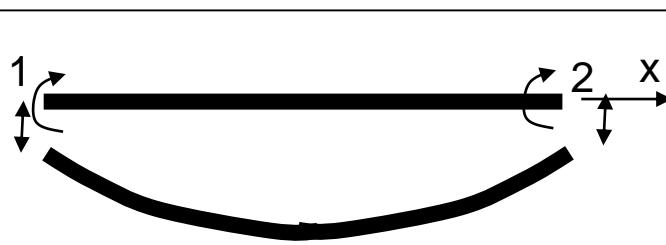


# Stiffness matrix from basics: beam elements


$$EI \frac{d^4 w}{dx^4} = p(x) = 0$$
$$w = \frac{1}{6}C_1 x^3 + \frac{1}{2}C_2 x^2 + C_3 x + C_4$$

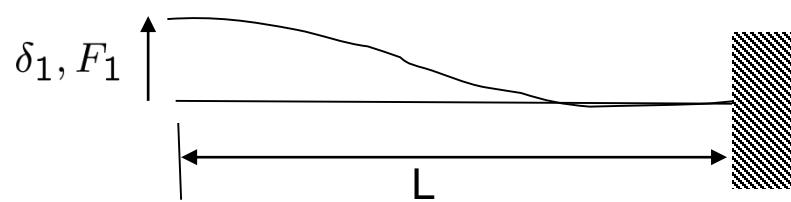
At  $x = 0$

$$\frac{dw}{dx} = 0 \Rightarrow C_3 = 0$$

At  $x = L$

$$\frac{dw}{dx} = 0, w = 0 \Rightarrow$$

$$C_1 L = -2C_2, C_4 = -\frac{1}{6}C_2 L^2$$



$$\begin{aligned}
\frac{d^2w}{dx^2} &= \frac{M(x)}{EI} \\
\frac{dM}{dx} &= V \\
\Rightarrow EI \frac{d^3w}{dx^3} &= V.
\end{aligned} \tag{1}$$

At  $x = 0$

$$\frac{d^3w}{dx^3} = \frac{F_1}{EI} \Rightarrow C_1 = \frac{F_1}{EI}$$

Thus,

$$K_{11} = \frac{F_1}{\delta_1} = \frac{F_1}{w(x=0)} = \frac{12EI}{L^3}$$

Intuitive but not easy

$$\mathbf{K}_L = \frac{EI}{L} \begin{pmatrix} 12/L^2 & -6/L & -12/L^2 & -6/L \\ 4 & 6/L & 2 & \\ & 12/L^2 & 6/L & \\ & & 4 & \end{pmatrix}$$



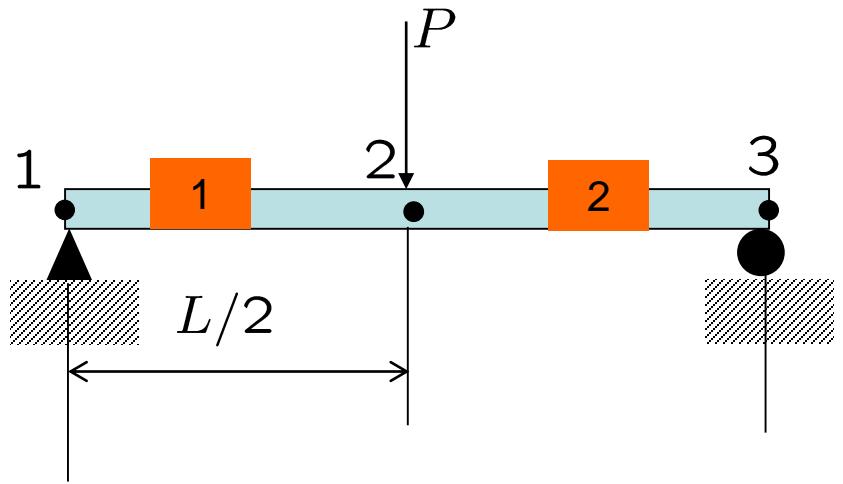
22



12



13



For each element

$$K_L^{1/2} = \frac{8EI}{L^3} \begin{pmatrix} 12 & 6(L/2) & -12 & 6(L/2) \\ 6L/2 & 4(L/2)^2 & -6(L/2) & 2(L/2)^2 \\ -12 & -6(L/2) & 12 & -6(L/2) \\ 6(L/2) & 2(L/2)^2 & -6(L/2) & 4(L/2)^2 \end{pmatrix}$$

The assembled global stiffness

$$\mathbf{K}_G = \frac{8EI}{L^3} \begin{pmatrix} 12 & 6(L/2) & -12 & 6(L/2) & 0 & 0 \\ 6L/2 & 4(L/2)^2 & -6(L/2) & 2(L/2)^2 & 0 & 0 \\ -12 & -6(L/2) & 12 + 12 & -6(L/2) - 6(L/2) & -12 & 6(L/2) \\ 6(L/2) & 2(L/2)^2 & -6(L/2) + 6(L/2) & 4(L/2)^2 + 4(L/2)^2 & -6(L/2) & 2(L/2)^2 \\ 0 & 0 & -12 & -6(L/2) & 12 & -6(L/2) \\ 0 & 0 & 6(L/2) & 2(L/2)^2 & -6(L/2) & 4(L/2)^2 \end{pmatrix}$$

The global force vector

$$\mathbf{F} = \begin{Bmatrix} 0 \\ 0 \\ -P \\ 0 \\ 0 \\ 0 \end{Bmatrix}$$

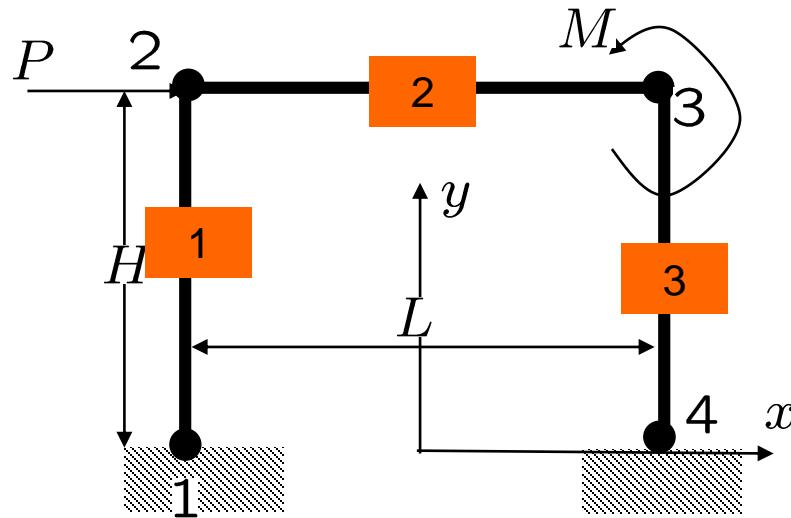
The global displacement vector

$$\mathbf{F} = \begin{Bmatrix} w_1 = 0 \\ \theta_1 \\ w_2 \\ \theta_2 \\ w_3 = 0 \\ \theta_3 \end{Bmatrix}$$

Finally we solve

$$\frac{8EI}{L^3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 4(L/2)^2 & -6(L/2) & 2(L/2)^2 & 0 & 0 \\ 0 & -6(L/2) & 12 + 12 & -6(L/2) - 6(L/2) & 0 & 6(L/2) \\ 0 & 2(L/2)^2 & -6(L/2) + 6(L/2) & 4(L/2)^2 + 4(L/2)^2 & 0 & 2(L/2)^2 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 6(L/2) & 2(L/2)^2 & 0 & 4(L/2)^2 \end{pmatrix} \begin{Bmatrix} w_1 \\ \theta_1 \\ w_2 \\ \theta_2 \\ w_3 \\ \theta_3 \end{Bmatrix} = \begin{Bmatrix} 0 \\ 0 \\ -P \\ 0 \\ 0 \\ 0 \end{Bmatrix}$$

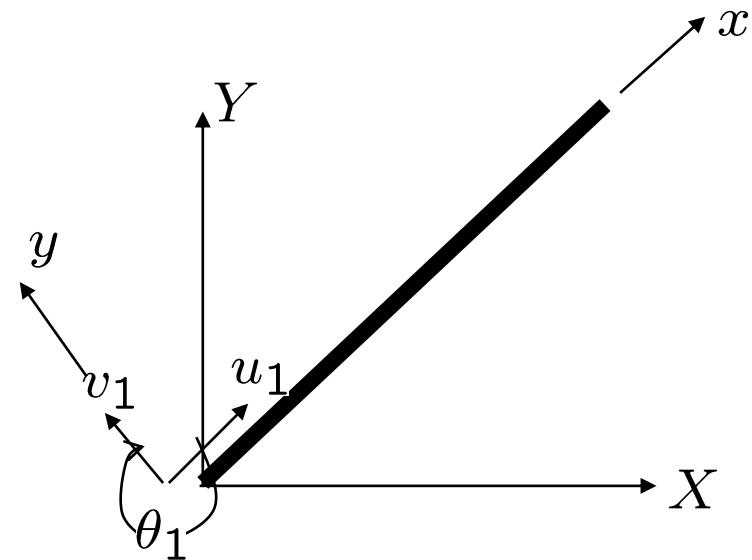
Beams in 2-d



The local stiffness for any element (say 2)

$$K_L^2 = \begin{pmatrix} AE/L & 0 & 0 & -AE/L & 0 & 0 \\ 0 & 12EI/L^3 & 6EI/L^2 & 0 & -12EI/L^3 & 6EI/L^2 \\ 0 & 6EI/L^2 & 4EI/L & 0 & -6EI/L^2 & 2EI/L \\ -AE/L & 0 & 0 & AE/L & 0 & 0 \\ 0 & -12EI/L^3 & -6EI/L^2 & 0 & 12EI/L^3 & -6EI/L^2 \\ 0 & -6EI/L^2 & -4EI/L & 0 & 6EI/L^2 & -2EI/L \end{pmatrix}$$

Transformation Matrix for a 2-d beam element



$$\begin{pmatrix} u_1 \\ v_1 \\ \theta_1 \end{pmatrix} = \begin{pmatrix} \cos \phi & \sin \phi & 0 \\ -\sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} U_1 \\ V_1 \\ \Theta_1 \end{pmatrix} = tU$$

$$u = \begin{pmatrix} t & 0 \\ 0 & t \end{pmatrix} U = TU$$

$$\mathbf{K}_G = \mathbf{T}^T \mathbf{K}_L \mathbf{T}$$

## Solving the equations

The stiffness matrix in the worst case is  $N \times N$ . We can take advantage of its sparsity, diagonal dominance and symmetry to design storage and solution methodologies. The basic problem is to start from

$$\mathbf{K}\mathbf{U} = \mathbf{F}$$

and determine the unknown displacements  $\mathbf{U}$  by

$$\mathbf{U} = \mathbf{K}^{-1}\mathbf{F}$$

We will look at a few numerical techniques for tackling this problem. In particular, we will look at direct methods. There are other iterative methods like conjugate gradient schemes that are suitable for FE calculations in large parallel computing environments.

# Basic Gauss elimination:an example

$$\left( \begin{array}{cccc} 18 & -6 & -6 & 0 \\ -6 & 12 & 0 & -6 \\ -6 & 0 & 12 & -6 \\ 0 & -6 & -6 & 12 \end{array} \right) \left\{ \begin{array}{c} U_1 \\ U_2 \\ U_3 \\ U_4 \end{array} \right\} = \left\{ \begin{array}{c} 60 \\ 0 \\ 20 \\ 0 \end{array} \right\}$$

Original problem

$$\left( \begin{array}{cccc} 18 & -6 & -6 & 0 \\ 0 & 10 & -2 & -6 \\ 0 & -2 & 10 & -6 \\ 0 & -6 & -6 & 12 \end{array} \right) \left\{ \begin{array}{c} U_1 \\ U_2 \\ U_3 \\ U_4 \end{array} \right\} = \left\{ \begin{array}{c} 60 \\ 20 \\ 40 \\ 0 \end{array} \right\}$$

1<sup>st</sup> elimination:  
Row2- Row 1\*(-6/18)  
Row3-Row1\*(-6/18)

$$\left( \begin{array}{cccc} 18 & -6 & -6 & 0 \\ 0 & 10 & -2 & -6 \\ 0 & 0 & 9.6 & -7.2 \\ 0 & 0 & -7.2 & 8.4 \end{array} \right) \left\{ \begin{array}{c} U_1 \\ U_2 \\ U_3 \\ U_4 \end{array} \right\} = \left\{ \begin{array}{c} 60 \\ 20 \\ 44 \\ 12 \end{array} \right\}$$

2<sup>nd</sup> elimination  
Row3-Row2\*(-2/10)  
Row4-Row2\*(-6/10)

$$\left( \begin{array}{cccc} 18 & -6 & -6 & 0 \\ 0 & 10 & -2 & -6 \\ 0 & 0 & 9.6 & -7.2 \\ 0 & 0 & 0 & 3 \end{array} \right) \left\{ \begin{array}{c} U_1 \\ U_2 \\ U_3 \\ U_4 \end{array} \right\} = \left\{ \begin{array}{c} 60 \\ 20 \\ 44 \\ 45 \end{array} \right\}$$

3<sup>rd</sup> elimination:  
Row4-Row3\*(-7.2/9.6)

To solve for  $u_1, \dots, u_4$ , we now *back substitute*:

$$u_4 = 45/3$$

$$u_3 = (44 + 7.2U_4)/9.6 = 15.83$$

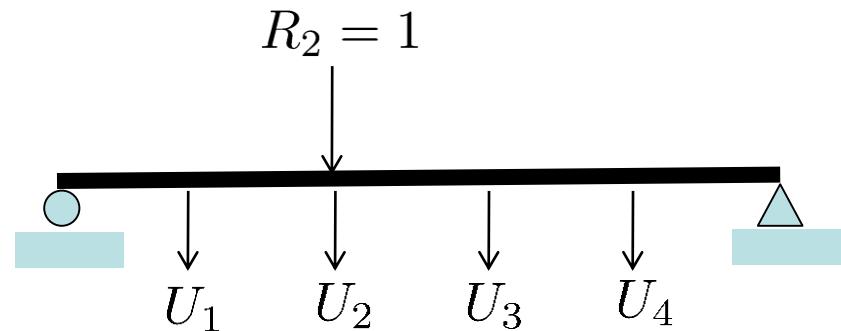
$$u_2 = (20 + 2U_3 + 6U_4)/10 = 14.17$$

$$u_1 = (60 + 6U_2 + 6U_3)/18 = 13.33$$

Note that, at any stage if we encounter a zero on the diagonal, the procedure will fail. It is advisable to resort to *pivoting* to take care of this issue.

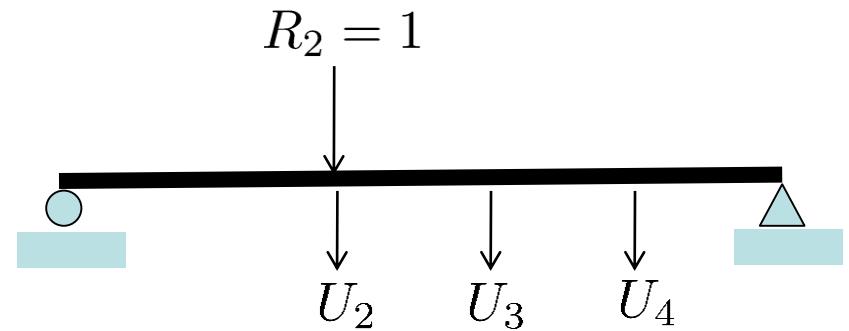
# A physical picture of the Gauss elimination procedure

To illustrate the Gauss elimination procedure, consider a simply supported beam as shown in the figure. We wish to solve for dof's  $U_1 \dots U_4$ . Assume that the final equation of equilibrium for this problem is



$$\begin{pmatrix} 5 & -4 & 1 & 0 \\ -4 & 6 & -4 & 1 \\ 1 & -4 & 6 & -4 \\ 0 & 1 & -4 & 5 \end{pmatrix} \begin{Bmatrix} U_1 \\ U_2 \\ U_3 \\ U_4 \end{Bmatrix} = \begin{Bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{Bmatrix}$$

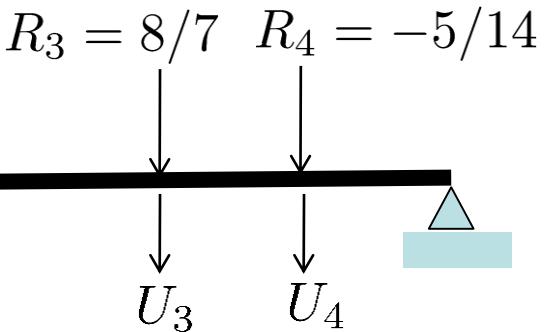
The first elimination will be equivalent to getting rid of the first dof  $U_1$  by modifying the load accordingly.



Thus the physical problem being solved after the first elimination is

$$\begin{pmatrix} \frac{14}{5} & -\frac{16}{5} & 1 \\ -\frac{16}{5} & \frac{29}{5} & -4 \\ 1 & -4 & 5 \end{pmatrix} \begin{Bmatrix} U_2 \\ U_3 \\ U_4 \end{Bmatrix} = \begin{Bmatrix} 1 \\ 0 \\ 0 \end{Bmatrix}$$

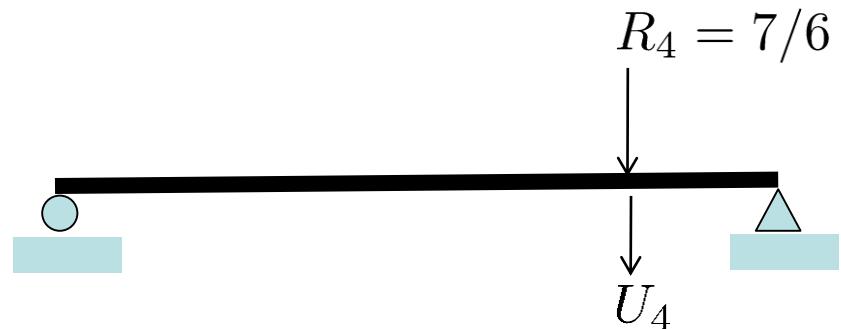
The physical problem now has 3 dof's.



The next elimination results in:

$$\begin{pmatrix} \frac{15}{7} & -\frac{20}{7} \\ -\frac{20}{7} & \frac{65}{14} \end{pmatrix} \begin{Bmatrix} U_3 \\ U_4 \end{Bmatrix} = \begin{Bmatrix} \frac{8}{7} \\ -\frac{5}{14} \end{Bmatrix}$$

Note now that a different loading has to be applied to maintain the equivalence with the original problem.



Finally we get

$$\frac{5}{6}U_4 = \frac{7}{6}.$$

Thus, from a physical standpoint, we can argue that we should not encounter a zero diagonal element while performing a Gauss elimination on a FE stiffness matrix. This is because the  $i$  th diagonal at any stage of elimination represents the stiffness of the structure when  $(i - 1)$  dof's have been released. This stiffness should be positive if the structure is stable.

# The numerics of Gauss elimination

In this chapter we will look at methods to solve the linear system of simultaneous equations

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where, the bold capitals denote matrices and bold small symbols denote vectors.

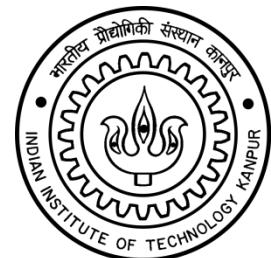
It is particularly easy to solve a  $n \times n$  *upper triangular system* like

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + a_{14}x_4 \cdots + a_{1n}x_n &= b_1 \\ a_{22}x_2 + a_{23}x_3 + a_{24}x_4 \cdots + a_{2n}x_n &= b_2 \\ a_{33}x_3 + a_{34}x_4 \cdots + a_{3n}x_n &= b_3 \\ &\vdots \\ a_{nn}x_n &= b_n \end{aligned}$$

Clearly,

$$x_n = \frac{b_n}{a_{nn}}.$$

*Applied Numerical Methods*



Given that  $a_{jj} \neq 0$ , for  $j = n - 1, \dots, 1$ , we can perform *backward substitution* to obtain:

$$x_j = \frac{1}{a_{jj}} \left( b_j - \sum_{k=j+1}^n a_{jk} x_k \right).$$

To count the number of operations, we assume that multiplication/division requires the longest time (though, division takes somewhat longer than multiplication in a computer) and addition/subtraction takes short times. In each step, we have  $n - j$  multiplications and 1 division. Thus, for the entire back substitution process, we have

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

multiplications/divisions. We also have  $n - j - 1$  additions and 1 subtraction per step, which leads to

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

addition/subtraction operations. The number of arithmetic operations in back substitution is:

$$\frac{n^2}{2} + O(n)$$



In the *Gauss elimination method*, a series of “elimination” steps transforms the system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  to an upper triangular system  $\mathbf{R}\mathbf{x} = \mathbf{c}$ , which is mathematically equivalent (may not be exactly equivalent due to round off errors) to the original system. Two operations *permutation of rows* and *scaling of rows* are performed in the elimination process. The process goes through the following steps:

- Form the augmented matrix

$$[\mathbf{A}|\mathbf{b}] .$$

- Set

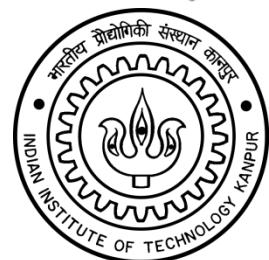
$$\mathbf{A}^0 = \mathbf{A}, \mathbf{b}^0 = \mathbf{b}$$

- For  $r \in 1, \dots, n$ , find  $a_{r1} \neq 0$ . Such an element must exist as otherwise  $\mathbf{A}$  will be singular. Permute 1st and  $r$  th rows. Let the new augmented matrix be

$$[\tilde{\mathbf{A}}^0|\tilde{\mathbf{b}}^0] .$$

- For  $j = 2, \dots, n$ , multiply row 1 by  $q_{j1}$  and subtract the result from  $j$  th row, where,

$$q_{j1} = \frac{\tilde{a}_{j1}^0}{\tilde{a}_{11}^0} \quad \text{Applied Numerical Methods} \quad b_j^1 = \tilde{b}_j^0 - q_{j1} \tilde{b}_1^0.$$



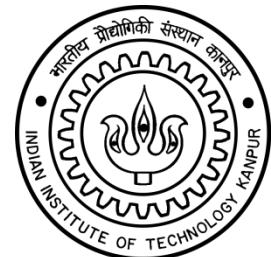
After the last step, the result is

19

$$[\mathbf{A}^1 | \mathbf{b}^1] = \left( \begin{array}{cccc|c} \tilde{a}_{11}^0 & \tilde{a}_{12}^0 & \dots & \tilde{a}_{1n}^0 & b_1^0 \\ 0 & a_{21}^1 & \dots & a_{2n}^1 & b_2^1 \\ 0 & \vdots & \vdots & \vdots & \vdots \\ 0 & a_{n1}^1 & \dots & a_{nn}^1 & b_n^1 \end{array} \right)$$

We can express the transformation  $[\mathbf{A}|\mathbf{b}] \rightarrow [\tilde{\mathbf{A}}^0|\tilde{\mathbf{b}}^0]$  as  $[\tilde{\mathbf{A}}^0|\tilde{\mathbf{b}}^0] = \mathbf{P}_1[\mathbf{A}|\mathbf{b}]$ , where  $\mathbf{P}$  is a *permutation matrix*. This matrix has the form

$$\mathbf{P} = \begin{pmatrix} \mathbf{1} & & & & & \mathbf{r} \\ & \mathbf{1} & & & & \\ & & \begin{matrix} 0 & 0 & 0 & \dots & 0 & 1 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots & \dots & \dots & 0 \\ & & & \ddots & & \vdots & \dots & 0 \\ & & & & & 0 & \dots & 0 \\ \mathbf{r} & & 1 & \dots & & & & \\ & & & & \vdots & & \ddots & \vdots \\ & & 0 & & & & & 1 \end{matrix} \end{pmatrix}$$



Similarly, we can represent  $[A^1|b^1] = G_1[\tilde{A}^0|\tilde{b}^0]$ , where  $G_1$  is called a *Frobenius matrix* and is given by:

$$\begin{pmatrix} 1 & & & \\ -q_{21} & 1 & & \\ \vdots & & \ddots & \\ -q_{n1} & \dots & & 1 \end{pmatrix}$$

It can be shown that  $\det P_1 = \det G_1 = 1$  and  $P_1^{-1} = P_1$ . The inverse of  $G_1$  is

$$G_1^{-1} = \begin{pmatrix} 1 & & & \\ q_{21} & 1 & & \\ \vdots & & \ddots & \\ q_{n1} & \dots & & 1 \end{pmatrix}$$

Clearly,  $Ax = b$  has the same solution as  $A^1x = b^1$  as

$$G_1 P_1 A x = G_1 P_1 b = b^1.$$

The series of steps

$$[\mathbf{A}|\mathbf{b}] \rightarrow [\mathbf{A}^1|\mathbf{b}^1] \rightarrow \cdots \rightarrow [\mathbf{A}^{n-1}|\mathbf{b}^{n-1}] \rightarrow [\mathbf{R}|\mathbf{c}],$$

leads to an upper triangular matrix  $\mathbf{R}$ , where

$$[\mathbf{R}|\mathbf{c}] = \mathbf{G}_{n-1}\mathbf{P}_{n-1} \dots \mathbf{G}_1\mathbf{P}_1[\mathbf{A}|\mathbf{b}],$$

is equivalent to the original system. Here,

$$\mathbf{P}_i = \begin{pmatrix} 1 & & & & \\ & \boxed{i} & & \boxed{r} & \\ & \ddots & & & \\ & & 1 & & \\ & & & \vdots & 1 \\ & & & & \vdots \\ & & & \boxed{r} & \\ & & & 1 & \dots & 0 \\ & & & & & 1 \\ & & & & & & \ddots \\ & & & & & & & 1 \end{pmatrix} \quad \mathbf{G}_i = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & -q_{i+1,i} & 1 \\ & & & \vdots & \ddots \\ & & & -q_{ni} & & 1 \end{pmatrix}$$

The element  $a_{r1} = \tilde{a}_{11}^0$  is called the *pivot element* and is chosen so that

$$|a_{r1}| = \max_{1 \leq j \leq n} |a_{j1}|.$$

Permutation of rows is called *column pivoting*. Though *total pivoting* where columns are also permuted is possible, it is expensive and rarely used.

We can perform an operation count for Gauss elimination. For simplicity, we will not count the steps required for pivoting. Thus,

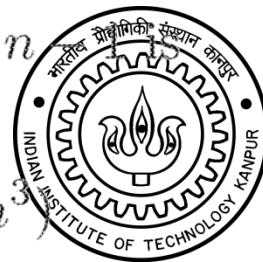
- The  $k$ th elimination step involves

$$a_{ij}^k = a_{ij}^{k-1} - \frac{a_{ik}^{k-1}}{a_{kk}^{k-1}} a_{kj}^{k-1}, \quad b_i^k = b_i^{k-1} - \frac{a_{ik}^{k-1}}{a_{kk}^{k-1}} b_k^{k-1}$$

for  $i, j = k \dots n$ .

- This involves  $n - k$  divisions to form  $a_{ik}^{k-1}/a_{kk}^{k-1}$
- followed by  $(n - k)(n - k + 1)$  multiplications.
- Total number of multiplications/divisions over  $k = 1, \dots, n$

$$\sum_{k=1}^{n-1} (n-k)(n-k+1) \frac{2n^3 + 3n^2 - 5n}{6} = O(n^3)$$



# Example

As an example, we solve the system

$$\begin{pmatrix} 3 & 1 & 6 \\ 2 & 1 & 3 \\ 1 & 1 & 1 \end{pmatrix} \begin{Bmatrix} x_1 \\ x_2 \\ x_3 \end{Bmatrix} = \begin{Bmatrix} 2 \\ 7 \\ 4 \end{Bmatrix}$$

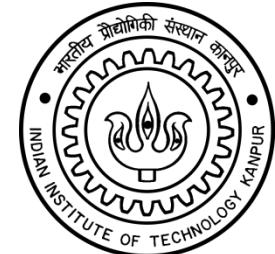
Step 1: pivoting

$$\left( \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 2 & 1 & 3 & 7 \\ 1 & 1 & 1 & 4 \end{array} \right)$$

Step 2: elimination

Instead of  $0_{q,j1}$  can be stored.

$$\left( \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 2/3 & 1/3 & -1 & 17/3 \\ 1/3 & 2/3 & -1 & 10/3 \end{array} \right)$$



Step 3: pivoting

$$\left( \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 1/3 & 2/3 & -1 & 10/3 \\ 2/3 & 1/3 & -1 & 17/3 \end{array} \right)$$

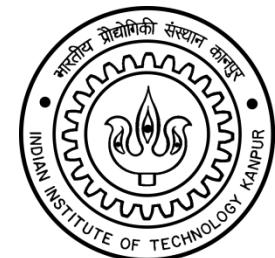
Step 4: elimination

$$\left( \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 1/3 & 2/3 & -1 & 10/3 \\ 2/3 & 1/2 & -1/2 & 4 \end{array} \right)$$

Step 5: back substitute

$$\begin{aligned} x_3 &= -8 \\ x_2 &= \frac{3}{2} \left( \frac{10}{3} - x_3 \right) = -7 \\ x_1 &= \frac{1}{3} (2 - x_2 - 6x_3) = 19 \end{aligned}$$

Solution also using the demo code `gauss_elim.m`



After the elimination is complete, the transformed matrix can be written as

25

$$\begin{pmatrix} 1 & 0 & 0 \\ 1/3 & 1 & 0 \\ 2/3 & 1/2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 & 6 \\ 0 & 2/3 & -1 \\ 0 & 0 & -1/2 \end{pmatrix} = \begin{pmatrix} 3 & 1 & 6 \\ 1 & 1 & 1 \\ 2 & 1 & 3 \end{pmatrix},$$

which is same as the matrix we started with (with two rows permuted). This is called a *LU* decomposition.

Solution also using the demo code LU.m

Formally,

$$L = \begin{pmatrix} l_{11} & & & 0 \\ l_{21} & 1 & & \\ \vdots & \ddots & \ddots & \\ l_{n1} & \dots & l_{n,n-1} & \end{pmatrix}, U = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ r_{22} & \dots & & r_{2n} \\ \ddots & & & \\ 0 & & & r_{nn} \end{pmatrix}$$

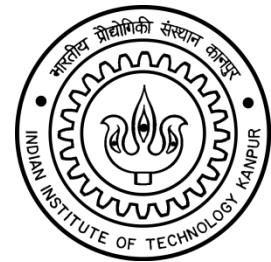
are the factor in a *LU* decomposition such that

$$LU = PA,$$

where 10/24/2018

Applied Numerical Methods

$$P = P_{n-1} \dots P_1.$$



In the  $i$  th elimination step, i.e.  $[A^{i-1}|b^{i-1}] \rightarrow [A^i|b^i]$ , we store the factors  $q_{i+1,i} \dots q_{n,i}$  of the Frobenius matrix  $G_i^{-1}$  in place of the zeros created. Thus, we work with the augmented matrix:

$$\left( \begin{array}{ccccccc|c} r_{11} & r_{12} & \dots & r_{1i} & r_{1,i+1} & \dots & r_{1n} & c_1 \\ \lambda_{21} & r_{22} & \dots & r_{2i} & r_{2,i+1} & \dots & r_{2n} & c_2 \\ \lambda_{31} & \lambda_{32} & \dots & r_{3i} & r_{3,i+1} & \dots & r_{3n} & c_3 \\ \vdots & \vdots & \ddots & & \vdots & & \vdots & \vdots \\ \lambda_{i1} & \lambda_{i2} & & r_{ii} & r_{i,i+1} & & r_{in} & c_i \\ \lambda_{i+1,1} & \lambda_{i+1,2} & & \lambda_{i+1,i} & a_{i+1,i+1}^i & \dots & a_{i+1,n}^i & b_{i+1}^i \\ \vdots & & & \vdots & \vdots & & \vdots & \vdots \\ \lambda_{n1} & \lambda_{n2} & \dots & \lambda_{ni} & a_{n,i+1}^i & \dots & a_{nn}^i & b_n^i \end{array} \right)$$

Here, the  $\lambda_{ki}$  are perturbations of  $q_{ki}$  and the matrix above finally leads to  $L$  and  $U$ . Once the decomposition is done, two sets of equations need to be solved to obtain the solution  $x$ :

$$Ly = Pb, \text{ and } Ux = y.$$

## Exercise

The importance of the pivoting process can be understood from the following example:

$$\begin{pmatrix} 10^{-4} & 1 \\ 1 & 1 \end{pmatrix} \begin{Bmatrix} x_1 \\ x_2 \end{Bmatrix} = \begin{Bmatrix} 1 \\ 2 \end{Bmatrix}$$

The exact solution to this is  $x_1 = 1.00010001$  and  $x_2 = 0.99989999$ . Let us first do this without pivoting. The elimination step leads to

$$\left( \begin{array}{cc|c} 0.1 \cdot 10^{-3} & 0.1 \cdot 10^1 & 0.1 \cdot 10^1 \\ 0 & -0.1 \cdot 10^5 & -0.1 \cdot 10^5 \end{array} \right).$$

After backward substitution, the solution is  $x_1 = 0, x_2 = 1$ . With pivoting, the two rows are first swapped leading to:

$$\left( \begin{array}{cc|c} 0.1 \cdot 10^1 & 0.1 \cdot 10^1 & 0.2 \cdot 10^1 \\ 0 & 0.1 \cdot 10^1 & 0.1 \cdot 10^1 \end{array} \right).$$

which gives the correct result  $x_1 = x_2 = 1$ .



Sometimes, even pivoting cannot lead us to the correct solution and a *scaling* is necessary. For example, multiplying the first row of the matrix in this example by a large number (say 20000), we get the system

$$\begin{pmatrix} 2 & 20000 \\ 1 & 1 \end{pmatrix} \begin{Bmatrix} x_1 \\ x_2 \end{Bmatrix} \begin{Bmatrix} 20000 \\ 2 \end{Bmatrix}$$

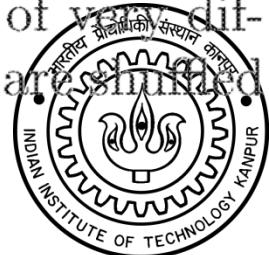
With or without pivoting, this system gives the incorrect solution  $x_1 = 0, x_2 = 1$ . However, we can apply a *equilibration step* to this system before solving it, by multiplying it with a diagonal matrix, i.e.

$$DAx = Db,$$

where

$$D_{ii} = \sum_{j=1}^n (|a_{ij}|)^{-1}.$$

An even better method for stabilising matrices that have elements of very different size is to apply a *total pivoting* where both rows and columns are ~~swapped~~ in a equilibration step.



# Computation of inverse by LU decomposition

In principle, the inverse of a square matrix  $\mathbf{A}$  can be computed by first LU decomposing  $\mathbf{PA}$ . Solution of the systems

$$\mathbf{Ly}^i = \mathbf{Pe}^i, \mathbf{Ux}^i = \mathbf{y}^i,$$

where,  $\mathbf{e}^i$  are the Cartesian basis vectors, yields the inverse  $\mathbf{A}^{-1} = [\mathbf{x}^1 \ \mathbf{x}^2 \ \dots \ \mathbf{x}^n]$ . The practical implementation is illustrated through an example.

# Exercise

Consider again the matrix

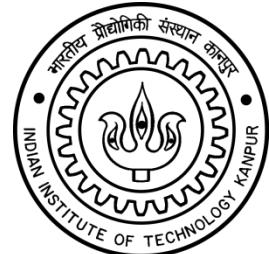
$$\begin{pmatrix} 3 & 1 & 6 \\ 2 & 1 & 3 \\ 1 & 1 & 1 \end{pmatrix}$$

Find its inverse. Step 1: pivoting

$$\left( \begin{array}{ccc|ccc} 3 & 1 & 6 & 1 & 0 & 0 \\ 2 & 1 & 3 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{array} \right)$$

Step 2: elimination

$$\left( \begin{array}{ccc|ccc} 3 & 1 & 6 & 1 & 0 & 0 \\ 0 & 1/3 & -1 & -2/3 & 0 & 1 \\ 0 & 2/3 & -1 & -1/3 & 1 & 0 \end{array} \right)$$



Step 3: Permute rows

$$\left( \begin{array}{ccc|ccc} 3 & 1 & 6 & 1 & 0 & 0 \\ 0 & 2/3 & -1 & -1/3 & 0 & 1 \\ 0 & 1/3 & -1 & -2/3 & 1 & 0 \end{array} \right)$$

Step 4: Next elimination

$$\left( \begin{array}{ccc|ccc} 3 & 1 & 6 & 1 & 0 & 0 \\ 0 & 2/3 & -1 & -1/3 & 0 & 1 \\ 0 & 0 & -1/2 & -1/2 & 1 & -1/2 \end{array} \right)$$

Step 4: Back substitution

$$\left( \begin{array}{ccc|ccc} 3 & 1 & 0 & -5 & 12 & -6 \\ 0 & 2/3 & 0 & 2/3 & -2 & 2 \\ 0 & 0 & -1/2 & -1/2 & 1 & -1/2 \end{array} \right)$$



Back substitution leads to the inverse matrix:

$$A^{-1} = \begin{pmatrix} -2 & 5 & -3 \\ 1 & -3 & 3 \\ 1 & -2 & 1 \end{pmatrix}$$



# Matrix norms and conditioning

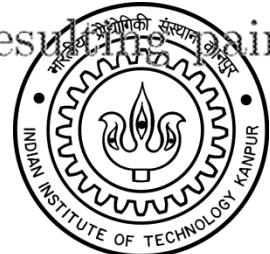
Assuming a basic exposure to linear algebra, we will quickly recapitulate some essential aspects. We have a vector space  $\mathbb{K}^n$  of all  $n$  dimensional vectors. That is, we deal with vectors like  $\mathbf{x} \in \mathbb{K}^n$  where,  $\mathbf{x}^T = \langle x_1 \ x_2 \ \dots \ x_n \rangle$ , where  $x_i$  are components of the vector with respect to a fixed Cartesian basis  $\mathbf{e}_i$ ,

$$\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i.$$

A mapping  $\|\cdot\| : \mathbb{K}^n \rightarrow \mathbb{R}$  is called a *vector norm* if it satisfies:

- *Definiteness*:  $\|\mathbf{x}\| \geq 0, \|\mathbf{x}\| = 0 \Rightarrow \mathbf{x} = 0, \forall \mathbf{x} \in \mathbb{K}^n$ .
- *Homogeneity*:  $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|, \alpha \in \mathbb{K}, \mathbf{x} \in \mathbb{K}^n$
- *Triangle inequality*:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|, \mathbf{x}, \mathbf{y} \in \mathbb{K}^n$

The norm can be defined over any vector space  $V$  over  $\mathbb{K}$ . The resulting pair  $\{V, \|\cdot\|\}$  is called a *normed vector space*



Examples of vector norms are the  $L_2$  norm

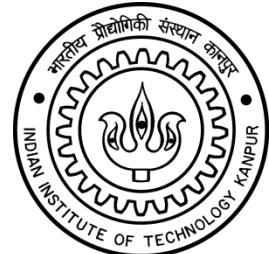
$$\|x\|_2 = \left( \sum_{i=1}^n x_i^2 \right)^{1/2},$$

or the more general  $L_p$  norm

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

More useful examples include the  $L_\infty$  norm and the  $L_1$  norm defined as

$$\|x\|_\infty = \max_{i=1,\dots,n} |x_i|, \|x\|_1 = \sum_{i=1}^n |x_i|.$$



A mapping  $(\mathbf{x}, \mathbf{y}) : \mathbb{K}^n \times \mathbb{K}^n \rightarrow \mathbb{K}$  is called a *scalar product* if it satisfies

- *Symmetry:*  $(\mathbf{x}, \mathbf{y}) = (\mathbf{y}, \mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{K}$
- *Linearity:*  $(\alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z}) = \alpha(\mathbf{x}, \mathbf{z}) + \beta(\mathbf{y}, \mathbf{z}), \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{K}^n, \alpha, \beta \in \mathbb{K}$
- *Definiteness:*  $(\mathbf{x}, \mathbf{x}) \in \mathbb{R}, (\mathbf{x}, \mathbf{x}) > 0 \quad \mathbf{x} \in \mathbb{K}^n / \{0\}.$

We will mostly use the Euclidean scalar product

$$(\mathbf{x}, \mathbf{y})_2 = \sum_{j=1}^n x_j y_j.$$

Two vectors are orthogonal if

$$(\mathbf{x}, \mathbf{y})_2 = 0.$$

A set of vectors in  $\{\mathbf{a}^1, \mathbf{a}^2 \dots \mathbf{a}^n\}$ , with  $\mathbf{a}^i \neq 0$  is called an *orthogonal basis* if  $(\mathbf{a}^i, \mathbf{a}^j)_2 = 0$  for  $i \neq j$ . Further, they are called an *orthonormal basis* if, additionally,  $(\mathbf{a}^k, \mathbf{a}^k) = 1$  for  $k = 1, \dots, n$ .



With an orthonormal basis in place, we can easily see that

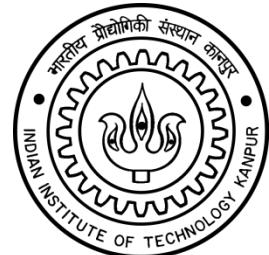
$$(\mathbf{x}, \mathbf{a}^i)_2 = \left( \sum \alpha_j \mathbf{a}^j, \mathbf{a}^i \right)_2 = \alpha_i.$$

This leads to the representation

$$\mathbf{x} = \sum_{i=1}^n (\mathbf{x}, \mathbf{a}^i)_2 \mathbf{a}^i.$$

Further, it can easily be seen that

$$\|\mathbf{x}\|_2^2 = (\mathbf{x}, \mathbf{x})_2 = \sum_{i=1}^n |(\mathbf{x}, \mathbf{a}^i)_2|^2.$$



Consider the vector space of all  $n \times n$  matrices  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . The norm induced by a vector norm  $\|\cdot\|$  is

$$\|\mathbf{A}\| = \max \left( \frac{\|\mathbf{Ax}\|}{\|x\|} : \|x\| \neq 0 \right),$$

for  $x \in \mathbb{K}^n$ . Through this norm, we are trying to determine the direction that is amplified most by  $\mathbf{A}$ .

A matrix norm satisfies

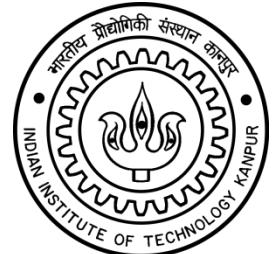
$$\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|x\|, x \in \mathbb{K}^n, \mathbf{A} \in \mathbb{K}^{n \times n}$$

and,

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|, \mathbf{A}, \mathbf{B} \in \mathbb{K}^{n \times n}.$$

Further, a *natural matrix norm* satisfies

$$\|\mathbf{I}\| = 1.$$



The natural matrix norm generated from the Euclidean norm is called the *spectral norm*. Note that it is not

$$\|A\| = \left( \sum_{j,k=1}^n |a_{jk}^2| \right)^{1/2},$$

as, for this norm,  $\|I\| = \sqrt{n}$ . If the eigenvalues of  $A^T A$  are  $\lambda$ , then the spectral norm is

$$\|A\|_2 = \max |\lambda|.$$

Now, we can do an error analysis of the linear system

$$Ax = b.$$

This implies that if the matrix  $A$  and vector  $b$  are faulty by  $\delta A$  and  $\delta b$  respectively, we wish to find the error  $\delta x$  in the solution of the system

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

with  $\tilde{A} = A + \delta A$ ,  $\tilde{b} = b + \delta b$  and  $\tilde{x} = x + \delta x$ .



We state the so-called *perturbation theorem* without proof. Let  $\mathbf{A} \in \mathbb{K}^{n \times n}$  be an invertible matrix such that  $\|\delta\mathbf{A}\| \leq \|\mathbf{A}^{-1}\|^{-1}$ , then, the perturbed matrix  $\tilde{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A}$  is also invertible and

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}\mathbf{A}}{1 - \text{cond}\mathbf{A}\|\delta\mathbf{A}\| \|\mathbf{A}\|} \left\{ \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right\},$$

where  $\text{cond}\mathbf{A} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ .

Further, it can be shown that for the spectral norm

$$\text{cond}_2 \mathbf{A} = \frac{|\lambda_{max}|}{|\lambda_{min}|},$$

where  $\lambda_{max}$  and  $\lambda_{min}$  are the largest and smallest eigenvalues of  $\mathbf{A}$ .

For cases where  $\text{cond}\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \ll 1$ , for  $\text{cond}\mathbf{A} \sim 10^s$ , and relative errors in  $\mathbf{A}$  and  $\mathbf{b} \sim 10^{-k}$ , the relative error in  $\mathbf{x}$  will be  $10^{s-k}$ .

