

Fall 2019

Course Summary

CSCC37 - Introduction to Numerical Algorithms for Computational Mathematics



Computer & Mathematical Sciences
UNIVERSITY OF TORONTO
S C A R B O R O U G H

Instructor: Richard Pancer
Email: pancer@utsc.utoronto.ca
Office: TBA
Office Hours: Mon 17:10 - 18:30

1 What is Numerical Analysis?

Physical Systems are typically transformed into **Mathematical Models** by engineers, chemists, physicists...

Mathematical Models can then have a closed form solution (rarely!), however most times there will only be a numerical solution which requires **Numerical Analysis**

2 Floating Point Arithmetic

2.1 Representation of non-negative integers

Decimal System (Base 10): $350 = (350)_{10} = 3 \cdot 10^2 + 5 \cdot 10^1 + 0 \cdot 10^0$

Binary System (Base 2): $350 = (101011110)_2 = 1 \cdot 2^8 + 0 \cdot 2^7 + \dots + 1 \cdot 2^1 + 0 \cdot 2^0$

In general, we have **Base b System**, for $b > 0, b \in \mathbb{N}$

$$x = (d_n d_{n-1} \dots d_0)_b = d_n \cdot b^n + d_{n-1} \cdot b^{n-1} + \dots + d_0 \cdot b^0$$

where $x \geq 0, x \in \mathbb{N}, 0 \leq d_i < b, i \in [1, n]$

Hexadecimal System (Base 16): Symbols $0, 1, \dots, 8, 9, A, B, C, D, E, F$

$$(8FA)_{16} = 8 \cdot 16^2 + F \cdot 16^1 + A \cdot 16^0 = 8 \cdot 16^2 + 15 \cdot 16^1 + 10 \cdot 16^0 = (2298)_{10}$$

Converting decimal to base b

Example, $(350)_{10}$ in base $b = 2$

Numerator	Denominator	Quotient	Remainder
350	2	175	0
175	2	87	1
87	2	43	1
43	2	21	1
21	2	10	1
10	2	5	0
5	2	2	1
2	2	1	0
1	2	0	1

Once we hit 0 in the **Quotients** column, we simply read the Remainder column from bottom up and we have our conversion. Thus $(350)_{10} = (101011110)_2$.

This algorithm is **safe** (will always terminate). Where **overflow** is the only potential problem.

2.2 Representation of real numbers

if $x \in \mathbb{R}$ then $x = \pm(x_I x_F)_b = \pm(d_n d_{n-1} \dots d_0 . d_{-1} d_{-2} \dots)_b$

Where x_I is the integral part and x_F is the fractional part. The sign (+ or -) is a single bit 0 or 1

x_I is the non-negative integer talked about above, and x_F can have infinite digits.

Example: $(0.77\dots)_{10} = (0.\bar{7})_{10} = 7 \cdot 10^{-1} + 7 \cdot 10^{-2} + \dots$

Binary System (Base 2): $(0.77\dots)_{10} = (.000110011\dots)_2 = (0.000\bar{1}1)_2 = 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} + \dots$

In general, we have **Base b System**, for $b > 0, b \in \mathbb{N}$

$$x_F = (.d_{-1} d_{-2} d_{-3} \dots)_b = d_{-1} \cdot b^{-1} + d_{-2} \cdot b^{-2} + \dots = \sum_{i=1}^{\infty} d_{-i} \cdot b^{-i}$$

We say that a **Binary Fraction** with n digits is terminating if it does not have any cycling and is finite, furthermore it has a terminating decimal representation.

Converting Decimal Fractions to base b **Example**, $(0.625)_{10}$ in base $b = 2$

Multiplier	Base	Product	Integral	Fraction
0.625	2	1.25	1	0.25
.25	2	0.5	0	0.5
0.5	2	1	1	0

Once we hit 0 in the **fractions** column, we simply read the Integral column from top down and we have our conversions. Thus $(0.625)_{10} = (.101)_2$

However you may not always hit 0 as this algorithm is **not safe** (won't always terminate)!

For example: What is $(0.1)_{10}$ in binary?

Multiplier	Base	Product	Integral	Fraction
0.1	2	0.2	0	0.2
0.2	2	0.4	0	0.4
0.4	2	0.8	0	0.8
0.8	2	1.6	1	0.6
0.6	2	1.2	1	0.2

We see that **0.2** occurs again as a fraction. This means that there must be a cycle! Thus $(0.1)_{10} = (0.000\bar{1}1)_2$

2.3 Machine Representation of Reals

Reals represented in computers as Floating Point numbers (FP for short).

A FP x in base b has the form: $x = (F)_b \cdot b^{(e)_b}$, where F is the fraction and has the form $F = \pm(d_1, d_2 \dots d_t)_b$ such F is also called the **mantissa**.

and $e = \pm(c_s c_{s-1} \dots c_1)_b$ is called the **exponent**

A FP number is **normalized** if $d_1 \neq 0$ unless $d_1 = d_2 = \dots = d_t = 0$

Significant Digits: (Of a nonzero FP) are the digits following and including the first nonzero digit. In the case of a normalized FP, all digits of the mantissa is significant (storing useful information).

Absolute value of mantissa is always ≥ 0 and < 1

Exponent is limited: $e \in [-M, M]$ where $M = (a_s a_{s-1} \dots a_1)_b$ with each $a_i = (b - 1)$

Largest FP number in absolute value defined by this system is $(.aa \dots a)_b \cdot b^{(aa \dots a)_b}$ where each $a = (b - 1)$

Smallest nonzero normalized FP number in absolute value is $(0.10 \dots 0)_b \cdot b^{-(aa \dots a)_b}$, where each $a = (b - 1)$

Non-normalized FPs allow us to get very very close to 0, consider $(.00 \dots 01)_b \cdot b^{-(aa \dots a)_b}$

$\mathbb{R}_b(t, s)$ denotes the set of all floating point numbers with t digit mantissa and s digit exponent.

And any $\mathbb{R}_b(t, s)$ is finite, while \mathbb{R} is infinite.

Overflow or **Underflow** occurs whenever a nonzero floating point number with absolute value outside these ranges must be stored on the computer.

Note that when the number is too close to zero it's called underflow. not when its largely negative.

Furthermore, \mathbb{R} is compact, where as $\mathbb{R}_b(t, s)$ is not. This is to say that between any two real numbers, there is an infinite number of reals in between. (infinite density).

A real number $x = \pm(x_I x_F)_b = \pm(d_k d_{k-1} \dots d_0 . d_{-1} d_{-2} \dots)_b$

We say we convert a real x to a fp number with the following notation, $x \in \mathbb{R} \rightarrow Fl(x) \in \mathbb{R}_b(t, s)$

We can represent this in $\mathbb{R}_b(t, s)$ by the following algorithm:

1. Normalizing the mantissa.

We shift the decimal point, and readjust with the exponent

$$x = \pm(d_k d_{k-1} \dots d_0 . d_{-1} d_{-2} \dots)_b = \pm(.D_1 D_2 \dots)_b \cdot b^{k+1}$$

2. chop or round the mantissa

(a) chopping - chop after digit t of the mantissa.

(b) rounding - chop after digit t then round D_t depending on whether $D_{t+1} \geq b/2$ or $D_{t+1} < b/2$.

A more efficient way to round is to add $b/2$ to D_{t+1} , then simply chop.

Example

Consider $Fl(3/2) \in \mathbb{R}_{10}(2, 4) = \begin{cases} +0.66 & \text{if chop} \\ +0.67 & \text{if round} \end{cases}$

Round off Error

Precisely, this is the difference between an $x \in \mathbb{R}$ and $FL(x) \in \mathbb{R}_b(t, s)$.

Typically measured relative to x as $\frac{x - FL(x)}{x} = \delta$ or $FL(x) = x(1 - \delta)$ where δ is the relative round off.

δ can be bound independently of x , $\delta < b^{1-t}$ for chopping normalized FPs.

$|\delta| < \frac{1}{2}b^{1-t}$ for rounding normalized FPs.

2.4 Machine Arithmetic

Let $x, y \in \mathbb{R}$ and $FL(x), FL(y) \in \mathbb{R}_b(t, s)$

Consider operations $\circ \in \{+, -, \times, /\}$, the computer gives $x \circ y \approx FL(FL(x) \circ FL(y))$

Example in $\mathbb{R}_{10}(2, 4)$

let $x = 2, y = 0.0000058 \in \mathbb{R}$ we have $x + y = 2.0000058$

then $FL(x) = (+0.20 \cdot 10^1)_{10}$ and $FL(y) = (+0.58 \cdot 10^{-5})_{10}$ Note that no RROs in $FL(x), FL(y)$

then $FL(x) + FL(y) = 0.20 \cdot 10^1 + 0.58 \cdot 10^{-5} = 0.20000058 \cdot 10^1$ stored temporarily.

finally, $FL(FL(x) + FL(y)) = (+0.20 \cdot 10^1)_{10}$

2.5 Machine Precision

Definition: The smallest **non-normalized** FP number eps such that $1 + eps > 1$. This number (eps) is referred to as machine epsilon.

This is important as this number $eps = \begin{cases} b^{1-t} & \text{chopping} \\ \frac{1}{2}b^{1-t} & \text{rounding} \end{cases}$

recall $\delta = \frac{x - FL(x)}{x}$ is the **RRO** where δ can be bound independently.

$0 \leq \delta \leq eps$ for chopping, and $|\delta| \leq eps$ for rounding.

How exactly is error propagated in each of $\circ \in \{+, -, \times, /\}$?

Let $x, y \in \mathbb{R}$, then $Fl(x) = x(1 - \delta)$, $Fl(y) = y(1 - \delta)$

Multiplication: $x \cdot y$, where computer gives $FL(FL(x) \cdot FL(y))$

We have $[x(x - \delta_x) \cdot y(1 - \delta_y)](1 - \delta_{xy}) = x \cdot y(1 - \delta_x)(1 - \delta_y)(1 - \delta_{xy}) \approx x \cdot y(1 - \delta_x - \delta_y - \delta_{xy}) = x \cdot y(1 - \delta)$

We say this is a good approximation as each δ is a small fraction bound by machine epsilon, there's a good chance that products of δ s will be so small that they are lower than such machine epsilon.

Where $|\delta| \leq 3 \cdot eps$, this is the worst case scenario.

Addition: $x + y$, where computer gives $FL(FL(x) + FL(y))$

We have $(x(1-\delta_x) + y(1-\delta_y))(1-\delta_{xy}) = x(1-\delta_x)(1-\delta_{xy}) + y(1-\delta_y)(1-\delta_{xy}) \approx x(1-\delta_x-\delta_{xy}) + y(1-\delta_y-\delta_{xy})$

$$= (x + y) \left[1 - \frac{x(1-\delta_x-\delta_{xy})}{x+y} + \frac{y(1-\delta_y-\delta_{xy})}{x+y} \right] = (x + y)[1 - \delta_+]$$

$$\text{Where } |\delta_+| \leq \left| \frac{x}{x+y} \right| 2 \cdot eps + \left| \frac{y}{x+y} \right| 2 \cdot eps = \frac{|x| + |y|}{|x+y|} 2 \cdot eps$$

$$\text{i.e. } |\delta_+| \leq \begin{cases} 2 \cdot eps & \text{x and y are the same sign} \\ 2 \cdot eps \frac{|x-y|}{|x+y|} & \text{x and y are different signs} \end{cases}$$

We see that when $x \approx -y$, the error bound could approach infinity. This is called subtractive cancellation.

Example of Subtractive Cancellation in $\mathbb{R}_{10}(3, 1)$ with rounding.

Compute $a^2 - 2ab + b^2$ with $a = 15.6, b = 15.7 \in \mathbb{R}$

first we have $fl(a) = +(0.156 \cdot 10^2)$ and $fl(b) = +(0.157 \cdot 10^2)$ with no initial round off.

In $\mathbb{R}_{10}(3, 1)$, $fl(a^2) \rightarrow fl(243.36) = +(0.243 \cdot 10^3)$

$fl(2ab) \rightarrow fl(489.84) = +(0.490 \cdot 10^3)$

$fl(b^2) \rightarrow fl(246.49) = +(0.246 \cdot 10^3)$

$fl(a^2 - 2ab + b^2) \rightarrow fl(243 - 490 + 246) = -(0.100 \cdot 10^1)$

But this is a problem as $a^2 - 2ab + b^2 = (a - b)^2$ which is positive.

3 Linear Systems

$A\vec{x} = \vec{b}$, $A \in \mathbb{R}^{n \times n}$, $\vec{x}, \vec{b} \in \mathbb{R}^n$, Given A, \vec{b} , calculate \vec{x}

General Solution Technique:

1. Reduce the problem to an equivalent one that is easier to solve.
2. Solve the reduced problem.

$$\text{Ex, given } \begin{cases} 3x_1 - 5x_2 = 1 \\ 6x_1 - 7x_2 = 5 \end{cases} \Leftrightarrow \begin{bmatrix} 3 & -5 \\ 6 & -7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \end{bmatrix} \Leftrightarrow \begin{bmatrix} 3 & -5 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \Rightarrow x_1 = 1, x_2 = 2$$

Generalize this to n equations in n unknowns.

Let matrix $A = [a_{ij}]$ where a_{11} refers to the top left element.

We have:

Equation 1: $a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$

Equation 2: $a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$

\vdots

Equation n: $a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n$

Reduce system to triangular form.

1. Assume that $a_{11} \neq 0$

multiply Equation 1 by $\frac{a_{21}}{a_{11}}$ and subtract that from Equation 2.

multiply Equation 1 by $\frac{a_{31}}{a_{11}}$ and subtract that from Equation 3.

\vdots

multiply Equation 1 by $\frac{a_{n1}}{a_{11}}$ and subtract that from Equation n.

This results in a equivalent system, where x_1 has been eliminated from Equations 2 to n.

Now we have:

Equation 1: $a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$

Equation $\hat{2}$: $0 + \hat{a}_{22}x_2 + \cdots + \hat{a}_{2n}x_n = \hat{b}_2$

\vdots

Equation \hat{n} : $0 + \hat{a}_{n2}x_2 + \cdots + \hat{a}_{nn}x_n = \hat{b}_n$

2. Assume that $\hat{a}_{22} \neq 0$

multiply Equation $\hat{2}$ by $\frac{\hat{a}_{32}}{\hat{a}_{22}}$ and subtract that from Equation $\hat{3}$.

multiply Equation $\hat{2}$ by $\frac{\hat{a}_{n2}}{\hat{a}_{22}}$ and subtract that from Equation \hat{n} .

\vdots

finally, at step $n - 1$, we assume $\tilde{a}_{n-1,n-1} \neq 0$, multiply equation $n - 1$ by $\frac{\tilde{a}_{n,n-1}}{\tilde{a}_{n-1,n-1}}$ and subtract from

Equation \tilde{n}

this will result in an upper triangular system.

Equation 1: $a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$

Equation $\hat{2}$: $0 + \hat{a}_{22}x_2 + \cdots + \hat{a}_{2n}x_n = \hat{b}_2$

\vdots

Equation \tilde{n} : $0 + 0 + \cdots + \tilde{a}_{nn}x_n = \tilde{b}_n$

Another way of looking at this problem is using matrix vector notation, looking to solve $A\vec{x} = \vec{b}$ where $A \in \mathbb{R}^{n \times n}$, $\vec{b}, \vec{x} \in \mathbb{R}^n$

Triangulating such A requires the following steps:

1. Eliminate the first column of A below a_{11}

$$L_1 A \vec{x} = L_1 \vec{b}, \text{ where } L_1 = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ -a_{21}/a_{11} & 1 & \cdots & 0 \\ -a_{31}/a_{11} & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -a_{n1}/a_{11} & 0 & \cdots & 1 \end{bmatrix}$$

This is very similar to the identity matrix, except the first column is filled with the multipliers used in the first step.

2. Eliminate the second column of $L_1 A$ below \hat{a}_{22}

$$L_2(L_1 A) \vec{x} = L_2(L_1 \vec{b}) = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & \hat{a}_{22} & \cdots & \hat{a}_{2n} \\ 0 & 0 & \hat{\hat{a}}_{33} & \hat{\hat{a}}_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \hat{\hat{a}}_{n3} & \hat{\hat{a}}_{nn} \end{bmatrix}, \text{ where } L_2 = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & -\hat{a}_{32}/\hat{a}_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & -\hat{a}_{n2}/\hat{a}_{22} & \cdots & 1 \end{bmatrix}$$

\vdots

We continue until we have $L_{n-1}L_{n-2} \cdots L_1 A \vec{x} = L_{n-1}L_{n-2} \cdots L_1 \vec{b}$

let $L_{n-1}L_{n-2} \cdots L_1 A = U$, where U is an upper triangular matrix. This becomes very easy to solve.

3.1 LU Factorization

We have $L_{n-1}L_{n-2}\dots L_1A = U \Leftrightarrow A = L_1^{-1}L_2^{-1}\dots L_{n-1}^{-1}U$

Lemma 1: If L_i is a Gauss Transform. Then L_i^{-1} exists and is also a Gauss Transform.

Lemma 2: If L_i and L_j are Gauss Transforms, and $i < j$, $L_iL_j = L_i + L_j - I$

Then $A = L_1^{-1}L_2^{-1}\dots L_{n-1}^{-1}U = LU$

Using $A = LU$ to solve $A\vec{x} = \vec{b}$

$$A\vec{x} = \vec{b} \Leftrightarrow LU\vec{x} = \vec{b}, \text{ let } \vec{d} = \vec{b} \text{ for a lower triangular } L\vec{d}$$

$$\text{Then } L\vec{d} = \vec{b} \Leftrightarrow U\vec{x} = \vec{d} \text{ for an upper triangular } U\vec{x}$$

We use LU factorization because its far cheaper, and if we have several linear systems with the same coefficient matrix A , but different right hand side, i.e. $A\vec{x} = \vec{b} = \vec{c} = \vec{d} \dots$, then we can use the same LU .

Example of LU factorization

$$\text{let } A = \begin{bmatrix} 2 & -1 & 1 \\ 2 & 2 & 2 \\ -2 & 4 & 1 \end{bmatrix}$$

$$L_1 = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, L_1A = \begin{bmatrix} 2 & -1 & 1 \\ 0 & 3 & 1 \\ 0 & 3 & 2 \end{bmatrix}$$

$$L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix}, L_2L_1A = \begin{bmatrix} 2 & -1 & 1 \\ 0 & 3 & 1 \\ 0 & 0 & 1 \end{bmatrix}, L_2L_1A = U \Leftrightarrow A = L_1^{-1}L_2^{-1}U$$

$$L_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}, L_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix}, L = L_1^{-1}L_2^{-1} = L_1^{-1} + L_2^{-1} - I = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ -1 & -1 & 1 \end{bmatrix}$$

3.2 Gaussian Elimination with Pivoting

From last lecture, we have $L_2P_2L_1P_1A = U$

$L_2P_2L_1P_1A \equiv L_2P_2L_1P_2P_2P_1A$, as $P_2P_2 \equiv I$

Then $L_2P_2L_1P_2P_2P_1A \equiv L_2(P_2L_1P_2)P_2P_1A$

Let $\tilde{L}_1 = P_2L_1P_2$

Claim: \tilde{L}_1 is a Gauss Transform.

$$\text{Example } P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, L_1 = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & 0 & 1 \end{bmatrix}, \text{ then } P_2L_1P_2 = \begin{bmatrix} 1 & 0 & 0 \\ l_{31} & 1 & 0 \\ l_{21} & 0 & 1 \end{bmatrix} \text{ (swapped multipliers)}$$

$$\text{Then } L_2\tilde{L}_1P_2P_1A = U \Leftrightarrow P_2P_1A = \tilde{L}_1^{-1}L_2^{-1}U \Leftrightarrow PA = LU$$

Now, to solve $A\vec{x} = \vec{b}$, given $PA = LU$

$$A\vec{x} = \vec{b} \Leftrightarrow PA\vec{x} = P\vec{b} \Leftrightarrow LU\vec{x} = P\vec{b}, \text{ note that } P\vec{b} \text{ is } \vec{b} \text{ with 2 elements swapped.}$$

Then $LU\vec{x} = \vec{b}$, let $\vec{d} = U\vec{x}$

$L\vec{d} = \vec{b}$ is easy to solve because L is lower triangular

$U\vec{x} = \vec{d}$ is easy to solve because U is upper triangular.

What happens if at some k -th stage, the diagonal and everything below it is 0?

For example $L_{k-1}\dots L_2L_1A$, we cannot divide by the k -th row, k -th column element because its 0, furthermore every element below it is 0.

Remember our goal originally is to make every element in the k -th column under k zero, we just continue.

This will result in a matrix U with a 0 along one of the diagonal entries.

Then we have a singular matrix U , but the factorization still stays.

While its possible for U to be singular, it's not for the matrix L .

As a combination of gauss transforms, they must be invertible and thus non-singular.

In the last step, we have $U\vec{x} = \vec{d}$, but U is singular. We could have no solution, or infinitely many solutions.

Example: $U\vec{x} = \vec{d} \rightarrow \begin{bmatrix} 2 & 5 & 4 \\ 0 & 0 & 1 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \Rightarrow 2x_3 = b_3 \text{ and } x_2 = b_2$

Now if $b_3 = 2b_2$, we have a free parameter x_2 . Otherwise, there are no solutions.

Now suppose during a numerical factorization of FPS, if all elements below the diagonal in column k and $[a_{kk}]$ at the k -th stage and are of magnitude $\leq \text{eps} \cdot \max_{j \in [1, k]} |u_{jj}|$, we call this **Numerical Singularity** (Near Singularity).

3.3 Complexity of Gaussian Elimination

Let matrix A be of order $n \times n$

We will count multiplication/addition pairs, i.e. $mx + b$ as a **Flop** (Floating Point Operation).

Lets begin with factorization. Computing the LU factorization,

1st stage is the zeroing out the first column, we have $(n-1)^2$ flops as we are adding multiples of row 1 to rows 2 to n .

2nd stage is the zeroing out the first column, we have $(n-2)^2$ flops as we are adding multiples of row 2 to rows 3 to n .

\vdots

$(n-1)$ -th stage is the zeroing out the first column, we have 1 flop.

Finally, we have $(n-1)^2 + (n-2)^2 + \dots + 1$ total flops.

$$\sum_{k=1}^n k^2 = \frac{n(n+1)(2n+1)}{6}, \text{ then } (n-1)^2 + (n-2)^2 + \dots + 1 = \frac{n(n-1)(2n)}{6} = \frac{n^3}{3} + O(n^2)$$

Computing the forward solve, $L\vec{d} = \vec{b}$ where L is a lower triangular matrix.

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ l_{21} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ l_{n1} & l_{n2} & \dots & 1 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \Leftrightarrow d_1 = b_1, d_2 = b_2 - l_{21}d_1, d_3 = b_3 - l_{31}d_1 - l_{32}d_2 \dots$$

$$\text{Total} = 0 + 1 + 2 + \dots + (n-1) = \frac{n^2}{2} + O(n) \text{ Flops}$$

$$\text{Back-solving is similar, resulting in another } \frac{n^2}{2} + O(n) \text{ Flops}$$

$$\text{Total Forward/Backward solve: } n^2 + O(n) \text{ Flops}$$

3.4 Round Off Error of Gaussian Elimination

Recall we have a factorization $PA = LU$ computed in a floating point system.

Because of machine round off error, we actually get $\hat{P}(A+E) = \hat{L}\hat{U}$, where $\hat{P}, \hat{L}, \hat{U}$ are the computed factors.

Then solving $A\vec{x} = \vec{b}$ becomes solving $(A+E)\vec{\hat{x}} = \vec{b}$, where $\vec{\hat{x}}$ is the computed solution.

Equivalently, let $E\vec{\hat{x}} = \vec{r}$, then $(A+E)\vec{\hat{x}} = \vec{b}$ becomes $A\vec{\hat{x}} + \vec{r} = \vec{b}$, or $\vec{r} = \vec{b} - A\vec{\hat{x}}$

Where \vec{r} is the residual. (We would like this to be $\vec{0}$).

We can show that if row partial pivoting is used,

$\|E\| \lesssim k \cdot \text{eps} \cdot \|A\|$ where k is not too large and depends on n .

And that $\|\vec{r}\| \lesssim k \cdot \text{eps} \cdot \|b\| \Leftrightarrow \frac{\|\vec{r}\|}{\|b\|} \lesssim k \cdot \text{eps}$

However, this does not mean that $\|\vec{x} - \hat{x}\|$ or $\frac{\|\vec{x} - \hat{x}\|}{\|x\|}$ is small.

Example: $\begin{bmatrix} 0.780 & 0.563 \\ 0.913 & 0.656 \end{bmatrix} \vec{x} = \begin{bmatrix} 0.217 \\ 0.254 \end{bmatrix}$, we know true solution is $\vec{x} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

Consider two computed solutions, $\hat{x}_1 = \begin{bmatrix} 0.999 \\ -1.001 \end{bmatrix}$, $\hat{x}_2 = \begin{bmatrix} 0.341 \\ -0.087 \end{bmatrix}$

$\vec{r}_1 = \vec{b} - A\hat{x}_1 = \begin{bmatrix} -0.001243 \\ -0.001572 \end{bmatrix}$, $\vec{r}_2 = \vec{b} - A\hat{x}_2 = \begin{bmatrix} -0.000001 \\ 0 \end{bmatrix}$

We see that $\frac{\|\vec{r}_2\|}{\|\vec{b}\|}$ is much smaller than $\frac{\|\vec{r}_1\|}{\|\vec{b}\|}$, yet we see that \hat{x}_2 is a terrible solution.

But why is $\frac{\|\vec{x} - \hat{x}_1\|}{\|x\|}$ so much smaller than $\frac{\|\vec{x} - \hat{x}_2\|}{\|x\|}$?

We need the relationship between relative error and relative residual.

We have $A\hat{x} = \vec{b} - \vec{r}$ and $A\vec{x} = \vec{b}$, subtract the top from the bottom yields:

$A(\vec{x} - \hat{x}) = \vec{r} \Leftrightarrow \vec{x} - \hat{x} = A^{-1}\vec{r}$

Taking the norm of both sides yields: $\|\vec{x} - \hat{x}\| = \|A^{-1}\vec{r}\| \leq \|A^{-1}\| \|\vec{r}\|$ (1)

We can take the original $\vec{b} = A\vec{x} \Leftrightarrow \|\vec{b}\| = \|A\vec{x}\| \leq \|A\| \|\vec{x}\|$ (2)

Combining (1) and (2): $\frac{\|\vec{x} - \hat{x}\|}{\|A\| \|\vec{x}\|} \leq \frac{\|A^{-1}\| \|\vec{r}\|}{\|\vec{b}\|} \Leftrightarrow \frac{\|\vec{x} - \hat{x}\|}{\|\vec{x}\|} \leq \frac{\|A\| \|A^{-1}\| \|\vec{r}\|}{\|\vec{b}\|}$

Where $\|A\| \|A^{-1}\| = \text{cond}(A)$

We can try to derive a lower bound for this, again we have equations $A\vec{x} = \vec{b}$ and $A\hat{x} = \vec{b} - \vec{r}$,

Subtracting the two yields $A(\vec{x} - \hat{x}) = \vec{r} \Leftrightarrow \|A(\vec{x} - \hat{x})\| = \|\vec{r}\|$

By triangular inequality, $\|A\| \|\vec{x} - \hat{x}\| \geq \|\vec{r}\|$ (1)

We take the original $A\vec{x} = \vec{b}$, rearrange: $\vec{x} = A^{-1}\vec{b} \Rightarrow \vec{x} \leq \|A^{-1}\| \|\vec{b}\|$ (2)

Combining equations (1) and (2) yields $\frac{\|\vec{x} - \hat{x}\|}{\|\vec{x}\|} \geq \frac{\|\vec{r}\|}{\|\vec{b}\| \|A\| \|A^{-1}\|}$

Finally, combining two bounds, we have $\frac{\|\vec{r}\|}{\|\vec{b}\| \text{cond}(A)} \leq \frac{\|\vec{x} - \hat{x}\|}{\|\vec{x}\|} \leq \text{cond}(A) \frac{\|\vec{r}\|}{\|\vec{b}\|}$

If $\text{cond}(A)$ is very large, problem is poorly conditioned. Small relative residual does not mean small relative error.

If the $\text{cond}(A)$ is not too large, the problem is well conditioned and the small relative residual is a reliable indicator of small relative error.

Conditioning is a continuous spectrum, how large is "very large" depends on context.

Recall in the previous example: $A = \begin{bmatrix} 0.780 & 0.563 \\ 0.913 & 0.656 \end{bmatrix}$, $A^{-1} = 10^6 \begin{bmatrix} 0.656 & -0.565 \\ -0.913 & 0.780 \end{bmatrix}$, where $10^6 = \frac{1}{\det(A)}$

Determinant formula: $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, $A^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ where $\det(A) = 0 \Rightarrow A^{-1}$ DNE.

$$\|A\|_\infty = 1.572, \|A^{-1}\| = 1.693 \cdot 10^6 \Rightarrow \text{cond}_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty = 2.66 \cdot 10^6$$

$$\frac{\|\vec{x} - \hat{x}\|}{\|\vec{x}\|} \leq 2.66 \cdot 10^6 \frac{\|r\|}{\|b\|}, \text{ i.e. relative error in } x \text{ could be as big as } 2.66 \cdot 10^6 \text{ times the relative residual.}$$

Thus this A is a poorly conditioned matrix. and relative residual is not a reliable indicator of relative error.

3.5 Iterative Refinement/Iterative Improvement

Can we improve \hat{x} ? One easy way is to improve mantissa length.

Want to solve $A\vec{x} = \vec{b}$, having already solved $(A + E)\hat{x} = \vec{b}$

Subtracting the two equations yields $A(\vec{x} - \hat{x}) = \vec{r}$, let $\vec{z} = \vec{x} - \hat{x}$, and we solve $A\vec{z} = \vec{r}$

Then $\vec{x} = \hat{x} + \vec{z}$, however this is a fallacy since we can not solve for \vec{z} exactly. thus we get \hat{z} , and not \vec{z} .

But if we were to get 2 leading digits correct in \hat{x} , then we will get 2 leading digits correct in \hat{z} , and then you will have 4 leading digits correct in $\hat{x} + \hat{z}$

The Algorithm

Compute $\hat{x}^{(0)}$ by solving $A\vec{x} = \vec{b}$ in a floating point system.

For $i = 0, 1, 2, \dots$ until solution is good enough.

 Compute the residual: $r^{(i)} = b - a\hat{x}^{(i)}$

 Solve $A\vec{z}^{(i)} = r^{(i)}$ For some $\hat{z}^{(i)}$

 Update $\hat{x}^{(i+1)} = \hat{x}^{(i)} + \hat{z}^{(i)}$

4 Non Linear Equations

Problem: $F : \mathbb{R} \rightarrow \mathbb{R}$

We need to solve $F(\tilde{x}) = 0$ for some $\tilde{x} \in \mathbb{R}$ (root finding problem)

Typically, we cannot find a "closed form" (analytic) solution to non linear F , which is why we need numerical techniques.

We can find iterative methods which generate an approximation $\hat{x}_k, k = 0, 1, 2, \dots$ such that $\hat{x}_k \rightarrow \tilde{x}$ as $k \rightarrow \infty$

4.1 Fixed Point Methods

Given a root finding problem, $F(\tilde{x}) = 0$, this is equivalent to a fixed point problem $\tilde{x} = g(\tilde{x})$

where $\tilde{x} = g(\tilde{x})$ is obtaining a fixed point from g .

Example: $F(x) = x - e^{-x}$, this is equivalent to finding the fixed point of $x = e^{-x}$

We can always use $g(x) = x - F(x)$ to obtain a working $g(x)$. Alternatively, we can use $g(x) = x - h(x)F(x)$ we refer to $x - F(x)$ to be the first form, $x - h(x)F(x)$ as the second form.

We often use the algebraically equivalent fixed point problem to solve the root finding problem because there is an obvious iteration algorithm to do so.

First form, $F(\tilde{x}) = 0 \Leftrightarrow \tilde{x} = g(\tilde{x})$

Second form, $F(\tilde{x}) = 0 \Rightarrow \tilde{x} = g(\tilde{x})$, however we can have a fixed point $\tilde{x} = g(\tilde{x})$ such that $F(\tilde{x}) \neq 0$.
(i.e., $h(\tilde{x}) = 0$ but $F(\tilde{x}) \neq 0$)

The advantage of the second form is that there is flexibility in designing $g(x)$, to make iteration converge faster.

4.2 Fixed Point Iteration

Start with an approximate solution \hat{x}_0 then iterate:

$$\hat{x}_{k+1} = g(\hat{x}_k), k = 0, 1, 2, \dots \text{ until convergence/failure.}$$

Example: Let $F(x) = x^2 + 2x - 3$ with exact roots $\tilde{x} = 1, -3$

Consider the fixed point iteration (FPI): $x_{k+1} = x_k + \frac{x_k^2 + 2x_k - 3}{x_k^2 - 5}$

for the fixed point problem $x = g(x) = x + \frac{x^2 + 2x - 3}{x^2 - 5}$

We see that this is the second form of the fixed point problem, where $h(x) = -\frac{1}{x^2 - 5}$
notice using this $h(x) \neq 0, \forall x \in \mathbb{R}$, this means we don't need to check after convergence.

Fixed Point Iteration of $\hat{x}_0 = 5 \rightarrow \hat{x}_k$ s converges to -3.

However, Fixed Point Iteration of $\hat{x}_0 = +5 \rightarrow \hat{x}_k$ s do not converge.

We can keep trying, but that would be a significant waste of time.

4.3 Fixed Point Theorem

If there is an interval $[a, b]$ such that $g(x) \in [a, b]$ for all $x \in [a, b]$

And if $\|g'(x)\| \leq L < 1, \forall x \in [a, b]$

Then $g(x)$ has a unique fixed point in $[a, b]$

Proof, start with any $\hat{x}_0 \in [a, b]$ and iterate with $\hat{x}_{k+1} = g(\hat{x}_k), k = 0, 1, 2, \dots$ then all $\hat{x}_k \in [a, b]$

Moreover, $x_{k+1} - x_k = g(x_k) - g(x_{k-1}) = g'(\eta_k)g(x_k - x_{k-1})$

For some $\eta_k \in [x_{k-1}, x_k] \subset [a, b]$ By MVT (Mean Value Theorem)

Thus, $|x_{k+1} - x_k| < |g'(\eta_k)||x_k - x_{k-1}| \leq L|x_k - x_{k-1}|$

Then, $L|x_k - x_{k-1}| \leq \dots \leq L^k|x_1 - x_0|$

Since $L < 1, |x_{k+1} - x_k| \rightarrow 0$ as $k \rightarrow \infty$

This means the fixed points are converging to some $\tilde{x} \in [a, b]$

to complete this proof, we must show:

1. $\tilde{x} = g(\tilde{x})$ (i.e., \tilde{x} is a fixed point!)
2. \tilde{x} is unique in $[a, b]$.

Using the Fixed Point Theorem, **Example:** $F(x) = x - e^x = 0$ and $g(x) = e^{-x}$ for $x = g(x)$

What is the maximum interval of guaranteed convergence of x_k ?

4.4 Rate of Convergence

Definition: If $\lim_{x_k \rightarrow \tilde{x}} \frac{|\tilde{x} - x_{k+1}|}{|\tilde{x} - x_k|^p} = c \neq 0$, we have **p-th** order convergence to fixed point \tilde{x}

Example: Table of absolute error of iterates, $|\tilde{x} - x_k|$

k	p = 1, c = 1/2	p = 2, c = 1
0	10^{-1}	10^{-1}
1	$5 \cdot 10^{-2}$	10^{-2}
2	$2.5 \cdot 10^{-2}$	10^{-4}
3	$1.25 \cdot 10^{-2}$	10^{-8}
4	$6.125 \cdot 10^{-3}$	10^{-16}

We see that despite the second column having $c = 1, p = 2$ converges much much faster.

4.5 Rate of Convergence Theorem

For the Fixed Point Iteration $x_{k+1} = g(x_k)$, if $g'(\hat{x}) = g''(\hat{x}) \cdots = g^{(p-1)}(\hat{x}) = 0$ but $g^{(p)}(\hat{x}) \neq 0$, then we have p -th order convergence.

Proof let $x_{k+1} = g(x_k) = g(\tilde{x} + (x_k - \tilde{x}))$

$$= g(\tilde{x}) + g'(x_k - \tilde{x})g'(\tilde{x}) + \frac{(x_k - \tilde{x})^2}{2!}g''(\tilde{x}) + \cdots + \frac{(x_k - \tilde{x})^{p-1}}{(p-1)!}g^{(p-1)}(\tilde{x}) + \frac{(x_k - \tilde{x})^p}{p!}g^{(p)}(\eta_k)$$

Where $\eta_k \in [\tilde{x}, x_k]$

However, if we have $g'(\hat{x}) = g''(\hat{x}) = \cdots = g^{(p-1)}(\hat{x}) = 0$

The Taylor expansion yields $x_{k+1} = g(x_k) = g(\hat{x}) + \frac{(x_k - \hat{x})^p}{p!}g^{(p)}(\eta_k) \Leftrightarrow \frac{x_{k+1} - \hat{x}}{(x_k - \hat{x})^p} = \frac{1}{p!}g^{(p)}(\eta_k)$

We see that as the $p \rightarrow \infty$, $x_k \rightarrow \hat{x}$, $\eta_k \in [\hat{x}, x_k] \Rightarrow \eta_k = \hat{x}$

We can rewrite this as $\lim_{x_k \rightarrow \hat{x}} \frac{|x_{k+1} - \hat{x}|}{|x_k - \hat{x}|^p} = \frac{1}{p!}g^{(p)}(\hat{x})$

Thus we see that when the p -th derivative of g at \hat{x} is not zero, we reach p -th order convergence.

We see that now we can use the second form of the fixed point iteration $g(x) = x - h(x)F(x)$, and pick such $h(x)$ so that the p -th derivative of g is not zero at \hat{x} to accelerate iteration.

4.6 Newton's Method

Recall Newton's Method of $F(x) = 0$ where $x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}$

Where it takes the root finding problem of $F(x) = 0 \Leftrightarrow x = g(x)$ with $g(x) = x - \frac{F(x)}{F'(x)}$

This is the second form of the fixed point problem with $h(x) = \frac{1}{F'(x)}$

Speed of Newton's Method

Suppose that $F'(\tilde{x}) = 0$ but $F''(\tilde{x}) \neq 0$

$$g'(x) = 1 - \frac{F'(x)F'(x) - F(x)F''(x)}{(F'(x))^2} = \frac{F(x)F''(x)}{(F'(x))^2} \text{ i.e., } g'(\tilde{x}) = 0 \text{ For any function } F.$$

By rate of convergence theorem, the newtons method has at-least quadratic convergence. since $g'(\tilde{x}) = 0$ for any function F .

Geometric Interpretation of Newton's Method

Want to solve $F(x) = 0$, at an initial guess x_k approximate (or model) $F(x)$ by a linear polynomial $p(x)$ that satisfies conditions: $p(x_k) = F(x_k)$ and $p'(x_k) = F'(x_k) \Rightarrow p_k(x) = F(x_k) + (x - x_k)F'(x_k)$

Then x_{k+1} is the root of $p_k(x)$, i.e. $p_k(x_{k+1}) = 0 \Rightarrow F(x_k) + (x_{k+1} - x_k)F'(x_k) = 0$

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

Will Newton's Method always Converge? No.

For example Newton's method on $F(x) = x - e^{-x}$

$$x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)} = x_k - \frac{x_k - e^{-x_k}}{1 + e^{-x_k}} = \frac{(1 + x_k)e^{-x_k}}{1 + e^{-x_k}}$$

4.7 Secant Method

Recall NM: $x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}$, notice that we have to supply the first derivative of F

We can approximate $F'(x_k)$ by $\frac{F(x_k) - F(x_{k-1})}{x_k - x_{k-1}}$

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

Another way to derive this is using the geometric interpretation.

However This is not a fixed point iteration of the form $x_{k+1} = g(x_k)$, as no function g could support two values.

We cannot directly use FPT or RCT to analyze this method. However, with some adjustment can prove rate of convergence is $p = \frac{1 + \sqrt{5}}{2} \approx 1.62$

Super Linear Convergence, perhaps not as fast as Newton ($p = 2$) ?

Newton's Method: $x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}$

Secant Method: $x_{k+1} = x_k - \frac{F(x_k)}{\frac{F(x_k) - F(x_{k-1})}{x_k - x_{k-1}}} = x_k - F(x_k) \frac{(x_k - x_{k-1})}{F(x_k) - F(x_{k-1})}$

But we see that newtons method requires 2 Function evaluations per step as F, F' are not the same.

And the secant method requires only 1 function evaluation. Although $F(x)$ needs to be evaluated, $F(x_{k-1})$ has been computed during the previous iteration!

Effect rate of convergence takes cost per iteration into account, as well as speed. Thus secant method is "effectively" faster than Newton's method.

4.8 Bisection Method

Oftentimes, if we start Newton's/Secant method with the wrong initial guess, they will not converge. (Wrong side of the hill)

Bisection Method however has guaranteed convergence, this is to say that it always finds a root, but slowly.

We need to find an $a \leq b$ such that $F(a) \leq 0 \leq F(b)$ or $F(b) \leq 0 \leq F(a)$

This means there is at least one root of F in $[a, b]$

Assume $F(a) \leq 0 \leq F(b)$

Loop until $b - a$ is small enough.

let $m = (a + b)/2$

If $F(m) \leq 0$ then let $a = m$, else $b = m$

Repeat for the interval $[m, b]$ or $[a, m]$

Then we have linear rate of convergence $p = 1, c = \frac{1}{2}$ with guaranteed convergence.

4.9 Hybrid Method

We can start off with bisection, then switch to a faster one later on.

Example: Combine Bisection and Newton's Method

4.10 System of Non-Linear Equations

Problem: Solve $F(\vec{x}) = \vec{0}$

Where $F(\vec{x}) = \begin{bmatrix} F_1(x_1, x_2) \\ F_2(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 4(x_1)^2 + 9(x_2)^2 - 16x_1 - 54x_2 + 61 \\ x_1x_2 - 2x_1 - 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Extending Newton's Method to Systems. Recall Newton's Method for a single non linear equation where

$F : \mathbb{R} \rightarrow \mathbb{R}$ and $x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}$

Then for $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ then $\vec{x}_{k+1} = \vec{x}_k - \frac{F(\vec{x}_k)}{F'(\vec{x}_k)}$

Need to Know: F' is the Jacobian matrix of F , i.e. $F' \in \mathbb{R}^{n \times n}$
 then $\vec{x}_{k+1} = \vec{x}_k - [F'(\vec{x}_k)]^{-1}F(\vec{x}_k)$ or $[F'(\vec{x}_k)](\vec{x}_{k+1} - \vec{x}_k) = -F(\vec{x}_k)$

This is simply in the form of some $A\vec{x} = \vec{b}$, a linear system!

This is very expensive! We can use the Pseudo Newton Method by holding the Jacobian matrix fixed for a few iterations, this means we can reuse the $PA = LU$ factorization. This is alright so long as we are still converging.

5 Approximation/Interpolation

Truncated Taylor Series: $p(x) = F(a) + F'(a)(x-a) + \dots + \frac{F^{(n)}(a)}{n!}(x-a)^n$

This is a polynomial because of the $(x-a)^i$ for $i = 0, 1, 2, \dots$

Then the error $e(x) = p(x) - F(x) = \frac{F^{(n+1)}(\eta)}{(n+1)!}(x-a)^{n+1}$

Other approximations:

1. Interpolation: Find a polynomial p such that $p(x_i) = F(x_i), i = 0, 1, 2, \dots$

F is the function we are trying to approximate, could simply be a set of data.

2. Least Squares: Finding a polynomial p such that $p(x)$ minimize $\|F - p\|_2 = \left(\int_a^b (F(x) - p(x))^2 dx \right)^{1/2}$

Other norms we can use for Least Squares approximation,

$$\|F - p\|_\infty = \max_{a \leq x \leq b} |F(x) - p(x)|$$

$$\|F - p\|_1 = \int_a^b |F(x) - p(x)| dx$$

5.1 Polynomial Interpolation

Consider \mathcal{P}_n : the set of all polynomials of degree $\leq n$

This is a function space, and requires a basis of $n+1$ functions. (\mathcal{P}_n is a function space of dimension $n+1$)

The most common basis is the monomial basis $\{x^i\}_{i=0}^n$

Weierstrass' Theorem (Existence)

If a function F is continuous on an interval $[a, b]$, then for any $\epsilon > 0, \exists p_\epsilon \ni \|F - p_\epsilon\|_\infty < \epsilon$

This is to say that for any continuous function on a closed interval $[a, b]$ there exists some polynomial that is as close to it as it can be.

Numerical Methods for Polynomial Interpolation

1. Vandermonde: Method of Undetermined Coefficients

Theorem For any sets: $\{x_i : i = 0, 1, \dots, n\}, \{y_i : i = 0, 1, \dots, n\}$ for distinct x_i s and undistinct y_i s

Then $\exists! p \in \mathcal{P}_n, \ni p(x_i) = y_i, i = 0, 1, \dots, n$

Proof: If such $p(x)$ exists, it must be possible to write it as a monomial space polynomial:

$$p(x) = \sum_{i=0}^n a_i x^i, \text{ Then we have } n+1 \text{ unknowns, } \{a_i\}, i = 0, 1, \dots, n$$

This can be converted into a matrix problem with $p(x_i) = y_i, i = 0, 1, \dots, n$

We can solve for the a_i s by using the Vandermonde matrix:

$$\begin{bmatrix} (x_0)^0 & (x_0)^1 & (x_0)^2 & \dots & (x_0)^n \\ (x_1)^0 & (x_1)^1 & (x_1)^2 & \dots & (x_1)^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ (x_n)^0 & (x_n)^1 & (x_n)^2 & \dots & (x_n)^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Now is the Vandermonde matrix non-singular?

We can return to polynomial problem, we are asking if its possible

for $p \in \mathcal{P}_n, \exists: p(x_i) = 0, i = 0, 1, 2, \dots, n$

But this is impossible as no polynomial in \mathcal{P}_n have $n + 1$ roots

Thus the polynomial is the x-axis with 0 coefficients.

$$\text{Or } \det(V) = (-1)^{n+1} \prod_{i=1}^n \prod_{j=0}^{i-1} (x_i - x_j)$$

The Vandermonde however does not lead to the best algorithm, it can be poorly conditioned, $PV = LU$ factorization then forward/backward solve will cost.

This will give the monomial basis $p(x)$

2. Lagrange Basis

Consider a simple interpolation problem $p(x_i) = y_i, i = 0, 1, \dots, n$

Consider the basis $\{l_i(x)\}$ where $l_i(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j}$ for $i = 0, 1, \dots, n$

$$\text{Note, } l_i(x) \in \mathcal{P}_n, \forall i \text{ and } l_i(x_j) = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

With this basis function, we can write out the interpolating polynomial for free.

$$\text{Where } p(x) = \sum_{i=0}^n l_i(x) y_i \Rightarrow p(x) \in \mathcal{P}_n \text{ and } p(x_i) = y_i, i = 0, 1, \dots, n$$

3. Newton (Divided Difference) Basis

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

We look for an interpolating $p(x)$ of the form:

$$p(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots + a_n(x - x_0)(x - x_1)(x - x_2) \dots (x - x_{n-1})$$

$$\text{Converting into a matrix, } \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & (x_1 - x_0) & 0 & \dots & 0 \\ 1 & (x_2 - x_0) & (x_2 - x_0)(x_2 - x_1) & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & (x_n - x_0) & \dots & \prod_{i=0}^{n-1} (x_n - x_i) & \dots \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

This is a lower triangular system, which means there is no factorization involved.

5.2 Divided Differences

Definition: $y[x_1] = y(x_1) = y_1$

$$y[x_{i+k}, \dots, x_i] = \frac{y[x_{i+k}, \dots, x_{i+1}] - y[x_{i+k+1}, \dots, x_i]}{x_{i+k} - x_i}$$

example: $y[x_2, x_1, x_0] = \frac{y[x_2, x_1] - y[x_1, x_0]}{x_2 - x_0}$

Newton's Polynomial: $p(x) = y[x_0] + (x - x_0)y[x_1, x_0] + \dots + (x - x_0)(x - x_1) \dots (x - x_{n-1})y[x_n, \dots, x_0]$
Then, $p(x) \in \mathcal{P}_n$ and $p(x_i) = y_i, i = 0, 1, \dots, n$

Proof by Induction, let $p(x) \equiv p_n(x)$

Basis: $n = 0, n = 1$

$$p_0(x) = y[x_0], p_0(x_0) = y[x_0] = y_0$$

$$p_1(x) = y[x_0] + (x_1 - x_0)y[x_1, x_0]$$

$$p_1(x_0) = y[x_0] = y_0$$

$$p_1(x_1) = y[x_0] + (x_1 - x_0) \left(\frac{y_1 - y_0}{x_1 - x_0} \right) = y[x_0] + y_1 - y_0 = y_1$$

Induction Hypothesis:

SKIPPING THIS FOR NOW NO TIME LMFAO

How are divided differences and derivatives related?

Example: $y[x_1, x_0] = \frac{y(x_1) - y(x_0)}{x_1 - x_0}$, but what happens when $x_1 \rightarrow x_0$

$$\text{Then } \lim_{x_1 \rightarrow x_0} y[x_1, x_0] = \lim_{x_1 \rightarrow x_0} \frac{y(x_1) - y(x_0)}{x_1 - x_0} = y'(x_0)$$

Provided the y' exists.

Example: $y[x_2, x_1, x_0] = \frac{y[x_2, x_1] - y[x_1, x_0]}{x_2 - x_0}$

$$\text{Then } \lim_{\substack{x_2 \rightarrow x_0 \\ x_1 \rightarrow x_0}} y[x_2, x_1, x_0] = \lim_{x_1 \rightarrow x_0} \frac{y[x_2, x_1] - y[x_1, x_0]}{x_2 - x_0} = \frac{y''(x_0)}{2!}$$

We can show that in general, $\lim_{\substack{x_k \rightarrow x_0 \\ x_{k-1} \rightarrow x_0 \\ \vdots \\ x_1 \rightarrow x_0}} y[x_k, x_{k-1}, \dots, x_0] = \frac{y^{(k)}(x_0)}{k!}$

This helps with Osculatory Interpolation (i.e., interpolation with derivatives).

Example: Find $p \in \mathcal{P}_4$ such that.

$$p(0) = 0$$

$$p(1) = 1, p'(1) = 1, p''(1) = 2$$

$$p(2) = 6$$

6 Numerical Integration

Problem: Approximate some $I(F) = \int_a^b F(x)dx$

recall Riemann Sums of a function F , $R(F) = \sum_{i=0}^{u-1} F(\eta_i)(x_{i+1} - x_i)$

Composite Rule: Apply a simple format on many short integrals, i.e. $\int_{x_1}^{x_2} F(x)dx$

Basic Formulas: Based off interpolation, instead of integrating we can use the interpolated polynomial. We approximate $F(x)$ by some $p(x)$ on $[a, b]$

Then $I(F) \approx Q(F) = I(P) = \int_a^b p(x)dx = \int_a^b [\sum_{i=0}^n F(x_i)l_i(x)] = \sum_{i=0}^n F(x_i) \int_a^b l_i(x)dx$

We arrive at the general form of quadrature: $Q(F) = \sum_{i=0}^n A_i F(x_i)$, $A_i = \int_a^b l_i(x)dx$

6.1 Theorem: Existence of Interpolatory Quadrature

Given any set of $n + 1$ distinct nodes, we can choose the weights A_i such that the quadrature rule is exact for all polynomials of degree $\leq n$, moreover A_i s are unique.

Proof: Since $Q(F)$ is linear, then for $\{x_i\}_{i=0}^n$ the polynomials can be divided into basis form.

Then for each x^k , $Q(x^k) = \sum_{i=0}^n A_i (x_i)^k = \frac{1}{k+1} (b^{k+1} - a^{k+1})$

This results in $n + 1$ equations with $n + 1$ unknowns.

$$\begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_0 & x_1 & x_2 & \dots & x_n \\ (x_0)^2 & (x_1)^2 & (x_2)^2 & \dots & (x_n)^2 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ (x_0)^n & (x_1)^n & (x_2)^n & \dots & (x_n)^n \end{bmatrix} \begin{bmatrix} A_0 \\ A_1 \\ A_2 \\ \vdots \\ A_n \end{bmatrix} = \begin{bmatrix} b-a \\ (b^2 - a^2)/2 \\ (b^3 - a^3)/3 \\ \vdots \\ (b^{n+1} - a^{n+1})/(n+1) \end{bmatrix}$$

We see that the first matrix is simply the transposed Vandermonde Matrix, thus by invertability, weights are unique.

6.2 Definition of Precision

If m is the highest natural number such that Q integrates all polynomials of degree $\leq m$, we say the precision of Q is m .

Midpoint Rule: $I(F) \approx Q(F) = M(F) = I(p_0) = (b-a)F(\frac{a+b}{2})$

Trapezoid Rule: $I(F) \approx Q(F) = T(F) = I(p_1) = \frac{b-a}{2}[F(a) + F(b)]$

Simpsons Rule: $I(F) \approx Q(F) = T(F) = I(p_2) = \frac{b-a}{6}[F(a) + F(\frac{b+a}{2}) + F(b)]$

6.3 Error in Interpolatory Quadrature

$E_n = I(F) - Q_n(F)$, where Q_n will integrate all polynomials of degree n

$$= \int_a^b [F(x) - p_n(x)] dx = \int_a^b \left(\frac{F^{(n+1)}(\eta)}{(n+1)!} \prod_{i=0}^n (x - x_i) \right) dx$$

Will interpolatory quadrature rules converge to $I(F)$ as $n \rightarrow \infty$?

Theorem: Let F be continuous on $[a, b]$ and $Q_n(F) = \sum_{i=0}^n A_i^{(n)} F(x_i)$

Then: $\lim_{n \rightarrow \infty} Q_n(F) = I(F) \iff \exists k \in \mathbb{R}, \exists: \sum_i |A_i^{(n)}| \leq k, \forall n \in \mathbb{N}$