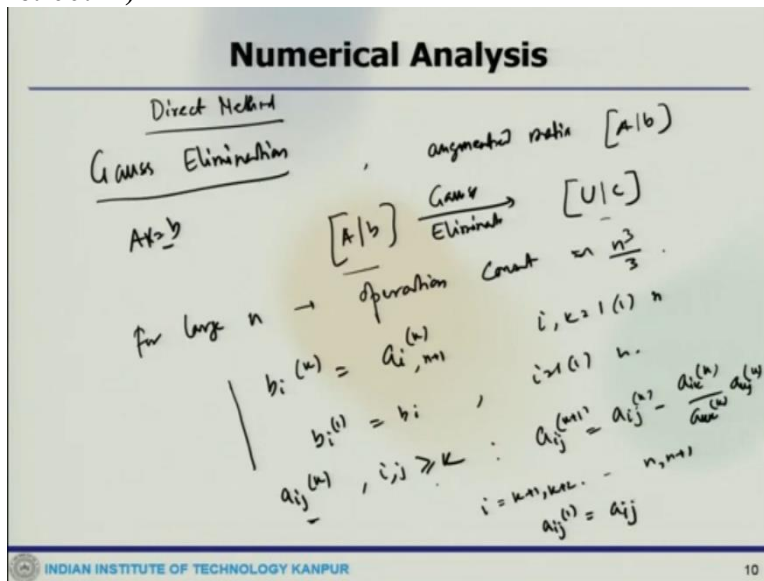


**Computational Science in Engineering**  
**Prof. Ashoke De**  
**Department of Aerospace Engineering**  
**Indian Institute of Technology, Kanpur**

**Lecture - 31**  
**Numerical Analysis**

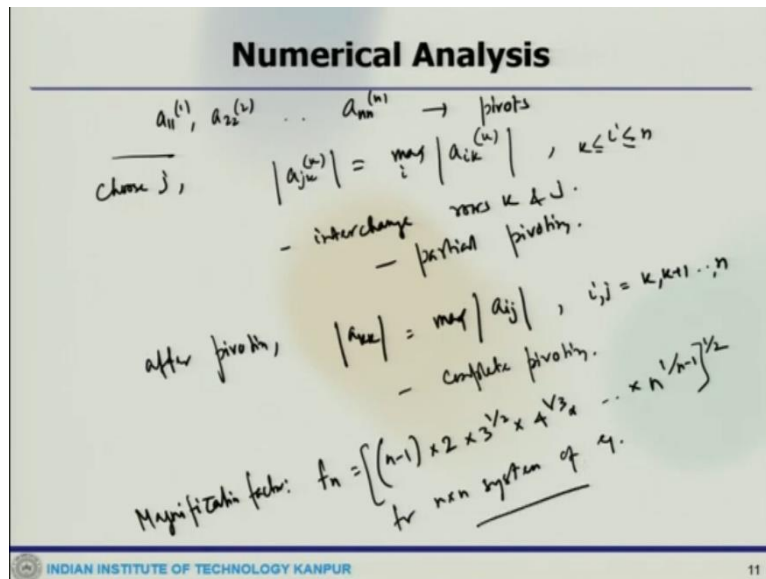
Okay, so let us continue the discussion of the linear system or the linear system equations  $Ax = b$  solution. So, we were talking about Gauss elimination.

(Refer Slide Time: 00:27)



And this is where we have stopped in the last session that when you do this elimination what you get finally this and the order of operation, would be like this and then you find out kind of like, you said this and for  $i, j \geq k$ , you can find out the elements like this. Here,  $i = k + 1$ , like that it will go to  $n, n + 1$  and  $a_{ij}^{(k+1)} = a_{ij}$ .

(Refer Slide Time: 01:08)



Now, the elements which are there like,

$$a_{11}^{(1)}, a_{22}^{(2)}, \dots, a_{nn}^{(n)}$$

these are called sort of your pivots, if you recall. To avoid division by zero or to round off error, partial pivoting is normally used. So, the pivot is chosen like this, let us say you can choose  $j$  and the smallest integer for which

$$|a_{jk}^{(k)}| = \max_i |a_{ik}^{(k)}|$$

where  $i$  lies between  $n$  and  $k$  and interchange rows  $k$  and  $j$ .

So, it is called, so there is an interchange rows  $k$  and  $j$ , so this is called partial pivoting. This is also we have discussed. Now, if at the  $k$  step we interchange both the rows and columns of the matrix, so that the largest number in magnitude in the remaining matrix is used as pivot that is after pivoting, so what we get like

$$|a_{kk}| = \max_{i,j} |a_{ij}|$$

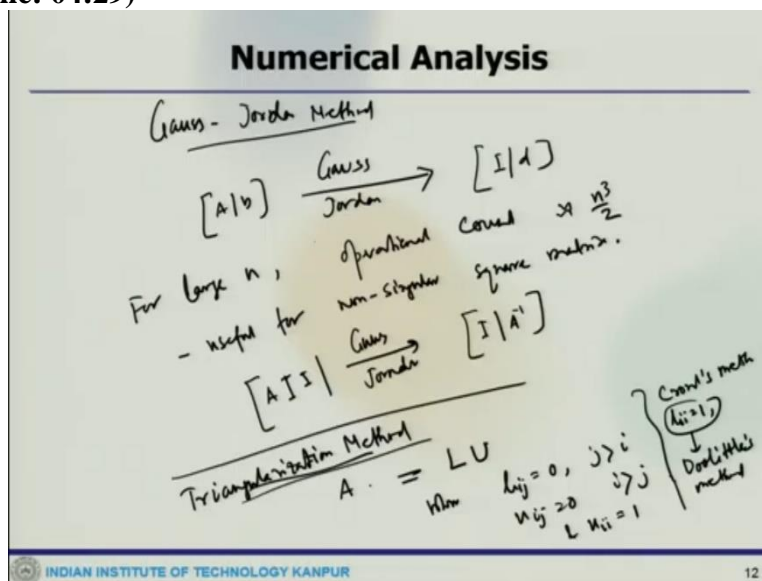
where  $i, j = k, k + 1, \dots, n$ . Now, this is called complete pivoting.

Now, here, we can interchange 2 columns and the position of the corresponding elements in this solution matrix is also change. Now, the complete pivoting is safe as errors are never magnified unreasonably. Secondly, the magnification factor is less than or equal to like, this is called magnification factor which is called

$$f_n = \left[ (n - 1) \times 2 \times 3^{\frac{1}{2}} \times 4^{\frac{1}{3}} \dots \times n^{\frac{1}{n-1}} \right]^{1/2}$$

So, this is what we get it for the magnification factor. So, this magnification factor actually reveals that the growth is within the limits. Even though the bound for the magnification factor in the case of partial pivoting cannot be given by this expression, it is known that the magnification error is almost eight times in most cases which is for complete pivoting. So, the complete pivoting actually doubles the cost, while the partial pivoting costs negligibly more than the Gauss elimination. So, the bottom line here is that, Gauss elimination with or without partial pivoting is same for diagonal dominant matrices.

(Refer Slide Time: 04:29)



Now, we look at Gauss Jordan method. So, what it does is that, now, we can start with the same augmented matrix like,  $[A|b]$ . So, the coefficient is reduced to diagonal. So, after this process what we get is that identity and different matrix. So, this means the elimination is done only with the equations below, here and the method is more expensive from the computation point of view compared to the gauss elimination method.

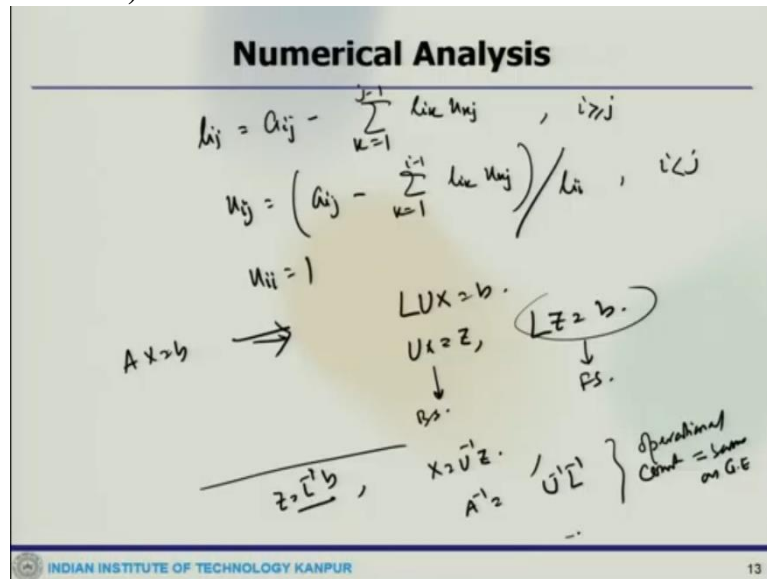
For large  $n$ , the operational count is,  $\frac{n^3}{2}$ , however, this method is useful in finding the inverse of a non-singular square matrix but this is useful for non singular square matrix. So that is important. So, what do we get like,  $[A|I]$ , now, if we use Gauss Jordan elimination here then, we get  $[I|A^{-1}]$ . So, now, we have to do some sort of like triangularization method. Here, what is that in this particular method?

You have the coefficient of the matrix  $A$  which is decomposed into the product of the lower triangular and upper triangular that means, this should be written as or decompose as the  $L U$

decomposition. So, this is what we have discussed as L U decomposition. So, where  $l_{ij} = 0$  for  $j > i$  and  $u_{ij} = 0$ , for  $i > j$  and  $u_{ii} = 1$ . So, this is also known as Crout's method.

Instead of  $u_{ii} = 1$ , if we take  $l_{ii} = 1$  then the method is called, so, instead of that, if we take this, the method is named as Doolittle's method. So, there are different ways one can handle the parameters. we uniquely determine lower triangular and upper triangular.

(Refer Slide Time: 07:36)



So, what we get

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj}$$

So, this is  $i \geq j$  and

$$u_{ij} = \frac{(a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj})}{l_{ii}}$$

which is  $i < j$  and  $u_{ii} = 1$ . So, the system  $AX = b$ , now, get it like,  $LUX = b$ . So, we rewrite that  $UX = Z$  and  $LZ = b$ . So, first we can do that from this equation, we using the forward substitution, so, we first find  $Z$  from here, forward first, so, this is forward substitution.

And you can find out  $Z$  and then we find out  $X$  from here is in the backward substitution. So that is how one can do or alternatively, one can do like first you find out  $L, Z$  from here and then you can find out the  $X$  like,

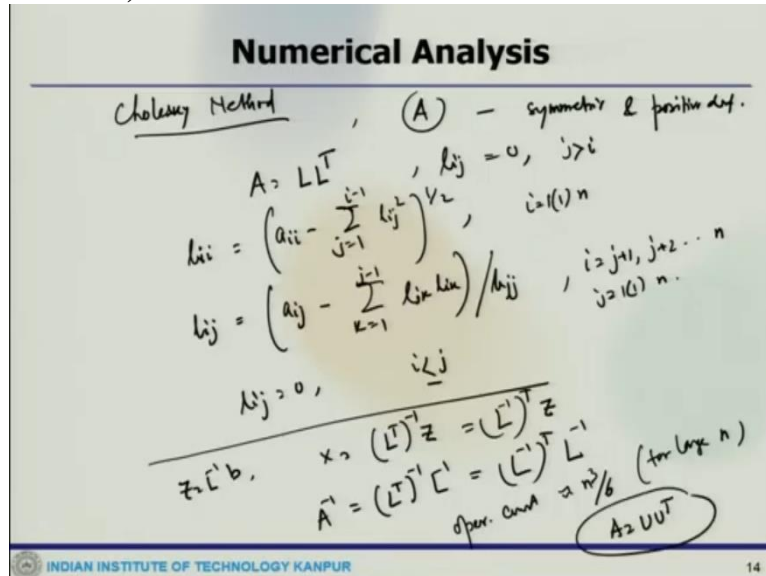
$$X = U^{-1}Z$$

and  $A$  inverse which can be obtained as that

$$A^{-1} = U^{-1}L^{-1}$$

So, this triangularization is used more often than the Gauss elimination and this case the operational count is also same as Gauss elimination process. LU decomposition is not always guaranteed for arbitrary matrices, so, decomposition is guaranteed when the A is positive definite, so that is important, otherwise this is not varying.

(Refer Slide Time: 09:41)



Now, there is another method which is called the Cholesky method or this is also known as square root method. If in the coefficient of A is symmetric and positive definite then, A can be decomposed as

$$A = LL^T$$

where  $l_{ij} = 0$  for  $j > i$ . The elements of L here given as

$$l_{ii} = \left( a_{ii} - \sum_{j=1}^{i-1} l_{ij}^2 \right)^{1/2}$$

where  $i$  goes from 1 to  $n$ . And

$$l_{ij} = \frac{\left( a_{ij} - \sum_{k=1}^{j-1} l_{jk} l_{ik} \right)}{l_{jj}}$$

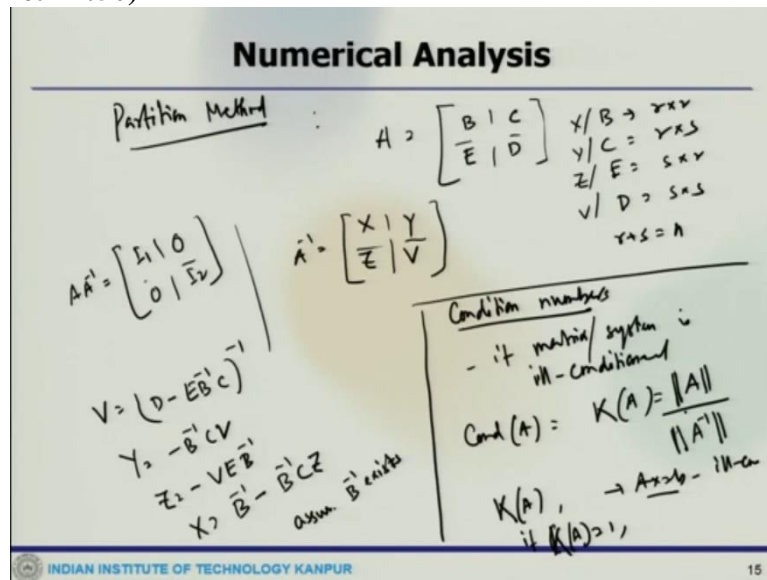
Where  $i = j + 1, j + 2, \dots, n$  where  $j$  goes 1 to  $n$  and  $l_{ij} = 0$  for  $i < j$ . Now, we have already seen that. So, this is now when we basically decompose the system. So, what we get here like, whenever we have done that, lower triangular and upper triangular, so, we get  $LZ = b$ , so, from there we can find out  $Z = L^{-1}b$  and  $X = (L^T)^{-1}Z = (L^{-1})^T Z$ .

So,

$$A^{-1} = (L^T)^{-1}L^{-1} = (L^{-1})^T L^{-1}$$

and this particular method in the Cholesky method, the operational count which is also like  $\frac{n^3}{6}$ , this is for large  $n$ , we can get this. Now, instead  $A = L^T L$ , we can also decompose like, similarly, we can decompose  $A = U U^T$ , this is also a possible decomposition one can do to find out the thing.

(Refer Slide Time: 12:50)



Now, similarly, you can have partition method. So, what this method does? That usually it is used to find the inverse of a large nonsingular square matrix by partitioning. So, let us say  $A$  can be partitioned like this

$$A = \begin{bmatrix} B & | & C \\ - & & - \\ E & | & D \end{bmatrix}$$

For  $B, C, E, D$  these are of the orders like  $B$  is of the orders of  $r \times r$ ,  $C$  is  $r \times s$ ,  $E$  is  $s \times r$ ,  $D$  is  $s \times s$ , where  $r + s = A$ . Similarly, we can partition  $A^{-1}$  which would be

$$A^{-1} = \begin{bmatrix} X & | & Y \\ - & & - \\ Z & | & V \end{bmatrix}$$

where  $X, Y, Z, V$  are the same orders as like. So, this would be  $X, Y, Z$  and  $V$  as the same orders like  $B, C, E, D$ .

Now, using the identity, what we can get? Like

$$A A^{-1} = \begin{bmatrix} I_1 & | & 0 \\ - & & - \\ 0 & | & I_2 \end{bmatrix}$$

and

$$V = (D - EB^{-1}C)^{-1}$$

$$Y = -B^{-1}CV$$

$$Z = -VEZ^{-1}$$

$$X = B^{-1} - B^{-1}CZ$$

where we have assumed that  $B^{-1}$  exists. So, if  $B^{-1}$  does not exist but  $D^{-1}$  exists then the equations can be modified suitably.

So, this is what it requires. Now, this is what partition method does. Then, whatever we say we talk about some important thing called the condition number or numbers, whatever. Sometimes one come across a system of equations which are very sensitive to round off errors that is, one gets different kinds of solution when the elements are rounded to different number of digits. In such cases, the system is called ill condition system.

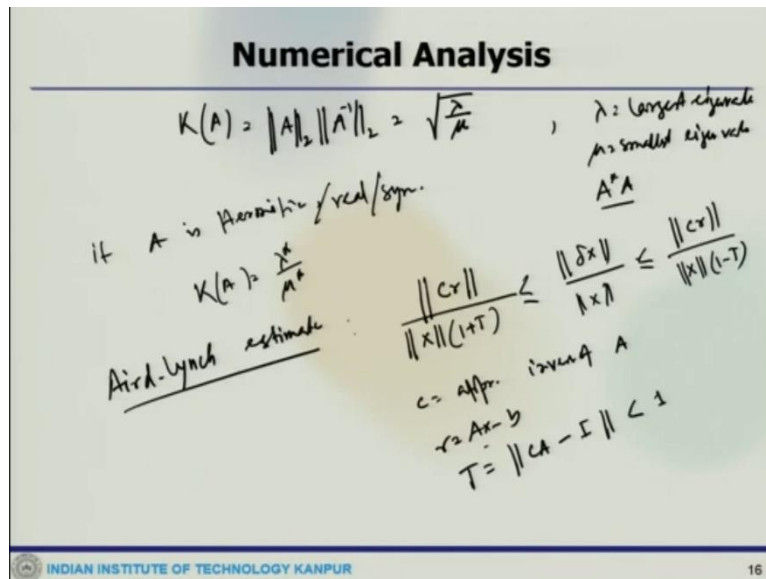
So, the measure of the ill condition ness is given by the value of condition number. So, if matrix or the system is ill conditioned then the condition number give the measure of that and how we define that,

$$Cond(A) = K(A) = \frac{\|A\|}{\|A^{-1}\|}$$

So, this is what we have already discussed, how to find out this number? The number is usually referred as a standard condition number.

Now, if  $K(A)$  is large then, the small changes in A or B produces those relative changes in X and the system of equation like  $Ax = b$  becomes ill conditions or if  $K(A) = 1$  then the system is ill condition and that time this norm is the spectral norm.

**(Refer Slide Time: 16:57)**



And then we can write that.

$$K(A) = \|A\|_2 \|A^{-1}\|_2 = \sqrt{\frac{\lambda}{\mu}}$$

where  $\lambda$  is the largest Eigen value,  $\mu$  is the smallest Eigen value, in modulus of  $A^*A$ , now if  $A$  is Hermitian or real or symmetric then, the condition number would be

$$K(A) = \frac{\lambda^*}{\mu^*}$$

So, this is the largest and smallest Eigen values in modulus of  $A$ . So, another important condition number is the Aird-Lynch estimate which is, Aird Lynch estimate.

So, these estimates give both the lower and upper bounds for the error magnification. So, we have estimated like

$$\frac{\|Cr\|}{\|X\|(1+T)} \leq \frac{\|\delta X\|}{\|X\|} \leq \frac{\|Cr\|}{\|X\|(1-T)}$$

where  $C$  is the appropriate inverse of  $A$ , usually the outcome of Gauss elimination process,  $r$  is  $Ax - b$  which is a residual vector and  $X$  is the computed solution and  $T = \|CA - I\| < 1$ . So, this is what you get for all this. Because the condition number is important to define the system, whether it is ill condition or you can get a solution for the given system.

**(Refer Slide Time: 19:18)**



**Numerical Analysis**

---

Iterative Methods

$$x^{(k+1)} = Hx^{(k)} + C \quad \text{--- (1)}$$

$\uparrow$       $\uparrow$       $\uparrow$   
 $H$       $x^{(k)}$       $C$

$k \rightarrow \infty, x^{(k)}$  converges to  $x = A^{-1}b$

---

if  $\|H\| < 1$ , w. i. t.  $\rho(H) < 1$   
 $A = L + D + U, (L + D + U)x = b \dots \text{--- (2)}$

$Ax = b$   
 $H =$  iteration matrix.  
 $C =$  column vector.

INDIAN INSTITUTE OF TECHNOLOGY KANPUR 17

Now, with that we go to some of that so, this is what happens when you have direct method, now, we will go to some iterative methods for finding the solution of  $Ax = b$ . In general, the iterative method of the solution of a system which is given as that  $Ax = b$ , is defined in the formula,

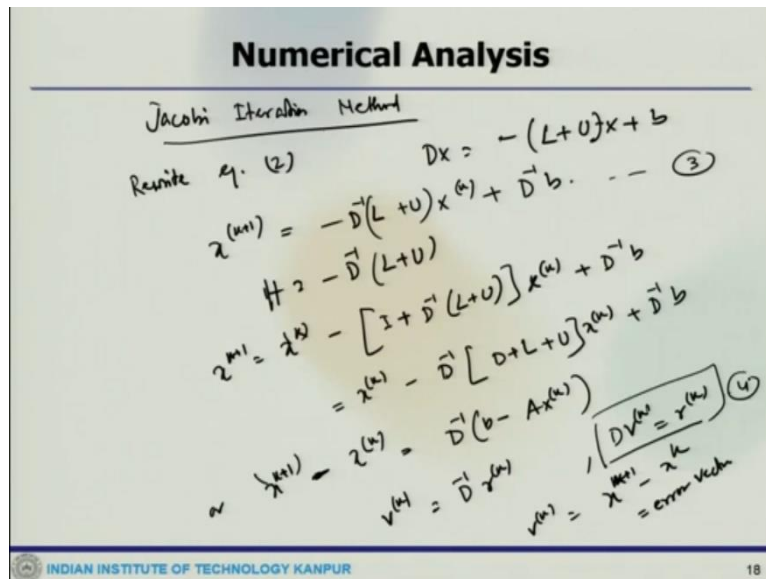
$$x^{(k+1)} = Hx^{(k)} + C$$

where the  $x^{(k+1)}$  and  $x^{(k)}$  are the approximation at  $(K+1)$ th and  $K$ th iteration. So, this is an approximation at  $(K+1)$  iteration, this is an approximation at  $K$ th iteration and  $H$  is called the iteration matrix.

So, depending on  $A$  and  $C$ , is the column vector, so this is column vector. So, in the limiting case when  $K$  tends to infinity  $\infty$ ,  $x^{(k)}$  converges to the exact solution like  $X = A^{-1}b$ . Now, there is one small quick theorem which says that, the iteration method, let us say, we say here, this iteration method which we have in the form of equation one for getting a solution of this  $Ax = b$ , converges to the exact solution for any initial vector if modulus is less than 1 or if and only if  $\rho(H) < 1$ .

So, let the coefficient matrix of  $A$  can be written as,  $A = L + D + U$  where  $L$  is lower triangular diagonal and strictly upper triangular parts, then what we can write that,  $(L + D + U)X = b$ , so, under this category of iterative method.

**(Refer Slide Time: 21:30)**



The one which we will first talk about is the Jacobi iteration method. So, which says that, so, we rewrite this equation, let us say we say, this one then we say this one is 2. So, we rewrite equation 2 as

$$DX = -(L + U)X + b$$

Now, we define an iterative process like,

$$x^{(k+1)} = -D^{-1}(L + U)x^{(k)} + D^{-1}b$$

And the iteration matrix is given as that

$$H = -D^{-1}(L + U)$$

now, this is called the Jacobi iteration method, the way the iteration is done.

So, we can rewrite that 3, like

$$x^{(k+1)} = x^{(k)} - [I + D^{-1}(L + U)]x^{(k)} + D^{-1}b$$

which one can write

$$x^{(k+1)} = x^{(k)} - D^{-1}[D + L + U]x^{(k)} + D^{-1}b$$

or what we can write like,

$$x^{(k+1)} - x^{(k)} = D^{-1}(b - Ax^{(k)})$$

So, this is what you can write or you can write that

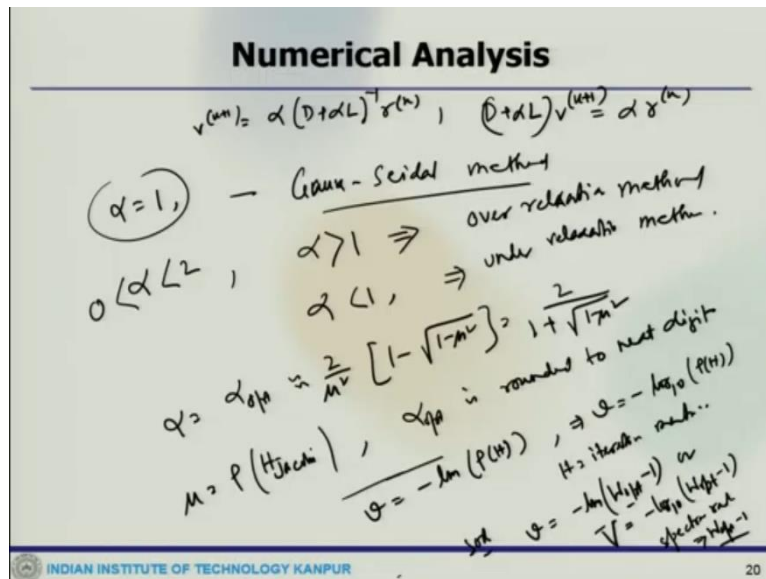
$$V^{(k)} = D^{-1}r^{(k)}$$

where,

$$DV^{(k)} = r^{(k)}$$

and  $V^{(k)} = x^{(k+1)} - x^{(k)}$  and this is called the error vector.

**(Refer Slide Time: 23:56)**



And  $r^{(k)} = b - Ax^{(k)}$  which is called a residual vector. So, now, if you see this guy, what do we have written here? This  $V^{(k)}$  equals to like this one,  $DV^{(k)} = r^{(k)}$ . So, equation 4, so, from the computational point this may be preferred as we are dealing with the errors and not the solution. So, this is from the point of view of the computation one can prefer this one. Now, you can have some other iterative methods like, called the Gauss-Seidel iteration method.

So, in terms of the error vector what we can write from this equation? What we can write that,

$$V^{(k+1)} = \alpha(D + \alpha L)^{-1} r^{(k)}$$

Or

$$(D + \alpha L)V^{(k+1)} = \alpha r^{(k)}$$

Here,  $\alpha = 1$ , these equations here, this particular equation it reduces to Gauss Seidel method. So, the SOR method essentially for  $\alpha = 1$ , it becomes Gauss Seidel method and the relaxation parameter actually satisfied a condition that  $\alpha$  lies between 0 to 2 when, alpha is greater than one. This method is called over relaxation method.

If  $\alpha$  is less than one this is called under relaxation method and the maximum convergence of a SOR is obtained like

$$\alpha = \alpha_{opt} \approx \frac{2}{\mu^2} [1 - \sqrt{1 - \mu^2}] = \frac{2}{1 + \sqrt{1 - \mu^2}}$$

where

$$\mu = \rho(H_{Jacobi})$$

and  $\alpha_{opt}$  is rounded to next digit. So, the rate of convergence of an iterative method is defined as

$$v = -\ln \rho(H)$$

also like, one can write,

$$v = -\log_{10} \rho(H)$$

where H is the iteration matrix.

Now, the spectral radius of the SOR method is  $(W_{opt} - 1)$  and its rate of convergence would be, for SOR rate of convergence would be

$$v = -\ln(W_{opt} - 1)$$

or

$$v = -\log_{10}(W_{opt} - 1)$$

so, this is the spectral radius here. The spectral radius is  $(W_{opt} - 1)$ , so that is how we get the rate of convergence like this. So, you see there, I mean, there are multiple variants of iterative methods, what one can use and we will look some of this more in the sort of in the next session. We will stop it here.