

**Numeric Methods in Civil Engineering**  
**Prof. Arghya Deb**  
**Department of Civil Engineering**  
**Indian Institute of Technology, Kharagpur**


**Lecture - 4**  
**Linear Systems – II**

Going to continue our discussion on linear systems, this is the second part of our lecture on linear systems. Last time we talked about singular value decomposition, when we talked about how to find the singular values of an  $m$  by  $n$  matrix  $A$ , which is of rank  $r$ .

(Refer Slide Time: 00:33)

**Singular Value Decomposition**

- Let  $A$  be a  $m \times n$  matrix of rank  $r$ . Then there exist:
- A  $m \times m$  orthogonal matrix  $U$  (i.e.  $U^T U = I$ )
- A  $n \times n$  orthogonal matrix  $V$  (i.e.  $V^T V = I$ )
- A  $r \times r$  diagonal matrix  $D$  with strictly positive elements, called the singular values of  $A$  such that:


$$A = U \Sigma V^T, \quad \Sigma = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}$$


If  $A$  is an  $m$  by  $n$  matrix of rank  $r$ . We can be assured that there exist a  $m$  by  $m$  orthogonal matrix  $U$ , orthogonal matrix naming that matrix, transpose of that matrix in product with itself, that is  $U$  transpose  $U$  is going to be the identity matrix. An  $n$  by  $n$  orthogonal matrix  $V$ , that is, which is also an orthogonal matrix, which means  $V$  transpose  $V$  is again equal to the identity matrix  $I$  and an  $r$  by  $r$  diagonal matrix  $D$ , which has got strictly positive elements and those positive elements are called the singular values of  $A$ , so all these matrices  $U$ ,  $V$  and  $D$ , when  $D$  is working as an expanded matrix  $\Sigma$ , they can be combined to form the matrix  $A$ . So,  $A$  is equal to  $U$  product,  $\Sigma$  product  $V$  transpose. So,  $U$  product  $\Sigma$  transpose  $V$  transpose is the singular value decomposition of the matrix  $A$ .

(Refer Slide Time: 01:53)

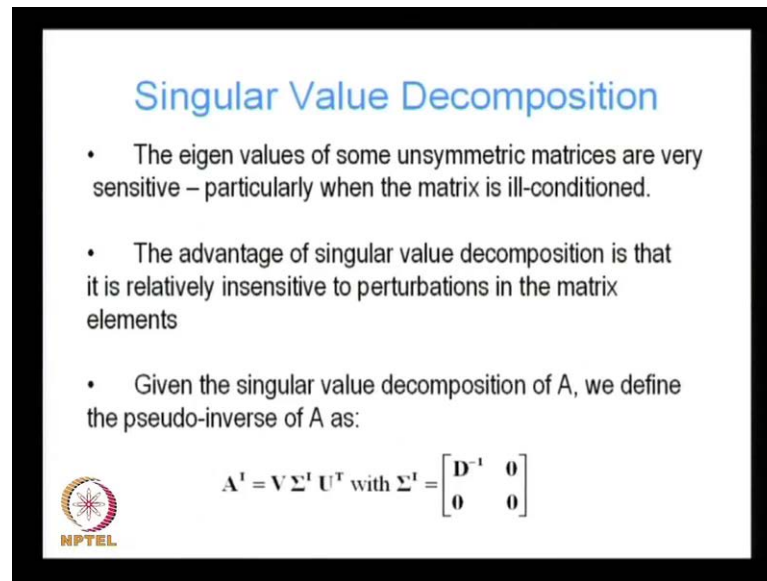
### Singular Value Decomposition

- Obviously in the most general case,  $\Sigma$  has to be a  $m \times n$  matrix
- If  $r = m$ , then  $\Sigma = \begin{bmatrix} \mathbf{D} & \mathbf{0} \end{bmatrix}$
- If  $r = n$ , then  $\Sigma = \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix}$
- If  $r = m = n$ , then  $\Sigma = \mathbf{D}$
- Finally,  $\mathbf{A}\mathbf{A}^T = \mathbf{U} \Sigma^2 \mathbf{U}^T$ ,  $\mathbf{A}^T \mathbf{A} = \mathbf{V} \Sigma^2 \mathbf{V}^T$




In the most general case, sigma has to be a m by n matrix; however, if the rank of the matrix is equal to m, that is the rank of the matrix is equal to the number of rows, in the matrix, then sigma is equals to D 0 right, however if the rank of the matrix is equal to n. The n is the column, number of columns in my original matrix, and then sigma is equal to D 0, where D is now in the form of a column vector. However, if sigma has rank r, which is equal to m, which is equal to n, which means that, A is the n by m matrix is a square matrix and it has got full rank then sigma is just equal to D. That is the matrix is no longer singular. It has got all its eigen values are positive real numbers. In that case, we can write. So, given the singular value decomposition of A, we can write AA transpose, as U sigma square U transpose similarly, by taking the transpose of that, we can write A transpose A is equal to V sigma square V transpose.

(Refer Slide Time: 03:20)



**Singular Value Decomposition**

- The eigen values of some unsymmetric matrices are very sensitive – particularly when the matrix is ill-conditioned.
- The advantage of singular value decomposition is that it is relatively insensitive to perturbations in the matrix elements
- Given the singular value decomposition of A, we define the pseudo-inverse of A as:

$$A^{\dagger} = V \Sigma^{\dagger} U^T \text{ with } \Sigma^{\dagger} = \begin{bmatrix} D^{-1} & 0 \\ 0 & 0 \end{bmatrix}$$



Why do we need singular value decomposition? The reason is, because many matrices particularly unsymmetric matrices their eigen values are extremely sensitive to small numerical perturbations. So, small numerical perturbations can lead to large change in the eigen values, such matrices are typically known as ill conditioned matrices, for such matrices singular value decomposition is an efficient way to find eigen values of that matrix.

Given the singular value decomposition of A, we defined the pseudo-inverse of A,  $A^{\dagger}$  note that, this is different from  $A^{-1}$ , where  $A^{-1}$  is the inverse. So,  $A^{\dagger}$  is the pseudo-inverse, the pseudo-inverse is again defined in terms of the orthogonal matrices V and U, as well as the diagonal matrix sigma. So, it can be written as the pseudo-inverse can be written as  $V \Sigma^{\dagger} U^T$  with  $\Sigma^{\dagger}$  is the pseudo-inverse of sigma, which is equal to  $D^{-1}$ , 0, 0, 0.

(Refer Slide Time: 04:37)

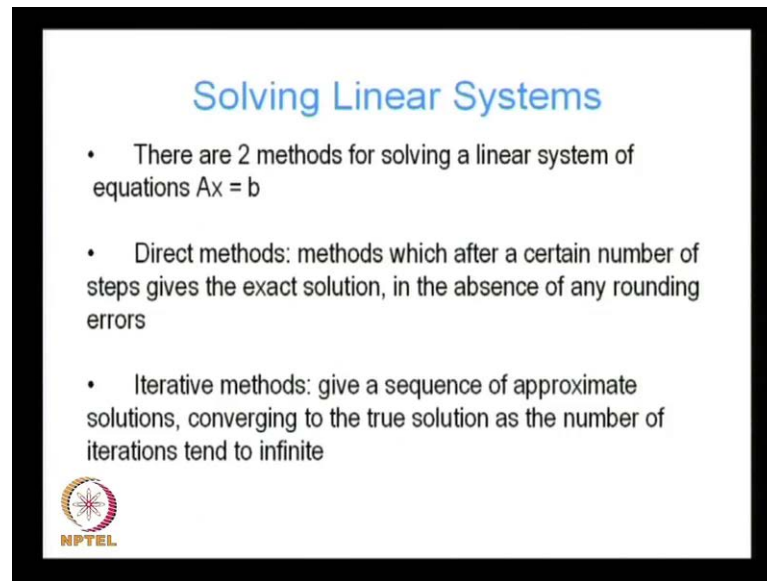
### Singular Value Decomposition

- Obviously in the most general case,  $\Sigma$  has to be a  $m \times n$  matrix
- If  $r = m$ , then  $\Sigma = \begin{bmatrix} \mathbf{D} & \mathbf{0} \end{bmatrix}$
- If  $r = n$ , then  $\Sigma = \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix}$
- If  $r = m = n$ , then  $\Sigma = \mathbf{D}$
- Finally,  $\mathbf{A}\mathbf{A}^T = \mathbf{U} \Sigma^2 \mathbf{U}^T$ ,  $\mathbf{A}^T\mathbf{A} = \mathbf{V} \Sigma^2 \mathbf{V}^T$




Let us recall, what was the form of sigma? So, sigma in its most general case had this form. So, the pseudo-inverse of sigma is D inverse, 0, 0, 0 in reality sigma does not have an inverse right, because this matrix has got 0 elements on its diagonal, which is singular matrix. It cannot be inverted, that is why we denote sigma I and call it pseudo-inverse rather than inverse. So, AI is the pseudo-inverse of A and we can write it as V sigma I U transpose with sigma I denoted as follows.

(Refer Slide Time: 05:30)



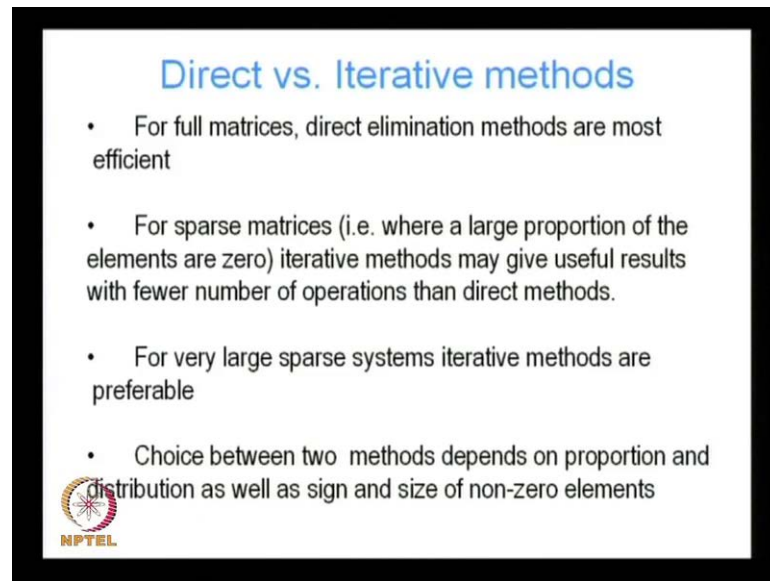
**Solving Linear Systems**

- There are 2 methods for solving a linear system of equations  $Ax = b$
- Direct methods: methods which after a certain number of steps gives the exact solution, in the absence of any rounding errors
- Iterative methods: give a sequence of approximate solutions, converging to the true solution as the number of iterations tend to infinite




So, that was our discussion on solving as on singular value decomposition. As we discussed, it is very useful for solving problems, which are ill conditioned coefficient matrices. Let us, now go back to a discussion on linear systems and talk about certain general procedures for solving linear systems. They are basically two main methods for solving linear systems, one method is known as the direct method, the other method is known as the iterative method. In direct methods, we are assured that after a certain number of steps we are going to get the exact solution, in the absence of any round out and in case of iterative methods, we get a sequence of approximate solutions and if we continue the sequence to infinite number of if we continue the sequence in infinite number of times, we are going to get the exact solution or the true solution.

(Refer Slide Time: 06:37)



**Direct vs. Iterative methods**

- For full matrices, direct elimination methods are most efficient
- For sparse matrices (i.e. where a large proportion of the elements are zero) iterative methods may give useful results with fewer number of operations than direct methods.
- For very large sparse systems iterative methods are preferable
- Choice between two methods depends on proportion and distribution as well as sign and size of non-zero elements

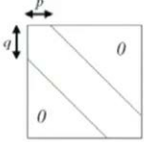



For, full matrices direct elimination methods are the most efficient; however, for sparse matrices, that is where a large proportion of the elements are 0, iterative methods make if useful results with the fewer number of operations than direct methods. For, very large sparse systems iterative methods are preferable. Therefore, what we get from this discussion is that, the decision on when to go for direct methods or when to choose iterative methods, typically depends on this sparsity of the matrix that is the number of non-zero elements in the matrix. So, depends on the proportion, distribution, as well as the sign and the size of non-zero elements.

(Refer Slide Time: 07:31)

### Sparsity and Band Width

- Whether a matrix is full or sparse can be characterized by a property known as the “band width” of the matrix.
- A matrix  $A$  for which  $a_{ij} = 0$  if  $j > i+p$  or  $i > j+q$  is called a banded matrix with band width  $w = p+q+1$
- The number of non-zero elements in any row/column  $< w$
- If  $A$  is  $n \times n$ , total number of non zero elements  $< n.w$




Whether, a matrix is full or sparse can be characterized by a property known as the band width of a matrix, a matrix  $A$  for which  $a_{ij}$  is equal to 0, if  $j$  is greater than  $i$  plus  $p$ , it means that, all if we consider a particular row, then all elements in that row, which are more than  $p$  elements away from the principle diagonal are 0.

Similarly, if  $i$  is greater than  $j$  plus  $q$ , which basically means that, all elements in a column, which are more than  $q$  elements away from the principle diagonal are 0. So, if  $A$  is a finite strip bounding the principle diagonal in the neighbourhood of the principle diagonal. The matrix as non-zero elements, while everywhere else; the elements of the matrix are 0 and in that case, we say that the band width of the matrix is  $W$ , the  $W$  is equal to  $p$  plus  $q$  plus 1. So, the band width of the matrix is  $p$  plus  $q$  plus 1 denoting the principle diagonal. So, this is the band width of a matrix. The number of this automatically gives us the total number of non-zero elements in the matrix, which has to be less than  $n$  dotted with  $W$ .

(Refer Slide Time: 09:11)

### Band width reduction

- Sometimes it is possible to re-label the unknown variables to reduce the band width of a matrix
- Example: Given the full system:
$$\begin{bmatrix} 1 & 0 & 3 \\ 13 & 6 & 5 \\ 0 & 2 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 12 \\ 3 \\ 4 \end{bmatrix}$$
- The original system has a full coefficient matrix



Sometimes, it is possible to re-label the unknown variables in a matrix, to reduce the band width of a matrix. Consider this 3 by 3 system it is clear, where this matrix has full band width that is, its band width is 3 for instance in the first row, in first row the third column is non-zero. So, it has got full band width. So, the band width of this matrix is 3. So, it has got a full coefficient matrix; however, we will see that by re-labelling the variables it is possible to reduce the band width of this matrix.

How is that possible? Well, if we re-label  $x_2$  as  $x_3$  and  $x_3$  as  $x_2$ , we get a band limited system with a band width of 2 instead of 3. Now, you can see that the maximum number of non-zero elements in a row is two elements away from the principle diagonal, this is the last element are 0, now similarly for the row for the column. So, this is my principle diagonal, we have one non-zero element away from the principle diagonal, beyond that we have 0 elements.

So, now the band width of the matrix is two rather than three, algorithms to reduce band width. It can be used for large matrices, this is because solution of band limited matrices is much less expensive than full matrices for example, there is the Cuthill-McKee algorithm, which is a widely used algorithm for the numbering, the equations in order to reduce the band width.




Similarly, there are more advanced versions of the Cuthill-McKee algorithm there are several other algorithms, which can be used for band width minimization. One thing we note is that even, if an original matrix is band limited its inverse may be full therefore, if we work on the original matrix and we use a band width minimiser and the original matrix and reduce its band width, but then we try to invert it we might end up with the full matrix. So, all work on band width minimization would be of no use, because eventually we have to use that inverted matrix to solve our system of equations and if the inverted matrix is full that is going to be as expensive. It is not going to reduce the computational expense at all.

(Refer Slide Time: 12:05)

### Solving Triangular systems

- Linear systems in which the coefficient matrix is triangular are particularly easy to solve
- Recall, an upper triangular system has the form:
 
$$\begin{aligned} u_{11}x_1 + u_{12}x_2 + \dots + u_{1n}x_n &= b_1 \\ u_{22}x_2 + \dots + u_{2n}x_n &= b_2 \\ &\dots \\ u_{nn}x_n &= b_n \end{aligned}$$
- Assuming diagonal elements  $u_{11}, u_{22}, \dots, u_{nn}$  are non-zero, system can be solved sequentially using back substitution:
 
$$x_n = \frac{b_n}{u_{nn}}, x_{n-1} = \frac{b_{n-1} - u_{n-1,n}x_n}{u_{n-1,n-1}}, \dots$$

$$x_i = \frac{b_i - \sum_{j=i+1}^n u_{ij}x_j}{u_{ii}}, \quad i = n, n-1, \dots, 1$$



Next, we would like to look at particular types of matrix, which are very useful with for solving linear equations, which give us a form, which is very useful for solving linear equations and these are known as triangular systems. Recall an upper triangular system has the following form:  $U_{11}x_1 + U_{12}x_2 + \dots + U_{1n}x_n = b_1$ , which is the first equation in my system. The second equation has  $U_{22}x_2 + U_{23}x_3 + \dots + U_{2n}x_n = b_2$  and so on. So, forth until the last equation is just  $U_{nn}x_n = b_n$ . So, it is clear that in upper triangular system have all elements below the principle diagonal 0, an all elements at or about the principle diagonal non-zero. So, this system can be solved very simply for instance. We can easily invert the last equation, we can


solve for  $x_n$  as  $b_n$  by  $u_{nn}$ . So, it just one division which gives  $x_n$ .

So, once I know  $x_n$  I can use the  $n$  minus 1th equation to solve for  $x_{n-1}$ , because  $n$  minus 1th equation has the form,  $u_{n-1, n-1} x_{n-1} + u_{n-1, n} x_n$  is equal to  $b_{n-1}$ . So, if we do  $x_n$  we can solve for  $x_{n-1}$  using this simple equation just by one division, one addition, and one multiplication, and one addition. Similarly, we can start from the end we go progressively backwards and we can use this general formula  $x_i$  is equal to  $b_i$  minus  $\sum_{j=i+1}^n u_{ij} x_j$ ,  $j$  is sum do where  $i$  plus 1 to  $n$  divided by  $u_{ii}$  for  $i$  is equal to  $n$  through one to solve this entire system. So, we start with the simplest equation and worked have a backwards to solve for the remaining variables.

(Refer Slide Time: 14:31)

### Solving Triangular systems

- A lower triangular system can be solved very similarly.
- Recall, a lower triangular system has the form:
 
$$\begin{aligned} l_{11}x_1 &= b_1 \\ l_{21}x_1 + l_{22}x_2 &= b_2 \\ &\dots\dots\dots \\ l_{n1}x_1 + l_{n2}x_2 + \dots + l_{nn}x_n &= b_n \end{aligned}$$
- Assuming diagonal elements  $l_{11}, l_{22}, \dots, l_{nn}$  are non-zero, system can be solve sequentially using forward substitution:
 
$$x_i = \frac{b_i - \sum_{j=1}^{i-1} l_{ij}x_j}{l_{ii}}, \quad i = 1, 2, \dots, n$$




Similarly, if we have a lower triangular system, that can also be solved very easily. The lower triangular system is one, which all elements at or below the principle diagonal is non-zero that all elements above the principle diagonal are 0. In this case, the last case was backwards substitution, because we solve the last equation and worked have a backwards. While in this case, we solve the first equation, which is easy enough to solve  $x_1$  is equal to  $b_1$  by  $l_{11}$  and work have a forward. We solve for  $x_1$ , we solve for  $x_2$ , we solve for  $x_3$  and so on. So, forth, but important thing to note is that, this assumes that the

diagonal elements  $l_{11}, l_{22}$  through  $l_{nn}$  are non-zero, only when these diagonal elements are non-zero, can the system be solved sequentially using either forward substitution or backward substitution.

(Refer Slide Time: 15:45)

### Operation count

- For both back substitution and forward substitution we need to perform ' $n$ ' divisions and
 
$$\sum_{i=1}^n (i-1) = \frac{1}{2}n(n-1) = \frac{1}{2}n^2 - n \approx \frac{1}{2}n^2 \text{ (for large } n\text{)}$$
 additions and multiplications
- The no. of operations is thus almost identical to the no. of operations required to multiply a vector by a triangular matrix
- One of the most widely used methods of direct solution, Gaussian elimination therefore works by reducing a "full" system to a triangular system



For, both back substitution and forward substitution, we need to perform  $n$  divisions. Why do we need to perform  $n$  divisions? So, that for each equation, we have to perform one division. So,  $n$  divisions and we have to perform for each equation  $\sum_{i=1}^n (i-1)$  additions or multiplications basically let us take, a look at our equations and I think that would make it clear. For instance, for solving the first equation, we just have to perform one division. For solving the second equation, we have to perform one multiplication and one subtraction and one division. Similarly, for solving the third equation, we have to perform one division, but then we have to perform two additions and subtractions followed by two multiplications and two additions and subtractions.

So, the total number of additions and multiplications, for the first equation, we have for  $i$  is equal to 1, we do not have any multiplications additions, so 0 additions and multiplications. For the second equation, we have  $i$  is equal to 2, so one addition and one multiplication. For the third equation, we have  $i$  is equal to 3, so two additions and two multiplications and so on. If, we keep adding the number of additions and

multiplications, we will find the total number of additions and multiplications is equal to  $1/2 n^2 - n$ .


And if a large systems,  $n^2$  is going to be much larger than  $n$ . So,  $1/2 n^2 - n$  can be written as approximately equal to  $1/2 n^2$ . Therefore, backward substitution and forward substitution, we need to perform  $1/2 n^2$  additions and multiplications. But since, and the number of divisions is  $n$ , but again  $n$  is comparatively much smaller than  $1/2 n^2$ . So, the total operation cost approximately is  $1/2 n^2$ . The number of operations  $1/2 n^2$  is thus almost identical to the number of operations required to multiply a vector by a triangular matrix. Why that is the triangular matrix has only  $1/2 n^2$  elements right. So, if I multiply a triangular matrix with the vector. I have to perform  $1/2 n^2$  operations. So, the number of operations is identical to the number of operations required to multiply a vector by a triangular matrix.

One of why do we talk, so much about triangular systems, because to solve an equations using Gaussian elimination, triangular systems are probably the most efficient way of solving it. So, the basic idea is that we try to reduce the system to a triangular system. So, the first part in any Gaussian elimination is attempts to reduce a full system to a triangular system and once it has reduced it to a triangular system it uses forward substitution or backward substitution to solve the system.

(Refer Slide Time: 19:27)

**Gaussian elimination**

- The idea is to eliminate the unknowns in a systematic way so that we end up with a triangular system
- Let us consider the full linear system:  
$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$
$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$
$$\dots\dots\dots$$
$$a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n$$
- Let us assume the coefficient matrix A is non-singular and hence the system has a unique solution.




So, the idea is to eliminate the unknowns in a systematic way. So, that we end up with a triangular system. Let us consider the full linear system, which is an n by n system with the coefficient matrix is of size n by n,  $a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$ . Similarly,  $a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$  and so on. So, forth let us also assume that, the coefficient matrix A is non singular and hence, the system has a unique solution. We recall our discussion from last lecture, when we talked about the conditions, which are necessary for a system of equations to be non singular.

(Refer Slide Time: 20:21)

**Gaussian elimination: reduction**

- Assume  $a_{11}$  is not equal to zero
- Eliminate  $x_1$  from the last  $n-1$  equations by subtracting from the  $i^{\text{th}}$  equation the multiple  $m_{i1}$  times the first equation where  $m_{i1}$  are given by:
$$m_{i1} = \frac{a_{i1}}{a_{11}}, \quad i = 2, 3, \dots, n$$
- Then the last  $n-1$  equations becomes:
$$\begin{aligned} a_{22}^{(2)}x_2 + \dots + a_{2n}^{(2)}x_n &= b_2^{(2)} \\ \dots & \\ a_{n2}^{(2)}x_2 + \dots + a_{nn}^{(2)}x_n &= b_n^{(2)} \end{aligned}$$



We also assume that  $a_{11}$  is not equal to 0, if we eliminate  $x_1$  from the last  $n$  minus 1 equations by subtracting from each equation. The multiple  $m_{i1}$  times the first equation, where  $m_{i1}$  is given by  $a_{i1}$  by  $a_{11}$  then, the last  $n$  minus 1 equations do not have any variable  $x_1$ . So, the size of the system gets reduce to  $n$  minus 1 by  $n$  minus 1.


Let's go back and take a look at a full system and see what, we are talking about. So, what we are saying is that, the first equation, we divide all the elements by  $a_{11}$  then, we add  $a_{21}$  times the first equation. We subtract  $a_{21}$  times the first equation, from the second equation. We subtract  $a_{31}$  times the first equation, from the third equation. We subtract  $a_{n1}$  times the first equation from the  $n$ th equation. If we do that the first coefficient always become 0, because we are dividing the first equation by  $a_{11}$ .

So, the coefficient of  $x_1$  in the first equation is 1. If, we multiply this by  $a_{21}$  and subtracted from the second equation, the coefficient of  $x_1$  is going to vanish. Similarly, we do the same thing for the third, fourth, fifth and so on so forth until the  $n$ th equations. If, we do this then, we are going to get the last  $n$  minus 1 equation and no longer going to have a  $x_1$  term. So, the coefficient to the  $x_1$  in the last  $n$  minus 1 equation, is going to become is 0 and we are going to get that equation. We are going to get  $n$  minus one equation in terms of  $n$  minus 1 variable  $x_2$  through  $x_n$ .

(Refer Slide Time: 22:28)

**Gaussian elimination: reduction**

- The new coefficients are given by:  
$$a_{ij}^{(2)} = a_{ij} - m_{i1}a_{1j}, \quad b_i^{(2)} = b_i - m_{i1}b_1, \quad i = 2, 3, \dots, n$$
- The reduced system now has  $n-1$  equations and  $n-1$  unknowns
- If  $a_{22}^{(2)} \neq 0$  we can also eliminate  $x_2$  from the last  $n-2$  equations by subtracting from the  $i^{\text{th}}$  equation the multiple  $m_{i2}$  times the second equation, where  $m_{i2}$  are given by:  
$$m_{i2} = \frac{a_{i2}^{(2)}}{a_{22}^{(2)}}, \quad i = 3, 4, \dots, n$$




So, the new coefficients are given by this, it is just in short form in symbolically. So, the reduced system now has  $n$  minus 1 equation and  $n$  minus 1 unknown. If, at the second stage after the first reduction, we are left with the  $n$  minus 1 system and if it turns out that this  $a_{22}$  is also non-zero then, we can continue with the same process throughout. So then, what we are going to do? We are going to divide this equation by  $a_{22}$ . So, the first coefficient is going to become 1 and then we are going to multiply this equation by the first coefficient of  $x_2$ , in the second equation and subtracted from the second equation.

So, in that case, the second equation is going to have 0 as the coefficient of  $x_2$ . Similarly, we will do the same thing for the third equation and we will continue for up till the  $n^{\text{th}}$  equation. So, if we do that, if we can also eliminate  $x_2$  from the last  $n$  minus 2 equations by subtracting from the  $i^{\text{th}}$  equation, the multiple  $m_{i2}$  times the second equation, where  $m_{i2}$  is given by  $a_{i2}^{(2)}$  by  $a_{22}^{(2)}$ , where  $i$  is equal to 2, 3, 4 through  $n$ . The subscript 2 denotes the step in the reduction process. So, 1 denote the first step, 2 denotes the second step.

(Refer Slide Time: 24:22)

### Gaussian elimination: reduction

- Note that the superscript (2) in  $a_{i2}^{(2)}$  and  $a_{22}^{(2)}$  indicate that these are the coefficients of the transformed system after the first elimination, not the coefficients of the original set of equations
- Again after elimination of  $x_2$  we get a reduced system of  $(n-2)$  equations in  $(n-2)$  unknowns, whose coefficients are given by:


$$a_{ij}^{(3)} = a_{ij}^{(2)} - m_{i2}a_{2j}^{(2)}, \quad b_i^{(3)} = b_i - m_{i2}b_2^{(2)}, \quad i = 3, 4 \dots n$$


These are the subscript 2 in a i2 2 and a 22 2 indicates that these are the coefficients of the transformed system, after the first elimination. Similarly, we can continue this process, for the thirds for the third step and so on through the end steps. Until, we get a system like this, which we will call is an upper triangular system.

(Refer Slide Time: 24:45)

### Gaussian elimination: reduction

- The elements  $a_{11}, a_{22}, a_{33}$ , which are used as the denominator in the scale factors  $m_{i1}, m_{i2}, m_{i3}$  are called the pivotal elements.
- If these elements are non-zero, the elimination can be continued until after  $(n-1)$  steps we get the single equation:
 
$$a_{nn}^{(n)} x_n = b_n^{(n)}$$
- Finally collecting the first equation from each step in the elimination, we get an upper triangular system:

$$\begin{aligned} a_{11}^{(1)} x_1 + a_{12}^{(2)} x_2 + \dots + a_{1n}^{(n)} x_n &= b_1^{(1)} \\ a_{22}^{(2)} x_2 + \dots + a_{2n}^{(2)} x_n &= b_2^{(2)} \\ \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots & \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \\ \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots & \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \end{aligned}$$





So, from the second equation, we have made sure that the coefficient of  $x_1$  is 0. From the third equation, we have made sure that the coefficient of  $x_1$  and  $x_2$  are 0. Similarly, until we have continued with that process until, we have reach the  $n$ th equation, where we have made sure that coefficients of  $x_1$  through  $x_{n-1}$  are 0 and the only variable with non-zero coefficients is  $x_n$ .

So, we have reduce this system to an upper triangular system, when we have got a single equation  $a_{nn}x_n = b_n$  again, we can solve this equation by back substitution. We solve for  $x_n$  substitute that in the  $(n-1)$ th equation to get  $x_{n-1}$ , substitute  $x_n$  and  $x_{n-1}$  in  $(n-2)$ th equation to get  $x_{n-2}$  and so on. So, forth until, we get all the variables  $x_1, x_2, x_3$  to through  $x_n$ .

(Refer Slide Time: 26:05)

### Gaussian elimination: reduction

- In the above  $a_{ij}^{(k)}$  and  $b_i^{(k)}$  refer to coefficients and right hand values for the original system.
- The right hand side of  $b$  is transformed exactly in the same way as the columns of  $A$
- Hence the description of the elimination is simplified if  $A$  is augmented to include  $b$  as its  $(n+1)$ th column i.e.

$$a_{i,n+1} = b_i$$



So, in the above  $a_{ij}$  and  $b_i$  refer to coefficients and right hand values for the original system,  $b_i$  refers to the right hand values, after the first step in the reduction procedure. Similarly,  $b_i^{(n)}$  would refer to the right hand values after the  $n$ th step in the reduction procedure. The right hand side of  $b$  is transformed exactly in the same way as the columns of  $A$ , hence the description of the elimination is simplified. If,  $A$  is augmented to include  $b$  as it is,  $(n+1)$ th column. Basically, we move the right hand side to the left hand side with of course, with the negative sign and that resulting system is called an

augmented system.

(Refer Slide Time: 27:07)

### Augmenting the no. of rows

- If it is required to solve several systems of equation with the same coefficient matrix A:  
$$Ax_1 = b_1, \quad Ax_2 = b_2, \dots, Ax_n = b_n$$
they can be treated similarly by augmenting A by the additional p columns  $b_1, b_2, \dots, b_p$
- Now after elimination 'p' triangular systems will need to be solved



So, if it is required to solve several systems of equations with the same coefficient matrix A, that is we are interested in solving several systems  $Ax_1$  equal to  $b_1$ ,  $Ax_2$  equal to  $b_2$  through  $Ax_n$  equal to  $b_n$ , that is the right hand side is changing, but the coefficient matrix remains the same. Then we can augment A by including the additional p columns  $b_1, b_2$  through  $b_p$  and increasing the size of the matrix from n by n to n by n plus pn by n plus p.


(Refer Slide Time: 27:57)

**Example: Gauss elimination**

Suppose we want to solve the following equation by Gauss elimination:

$$\begin{aligned}2x - 3y - z + 3w &= 4 \\4x - 4y - z + 11w &= 4 \\2x - 5y - 2z - w &= 9 \\2y + z + 4w &= -5\end{aligned}$$

The augmented system:

$$\begin{bmatrix} 2 & -3 & -1 & 3 & 4 \\ 4 & -4 & -1 & 11 & 4 \\ 2 & -5 & -2 & -1 & 9 \\ 0 & 2 & 1 & 4 & -5 \end{bmatrix}$$


So, after elimination, we will have to solve  $p$  triangular systems for each of the right hand sides  $b_1$  through  $b_p$ . Let us, look at a simple example, suppose we want to solve the system of equations with four equations and four unknowns and we want to solve it using gaussian elimination. The augmented system is given by this, where we have moved right hand side to the left and we haven't at and we are going to operate on this using gaussian elimination, with first step will be the reduction and then, we will use back substitution to solve this system.

So, in the first step, we want to eliminate the coefficients of  $x$  in the second equation right. So, this  $4x$  we want to get rid of this  $4x$ . So, what do we do? We divide the first equation by 2, because 2 is the coefficient of  $x$  in the first equation, we multiply the after dividing the first equation by 2, and we multiply the first equation by four and subtract that from the second equation. Similarly, we multiply the first equation by 2 and subtract it from the third equation and so on.

(Refer Slide Time: 29:20)


**Example: Gauss elimination**

In the first step we want to eliminate the coefficient of  $x$  from the last 3 equations

The multipliers are: equation (2)  $m_{21} = a_{21}/a_{11} = 4/2 = 2$   
equation (3)  $m_{31} = a_{31}/a_{11} = 2/2 = 1$   
equation (4)  $m_{41} = a_{41}/a_{11} = 0/2 = 0$

The reduced system at the end of step 1 is obtained by subtracting:

- $m_{21}$  x equation (1) from equation (2)
- $m_{31}$  x equation (1) from equation (3)
- $m_{41}$  x equation (1) from equation (4)




So, if we do that the multipliers for equation 2 is  $m_{21}$  is equal to  $a_{21}$  by  $a_{11}$ , which is 4 by 2 is equal to 2. The multiplier for equation 3 is  $a_{31}$  by  $a_{11}$ . So, 2 by 2 that is equal to 1. And the multiplier for equation 4, where the coefficient of  $x$  is 0 is nothing but, 0 by 2 which is equal to 0. So, if we subtract  $m_{21}$  times equation 1 from equation 2,  $m_{31}$  times equation 1 from equation 3 and  $m_{41}$  times equation 1 from equation 4.

(Refer Slide Time: 30:10)

**Example: Gauss elimination**

The reduced system at the end of the first step:

$$\begin{bmatrix} 2 & -3 & -1 & 3 & 4 \\ 0 & 2 & 1 & 5 & -4 \\ 0 & -2 & -1 & -4 & 5 \\ 0 & 2 & 1 & 4 & -5 \end{bmatrix}$$


Then, we are going to get this system. This is the reduced system at the end of the first step then, you can see that the second equation. Now, have 0 coefficients it has got a 0 coefficient as on the first column. So, the coefficient of  $x_1$  is actually equal to 0.

(Refer Slide Time: 30:33)

**Example: Gauss elimination**

In the second step we want to eliminate the coefficient of  $y$  from the last 2 equations


The multipliers are: equation (3)  $m_{32} = \frac{a_{32}^{(2)}}{a_{22}^{(2)}} = -2/2 = -1$

equation (4)  $m_{42} = \frac{a_{42}^{(2)}}{a_{22}^{(2)}} = 2/2 = 1$

The reduced system at the end of step 1 is obtained by subtracting:

$m_{32} \times$  equation (2) from equation (3)

$m_{42} \times$  equation (2) from equation (4)



In the second step, we want to eliminate the coefficient of  $y$  from the last two equations. So, this is the coefficient of  $y$ . So, I want to eliminate this, from the last two equations. So, what do I do? I divide the second equation by 2 and then, I multiply the second equation by minus 2 and added to the third equation. I multiply the second equation by 2 and add it to the fourth equation and if I do that. So, basically my multipliers, now are a  $32/2$  by a  $22/2$ . So, minus 2 by 2, this is equal to minus 1, for equation 3. For equation 4 my multipliers are  $m_{42}$  is equal to a  $42/2$  by a  $22/2$ , which is 2 by 2, which is equal to 1. And then, I obtain my reduce system by multiplying equation 2 by  $m_{32}$  and subtracting it from equation 3 and multiply equation 2 by  $m_{42}$  and subtracting it from equation 4, then after, I do that.

(Refer Slide Time: 31:55)


**Example: Gauss elimination**

The reduced system at the end of the second step:

$$\begin{bmatrix} 2 & -3 & -1 & 3 & 4 \\ 0 & 2 & 1 & 5 & -4 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & -1 & -1 \end{bmatrix}$$

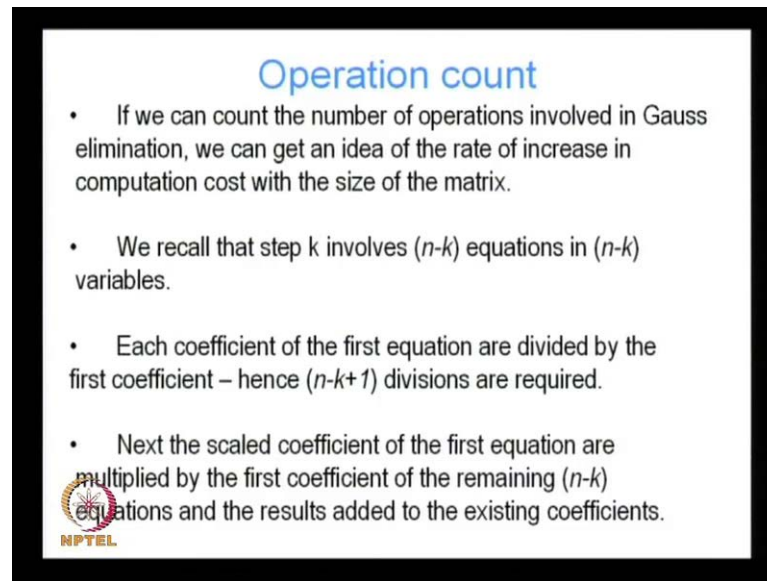
↴

The system is triangular, and can be solved.




This is my system, this is my resultant system and as you can see, this system is triangular right. If you look at my coefficient matrix, it is a triangular system and this can be easily solved using back substitution. Actually for, if this matrix were different we might have had to do a third step, but the structure of the matrix is such that even after, two steps of reduction we get a fully triangular system and we can go ahead with the back substitution and can solve the system.

(Refer Slide Time: 32:42)



**Operation count**

- If we can count the number of operations involved in Gauss elimination, we can get an idea of the rate of increase in computation cost with the size of the matrix.
- We recall that step  $k$  involves  $(n-k)$  equations in  $(n-k)$  variables.
- Each coefficient of the first equation are divided by the first coefficient – hence  $(n-k+1)$  divisions are required.
- Next the scaled coefficient of the first equation are multiplied by the first coefficient of the remaining  $(n-k)$  equations and the results added to the existing coefficients.

 NPTEL

Let us do an operation count, basically we want to count the number of operations involved in gauss elimination. Why do we want to do this, we want to get an idea of the rate of increase in computational cost with the size of the matrix. So, as  $n$  increases how quickly does the computational? What is the weight of increase of the computational cost? That is, what we want to find out, to do that we have to count the number of operations in each step.


We recall that step  $a, k$  involves  $n$  minus  $k$  equations in  $n$  minus  $k$  variables for instance, when after the first step we have  $n$  minus one equations and we have  $n$  minus one variables. The second step has  $n$  minus two equations,  $n$  minus two variables and so on. So, forth each coefficient of the first equation are divided by the first coefficient hence  $n$  minus  $k$  plus 1 division are required, because there are  $n$  minus  $k$  terms and then, there is a right hand side also. So, the  $n$  minus  $k$  plus 1 terms and each term in that first equation has to be divided by the first coefficient. So, the  $n$  minus  $k$  plus 1 division are **required**.

Next, the scaled coefficients of the first equation are multiplied by the first coefficient of the remaining  $n$  minus  $k$  equations, and the results are added to the existing coefficients.

(Refer Slide Time: 34:30)

### Operation count

- Hence  $(n-k).(n-k+1)$  multiplications and additions are required.
- For large values of  $n$  the cost the  $(n-k+1)$  divisions are negligible compared to the cost of the  $(n-k).(n-k+1)$  multiplications and additions.
- Hence neglecting the divisions, the total number of operations is approximately:

$$\sum_{k=1}^{n-1} (n-k)(n-k+1) = \frac{1}{3}n^3 + \frac{1}{2}n^2(p-1) \approx \frac{1}{3}n^3 \text{ (for large } n)$$


So, what does that mean; that means, there are  $n$  minus  $k$  times  $n$  minus  $k$  plus 1 multiplications and additions, why the first equation, we have already, we disregard the first equations. So, now, the  $n$  minus  $k$  equations, but for each  $n$  minus  $k$  equation we have to do  $n$  minus  $k$  plus 1 multiplications and additions. So,  $n$  minus  $k$  times  $n$  minus  $k$  plus 1 multiplication and additions are required.

For, large values of  $n$  the cost of the  $n$  minus  $k$  plus 1 divisions are negligible, compared to the cost of the  $n$  minus  $k$  times  $n$  minus  $k$  plus one multiplications and additions, why is that, because this involves  $n$  square terms. So, the cost of this is negligible compared to the cost of the  $n$  minus  $k$  plus 1 divisions hence, neglecting the divisions the total number of operations is  $n$  minus  $k$  times  $n$  minus  $k$  plus 1, but remember this has to be performed for each step in the reduction and there are  $n$  minus 1 steps. So, we sum over  $k$  is equal to 1 to  $n$  minus 1 and if we do that, we get one third  $n$  cube plus 1 by 2  $n$  square  $p$  minus 1  $p$  is a. You remember is the number of right hand sides the sake of simplicity, if we assume that the number of right hand sides is 1. So, we have 1 by 3  $n$  cube approximately 1 by 3  $n$  cube operations for large  $n$ .

Even if,  $p$  is not 1 by 3  $n$  cube would be much greater than 1 by 2  $n$  square times  $p$  minus 1 again, because  $n$  cube is much greater than  $n$  square. So, the total number of operations



is about  $1/3 n^3$  for large values of  $n$ .

So, once we have obtained a triangular system, to solve the triangular system, it takes  $n$  divisions and as we saw earlier. It takes  $1/2 n(n-1)$  times  $n$  addition or multiplications, which is approximately equal to  $1/2 n^2$  square operations,  $1/2 n^2$  square additions and multiplications hence, the back substitutions for large systems requires  $1/2 n^2$  square additional operations.

Since, for large systems  $1/3 n^3$  recalled, that the from our previous slide  $1/3 n^3$  cube is the total cost of the reduction,  $1/3 n^3$  is always going to be much larger than  $1/2 n^2$  squared for large systems. So, the major cost in Gaussian elimination lies in the reduction to the triangular system. So, if you want to speed up the cost of gaussian elimination. We have to find out cheaper ways to do the reduction, we have to cut down the cost of the reduction and there are several ways of doing that, we are going to talk about some of them later on.

But recall our discussion of Gaussian elimination. We always impose the condition, that pivot the term by, which we divide those equations. The term by which, we normalize the equation for instance. Let us go back and look at our equations again. So, these pivots right a  $a_{22}$  or for instance in this case a  $a_{11}$ , these pivots have to be non-zero. If they become 0 then, this procedure stops. We cannot proceed with a reduction procedure. So, the prerequisite for Gaussian elimination to work without any interchange of rows or columns is that these pivots must be non-zero. So, from the algorithm for Gaussian elimination, we have considered. So, far it is evident that if any of the terms  $a_{11}$ ,  $a_{22}$  or  $a_{33}$  all of these terms are the pivots right. The terms by which we normalize the equation if these are 0 then, our algorithm is going to fail.

Hence, all practical applications of Gaussian elimination require some pivoting, because we cannot be assured, that these terms are not going to be 0 to an understand how the pivoting works. Let us consider the system, the simple system  $3 \times 3$  system in 3 unknowns,  $x_1 + x_2 + x_3$  is equal to 1,  $x_1 + x_2 + 2x_3$  is equal to 2 and  $x_1 + 2x_2 + 2x_3$  is equal to 1, after the first step of the elimination, we get  $0x_2 + 1x_3$  is equal to 1 and  $1x_2 + x_3$  is equal to 0.

Since, the coefficient of  $x^2$  in the first equation is 0, this algorithm is going to break down, because the pivot is 0; however, the problem is resolved, if we interchange the rows in the above equation. So, if we interchange the rows in the above equation since, the coefficient of  $x^2$  in the second equation is positive. We can again proceed with the reduction.

So, in general suppose at step  $k$  in our reduction process we find a  $a_{kk}$  is equal to 0. In that case, there will always be some other element in the  $k$ th column suppose  $a_{rk}$ , which is non-zero go back and look at our previous example, in this case  $a_{22}$  is 0, but  $a_{32}$  is no longer is not 0. So, that is why we can interchange these two rows to make sure, that the pivot is no longer 0.

Why are we assured that, there must be at least one element in the  $k$ th column, which is not 0, because if all the elements in the  $k$ th column is 0, this implies that the first  $k$  columns are linearly dependent, because we have formed the  $k$ th column by taking scalar multiples of the first  $k-1$  columns, so if the  $k$ th column, if all the terms in the  $k$ th column as 0. It only means that the first  $k-1$  columns are linearly dependent, recall our discussion of linear dependence and linear independence from our previous lecture.

So, if  $k-1$  vectors, if we can add  $k-1$  columns that is,  $k-1$  vectors. If, we multiply each column by scalar and we add them together and the result is 0. We say that those first  $k-1$  columns are linearly dependent, they are not linearly independent and if they are linearly dependent, what does that mean; that means, that coefficient matrix  $A$  is singular, because its rank is less than  $n$ .

But since, we have assumed that the matrix is non-singular; that means, it has got full rank that is the  $k$ th column cannot be obtained as a linear combination of the first  $k-1$  columns. So, we are assured that, at least some elements in the  $k$ th column must be non-zero otherwise, the matrix is singular. Since at least, one of the elements in the  $k$ th column is non-zero, we can always interchange that row with the  $k$ th row to restore not to make sure, that a pivot is non zero and we can continue with reduction proceeding.

Since, a  $a_{rk}$  is non-zero. We can interchange rows  $r$  and  $k$  and proceed with elimination thus it follows that any non-singular system of equations can be reduced due to triangular form by gaussian elimination, accompanied by row interchanges. If the system is singular, there is no assurance this can happen because we can get eventually a row, which has got 0 entries all 0 entries and in that case we can no longer get a non-zero pivot.

However, it is not always true that only, if the pivot term is 0, it is necessary to do pivoting. If the pivot term, becomes really small it becomes sufficiently close to 0 that if we divide the equation by that pivot term, we are going to get numerical instability. In that case, also it might be necessary to do pivoting. So, it may be necessary to perform row interchanges not only when the pivotal element is exactly 0, but also when it is nearly 0 in order to improve numerical instability to see the need for this we consider the previous system the previous system which we considered earlier, but in this case we change the coefficient of  $a_{22}$  from one to 1.0001.

Let us take a look at our previous system. So, here the coefficient of  $a_{22}$  was 1. So, instead of change if keeping it 1, let us change it to 1.0001 and see what happens, if we do that, we will see that our triangular system after Gaussian elimination is going to become something like this, and if we do that substitution and using round off after three decimal. So, we will get  $x_1$  equal to 0,  $x_2$  is equal to 0,  $x_3$  is equal to 1; however, the true solution if we use four decimals, if we round to four decimals is actually  $x_1$  equal to 1,  $x_2$  equal to minus 1.0001,  $x_3$  is equal to 1.0001.

So, we can see by using a nearly 0 element as the pivot our system of equations has got a totally erroneous solution, why because when we use the pivot 1.0001. Our system became ill conditioned. So, minor perturbations give totally different solutions. So, if instead of using as near 0 pivot if we did pivoting if we interchange the rows to get sufficiently large pivot, we would have got the true solution.

So, it is very important to do pivoting not only when the pivotal element is fully 0, but when the pivotal element is nearly 0. It is clear that without interchanges Gaussian elimination is going to be inherently unstable and will give rise to meaningless results;

however, by interchanging rows 2 and 3 in our in rows 2 and 3 in this equation we could and again using three decimal points in our computations, we could get  $x_1$  is equal to 1 and  $x_2$  is equal to minus 1.0001, which is almost exactly the true solution thus it is clear that interchange of rows is essential in order to obtain accurate solutions in such situations.

Now, let us talk about two particular types of pivoting first, we want to talk about partial pivoting. If, we interchange of rows is extended to make sure that, we interchange the  $k$ th row with the row that has the largest element in the  $k$ th column, it is known as partial pivoting, this is clear if we look at this little picture. Suppose, we are looking at this system we have reduced it up to here. So, it is a triangular system up to here and this is our and we have here left with this system and we want to reduce it further and we look at the pivotal element here and we find that the pivotal element is either 0 or sufficiently close to 0, so that we need to do pivoting.

So, what do we do? We look at all the elements in the shaded region that is basically we look at all the elements in this column and we select the element which has got the largest non-zero magnitude, which has got the largest absolute value and then we interchange this row with the top row. So, that our pivot we get the largest with the best estimate we can of the pivot.

So, in order to avoid ill conditioning and then, we proceed with rest of our operations. So, this the partial pivoting basically, requires that we choose  $r$  here, such that  $a_{rk}$  is the maximum in that column, over  $k$  equal to one through  $n$  this term is the maximum in this column and after we find the maximum element in this column we interchange this row with our  $k$ th row.

So, this is known as partial pivoting, we also have something known as complete pivoting. In complete pivoting, we not only require interchange the rows, we also interchange the rows and columns instead of restricting to the  $k$ th column to find the element to interchange with the current  $a_{kk}$  the search is extended to other columns also.

So,  $a_{kk}$  is replaced by the element with the largest magnitude in the partition basically let us look at our picture in the next slide. So, now, again suppose this is my near 0 element or 00 elements. So, instead of finding the maximum element in this column, we find the element which has the maximum magnitude in this entire partition, that is we look over all the rows and columns in this partition in the shaded region and find the maximum element and then, we interchange this column with that column and this row with that row.

So, basically instead of restricting to the  $k$ th column, we find the element to interchange with the current  $a_{kk}$  the search is extended to other columns also. So, extend the search to other columns and  $a_{kk}$  is replaced by the element with the largest magnitude in the entire partition. So, it is replaced by the largest element in the entire partition.

Thus complete pivoting requires that we choose  $r$  and  $s$  as this smallest integers, which satisfy  $a_{rsk}$  is equal to maximum of  $a_{ij}$ ,  $k$  is lesser than equal to  $i$ ,  $j$  is lesser than equal to  $k$  and interchange the rows  $r$  and  $s$  and the columns  $k$  and  $s$ . Complete pivoting therefore, requires searching over the  $n$  minus  $k$  plus 1 elements of each of the  $n$  minus  $k$  plus 1 columns of the partition. To determine, the pivot thus complete pivoting is much more expensive than partial pivoting, because we have to search over all these elements in the entire partition.

In practice; however, partial pivoting is usually sufficient to prevent numerical instability and is commonly adopted. Thank you, we will continue with our discussion of Gaussian elimination in particular, we will look at some particular algorithms for doing Gaussian elimination for instance lu decomposition, crowd factorization, etcetera, which can further reduce the cost of Gaussian elimination by making it more efficient and more compact.