2)}\|_\infty && (g \text{ is a contraction on } D) \\
&\leq L^{k-1}\|\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(0)}\|_\infty && \text{(induction)}
\end{aligned}$$

Now, fix positive integers $m, k$ such that $m > k$. Then

$$\begin{aligned}
\|\boldsymbol{x}^{(m)} - \boldsymbol{x}^{(k)}\|_\infty &= \|\boldsymbol{x}^{(m)} - \boldsymbol{x}^{(m-1)} + \boldsymbol{x}^{(m-1)} + \cdots + \boldsymbol{x}^{(k-1)} - \boldsymbol{x}^{(k)}\|_\infty \\
&\leq \|\boldsymbol{x}^{(m)} - \boldsymbol{x}^{(m-1)}\|_\infty + \cdots + \|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\|_\infty && \text{(triangle inequality)} \\
&\leq (L^{m-1} + \cdots + L^k)\|\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(0)}\|_\infty && (g \text{ a contraction}) \\
&= L^k(L^{m-k-1} + \cdots + 1)\|\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(0)}\|_\infty \\
&\leq L^k \frac{1}{1-L}\|\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(0)}\|_\infty && \text{(geometric series)}
\end{aligned}$$

Since $L \in (0,1)$, $\lim_{k\to\infty} L^k = 0$. Therefore, $\boldsymbol{x}^{(k)}$ is a Cauchy sequence in $\mathbb{R}^n$, that is for all $\epsilon > 0$, there exists a $k_0$ such that

$$\|\boldsymbol{x}^{(m)} - \boldsymbol{x}^{(k)}\|_\infty < \epsilon \quad \forall m, k \geq k_0 \tag{77}$$

Any Cauchy sequence in $\mathbb{R}^n$ is convergent in $\mathbb{R}^n$. Thus, there exists some $\boldsymbol{\xi} \in \mathbb{R}^n$ such that $\boldsymbol{\xi} = \lim_{k\to\infty} \boldsymbol{x}^{(k)}$.

$\boldsymbol{\xi}$ is indeed the fixed point: Since $g$ satisfies the Lipschitz condition on $D$, $g$ is continuous on $D$. Therefore,

$$\boldsymbol{\xi} = \lim_{k\to\infty} \boldsymbol{x}^{(k+1)} = \lim_{k\to\infty} g(\boldsymbol{x}^{(k)}) = g\left(\lim_{k\to\infty} \boldsymbol{x}^{(k)}\right) = g(\boldsymbol{\xi}) \tag{78}$$

therefore $\boldsymbol{\xi}$ is a fixed point of $g$, and observe that $\boldsymbol{\xi} \in D$ since $D$ is closed. $\qquad\square$

**Definition 18.** (Jacobian) Let $g = (g_1, \ldots, g_n)^T : \mathbb{R}^n \to \mathbb{R}^n$ be a function defined and

continuous in an (open) neighborhood of $\boldsymbol{\xi} \in \mathbb{R}^n$. Suppose the first partial derivatives of each $g_i$ exist at $\boldsymbol{\xi}$. The Jacobian matrix $J_g(\boldsymbol{\xi})$ of $g$ at $\boldsymbol{\xi}$ is the $n \times n$ matrix with elements

$$J_g(\boldsymbol{\xi})_{ij} = \frac{\partial g_i}{\partial x_j}(\boldsymbol{\xi}) \tag{79}$$

**Theorem 15.** Let $g = (g_1, \ldots, g_n)^T : \mathbb{R}^n \to \mathbb{R}^n$ be a function defined and continuous on a closed set $D \subset \mathbb{R}^n$. Let $\boldsymbol{\xi} \in D$ be a fixed point of $g$. Suppose the first partial derivatives of each $g_i$ are defined and continuous in some (open) neighborhood $N(\boldsymbol{\xi}) \in D$ of $\boldsymbol{\xi}$, with

$$\|J_g(\boldsymbol{\xi})\|_\infty < 1 \tag{80}$$

Then there exists $\epsilon > 0$ such that $g(\bar{B}_\epsilon(\boldsymbol{\xi})) \subset \bar{B}_\epsilon(\boldsymbol{\xi})$, and the sequence $x^{(k+1)} = g(x^k)$ converges to $\boldsymbol{\xi}$ for all $x^{(0)} \in \bar{B}_\epsilon(\boldsymbol{\xi})$ (in other words, the sequence converges to $\boldsymbol{\xi}$ as long as $x^{(0)}$ is close enough to $\boldsymbol{\xi}$).

**Example 6.**

## Newton's Method

**Definition 19.** (Newton's Method) The sequence defined by

$$x^{(k+1)} = x^{(k)} - [J_f(x^{(k)})]^{-1} f(x^{(k)}) \tag{81}$$

where $x^{(0)} \in \mathbb{R}^n$, is called Newton's method.

**Theorem 16.** Suppose $f(\boldsymbol{\xi}) = 0$, that in some (open) neighborhood $N(\boldsymbol{\xi})$ of $\boldsymbol{\xi}$, where $f$ is defined and continuous, all the second-order partial derivatives of $f$ are defined and continuous, and that the Jacobian matrix $J_f(x^{(k)})$ of $f$ at the point $\boldsymbol{\xi}$ is nonsingular. Then the sequence defined by Newton's method converges to $\boldsymbol{\xi}$ provided that $x^{(0)}$ is sufficiently close to $\boldsymbol{\xi}$.

# 6 Eigenvalues of Eigenvectors of a symmetric matrix

Another matrix decomposition is

$$A = X \Lambda X^{-1} \tag{82}$$

where $X$ is a matrix of the eigenvectors and $\Lambda$ is a diagonal matrix with the eigenvalues.

## 6.1 Why we use iteration to calculate eigenvalues/eigenvectors

We call $\lambda$ an eigenvalue and $x \neq 0$ an eigenvector of $A$ if $Ax = \lambda x$. Thus, $(Ax - \lambda I) = 0$. Therefore, $x \in Null(A - \lambda I)$. Since $x \neq 0$, $A - \lambda I$ has a non-trivial nullspace, so we

must have $\det(A - \lambda I) = 0$. This suggests a way to transform an eigenvalue finding problem to a root finding problem. Define

$$\rho(\lambda) = \det(A - \lambda I) \tag{83}$$

Recall that the determinant of a matrix is the product of its eigenvalues. If $A$ is a $n \times n$ real, symmetric matrix, then $\rho(\lambda)$ is an $n$-th order polynomial in $\lambda$, whose roots are the eigenvalues of $A$.

**Theorem 17.** (Abel(-Ruffini) Theorem, or "No-go Theorem") There is no algebraic solution (that is, a solution expressed in terms of radicals) to general polynomial equations of degree five or higher with arbitrary coefficients.

Therefore, there is no finite-number operation procedure that provides an eigenvalue decomposition.

## 6.2  Power Iteration

**Find the biggest eigenvalue/vector.**

---
**Algorithm 2** Power Iteration

---
**Require:** $v^{(0)} = $ some vector with $\|v^{(0)}\| = 1$
  1: **for** $k \leftarrow 1, 2, \ldots$ **do**
  2:      $w \leftarrow A v^{(k-1)}$                          $\triangleright$ Apply $A$
  3:      $v^{(k)} \leftarrow w / \|w\|$                     $\triangleright$ Normalize
  4:      $\lambda^{(k)} \leftarrow (v^{(k)})^T A v^{(k)} = \langle v^{(k)}, A v^{(k)} \rangle$      $\triangleright$ Rayleigh Quotient
  5: **end for**

---

**Theorem 18.** (Convergence of Power Iteration) Suppose $|\lambda_1| > |\lambda_2| \geq \ldots \geq |\lambda_n|$ and $q_1^T v^{(0)} \neq 0$. Then the iterates of power iteration satisfy

$$\|v^{(k)} - (\pm q_1)\| = \mathcal{O}\left( \left| \frac{\lambda_2}{\lambda_1} \right|^k \right) \qquad \text{(error of eigenvector)}$$

$$|\lambda^{(k)} - \lambda_1| = \mathcal{O}\left( \left| \frac{\lambda_2}{\lambda_1} \right|^{2k} \right) \qquad \text{(error of eigenvalue)}$$

*Proof.* Convergence of eigenvector: Write $v^{(0)} = v$ as a linear combination of the orthonormal eigenvectors $q_i$:

$$v = c_1 q_1 + \cdots + c_n q_n \tag{84}$$

$v^{(k)}$ is a scalar multiple of $A^k v^{(0)}$. Therefore

$$v^{(k)} = \alpha_k A^k v^{(0)} \qquad\qquad (\alpha_k \text{ a normalization constant})$$

$$= \alpha_k \left( \sum_{i=1}^n \lambda_i^k a_i q_i \right)$$

$$= \alpha_k \lambda_1^k \left( c_1 q_1 + c_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k q_2 + \ldots + c_n \left( \frac{\lambda_n}{\lambda_1} \right)^k q_n \right)$$

We can choose $\alpha_k$ such that $\alpha_k \lambda_1^k$ is 1. Therefore, $c_1 q_1$ is dominating (as long as $c_1 \neq 0$). The other terms are of order $\mathcal{O} \left( \left| \frac{\lambda_2}{\lambda_1} \right|^k \right)$.

Convergence of eigenvalue: see proposition below.

$\square$

**Theorem 19.** (Error of Rayleigh Quotient) Let $x_1$ be the eigenvector that corresponds to the largest (in absolute value) eigenvalue. If $\|x - x_1\| = \mathcal{O}(\epsilon)$, then

$$\left| \frac{\langle x, Ax \rangle}{\langle x, x \rangle} - \lambda_1 \right| = \mathcal{O}(\epsilon^2) \tag{85}$$

*Proof.* TODO. $\square$

## 6.3 Inverse Iteration

**Find the smallest eigenvalue/vector.**

## 6.4 Simultaneous Iteration

**Obtain the full set of eigenvalues/vectors simultaneously.**

---
**Algorithm 3** Simultaneous Iteration

---
**Require:** $Q^{(0)} = V = I$, a list of vectors $V$, which we choose to be the identity
1: **for** $k \leftarrow 1, 2, \ldots$ **do**
2: $\quad Z \leftarrow A \underline{Q}^{(k-1)}$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \triangleright$ Apply $A$
3: $\quad Z \leftarrow \underline{Q}^{(k)} R^{(k)}$ $\qquad\qquad\qquad\qquad\qquad\qquad \triangleright$ QR factorization of $Z$
4: $\quad A^{(k)} \leftarrow (\underline{Q}^{(k)})^T A \underline{Q}^{(k)}$ $\qquad\qquad\qquad \triangleright A_{ii}^{(k)} = \langle q_i^{(k)}, A q_i^{(k)} \rangle$
5: **end for**

---

Intuitively,

$$A^K \cdot V = \left[ \sum_i \lambda_i^k c_{1i} \tilde{q}_i \;\; \middle| \;\; \sum_i \lambda_i^k c_{2i} \tilde{q}_i \;\; \middle| \;\; \cdots \right] \tag{86}$$

The first column vector will converge to $\tilde{q}_1$. The second vector will converge to $\tilde{q}_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)\tilde{q}_2$.

$\underline{Q}^{(k)}$ will converge to the matrix of eigenvectors:

$$X = \begin{bmatrix} x_1 & | & x_2 & | & \cdots & | & x_n \end{bmatrix} \tag{87}$$

$\underline{A}^{(k)}$ will converge to a diagonal matrix containing the eigenvalues.

## 6.5 Shifted Power Iteration

**Find the eigenvalue close to a specific number.**

## 6.6 QR Algorithm

The $QR$ can be viewed as a stable procedure for computing $QR$ factorizations of the matrix powers $A, A^2, A^3, \dots$

---

**Algorithm 4** QR Algorithm (without shifts)

---

**Require:** $A^{(0)} = A$
1: **for** $k \leftarrow 1, 2, \dots$ **do**
2: $\quad Q^{(k)}R^{(k)} \leftarrow A^{(k-1)}$ $\qquad\qquad\qquad\qquad$ ▷ $QR$ factorization of $A^{(k-1)}$
3: $\quad A^{(k)} \leftarrow R^{(k)}Q^{(k)}$ $\qquad\qquad\qquad$ ▷ Recombine factors in reverse order
4: **end for**

---

## 6.7 Simultaneous Iteration equivalent to QR Algorithm

The $QR$ algorithm is equivalent to simultaneous iteration applied to a full set of initial vectors, namely, $\hat{Q}^{(0)} = I$. Summary of each algorithm:

**Simultaneous Iteration**

$$
\begin{aligned}
\underline{Q}^{(0)} &= I && \text{(initial condition)} \\
Z &= A\underline{Q}^{(k-1)} && \text{(apply } A\text{)} \\
Z &= \underline{Q}^{(k)}R^{(k)} && \text{(resemblance of normalization, } QR \text{ factorization of } Z\text{)} \\
A^{(k)} &= (\underline{Q}^{(k)})^T A\underline{Q}^{(k)} && \text{(resemblance of Rayleigh quotient)}
\end{aligned}
$$

## QR Algorithm

$$A^{(0)} = A \qquad \text{(initial condition)}$$
$$A^{(k-1)} = Q^{(k)} R^{(k)} \qquad \text{(compute } QR \text{ factorization)}$$
$$A^{(k)} = R^{(k)} Q^{(k)} \qquad \text{(reverse order of factors)}$$
$$\underline{Q}^{(k)} = Q^{(1)} Q^{(2)} \cdots Q^{(k)} \qquad \text{(definition of } \underline{Q}^{(k)})$$

and
$$\underline{R}^{(k)} = R^{(k)} R^{(k-1)} \cdots R^{(1)} \qquad \text{(definition of } \underline{R}^{(k)})$$

**Theorem 20.** (Equivalence of Simultaneous Iteration and the QR Algorithm) Simultaneous Iteration and the QR Algorithm generate identical sequences of matrices $\underline{R}^{(k)}, \underline{Q}^{(k)}, A^{(k)}$. Both give

$$(a) : A^{(k)} = \underline{Q}^{(k)} \underline{R}^{(k)} \qquad (QR \text{ factorization of the } k\text{th power of } A)$$
$$(b) : A^{(k)} = (\underline{Q}^{(k)})^T A \underline{Q}^{(k)} \qquad \text{(projection)}$$

*Proof.* By induction on $k$ (number of iterations). The base case $k = 0$ is trivial.

1. *QR gives* $(a)$: Assume $A^{(k-1)} = \underline{Q}^{(k-1)} \underline{R}^{(k-1)}$. The inductive hypothesis for $(b)$ gives that $A^{(k-1)} = (\underline{Q}^{(k-1)})^T A \underline{Q}^{(k-1)}$ or that $\underline{Q}^{(k-1)} A^{(k-1)} = A \underline{Q}^{(k-1)}$. Then

$$A^{(k)} = A A^{(k-1)} \qquad \text{(decompose to use inductive hypothesis)}$$
$$= A \underline{Q}^{(k-1)} \underline{R}^{(k-1)} \qquad \text{(inductive hypothesis)}$$
$$= \underline{Q}^{(k-1)} A^{(k-1)} \underline{R}^{(k-1)} \qquad \text{(inductive hypothesis from } (b))$$
$$= \underline{Q}^{(k-1)} R^{(k-1)} Q^{(k-1)} \underline{R}^{(k-1)} \qquad \text{(from algorithm)}$$
$$= \underline{Q}^{(k)} \underline{R}^{(k)} \qquad \text{(from definitions of } \underline{Q}^{(k)}, \underline{R}^{(k)})$$

2. *QR gives* $(b)$: Assume $A^{(k-1)} = (\underline{Q}^{(k-1)})^T A \underline{Q}^{(k-1)}$. From the relationship $A^{(k-1)} = Q^{(k)} R^{(k)}$ and the fact that $Q^{(k)}$ is orthogonal, we can apply $(Q^{(k)})^T$ to both sides (on the left) to get that $(Q^{(k)})^T A^{(k-1)} = R^{(k)}$. Then

$$A^{(k)} = R^{(k)} Q^{(k)}$$
$$= (Q^{(k)})^T A^{(k-1)} Q^{(k)}$$
$$= (Q^{(k)})^T (\underline{Q}^{(k-1)})^T A \underline{Q}^{(k-1)} Q^{(k)} \qquad \text{(inductive hypothesis)}$$
$$= (\underline{Q}^{(k)})^T A \underline{Q}^{(k)} \qquad \text{(definition of } \underline{Q}^{(k)})$$

$\square$