

**Numerical Analysis.**  
**Professor R. Usha.**  
**Department of Mathematics.**  
**Indian Institute of Technology, Madras.**  
**Lecture-48.**  
**Practical Problems.**

(Refer Slide Time: 0:45)

Example 1.  $A = \begin{pmatrix} 1 & 2 \\ 1 & -1 \end{pmatrix}$  In this case, we can compute the eigen values:  $|A - \lambda I| = 0 \Rightarrow \begin{vmatrix} 1-\lambda & 2 \\ 1 & -1-\lambda \end{vmatrix} = 0$   
 $-(1-\lambda)(1+\lambda) - 2 = 0 \Rightarrow -1 + \lambda^2 - 2 = 0 \Rightarrow \lambda = \pm\sqrt{3}$ .

From the rows of matrix A,  $|\lambda - 1| \leq 2 \rightarrow I$   
and  $|\lambda + 1| \leq 1 \rightarrow II$

(\*)  $\rightarrow$  eigenvalues obtained through actual calculation.

Good morning everyone, in the last class we learned how to compute Gerschgorin bounds which help us to specify the location of the eigenvalues of a given matrix A. So now we shall consider some examples and see how we can obtain the Gerschgorin bounds. Let us start with a simple example of a 2 by 2 matrix A. So here the matrix A has entries 1, 2, 1, -1 and we can actually compute the eigenvalues in this case. The eigenvalues are given by determinant of A - lambda I equal to 0. So if you compute the determinant of this matrix and equate it to 0 and solve the corresponding characteristic polynomial for lambda, then you get lambda to be + or - root 3.

So you know that one of the eigenvalues of matrix A is + root 3 and the other eigenvalue is - root 3. Now let us apply the Gerschgorin's theorem results and see what information that we get about the location of these 2 eigenvalues. So the 2 eigenvalues + - root 3 have actually been marked on the real line, they are real eigenvalues, I have marked them here by means of stars. This is lambda equal to - root 3 and this is lambda equal to + root 3. Now let us use Gerschgorin's theorem. So from the rows of matrix A I consider the 1<sup>st</sup> row, this is A<sub>11</sub> and Gerschgorin's theorem tells us that modulus of lambda - A<sub>11</sub> which is 1 is less than or equal to sum of the absolute values of the other entries in that row.

So here we have only one element which is 2, so the 1<sup>st</sup> result else that absolute value of  $\lambda - 1$  is less than or equal to 2, this gives us one circle and some information about an eigenvalue of the matrix A can be obtained from here. Then taking the 2<sup>nd</sup> row we have  $\lambda - (-1)$  which is  $A_{22}$  that this  $\lambda + 1$  in absolute value must be less than or equal to the sum of the absolute values of the other entries, here in this case it is 1, so modulus of  $\lambda + 1$  is equal to 1 is going to be the 2<sup>nd</sup> Gerschgorin circle that we have.

So we observe that for the matrix A we have considered, there are now 2 discs in the complex plane and the Centre of the discs are  $A_{11}$  and  $A_{22}$  which appear in the matrix A. So each centred on one of the diagonal entries of the matrix A.

(Refer Slide Time: 3:21)

For the matrix  $A_{2 \times 2}$ , there are two disks in the complex plane, each centered on one of the diagonal entries of the matrix A.

From the theorems we know that every eigenvalue must lie within one of these circular disks. However, it does not say that each disk has an eigenvalue.

Example 2  $A = \begin{pmatrix} 1 & -1 \\ 2 & -1 \end{pmatrix}$

$$|A - \lambda I| = \begin{vmatrix} 1-\lambda & -1 \\ 2 & -1-\lambda \end{vmatrix} = 0$$

$$\Rightarrow -(1-\lambda)(\lambda+2) = 0$$

$$\Rightarrow \lambda = -1, -2$$

has an eigenvalue.

Example 2  $A = \begin{pmatrix} 1 & -1 \\ 2 & -1 \end{pmatrix}$

$$|A - \lambda I| = \begin{vmatrix} 1-\lambda & -1 \\ 2 & -1-\lambda \end{vmatrix} = 0$$

$$\Rightarrow -(1-\lambda)(\lambda+2) = 0$$

$$\Rightarrow \lambda = -1, -2$$

The circular disks are  $|\lambda - 1| \leq 1 \rightarrow \text{I}$   
 $|\lambda + 1| \leq 2 \rightarrow \text{II}$

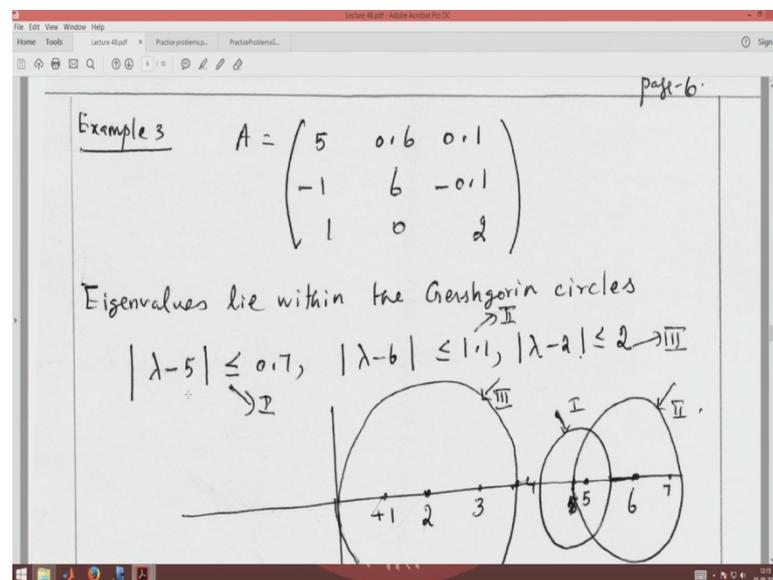
(\*  $\rightarrow$  eigenvalues obtained through actual calculation)

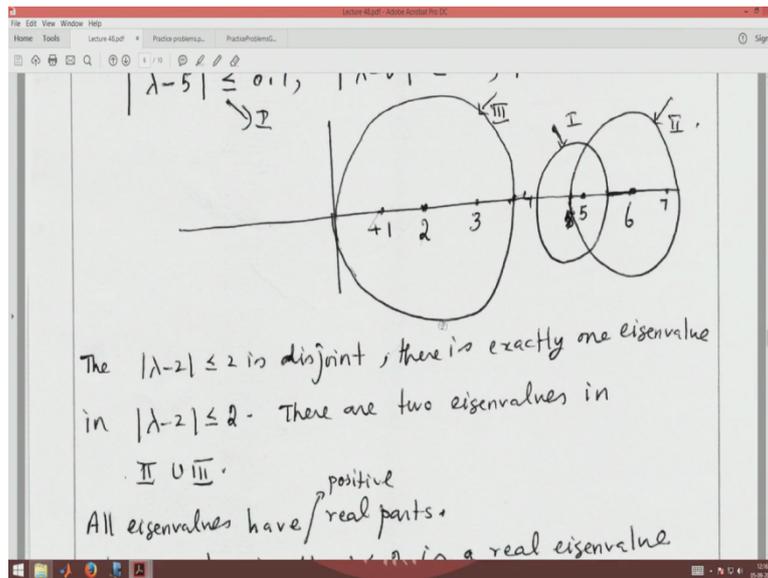
We observe that all the eigenvalues lie within the circular disk  $|\lambda + 1| \leq 2$  defined by the second row.

And from the theorems we know that every eigenvalue must lie within one of these circles and however we do not have any information about the fact that the each disc contains an eigenvalue. So this information is not obtained from Gerschgorin theorem. So let us see another example. So in the 2<sup>nd</sup> example, again I consider a 2 by 2 matrix, as before we can consider the eigenvalues of this matrix A, they can be actually computed by solving the characteristic polynomial even by determinant of A - lambda I equal to 0 and we observe that the eigenvalues are given by lambda equal to + - I.

So those eigenvalues are actually plotted here, this is lambda equal to i had this is lambda equal to - i, they are marked by star here in the complex plane. Now let us see what the Gerschgorin's theorems tell us. The Gerschgorin's theorems tell us that the circles within which the eigenvalues + or - i given by modulus of lambda - A<sub>11</sub> where A<sub>11</sub> is 1 is less than or equal to modulus of -1 that is 1. And modulus of lambda - of -1, that is lambda of modulus of lambda +1 is less than or equal to 2 and that gives us the 2<sup>nd</sup> circle. So 1 and 2 are Gerschgorin's circles and the eigenvalues i and - i lie in the union of these 2 circles is what is given by Gerschgorin's theorem.

(Refer Slide Time: 5:27)





The  $|\lambda - 2| \leq 2$  is disjoint, there is exactly one eigenvalue in  $|\lambda - 2| \leq 2$ . There are two eigenvalues in  $II \cup III$ . All eigenvalues have positive real parts. The eigenvalue in  $|\lambda - 2| \leq 2$  is a real eigenvalue since complex eigenvalues must occur in conjugate pairs. We now apply Gerschgorin theorem to  $A^T$ . Eigenvalues must lie within the Gerschgorin circles  $|\lambda - 5| \leq 2$ ,  $|\lambda - 6| \leq 0.6$  and  $|\lambda - 2| \leq 0.2$

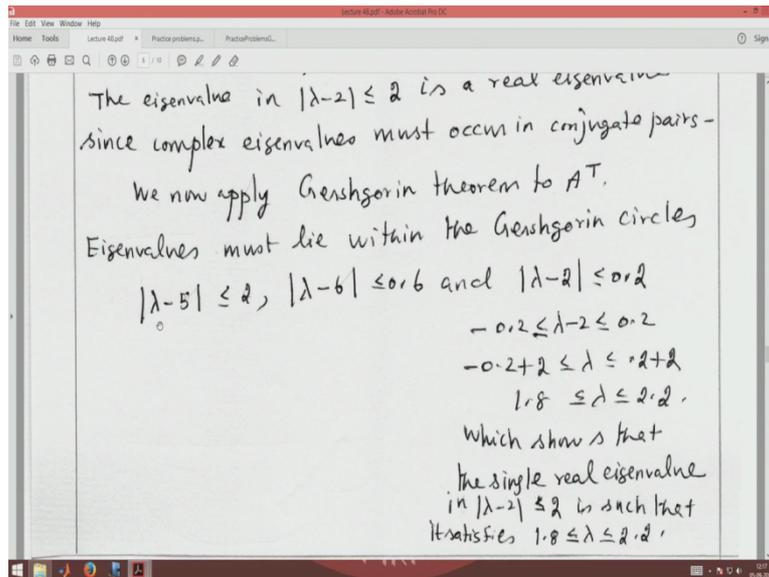
$-0.2 \leq \lambda - 2 \leq 0.2$   
 $-0.2 + 2 \leq \lambda \leq 2 + 0.2$

And we observe in this case that all the Eigen values, namely  $+i$  and  $-i$  lie within the circular disk given by modulus of lambda +1 less than or equal to 2 and that is defined by the 2<sup>nd</sup> row. Right, let us now consider another example. Namely consider a 3 by 3 matrix A and here Gerschgorin's Circle theorems tell us that all the eigenvalues lie within the Gerschgorin circles given by mod lambda -5 is less than equal to the sum of the absolute values of the other 2 entries in the row and that is 0.7. Then modulus of lambda -6 which is this is less than or equal to the sum of the absolute values of the other 2 entries that is 1.1 and finally modulus of lambda -2 is less than equal to 1 +0 and therefore we have the Gerschgorin's circles are given by these.

So finally we see that the circle which has its centre at 2 and radius 2 is disjoint from the other 2 circles which clearly tells us that there is exactly one eigenvalue in this circle which is

disjoint from the others. And there are 2 more eigenvalues and these 2 Eigen values lie in the union of other 2 circles which we have given. And we also observe that all the Eigen values are real because they have a positive real part. And the 2<sup>nd</sup> thing is therefore the circle, modulus of lambda -2 less than equal to 2 is such that it has an eigenvalue which is a real eigenvalue since complex eigenvalues, if they occur, they must after in conjugate pairs.

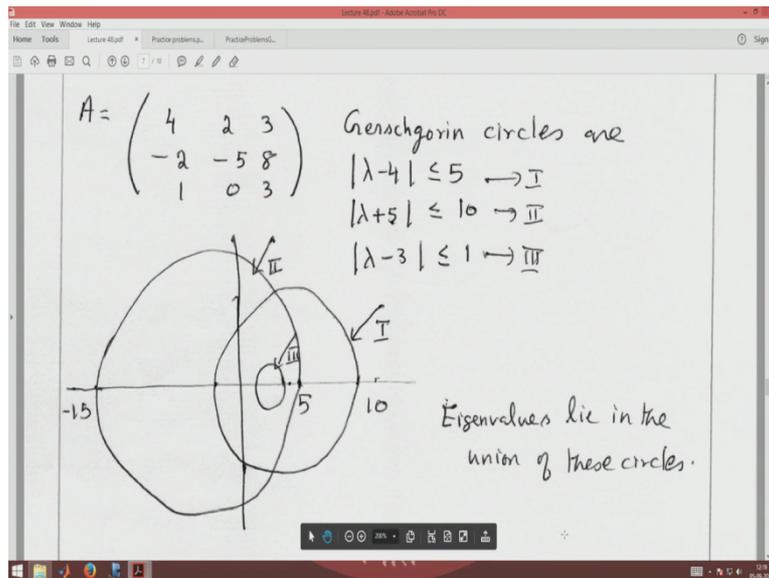
(Refer Slide Time: 7:38)



And so this circle, the disjoint circle encloses one and only one eigenvalue which is the real eigenvalue. Now to get more information, we apply Gerschgorin's theorems to the transpose of the matrix A, that is A transpose. Then we observe that eigenvalues must lie within the Gerschgorin circles which are given by modulus of lambda -5 less than or equal to the sum of the absolute values in the 1<sup>st</sup> column, that is 2 and modulus of lambda -6 is less than or equal to the sum of the absolute values in the 2<sup>nd</sup> column which is 0.6 and finally if you use the 3<sup>rd</sup> column of matrix A, that gives you mod lambda -2 is less than or equal to 0.2.

Right. So you observe that this is the disjoint circle with Centre at the point 2 and this result which is obtained from the columns of P tell us that that single eigenvalue which lies in the disjoint circle must be such that lambda must be greater than or equal to 1.8 and less than or equal to 2.2 and so lambda lies in the interval 1.8 to 2.2. So it clearly gives you an interval within which this single eigenvalue lie within that disjoint circle. So we are able to get the location of the eigenvalues of a given matrix A by applying Gerschgorin circle theorem in terms of rows and in terms of the columns of the matrix A or the rows of the corresponding matrix which is A transpose.

(Refer Slide Time: 9:10)



So we now consider another example in which say  $A$  is again a 3 by 3 matrix. Immediately we can write down the Gerschgorin circles, what are they, they are given by modulus of  $\lambda - 4$  less than or equal to  $2+3$  which is 5, modulus of  $\lambda - (-5)$  that is  $\lambda + 5$  is less than or equal to modulus of  $-2+8$  which is less than 10 and modulus of  $\lambda - 3$  is less than or equal to modulus  $1+0$  that is one. So these are the 3 circles and the eigenvalues of this given matrix lie within the union of these 3 circles. So you just plot these circles because you know the Centre of the circles are going to be respectively 4, 5 and 3, I mean 4, -5 and 3.

And the radii of the circles are going to be 5, 10 and 1 respectively, so plot the circles in the complex plane and you know the eigenvalues of the given matrix lie within the union of these 3 circles in this case.

(Refer Slide Time: 10:22)

$$A = \begin{pmatrix} 15 & -3 & 1 \\ 1 & -2 & 1 \\ 1 & 6 & 1 \end{pmatrix}$$

The eigenvalues of  $A$  lie within the union of three disks  $|z-15| \leq 4$ ,  $|z+2| \leq 2$ ,  $|z-1| \leq 7$ .

The circle centered on  $z=15$  is entirely disjoint from circles II & III, and  $\therefore$  must contain one eigenvalue and the other two eigenvalues lie in the union of circles II and III.

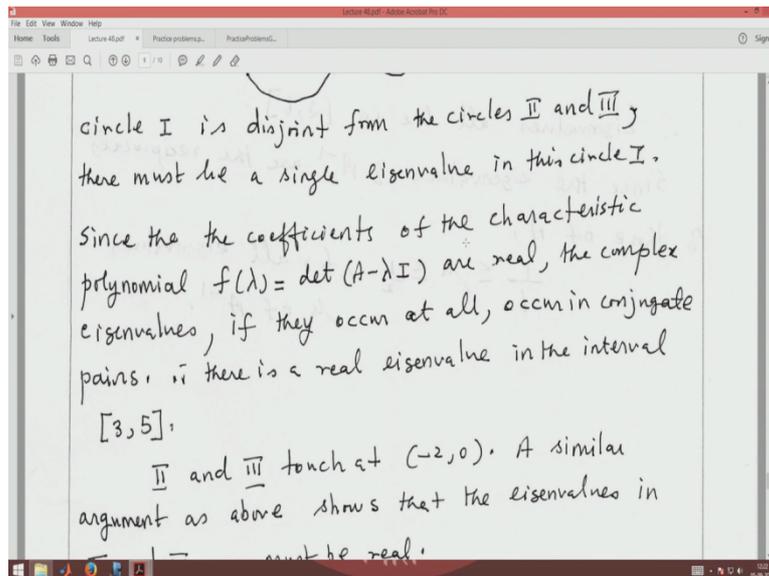
So let us now take another example. So again  $A$  is a 3 by 3 matrix and it is clear that the Eigen values of the matrix  $A$  lie within the union of the circles mod  $Z -15$  less than or equal to 4, mod  $Z +2$  less than or equal to 2 and mod  $Z -1$  less than or equal to 7. And the one with Centre at 15 and radius is equal to 4 is going to be a circle which is the joint from the other 2 circles. So as before we argue and say and conclude that that circle which is disjoint from the other circles must contain one eigenvalue and the other 2 Eigen values lie in the union of the other 2 circles namely mod  $Z + 2$  less than or equal to 2 and mod  $Z -1$  less than or equal to 7.

(Refer Slide Time: 11:35)

$$A = \begin{pmatrix} 4 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & 1 & -4 \end{pmatrix}$$

Gershgorin's circles are  $|\lambda-4| \leq 1$ ,  $|\lambda| \leq 2$ ,  $|\lambda+4| \leq 2$ .

Circle I is disjoint from the circles II and III, there must be a single eigenvalue in this circle I.



So the Gerschgorin's theorem is very very useful in understanding the location of the eigenvalues of a given matrix A. So let us look into another example, suppose A is given by these 3 by 3 matrix, we can write down the Gerschgorin circles, so  $\text{mod } \lambda - 4$  is less than equal to 1,  $\text{mod } \lambda - 0$ , that is  $\text{mod } \lambda$  less than equal to 2 and  $\text{mod } \lambda +$  for is lesser equal to 2, so these are the rather circles which we have plotted here. And we observe that this circle 1 is disjoint from the other 2 circles and therefore it must contain a single eigenvalue within it.

And we also observe that if we write down the characteristic polynomial, the coefficients of the characteristic polynomial given by determinant of  $A - \lambda I$  equal to 0 are real and hence the complex values which are eigenvalues, if they occur at all, then they can in conjugate pairs. Therefore there is a single real eigenvalue within this interval namely 3 to 5 because this is a disjoint circle, right, it single eigenvalue lies within it and radius of the circle is 1, so that single real, single eigenvalue is real and that real eigenvalue lies between 3 and 5. So we know the interval within which this real eigenvalue lies.

(Refer Slide Time: 13:20)

pairs, if there is a real eigenvalue  $\lambda$  in  $[3, 5]$ ,  
 $\text{II}$  and  $\text{III}$  touch at  $(-2, 0)$ . A similar argument as above shows that the eigenvalues in  $\text{II}$  and  $\text{III}$  ~~are~~ must be real.  
 $\lambda = -2$  is not an eigenvalue, since

$$|A - \lambda I| = \begin{vmatrix} 4+2 & 1 & 0 \\ 1 & 2 & -1 \\ 1 & 1 & -2 \end{vmatrix} = 6(-4+1) - 1(-2+1) \\ = 6(-3) - 1(-1) \\ = -18 + 1 = -17 \neq 0$$

There is one eigenvalue in  $[-6, -2)$  and one in  $(-2, 2]$ .

$\text{II}$  and  $\text{III}$  ~~are~~ must be real.  
 $\lambda = -2$  is not an eigenvalue, since

$$|A - \lambda I| = \begin{vmatrix} 4+2 & 1 & 0 \\ 1 & 2 & -1 \\ 1 & 1 & -2 \end{vmatrix} = 6(-4+1) - 1(-2+1) \\ = 6(-3) - 1(-1) \\ = -18 + 1 = -17 \neq 0$$

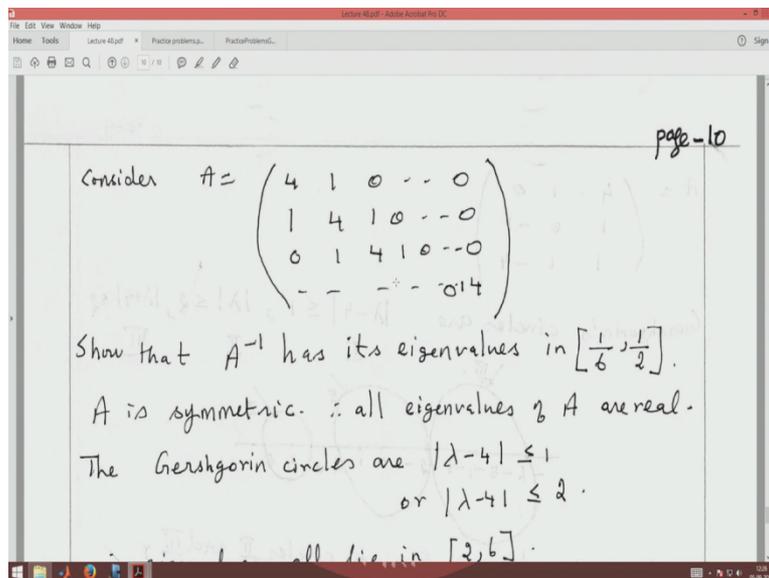
There is one eigenvalue in  $[-6, -2)$  and one in  $(-2, 2]$ .  
 $\therefore A$  has one eigenvalue in each of the intervals,  $[-6, -2)$ ,  $(-2, 2]$ ,  $[3, 5]$ .

And then the other 2 circles, so the other 2 Eigen values lie in the union of these 2 circles is what we know. But we observe that these 2 circles meet each other or touch each other at the point -2. So we would like to see whether -2 is an eigenvalue of the matrix A. How do you check whether lambda equal to lambda 0 is an eigenvalue of the matrix A? The just have to see whether it is the root of the characteristic polynomial, determinant of A - lambda I equal to 0, that is what you have to check. So let us compute determinant of lambda - A - lambda I taking lambda to be -2, so if you write down the determinant and expand and you observe that that turns out to be -17 which is different from 0 and therefore the characteristic polynomial does not have lambda equal to -2 as its root.

And therefore lambda equal to -2 is not an eigenvalue of the matrix A is what we conclude. What does it mean, we said earlier that the other 2 Eigen values lie in the union of those 2 circles and the point where they touch each other, namely -2 is not an Eigen values, all right. And therefore the other 2 Eigen values must be such that each one must lie in each of these circles and that is the conclusion that we get from our discussion. So there is one eigenvalue in the interval -6 to -2 and another eigenvalue which is in -2 to -2. Why does that happen, so let us look at the circles, so this is a circle with Centre at -4 and radius 2.

And therefore it passes through the points -6 in -2, so one of the eigenvalues lies between -6 and -2. And then you have another circle whose centre is the origin and radius is 2 units, and there is an eigenvalue within this circle and it is a real eigenvalue and therefore it lies between -2 and 2. So the matrix A is such that it has one real Eigen value between -6 and -2, another real eigenvalue between -2 and 2 and 3<sup>rd</sup> real eigenvalue between 3 and 5. And so we conclude that A has one eigenvalue in each of these intervals -6 to -2, -2 to 2 and 3 to 5. So information about the location of the eigenvalues can be obtained by writing out the Gerschgorin's circles within which the eigenvalues are located and arguing out depending upon whether these circles are disjoint or the Eigen values lie within the union of these circles and so on.

(Refer Slide Time: 16:14)

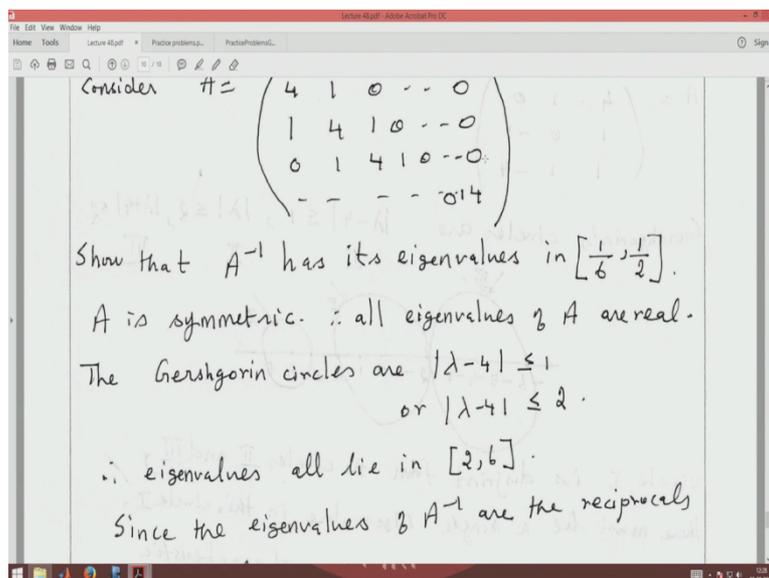


So now we have another problem where A is a matrix say it is an n Cross n matrix having entries given by these. Namely along the diagonals you have 4 appearing throughout and just above the diagonal, entries have values 1 and just below the angle, the entries are again 1, the rest of the entries are all 0. So you have an N cross N matrix where a I J are given by the

entries which are specified here. You are asked to show that the inverse of this matrix  $A$  has its eigenvalues in the interval  $1/6$  to  $1/2$ . So you need some information about the location of the eigenvalues of the inverse of this matrix  $A$ .

We have already seen while doing the Gerschgorin theorems and other details, we have remarked that if  $\lambda$  is an eigenvalue of  $A$ , then  $1/\lambda$  is an eigenvalue of  $A$  inverse, we have actually shown this result. So now we want to show that  $A$  inverse has its eigenvalues interval  $1/6$  to  $1/2$ . So let us  $1^{\text{st}}$  data mean the interval within which eigenvalues of  $A$  lie and then we can use this result and make a conclusion about the location of eigenvalues of the matrix  $A$  inverse.  $1^{\text{st}}$  of all we observe that  $A$  is a symmetric matrix, a  $II$  is a  $JI$ , so what do you immediately conclude, all its eigenvalues are, right, are going to be real.

(Refer Slide Time: 18:40)

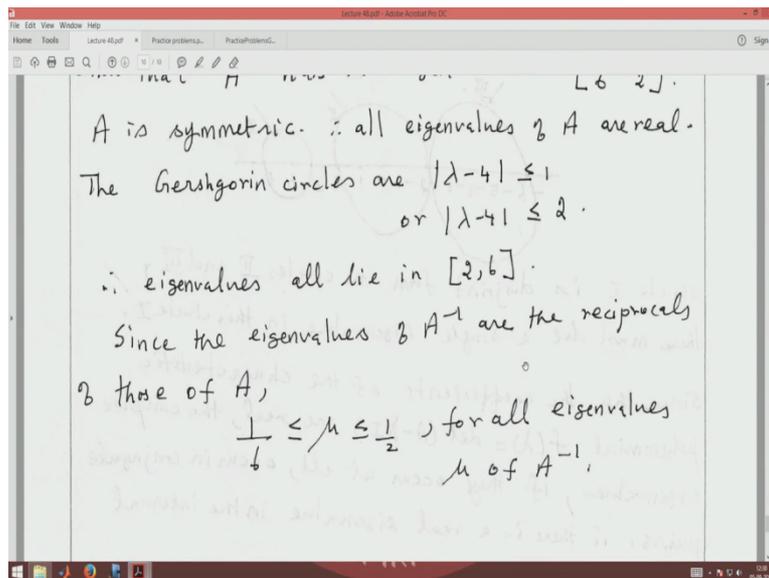


Namely all the Eigen values of matrix  $A$  are going to be real. That is the  $1^{\text{st}}$  information that we get by looking at what the matrix  $A$  and matrix  $A$  is symmetric. Now let us apply Gerschgorin's circles theorem to the rows of this matrix  $A$ . So every row has 4 along its diagonal and therefore the Gerschgorin circles will be given by modulus of  $\lambda - 4$  is less than or equal to the sum of absolute values of the entries in that row. And we observe that the  $1^{\text{st}}$  row has entry to the right of 4, so the  $1^{\text{st}}$  row will give us the Gerschgorin circle to be modulus  $\lambda - 4$  less than or equal to 1.

Similarly the last row will give us modulus  $\lambda - 4$  is less than or equal to 1 because the rest of the entries are all zeros. Now apply Gerschgorin's theorem to all the other rows, starting from

the 2<sup>nd</sup> row to the N -1th row. We observe that the diagonal entries are 4 and you have on the either side of the diagonal entry, value 1 appears to the right of it and value 1 appears to the left of it. And hence we get modulus of lambda -4 less than or equal to 2 to be the Gerschgorin circles. So what does this mean, if you write out this, this tells you -2 less than equal to lambda -4 less than or equal to 2 and therefore that immediately tells you that all the eigenvalues lie in the interval 2 to 6.

(Refer Slide Time: 20:16)



So what are these eigenvalues? They are the eigenvalues of matrix A. So Gerschgorin circle theorems tell us that all the eigenvalues of A are real and all the eigenvalues of the matrix A lie within the interval 2 to 6. So now we make use of the result that we have shown. What is it, eigenvalues of A inverse a reciprocal of eigenvalues of A. So A has eigenvalues lying in the interval 2 to 6 and therefore if mu denotes the Eigen values of A inverse, then mu has to lie between the reciprocal of 2 and reciprocal of 6, namely mu has to lie between 1 by 6 and 1 by 2.

So now lies between 1 by 6 and 1 by 2 and therefore the conclusion is that all the eigenvalues mu of A inverse lie in the interval 1 by 6 to 1 by 2 and that is what we are asked to show. Right, so we have made use of the 2 results, what are the 2 results, we said that if A is a symmetric matrix, all eigenvalues are real and that secondly we use the fact that eigenvalues of A inverse are given by the reciprocal of eigenvalues of A, so we 1<sup>st</sup> computed the eigenvalues of A and then we used the result and got some information about the location of the eigenvalues of the matrix A inverse.

So these are some examples in Gerschgorin theorems and location of eigenvalues of a given matrix. Now we move onto some examples where we discuss the computation of norm of a matrix and then some problems from the error analysis by the direct methods of solving a system of equation  $AX$  is equal to  $B$ .

(Refer Slide Time: 21:58)

TRUE OR FALSE

The matrix  $A = \begin{pmatrix} 100 & -200 \\ -200 & 401 \end{pmatrix}$  is ill-conditioned.

$\|A\| = \max(300, 601) = 601$ .

$A^{-1} = \frac{1}{\det A} \begin{pmatrix} 401 & 200 \\ 200 & 100 \end{pmatrix} = \begin{pmatrix} 4.01 & 2 \\ 2 & 1 \end{pmatrix}, |A| = \det A = 100$ .

$\|A^{-1}\| = \max(6.01, 3) = 6.01$

$\therefore \kappa(A) = \|A\| \|A^{-1}\| = 601 (6.01) = 3612$

So suppose say I give you a matrix  $A$  which is this and I make a statement, I make the statement that the matrix  $A$  given by this is ill conditioned and I ask you to tell me whether this statement is true or false. It is just not enough if you would say that the statement is true or the statement is false, you have to justify your statement whether it is a true statement or a false statement, then only your answer will be correct and you have justified your conclusion. So how are we going to conclude whether the statement is true or false, how do I say or how do I find out whether the matrix is an ill conditioned matrix or not, I compute what is its condition number.

If the condition number is very large, then we conclude that the matrix  $A$  is ill conditioned. So I compute the condition number. And what is the condition number, condition number is norm  $A$  into norm  $A$  inverse. So given a matrix  $A$ , I compute its inverse and then I compute norm of  $A$ , then norm of  $A$  inverse and then find what  $K$  of  $A$ , the condition number of  $A$  is. So I compute norm  $A$  using maximum column sum norm. I have maximum of  $100 +$  modulus of  $-200$ , so  $300$ , then modulus of  $-200 + 401$ , that is  $601$ , so norm  $A$  is  $601$ .

So with respect to the maximum columns sum Norm, Norm  $A$  is  $601$ . Now I compute norm  $A$  inverse with respect to the maximum column sum, so norm  $A$  inverse is going to be equal to

maximum of the some of the absolute values of the entries in the 1<sup>st</sup> column which is 6.01 and the sum of the entries in the 2<sup>nd</sup> column after taking the absolute value, right and that gives you 3 and so it is 6.01. And therefore norm A inverse is 6.01. So now compute the condition number of A, it is norm A into norm A inverse, so that is 601 into 6.01 and it turns out to be 3612 which is very much greater than the value 1, right.

For a well conditioned system, the condition number must be close to value 1 greater than 1 and here it turns out to be 3612 and therefore A is an ill conditioned matrix, so the statement that has been given, namely the matrix A is ill conditioned is a true statement. So whenever you are given some statement and we are asked to check whether the statement is true or false, right, in order to conclude what the answer is, you must work out the details and based on the results that you get, you must conclude whether the statement is true or false, in which case you have justified your conclusion.

(Refer Slide Time: 25:23)

The image shows a handwritten mathematical derivation on a digital whiteboard. The text is as follows:

Example. Let  $\epsilon > 0$ .

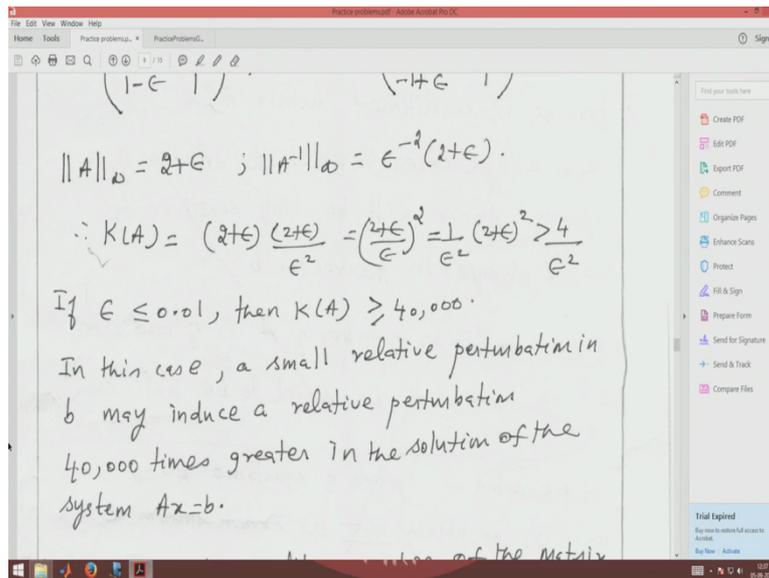
$$A = \begin{pmatrix} 1 & 1+\epsilon \\ 1-\epsilon & 1 \end{pmatrix}, \quad A^{-1} = \epsilon^{-2} \begin{pmatrix} 1 & -1-\epsilon \\ -1+\epsilon & 1 \end{pmatrix}$$

$$\|A\|_{\infty} = 2+\epsilon; \quad \|A^{-1}\|_{\infty} = \epsilon^{-2}(2+\epsilon).$$

$$\therefore K(A) = \frac{(2+\epsilon)(2+\epsilon)}{\epsilon^2} = \frac{(2+\epsilon)^2}{\epsilon^2} = \frac{1}{\epsilon^2}(2+\epsilon)^2 > \frac{4}{\epsilon^2}$$

If  $\epsilon \leq 0.01$ , then  $K(A) \geq 40,000$ .

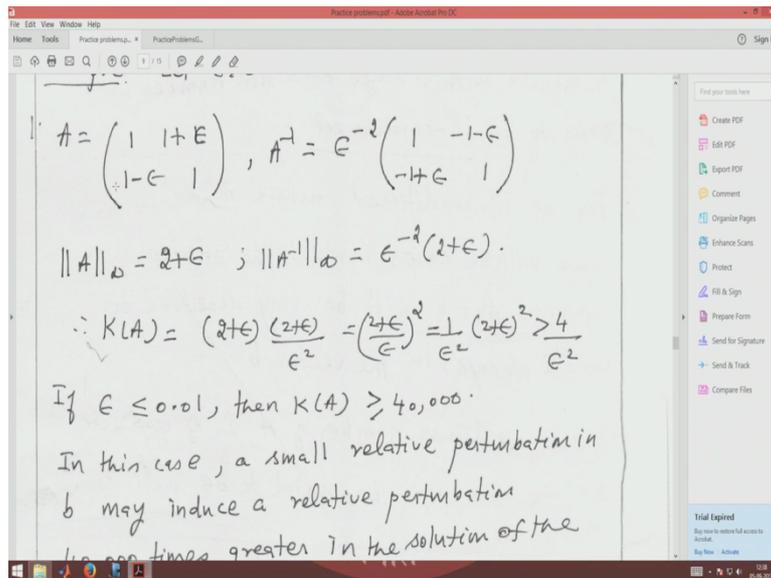
In this case a small relative perturbation in



So let us now consider an example where we want to perform some error analysis. Suppose you are given  $\epsilon$  greater than 0 and a matrix  $A$  having entries 1,  $1 + \epsilon$  in the 1<sup>st</sup> row and  $1 - \epsilon$  and 1 in the 2<sup>nd</sup> row. Immediately I compute  $A$  inverse and that is given there. And this time I would like to obtain norm of  $A$  using the infinity norm. So what is, the condition number is  $\|A\|_\infty$  into  $\|A^{-1}\|_\infty$  and that turns out to be  $2 + \epsilon$  the whole square by  $\epsilon$  square. So this is surely greater than 4 by  $\epsilon$  square.

If suppose in the matrix that is given, I take my  $\epsilon$  which is positive to be such that it is less than or equal to 0.01, then I observe that in that case  $K$  of  $A$  is greater than 4 by 0.01 the whole square and that will turn out to be 40,000. So  $K$ , condition number of  $A$  will turn out to be greater than or equal to 40,000 which is a very very large quantity.

(Refer Slide Time: 27:10)



$$A = \begin{pmatrix} 1 & 1+\epsilon \\ -1-\epsilon & 1 \end{pmatrix}, \quad A^{-1} = \epsilon^{-2} \begin{pmatrix} 1 & -1-\epsilon \\ -1-\epsilon & 1 \end{pmatrix}$$

$$\|A\|_{\infty} = 2+\epsilon; \quad \|A^{-1}\|_{\infty} = \epsilon^{-2}(2+\epsilon).$$

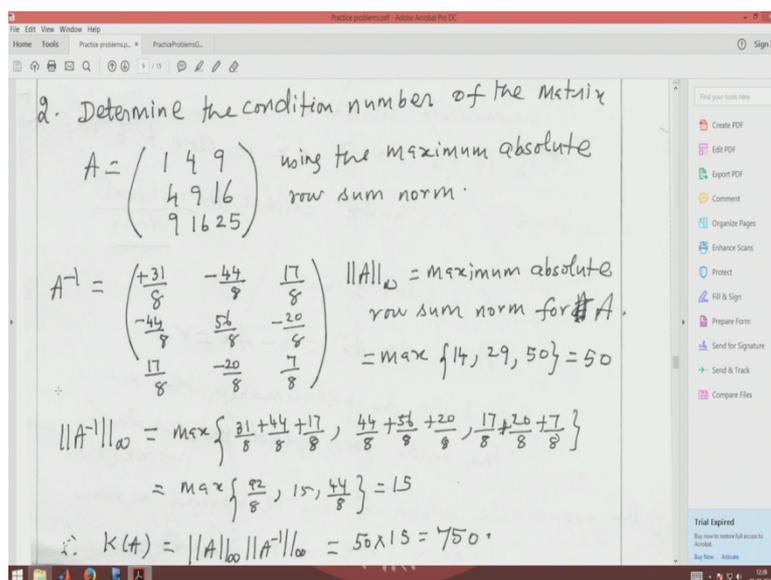
$$\therefore K(A) = \frac{(2+\epsilon)(2+\epsilon)}{\epsilon^2} = \frac{(2+\epsilon)^2}{\epsilon^2} \geq \frac{4}{\epsilon^2}$$

If  $\epsilon \leq 0.01$ , then  $K(A) \geq 40,000$ .

In this case, a small relative perturbation in  $b$  may induce a relative perturbation 40,000 times greater in the solution of the given system  $AX$  is equal to  $B$ .

Suppose I give a small, suppose I give a small change in one of the entries in the matrix  $A$ , namely I add Epsilon to this entry and subtract Epsilon from this entry, so if I perturbed entries in the given matrix  $A$ , then what happens to the condition number. In that case the condition number becomes very very large and so this says that a relative perturbation in the right-hand side may induce a relative perturbation 40,000 times greater in the solution of the given system  $AX$  is equal to  $B$ . So change in some entries in  $A$  will give the corresponding changes in the solution and that is going to be such that it is going to be 40,000 times greater in the solution of the given system  $AX$  is equal to  $B$ .

(Refer Slide Time: 28:01)



2. Determine the condition number of the matrix  $A = \begin{pmatrix} 1 & 4 & 9 \\ 4 & 9 & 16 \\ 9 & 16 & 25 \end{pmatrix}$  using the maximum absolute row sum norm.

$$A^{-1} = \begin{pmatrix} \frac{31}{8} & -\frac{44}{8} & \frac{17}{8} \\ -\frac{44}{8} & \frac{56}{8} & -\frac{20}{8} \\ \frac{17}{8} & -\frac{20}{8} & \frac{7}{8} \end{pmatrix}$$

$$\|A\|_{\infty} = \max\{14, 29, 50\} = 50$$

$$\|A^{-1}\|_{\infty} = \max\left\{\frac{31}{8} + \frac{44}{8} + \frac{17}{8}, \frac{44}{8} + \frac{56}{8} + \frac{20}{8}, \frac{17}{8} + \frac{20}{8} + \frac{7}{8}\right\}$$

$$= \max\left\{\frac{92}{8}, 15, \frac{44}{8}\right\} = 15$$

$$\therefore K(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty} = 50 \times 15 = 750$$

So let us now consider what happens if you want to compute the condition number of a matrix A which is a 3 by 3 matrix using the maximum absolute by row sum Norm. So A is given, so you can compute A inverse using Gauss Jordan technique, now you have learnt that. So you know how to compute the inverse of a given matrix using Gauss Jordan technique and now you compute norm A infinity. What is it, it is maximum absolute row sum Norm for A. So it is a maximum of the sum of these entries in each of these rows after taking the absolute value. And that turns out to be maximum of 14, 29 and 50 and that is 50.

Now compute norm A inverse infinity, here again I want to use the maximum row sum and so compute the sum of the entries in absolute values in each of the rows and from among them choose that which is the maximum that turns out to be 15. So norm A inverse infinity is 15, so the condition, condition number of given matrix A is K of A given Norm A infinity into A inverse infinity which is 50 from here and 15 from norm A inverse infinity, that gives you 750. So we can say that the matrix A is an ill conditioned matrix.

(Refer Slide Time: 29:34)

Let  $A = \begin{pmatrix} 2 & 1 \\ 2 & 1.01 \end{pmatrix}$ , with respect to  $\|\cdot\|_2$ ,

a) the condition number of A is \_\_\_\_\_.

b) the matrix is \_\_\_\_\_.

Now  $\|A\|_2 = 3.165$

$\|A^{-1}\|_2 = 158.27$

a)  $\therefore \kappa(A) = \|A\|_2 \|A^{-1}\|_2 = 500.974$ .

b) the matrix is ill-conditioned.

Let  $\epsilon > 0$ ,  $A = \begin{pmatrix} 1 & 1+\epsilon \\ 1 & 1-\epsilon \end{pmatrix}$ ,  $A^{-1} = \begin{pmatrix} 1 & 1 \\ 1 & -1-\epsilon \end{pmatrix}$

Now I consider another example. Taking A to be a 2 by 2 matrix and I want you to suit the condition number of A with respect to the 2 norms. So when I work out the 2 norms for the matrix A, that turns out to be 3.165. I also require the 2 norms for A inverse because I am asked to compute the condition number of A. So norm A inverse 2 turns out to be 158.27. Now that you know the definition of 2 norms for a matrix A, you can compute A inverse and work out the details and show that norm A inverse 2 is 158.27. And therefore the condition number is norm A into norm A inverse with respect to the 2 norms and so it is the product of these 2 values and that turns out to be this.

So you have actually computed it and you are asked to fill in the blanks here, so write out the results here that it is 500.974. And you observe that this is a large quantity and therefore the, you are asked to make some conclusion about what the matrix is. And you know since the condition number is large, the matrix is ill conditioned, so you fill in the blanks here by saying that the matrix is an ill conditioned matrix.

(Refer Slide Time: 31:05)

3. Compute the condition number of the following matrix relative to  $\|\cdot\|_\infty$ .

$$A = \begin{pmatrix} \frac{1}{2} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{4} \end{pmatrix}; \|A\|_\infty = \max \left\{ \left| \frac{1}{2} \right| + \left| \frac{1}{3} \right|, \left| \frac{1}{3} \right| + \left| \frac{1}{4} \right| \right\}$$

$$= \max \left\{ \frac{5}{6}, \frac{7}{12} \right\}$$

$$= \max \{ 0.83333333, 0.58333333 \}$$

$$= 0.83333333.$$

$$|A| = \frac{1}{8} - \frac{1}{9} = \frac{9-8}{72} = \frac{1}{72}.$$

$$= 0.83333333.$$

$$|A| = \frac{1}{8} - \frac{1}{9} = \frac{9-8}{72} = \frac{1}{72}.$$

$$A^{-1} = 72 \begin{pmatrix} \frac{1}{4} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{1}{2} \end{pmatrix} = \begin{pmatrix} 18 & -24 \\ -24 & 36 \end{pmatrix}$$

$$\|A^{-1}\|_\infty = \max \{ |18| + |-24|, |-24| + |36| \}$$

$$= \max \{ 42, 60 \} = 60$$

$$\therefore \text{condition number } K(A) = \|A\|_\infty \|A^{-1}\|_\infty$$

$$= (0.83333333)(60)$$

$$= 49.999998$$

$$= \max\{42, 60\} = 60$$

$$\therefore \text{condition number } K(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty}$$

$$= (0.3333333)(60)$$

$$= 49.999998$$

$$\approx 50$$

$$K(A) \approx 50$$

Compute the condition number of the following matrix relative to infinity norm. So A is a 2 by 2 matrix, norm A infinity is maximum row sum, so compute the row sum, namely the 1<sup>st</sup> row sum in absolute value is and the 2<sup>nd</sup> row sum and you conclude that norm A infinity turns out to be this. Now what is the determinant of A, that turns out to be 1 by 72. Get A inverse and workout what norm A inverse infinity is. Again maximum row sum Norm and so we have maximum of 42 and 60 to give you norm A inverse infinity and that turns out to be 60. So what is the condition number? The condition number is norm A infinity into norm A inverse infinity.

We have computed each of these, so the product of this gives you norm A A infinity into norm A inverse infinity and that is 50 and you are asked to get what the condition number is. So the condition number of A is 50. So now we know how to compute norm of a vector or norm of a matrix associated with the norm of a vector which is given. So let us work out some details about the error analysis for direct method.

(Refer Slide Time: 32:44)

d) The following linear system  $Ax=b$  has  $\bar{x}$  as the actual solution and  $\tilde{x}$  as an approximate solution. Compute  $\|x-\tilde{x}\|_\infty$  and  $\|A\tilde{x}-b\|_\infty$ .

$$\frac{1}{2}x_1 + \frac{1}{3}x_2 = \frac{1}{63} \quad x = \left(\frac{1}{7}, -\frac{1}{6}\right)^T$$

$$\frac{1}{3}x_1 + \frac{1}{4}x_2 = \frac{1}{168} \quad \tilde{x} = (0.142, -0.166)^T$$

$$x - \tilde{x} = \left\| \frac{1}{7} - 0.142, -\frac{1}{6} + 0.166 \right\|_\infty$$

$$\|x - \tilde{x}\|_\infty = \max \left\{ \left| \frac{1}{7} - 0.142 \right|, \left| -\frac{1}{6} + 0.166 \right| \right\}$$

$$= \max \{ 0.00085714, -0.00066667 \}$$

$$x - \tilde{x} = \left\| \frac{1}{7} - 0.142, -\frac{1}{6} + 0.166 \right\|_\infty$$

$$\|x - \tilde{x}\|_\infty = \max \left\{ \left| \frac{1}{7} - 0.142 \right|, \left| -\frac{1}{6} + 0.166 \right| \right\}$$

$$= \max \{ 0.00085714, -0.00066667 \}$$

$$= 0.00085714$$

$$A\tilde{x} = \begin{pmatrix} \frac{1}{2} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{4} \end{pmatrix} \begin{pmatrix} 0.142 \\ -0.166 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(0.142) - \frac{1}{3}(0.166) \\ \frac{1}{3}(0.142) + \frac{1}{4}(-0.166) \end{pmatrix}$$

$$= \begin{pmatrix} 0.071 - 0.05533333 \\ 0.01566667 \end{pmatrix} = \begin{pmatrix} 0.01566667 \\ 0.01566667 \end{pmatrix}$$

So the problem is, the following linear system  $AX$  is equal to  $B$  has  $X$  as the actual solution and  $X$  tilde is an approximate solution. So you are asked to compute the magnitude of the absolute error and also compute what is the residual a namely norm of  $A X$  tilde -  $B$  with respect to the infinity norm. So you are given a system of 2 equations in 2 unknowns. And you observe that these coefficients have 1 by 3, 1 by 3, 1 by 4 are the same as the examples which we have considered earlier. And  $X$  is equal to 1 by 7 - 1 by 6 transpose is the actual solution and  $X$  tilde which is this is an approximate solution.

So solve this system say by applying Gauss elimination method with partial pivoting and get the solution correct to 3 decimal places and that is what your  $X$  tilde. So now you want to compute the error, what is the error, it is  $X - X$  tilde. So  $X$  is 1 by 7,  $X$  tilde, 1<sup>st</sup> component is

0.142. 2<sup>nd</sup> component is -1 by 6 and the 2<sup>nd</sup> component of X tilde is -0.166, so it is this quantity. And you want norm X - X tilde infinity. So it is the maximum of the 1<sup>st</sup> entry and the 2<sup>nd</sup> entry. So maximum of absolute value of the 1<sup>st</sup> entry and the absolute value of the 2<sup>nd</sup> entry, so compute this.

(Refer Slide Time: 35:05)

$$A\tilde{x} = \begin{pmatrix} \frac{1}{2} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{4} \end{pmatrix} \begin{pmatrix} 0.142 \\ -0.166 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(0.142) - \frac{1}{3}(0.166) \\ \frac{1}{3}(0.142) + \frac{1}{4}(-0.166) \end{pmatrix}$$

$$= \begin{pmatrix} 0.071 - 0.05533333 \\ 0.06473333 - 0.0415 \end{pmatrix} = \begin{pmatrix} 0.01566667 \\ 0.00583333 \end{pmatrix}$$

$$A\tilde{x} - b = \begin{pmatrix} 0.01566667 \\ 0.00583333 \end{pmatrix} - \begin{pmatrix} \frac{1}{63} \\ \frac{1}{168} \end{pmatrix}$$

$$= \begin{pmatrix} 0.01566667 \\ 0.00583333 \end{pmatrix} - \begin{pmatrix} 0.01587302 \\ 0.00595238 \end{pmatrix}$$

$$= \begin{pmatrix} -0.00020635 \\ -0.00011905 \end{pmatrix}; \quad \|A\tilde{x} - b\|_{\infty} = \max\{|-0.00020635|, |-0.00011905|\} = 0.00020635$$

4(b) Compute  $K(A)$  relative to  $\|\cdot\|_{\infty}$

So you get this to be maximum of this and this which turns out to be the value 0.00085714, so that is the absolute error in the solution X of the system AX is equal to B. So you would like to find what is the residual, so compute A into X tilde. So A matrix is given, X tilde is given, so work out what is A into X tilde which turns out to be this vector. So what is the residual, residual is AX tilde - B. Why is it a residual, you want X to satisfy the equation AX

is equal to B, so  $AX - B$  must be 0 but  $X$  tilde is an approximate solution, so  $AX$  tilde - B will not be equal to 0, it gives you residual, and so you call r as  $AX$  tilde - B.

Are you computed  $AX$  tilde as this, so that vector - you know what the right-hand side vector is, so substitute that and then evaluate what is the difference and compute norm  $AX$  tilde - B, that is the maximum of the absolute values of the entries that you have written here and that turns out to be 0.00020635. So you will be in a position to compute the absolute error which is norm  $X - X$  tilde as well as the norm of the residual r which is norm  $AX$  tilde - B with respect to any norm that will be asked. In this example you had been asked to get these values using the infinity norm and you have computed them with respect to infinity norm.

(Refer Slide Time: 36:39)

Q(b) Compute  $K(A)$  relative to  $\|\cdot\|_\infty$   
 and  $K(A) \frac{\|b - Ax\|_\infty}{\|A\|_\infty}$

$$K(A) = \|A\|_\infty \|A^{-1}\|_\infty,$$

We have already computed  $K(A)$  relative to  $\|\cdot\|_\infty$   
 $K(A) = 50$  Also  $\|b - Ax\|_\infty = 0.00020635$

we also have computed  $\|A\|_\infty = 0.8333333$

$$\therefore \frac{K(A) \|b - Ax\|_\infty}{\|A\|_\infty} = \frac{(50)(0.00020635)}{0.8333333}$$

$$= \frac{0.0103175}{0.8333333}$$

$$= 0.012381$$

And now you are asked to compute the condition number relative to infinity norm and also the condition number times the residual by norm A with respect to infinity norm. So you compute what is condition number, it is norm A Infinity into norm A inverse infinity. We already have computed the condition number earlier when we worked out the example, it turned out to be 50. And now we have computed what is the residual, namely norm B - AX tilde with respect to infinity norm and that is 0.00020635. And we also have computed norm A with respect to infinity norm which is this so we simply have to substitute these values which appear on the left-hand side, namely these and then simplify and that gives you the value of this namely condition number multiplied by norm of the residual by norm A infinity to be equal to this.

(Refer Slide Time: 38:05)

Suppose  $\tilde{x}$  is an approximation to the solution  $B Ax=b$ ,  $A$  is a nonsingular matrix, and  $\tilde{r}$  is the residual vector for  $\tilde{x}$ . Then, show that

(i) for any natural norm,  

$$\|x - \tilde{x}\| \leq \|r\| \cdot \|A^{-1}\|.$$

(ii) and if  $x \neq 0$ , and  $b \neq 0$ ,  

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|r\|}{\|b\|}.$$

Since residual  $r = b - A\tilde{x}$   

$$= Ax - A\tilde{x} \quad (-: Ax=b)$$

$$= A(x - \tilde{x})$$

$$\therefore A^{-1}r = A^{-1}A(x - \tilde{x}) \quad (-: A \text{ is nonsingular, } A^{-1} \text{ exists})$$

$$= x - \tilde{x}$$

$$\therefore \|x - \tilde{x}\| = \|A^{-1}r\|$$

$$\leq \|A^{-1}\| \|r\| \quad (-: \|AB\| \leq \|A\| \|B\|)$$

which proves result (i).

Now, since  $b = Ax$ ,  $\|b\| = \|Ax\| \leq \|A\| \|x\|$

If we know how to compute the matrix norm given the vector norm, then we will be in a position to compute the absolute error and also the relative error. So that is what is given in the next example. Suppose that  $\tilde{x}$  is an approximation to the solution of  $Ax = b$ ,  $A$  is a non-singular matrix and  $r$  is the residual vector for  $\tilde{x}$ , that is  $r = Ax - b$ . Then you are asked to show that for any natural norm, namely infinity norm or 2 norm, you should be able to show  $\|x - \tilde{x}\| \leq \|A^{-1}\| \|r\|$ .

And in addition, if  $x$  is different from 0, is a nonzero vector, and  $b$  is also a nonzero vector, then show that the relative error in our computation, what is the relative error, it is change in  $x$  with respect to  $x$ , so it is  $\frac{\|x - \tilde{x}\|}{\|x\|}$ , that must be less than or equal to the condition number which is  $\|A\| \|A^{-1}\|$  multiplied by the norm of the residual by the norm of the right-hand side vector, this is what we have to show. So these 2 results clearly tell us the bound on the absolute error, namely this is the right-hand side, your absolute error cannot exceed the value on the right-hand side.

Your relative error in the computation of solution of a system  $Ax = b$  cannot exceed the value on the right-hand side, what is it, it is the condition number of the matrix  $A$  times the norm of the residual by norm of the right-hand side vector  $b$ . So let us work out the details. Now residual is  $r$ , which is  $b - Ax$ . But what do you know about  $b$ ,  $b = Ax$ , so  $Ax - Ax$ , that is  $A(x - \tilde{x})$ . So I pre-multiply both sides by  $A^{-1}$  because it is given  $A$  is a nonsingular matrix, so  $A^{-1}$  exists and therefore I multiply by  $A^{-1}$ .

(Refer Slide Time: 40:50)

The image shows a whiteboard with handwritten mathematical derivations. At the top, it states  $\|r\| \leq \|A\| \|x - \tilde{x}\|$  (labeled as result 1) and  $\|r\| \leq \|A\| \|b\|$  (labeled as result 2). Below this, it says "which proves result (1)". Then, it states "Now, since  $b = Ax$ ,  $\|b\| = \|Ax\| \leq \|A\| \|x\|$ ". This leads to the inequality  $\frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$ . Next, it says "Consider (1):  $\|x - \tilde{x}\| \leq \|A^{-1}\| \|r\|$ ". This leads to the final inequality  $\frac{\|x - \tilde{x}\|}{\|x\|} \leq \frac{\|A^{-1}\| \|r\| \|A\|}{\|b\|} = \frac{\|A\| \|A^{-1}\| \|r\|}{\|b\|}$ , which is labeled as "which proves (2)".

So we have  $\|A^{-1}\|_r$  to be  $\|A^{-1}\|_r$  into  $\|A^{-1}\|_r \|X - \tilde{X}\|_r$ , so that will be  $\|X - \tilde{X}\|_r$ . And therefore  $\|X - \tilde{X}\|_r$  will be  $\|A^{-1}\|_r$  but we have shown already  $\|A^{-1}\|_r \|B\|_r$ , where  $A$  and  $B$  are matrices will be less than or equal to  $\|A\|_r \|B\|_r$ , so let us make use of that result. So  $\|A^{-1}\|_r$  is less than or equal to  $\|A^{-1}\|_r \|r\|_r$ . So this proves our 1<sup>st</sup> result, that is what we have to show, namely  $\|X - \tilde{X}\|_r$  is less than or equal to  $\|r\|_r \|A^{-1}\|_r$  and we have shown that result 1.

Now, we are, we know that  $B = AX$ , so  $\|B\|_r = \|AX\|_r$ , that is less than or equal to  $\|A\|_r \|X\|_r$  and therefore  $\|X\|_r$  is that an equal to  $\|B\|_r$  divided by  $\|A\|_r$ . Now consider  $\|X - \tilde{X}\|_r$  less than or equal to  $\|A^{-1}\|_r \|r\|_r$ . And therefore divided by  $\|X\|_r$  Throughout, so you get  $\|X - \tilde{X}\|_r / \|X\|_r$  is  $\|A^{-1}\|_r \|r\|_r / \|X\|_r$  and what about  $\|X\|_r$ , we know  $AX = B$ , so I also know  $\|B\|_r$  is less than or equal to  $\|A\|_r \|X\|_r$  and so I can replace  $\|X\|_r$  by a quantity which is a miracle to norm that able  $\|B\|_r$ . So I replace this  $\|X\|_r$  by this and this give me  $\|A^{-1}\|_r \|r\|_r \|A\|_r / \|B\|_r$  and that is what we have to show here, namely  $\|A^{-1}\|_r \|A\|_r \|r\|_r / \|B\|_r$ .

We can immediately get the bound on the absolute error and the relative error for a given problem when we get an approximate solution using the direct methods, namely Gauss elimination method or any of the direct methods that we have studied. So given a system  $AX = B$  and if you are if you apply direct methods to solve the system, then you can immediately specify the amount of absolute error and relative error that are involved in your computations by using the bounds which are given by these examples. What is that we have done all along in this course?

We have focused our attention on 5 important topics which are very useful to us in obtaining approximation of a function whose values are specified at discrete points. Namely, we discussed polynomial approximation and then when the information is given at a set of equally spaced points, we used Newton's forward interpolation polynomial or backward interpolation polynomial and obtained approximate polynomial, namely the interpolating polynomial that interpolates the function at a certain of given point. If the information is given at a set of arbitrarily located points, then we had Lagrange interpolation method and divided difference method to compute the interpolating polynomial that interpolates the given data and we have discussed error in interpolation and also the bound on the error in interpolation.

(Refer Slide Time: 44:12)

The truncation error in the four-point backward difference formula used to approximate the second derivative

$$f''(x_i) = \frac{1}{h^2} [f(x_i) - 5f(x_{i-1}) + 4f(x_{i-2}) - f(x_{i-3})]$$

is

A)  $O(h)$     B)  $O(h^2)$     C)  $O(h^3)$     D)  $O(h^4)$

RHS =  $\frac{1}{h^2} [f(x_i) - 5f(x_{i-1}) + 4f(x_{i-2}) - f(x_{i-3})]$

collect coefficient

collect coefficient

$$f_i [0]$$

$$f_{i-1} [0]$$

$$f''_i [1 + \frac{5}{2} - \frac{11}{2} + \frac{9}{2}] = f''_i [0] \checkmark$$

$$f'''_i [\frac{5}{6} - \frac{32}{6} + \frac{27}{6}] = 0 \cdot f'''_i$$

$$\frac{h^4 f^{(4)}_i}{h^2} [\frac{-5}{24} + \frac{64}{24} - \frac{81}{24}] \neq 0 \therefore TE \leq O(h^2)$$

$$\frac{h^4 f^{(4)}_i}{h^2} [\frac{-5}{24} + \frac{64}{24} - \frac{81}{24}] \neq 0 \therefore TE \leq O(h^2)$$

Ans: B

---

page-1.

1. Choose the appropriate answer.

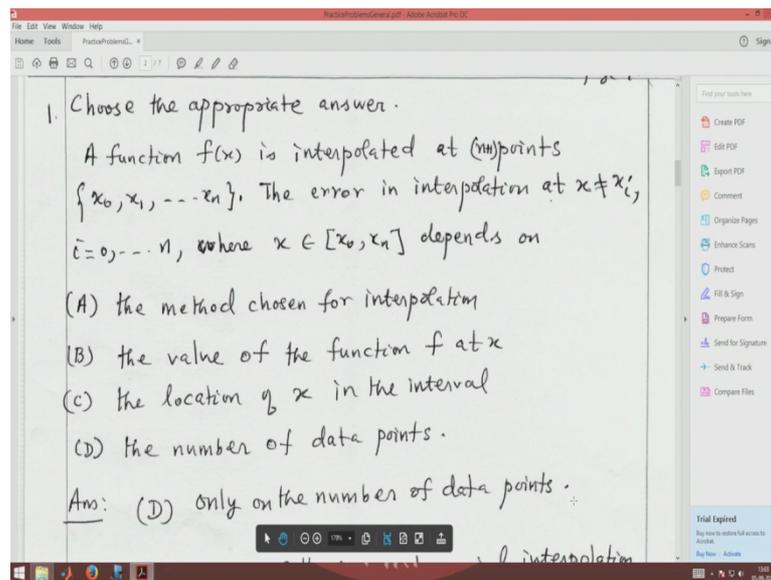
A function  $f(x)$  is interpolated at  $(n+1)$  points in an interval  $[a, b]$  at  $x_0, x_1, \dots, x_n$ .

Then we considered the topic on numerical differentiation. So then we derived finite difference approximations to various ordered derivatives which can be given in terms of the values of the function at a certain set of points. So let us consider an example here. Suppose say I give you a formula which approximates the 2<sup>nd</sup> derivative,  $f''(x)$  is given by this and I asked you the following. Show that the truncation error in the 4 point backward difference formula which is given by this, that is used to approximate the 2<sup>nd</sup> derivative is, is it order of  $h$ , order of  $h^2$ , order of  $h^3$  or order of  $h^4$ .

I ask you to choose the correct answer from the 4 choices which I have given here. Unless you work out the details, you will not be able to give the results. So you start with the right-hand side, expand using Taylor's theorem and get the coefficient of  $f$ ,  $f'$ ,  $f''$ ,  $f'''$  and the 4<sup>th</sup> derivative of  $f$  and so on. And then that is equal to  $f''$  which appears on the left-hand side. So collect the coefficients of  $f$ , you will see it is 0,  $f'$ , that is also 0,  $f''$  also turns out to be 0,  $f'''$  also turns out to be 0.

On the other hand the 4<sup>th</sup> derivative of  $f$  coefficient turns out to be different from 0 and so you have right-hand side has  $h^4$  into 4<sup>th</sup> derivative at  $i$  divided by  $h^2$  because the formula had  $1/h^2$  in it. So we observe that you have the 4<sup>th</sup> derivative term appearing, multiplied by  $h^4$  by  $h^2$  and that is a nonzero term, so it is  $h^2$  times the 4<sup>th</sup> derivative. And so that is the 1<sup>st</sup> term which is a nonzero term that appears in this formula when you approximate the 2<sup>nd</sup> derivative by means of the formula on the right-hand side. So you conclude that the truncation error is of order of  $h^2$ .

(Refer Slide Time: 46:24)



And we have also discussed after numerical differentiation some results on numerical integration when we derived the closed type formulas, which are Newton Cotes type from last, namely the trapezoidal rule and Simpson's rule and we derived certain integration methods which are exact of polynomial is of degree  $N$  when we use  $N$  nodes in the interval say  $A, B$  which is a closed interval.

Then we said that if we can sacrifice the requirement, that these  $N$  nodes are equally spaced and look for methods which are quadrature methods, such that they are exact for polynomials of degree greater than  $N$ , when  $N$  nodes are used within an interval of the form  $A, B$  in particular an interval of the form  $-1$  to  $1$ , then we said that we obtain open type integration methods where we have used properties of the Legendre polynomials and the integration methods turn out to be methods which are exact for polynomial is of degree less than or equal to  $2N - 1$  where the function values at only  $N$  nodes are required.

And these  $N$  nodes turn out to be the zeros of the Legendre polynomial of degree  $N$  and these formulas were of the open type. The notes were all within the interval  $-1$  to  $1$ . And then we moved on to the solution of ordinary differential equations where we derived some single step methods like Taylor's series method, Euler smothered, modified Euler's smothered and Runge Kutta methods of order 2 and 4 and solved initial value problem using any of these methods. Then we said that it is possible to consider iterative type of methods where we have a predictor to predict the solution and then a corrector which would recorrect the solution and so that we can obtain solution correct to the desired degree of accuracy.

And these methods involved a number of points close to the point at which the initial condition is specified. So that information at 4 points were used to compute the solution at the 5<sup>th</sup> point. And these were multistep methods which we used as predictor corrector methods and we had solved an initial value problem using Milne's predictor character and Adam Bashforth's predictor character methods. Then we moved on to solution of nonlinear equations of the form  $F(X) = 0$ .

So we discussed 2 types of methods, namely the enclosure methods and then the fixed point iteration methods and discussed the order of convergence of these methods and we learnt some techniques of solving equations of the form  $F(X) = 0$ . And finally we moved on to the solution of a system of algebraic equations and then discussed the direct methods and iterative methods. And the direct methods included the decomposition methods of Doolittle, Crout and Cholesky and Gauss elimination method incorporating partial pivoting and then we discussed the iterative methods such as Gauss Seidel method and Gauss Jacobi method.

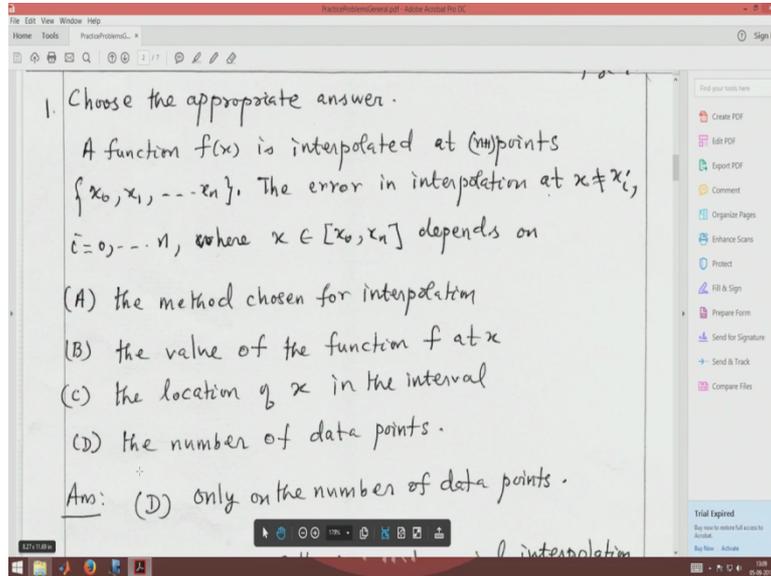
But then we wanted to see how in the, in the case of direct methods, we can have some information about the error that has been committed by us. So we performed error analysis and so we introduced norm of a vector, norm of a matrix associated with the given norm for a vector and we have shown how to get a bound on the absolute error and the relative error when we know the condition number of a given matrix  $A$  which appears as the coefficient matrix in the system  $AX = B$ . And when we know what the norm of the residual which is  $\|AX - B\|$  where  $X$  is an approximate solution  $B$  is the right-hand side vector that appears in the system  $AX = B$ .

And finally we moved to the matrix eigenvalue problems and computed the most dominant eigenvalue by power method and discussed the Gerschgorin bounds using Gerschgorin circle theorem and Brauer's theorem. And we have illustrated all these methods which we have learnt in this course by means of some examples, problems and you have solved number of assignment problems, you also have looked into a number of problems which we have given to you for practice and finally you are ready to appear for the end semester examinations.

The type of problem that you will get in your end semester will be such that you will have questions where you may have to choose the appropriate answer with proper justification, you may have to choose whether a given statement is true or false, you may have to fill in certain blanks where you may have to work out the details 1<sup>st</sup> and then fill in the results in the

blank and there will be also be problems where you will have to completely work out the details and present the solution. Some questions similar to what may come in your examination, I have included a few examples here. So let us quickly go through them.

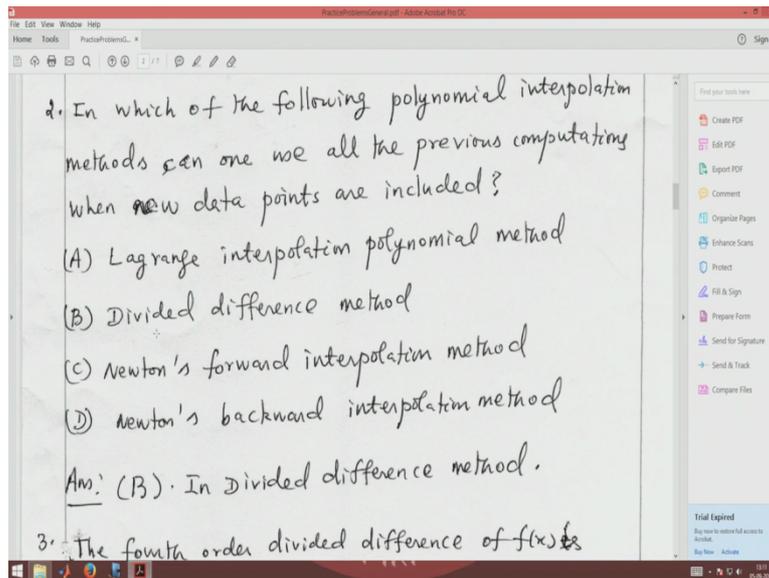
(Refer Slide Time: 52:31)



So you have the 1<sup>st</sup> question to be, choose the appropriate answer. So you are given a function  $F$  of  $X$  is interpolated at  $N + 1$  points,  $X_0$  to  $X_N$ . The error in interpolation at any point  $X$  which does not, which does not coincide with any of the  $X_i$  for  $i$  is equal to 1 to  $N$ , depends on what is the question. The answer is, is it a dependent on the method that is chosen for interpolation, surely not, it is clear. Does it depend on the value of the function  $F$  at  $X$ , no, you do not even know what the value of  $F$  at  $X$ . Thirdly is it dependent on the location of extremely interval, no, not at all,  $X$  can be anywhere in the interval  $X_0$  to  $X_N$ .

Is it dependent on the number of data points, yes, of course. If you choose 3 points from the data, you can give an interpolating polynomial of degree 2. If you choose 10 points from the data, you can give and you can get an interpolating polynomial of degree  $N$ . So the error in interpolation depends on the number of data points and not on any of the other 3 statements which are given. So you conclude that the result is only on the number of data points and Mark D to be the correct result and giving reason for that.

(Refer Slide Time: 54:02)



Then you have the another question which comes in interpolation again. In which of the following polynomial interpolation methods can one use all the previous all the previous computations where new data points are included. We have already discussed this in the class and therefore I straightaway go to the answer, we know that the use in Divided difference method, all the computations that have been used earlier because suppose say we are given 3 points and we are asked to get the Divided difference interpolation polynomial. Then we would be getting the interpolating polynomial of degree 2 that interpolates the setup given data which are 3 in number.

But now suddenly I have another information coming about this function, so another point and the value of that function at that point is given, now I have 4 points, so I can get an interpolating polynomial of degree 3 which I call as  $P_3$ , then what is  $P_3$ ,  $P_3$  is simply  $P_2$  of  $X$ , namely the quadratic interpolation polynomial that we already have computed. So  $P_2$  of  $X$  + an extra term, that extra term which arises due to the new information that is given to us. Right. And you know what the extra term should be, so if  $X_0, X_1, X_2$  are the previous point and  $X_3$  is the new point, then the extra term is going to be  $+A$  times  $X - X_0, X - X_1, X - X_2$ .

How do you determine  $A$  so your  $P_3$  of  $X$  is  $P_2$  of  $X$  +  $A$  times  $X - X_0, X - X_1, X - X_2$ . How do you determine  $A$ ? This  $P_3$  must be such that it should interpolate at the new point that is given. So substitute the  $X$  value as the  $X$  coordinate of the point, new point that is given and the  $Y$  value as a function value that is specified and determine  $A$  and you now have a cubic polynomial which interpolate the given data where you have used the computation of  $P_2$  of  $X$  which you have already done. So it is in the divide and difference formula that you are able to

use the previous computations. So you write the result as B, that is so in Divided difference method.

(Refer Slide Time: 56:30)

4. Newton-Raphson method is applied to solve  $f(x) = x^2 - 4 = 0$  and the computations are started with an initial approximation  $p_0$  close to 2, then

$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^2}$  is

(A)  $\frac{1}{2}$  (B) 1 (C)  $\frac{1}{4}$  (D) 0

Ans: (C)  $\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^2} = \frac{1}{2} \frac{f''(2)}{f'(2)} = \frac{1}{2} \frac{(2)}{4} = \frac{1}{4}$

$f(x) = x^2 - 4$  since  $x=2$  is a simple root and N-R method

$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^2}$  is

(A)  $\frac{1}{2}$  (B) 1 (C)  $\frac{1}{4}$  (D) 0

Ans: (C)  $\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^2} = \frac{1}{2} \frac{f''(2)}{f'(2)} = \frac{1}{2} \frac{(2)}{4} = \frac{1}{4}$

$f(x) = x^2 - 4$  since  $x=2$  is a simple root and N-R method has quadratic convergence

$f'(x) = 2x, f''(x) = 2$

and  $\lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n^2} = \frac{1}{2} \frac{f''(p)}{f'(p)}$

or  $e_{n+1} = C e_n^2$

So you have another problem on Divided difference which you can try to work out. Now if I consider examples in say methods of solving a non-linear equation of the form  $F$  of  $X$  is equal to 0 and I say apply Newton Raphson method to solve this equation. And the computations are started with an initial approximation  $P_0$  close to 2, that is an initial approximation to a root of the equation. And you are asked to find what is limit  $N$  tending to infinity of modulus of  $X_{N+1} - P$  by modulus of  $X_N - P$  the whole squared. Immediately you have recognised that this is your error at the  $N + 1$ th step, this is your error at the  $N$ th step.

So you are asked to get what is limit as N tending to infinity of EN +1 by EN square, is it any one of these values. Unless you work out, you cannot give the answer, so let us work out the details, what is it. I start with EN +1 by modulus of EN square. We have already done error analysis and we have shown, when we apply Newton Raphson method, EN +1 by EN to the power of Alpha where P is an approximation is nothing but half of the 2<sup>nd</sup> derivative of F at the root which is P divided by the 1<sup>st</sup> derivative of the function at the root. So P is given, because it says close to 2, that information is given to us.

So P is 2, and therefore we take this value, we know the function, how do we know the function, we are given F of X is equal to X square -4. Compute the 2<sup>nd</sup> derivative, compute the 1<sup>st</sup> derivative, substitute these derivatives, evaluate it at 2 and give the values here and that gives you 1 by 4. So immediately you know that the answer C is correct because you have been able to show that it is equal to 1 by 4. Because the Newton Raphson method is such that the error at the N +1th step is given by the C into the error at the Nth step squared whereas C, the asymptotic error constant is given by this information and that result has been used.

(Refer Slide Time: 59:14)

5. Newton-Raphson method for finding the square root of a real number R is

A)  $x_{n+1} = x_n - \frac{R}{2}$

B)  $x_{n+1} = x_n + \frac{R}{2}$

C)  $x_{n+1} = \frac{1}{2} \left[ x_n + \frac{R}{x_n} \right]$

D)  $x_{n+1} = \frac{1}{2} \left[ 3x_n - \frac{R}{x_n} \right]$

Ans = C  $x_{n+1} = \frac{1}{2} \left( x_n + \frac{R}{x_n} \right)$

$f(x) = x^2 - R$   
 $f'(x) = 2x$   
 $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} = x_n - \frac{(x_n^2 - R)}{2x_n}$   
 $= \frac{2x_n^2 - x_n^2 + R}{2x_n} = \frac{x_n^2 + R}{2x_n}$   
 $= \frac{1}{2} \left[ x_n + \frac{R}{x_n} \right]$

So all the results that we have shown in that there are very important in the sense that those results like what we have now shown, those results can be used to solve some examples and that has been illustrated above. Let us consider this example, Newton Raphson method for finding the square root of a real number R. What is it, is it any one of these, is the question. Yes, take the new trend of their method, take the function you want to get the square root of a real number R, so what is equation that we want to solve, X square - - R equal to 0 is equation, call that is F of X, so find F dash of X, substitute in XN +1 equal to XN - F of X N

by  $f'(x)$  to  $\frac{1}{xN}$ , that is what appears here, substitute, simplify and you end up and you observe that that appears here wherein this half of  $xN + R$  by  $xN$ . So you Mark C as the correct answer.

(Refer Slide Time: 60:04)

One wants to compute the reciprocal of a positive real number  $N$  using Newton-Raphson method, starting with an initial approximation. The sequence of iterates are generated from

(A)  $x_{n+1} = x_n(1 - Nx_n)$  ; (B)  $x_{n+1} = x_n(Nx_n - 1)$   
 (C)  $x_{n+1} = x_n(2 - Nx_n)$  ; (D)  $x_{n+1} = x_n(Nx_n - 2)$ .

Ans: [C]  $x_{n+1} = x_n(2 - Nx_n)$ .

$f(x) = N - \frac{1}{x} = 0, N > 0$   
 $f'(x) = \frac{1}{x^2}$  ∴ NRM is given by  $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

$$= x_n - \frac{N - \frac{1}{x_n}}{\frac{1}{x_n^2}}$$

$$= x_n \left[ \frac{1 - N + \frac{1}{x_n}}{\frac{1}{x_n}} \right] = \frac{2 - Nx_n}{\frac{1}{x_n}}$$

$$= 2x_n - Nx_n^2$$

$$= x_n(2 - Nx_n)$$

Similarly I have another example where we say the reciprocal a positive real number  $N$  using Newton Raphson method starting with an initial approximation. So the sequence of iterates are generated from which one of these? So we are asked to do by Newton Raphson method, you want to get the reciprocal of a number, so the equation that you want to solve is actually equal to 1 by  $N$  or the equation is  $N - 1/x = 0$ , compute  $f'(x)$ , substitute in Newton Raphson method,  $x_{n+1}$  equal to  $x_n - f(x_n)/f'(x_n)$  and that gives you this and that appears as one of the results, namely C and so you can fill in and say the answer is C.

(Refer Slide Time: 60:52)

Start with reason whether the following statement is TRUE or FALSE

1. A nonlinear equation  $f(x)=0$  is solved by the bisection method. If the number of steps it takes to get an error  $\epsilon$  is  $n$  with initial interval  $[a, b]$  and when the length of the initial interval is divided by 3, the number of steps it takes to get an error  $\epsilon$  is  $n_1$ , then  $n_1 > n$ .

Ans: FALSE.

We know that  $e \leq 2^{-n}(b-a) \rightarrow (1)$   
and it is given that  $e \leq 2^{-n_1} \frac{(b-a)}{3} \rightarrow (2)$

Ans: FALSE.

We know that  $e \leq 2^{-n}(b-a) \rightarrow (1)$   
and it is given that  $e \leq 2^{-n_1} \frac{(b-a)}{3} \rightarrow (2)$

This gives, dividing (1) by (2) and taking logarithm with respect to base 2 that

$$n = n_1 + 1.585 \quad (\because \log_2 3 = 1.585)$$

$\therefore n > n_1$ . The result shows that it takes 2 fewer iterations when the interval is divided by 3 to get an error  $\epsilon$ .

2. Suppose  $f$  is a continuous function on  $[-1, 2]$  and

and it is given that  $e \leq 2^{-n_1} \frac{(b-a)}{3} \rightarrow (2)$

This gives, dividing (1) by (2) and taking logarithm with respect to base 2 that

$$n = n_1 + 1.585 \quad (\because \log_2 3 = 1.585)$$

$\therefore n > n_1$ . The result shows that it takes 2 fewer iterations when the interval is divided by 3 to get an error  $\epsilon$ .

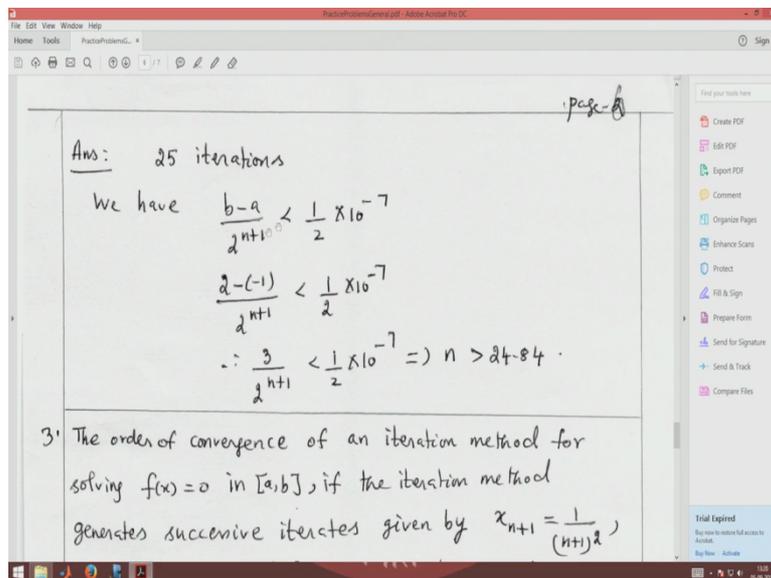
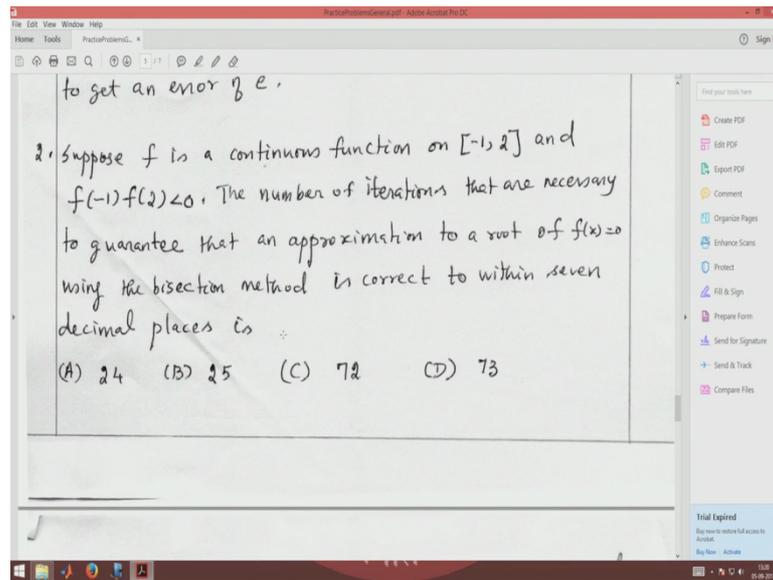
2. Suppose  $f$  is a continuous function on  $[-1, 2]$  and  $f(-1)f(2) < 0$ . The number of iterations that are necessary to guarantee that an approximation to a root of  $f(x)=0$  is correct to within seven

Then we have a desert which says states with reason whether the following statement is true or false. So the statement is, a non-linear equation  $F$  of  $X$  is equal to  $0$  is solved by bisection method. So it is a direct method for solving an equation of the form  $F$  of  $X$  is equal to  $0$ , it belongs to the class of enclosure method. Now if the number of steps to be taken, so that you want to get an error of  $E$  and that is  $N$  steps, with initial interval  $A, B$ . And you divide the interval in such a way that the length of the initial interval is divided by  $3$ . And the number of steps that you need to commit the same error  $E$  is say  $N_1$ .

Then the statement says  $N_1$  has to be greater than  $N$ , is your statement correct is the question. So do not worry about what you should write, start writing down what is given. What is given, we know that the error is in bisection method given by the length of the interval by  $2$  to the power of  $N + 1$  at the end of  $N$  steps. So you start, then what do you do, you divide the interval by  $3$ , so what is the length of the interval now, it is  $B - A$  by  $3$ . So what is the error at step, it is  $B - A$  by  $3$  divided by  $2$  power, now the number of steps is  $N_1$ , so  $2$  power  $N_1 + 1$ .

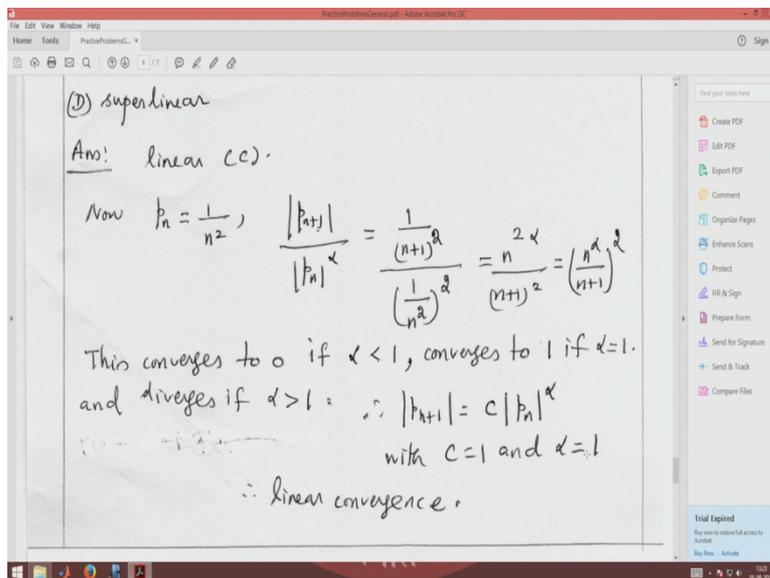
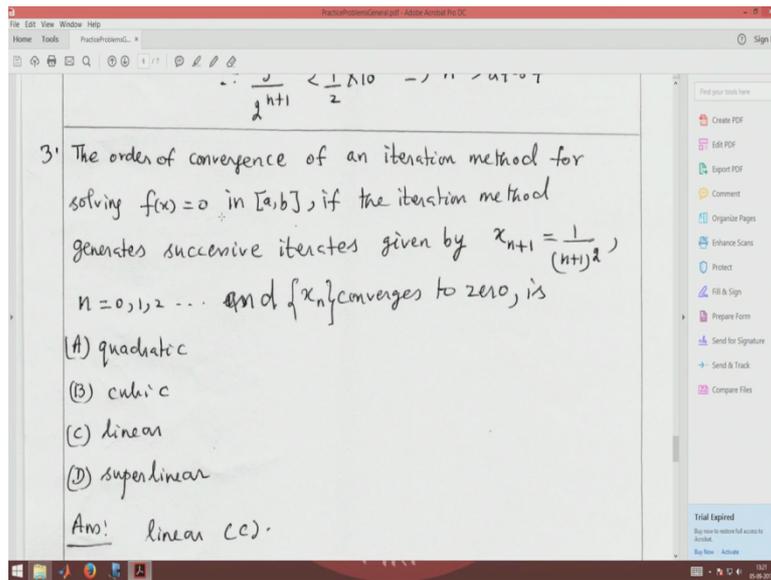
Now what is the error that you commit, it is the same  $e$ . So I have written whatever that is given in the problem, now I divide one by the other. And I see that when I take logarithm with respect to base  $2$ ,  $N$  turns out to be  $N_1 +$  some quantity. So  $N$  is bigger than  $N_1$ , so  $N$  is greater than  $N_1$ . The question says  $N_1$  is greater than  $N$ , no, the statement is therefore false, why because of all this. So you have to justify by giving the details of why you are saying that the statement is false.

(Refer Slide Time: 63:09)



Then suppose say  $F$  is a continuous function on the interval  $-1$  to  $2$ . And  $F$  of  $-1$  into  $F$  of  $-2$  is less than  $0$ , what does it mean, the root lies in  $-1$  to  $2$ . The number of iterations that are necessary to guarantee that an approximation to a root of  $F$  of  $X$  equal to  $0$  using bisection method is correct to within  $7$  decimal places is what, which one of them is correct. So you have a formula which says  $B - A$  by  $2$  power  $N + 1$ , you want it to be less than half of  $10$  to the power of  $-7$  because correct to within  $7$  decimal places is what your requirement is. Simplify and you get  $N$  to be bigger than  $24.84$  and you observe that  $N$  greater than  $N$  is equal to  $25$  is what appears here.

(Refer Slide Time: 64:01)



So now, the order of convergence of an iteration method for solving  $F$  of  $X$  is equal to 0 in  $A, B$ , if the iteration method generates successive iterates given by  $X_{N+1}$  equal to  $1$  by  $N+1$  the whole square and  $X_N$  converges to 0 is what. So you are asked to get the order of convergence of a method. You are given a sequence that generates, sequence which generates the successive iterates and which converges to 0. So what is the order? So  $P_n$  is given, how do you compute the order, you must get what is  $E_{N+1}$  divided by  $E_N$  power Alpha, that must be an asymptotic error constant times, right,  $E_N$ , asymptotic error constants, which is  $C$ .

So I repeat again, you have to show that  $E_{N+1}$  is equal to  $E_N$  to the power of Alpha multiplied by  $C$ . Alpha gives you the order of convergence of the method and  $C$  give you the asymptotic error constant. So let us see what is this in this example. Now the sequence

converges to 0, so it is  $P_{N+1} - 0$  is  $E_{N+1}$ ,  $P_N - 0$  is  $E_N$ . So I consider  $\text{mod } E_{N+1}$  by  $\text{mod } E_N$  power  $\alpha$ . So it is given  $1$  by  $N+1$  the whole square by  $1$  by  $N$  square the whole square. So that is  $N$  to the  $1$  by  $N$  square to the power of  $P_N$  power  $\alpha$ . So this must be  $\alpha$ .

So that gives you  $N$  power  $2\alpha$  by  $N+1$  the whole square, this is  $\alpha$ , right. You have  $N$  power  $\alpha$  by  $N+1$  the whole square. Now when will this quantity converge to 0? That will converge to 0 if  $\alpha$  is less than 1, it will converge to 1 if  $\alpha$  is 1 and it will diverge if  $\alpha$  is greater than 1. And you want the sequence to converge, so  $\alpha$  must be such that it converges to 1 if  $\alpha$  is equal to 1. So modulus of  $P_{N+1}$  must be  $C$  into modulus of  $P_N$  to the power of  $\alpha$ , with  $C$  is equal to 1 and  $\alpha$  equal to 1 and so you conclude that the method converges and its order of convergence is 1, so it has linear convergence.

I have illustrated some examples from the different topics that we have discussed in our course so that one can expect questions of this type also in your examination. So I hope you have enjoyed this course and you have had learnt some good portion of the topics from numerical analysis from this course and I wish you all the best in your career and in future. Thank you very much.