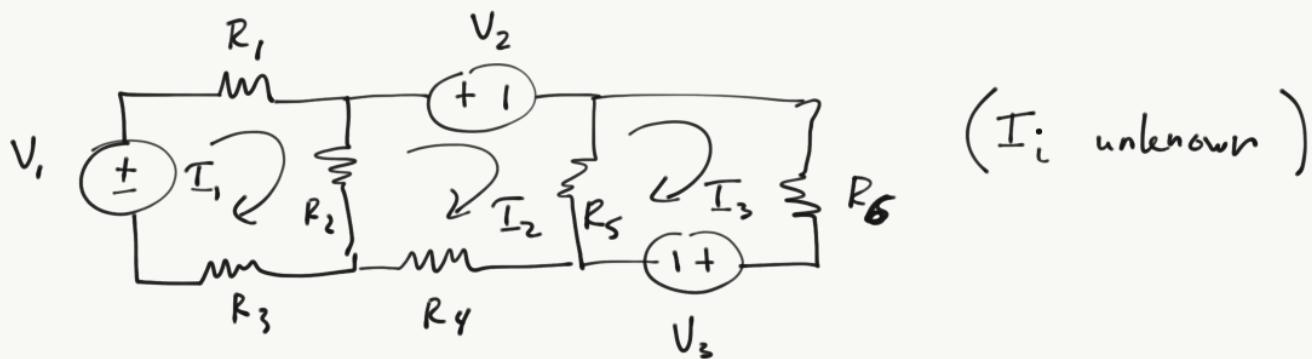


# Linear Algebra

- Sets of linear equations play a key role in many physical systems → these can be represented as matrices
- For example analysis of circuits:



Applying Kirchoff's voltage law around each loop:

$$V_1 - I_1 R_1 - (I_1 - I_2) R_2 - I_1 R_3 = 0$$

$$-V_2 - (I_2 - I_3) R_5 - I_2 R_4 - (I_2 - I_1) R_2 = 0$$

$$- (I_3 - I_2) R_5 - I_3 R_6 - V_3 = 0$$

Rearranging, we find :

$$(-R_1 - R_2 - R_3) I_1 + R_2 I_2 + 0 = -V_1$$

$$+ R_2 I_1 (-R_5 - R_4 - R_2) I_2 + R_5 I_3 = V_2$$

$$+ 0 + R_5 I_2 (-R_5 - R_6) I_3 = V_3$$

- More generally a system of 3 eqns. & 3 unknowns :

$$A_{11}x_1 + A_{12}x_2 + A_{13}x_3 = b_1$$

$$A_{21}x_1 + A_{22}x_2 + A_{23}x_3 = b_2$$

$$A_{31}x_1 + A_{32}x_2 + A_{33}x_3 = b_3$$

- These can be represented in matrix form as :

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

(note compactness!)

$A_{ij}$  are elements of  $\underline{\underline{A}}$   
 $x_j$  are elements of  $\underline{x}$   
 $b_i$  are elements of  $\underline{b}$

- Referring back to the system above & enforcing equivalence with matrix form, it is clear that matrix multiplication is defined by :

$$(\underline{\underline{A}} \underline{x})_i \equiv \sum_j A_{ij}x_j \rightarrow \text{more generally for two matrices}$$

$$(\underline{\underline{AB}})_{ij} \equiv \sum_k A_{ik}B_{kj}$$

matrix multiplication

- For multiplication to make sense matrices must be conformable, viz. of size:  $(m \times p) \cdot \underbrace{(p \times n)}$

inner dimensions equal

$m \times n$

(3)

- Mathematically speaking, a system of equations could fall into one of four categories :
  - 1) no solution (inconsistent)
  - 2) trivial solution  $\underline{x} = \underline{0}$
  - 3) infinite solutions conforming to some constraint
  - 4) unique solution
- For physical situations we expect 4) !!!
- Linear systems play a large role in numerical analysis - networks, fitting & estimation of parameters, interpolation, & ODEs/PDEs
- Special Matrices :

$$1) \text{ column vector } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

$$2) \text{ row vector } \underline{y} = [y_1 \ y_2 \dots]$$

by notational default we assume  $\underline{x}$  is column &  
 $\underline{x}^T$  is row, unless otherwise stated ..

3) square matrix has size  $n \times n$

$$4) \text{ diagonal matrix } \underline{\underline{D}} = \begin{bmatrix} d_{11} & & & \\ & d_{22} & & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$$

11

identity matrix  $\underline{\underline{I}} = \begin{bmatrix} 1 & & & \\ & 1 & 0 & \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix}$  (4)

upper triangular  $\underline{\underline{U}} = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \dots \\ & u_{22} & u_{23} & \\ 0 & & \ddots & \end{bmatrix}$

lower triangular  $\underline{\underline{L}} = \begin{bmatrix} l_{11} & & & \\ & l_{21} & l_{22} & \\ & l_{31} & l_{32} & l_{33} & \dots \end{bmatrix}$

banded  $\underline{\underline{B}} = \begin{bmatrix} b_{11} & 0 & b_{13} & 0 \\ 0 & b_{22} & 0 & b_{24} \\ b_{31} & 0 & b_{33} & 0 \\ 0 & b_{42} & 0 & b_{44} \end{bmatrix}$

- Matrix operations all follow from linear equation properties

$$(\underline{\underline{A}} + \underline{\underline{B}})_{ij} = A_{ij} + B_{ij} ; (\underline{\underline{A}} - \underline{\underline{B}})_{ij} = A_{ij} - B_{ij}$$

$\underline{\underline{A}}, \underline{\underline{B}}$  must have same size!

$$(\alpha \underline{\underline{A}})_{ij} = \alpha A_{ij}$$

$$\underline{\underline{A}}(\underline{\underline{B}} + \underline{\underline{C}}) = \underline{\underline{AB}} + \underline{\underline{AC}} ; \quad \boxed{\underline{\underline{AB}} \neq \underline{\underline{BA}}}$$

(5)

- Matrix inverse  $\underline{\underline{A}}^{-1}$  defined s.t. :

$$\underline{\underline{A}} \underline{\underline{A}}^{-1} = \underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}}$$

- Solutions to  $\underline{\underline{A}} \underline{\underline{B}} = \underline{\underline{C}}$  for  $\underline{\underline{A}}$  :

$$\underline{\underline{A}} = \underline{\underline{B}}^{-1} \underline{\underline{C}}$$

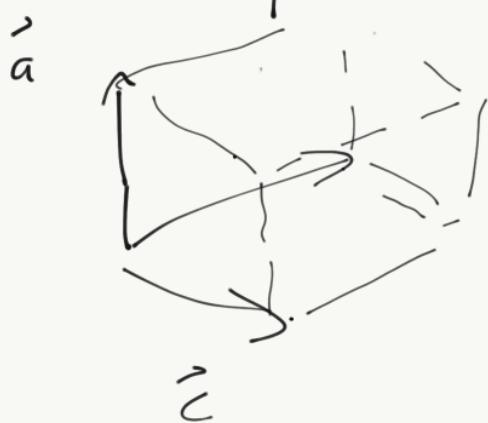
similarly  $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}} \rightarrow \boxed{\underline{\underline{x}} = \underline{\underline{A}}^{-1} \underline{\underline{b}}}$

- Determinants of matrices can be thought of as analogous to concept of "volume" or magnitude

$$\underline{\underline{A}} = \begin{bmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{bmatrix} \quad \vec{a} = (a_1, a_2, a_3) \\ \vec{b} = \dots \text{ etc.}$$

$\det(\underline{\underline{A}}) \equiv |\underline{\underline{A}}|$  is related to the triple product

$\vec{a} \cdot (\vec{b} \times \vec{c})$  which is volume spanned by  $\vec{a}, \vec{b}, \vec{c}$



if  $|\underline{\underline{A}}| = 0$

then  $\vec{a}$  lies  
in plane of  $\vec{b}, \vec{c}$

linearly dependent

- More generally we can define  $|\underline{A}|$  by: ⑥

$$|\underline{A}| = \sum_{j=1}^n A_{ij} \text{ cof}(A_{ij}) = \sum_{j=1}^n A_{ij} (-1)^{i+j} M_{ij}$$

(i arb)     $\downarrow$     cofactor of  $A_{ij}$      $\downarrow$     minor of  $A_{ij}$

$M_{ij}$  is the minor defined by evaluation det of augmented matrix, without row  $i$ , column  $j$

- Note that  $\det(\text{scalar}) = \text{scalar}$ .

- Consider a  $2 \times 2$  determinant

$$|\underline{A}| = \begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = A_{11} (-1)^{1+1} M_{11} + A_{12} (-1)^{1+2} M_{12}$$

$$M_{11} \sim \begin{vmatrix} \cancel{A_{11}} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = \det(A_{22}) = A_{22}$$

$$M_{12} \sim \begin{vmatrix} A_{11} & \cancel{A_{12}} \\ A_{21} & \cancel{A_{22}} \end{vmatrix} = \det(A_{21}) = A_{21}$$

$$\boxed{|\underline{A}| = A_{11} A_{22} - A_{12} A_{21}} \quad (2 \times 2 \text{ det})$$

- Can also expand cofactors along column:

$$\boxed{|\underline{A}| = \sum_{i=1}^n A_{ij} \text{ cof}(A_{ij})} \quad (j \text{ fixed})$$

-  $3 \times 3$  determinants proceed similarly : (7)

$$|A| = A_{11} (-1)^{1+1} M_{11} + A_{12} (-1)^{1+2} M_{12} + A_{13} (-1)^{1+3} M_{13}$$

$$M_{11} \sim \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{vmatrix} = A_{22} A_{33} - A_{23} A_{32}$$

$$M_{12} = A_{21} A_{33} - A_{23} A_{31} ; \quad M_{13} = A_{21} A_{32} - A_{22} A_{31}$$

$$\begin{aligned} |A| &= A_{11} (A_{22} A_{33} - A_{23} A_{32}) - A_{21} (A_{21} A_{33} - A_{23} A_{31}) \\ &\quad + A_{13} (A_{21} A_{32} - A_{22} A_{31}) \end{aligned}$$

- Evaluating an  $n \times n$  determinant requires some adds  
adds &  $n!$  terms, each containing  $(n-1)$   
multiplications, so it scales like:

$$\mathcal{O}((n-1)n!) \quad \begin{matrix} \text{(cofactor expansion} \\ \text{of det)} \end{matrix}$$

## Elimination Methods

(8)

- Cramer's rule is a direct method & useful for establishing a baseline approach for comparing efficiency of various methods:

$$\underline{\underline{A}} \underline{x} = \underline{b} ; \quad x_j = \frac{|\underline{\underline{A}}^j|}{|\underline{\underline{A}}|} ; \quad (\underline{\underline{A}}^j \text{ is } \underline{\underline{A}} \text{ with the } j^{\text{th}} \text{ column replaced by } \underline{b})$$

$$- \text{ e.g. } x_1 = \frac{\begin{vmatrix} b_1 & A_{12} \\ b_2 & A_{22} \end{vmatrix}}{A_{11}A_{22} - A_{12}A_{21}} ; \quad ; \quad \frac{\begin{vmatrix} A_{11} & b_1 \\ A_{21} & b_2 \end{vmatrix}}{A_{11}A_{22} - A_{12}A_{21}} = x_2$$

$$x_1 = \frac{b_1 A_{22} - A_{12} b_2}{A_{11}A_{22} - A_{12}A_{21}} ; \quad ; \quad x_2 = \frac{A_{11} b_2 - b_1 A_{21}}{A_{11}A_{22} - A_{12}A_{21}}$$

- Dealing with larger systems is problematic. For  $\underline{\underline{A}} \sim n \times n$  we must evaluate  $n+1$  determinants.
- Each determinant has  $n!$  products (all containing  $n$  items multiplied) So  $(n!)(n-1)$  multiplications per det.
- Total multiplications :

$$|\underline{\underline{M}}| = \mathcal{O}((n+1)(n-1)n!) \\ = \mathcal{O}((n-1)(n+1)!)$$

(9)

- Thus, Cramer's rule is too computationally expensive to practically use. (Slow IN MATLAB).
- Elimination methods address this shortcoming; these operate very much like a normal solution method by hand.
  - 1) perform algebraic substitutions until value of one variable can be found (fwd elimination)
  - 2) substitute this variable back into simplified eqns successively to find other vars.  
(back substitution)
- Elementary row operations (allow elimination steps)
  - 1) Any row can be multiplied by scalar (scaling)
  - 2) Rows can be interchanged (pivoting)
  - 3) Rows may be linearly combined with other rows.  
(elimination)
- None of these row operations change the solution to the system  $\underline{A}\underline{x} = \underline{b}$
- Connection of row operations to matrix multiplication, e.g.:

$$1) \underline{E}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}; \underline{E}_1 \underline{A} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ \alpha A_{21} & \alpha A_{22} & \alpha A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}$$

$$2) \underline{E}_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \underline{E}_2 \underline{A} = \begin{bmatrix} A_{21} & A_{22} & A_{23} \\ A_{11} & A_{12} & A_{13} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}$$

$$\underline{\underline{E}}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \alpha & 0 & 1 \end{bmatrix}; \quad \underline{\underline{E}}_3 \underline{\underline{A}} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ \alpha A_{11} + A_{31} & \alpha A_{12} + A_{32} & \alpha A_{13} + A_{33} \end{bmatrix}$$

(10)

- Because of the connection of row ops to multiplication we can conceive as an inverse as a sequence of row ops.

$$\left. \begin{array}{l} \underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}} \\ \left( \prod_i \underline{\underline{E}}_i \right) \underline{\underline{A}} = \underline{\underline{I}} \end{array} \right\}$$

so its possible to  
find an inverse  
hence solution from  
row operations!!!

- So if we find row ops that reduce  $\underline{\underline{A}}$  to  $\underline{\underline{I}}$  we can apply those ops to  $\underline{\underline{I}}$  to get  $\underline{\underline{A}}^{-1}$ :

$$\underline{\underline{A}}^{-1} = \left( \prod_i \underline{\underline{E}}_i \right) \cdot \underbrace{\left( \prod_i \underline{\underline{E}}_i \right) \underline{\underline{I}}}_{\underline{\underline{I}}}$$

NOTE: Matrix multiplication is associative!

$$\underline{\underline{A}}(\underline{\underline{B}}\underline{\underline{C}}) = (\underline{\underline{A}}\underline{\underline{B}})\underline{\underline{C}}$$

$$\begin{array}{l}
 x_1 + 4x_2 + 2x_3 = 15 \\
 3x_1 + 2x_2 + x_3 = 10 \\
 2x_1 + x_2 + 3x_3 = 13
 \end{array} \rightarrow \left[ \begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 3 & 2 & 1 & 10 \\
 2 & 1 & 3 & 13
 \end{array} \right] \quad (11)$$

elim. from row 1 :

$$\left[ \begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 3 & 2 & 1 & 10 \\
 2 & 1 & 3 & 13
 \end{array} \right] \quad \begin{array}{l}
 R_2 - 3R_1 \\
 R_3 - 2R_1
 \end{array} \quad \begin{array}{c}
 \\ \\
 \cdots \cdots \cdots
 \end{array}$$

elim. from row  $R_2$  :

$$\left[ \begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 0 & -10 & -5 & -35 \\
 0 & -7 & -1 & -17
 \end{array} \right] \quad \begin{array}{l}
 R_3 - \frac{7}{10}R_2
 \end{array} \quad \begin{array}{c}
 \\ \\
 \begin{array}{l}
 -1 - \frac{7}{10}(-5) \\
 = \frac{5}{2}
 \end{array} \\
 -17 - \frac{7}{10}(-35)
 \end{array}$$

final state following elim :

$$\left[ \begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 0 & -10 & -5 & -35 \\
 0 & 0 & \frac{5}{2} & \frac{15}{2}
 \end{array} \right] \quad \begin{array}{c}
 \\ \\
 \begin{array}{l}
 -170 + 245 \\
 \hline
 10 \\
 = \frac{75}{10} = \frac{15}{2}
 \end{array}
 \end{array}$$

back substitution :

$$x_3 = \left(\frac{15}{2}\right) \left(\frac{2}{5}\right) = 3$$

$$x_2 = \left(-35 + 5(3)\right) \left(-\frac{1}{10}\right) = -20 \left(-\frac{1}{10}\right) = 2$$

$$x_1 = (15 - 2(3) - 4(2)) = 1$$

- Two modifications to standard elimination schemes are advisable : pivoting & slicing. (12)
- The pivot element for an elimination step is the element being used for elimination. e.g.

pivot  $\left[ \begin{array}{ccc|c} 3 & 2 & 1 & 10 \\ 2 & 1 & 3 & 13 \\ 1 & 4 & 2 & 15 \end{array} \right]$   $R_2 - \frac{2}{3}R_3$   $R_3 - \frac{1}{3}R_3$  elimination multiplier

- Interchanging rows (reordering eqns.) doesn't change solution - so one can switch pivots with impunity
- For reasons of numerical accuracy, it is always best to scale each row by largest element in that row & then switch to largest pivot.

$$\left[ \begin{array}{ccc|c} 1 & 4 & 2 & 15 \\ 3 & 2 & 1 & 10 \\ 2 & 1 & 3 & 13 \end{array} \right] \quad \begin{array}{l} \text{pivot: } \frac{1}{15} \rightarrow \text{swap} \\ \text{pivot: } \frac{3}{10} \\ \text{pivot: } \frac{2}{13} \end{array}$$

↓

$$\left[ \begin{array}{ccc|c} 3 & 2 & 1 & 10 \\ 1 & 4 & 2 & 15 \\ 2 & 1 & 3 & 13 \end{array} \right] \quad \begin{array}{l} \text{do elimination, then reorder} \\ \text{again for largest scaled pivot} \end{array}$$

→ defines Gaussian elimination

- Gauss-Jordan elimination is based on the fact that one may compute inverse as a sequence of row ops on  $\underline{\underline{I}}$ :

$$\underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}} \rightarrow \left( \prod_i \underline{\underline{E}}_i^{-1} \right) \underline{\underline{A}} = \underline{\underline{I}}$$

$$\therefore \left( \prod_i \underline{\underline{E}}_i^{-1} \right) \underline{\underline{I}} = \underline{\underline{A}}^{-1}$$

$$\left[ \begin{array}{ccc|ccc} 1 & 4 & 2 & 1 & 0 & 0 \\ 3 & 2 & 1 & 0 & 1 & 0 \\ 2 & 1 & 3 & 0 & 0 & 1 \end{array} \right] = \left[ \begin{array}{c|c} \underline{\underline{A}} & \underline{\underline{I}} \end{array} \right] \quad \begin{aligned} R_2 &- 3R_1 \\ R_3 &- 2R_1 \end{aligned}$$

- See book examples & homework.
- Generally elimination methods of all types are more efficient than solving a system via Cramer's Rule; which scaled like

$$\mathcal{O}(n(n-1)(n+1)!) \quad \underbrace{\text{Cramer's Rule}}$$

- E.g. Gauss elimination scales like:  $\mathcal{O}\left(\frac{n^3}{3} - \frac{n}{3}\right)$

## Determinants

(14)

- Determinants can be evaluated efficiently using elimination methods.
- Consider an upper triangular matrix resulting from elimination:

$$\underline{U} = \begin{bmatrix} U_{11} & U_{12} & U_{13} & \dots \\ 0 & U_{22} & U_{23} & \dots \\ 0 & 0 & U_{33} & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}$$

$\det(\underline{U}) = (\text{expand along } 1^{\text{st}} \text{ column})$

$$= U_{11} (-1)^{1+1} \left| \underline{M}_{11} \right| = U_{11} (-1)^{1+1} \left[ U_{22} (-1)^{2+1} \left| \underline{M}_{22} \right| \right]$$

$$= \dots = \prod_i U_{ii} \quad (\text{expand along } 1^{\text{st}} \text{ column of } \underline{M}_{ii})$$

- So evaluating dets of triangular matrices is trivial
- Performing simple elimination does not change the determinant, thus we can do fwd elimination then evaluate det from result.
- This is more efficient than the  $(n-1)n!$  operations implied by cofactor expansions.
- Pivoting changes value of determinants in predictable ways:

- 1) scaling row by a const also scales det by const
- 2) swapping rows flips sign of det.

## LU factorization

(15)

- One can factor a matrix in infinite ways; one useful way is to factor into upper & lower triangular matrices:

$$\underline{\underline{A}} = \underline{\underline{L}} \underline{\underline{U}}$$

- This can be used to solve  $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}}$  as follows:

$$\underline{\underline{L}} \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{b}} ; \quad \underline{\underline{L}}^{-1} \underline{\underline{L}} \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{L}}^{-1} \underline{\underline{b}}$$

$$\underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{L}}^{-1} \underline{\underline{b}}$$

- Defining  $\underline{\underline{b}}' = \underline{\underline{L}}^{-1} \underline{\underline{b}}$  we have the following eqns:

$$(1) \quad \underline{\underline{L}} \underline{\underline{b}}' = \underline{\underline{b}} \quad (2) \quad \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{b}}'$$

- Thus we can view solving  $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}}$  as solving (1), (2) above
- This is advantageous since triangular systems can be solved easily using forward/backward substitution. Also for systems with multiple RHS,  $\underline{\underline{L}}, \underline{\underline{U}}$  need only be computed once!
- An example with multiple RHS would be finding an inverse via LU factorization:

find  $\underline{\underline{A}}^{-1}$  s.t.  $\underline{\underline{A}} (\underline{\underline{A}}^{-1}) = \underline{\underline{I}}$  can be viewed as

76

$$\left. \begin{array}{l} \underline{\underline{A}} \underline{\underline{x}}_1 = \underline{\underline{b}}_1 \\ \underline{\underline{A}} \underline{\underline{x}}_2 = \underline{\underline{b}}_2 \\ \underline{\underline{A}} \underline{\underline{x}}_3 = \underline{\underline{b}}_3 \end{array} \right\} \quad \underline{\underline{A}}^{-1} = \begin{bmatrix} \underline{\underline{x}}_1 & \underline{\underline{x}}_2 & \underline{\underline{x}}_3 \end{bmatrix}$$
$$\underline{\underline{b}}_1 = \underline{\underline{I}}_{i1} \quad (\text{first column of } \underline{\underline{I}})$$
$$\underline{\underline{b}}_2 = \underline{\underline{I}}_{i2} \quad (2^{\text{nd}} \quad " \quad " \quad ")$$
$$\underline{\underline{b}}_3 = \underline{\underline{I}}_{i3} \quad (3^{\text{rd}} \quad " \quad " \quad ")$$

## Tri diagonal Systems

(17)

- Consider the system :

$$\underline{\underline{T}} \underline{x} = \underline{b};$$

$$\underline{\underline{T}} = \begin{bmatrix} T_{11} & T_{12} & & & \\ T_{21} & T_{22} & T_{23} & & \\ & T_{23} & T_{33} & T_{34} & \\ & & T_{34} & T_{44} & T_{45} \\ & & & \ddots & \\ & & & & T_{i,i-1} & T_{ii} & T_{i,i+1} \\ & & & & & \ddots & \\ & & & & & & T_{n,n-1} & T_{nn} \end{bmatrix}$$

- Elimination applied to this system is particularly simple. E.g. only one element in row 2, to be eliminated using row 1. (column 1)

$$R_2 \rightarrow R_2 - \frac{T_{21}}{T_{11}} R_1; \text{ viz. } \overline{T_{2j}} = T_{2j} - \frac{T_{21}}{T_{11}} T_{1j} \quad j \in \{2, 3\}$$

- Now full elimination can proceed downward to subsequent rows...

- We must also alter RHS from elim step : (18)

$$b_2 = b_2 - \frac{T_{21}}{T_{11}} b_1$$

- These two steps can be encoded generally as ( $j \leq n$ )

$$1) \quad T_{i,j} = T_{i,j} - \left( \frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,j}$$

$$(i \in \{1, n\}; j \in \{i-1, i+1\})$$

there are two of these to perform :

$$T_{i,i} = T_{i,i} - \left( \frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,i}$$

$$T_{i,i+1} = T_{i,i+1} - \left( \frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,i+1}$$

- 2) Following by RHS :

$$b_i = b_i - \left( \frac{T_{i,i-1}}{T_{i-1,i-1}} \right) b_{i-1}$$

- Following these two steps  $\underline{T}$  has been triangulated

$\underline{T} \rightarrow \underline{T}'$ ; we can now perform backsubstitution...

$$x_n = \frac{b_n}{T_{nn}}; \quad x_i = \frac{1}{T_{ii}} (b_i - T_{i,i+1} x_{i+1})$$

# Condition Number

19

- Elimination methods can be sensitive to round-off error which propagates thru subsequent steps ...
- One way to measure this sensitivity is thru the matrix condition number which is computed using norms
- A norm is a measure of magnitude & satisfies the constraints:

$$a) \|\underline{A}\| > 0$$

$$b) \|\underline{A}\| = 0 \text{ iff } \underline{A} = \underline{0}$$

$$c) \|\alpha \underline{A}\| = |\alpha| \|\underline{A}\|$$

$$d) \|\underline{A} + \underline{B}\| = \|\underline{A}\| + \|\underline{B}\|$$

$$e) \|\underline{A}\underline{B}\| \leq \|\underline{A}\| \|\underline{B}\| \quad (\underbrace{\text{Schwarz}}_{\text{inequality}})$$

- Commonly used vector norms:

$$\|\underline{x}\|_1 = \sum_i |x_i| \quad L^1 - \text{norm}$$

$$\|\underline{x}\|_2 = \sqrt{\sum_i x_i^2} \quad L^2 - \text{norm} \quad (\text{Euclidean})$$

$$\|\underline{x}\|_\infty = \max(|x_i|) \quad L^\infty - \text{norm}$$

- Matrix norms can be similarly defined : (20)

$$\|\underline{A}\|_1 = \max_j \left\{ \sum_i |A_{ij}| \right\} \quad (\text{max col. sum})$$

$$\|\underline{A}\|_\infty = \max_i \left\{ \sum_j A_{ij} \right\} \quad (\text{max row sum})$$

$$\|\underline{A}\|_2 = \min(\lambda_{\cdot\cdot}) \quad / \begin{matrix} \text{min eigenvalue,} \\ \text{spectral norm} \end{matrix}$$

$$\|\underline{A}\|_e = \sqrt{\sum_{i,j} A_{ij}^2} \quad (\text{Euclidean norm})$$

- The condition number measures sensitivity to small changes (e.g. round-off) in system.
- Considering the system :

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

we know  $\|\underline{b}\| \leq \|\underline{A}\| \|\underline{x}\|$

- Let us examine system perturbed by  $\underline{\delta b}$

$$\underline{A}(\underline{x} + \underline{\delta x}) = \underline{b} + \underline{\delta b} = \underline{A}\underline{x} + \underline{A}\underline{\delta x} = \underline{b} + \underline{\delta b}$$

- Subtract original equation :

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

$$\Rightarrow \underline{\delta x} = \underline{A}^{-1} \underline{\delta b}$$

thus  $\|\underline{\delta x}\| \leq \|\underline{A}^{-1}\| \|\underline{\delta b}\|$

- Multiplying the two Schwarz-inequalities : (21)

$$\|\underline{b}\| \|\underline{\delta x}\| \leq \|\underline{A}\| \|\underline{x}\| \|\underline{A}^{-1}\| \|\underline{\delta b}\|$$

$$\|\underline{b}\| \|\underline{\delta x}\| \leq \|\underline{x}\| \|\underline{\delta b}\| (\|\underline{A}\| \|\underline{A}^{-1}\|)$$

$$\frac{\|\underline{\delta x}\|}{\|\underline{x}\|} \leq \frac{\|\underline{\delta b}\|}{\|\underline{b}\|} C(\underline{A})$$

$$C(\underline{A}) \equiv \|\underline{A}\| \|\underline{A}^{-1}\| \quad (\text{condition number})$$

- A large condition number  $\Rightarrow$  modest  $\|\underline{\delta b}\|/\|\underline{b}\|$  can lead to large  $\|\underline{\delta x}\|/\|\underline{x}\|$  (BAD!!)
- $C(\underline{A}) \sim 1$  is best.

## Iterative Methods

(22)

- One drawback of elimination methods is the accumulation of error through successive steps.
- Iterative methods can address this issue, but require a matrix to be diagonally dominant, i.e.

$$\boxed{\exists i \text{ s.t. } |A_{ii}| \geq \sum_{j, j \neq i} |A_{ij}|}$$

- Many Systems will not meet this requirement...
- Iterative methods start with an initial guess and then refine that guess repeatedly until convergence is achieved.
- Consider:  $\underline{A}\underline{x} = \underline{b}$ , for some guess  $\underline{x}^{(0)}$

the residual (error in a sense) is  $\underline{b} - \underline{A}\underline{x}^{(0)}$

$$\boxed{R_i^{(k)} = b_i - \sum_j A_{ij}x_j^{(k)}} \quad \begin{matrix} \text{residual on } i^{\text{th}} \\ \text{var } \in \text{ iteration} \\ k \end{matrix}$$

- When the residual is small the iterated solution is near the true solution.
- When the solution  $\underline{x}^{(k)}$  does not change appreciably between iterations, it has converged.

- Jacobi iteration involves a correction to each variable based on residual : (23)

$$x_i^{(k+1)} = x_i^{(k)} + \frac{r_i^{(k)}}{A_{ii}}$$

- The basis of this idea is that, given some guess for all vars. (sans  $i$ ) we can directly compute the variable  $i$  :

$$\begin{aligned} A_{ii}x_i^{(k+1)} &\approx b_i - \sum_{j \neq i} A_{ij}x_j^{(k)} \\ &= A_{ii}x_i^{(k)} + b_i - \sum_j A_{ij}x_j^{(k)} \\ \therefore x_i^{(k+1)} &= x_i^{(k)} = \frac{1}{A_{ii}} \left( b_i - \sum_j A_{ij}x_j^{(k)} \right) \end{aligned}$$


---

- letting  $\underline{\Delta x} = \underline{x}^{(k+1)} - \underline{x}^{(k)}$  we can define convergence as :  $|\underline{\Delta x}| \leq \varepsilon$  for some iteration step
- Relative error can also be used for convergence but becomes problematic if a solution is near zero.

- Gauss-Seidel iteration uses the partially updated solution to determine  $R_i^{(k)}$ :

$$R_i^{(k)} = b_i - \sum_{j < i} A_{ij} x_j^{(k+1)} - \sum_{j \geq i} A_{ij} x_j^{(k)}$$

$$x_i^{(k+1)} = x_i^{(k)} + \frac{R_i^{(k)}}{A_{ii}} \quad (\text{same})$$

- Convergence is almost always faster than Jacobi iteration & it's "simpler" to code.
- Successive over-relaxation (SOR) is a further refinement to speed convergence.

$$x_i^{(k+1)} = x_i^{(k)} + \omega \frac{R_i^{(k)}}{A_{ii}} \quad (\omega < 2)$$

- The relaxation parameter ( $\omega$ ) must be chosen essentially by hand & some choices will make convergence slower...