

Matrix Solvers
Prof. Somnath Roy
Department of Mechanical Engineering
Indian Institute of Technology, Kharagpur

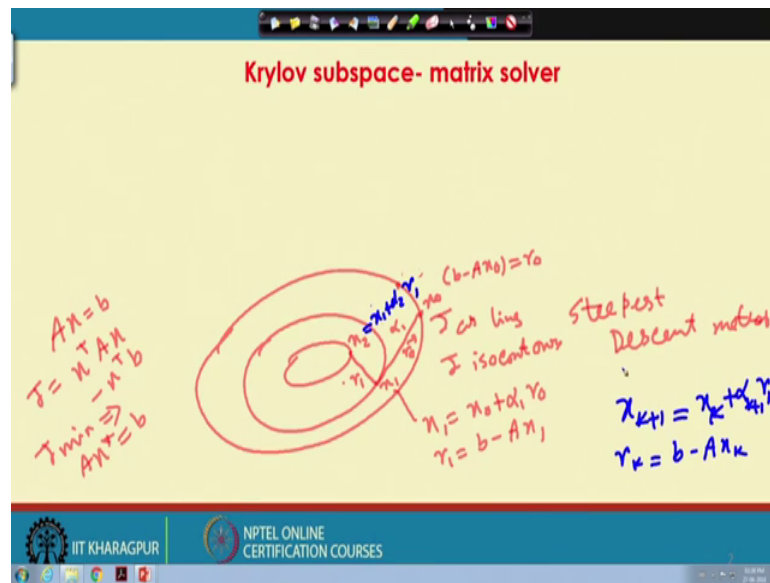
Lecture – 46
Conjugate Gradient Methods

Welcome we are discussing about Krylov subspace methods for solving Matrix equations iteratively. So, today we will discuss about Conjugate Gradient Method which is a very first Krylov subspace method for solving symmetric matrix equations. The classes before we have discussed about Arnoldi's method with through which using something like a modified Gram Schmidt method we can get orthogonal basis for Krylov subspace. And then we have also seen how to apply this for solving matrix equation using full orthogonal method or FOM.

And we also looked about different variants of FOM and then further looked into Lanczos algorithm in which, FOM can be converted for symmetric matrix equation which essentially gives something like a tri diagonal matrix system in which, direct solution like TDMA type algorithm can be utilized for faster solution. And now, we will continue this discussion for conjugate gradient method.

Before going into that, I will again quickly review few of the important things in Krylov subspace method, as well as full orthogonal method or and Arnoldi's method, which will give a base for continuing this discussion to conjugate gradient method.

(Refer Slide Time: 01:49)



So, if we look into the matrix solver solution equations for example, when we discussed about steepest descent algorithm, which is the building block of any projection method. And Krylov subspace method is just an extension of general projection method whether, the where the projection space is not one dimensional rather multi dimensional.

So, if you look into the steepest descent algorithm for example, we have a we have to solve Ax is equal to b and we define a function J is equal to x transpose A x minus x transpose b , which has to be minimized J min will imply, Ax star is equal to b . This has to be solved so, we get the exact solution of Ax is equal to b .

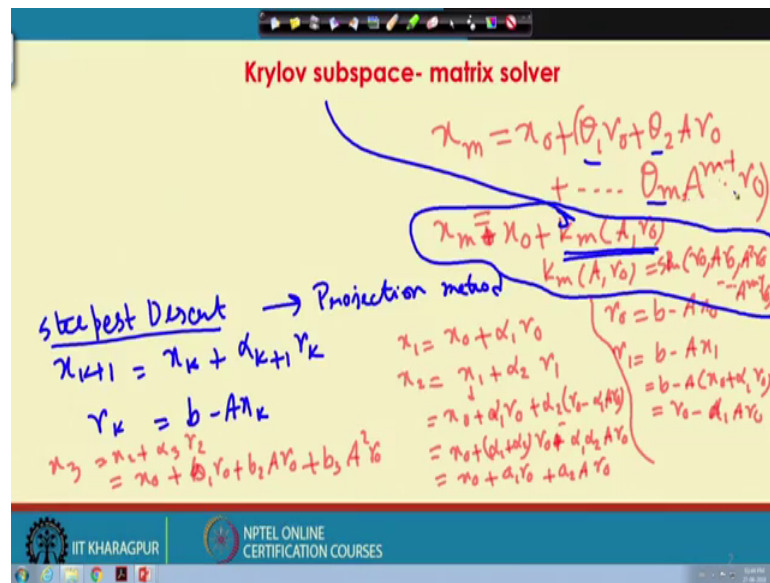
So, we start with the iso contours of J and this is basically, steepest descent what I am discussing here. So, these are the J constant lines or we call them J isocontours, these are the lines over which J is constant. So, we will start from any value x_0 and evaluate what is b minus Ax_0 that is equal to r_0 . And then I will move along r_0 , I will keep on moving along r_0 , till just one second this r_0 has to be orthogonal to it.

So, let us select the x_0 here. This is just for convenience of explaining that this x_0 here and I will keep on moving along r_0 , say up to a distance α_1 . Here, the r_0 vector will be tangential to the new iso contour and I have to change the search direction. I have to evaluate so, this is x_1 where, x_1 is equal to x_0 plus $\alpha_1 r_0$ for example. And I will evaluate r_1 is equal to b minus Ax_1 and then I will move along r_1 to find out x_2

where, x_2 is equal to I will write x_2 is equal to b minus sorry, x_2 is equal to x_1 plus $\alpha_2 r_1$ something like that.

So, at each step x will be updated as a function as x_{k+1} is equal to x_k plus $\alpha_{k+1} r_k$ something like this. And my r_k will be defined as $b - Ax_k$. So, now if we look into this update into little detail or just work it, work out the updates for first 2-3 steps.

(Refer Slide Time: 05:19)



So, I rather let us this I keep with the relations that, x_{k+1} is equal to x_k plus $\alpha_{k+1} r_k$ and r_k is equal to $b - Ax_k$. Exactly, this is what we have done in steepest descent, which is building block of projection method. So, this later we will give as what are the general projection methods. Now, I can write that let us use another pen, x_1 is equal to x_0 plus $\alpha_1 r_0$. Then, x_2 will be x_1 plus $\alpha_2 r_1$.

Now, what is x_1 ? x_1 is nothing but x_0 plus $\alpha_1 r_0$. And what is r_1 ? r_0 is equal to $b - Ax_0$, r_1 is equal to $b - Ax_1$ is equal to $b - A(x_0 + \alpha_1 r_0)$, which is $b - Ax_0 - \alpha_1 A r_0$. So, I can substitute here that x_1 is x_0 plus $\alpha_1 r_0$ and this is plus α_2 into r_1 minus $\alpha_1 A r_0$. So, this is x_0 plus $\alpha_1 r_0$ plus $\alpha_2 r_1$ plus $\alpha_2 (-\alpha_1 A r_0)$.

Now, can we see that this can be written as also x_0 plus say; $\alpha_1 r_0$ plus $\alpha_2 A r_0$. Similarly, if I write expression for x_3 , this will be x_2 plus $\alpha_3 r_2$. And now, I will

substitute calculate the values of r^2 and x^2 and substitute it, here I will again get x_0 plus A say, $b_1 r_0$ plus $b_2 A r_0$ plus $b_3 A^2 r_0$ so on. So, I will get a general expression for m th order or m th iteration of x as, let us let me wipe this part out so that, or let me write down the general expression that, x_m after m th iteration is equal to x_0 plus say, $\theta_1 r_0$ plus $\theta_2 A r_0$ plus and it goes up to $\theta_{m+1} A^m r_0$. x^3 is related with β^3 so, θ_3 rather $\theta_3 A^2 r_0$.

So, the entire update of x_0 becomes a function of r_0 plus $A r_0$ plus $A^2 r_0$ plus $A^3 r_0$ plus $A^4 r_0$ plus $A^5 r_0$ plus $A^6 r_0$ plus $A^7 r_0$ plus $A^8 r_0$ plus $A^9 r_0$ plus $A^{10} r_0$ plus $A^{11} r_0$ plus $A^{12} r_0$ plus $A^{13} r_0$ plus $A^{14} r_0$ plus $A^{15} r_0$ plus $A^{16} r_0$ plus $A^{17} r_0$ plus $A^{18} r_0$ plus $A^{19} r_0$ plus $A^{20} r_0$ plus $A^{21} r_0$ plus $A^{22} r_0$ plus $A^{23} r_0$ plus $A^{24} r_0$ plus $A^{25} r_0$ plus $A^{26} r_0$ plus $A^{27} r_0$ plus $A^{28} r_0$ plus $A^{29} r_0$ plus $A^{30} r_0$ plus $A^{31} r_0$ plus $A^{32} r_0$ plus $A^{33} r_0$ plus $A^{34} r_0$ plus $A^{35} r_0$ plus $A^{36} r_0$ plus $A^{37} r_0$ plus $A^{38} r_0$ plus $A^{39} r_0$ plus $A^{40} r_0$ plus $A^{41} r_0$ plus $A^{42} r_0$ plus $A^{43} r_0$ plus $A^{44} r_0$ plus $A^{45} r_0$ plus $A^{46} r_0$ plus $A^{47} r_0$ plus $A^{48} r_0$ plus $A^{49} r_0$ plus $A^{50} r_0$ plus $A^{51} r_0$ plus $A^{52} r_0$ plus $A^{53} r_0$ plus $A^{54} r_0$ plus $A^{55} r_0$ plus $A^{56} r_0$ plus $A^{57} r_0$ plus $A^{58} r_0$ plus $A^{59} r_0$ plus $A^{60} r_0$ plus $A^{61} r_0$ plus $A^{62} r_0$ plus $A^{63} r_0$ plus $A^{64} r_0$ plus $A^{65} r_0$ plus $A^{66} r_0$ plus $A^{67} r_0$ plus $A^{68} r_0$ plus $A^{69} r_0$ plus $A^{70} r_0$ plus $A^{71} r_0$ plus $A^{72} r_0$ plus $A^{73} r_0$ plus $A^{74} r_0$ plus $A^{75} r_0$ plus $A^{76} r_0$ plus $A^{77} r_0$ plus $A^{78} r_0$ plus $A^{79} r_0$ plus $A^{80} r_0$ plus $A^{81} r_0$ plus $A^{82} r_0$ plus $A^{83} r_0$ plus $A^{84} r_0$ plus $A^{85} r_0$ plus $A^{86} r_0$ plus $A^{87} r_0$ plus $A^{88} r_0$ plus $A^{89} r_0$ plus $A^{90} r_0$ plus $A^{91} r_0$ plus $A^{92} r_0$ plus $A^{93} r_0$ plus $A^{94} r_0$ plus $A^{95} r_0$ plus $A^{96} r_0$ plus $A^{97} r_0$ plus $A^{98} r_0$ plus $A^{99} r_0$ plus $A^{100} r_0$.

So, this is: what is the basic definition of Krylov subspace; what we have used earlier. So, now, if we look into the definition how x is updated and when updating x for example, when we looked into steepest descent x is always updated along r . And in that way, we made it also a point that the new residual vector r should be orthogonal to the previous residual vector or also to the functional subspace iso contour of J .

So, there is a constraint on how x will be updated and based on which the values of α are calculated. So, x is a x is any way updated in this particular plane, which is space or which is Krylov subspace, this is the Krylov subspace, x is always updated in along on Krylov subspace. And what are the coefficients through which, these updates will be there, what will be θ_1 θ_2 θ_m , that will come from the projection constraint that the update should be orthogonal to or the residual should be orthogonal to a particular plane.

(Refer Slide Time: 11:20)

Krylov subspace- matrix solver

General Projection Method

The projection method seeks an approximate solution x_m from an affine subspace $x_0 + K_m$ by imposing condition $b - Ax_m \perp L_m$, L_m is another subspace of dimension m , x_0 is the initial guess.

In the case of Krylov subspace methods, $K_m = K_m(A, r_0)$, $r_0 = (b - Ax_0)$ in R^n

$$K_m = \text{span}\{r_0, Ar_0, A^2 r_0, \dots, A^{m-1} r_0\}$$

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES

So, if we now look into the recap what we have seen in Krylov space, which is we will also write that this is a general projection method. And, we have also seen that, this update is updating x in Krylov subspace we will make it converge into the exact solution x^* and the residual as we update x and (Refer Time: 11:50) Krylov subspace also the residual will also converge to 0.

The projection methods seeks an approximate solution x_m from an affine subspace of $x_0 + K_m$, I have discussed: what is the definition of affine subspace in last few classes. By imposing the condition that $b - x_m$ is a this particular perpendicular to L_m where, L_m is another subspace of dimension f and x_0 is the initial guess in case of so, this is the statement for general projection method. In case of Krylov subspace methods, K_m is given by K_m of a r_0 where r_0 is $b - a x_0$ and K_m is a space spanned by $r_0, A r_0, A^2 r_0$ up to the power A to the power $m - 1 r_0$.

So, what we are thinking as multiple iterative steps or multiple search directions in a steepest descent algorithm is basically, the multiple bases vectors through which x_m is updated x_0 is updated as x_m and these are the multiple base several bases vectors of the Krylov subspace. And we have seen what is the grade of why with respect to A what can be the maximum dimension of Krylov subspace etcetera in last few classes and have shown that, this is a convergent method, this the in general Krylov subspace method should converge to the right solution based on any (Refer Time: 13:33) value of x_0 or

any starting residual $r \neq 0$. So, our goal is to develop efficient solvers using Krylov subspace method.

(Refer Slide Time: 13:36)

Krylov subspace methods

The different versions of Krylov subspace methods arise from different choice of L_m and from the way in which the system is preconditioned. Two broad choices for L_m give rise to the best known techniques:

$L_m = K_m$ FOM or Full orthogonal
 $L_m = AK_m$ GMRES, MINRES → oblique projection

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES

The different values of versions of Krylov subspace methods arise from different choices of L_m ; that means, the space on which to which the residual should be orthogonal and from the way in which the system is preconditioned. This is a recapitulation slide of previous lectures. There are 2 broad choices for L_m which, gives based on techniques: one is L_m is equal to K_m , which is an orthogonal method and another is L_m is equal to AK_m , which is an oblique projection.

And now, in this particular session we are concentrating on the full orthogonal method. We have earlier discussed about full orthogonal method also and we will see how for symmetric matrix we can generate very efficient solver from full orthogonal method.

(Refer Slide Time: 14:42)

Arnoldi-Modified Gram-Schmidt

At each step the algorithm multiplies Arnoldi vector v_j by A and then orthonormalizes the resulting vector w_j against all previous v_j 's by a standard Gram-Schmidt procedure

1. Choose a vector v_1 of norm 1
2. For $j = 1, 2, \dots, m$ Do:
3. Compute $w_j := Av_j$
4. For $i = 1, \dots, j$ Do:
5. $h_{ij} = (w_j, v_i)$
6. $w_j := w_j - h_{ij}v_i$
7. EndDo
8. $h_{j+1,j} = \|w_j\|_2$. If $h_{j+1,j} = 0$ Stop
9. $v_{j+1} = w_j/h_{j+1,j}$
10. EndDo

This method results in to v and w , bases V and W for K_m and L_m

h_{ij} are elements of a Hessenberg matrix \bar{H}_m

IIT KHARAGPUR NPTEL ONLINE CERTIFICATION COURSES

And Arnoldi modified Gram Schmidt is an algorithm through which, we can get the bases vectors of the this is bases vectors of v , the bases vectors of the Krylov subspace which are the base v 's are the bases vectors of Krylov subspace. We can generate it through Arnoldi modified Gram Schmidt method that, you start with any vector v_1 and which is basically, your when you are using full orthogonal method this will be $r=0$. And then, take a product of v with respect to the A matrix Av , from v you subtract Av , with v is the bases of Krylov subspace Av is the next bases of Krylov subspace so, you orthogonalize v and Av using something like, Gram Schmidt error in find out the orthonormal bases vectors.

In that way you get another matrix which is h_{ij} ; which is the product between the new bases of Krylov subspace space and the older bases of Krylov subspace, before orthogonalization of the new bases. And this h_{ij} forms a Hessenberg matrix \bar{H}_m . And this Hessenberg matrix means, sorry this is \bar{H}_m forms a Hessenberg matrix \bar{H}_m . The Hessenberg matrix is basically, an upper triangular matrix plus 1 sub diagonal term into it. And this is how we get h_{ij} because, it goes up to $h_{j+1,j}$ so, it goes up to the diagonal term and plus one sub diagonal term in it.

(Refer Slide Time: 16:39)

Residual vector from Arnoldi's method

$$r_m = b - Ax_m = b - A(x_0 + V_m y_m)$$

$$= b - Ax_0 - AV_m y_m$$

$$= r_0 - AV_m y_m$$

$$= \beta_1 - AV_m y_m$$

So: $r_m = \beta_1 - (V_m H_m + w_m e_m^T) y_m$

$$= \beta_1 - V_m H_m y_m - h_{m+1,m} v_{m+1} e_m^T y_m$$

$$\Rightarrow r_m = -h_{m+1,m} e_m^T y_m v_{m+1}$$

Now: $AV_m = V_m H_m + w_m e_m^T$
 $w_m = h_{m+1,m} v_{m+1}$

Further: $y_m = H_m^{-1}(\beta e_1)$
 $\Rightarrow H_m y_m = \beta e_1 \Rightarrow V_m H_m y_m = V_m \beta e_1 = \beta v_1$

Handwritten notes:
 $x_m = x_0 + K_m(A, x_0)$
 $V_m \rightarrow$ orthonormal basis of K_m
 FOM: $y_0 = ||r_0|| v_1 = \beta_1 v_1$
 Arnoldi
 Full orthogonal method:
 $L_m = K_m \& V_{m+1} K_m$

Now, we look into the residual vector from Arnoldi's method, what is a residual vector at one particular iterative step on Arnoldi's method. And this is r_m is equal to $b - Ax_m$ and x_m is x belongs to, so we have seen it earlier that x is equal to $x_0 + K_m(A, x_0)$ and V_m is orthonormal bases of K_m . So, we can write x_m is equal to $x_0 + V_m$ multiplied by sorry they K_m of x_0 multiplied by y . So, x_m is x_0 plus the Krylov subspace of A and r_0 and the orthonormal bases of this Krylov subspace is obtained as V_m from Arnoldi method. So, this comes from Arnoldi, from Arnoldi's method. So, you can write x_m is equal to $x_0 + V_m$ into y .

So, if substituted x_m is equal to $x_0 + V_m$ into y V_m into y_m (Refer Time: 18:03) and y_m 's are the coefficients which are multiplied with different bases vectors of Arnoldi's of Krylov subspace to get a general expression of a general vector in Krylov subspace. So, x_m is equal to x_0 plus a general vector in Krylov subspace and we have to find out what y_m is that so that, we can satisfy the equation come to it later.

And so this is $b - Ax_0 + V_m y_m$ $b - Ax_0$ is equal to r_0 . So, $r_0 - AV_m y_m$, which is r_0 is again we have if we go to the maybe the previous slide. So, we started with not here, but if we if we not in the previous slide if we see, in full orthogonal method right in FOM we started with the first bases of Krylov subspace is capital V is equal to mod of r_0 into v_1 and which I have written as β_1 into v_1 .

So, β is mod of r_0 which is the first this magnitude of the unit vector first case of first residual in Krylov subspace are 0 in Krylov space. So, βv_1 is equal to and we can write and this is so, we can write that V is equal to or V is equal rather V not should not write that we star put r_0 is mod r_0 into v_1 , v_1 is an unit bases you first unit bases first bases vectors of the Krylov subspace which is an unit vector this is βv_1 .

So, r_0 is equal to βv_1 I substitute r_0 is equal to βv_1 and the m th ordered residual is βv_1 minus $A V^m y^m$. Now, from Arnoldi's method we got these identities, this comes again from Arnoldi. If we look into the Arnoldi's algorithm we have just shown in the last slide, this these relations are evident from there that $A V^m$ is equal to $V^m H^m$ plus $w^m e^m$ transpose, w^m is h^m plus $1 h^m v^m$. So, now in $A V^m$, I will substitute $A V^m$ by this particular quantity. So, r^m is equal to βv_1 minus $V^m H^m$ plus $w^m e^m$ transpose y^m . And this is βv_1 minus $V^m H^m y^m$ minus h^m plus $1 v^m e^m$ transpose.

Now, what is y^m and that is what we are trying to find out what is the value of y^m . Because, if we know y^m we have defined the Krylov subspace, Arnoldi's method will give us the orthonormal bases of Krylov subspace. If we can find y^m , we will multiply it with different bases vectors and find out what is the update in Krylov subspace which is to be added with initial guess x_0 so that I get the final solution. So, now this update is obtained, if we can remember steepest descent method, α is the amount of distance that I should go along one particular direction vector r . And, this is obtained from the relation that the residual is orthogonal to particular subspace.

Similarly, full orthogonal method uses the fact that the residual is orthogonal to L^m which is same as K^m , the residual is also orthogonal to the prior surface only. And using that we can get, so this is using the fact that full orthogonal method or Arnoldi's method for solving linear equation that gives L^m is equal to K^m and r^m is orthogonal to r^m plus 1 rather is orthogonal to K^m .

So, using that we can find out we have shown it in few lectures before, y^m is H^m inverse βe_1 . So, $H^m Y^m$, H^m is the Hessenberg matrix the a part of the Hessenberg matrix which we have obtained in the last shown in the last slide. So, $H^m y^m$ is βe_1 multiply both side by V^m , $V^m H^m y^m$ which is the this term is V^m into βe_1 , e_1 is an unit first unit vector $1 0 0$ which is multiplied with V^m will give the first vector in

the of the orthonormal bases (Refer Time: 23:26) beta v 1. So, I can write that, this is equal to beta V 1 utilizing this factor and they will cancel out. So, I will get that r m is minus h m plus 1 e m transpose y m v m plus 1.

This is a very important identity from the Arnoldi's method or from full orthogonal method that, the residual is orthogonal to the next bases in the Krylov subspace, residual at one particular m level is orthogonal to the next bases function of Krylov subspace. Sorry, residual is not orthogonal, residual is parallel, the residual is along the next bases of Krylov subspace. So, residual will be orthogonal to the previous bases I am sorry. Residual r m is along v m plus 1 where, v m plus 1 is the m plus 1th bases of Krylov subspace; mth residual is in the direction of m plus 1th bases of Krylov subspace that is the finding.

(Refer Slide Time: 24:45)

Residual vector from Arnoldi's method

$$r_m = -h_{m+1,m} e_m^T y_m v_{m+1}$$

So: 1. Residual vector is in the direction of v_{m+1} .
2. Residual vectors are orthogonal to each other

Arnoldi modified Gram-Schmidt
 $u_1 \perp u_2 \dots u_m \perp u_{m+1}$
 $r_m = a_m v_{m+1}$
 $r_{m-1} = a_{m-1} v_m$
 $r_m \perp r_{m-1}$

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES

So, we start with r m is equal to minus h m plus 1 m e m transpose y m v m plus 1. So, the residual vector is in the direction of v m plus 1 and residual vectors are orthogonal to each other. Why? This is simply due to the fact that v 1 is orthogonal to v 2 is v m is orthogonal to v m plus 1, which comes due to Arnoldi method is modified by Gram Schmidt.

Gram Schmidt method will give a set of orthonormal vectors. So, these v v vectors are orthogonal to each other. Therefore, once I see residual r m or r m is equal to something say a m v m plus 1 and r m minus 1 is equal to a m minus 1 v m right, that is the

founding here; v_{m+1} and v_m are orthogonal to each other. So, r_m and r_{m-1} ; that means, r_m is orthogonal to r_{m-1} and so on. So, residual vectors are orthogonal to each other.

(Refer Slide Time: 26:08)

Residual vector from Arnoldi's method

$$r_m = -h_{m+1,m} e_m^T y_m v_{m+1}$$

So: 1. Residual vector is in the direction of v_{m+1} .
2. Residual vectors are orthogonal to each other

For symmetric matrix, $A: H_m = T_m \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_m & \\ & & & & \beta_m & \\ & & & & & \eta_1 & \\ & & & & & & \eta_2 & \\ & & & & & & & \ddots & \\ & & & & & & & & \eta_{m-1} & \\ & & & & & & & & & \eta_m \end{pmatrix}$

We get $r_m = -\beta_{m+1} e_m^T y_m v_{m+1}$

For a now for a symmetric matrix A : H_m or the part of that Hessenberg matrix is a tri diagonal matrix. We have shown it earlier if A becomes symmetric H_m is a Hessenberg matrix, so it has a it is a upper triangular matrix with one sub diagonal if A is symmetric, so H_m is also symmetric. So, only the sub diagonal and super diagonal exists and it becomes a tri diagonal matrix. And we can express this tri diagonal matrix as product of two bi diagonal matrices; one is a lower triangular form, another is an upper triangular form. And if we use this form these betas are comes in terms of the h here.

So, r_m becomes minus beta $m+1$ $e_m^T y_m v_{m+1}$ for symmetric matrix A . If we can tri diagonalize the Hessenberg matrix h_m for symmetric matrix A , we can directly calculate from that form of the Hessenberg matrix, we can directly calculate what is the coefficient which will be multiplied with the m transpose $y_m v_{m+1}$ to get the residual vector.

So, there is a relationship between the residual vector and the tridiagonal matrix which directly comes here. So, with will I will stop here in this session and with this particular idea that our residual vectors are at m th level is along the base $m+1$ in bases of Krylov subspace and residual vectors are orthogonal to each other. We will see how this

we got an expression for the residual vector for symmetric matrices, how this can be substituted into the Lanczos algorithm or direct Lanczos algorithm where, we are doing this TDMA type of calculation and we can get a faster solution method.

Thank you.