

$$E = \frac{f^{(2n)}(\xi)}{(2n)!} \int_a^b \left( \prod_{i=0}^{n-1} (x - x_i) \right)^2 w(x) dx, \quad (108)$$

where  $\xi \in (a, b)$ .

## § 4.5 Romberg Integration

Idea: use Richardson's extrapolation to increase accuracy of composite trapezoidal rule.  
 → get high accuracy method from multiple applications of a low accuracy method.

Recall composite trapezoidal rule:

$$\int_a^b f(x) dx = \frac{h}{2} \left( f(x_0) + f(x_m) + 2 \sum_{j=1}^{m-1} f(x_j) \right) - \frac{(b-a)}{12} h^2 f''(\xi)$$

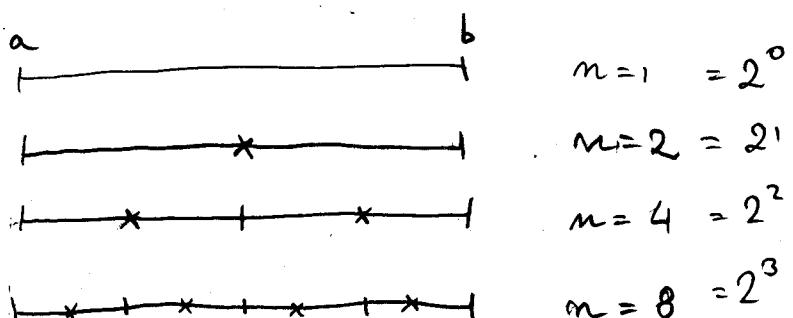
where  $\xi \in (a, b)$   $h = \frac{b-a}{m}$  and  $x_j = a + jh$ .

Romberg integration needs to evaluate composite trapezoidal rule

for  $n = 2^k$ ,

$k = 0, 1, \dots, M$

without superfluous function evaluations



\* = new function evaluations  
 | = function evaluation available from previous "level".

(109)

Thus composite trapz rule becomes with  $h_k = \frac{2^{-k}}{2^k-1} (b-a)$ :

$$\int_a^b f(x) dx = \frac{h_k}{2} \left[ f(a) + f(b) + 2 \sum_{i=1}^{2^k-1} f(a + i h_k) \right] - \frac{(b-a)}{12} h_k^2 f''(\xi_k)$$

$$= R_{k,0} + \text{error}$$

where  $\xi_k \in (a, b)$ .

So:  $R_{0,0} = \frac{h_0}{2} (f(a) + f(b)) = \frac{b-a}{2} (f(a) + f(b))$

$$R_{1,0} = \frac{h_1}{2} (f(a) + f(b) + 2f(a+h_1))$$

$$= \frac{1}{2} R_{0,0} + h_1 f(a+h_1)$$

$$R_{2,0} = \frac{1}{2} R_{1,0} + h_2 (f(a+h_2) + f(a+3h_2))$$

$$R_{k,0} = \frac{1}{2} R_{k-1,0} + h_k \sum_{j=1}^{2^{k-1}} f(a + (2j-1)h_k)$$

evaluation only at odd numbered nodes.

It can be shown that:

$$\textcircled{1} \quad \int_a^b f(x) dx = R_{k,0} + K_1 h_k^2 + K_2 h_k^4 + K_3 h_k^6 + \dots$$

$$\textcircled{2} \quad \int_a^b f(x) dx = R_{k+1,0} + K_1 h_{k+1}^2 + K_2 h_{k+1}^4 + K_3 h_{k+1}^6 + \dots$$

$$= R_{k+1,0} + K_1 \frac{h_k^2}{4} + K_2 \frac{h_k^4}{16} + K_3 \frac{h_k^6}{64} + \dots$$

Do Richardson's extrapolation trick to cancel out leading term in error:

$$4 \times ② - ① :$$

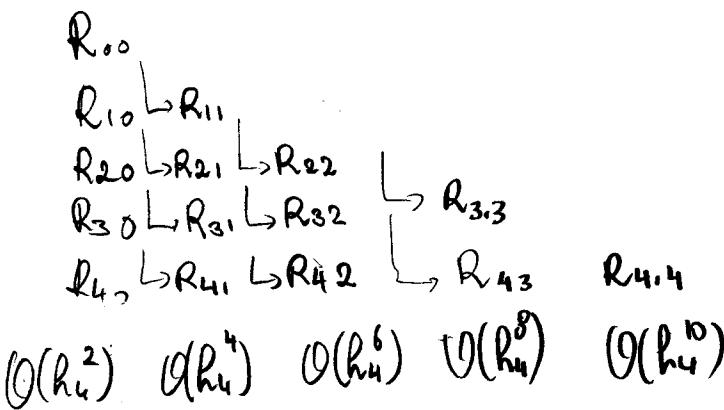
$$3 \int_a^b f(x) dx = 4R_{k+1,0} - R_{k,0} - \frac{3}{4} h_k^4 - \frac{15}{16} h_k^6 - \dots$$

$$\int_a^b f(x) dx = \underbrace{\frac{4R_{k+1,0} - R_{k,0}}{4 - 1}}_{-\frac{1}{4} h_k^4 - \frac{5}{16} h_k^6 - \dots} = R_{k+1,1} = O(h_k^4) \text{ approx}$$

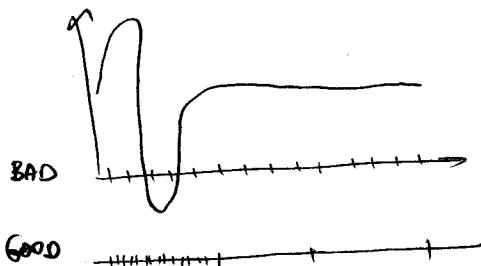
If we keep doing same thing:

$$R_{k+1,j} = \frac{4^j R_{k+1,j-1} - R_{k,j-1}}{4^j - 1} = O(h_k^{2(j+1)})$$

Organizing work in a table we get: (for example)



## § 4.6 Adaptive quadrature methods



Using a uniform partition of an interval for a quadrature rule for approximating the integral of a function is not the most efficient way of doing things!

Obviously we can do better if we adapt the # of points to put more points where they are needed the most  $\Rightarrow$  principle behind adaptive methods.

key ingredient: we need a way that tells us how good our quadrature is working. We could do a number of things: for example using more nodes to estimate error term or using a more accurate method. Here we choose to refine # of points to estimate error.

On some interval  $[a, b]$ :

$$(*) \int_a^b f(x) dx = \underbrace{\frac{b-a}{6} (f(a) + f(\frac{a+b}{2}) + f(b))}_{\text{for some } \xi \in (a, b)} - \frac{h^5}{90} f^{(4)}(\xi) \quad \text{where } h = \frac{b-a}{2}$$

Of course we could estimate error if we knew  $f^{(4)}$  but this is asking too much from end user!

Instead we apply Simpson's rule again on two subintervals of  $[a, b]$ :

$$(**) \int_a^b f(x) dx = S(a, \frac{a+b}{2}) + S(\frac{a+b}{2}, b) - \frac{h^5}{25920} f^{(4)}(\xi_1) - \frac{h^5}{25920} f^{(4)}(\xi_2)$$

$$= \quad " \quad - \frac{h^5}{90} \left(\frac{1}{16}\right) f^{(4)}(\tilde{\xi})$$

where  $\xi_1 \in (a, \frac{a+b}{2})$ ,  $\xi_2 \in (\frac{a+b}{2}, b)$  and  $\tilde{\xi} \in (a, b)$ .

Now assume  $f^{(4)}$  does not vary too much on  $(a, b)$  so:

(11)

$$f^{(4)}(\xi) \approx f^{(4)}(\tilde{\xi})$$

so to estimate the error we simply do  $(\star\star) - (\star)$  to get:

$$\frac{h^5}{90} f^{(4)}(\xi) \approx -\frac{16}{15} \left[ S(a, \frac{a+b}{2}) + S(\frac{a+b}{2}, b) - S(a, b) \right]$$

Plugging into  $(\star\star)$  and bounding integration error:

$$\left| \int_a^b f(x) dx - S(a, \frac{a+b}{2}) - S(\frac{a+b}{2}, b) \right| \leq \frac{1}{15} \left| S(a, b) - S(a, \frac{a+b}{2}) - S(\frac{a+b}{2}, b) \right|$$

So if we want  $(\star\star)$  to be accurate to within  $\epsilon$  we need:

$$\left| S(a, b) - S(a, \frac{a+b}{2}) - S(\frac{a+b}{2}, b) \right| < 15\epsilon$$

To be conservative (and account for some changes in  $f^{(4)}$ ) we can require  $10\epsilon$  bound.

Algorithm: the book implements a recursive procedure using a for loop.  
This is like writing a tree traversal algorithm with for loops: very complicated!  $\rightarrow$  use recursion: i.e. a function that calls itself.

- △ Using recursion is probably less efficient than running code in the book but it's much easier to understand!
- △ Also some programming languages are limited in # of recursion levels that are allowed or if you run a C program you may get "stack overflow" or other messages of the kind.
- △ It may be good to impose a limit on # of recursions in "industrial" code as infinite recursions can be nastier than infinite loops!

function  $r = \text{adapt Simpson}(a, b, f, \epsilon)$

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

$$S_1 = \frac{b-a}{12} (f(a) + 4f\left(a + \frac{b-a}{4}\right) + f\left(\frac{a+b}{2}\right))$$

$$S_2 = \frac{b-a}{12} (f\left(\frac{a+b}{2}\right) + 4f\left(a + \frac{3(b-a)}{4}\right) + f(b))$$

If  $|S - S_1 - S_2| < \epsilon$

$$r = S_1 + S_2 ; \quad \because \text{we are happy with approx}$$

else

$$\rightarrow r = \text{adapt Simpson}\left(a, \frac{a+b}{2}, f, \epsilon/2\right) + \text{adapt Simpson}\left(\frac{a+b}{2}, b, f, \epsilon/2\right)$$

- recursive calls of adaptive Simpson code for left and right halves of intervals. Requiring that error be  $\epsilon/2$  ensures the total error from left + right subint will not exceed  $\epsilon$  when summed up.
- instead of accumulating only integrals we can also accumulate e.g. points of discretization.

# LU and Cholesky factorizations

BLAS 1,2,3  
LAPACK

(14)

Consider the linear system

$$Ax = b, \text{ where } A \in \mathbb{R}^{n \times n}, x \in \mathbb{R}^n, b \in \mathbb{R}^n.$$

Easiest system to solve is when  $A = \text{diag}(a_1, a_2, \dots, a_m)$  i.e.

$$A = \begin{bmatrix} a_1 & & & \\ & a_2 & & \\ & & a_3 & \\ & & & \ddots \\ & & & a_m \end{bmatrix}$$

then: (if the  $a_i \neq 0$ ):

$$x = A^{-1}b = \begin{bmatrix} b_1/a_1 \\ b_2/a_2 \\ \vdots \\ b_m/a_m \end{bmatrix} = b \cdot 1/a \quad \text{in Matlab notation}$$

Here is another system that is convenient:

lower triangular:

$$\begin{pmatrix} \Delta \\ A \end{pmatrix} \begin{bmatrix} x \\ b \end{bmatrix} = \begin{bmatrix} ? \\ ? \end{bmatrix}$$

$$A = \begin{bmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ a_{21} & a_{22} & 0 & \cdots & 0 \\ a_{31} & a_{32} & a_{33} & 0 & \cdots \\ \vdots & & & & \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mm} \end{bmatrix}$$

How do we solve such a system?

$$x_1 = b_1 / a_{11}$$

$$x_2 = (b_2 - a_{21}x_1) / a_{22}$$

$$x_3 = (b_3 - a_{31}x_1 - a_{32}x_2) / a_{33}$$

$$\vdots \quad \sum_{j=1}^{i-1} \\ x_i = (b_i - \sum_{j=1}^{i-1} a_{ij}x_j) / a_{ii} \quad i \leq n$$

This is called forward substitution as we compute  $x_i$  "forward" (direction of increasing  $i$ ). Writing this as a for loop:

for  $i = 1 \dots n$

$$| \quad x_i = (b_i - \sum_{j=1}^{i-1} a_{ij}x_j) / a_{ii}$$

Yet another way of writing same algorithm but using Matlab:

for  $i = 1 \dots n$

$$| \quad x_i = (b_i - A(i, 1:i-1) \cdot x(1:i-1)) / a_{ii}$$

end.

In a similar way upper triangular ( $\Delta$ ) systems are easy to solve:

The algorithm is same but we walk our way backwards from  $x_n$  to  $x_1$ :

for  $i = n : -1 : 1$

$$| \quad x_i = b_i - \sum_{j=i+1}^n a_{ij}x_j$$

(backward substitution)

Obviously if we permute rows of system we can still solve easily. For example consider permutation

$p_1, p_2, \dots, p_m$  (which we can put in a single vector  $p$ )

if the matrix  $[a_{p_i, j}]_{i=1, n, j=1 \dots n}$  is lower (or upper) triangular

then we can perform forward substitution for permuted systems:

for  $i = 1 \dots n$

$$x_i = \left( b_{p_i} - \sum_{j=1}^{i-1} a_{p_i, j} x_j \right) / a_{p_i, i} \quad \underline{m^2 \text{ flops}}$$

or in Matlab notation:

for  $i = 1 : n$ ,

$$x(p_i) = (b_{p_i} - a(p_i), 1:i-1) x) / a_{p(i), i}.$$

Similarly backwards substitution become:

for  $i = n : -1 : 1$ ,

$$x_i = b_{p_i} - \sum_{j=i+1}^n a_{p_i, j} x_j / a_{p_i, i}.$$

LU decomposition:

$$A = L U$$

$\triangleleft$  not all matrices admit LU decomp



If we have LU decom position available:

(117)

solving  $Ax = b$   
 $\Rightarrow \begin{cases} \text{Solving } Lz = b \text{ for } z. \\ \text{Solving } Ux = z \text{ for } x \end{cases}$  { relatively easy to do}

Derivation of LU factorization:

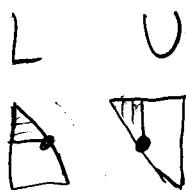


$$(L \cdot U)_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj}$$

$$= \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj}$$

For either values  $l_{kk}$  or  $u_{kk}$ , say  $u_{kk} = 1 \Rightarrow U$  unit upper triangular  
 $\equiv$  Gauss's factor.

$l_{kk} = 1 \Rightarrow L$  unit lower triangular.  
 $\equiv$  Doolittle factor.



Assume we have already computed first  $k \times k$  minor of  $L, U$ :

$$a_{kk} = \sum_{j=1}^{k-1} l_{kj} u_{jk} + l_{kk} u_{kk} \quad \text{let } \boxed{l_{kk} = 1} \quad \text{Doolittle}$$

$$\Rightarrow u_{kk} = a_{kk} - \sum_{j=1}^{k-1} l_{kj} u_{jk}$$

compare to A which has  $n^2$ .

Note: it makes sense to have to specify one of diagonals.

$L, U$  have subtotal  $\frac{n(n+1)}{2} + \frac{n(n+1)}{2}$  elements

$$\left\{ \begin{array}{l} a_{kj} = \sum_{s=1}^{k-1} l_{ks} u_{sj} + \underline{l_{kk}} u_{kj} \quad k+1 \leq j \leq n \\ a_{ik} = \sum_{s=1}^{k-1} l_{is} u_{sk} + \underline{l_{ik}} u_{ik} \quad k+1 \leq i \leq m \\ \\ \Rightarrow u_{kj} = (a_{kj} - \sum_{s=1}^{k-1} l_{ks} u_{sj}) \quad k \leq j \leq m \\ l_{ik} = (a_{ik} - \sum_{s=1}^{k-1} l_{is} u_{sk}) / u_{ik} \quad k+1 \leq i \leq m \end{array} \right.$$

### LW Algorithm

for  $k = 1 : m$ ,

$$l_{kk} = 1$$

for  $j = k : m$

$$| \quad u_{kj} = a_{kj} - \sum_{s=1}^{k-1} l_{ks} u_{sj}$$

end

for  $i = k+1 : m$

$$| \quad l_{ik} = (a_{ik} - \sum_{s=1}^{k-1} l_{is} u_{sk}) / u_{ik}$$

end

} can be done with  
BLAS

$$\frac{2}{3} n^3$$

note: algorithm can be done in place overwriting matrix A.



$$A = L \cup U$$

(L) (U)

L = unit lower triangular  
U = upper triangular

Theorem: If all  $n$ -leading principal minors of  $A \in \mathbb{R}^{n \times n}$  are non-singular than  $A$  has a LU-decomposition.

$k$   $\xrightarrow{\quad}$   $k$ -th leading minor      (this is just a sufficient condition)



In many applications  $A$  is symmetric positive definite: (S.p.d.)

meaning:  $A = A^T$  and

$$\forall x \neq 0 \quad x^T A x > 0$$

This means  $A$  is non-singular (why?)

$$Ax = 0 \Rightarrow x^T Ax = 0 \Rightarrow x = 0.$$

Does  $A$  have an LU factor? Yes since all principal minors of  $A$  are S.p.d. (why?)

$$(x_1, x_2, \dots, x_k, 0, \dots, 0) \quad A \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \\ 0 \\ \vdots \\ 0 \end{pmatrix} = (x_1, x_2, \dots, x_k) A_k \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}$$

where  $A_k$  =  $k$ -th principal minor of  $A$

$$= A(1:k, 1:k)$$

only way this can happen

$$\left. \begin{array}{l} A = LU \\ A^T = U^T L^T \end{array} \right\} \Rightarrow \underbrace{U L^{-T}}_{(\nabla)(\nabla)} = \underbrace{L^{-1} U^T}_{(\Delta)(\Delta)} = D$$

note  $(\nabla)^{-1} = \nabla$      $(\nabla)(\nabla) = (\nabla)$

$(\Delta)^{-1} = \Delta$      $(\Delta)(\Delta) = (\Delta)$

$$A = \underbrace{L D L^{-T}}_{D} L^T \quad \text{D must have positive entries since:}$$

$$D = L^{-1} A L^{-T}$$

$$x^T D x = x^T L^{-1} A L^{-T} x$$

$$= (L^{-T} x)^T A (L^{-T} x) > 0 \text{ when } x \neq 0.$$

$$D^{\frac{1}{2}} = \begin{pmatrix} \sqrt{d_{11}} & & \\ & \sqrt{d_{22}} & \\ & & \sqrt{d_{nn}} \end{pmatrix}$$

thus:  $A = L D^{\frac{1}{2}} D^{\frac{1}{2}} L^T = \tilde{L} \tilde{L}^T$  where  $\tilde{L} = LD^{\frac{1}{2}}$

Theorem: If  $A$  is spd then it has a unique Cholesky factorization

$$A = \tilde{L} \tilde{L}^T$$

(Q)(7)

Algorithm: Look at matrix matrix product:

$$a_{kk} = \sum_{t=1}^{k-1} l_{kt}^2 + l_{kk}^2$$

$$l_{kk} = \left( a_{kk} - \sum_{s=1}^{k-1} l_{ks}^2 \right)^{\frac{1}{2}}$$

$$\text{for } k=1..n$$

$$l_{kk} = \left( a_{kk} - \sum_{s=1}^{k-1} l_{ks}^2 \right)^{\frac{1}{2}}$$

$$\text{for } i=k+1..n$$

$$l_{ik} = \left( a_{ik} - \sum_{s=1}^{k-1} l_{is} l_{ks} \right) / l_{kk}$$

# Gaussian Elimination

(12)

Example:

$$A \mathbf{x} = A^{(1)} \mathbf{x} = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 12 & -8 & 6 & 10 \\ 3 & -13 & 9 & 3 \\ -6 & 4 & 1 & -18 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 34 \\ 27 \\ -38 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 12 & 3 & 1 & 0 \\ -1 & -\frac{1}{2} & 2 & 1 \end{bmatrix}$$

$$A^{(2)} \mathbf{x} = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & -12 & 8 & 1 \\ 0 & 2 & 3 & -14 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ 21 \\ -26 \end{bmatrix}$$

$$A^{(3)} \mathbf{x} = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 4 & -13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -21 \end{bmatrix}$$

$$A^{(4)} \mathbf{x} = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 0 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -3 \end{bmatrix}$$

Ans:

$$\mathbf{x} = \begin{bmatrix} 1 \\ -3 \\ -2 \\ 1 \end{bmatrix}$$

$$A^{(3)} = U$$

formally:

$$A = A^{(1)} \rightarrow A^{(2)} \rightarrow A^{(3)} \rightarrow \dots \rightarrow A^{(n)} = U$$

(all gears used)

at each step:

$$A^{(k)} = \left[ \begin{array}{c|c} \text{triangle} & L \\ \hline 0 & \begin{matrix} x \\ x \\ x \\ x \end{matrix} \end{array} \right] \xrightarrow{k \rightarrow} A^{(k+1)} = \left[ \begin{array}{c|c} \text{triangle} & L \\ \hline 0 & \begin{matrix} x \\ x \\ x \\ x \\ 0 \\ \vdots \\ 0 \end{matrix} \end{array} \right]$$

} unchanged  
} updated

thus to obtain  $A^{(k+1)}$  we need to introduce zeros in column marked by  $x$ ... and we do not need to change  $A^{(k)}$  ( $1:k, 1:k$ )

new entries are:

$$a_{ij}^{(k+1)} = \begin{cases} a_{ij}^{(k)} & \text{if } i \leq k \\ a_{ij}^{(k)} - \frac{a_{ik}^{(k)} a_{kj}^{(k)}}{a_{kk}^{(k)}} & \text{if } i \geq k \text{ and } j \geq k+1 \\ 0 & \text{if } i \geq k+1 \text{ and } j \leq k \end{cases}$$

and  $L$  is defined by:

$$l_{ik} = \begin{cases} a_{ik}^{(k)} / a_{kk}^{(k)} & \text{if } i \geq k+1 \\ 1 & \text{if } i = k \\ 0 & \text{if } i \leq k-1 \end{cases}$$

Theorem: If the pivot  $a_{kk}^{(k)} \neq 0$  then  $A = LU$

When can this break down?

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

can't add multiple of (1st eq) to introduce new zeros in (2nd eq).

even in situations like:

$$\begin{bmatrix} \varepsilon & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

↓ Gaussian Elim:

$$\begin{bmatrix} \varepsilon & 1 \\ 0 & 1-\varepsilon^{-1} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2-\varepsilon^{-1} \end{bmatrix}$$

$$x_2 = \frac{(2-\varepsilon^{-1})}{(1-\varepsilon^{-1})} \approx 1$$

$$x_1 = (1-x_2)\varepsilon^{-1} \approx 0$$

if  $\varepsilon$  is very small.  
say  $\varepsilon <$  machine epsilon

however true sol is:

$$x_2 = \frac{2\varepsilon-1}{\varepsilon-1} \approx 1$$

$$x_1 = \frac{(1-x_2)}{\varepsilon} = \frac{1}{\varepsilon} \left( \frac{\varepsilon-1}{\varepsilon-1} - \frac{(2\varepsilon-1)}{\varepsilon-1} \right)$$

$$= \frac{1}{\varepsilon} \frac{-\varepsilon}{\varepsilon-1} = \frac{1}{1-\varepsilon} \approx 1$$

The problem is not only with smallness of  $a_{11}$ , but so how small is  $a_{11}$  compared to all other elements in its row: (124)

$$\begin{pmatrix} 1 & \varepsilon^{-1} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \varepsilon^{-1} \\ 2 \end{pmatrix} \quad (\text{we just replaced by } \varepsilon \text{ first row})$$

GE

$$\begin{pmatrix} 1 & \varepsilon^{-1} \\ 0 & 1-\varepsilon^{-1} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \varepsilon^{-1} \\ 2-\varepsilon^{-1} \end{pmatrix}$$

sol:

$$x_2 = \frac{2-\varepsilon^{-1}}{1-\varepsilon^{-1}} \approx 1$$

$$x_1 = \varepsilon^{-1} - \varepsilon^1 x_2 \approx 0 \quad \text{which is again wrong!}$$

If we simply exchange rows we get no problem!

$$\begin{pmatrix} 1 & 1 \\ \varepsilon & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

GE

$$\begin{pmatrix} 1 & 1 \\ 0 & 1-\varepsilon \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1-2\varepsilon \end{pmatrix} \Rightarrow x_2 = \frac{1-2\varepsilon}{1-\varepsilon} \approx 1$$

$$x_1 = 2 - x_2 \approx 1$$

We usually need to reorder rows to get LU factorization. Assume we have a permutation  $P$  of  $\{1, 2, \dots, n\}$  so that if we carry operations on rows in this order we do not encounter these problems then:

algorithm becomes:

```

for k = 1 : n-1
  for i = k+1 : n
    |    $\gamma = \underline{a_{pik}} / \underline{a_{pk}}$ 
    |    $a_{pik} = 0$ 
    |   for j = k+1 .. n
    |     |    $a_{pj} = a_{pj} - \gamma a_{pk}$ 
  
```

Gaussian elimination w/ scaled row pivoting:

and result is factor:

$$PA = LU$$

where  $P$  = permutations matrix:  $uP\bar{A} = A$  with rows permuted

$$\text{To solve } Ax = b \Leftrightarrow PAx = Pb$$

$$\Leftrightarrow \underbrace{LUx}_{z} = Pb$$

$$1) \text{ solve } L\underline{z} = Pb$$

$$2) \text{ solve } Ux = z$$

In the factorization phase we simply "choose" for pivot the one that is relatively higher w.r.t rest of its row and we permit it to pivot position. i.e. find  $p_i$  s.t.

$$\frac{|a_{p_1,1}|}{\max_{1 \leq i \leq n} |a_{p_1,i}|} > \frac{|a_{j,1}|}{\max_{1 \leq i \leq m} |a_{j,i}|}$$

Another class of matrices where we can dispense from doing pivoting is:

diagonally dominant matrices:

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}| \quad \text{diagonal element is larger than all other elements in a row.}$$

Thm: Every diag dominant matrix is non-singular and has an L U factor.

Tridiagonal matrix: Assuming no pivoting is needed, it is very efficient to do Gaussian elimination on such matrices.

$$\begin{bmatrix} d_1, c_1 \\ a_1, & \ddots \\ & \ddots & \ddots & c_{n-1} \\ & & a_{n-1}, d_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

in matlab:  $A = \text{diag}(a, -1) + \text{diag}(d, 0) + \text{diag}(c, 1);$

Storage for A = 3 vectors.

$$\begin{bmatrix} x & x & & \\ x & x & x & \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix} \rightarrow \begin{bmatrix} x & x & & \\ 0 & x & x & \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix} \rightarrow \begin{bmatrix} x & x & x & & \\ 0 & x & x & x & \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix} \rightarrow \dots \begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix}$$

first eq:

$$d_2 = d_2 - \frac{a_{12}c_1}{d_1}$$

$$b_2 = b_2 - \frac{a_{12}b_1}{d_1}$$

thus:

$$\left\{ \begin{array}{l} \text{for } i=2 \dots n \\ d_i = d_i - \frac{a_{i-1}c_{i-1}}{d_{i-1}} \\ b_i = b_i - \frac{a_{i-1}b_{i-1}}{d_{i-1}} \end{array} \right\} \quad GE.$$

$$x_n = b_n/d_n$$

$$\text{for } i = n-1; -1; 1$$

$$| x_i = (b_i - c_i x_{i+1}) / d_i$$

Backward  
substitution

## Norms and the analysis of errors.

Doing LU factor + solution of triangular systems is  $\mathcal{O}(n^3)$  operations.  
 How do we know whether error accumulates so much that we cannot trust solution?  
 We need a tool for measuring error in  $\mathbb{R}^n$ : norm.

Def (norm):  $\|x\|$  is a norm, if:

$$\|x\| \geq 0 \quad \forall x \in \mathbb{R}^n$$

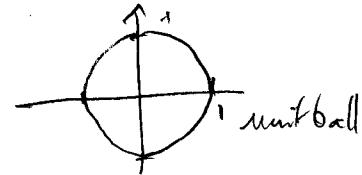
$$\|x\|=0 \Leftrightarrow x=0$$

$$\|\alpha x\| = |\alpha| \|x\| \quad \forall \alpha \in \mathbb{R}, x \in \mathbb{R}^n$$

$$\|x+y\| \leq \|x\| + \|y\| \quad \text{triangle inequality.}$$

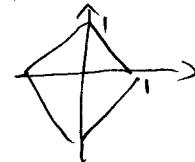
$$\|x\|_2 = \left( \sum_{i=1}^m x_i^2 \right)^{\frac{1}{2}}$$

Euclidean



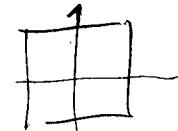
$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

l1 norm



$$\|x\|_\infty = \max_{i=1 \dots n} |x_i|$$

l infinity norm



Induced matrix norms: Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^n$  then:

$$\|A\| = \sup_{\|u\|=1} \|Au\| = \sup_{u \neq 0} \frac{\|Au\|}{\|u\|}$$

$\|\cdot\|$  satisfies all axioms of a norm and so is a norm.

Moreover induced matrix norms satisfy:

$$\|I\| = 1$$

$$\|AB\| \leq \|A\| \|B\|$$

An important case is

$$\|A\|_2^2 = \sup_{x \neq 0} \frac{\|Ax\|^2}{\|x\|^2} = \sup_{x \neq 0} \frac{x^T A^T A x}{x^T x} = \lambda_{\max}(A^T A)$$

$$\Rightarrow \|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$$

"largest eigenvalue of  $A^T A$ ".

Condition number

Suppose  $\tilde{b}$  is a perturbation of  $b$ , If  $x, \tilde{x}$  satisfy

$$Ax = b$$

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

then when are  $x, \tilde{x}$  close?

$$\begin{aligned}\|x - \tilde{x}\| &= \|A^{-1}b - A^{-1}\tilde{b}\| = \|A^{-1}(b - \tilde{b})\| \\ &\leq \|A^{-1}\| \|b - \tilde{b}\|\end{aligned}$$

What about relative error?

$$\begin{aligned}\|x - \tilde{x}\| &\leq \|A^{-1}\| \|b - \tilde{b}\| = \|A^{-1}\| \frac{\|Ax\|}{\|b\|} \|b - \tilde{b}\| \\ &\leq \|A^{-1}\| \|A\| \|x\| \frac{\|b - \tilde{b}\|}{\|b\|}\end{aligned}$$

$$\Leftrightarrow \frac{\|x - \tilde{x}\|}{\|x\|} \leq \underbrace{\|A^{-1}\| \|A\|}_{K(A)} \frac{\|b - \tilde{b}\|}{\|b\|},$$

$$K(A) = \text{condition number of } A = \|A^{-1}\| \|A\|$$

= measures how sensitive is result of linear system to perturbation

= depends on vector norm

$$\geq 1$$

There are ways of estimating  $K(A)$  efficiently to have an idea on how much we can trust solution.

Here is an example of induced matrix norm.

$$\|A\|_{\infty} = \sup_{\|u\|_{\infty}=1} \|Au\|_{\infty}$$

$$= \sup_{\|u\|_{\infty}=1} \max_{1 \leq i \leq n} |(Au)_i|$$

$$= \sup_{\|u\|_{\infty}=1} \max_{1 \leq i \leq n} \sum_{j=1}^m |a_{ij} u_j|$$

$$= \max_{1 \leq i \leq n} \sup_{\|u\|_{\infty}=1} \quad \text{sup attained when } u_j = \operatorname{sgn} a_{ij}$$

$$= \max_{1 \leq i \leq n} \underbrace{\sum_{j=1}^m |a_{ij}|}_{l_1 \text{ norm of row } i}$$

$$A = \begin{bmatrix} 1 & 1+\varepsilon \\ 1-\varepsilon & 1 \end{bmatrix} \quad A^{-1} = \varepsilon^{-2} \begin{bmatrix} 1 & -1-\varepsilon \\ -1+\varepsilon & 1 \end{bmatrix}$$

$$\|A\|_{\infty} = 2 + \varepsilon$$

$$\|A^{-1}\|_{\infty} = \varepsilon^{-2} (2 + \varepsilon)$$

$$\mathcal{F}(A) = \left(\frac{2+\varepsilon}{\varepsilon}\right)^2 \quad \text{if } \varepsilon \leq 0.01 \text{ then } \mathcal{F}(A) > 40000$$

# Neumann Series and Iterative refinement

131

notion of convergence in  $\mathbb{R}^n$ :

$$\lim_{k \rightarrow \infty} \|v^{(k)} - v\| = 0$$

which norm? In  $\mathbb{R}^n$  it doesn't matter! This is because all norms are equivalent in a finite dimensional space:

$\alpha \|x\| \leq \|x\| \leq \beta \|x\|$ , where  $\alpha, \beta > 0$  are constant  
and  $\|x\|$  and  $\|x\|$  are any two norms in  $\mathbb{R}^n$ .

Also  $\mathbb{R}^n$  is complete: any sequence satisfying Cauchy Criterion converges:

$$\left\{ \forall \epsilon > 0 \exists N \text{ s.t. } \forall i, j \geq N \quad \|v^{(i)} - v^{(j)}\| < \epsilon \right\}$$

## Theorem (Neuman Series)

let  $A \in \mathbb{R}^{n \times n}$  & such that  $\|A\| < 1$ , then  $I - A$  is invertible  
and:

$$(I - A)^{-1} = \sum_{k=0}^{\infty} A^k$$

Proof:

Assume  $I - A$  is singular (for contradiction):

$\exists x \in \mathbb{R}^n$  s.t.  $(I - A)x = 0$ , take  $x$  s.t.  $\|x\| = 1$ :

$$1 = \|x\| = \|Ax\| \leq \|A\| \|x\| = \|A\|$$

contradiction!

We need to show

$$\sum_{k=0}^{\infty} A^k \rightarrow (I - A)^{-1}$$

$$(I-A) \sum_{k=0}^m A^k = \sum_{k=0}^m A^k - A^{k+1} = \underbrace{A^0}_I - A^{m+1}$$

Now:  $\|A\| \leq \|A\|^{m+1} \rightarrow 0$  [By property of induced matrix norm:  $\|AB\| \leq \|A\|\|B\|$ ]

thus

$$\rightarrow \left\| (I-A) \sum_{k=0}^m A^k - I \right\| \leq \|A\|^{m+1} \rightarrow 0. \quad \text{QED.}$$

Note: This theorem is very intuitive if you know geometric series:

$$\frac{1}{1-x} = 1 + x + x^2 + \dots \text{ when } |x| < 1.$$

For example:

$$\|(I-A)^{-1}\| \leq \sum_{k=0}^{\infty} \|A^k\| \leq \sum_{k=0}^{\infty} \|A\|^k = \frac{1}{1-\|A\|}$$

Note: This theorem is very useful for the theory, however there are more powerful methods for finding inverse of a matrix.

Theorem (Generalization of Neumann series)

(due to St.)

If  $A$  and  $B$  are  $n \times n$  matrices s.t.  $\|I - \tilde{AB}\| < 1$ , then

$A, B$  are invertible and:

$$A^{-1} = B \sum_{k=0}^{\infty} (I - AB)^k$$

$$B^{-1} = \left[ \sum_{k=0}^{\infty} (I - AB)^k \right] A$$

Proof:  $I - \sum_{k=0}^{\infty} (I - AB)^k = AB$  is invertible and.

$$(AB)^{-1} = \sum_{k=0}^{\infty} (I - AB)^k$$

$$\Rightarrow A^{-1} = B B^{-1} A^{-1} = B \sum_{k=0}^{\infty} (I - AB)^k$$

$$B^{-1} = B^{-1} A^{-1} A = \left( \sum_{k=0}^{\infty} (I - AB)^k \right) A.$$

### Iterative refinement

Let  $x^{(0)}$  be an approx. sol to

$$Ax = b$$

then:

$$\begin{aligned} x &= x^{(0)} + \underbrace{A^{-1}(b - Ax^{(0)})}_{\text{error vector}} - x^{(0)} \\ &= e^{(0)} \end{aligned}$$

let  $r^{(0)} = b - Ax^{(0)}$  = residual vector:

$$Ae^{(0)} = r^{(0)}$$

### Algorithm

$x^{(0)}$  given

for  $k = 0, 1, \dots$

$$\begin{cases} r^{(k)} = b - Ax^{(k)} \\ Ae^{(k)} = r^{(k)} \\ x^{(k+1)} = x^{(k)} + e^{(k)} \end{cases}$$

can improve precision of solution, even if we use a "cheap" version of LU factorization that does not compute facts completely.

(How to look at this method?)

(134)

$$x^{(0)} = Bb \quad \text{where } B = \text{"approx inverse of } A \text{"}$$

$$x^{(k+1)} = x^{(k)} + B(b - Ax^{(k)}) , k \geq 0.$$

$B = \text{"approx inverse of } A \text{"} \text{ means } \|I - AB\| \leq 1$

$$A^{-1} = B \sum_{k=0}^{\infty} (I - AB)^k$$

thus  $x = B \sum_{k=0}^{\infty} (I - AB)^k b \quad (\star)$

solves system  $Ax = b$

Theorem Iterative refinement iterates are:

$$x^{(m)} = B \sum_{k=0}^m (I - AB)^k b \quad (m \geq 0)$$

and because of  $(\star)$ :

$$x^{(m)} \rightarrow x \text{ as } m \rightarrow \infty.$$

Proof: By induction on  $m$

$$\underline{m=0} \quad x^{(0)} = Bb$$

Assume  $m$ th case is true then:

$$x^{(m+1)} = x^{(m)} + B(b - Ax^{(m)})$$

$$= B \sum_{k=0}^m (I - AB)^k b + Bb - AB \sum_{k=0}^m (I - AB)^k b$$

$$= B \left( b + \underbrace{(I - AB) \sum_{k=0}^m (I - AB)^k b}_{\sum_{k=0}^m (I - AB)^{k+1} b} \right) = B \sum_{k=0}^{m+1} (I - AB)^k b$$