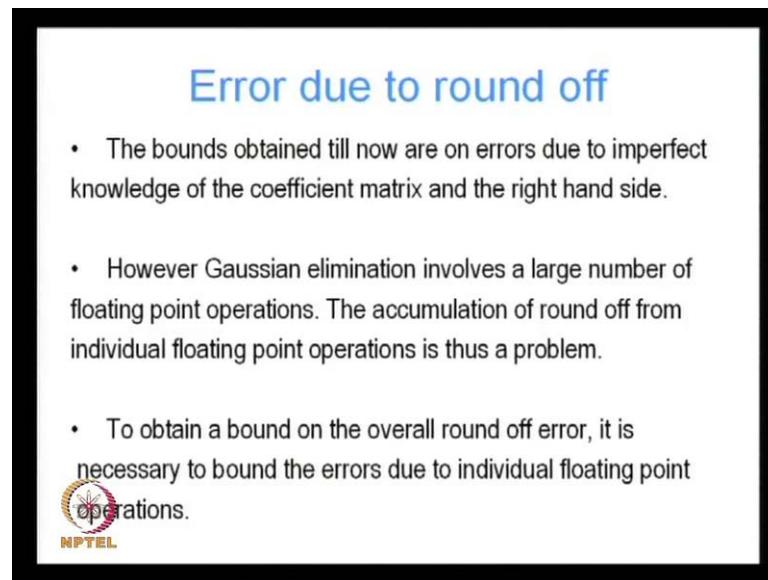


Numerical Methods in Civil Engineering
Prof. Arghya Deb
Department of Civil Engineering
Indian Institute of Technology, Kharagpur

Lecture - 7
Error Bounds and Iterative Methods for Solving Linear Systems


In the seventh lecture in our series in numerical methods in civil engineering, we are going to continue our discussion on error bounds for direct methods in addition. We are going to introduce iterative methods for solving linear systems.

(Refer Slide Time: 00:31)



Error due to round off

- The bounds obtained till now are on errors due to imperfect knowledge of the coefficient matrix and the right hand side.
- However Gaussian elimination involves a large number of floating point operations. The accumulation of round off from individual floating point operations is thus a problem.
- To obtain a bound on the overall round off error, it is necessary to bound the errors due to individual floating point operations.

 NPTEL

The bounds we obtained till now on the errors due to our on the errors on the solution of linear systems, were due to imperfect knowledge of the coefficient matrix, and the right hand side. Basically what we what we found was that if the coefficient matrix has an error δa or the right hand side has an error δb , we try to find bounds on the errors that will be introduced to the solution x due to the errors in a and b , due to the error δa in a , and the error δb in b that is going to result in solutions to errors to the solution x . And we try to find the bounds on those errors and what we found was that if the matrix is depending on the condition number, if the matrix is well conditioned then errors due to in the coefficient matrix, and the errors due to the right hand side will be bounded, and it is going the magnitude of the bound is going to depend on the condition number of the matrix of the coefficient matrix.

Today, we are going to talk about errors in Gaussian elimination, because of the large number of floating point operations, which take place during Gaussian elimination the accumulation of round off errors due to these floating point operations lead to overall errors in the solution, and is therefore a problem to obtain a bound on the overall round off error, it is necessary to bound the errors due to individual floating point operations. Since the Gaussian elimination process is a combination of a large number of floating point operations, we have to find bounds on the error of individual floating point operations in order to find a bound on the error of overall error, due to in the solution due to round off during Gaussian elimination.

(Refer Slide Time: 02:45)


Bound on round off error

- We consider the floating decimal representation of a number a :

$$a = m \cdot 10^q, \text{ where } .1 \leq m \leq 1, q \text{ is an integer.}$$
- In the computer a is represented by :

$$\bar{a} = \bar{m} \cdot 10^q, \text{ where } \bar{m} \text{ is the approximate representation of the mantissa } m, \text{ rounded off to 't' decimals, given that the precision of the machine is 't' decimal digits.}$$
- Hence the bound on the absolute error in the mantissa is:

$$|\bar{m} - m| \leq \frac{1}{2} 10^{-t}$$



We consider the floating point decimal representation of a number suppose we have a floating point number a how does the computer represent that number internally, well the computer almost always represents the number using this format. It represents it as some number m and times 10 depending on the base of the computer, if 10 is the base then that will be 10 there if instead of the instead of 10, it is some other number some other base then that is going to be the appropriate base. So, 10 to the power q where m lies between 0.1 and ones m m is a number between 0.1 and 1 and q is an integer. So, m is known as the mantissa and q is known as the exponent. So, this is how the computer represents any floating point number a, but in reality the computer cannot represent m up to infinite precision.

Because the computer has finite precision for instance if the computer rounds off floating point numbers to t decimals, then a will not actually be the computer will not actually be storing a , it will be storing an approximation to a which is denoted by \bar{a} and the approximation arises; because m is only the computer only stores an approximate representation of the mantissa m \bar{m} , which is rounded off to t decimals given that the precision of the machine is t decimal digits. So, \bar{m} is the approximation of m which is stored in the computer. Hence the bound on the absolute error in the mantissa m is given by $\bar{m} - m$, which is less than or equal to half 10 to the power minus t this we have obtained in a previous lecture. So, the error is bounded the round of error is bounded by half into 10 to the power minus t .


(Refer Slide Time: 04:59)

Bound on round off error

- Hence the bound on the relative error in ' a ' is:

$$\left| \frac{\bar{a} - a}{a} \right| \leq \left| \frac{\bar{m} \cdot 10^q - m \cdot 10^q}{m \cdot 10^q} \right| = \frac{|\bar{m} - m| \cdot 10^q}{|m| \cdot 10^q} \leq \frac{\frac{1}{2} \cdot 10^{-t}}{10^{-1}} = \frac{1}{2} \cdot 10^{1-t}$$

- This limiting value on the relative error due to round off is denoted as the machine unit (u).
- The relative error due to round off in a floating point number cannot exceed the machine unit of the computer.



That was the bound on the absolute error in the mantissa, if we look at the relative error in a that is given by $\bar{a} - a$ by a which is by definition the relative error. So, this has got to be less than or equal to $\bar{m} \cdot 10^q - m \cdot 10^q$ dotted with 10^q divided by $m \cdot 10^q$. That is simply because $\bar{m} - m$ is less than or equal to half 10 to the power minus t sorry. So, we can write that as this is the $\bar{m} - m$ dotted with 10^q divided by $m \cdot 10^q$ is lesser than or equal to half 10 to the power minus t 10 to the power 10 to the power q 10 to the power q cancels out, this is we know is lesser than or equal to half 10 to the power minus t , and since mod of m if we go back to the

previous slide and see m is lesser than or equal to point one and greater than or equal to point one less than or equal to one.

So, if we if we replace mod of m by point one this has always got to be greater than, because we are taking the smallest possible value of m in the denominator. So, this thing has to be less greater than or equal to this thing this thing has to be lesser than or equal to this thing, and this gives me a bound half in to 10 to the power one minus t . A correction this is this should not be lesser than or equal to this is actually equal to right. So, mod of a bar minus a by a is equal to m bar dotted with 10 to the power q minus m 10 to the power q divided by this which is equal to this, but this is lesser than or equal to this because on in the denominator, we have replace mod of m by its smallest possible value which is 10 to the power minus 1. So, we get this bound on the relative error in a .

This limiting value on the relative error due to round off is denoted as the machine unit as you can see this is totally dependent on the machine precision t . So, the relative error due to round off in a floating point number cannot exceed the machine unit of the computer denoted as this whole thing denoted as u .

(Refer Slide Time: 07:51)


Bound on round off error

- Next let us consider the effect of operating on two floating point numbers x and y , denoting the true result of the operation as $x \text{ op } y$ and the floating point representation as $fl(x \text{ op } y)$
- Since the relative error due to round off is bounded by u :

$$|fl(x \text{ op } y) - x \text{ op } y| \leq |x \text{ op } y| u$$
- Hence there exists a number δ with $|\delta| \leq u$, such that:

$$|fl(x \text{ op } y) - x \text{ op } y| = |x \text{ op } y| \delta$$

$$\therefore fl(x \text{ op } y) = x \text{ op } y(1 + \delta) \quad (**)$$



Next let us consider the effect of operating on two floating point numbers x and y , and let as denote the true result of the operation as $x \text{ op } y$, and the floating point representation as float of $x \text{ op } y$. So, $x \text{ op } y$ is the true solution and this is how the computer is going to represent this.

So, the relative error due to round off is bounded by u , we can say that $|x \text{ op } y - \text{fl}(x \text{ op } y)| \leq |x \text{ op } y| u$; which is exactly what we get if we replace $|x \text{ op } y|$ by $\text{fl}(x \text{ op } y)$, and $\text{fl}(x \text{ op } y)$ is the true solution by $x \text{ op } y$ right. So, this is what we get this is the bound on the relative error now; because this is less than $|x \text{ op } y| u$ we can be sure that, there is a number δ with $|\delta| \leq u$ such that this becomes equal to $|x \text{ op } y| u$. So, there must exist a number for which this becomes equal to that right, and this gives rise to this expression $\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta)$. So, this is how the computer represents $x \text{ op } y$, this is the true solution and the true solution times this error $1 + \delta$ is the floating point representation of the result of the operation of $x \text{ op } y$.

(Refer Slide Time: 09:45)


Bound on round off error

- Recall that in the k^{th} step of the Gaussian elimination of a symmetric matrix, the elements are transformed as:

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} = a_{ij}^{(k)} - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} a_{kj}^{(k)}, i, j = k+1 \dots n$$
- Because of round off the computed values of each of the quantities on the right hand side (denoted by an overbar) will differ from their true values and result in additional errors:

$$\bar{m}_{ik} = \frac{\bar{a}_{ik}^{(k)}}{\bar{a}_{kk}^{(k)}} (1 + \delta_1) \quad (*) \quad |\delta_1| \leq u$$

$$\bar{a}_{ij}^{(k+1)} = [\bar{a}_{ij}^{(k)} - \bar{m}_{ik} \bar{a}_{kj}^{(k)}] (1 + \delta_2) (1 + \delta_3) \quad (**) \quad |\delta_2|, |\delta_3| \leq u$$



Next let us go back to the Gaussian elimination, and let us recall that in the k -th step of the Gaussian elimination of symmetric matrix the elements are transformed as the following $a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)}$, where m_{ik} is the multiplier times $a_{kj}^{(k)}$ at the k -th iterations k -th step; which is equal to $a_{ij}^{(k)} - a_{ik}^{(k)} \frac{a_{kj}^{(k)}}{a_{kk}^{(k)}}$, and you can see that we have summing we have j goes from $k+1$ to n , where we have taken advantage of the fact that the matrices symmetric, because of round off the computed values of each of the quantities on the right hand side will differ from this true values and result in additional errors. So, instead of $a_{ij}^{(k+1)}$ what is actually stored in the computer is $\bar{a}_{ij}^{(k+1)}$, where $\bar{a}_{ij}^{(k+1)}$ includes the floating point

errors right it includes the round off errors. So, it is the floating point approximation right.

So, $\bar{a}_{ij}^{(k+1)}$ is equal to $\bar{a}_{ij}^{(k)} - m_{ik} \bar{a}_{kj}^{(k)}$. So, here when we compute m_{ik} from $\bar{a}_{ik}^{(k)}$ divided by $\bar{a}_{kk}^{(k)}$. Since this is the floating point operation, we introduce certain errors certain round off errors right, and this round off error is denoted by $1 + \delta_1$. So, this is the error in the computation of m_{ik} then we subtract $m_{ik} \bar{a}_{kj}^{(k)}$ from $\bar{a}_{ij}^{(k)}$. So, this operation the subtraction operation introduces additional floating point error, which is denoted by $1 + \delta_2$ and then on top of let me let me take a step back.

So, this subtraction operation introduces additional floating point error which is given by $1 + \delta_3$, and this operation $m_{ik} \bar{a}_{kj}^{(k)}$ introduces floating point of error which is given by $m_{ik} \bar{a}_{kj}^{(k)} (1 + \delta_2)$. So, this operation introduces of floating point error which is given by $1 + \delta_2$ this subtraction operation introduces a floating point error, which is given by $1 + \delta_3$ and this division operation introduces a floating point error which is given by $1 + \delta_1$; and we are guaranteed that each of these errors it must be less than or equal to the machine precision.


(Refer Slide Time: 12:53)

Bound on round off error

- We assume that it is possible to get exactly the same transformed values at the end of step 'k', by performing **exact** computations on perturbed values of $a_{ij}^{(k)}$
- The perturbations are only applied to the elements belonging to rows that are transformed by the step k (rows with index > k)

$$\hat{a}_{ij}^{(k)} = \bar{a}_{ij}^{(k)} + \varepsilon_{ij}^{(k)}$$

$$\therefore \bar{a}_{ij}^{(k+1)} = \hat{a}_{ij}^{(k)} - \frac{\hat{a}_{ik}^{(k)}}{\bar{a}_{kk}^{(k)}} \bar{a}_{kj}^{(k)} = \bar{a}_{ij}^{(k)} + \varepsilon_{ij}^{(k)} - \bar{m}_{ik} \bar{a}_{kj}^{(k)} (***)$$



$$\text{where } \bar{m}_{ik} = \frac{\bar{a}_{ik}^{(k)} + \varepsilon_{ik}^{(k)}}{\bar{a}_{kk}^{(k)}}$$

Now, we say that if we do these operations, then we are going to get these floating point errors right. So, instead of that we say that let us see if we can get exactly the same

transform values of the end of step k ; by performing the exact computations on perturbed values of a_{ij} . So, the idea is like this.

So, we operate on the values that are stored in the computer and we operate we go through this operation this operation, and we end up with additional errors due to the which are governed by $\delta_1 \delta_2 \delta_3$. So, now, we are saying we are saying that instead of working with the instead of assuming that there are instead of assuming that there are operating on two float no sorry. I am going back I have instead of getting instead of performing on the exact values a_{ijk} , what I am saying is that I am going to operate on the some perturbed values on the perturbed values; and I am going to going to go through the exact same transformation equations, but I will assume that I am not introducing any additional floating point errors during the transformations right.

So, I am going to operate on some perturbed values instead of operating on the exact a_{ijk} a_{ij} superscript k . I am going to operate on certain perturbed values, but I will assume that after by I will get the same results by performing exact operations on the perturbed values. So, the perturbations; however, only applied to the elements belonging to the rows that are transformed by the step k that is rows with index greater than k . So, basically I am saying that instead of operating on a_{ijk} . I will operate on $\bar{a}_{ijk} + \epsilon_{ijk}$. So, now I am get I am going to get $\bar{a}_{ijk} + 1$ is equal to $\hat{a}_{ijk} - \hat{a}_{ik} / \bar{a}_{kk} \times \bar{a}_{kj}$ \bar{a}_{k} these I have already be obtained. So, they have superscript k right. So, they have already been obtained from the previous step in the Gaussian elimination right.

So, these values are known now what I am saying is that the values that I am going to transform right, I am going to I am not going to operate on those values themselves. I am going to operate on those values plus some perturbed values, and I am going to assume that my floating point operations. I am not going to introduce any errors. So, I am going to get the exact solution by operating on these perturbed values perturbed values and I am going to get the exact same solution. I hope that is cleared but basically the idea is that instead of looking at the effect of the floating point operations, we say that we are trying to get this same solution by considering perturbations in my original matrix components.

So, I get some value after my operations right those operations, typically include floating point operations operating on the original numbers original numbers that was there on the matrix. So I say that instead of operating on the original numbers on the matrix, I will operate on some original numbers plus perturbations on some perturb numbers, but during the operations I will not introduce any floating point operation any floating point errors, and I want to get the same values as a result of this operation. So, the idea is that instead of we are transferring the problem to the perturbations. So, we want to find what perturbations in my original system will give me the same errors same floating point errors, as I would get during round off right. So, so what changes should I make to the my original matrix elements, in order to get the same error which I would have got due to round off.

(Refer Slide Time: 18:15)

Bound on round off error


- Hence from (*):
$$\varepsilon_{ik}^{(k)} = \bar{a}_{ik}^{(k)} \delta_1$$

Rewriting (**) as :

$$\bar{m}_{ik} \bar{a}_{kj}^{(k)} = \frac{\bar{a}_{ij}^{(k)} - \bar{a}_{ij}^{(k+1)}}{(1 + \delta_2)} \text{ and substituting in (***) we get :}$$

$$\varepsilon_{ij}^{(k)} = \bar{a}_{ij}^{(k+1)} (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1}) - \bar{a}_{ij}^{(k)} (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1})$$

$$\therefore \left| \varepsilon_{ij}^{(k)} \right| = \left| \begin{array}{l} \bar{a}_{ij}^{(k+1)} (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1}) \\ - \bar{a}_{ij}^{(k)} (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1}) \end{array} \right|$$

$$\leq \max \left(\left| \bar{a}_{ij}^{(k+1)} \right|, \left| \bar{a}_{ij}^{(k)} \right| \right) \left| (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1}) - (1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1}) \right|$$


From this expression \bar{m}_{ik} is equal to \bar{a}_{ik} by $\bar{a}_{kj}^{(k+1)}$ plus δ_1 . We therefore, get ε_{ik} is equal to \bar{a}_{ik} times δ_1 how do we get this well we compare this expression this expression with sorry, this expression this expression with this expression right, when we compare this expression with this expression we get ε_{ik} is equal to \bar{a}_{ik} times δ_1 . Then we can rewrite this expression this expression as this is just a question of substitution right we are going to substitute those values here and we get this and finally, substituting all this in this expression we are going to get finally.

This expression which you can see gives me an expression for the perturbation epsilon i j k, which gives me an expression for the perturbation epsilon i j k. So, that is the basic purpose of this exercise the basic purpose of this exercise is to try to find bounds on the perturbation epsilon i j k, because we know that the perturbations are equivalent the end result of the perturbations is going to be the same errors, which would have a proved if I had the floating point errors right. So, instead of finding try to find bounds on the floating point errors themselves, I am going to try to find bounds on the perturbations. So, this is just I am transferring the problem transferring the problem of finding the bounds on the perturbations, because I am saying the end result of the perturbations is equivalent to the floating point errors in the operations.

So, we get this and then if we take bounds on both sides we get an expression like this and this must be lesser than or equal to. So, this is a bar i j k plus one times this term minus a bar i j k times this term. So, this has to be lesser than or equal to maximum of this and this times this term right, because this is maximum this and this must be larger than this minus this times this minus this times this must be lesser than the maximum of this and this times this minus this right. So, this is lesser than that.


(Refer Slide Time: 21:08)

Bound on round off error

- Then:
$$\left| e_{ij}^{(k)} \right| \leq \max \left(\left| \bar{a}_{ij}^{(k+1)} \right| \left| \bar{a}_{ij}^{(k)} \right| \right) \left\{ \left| 1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1} \right| + \left| 1 - (1 + \delta_2)^{-1} \right| \right\}$$
- Since $1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1} \approx 1 - (1 - \delta_3)(1 - \delta_2) \approx \delta_2 + \delta_3$

$$\left| 1 - (1 + \delta_3)^{-1} (1 + \delta_2)^{-1} \right| \approx |\delta_2 + \delta_3| \leq |\delta_2| + |\delta_3| \leq 2u$$
- Similarly since $1 - (1 + \delta_2)^{-1} \approx 1 - (1 - \delta_2) \approx \delta_2$

$$\left| 1 - (1 + \delta_2)^{-1} \right| \leq |\delta_2| \leq u$$
- We finally get:
$$\left| e_{ij}^{(k)} \right| \leq \max \left(\left| \bar{a}_{ij}^{(k+1)} \right| \left| \bar{a}_{ij}^{(k)} \right| \right) (3u) \quad [i, j > k]$$



So, we get the same expression I have written here. So, mod of epsilon i j is lesser than or equal to max of this times this and then we know that since delta 1 delta 2 delta 3 are less than the machine precision u. So, these must be small numbers. So, I can do a

binomial expansion of this and if I do that I can write $1 + \delta_3 - 1$ as $1 - \delta_3 + \delta_2 - 1$ as $1 - \delta_2$ and this is approximately equal to $\delta_2 + \delta_3$. Where I have ignored terms which involved $2\delta_2$ is δ_2 times δ_3 . Similarly. So, if I take bounds on that I get this is approximately equal to $\delta_2 + \delta_3$ which is lesser than or equal to $\delta_2 + \delta_3$. And since both δ_2 and δ_3 are lesser than u this must be less than two times the machine precision.

Similarly, the second term $1 + 1 + \delta_2$ inverse I can write it as $1 - 1 - \delta_2$, again using binomial expansion taking into account the fact that δ_2 has is very small much smaller than 1, which is going to give me approximately δ_2 . So, again I take bounds on that this going to be lesser than or equal to δ_2 and again δ_2 is lesser than the machine precision. So, that is going to be less than u . So, we finally, get δ_{ij} to the power not to the power δ_{ij} at the k -th step is lesser than or equal to $2 \max$ of this times 3 times u where u is the machine precision.

(Refer Slide Time: 23:12)

Bound on round off error

- Recalling (***) ,

$$\bar{a}_{ij}^{(k+1)} = \bar{a}_{ij}^{(k)} + \varepsilon_{ij}^{(k)} - \bar{m}_{ik} \bar{a}_{kj}^{(k)} \quad (i, j > k)$$


Summing for $k = 1, 2, \dots, r$ where $r = \min(i-1, j)$, we get :

$$\sum_{k=1}^r \bar{a}_{ij}^{(k+1)} - \sum_{k=1}^r \bar{a}_{ij}^{(k)} = e_{ij} - \sum_{k=1}^r \bar{m}_{ik} \bar{a}_{kj}^{(k)} \quad \text{where } e_{ij} = \sum_{k=1}^r \varepsilon_{ij}^{(k)}$$

This can be written as :

$$\bar{a}_{ij}^{(1)} = a_{ij} = \bar{a}_{ij}^{(r+1)} + \sum_{k=1}^r \bar{m}_{ik} \bar{a}_{kj}^{(k)} - e_{ij} \quad (****)$$

since no computation and hence no round off in $\bar{a}_{ij}^{(1)} = a_{ij}$



Let us recall this equation this equation, which was my update formula for Gaussian elimination and which says that $\bar{a}_{ij}^{(k+1)}$ is equal to $\bar{a}_{ij}^{(k)} + \varepsilon_{ij}^{(k)} - \bar{m}_{ik} \bar{a}_{kj}^{(k)}$ if we sum this expression for k is equal to 1 to r , where r is the minimum of $i - 1$ and j being the row index j being the column index. We are going to get $\sum_{k=1}^r \bar{a}_{ij}^{(k+1)} - \sum_{k=1}^r \bar{a}_{ij}^{(k)} = e_{ij} - \sum_{k=1}^r \bar{m}_{ik} \bar{a}_{kj}^{(k)}$ where $e_{ij} = \sum_{k=1}^r \varepsilon_{ij}^{(k)}$

j , where e_{ij} is basically I have summing this term from k equal to one to $r = i - j$ minus $\sum_{k=1}^{r} \bar{m}_{ik} \bar{a}_{kj}$.

This everything is going to cancel except the r for except the term which is going to be for k equal to r , which is going to give me a \bar{a}_{ij}^{r+1} and the term which involves one \bar{a}_{ij}^1 . So, I have a $\bar{a}_{ij}^{r+1} - \bar{a}_{ij}^1$ the rest of the terms are going to cancel, the rest of the terms from this first term is going to cancel the rest of the terms from this second term. So, we are left with a $\bar{a}_{ij}^{r+1} - \bar{a}_{ij}^1$, but a \bar{a}_{ij}^1 is going to be a y_j ; because that is the first that this that is the first the first step and the first step there are no round off errors. So, a \bar{a}_{ij}^1 is going to be equal to a y_j . So, we can get an expression like a y_j is equal to a \bar{a}_{ij}^{r+1} bringing changing the sides right bringing this to the left hand side. So, we get a $\bar{a}_{ij}^{r+1} - y_j = \sum_{k=1}^{r} \bar{m}_{ik} \bar{a}_{kj} - e_{ij}$.

(Refer Slide Time: 25:35)


Bound on round off error

- For terms at or above principal diagonal i.e. $j \geq i, r = (i - 1)$
 $\therefore \bar{a}_{ij}^{(r+1)} = \bar{a}_{ij}^{(i)}$
- For terms below the principal diagonal, $i > j, r = j$. Hence
 $\bar{a}_{ij}^{(r+1)} = \bar{a}_{ij}^{(j+1)} = 0$
- Hence we can write (****) as:

$$a_{ij} = \bar{a}_{ij}^{(i)} + \sum_{k=1}^{i-1} \bar{m}_{ik} \bar{a}_{kj}^{(k)} - e_{ij} \quad j \geq i$$

$$= 0 + \sum_{k=1}^j \bar{m}_{ik} \bar{a}_{kj}^{(k)} - e_{ij} \quad i > j$$

Assuming $\bar{m}_{ii} = 1, a_{ij} = \sum_{k=1}^p \bar{m}_{ik} \bar{a}_{kj}^{(k)} - e_{ij}$ where $p = \min(i, j)$



For terms at or above the principal diagonal j is greater than or equal to i and since r is equal to minimum of $i - 1$ and j if j is greater than i . So, r must be equal to $i - 1$, because r is minimum of $i - 1$ and j so, r is equal to $i - 1$. So, in that case a \bar{a}_{ij}^{r+1} is going to be a \bar{a}_{ij}^1 , because r is equal to $i - 1$. So, this gives me a $\bar{a}_{ij}^{r+1} - \bar{a}_{ij}^1$ is equal to a $\bar{a}_{ij}^1 - \bar{a}_{ij}^1$ for terms below the principal diagonal the row index is going to be greater than the column index i is going to be greater than j therefore, minimum of $i - 1$ and j is going to be j . So, r is going to be j hence in that case a \bar{a}_{ij}^{r+1}

plus 1 is equal to $\bar{a}_{jj} + 1$, which we know from our Gaussian elimination is going to be 0. So, beyond the j -th step the \bar{a}_{ij} is going to be 0 the terms which are below the principal diagonal are going to be 0.

Hence we can write this previous expression this previous expression we can write it as \bar{a}_{ij} which is equal to $\bar{a}_{ij} + 1$ is equal to $\bar{a}_{ij} + 1$ plus this term, which does not change this is true for j greater than or equal to i and this is equal to zero plus this term when i is greater than j . So, basically I have split it up if split this equation in to two parts one for j greater than or equal to i and 1 for i greater than j and if we assume that \bar{m}_{ii} is equal to 1, I can put this term inside the summation and change this index from k equal to one to $i - 1$ to k equal to 1 to i .

So, in that case we can write \bar{a}_{ij} is equal to $\sum_{k=1}^{i-1} \bar{m}_{ik} \bar{a}_{kj} - e_{ij}$ provided \bar{m}_{ii} is equal to 1 that is true, then we can write combine these two equations together to write \bar{a}_{ij} is equal to $\sum_{k=1}^p \bar{m}_{ik} \bar{a}_{kj} - e_{ij}$, where p is equal to minimum of i and j if j is greater than i then p is going to be i when i is greater than j then p is going to be j . So, when p when j is greater than i , I am going to recover this first equation, when i is greater than j . I am going to recover the second equation. So, i finally get \bar{a}_{ij} is equal to $\sum_{k=1}^p \bar{m}_{ik} \bar{a}_{kj} - e_{ij}$.

(Refer Slide Time: 28:45)


Bound on round off error

- In the above equation, only components \bar{m}_{ik} at or below the principal diagonal or components $\bar{a}_{kj}^{(k)}$ at or above the principal diagonal are involved
- Assuming the components of \bar{m}_{ik} above principal diagonal and components of $\bar{a}_{kj}^{(k)}$ below diagonal are zero, we can write:

$$A = \bar{L}\bar{U} - E$$

$$E = (e_{ij}), \bar{L} = \begin{bmatrix} \cdot & \cdot & 0 \\ \bar{m}_{i1} & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}, \bar{U} = \begin{bmatrix} \cdot & \cdot & \bar{a}_{kj}^{(k)} \\ 0 & \cdot & \cdot \end{bmatrix}$$

Thus the computed matrices $\bar{L} = (\bar{m}_{ik})$ and $\bar{U} = (\bar{a}_{kj}^{(k)})$ are the exact triangular factors of the matrix $A + E$



So, in the above equation we can see that only components m_{ik} are needed when they are m_{ik} terms, we only involve the m_{ik} terms which are below the principal diagonal and a_{kj} terms, which are above the principal diagonal we can see because k is equal to 1 to p and p is equal to minimum of i, j , here p is the second index and k is and k can be as high as only as p and p is bounded by this. So, this term only involves a term, which are at or below the principal diagonal while here this term. Since the first index is k and k can be only as high as p and p is bounded by minimum of i, j . So, this term only involves terms which are above the principal diagonal. So, in that because of that this equation only involves components of m_{ik} at or below the principal diagonal and components of a_{kj} at or above the principal diagonal.

Since we do not use components of m which are which are above the principal diagonal and components of a , which are below the principal diagonal we can assume them to be 0 in which case this expression this expression is a $L U$ decomposition this is the lower triangular matrix that is an upper triangular matrix. So, in that case we are going to get, a is equal to $L U$ minus E where E is the matrix is components. I given by E_{ij} let us go back to the previous slide this components i given by E_{ij} and i bar, what we just discussed from is made of m is components are m_{ik} and u bar has components a_{kj} .

Thus the computed matrices L bar and U bar are the exact triangular factors of the matrix a plus e . So, LU we thought was the triangular decomposition of a right, but because of round off errors, we are going to get L bar U bar and L bar U bar is the exact decomposition of a right a plus E it is the exact decomposition of a plus E . So, now we are transfer the error into a perturbation in a right. So, a the original matrix a plus some perturbation matrix E is going to give me L bar U bar L bar U bar includes the effects of all the floating point errors.

(Refer Slide Time: 32:07)


Bound on LHS

- Recalling that

$$e_{ij} = \sum_{k=1}^r \varepsilon_{ij}^{(k)} \quad r = \min(i-1, j) \text{ is the sum of } \min(i-1, j)$$

quantities and $\left| \varepsilon_{ij}^{(k)} \right| \leq 3.u \cdot \max_k \left(\left| \bar{a}_{ij}^{(k)} \right|, \left| \bar{a}_{ij}^{(k+1)} \right| \right)$

Hence, $\left| e_{ij} \right| \leq 3.u \min(i-1, j) \max_k \left(\left| \bar{a}_{ij}^{(k)} \right| \right)$



So, recall that we have defined e_{ij} is the sum of all the $\varepsilon_{ij}^{(k)}$ is that is when what we defined exactly. Here when we defined E_{ij} equal to $\sum \varepsilon_{ij}^{(k)}$ and r is equal to minimum of $i-1$ and j . So, it we can write this as the sum of minimum $i-1$ and j quantities and we also remember that $\varepsilon_{ij}^{(k)}$ has this bound which we obtained earlier. So, we can say that $\left| \varepsilon_{ij}^{(k)} \right| \leq 3.u \cdot \max_k \left(\left| \bar{a}_{ij}^{(k)} \right|, \left| \bar{a}_{ij}^{(k+1)} \right| \right)$ then we get a bound on e_{ij} which is three u times are times are minimum of $i-1$ and j , because this sum is over r and r is minimum of $i-1$ and j . So, we are taking this value which is the maximum value of r right r is bounded r is given by this.

So, this value times the maximum of all these sums right this is the maximum of $\max_k \left(\left| \bar{a}_{ij}^{(k)} \right|, \left| \bar{a}_{ij}^{(k+1)} \right| \right)$ and that is maximum of $\left| \bar{a}_{ij}^{(k)} \right|$ over all the k s. So, we are taking the maximum possible. So, we have a sum of sum of $\varepsilon_{ij}^{(k)}$ which involve terms like this right maximum of this and this. So, we are taking we are looking at all the sums and we are saying that we are going to take of all the terms which comprise the sums we are going to take the maximum value that this maximum of all those values. So, that maximum value times $\min(i-1, j)$ is going to be an upper bound on $\left| e_{ij} \right|$.

So, that. So, that is just the absolute upper bound because we this term comprises a number of sums. So, we are taking the largest term in that sum we are taking the largest term in that sum and multiplying it multiplying that by the number of terms in the sum. So, that is going to be greater than or equal to mod of e i j. So, again let me repeat we this term comprise is a sum of a number of terms right we are taking the largest term in that sum and multiplying it where the number of terms in that sum, and we are saying that has to be a bound on the left hand side which is the sum of all those terms.


(Refer Slide Time: 35:25)

Bound on rhs

- Also since the original system is of the form $Ax = b$ and the rhs is transformed in exactly the same way as the columns of A:
 we have : $b_i = \sum_{k=1}^i \bar{m}_{ik} \bar{b}_k^{(k)} - c_i$ where $c_i = \sum_{k=1}^{i-1} \varepsilon_i^{(k)}$
 and $|\varepsilon_i^{(k)}| \leq 3.u. \max(\bar{b}_i^{(k)}, \bar{b}_i^{(k+1)})$
 Thus $|c_i| \leq 3.u.(i-1) \max_k(\bar{b}_i^{(k)})$
- Introducing the intermediate vector $\mathbf{y} = \bar{U}\mathbf{x}$ we solve the system $\bar{L}\bar{U}\mathbf{x} = \mathbf{b} + \mathbf{c}$:

$$\bar{L}\mathbf{y} = \mathbf{b} + \mathbf{c}$$

$$\bar{U}\mathbf{x} = \mathbf{y}$$



So, that was the bound on the left hand side let us consider, a bound on the right hand side since the original system is of the form $a x$ is equal to b , and the right hand side is transformed in a it remember we recall from Gaussian elimination the right hand side is transformed in exactly the same way as the left hand side right. So, the right hand side because the right hand side is transformed in exactly the same way as the columns of a , we can develop a bound on the error in the right hand side, because we are performing operations on the right hand side every time at every step, we are operating on the left hand side as well as in the right hand side.

So, when we when we operate on the right hand side we introduce floating point errors on the right hand side, and these are the bounds on the similar to the errors we get on the bounds on the errors, we get on the bounds we get on the errors in the left hand side we can similarly get bounds on error on the right hand side in a very similar fashion, and

they have a similar form. So, mod of c_i is the sum of the errors on the right hand side and that is bounded by something like this which is very similar to the bound on the error in the right hand side.

So, introducing the intermediate vector y is equal to $U \bar{x}$ we solve the system $L \bar{y} = U \bar{x} - b$ plus c right now $L \bar{y} = U \bar{x} - b$ plus c , why do we have this c because c is because of the round off errors in b right. So, this is the actual system that we are solving and introducing an intermediate vector y is equal to $U \bar{x}$ we get $L \bar{y} = U \bar{x} - b$ plus, and then we solve for y here we solve for y from this equation put y on the right hand side and finally, we solve for x $U \bar{x} = y$ equal to y .

(Refer Slide Time: 37:40)

Bound on rhs


- From the bounds:

$$|e_{ij}| \leq 3u \cdot \min(i-1, j) \max_k(\bar{a}_{ij}^{(k)})$$

$$|c_i| \leq 3u \cdot (i-1) \max_k(\bar{b}_i^{(k)})$$

it is clear that it is important to limit the values of $|\bar{a}_{ij}^{(k)}|$ and $|\bar{b}_i^{(k)}|$ and the goal of any pivoting strategy should be to limit their growth.

- The size of the multipliers $|\bar{m}_{ik}^{(k)}|$ is seen to have no effect on the magnitude of the round off errors during Gauss elimination




So, summarizing we have obtained bounds on e_{ij} and on the right hand side c_i and these bounds have the following form. Now, it is clear that if these errors are going to be small then these bounds this right these bounding values must be small two right. So, in order to bound these errors these terms $\bar{a}_{ij}^{(k)}$ and $\bar{b}_i^{(k)}$ must be small. So, it is very important that the pivoting strategy whatever pivoting strategy we adopt should limit the values of mod of $\bar{a}_{ij}^{(k)}$ and $\bar{b}_i^{(k)}$, but what is interesting to note is that these bounds do not involve the multipliers $\bar{m}_{ik}^{(k)}$, thus the magnitude of the multipliers has no effect on the magnitude of the round off errors during Gaussian elimination; because these terms these terms do not involve the multipliers they only

involved coefficients, and the right hand side the approximations to coefficients on the right hand side it do not involve the multipliers.

(Refer Slide Time: 39:04)

Error due to back substitution

- We have thus far obtained bounds on the round off error in the elimination step i.e. we have obtained bounds on \mathbf{E} where $\mathbf{E} = \overline{\mathbf{L}\mathbf{U}} - \mathbf{A}$ and \mathbf{c} where $\mathbf{c} = \overline{\mathbf{b}} - \mathbf{b}$
- However the full solution of the linear system also involves a back substitution step following Gaussian elimination
- For most matrices occurring in practice, the round off error accumulating due to back substitution is negligible compared to that during elimination. Thus the errors \mathbf{E} and \mathbf{c} during the elimination step limit the accuracy of the computed soln $\mathbf{Ax}=\mathbf{b}$



We have thus far obtained bounds on the round off error in the elimination step that is we have obtained bounds on \mathbf{e} where \mathbf{e} is equal to $\overline{\mathbf{L}\mathbf{U}}$ minus \mathbf{a} . So, $\overline{\mathbf{L}\mathbf{U}}$ is the approximation to $\mathbf{L}\mathbf{U}$ $\mathbf{L}\mathbf{U}$ is the exact decomposition of \mathbf{A} . So, $\overline{\mathbf{L}\mathbf{U}}$ minus \mathbf{a} is equal to \mathbf{e} write that \mathbf{e} matrix and we have obtained Boundson the terms of the \mathbf{e} matrix each term of the \mathbf{e} matrix is E_{ij} . We obtained Boundson the error matrix in the in the in the coefficient and we have also obtained errors on \mathbf{c} bounds on \mathbf{c} where \mathbf{c} is the error in the right hand side.

So, that. So, we have obtained bounds both on the coefficient matrix as well as on the right hand side errors due to round off, we know now know that those error due to round off cannot exceed these bounds right. So, those errors due to round off cannot exceed these bounds we have found bounds on those errors. However the full solution of the linear system also involves the back substitution following Gaussian elimination. So, this is the bound on the errors due to Gaussian elimination, there will be additional errors due to back substitution which follows Gaussian elimination.

However, for most matrices that occur practice the round off error due to Gaussian due to back substitution is negligible compared to the round off error, due to Gaussian elimination thus the errors \mathbf{E} and \mathbf{c} during the elimination step limit the accuracy of the

computed solution $Ax = b$. So, the most of the errors round off errors occur during the elimination step and if we obtained bounds on those errors whatever additional errors have proved during back substitution.

They are going to be negligibly small compared to they are still going to be bounded by bound which we have obtained during for Gaussian elimination. So, let us look again at those bounds mod of a_{ij} lesser than or equal to three times machine precision U times minimum 5×10^{-16} then maximum of $\bar{a}_{ij}^{(k)}$, overall the iterations k similarly mod of c_i less than or equal to three times machine precision u times $I - 1$ times maximum of $\bar{b}_i^{(k)}$ over all the iterations k these are basically. What are known as posteriori bounds why because we cannot predict those bound before, we have actually computed these values right because these values are the maximum over all the all the steps in the Gaussian elimination.

So, unless we have actually compute unless we have actually gone through the steps of the Gaussian elimination we do not know what these max values are we do not know what $\max_k \bar{a}_{ij}^{(k)}$ is or $\max_k \bar{b}_i^{(k)}$ is. So, these are posteriori bounds posteriori bounds and these have to be these can only be obtained after we have perform the Gaussian elimination. However it is desirable to obtained a priori bounds that is we want to know what my error will be I want to know what my bounds in my what the bounds in my error will be before I have actually started doing the computations right. So, a priori bounds are always more useful than posteriori bounds.


So, let us see how we can n obtained a priori bounds on the error in Gaussian elimination it has been shown that a priori bounds can only be obtained if mod of m_{ik} won not go into the derivation for that, but let us take it for granted it can be there are references and the references I have referred to in this course you can find discussions on this. So, it can be shown that if mod of m_{ik} is always lesser than or equal to 1 then we can obtained a priori bounds. But when is mod of m_{ik} will going to less than or equal to one when are those multipliers going to be less than or equal to one they are going to be less than or equal to one. Only when we have performed pivoting when at least we have performed partial pivoting right at every step, we choose the largest element in a column as the pivot only then can be assured that the multipliers, that we are going to use are always going to be lesser than or equal to one and in that case we can obtained a priori bounds.

(Refer Slide Time: 44:17)

A priori bounds

- It can be shown that if partial pivoting is performed at each step the bound on the error in A:
$$|e_{ij}| \leq 2u \cdot \min(i-1, j) \max_{i,j,k} (\bar{a}_{ij}^{(k)})$$
- The last factor on the right is computed as $\max_k (\max_{i,j} |\bar{a}_{ij}^{(k)}|)$
- If we can bound this computed value in terms of some a-priori values, then we can obtain an a-priori bound on the round off error.

Let us denote $g_{ij} = \frac{\max_{i,j,k} |\bar{a}_{ij}^{(k)}|}{\max_{i,j} |a_{ij}|}$



So, in the in case we perform partial pivoting then we can show that mod of e i j now be it changes. So, it the bound changes the bound becomes narrow instead of this lesser than or equal three u times this, now we are going to get this lesser than or equal two u times minimum of i minus one j, but most interesting is this part. Let's look at what this was this was maximum of a bar i j k over all the steps k ah right. So, now, instead of that I have maximum of a bar i j k over i j k. So, not only is the maximum over the steps k it is over all the elements of the matrix right it is over all the elements of the matrix i j k.

So, now I am saying that mod of e i j is lesser than or equal to the largest value in the coefficient matrix not only over the steps, but over each element not it is not only concerned with each element it is it is the maximum of over all the elements in that matrix it is a maximum of i j as well as k. So, it is not only a maximum of over the steps, it is also a maximum over the row and column indices. So, basically it is the maximum absolute maximum element in that matrix over all the steps. It is the maximum element in that matrix over all the steps. So, the last factor on the right is computed as maximum of i j mod of a bar i j k. So, at each step k we find the largest element in the matrix and then we find the largest over all the steps, so maximum of k over i and j. So, if we can bound this computed value in terms of some a-priori values then we can obtain an a-priori bound on the round off error.

So, if we can we can bound this in terms of some a-priori values in terms of some values which are known which are known before I do my Gaussian elimination typically some values, which are part of my original matrix a right using some components of my original matrix a if I can obtain a bound on this I am going to get in a-priori bound on my error due to Gaussian elimination to do, that is to do that let us denote g_n is equal to maximum of $|a_{ij}|$ divided by maximum of $|a_{ij}|$.


(Refer Slide Time: 47:22)

A priori bounds

- Then the error matrix is bounded by:

$$|E| \leq 2 \cdot u \cdot g_n \max_{i,j} |a_{ij}| \begin{bmatrix} 0 & 0 & 0 \dots & 0 & 0 \\ 1 & 1 & 1 \dots & 1 & 1 \\ 1 & 2 & 2 \dots & 2 & 2 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 2 & 3 \dots n-2 & n-2 & n-2 \\ 1 & 2 & 3 \dots n-1 & n-1 & n-1 \end{bmatrix}$$
- The elements of the matrix on the right hand side are each given by $\min(i-1, j)$. Hence the maximum norm of the matrix is given by the sum of the elements in the last row:

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



If we do that we can write this expression we can re write this expression as mod of e mod of e lesser than or equal two u two u remains the same maximum of $|a_{ij}|$. I am going to replace this by g_n times maximum of $|a_{ij}|$. So, that is going to be g_n times maximum of $|a_{ij}|$ and what is this matrix this matrix is nothing, but minimum of $i-1$ and j . So, at for each element I have computed I have obtained this element by taking minimum of the row index minus one and the column index. So, the row index here is one the column index is one. So, I have one minus one is 0 0 and 1 minimum of 0 and 1. So, that is going to be 0. Similarly, I have computed all the rest of the terms in this matrix right which is nothing, but minimum of $i-1$ and j elements of the matrix on the right hand side or each given by minimum of $i-1$ and j .

Hence the matrix norm of the maximum norm of the matrix is given by the sum of the elements in the last row recall what is the maximum norm it is the sum of the sum of the elements in each row maximum of that right. So, I have taken. So, that. So, that the my

last row is going to contribute that and that is going to be 1 plus 2 plus 3 plus n minus 1 plus n minus 1. I have group these terms together you can see these are the sum of the first n natural numbers which is given by half n n plus one. So, this is the half n n plus one minus one. So, we get mod of E infinity why is why have we taken e infinity because we have computed the infinite norm of this matrix right.

(Refer Slide Time: 49:33)

A priori bounds


- Hence:

$$|E|_{\infty} \leq 2.u.g_n \max_{i,j} |a_{ij}| \left\{ \frac{1}{2}n^2 + \frac{1}{2}n - 1 \right\}$$
- By slightly refining the estimate of E, it can be shown that

$$|E|_{\infty} \leq 2.u.g_n \max_{i,j} |a_{ij}| \left\{ \frac{1}{2}n^2 \right\} = n^2 g_n .u \max_{i,j} |a_{ij}|$$

But $\max_y |a_{ij}| \leq |A|_{\infty} \left[|A|_{\infty} = \max_i \sum_{j=1}^n |a_{ij}| \right]$

$$|E|_{\infty} \leq n^2 g_n .u .|A|_{\infty} \quad (*)$$



So, mod of e infinity is lesser than or equal to 2 u g n times max i j over mod of i j this term remains the same and then the infinite norm of this matrix infinite norm of this matrix, which we just calculated to be half of n square plus half of n minus 1 by slightly refining the estimate of E, which I am not going to go into again it is slight change actually we can show that mod of e infinity is lesser than or equal to 2 u g n max i j mod of i j times half of n square.

So, it is actually here we have obtained half of n square plus half of n minus 1 actually it is less than half of n square this is a sharper bound right it can be shown that we can write it like this and then 2 2 cancels out we get n square g n u maximum of i j mod of i a i j, but maximum of a mod of a i j i j is lesser than or equal the infinite norm of my original coefficient matrix why because the infinite norm is defined as sum of over all the columns sum of sum of each row sigma j equal to one to n mod of a i j and then the maximum of that right. So, maximum of i j mod of a i j is going to be bounded by the infinite norm of a right.


So, I can replace this on the right hand side by the infinite norm of a because this is greater than that. So, I finally, get this mod of e infinity is lesser than or equal to n square g n u times a infinite norm of a, but we still have this factor g n right what is this factor g n.

(Refer Slide Time: 51:38)

A priori bounds

- Hence, if \bar{L} and \bar{U} are the computed triangular factors with **partial or complete pivoting**, then $E = \bar{L}\bar{U} - A$ the matrix of round off errors is bounded as given in (*)
- The bound on E is satisfactory unless the ratio g_n is large.
- It becomes an a priori bound for an a priori estimate of g_n . For partial pivoting it has been shown $g_n \leq 2^{n-1}$. For complete pivoting $g_n \leq 1.8^{.25 \ln n}$ which gives a narrower bound for large n

Even for partial pivoting g_n rarely exceeds 8.




Let's recall g_n we define to be this we define this factor g_n to be this and it can be it has been it can by just by slowing a large number of problems finding out the typical values for g_n , it has been shown, that for partial pivoting g_n is lesser than or equal to 2 to the power n minus 1, where n is the dimension of the matrix and for complete pivoting g_n is lesser than or equal to 1 point 8 to the power point 25 l n, which gives a narrower bound, but even if we do just partial pivoting this bound is sufficiently narrow because g_n rarely exceeds eight right.

So, the bound on e is satisfactory, because g_n is usually quite small g_n rarely exceeds 8 g_n rarely exceeds 8 so in that case, we have we have obtained bounds on e bar on E \bar{L} \bar{U} bar and \bar{U} bar are the computed triangular factors with partial or complete pivoting why do we say partial or completing pivoting, because we have try to find a a-priori bound a-priori bound right, and we have mentioned that we can only obtain an a-priori bound and all those multipliers are less than 1. So, partial pivoting is necessary so this bound is satisfactory unless the ratio g_n is large.

(Refer Slide Time: 53:38)

The effect of scaling

- It has been seen that scaling of the unknowns or of the coefficient matrix has no effect on accuracy of the computed solution except by affecting the choice of pivots.
- Thus a suitable scaling, by increasing the size of the pivots may improve the stability of a solution or vice versa.
- Pivoting, as is clear from the error bound obtained earlier, does not affect the accuracy of the solution.




So, that was. So, we have found a bound on the Gaussian error in Gaussian elimination due to round off some people have often try to improve the solution improve the results of Gaussian elimination they often try scaling, basically they either scale the coefficient matrix or they scale the right hand side. And they assumed they sometimes think that before because their doing scaling they somehow going to get better solutions they are going to get smaller errors, but it has been seen that scaling the unknowns or the coefficient matrix has no effect on the accuracy of the computed solution except by affecting the choice of the pivots. So, it can improve the stability of the solution because pivot is going to the size of the pivot is going to determine the stability of the solution. So, if the pivots are large we are going to get more stable solution the pivots are small the error this solution is going to be less table, but it is not going to effect the accuracy of the solution the accuracy of the solution is not going to be effected by scaling.

(Refer Slide Time: 54:55)

Iterative improvement of soln

- When the coefficient matrix A is ill conditioned, but not excessively so, the Gauss elimination solution can be improved (errors reduced) by performing iterations.
- Recall, that extremely ill conditioned matrices may lead to incorrect solutions even though the residual is small. However if the residual vector is not too small, it is possible to improve the solution by iteration for a moderately ill-conditioned matrix.
- Let $\mathbf{r} = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}}$ be the residual for computed solution $\bar{\mathbf{x}}$
Then $\mathbf{A}(\mathbf{x} - \bar{\mathbf{x}}) = \mathbf{r}$



So, that was that was discussion on the error estimate due to round off during Gaussian elimination. Next from next lecture onwards we are going to talk about iterative methods for solving linear systems, but before we talk about iterative methods for instance Gauss Seidel iterations or Jacobi iterations. We will briefly talk about a simple technique to improve the solution of direct methods improve the solution of obtain using direct methods such as Gaussian elimination by some iterations. So, we solve the problem using Gaussian elimination, but then we can improve the solution by doing some simple iterations. So, we will talk about that first before going into directly into iterative methods for solving Gaussian for solving linear systems.

And we will find that these iterative methods are particularly suited for problems which have coefficient matrix which are very spars right which are very spars because Gaussian elimination, we know one of the problem with Gaussian elimination is that it destroys sparseness of this system right. So, so we want if the matrix is really spars we want to use iterative methods because those methods preserve this sparseness of this system. So, we can take advantage of this sparseness in reducing the number of mathematic mathematical operations as well as the storage required.

Thank you.