

**UC Berkeley**  
**SEMM Reports Series**

**Title**

A Discussion on Using Lanczos Vectors and Ritz Vectors for Computing Dynamic Responses

**Permalink**

<https://escholarship.org/uc/item/68q3g2sc>

**Authors**

Chen, Harn-Ching

Taylor, Robert

**Publication Date**

1988-08-01

REPORT NO.  
UCB/SEMM-88/15

**STRUCTURAL ENGINEERING  
MECHANICS AND MATERIALS**

**A DISCUSSION ON USING  
LANCZOS VECTORS AND RITZ VECTORS  
FOR COMPUTING DYNAMIC RESPONSES**

**BY**

**HARN C. CHEN and ROBERT L. TAYLOR**

AUGUST 1988

**DEPARTMENT OF CIVIL ENGINEERING  
UNIVERSITY OF CALIFORNIA  
BERKELEY, CALIFORNIA**

# **A discussion on using Lanczos vectors and Ritz Vectors for computing dynamic responses**

**Harn C. Chen and Robert L. Taylor**

*Department of Civil Engineering  
University of California, Berkeley, CA 94720, USA*

## **ABSTRACT**

The Lanczos vectors and the Ritz vectors have been used for computing the dynamic response of linear structures. Although the procedures of using these two sets of vectors appear similar to the procedure of using the eigenvectors to find an approximate solution, the fundamental mechanisms of the three are different. We compare the three sets of vectors in detail to show some of the important differences in the hope that this comparison will be helpful to the use of the Lanczos vectors or the Ritz vectors for computing dynamic responses.

## **CONTENTS**

1. Introduction
  2. Reduced System
  3. Eigenvectors
  4. Lanczos Vectors
  5. Ritz Vectors
  6. Recommendations
- References

## INTRODUCTION

The mode superposition method has been used frequently in the analysis of dynamic response of linear structures. The standard procedure of the method consists in finding the eigenvectors of the system and using them to transform the coupled equations of motion into an equivalent set of decoupled equations. These decoupled equations can then be solved individually and their solutions are superposed to give the response of the dynamic system. A major advantage of this method is the fact that often only the first few decoupled equations are needed to give a satisfactory approximation, thus resulting in considerable savings in the computational effort.

Recently, Wilson *et al*<sup>1</sup> proposed using a set of load-dependent Ritz vectors to solve a class of problems where the applied load  $\mathbf{f}(t)$  is of the form

$$\mathbf{f}(t) = \hat{\mathbf{f}} \epsilon(t) \quad (1)$$

in which  $\hat{\mathbf{f}}$  is a space vector and  $\epsilon(t)$  a time function. In this class of problems, the Ritz vectors generated by choosing the static deflection shape as the starting vector produce better approximations than the same number of eigenvectors. In addition, these Ritz vectors are less expensive to generate than the eigenvectors. These advantages make the Ritz vectors attractive for solving a variety of problems<sup>2-5</sup>. In the mean time, Nour-Omid and Clough<sup>6</sup> gave a detailed description on how to efficiently generate a set of Lanczos vectors for the dynamic analysis. In fact, the Lanczos vectors and the load-dependent Ritz vectors are equivalent since they span the same subspace, a special Krylov subspace. A Krylov subspace has been used to find eigenpairs of a large system for a long time; however, its use to obtain a reduced system for the dynamic response analysis is new.

Although the procedure of using the Lanczos vectors or the Ritz vectors appears quite the same as the procedure of using the eigenvectors to find the dynamic response, the fundamental mechanisms of the three sets of vectors are different. By contrasting the three sets in detail, we find out the following important points. First, it is not appropriate, at least physically, to use the term, "participation factors", for the right-hand side of the reduced system resulted from the Lanczos vectors or the Ritz vectors. Second, the error in the representation of the space vector  $\hat{\mathbf{f}}$  should be distinguished from the error in the

response of the system. When the eigenvectors are used, the error in representing  $\hat{\mathbf{f}}$  has a direct relation to the error in the response and therefore the error in representing  $\hat{\mathbf{f}}$  is usually used as an estimate of the error in the response; however, this relation does not carry over when the Lanczos vectors or the Ritz vectors are used. In addition, using the *dynamic* deflection shape as the starting vector for generating the Lanczos vectors or the Ritz vectors can further accelerate the convergence of the superposition process as compared with using the static deflection shape. To demonstrate these important points, we first present a theoretical background and then examine the three sets of vectors one by one in the following.

## REDUCED SYSTEMS

The differential equations of motion for a discretized model of structural systems can be expressed by

$$\mathbf{M} \ddot{\mathbf{u}}(t) + \mathbf{C} \dot{\mathbf{u}}(t) + \mathbf{K} \mathbf{u}(t) = \mathbf{f}(t) \quad (2)$$

where  $\mathbf{M}$ ,  $\mathbf{C}$ , and  $\mathbf{K}$  are, respectively, the  $n \times n$  mass, damping, and stiffness matrices, and  $\ddot{\mathbf{u}}(t)$ ,  $\dot{\mathbf{u}}(t)$ , and  $\mathbf{u}(t)$  are the  $n \times 1$  acceleration, velocity, and displacement vectors. In practical analysis of complicated systems, the order  $n$  of the equations of motion is usually very large. Therefore, how to efficiently obtain a satisfactory approximate solution becomes a major concern. A commonly used approach in this regard is to reduce the original set of equations to a much smaller set and to find the approximate solution by solving the reduced set. In mathematical terms this approach can be interpreted as a *projection* process. We summarize the essential ingredients below.

A set of  $n$  vectors  $\mathbf{Y} = [ \mathbf{y}_1, \dots, \mathbf{y}_n ]$  is *orthonormal* if they satisfy the condition  $\mathbf{Y}^T \mathbf{Y} = \mathbf{I}$ , where  $\mathbf{I}$  is the  $n \times n$  identity matrix. Any set of  $n$  orthonormal vectors can be chosen as a *basis* for expressing the solution  $\mathbf{u}(t)$  since it contains  $n$  components. That is, the exact solution to equation (2) can be expressed in a chosen basis  $\mathbf{Y}$  as

$$\mathbf{u}(t) = \sum_{j=1}^n \mathbf{y}_j x_j(t) = \mathbf{Y} \mathbf{x}(t) \quad (3)$$

in which the basis  $\mathbf{Y}$  serves to transform from the generalized coordinates  $\mathbf{x}(t)$  to the

geometric coordinates  $\mathbf{u}(t)$ . Theoretically, only linear independence is required for  $\mathbf{Y}$  to be a basis, but an orthonormal set is convenient for computation. Since the solution is exact no matter which basis is used, any basis can be used to express the solution. In practice, however, we are interested in finding a satisfactory approximation  $\mathbf{u}_m(t)$  by using only a small subset of  $\mathbf{Y}$ , i.e.,

$$\mathbf{u}_m(t) = \sum_{j=1}^m y_j x_j(t) = \mathbf{Y}_m \mathbf{x}_m(t) \quad (4)$$

Therefore, we have to choose a set of basis vectors such that the first few of them,  $\mathbf{Y}_m$  with  $m$  being much smaller than  $n$ , will produce a solution with satisfactory accuracy.

Geometrically, the  $\mathbf{u}_m(t)$  is a projection of  $\mathbf{u}(t)$  onto the subspace  $\text{span}[\mathbf{Y}_m]$ . Different  $\mathbf{u}_m(t)$  can be obtained by using different  $\mathbf{x}_m(t)$ . In practice, an *orthogonal* projection method is frequently used to find the approximate solution  $\mathbf{u}_m(t)$ . The requirement of the orthogonal projection method is that the residual vector resulting from the approximate solution  $\mathbf{u}_m(t)$

$$\mathbf{r}_m(t) = \mathbf{M} \ddot{\mathbf{u}}_m(t) + \mathbf{C} \dot{\mathbf{u}}_m(t) + \mathbf{K} \mathbf{u}_m(t) - \mathbf{f}(t) \quad (5)$$

be orthogonal to the basis  $\mathbf{Y}_m$ , i.e.

$$\mathbf{Y}_m^T \mathbf{r}_m(t) = 0 \quad (6)$$

Substituting in turn (5), (4) and its time derivatives into equation (6), we obtain an  $m \times m$  reduced system :

$$\mathbf{M}_m^* \ddot{\mathbf{x}}_m(t) + \mathbf{C}_m^* \dot{\mathbf{x}}_m(t) + \mathbf{K}_m^* \mathbf{x}_m(t) = \mathbf{f}_m^*(t) \quad (7)$$

which is actually the orthogonal projection of the original system, represented by (2), onto the subspace  $\mathbf{Y}_m$ . The projected mass, damping, and stiffness matrices and the force vector are given by the following

$$\mathbf{M}_m^* = \mathbf{Y}_m^T \mathbf{M} \mathbf{Y}_m \quad (7a)$$

$$\mathbf{C}_m^* = \mathbf{Y}_m^T \mathbf{C} \mathbf{Y}_m \quad (7b)$$

$$\mathbf{K}_m^* = \mathbf{Y}_m^T \mathbf{K} \mathbf{Y}_m \quad (7c)$$

$$\mathbf{f}_m^*(t) = \mathbf{Y}_m^T \mathbf{f}(t) \quad (7d)$$

That is, we solve the reduced system represented by (7) to find  $\mathbf{x}_m(t)$  and then use (4) to construct an approximate solution  $\mathbf{u}_m(t)$  to the original system represented by (2).

It is important to note that in the orthogonal projection method the quality of the approximate solution  $\mathbf{u}_m(t)$  depends entirely on the chosen set  $\mathbf{Y}_m$  since  $\mathbf{x}_m(t)$  is the solution of the reduce system, which is completely determined by  $\mathbf{Y}_m$ . In the following we compare two sets of subspace : the eigenspace spanned by the least dominant set of eigenvectors, and the Krylov subspace represented by  $\text{span}[\mathbf{b}, \mathbf{D}\mathbf{b}, \mathbf{D}^2\mathbf{b}, \dots, \mathbf{D}^{m-1}\mathbf{b}]$  with  $\mathbf{D}=\mathbf{K}^{-1}\mathbf{M}$  and  $\mathbf{b}=\mathbf{K}^{-1}\hat{\mathbf{f}}$ . An interesting physical interpretation for this Krylov subspace was given in Reference 1.

In practical computation the projected mass matrix  $\mathbf{M}_m^*$  is usually scaled to be the identity matrix  $\mathbf{I}_m$ . To be consistent, the definition of length (norm) and orthogonality should be generalized such that they are in terms of the  $\mathbf{M}$ -weighted inner product, defined by  $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{M}} = \mathbf{u}^T \mathbf{M} \mathbf{v}$ , rather than the conventional inner product  $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{v}$ . That is, the  $\mathbf{M}$ -norm of a vector  $\mathbf{v}$  is  $\|\mathbf{v}\|_{\mathbf{M}} = (\mathbf{v}^T \mathbf{M} \mathbf{v})^{1/2}$  and a set of vectors  $\mathbf{Y}_m$  are  $\mathbf{M}$ -orthonormal if  $\mathbf{Y}_m^T \mathbf{M} \mathbf{Y}_m = \mathbf{I}_m$ . For simplicity, we assume in this study that the damping matrix is proportional and can be expressed as

$$\mathbf{C} = a \mathbf{M} + b \mathbf{K} \quad (8)$$

with  $a$  and  $b$  being constants. Under such assumption, the projected system given by (7) can always be decoupled by the eigenvectors of the pencil  $(\mathbf{M}^*, \mathbf{K}^*)$ .

## EIGENVECTORS

In this section, we discuss using the eigenvectors, i.e., the solutions of  $\omega_j^2 \mathbf{M} \phi_j = \mathbf{K} \phi_j$ , as a basis to express the solution. It is convenient to scale the eigenvectors  $\Phi_m = [\phi_1, \dots, \phi_m]$  such that they satisfy the  $\mathbf{M}$ -orthonormality and  $\mathbf{K}$ -orthogonality :

$$\Phi_m^T \mathbf{M} \Phi_m = \mathbf{I}_m \quad \Phi_m^T \mathbf{K} \Phi_m = \Omega_m^2 \quad (9)$$

where  $\mathbf{I}_m$  is the identity matrix of dimension  $m$  and  $\Omega_m$  is a diagonal matrix of dimension  $m$  with elements  $\omega_j$ .

In order to take advantage of this  $\mathbf{M}$ -orthonormality and  $\mathbf{K}$ -orthogonality, we first pre-multiply (2) by  $\mathbf{M}^{-1}$  to obtain a new equation :

$$\ddot{\mathbf{u}}(t) + (a \mathbf{I} + b \mathbf{M}^{-1} \mathbf{K}) \dot{\mathbf{u}}(t) + \mathbf{M}^{-1} \mathbf{K} \mathbf{u}(t) = \mathbf{M}^{-1} \hat{\mathbf{f}} \boldsymbol{\epsilon}(t) \quad (10)$$

Here the operator  $\mathbf{M}^{-1} \mathbf{K}$  is unsymmetric but is symmetric with respect to the  $\mathbf{M}$ -weighted inner product since

$$(\mathbf{M}^{-1} \mathbf{K} \mathbf{u})^T \mathbf{M} (\mathbf{v}) = \mathbf{u}^T \mathbf{K} \mathbf{v} = (\mathbf{u})^T \mathbf{M} (\mathbf{M}^{-1} \mathbf{K} \mathbf{v})$$

Therefore, we can use the  $\mathbf{M}$ -orthogonal projection method to find an approximate solution to equation (10). By substituting  $\mathbf{u}_m(t) = \Phi_m \mathbf{v}(t)$  and its time derivatives into (10), and requiring that the resulting residual be  $\mathbf{M}$ -orthogonal to  $\Phi_m$ , we obtain

$$\begin{aligned} \Phi_m^T \mathbf{M} \Phi_m \ddot{\mathbf{v}}(t) + (a \Phi_m^T \mathbf{M} \Phi_m + b \Phi_m^T \mathbf{K} \Phi_m) \dot{\mathbf{v}}(t) \\ + \Phi_m^T \mathbf{K} \Phi_m \mathbf{v}(t) = \Phi_m^T \mathbf{M} \mathbf{M}^{-1} \hat{\mathbf{f}} \boldsymbol{\epsilon}(t) \end{aligned} \quad (11)$$

or, after using (9),

$$\ddot{\mathbf{v}}(t) + (a \mathbf{I}_m + b \Omega_m^2) \dot{\mathbf{v}}(t) + \Omega_m^2 \mathbf{v}(t) = \Phi_m^T \hat{\mathbf{f}} \boldsymbol{\epsilon}(t) \quad (12)$$

which is a set of decoupled equations. Introducing the notation  $2\xi_j \omega_j$  for  $a + b\omega_j^2$ , we can write the  $j^{\text{th}}$  equation of the set (12) as

$$\ddot{v}_j(t) + 2 \xi_j \omega_j \dot{v}_j(t) + \omega_j^2 v_j(t) = \phi_j^T \hat{\mathbf{f}} \boldsymbol{\epsilon}(t) \quad (13)$$

where the participation factor  $\phi_j^T \hat{\mathbf{f}}$  is generally used to measure the extent to which the  $\phi_j$  participates in synthesizing the total load on the system. We note from the above derivation that the participation factor  $\phi_j^T \hat{\mathbf{f}}$  is the  $\mathbf{M}$ -weighted inner product of the basis vector  $\phi_j$  and the right-hand side vector  $\mathbf{M}^{-1} \hat{\mathbf{f}}$  in (10). If we denote  $\phi_j^T \hat{\mathbf{f}}$  by  $p_j$ , then it follows that

$$\hat{\mathbf{f}} = \sum_{j=1}^n \mathbf{M} \phi_j p_j = \mathbf{M} \Phi \mathbf{p} \quad (14a)$$

or, after pre-multiplying both sides by  $\mathbf{M}^{-1}$ ,

$$\mathbf{M}^{-1} \hat{\mathbf{f}} = \sum_{j=1}^n \phi_j p_j = \Phi \mathbf{p} \quad (14b)$$

From this expression, we see that the participation factor  $p_j = \phi_j^T \hat{\mathbf{f}}$  is the component of  $\mathbf{M}^{-1} \hat{\mathbf{f}}$  in terms of the basis  $\phi_j$ . The  $\mathbf{M}$ -norm of the  $\mathbf{M}^{-1} \hat{\mathbf{f}}$  is

$$\| \mathbf{M}^{-1} \hat{\mathbf{f}} \|_{\mathbf{M}} = [ (\mathbf{M}^{-1} \hat{\mathbf{f}})^T \mathbf{M} (\mathbf{M}^{-1} \hat{\mathbf{f}}) ]^{1/2}$$



$$= [ \hat{\mathbf{f}}^T \mathbf{M}^{-1} \hat{\mathbf{f}} ]^{1/2} = [ \mathbf{p}^T \mathbf{p} ]^{1/2} = [ \sum_{j=1}^n p_j^2 ]^{1/2} \quad (15)$$

Therefore, we can scale the  $p_j$  by this norm to obtain the *normalized participation factor*

$$\frac{p_j}{(\sum_{j=1}^n p_j^2)^{1/2}} = \frac{\phi_j^T \hat{\mathbf{f}}}{[ \hat{\mathbf{f}}^T \mathbf{M}^{-1} \hat{\mathbf{f}} ]^{1/2}} \quad (16)$$

which actually is equal to the cosine of the angle between the force vector  $\hat{\mathbf{f}}$  and  $\phi_j$ . Since the square sum of these cosines is equal to one, we can assess how adequately the load vector  $\hat{\mathbf{f}}$  is represented in  $\Phi_m$  by examining the quantity :

$$1 - \frac{\sum_{j=1}^m (\phi_j^T \hat{\mathbf{f}})^2}{(\hat{\mathbf{f}}^T \mathbf{M}^{-1} \hat{\mathbf{f}})} = \frac{\sum_{j=m+1}^n p_j^2}{\sum_{j=1}^n p_j^2} \quad (17)$$

When this quantity is equal to 0, the force vector  $\hat{\mathbf{f}}$  is completely represented.

The participation factor is used to measure how adequately the spatial distribution of the load  $\mathbf{f}(t)$  is represented. To assess how accurate the approximate solution is, we examine the *error vector*  $\mathbf{e}_m(t)$ , defined as the difference between the exact solution  $\mathbf{u}(t)$  and the approximate solution  $\mathbf{u}_m(t)$ , i.e.

$$\mathbf{e}_m(t) = \mathbf{u}(t) - \mathbf{u}_m(t) = \sum_{j=m+1}^n \phi_j v_j(t) \quad (18)$$

The  $\mathbf{M}$ -norm of this error vector is given by

$$\|\mathbf{e}_m(t)\|_{\mathbf{M}} = [ \mathbf{e}_m^T(t) \mathbf{M} \mathbf{e}_m(t) ]^{1/2} = [ \sum_{j=m+1}^n v_j^2(t) ]^{1/2} \quad (19)$$

which gives the *absolute error* of the approximation. We can divide the absolute error by the  $\mathbf{M}$ -norm of the exact solution vector, i.e.

$$\|\mathbf{u}(t)\|_{\mathbf{M}} = [ \mathbf{u}^T(t) \mathbf{M} \mathbf{u}(t) ]^{1/2} = [ \sum_{j=1}^n v_j^2(t) ]^{1/2} \quad (20)$$

to obtain the *relative error*

$$\frac{\|\mathbf{e}_m(t)\|_{\mathbf{M}}}{\|\mathbf{u}(t)\|_{\mathbf{M}}} = \frac{[ \sum_{j=m+1}^n v_j^2(t) ]^{1/2}}{[ \sum_{j=1}^n v_j^2(t) ]^{1/2}} \quad (21)$$

We see that the expression of the relative error has the same structure as the expression in (17) except the square root in (21). Moreover, the  $p_j$  and  $v_j(t)$  are directly related through the solution to equation (13)

$$v_j(t) = \frac{p_j}{\omega_j^2} D_j(t) \quad (22)$$

where the *dynamic load factor*  $D_j(t)$  is given by

$$D_j(t) = \frac{\omega_j^2}{\tilde{\omega}_j} \int_0^t e^{-\xi_j \omega_j (t-\tau)} \sin \tilde{\omega}_j (t-\tau) \epsilon(\tau) d\tau \quad (23)$$

with  $\tilde{\omega}_j = \omega_j (1 - \xi_j^2)^{1/2}$ . Equation (22) simply indicates that the solution  $v_j(t)$  is equal to the static response  $p_j / \omega_j^2$  multiplied by the dynamic load factor  $D_j(t)$ . The fact that the expression of  $v_j(t)$  has a factor  $1 / \omega_j^2$  justifies that the higher modes are, in general, less important than the lower modes when the eigenvectors are used to find the solution. It follows from (21) and (22) that one can safely terminate the superposition process when the quantity in (17) is close to 0 unless there exists such a higher mode whose  $\xi_j$  is very small and  $\omega_j$  is very close to the frequency content of the loading so that  $D_j(t)$  is very large.

## LANCZOS VECTORS

We can apply the Gram-Schmidt orthonormalization process to the columns of the Krylov subspace  $\text{span}[\mathbf{b}, \mathbf{D}\mathbf{b}, \mathbf{D}^2\mathbf{b}, \dots, \mathbf{D}^{m-1}\mathbf{b}]$ , the results is a set of Lanczos vectors. A detailed theoretical description and numerical analysis on this subject can be found in Reference 7. In the following, we summarize the steps used to derive the Lanczos vectors for the convenience of discussion. In practical computation these steps have to be rearranged and reinforced by a reorthogonalization scheme to prevent the loss of orthogonality due to the roundoff errors.

Given an arbitrary vector  $\mathbf{b}$ , we normalized it to obtain the first Lanczos vector as

$$\beta_1 = (\mathbf{b}^T \mathbf{M} \mathbf{b})^{1/2} \quad (24a)$$

$$\mathbf{q}_1 = \mathbf{b} / \beta_1 \quad (24b)$$

and compute the rest of the Lanczos vectors,  $j = 1, \dots, m-1$ , by

$$\beta_{j+1} \mathbf{q}_{j+1} = \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{q}_j - \alpha_j \mathbf{q}_j - \beta_j \mathbf{q}_{j-1} \quad (24c)$$

where  $\mathbf{K}_\sigma = \mathbf{K} - \sigma \mathbf{M}$  with  $\sigma$  being an appropriate shift,  $\mathbf{q}_0 = \mathbf{0}$ , and

$$\alpha_j = \mathbf{q}_j^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{q}_j \quad (24d)$$

$$\beta_{j+1} = \mathbf{q}_{j+1}^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{q}_j \quad (24e)$$

After these steps, we have  $m$  Lanczos vectors  $\mathbf{Q}_m = [ \mathbf{q}_1, \dots, \mathbf{q}_m ]$  satisfying the matrix form of the three-term recurrence formula :

$$\beta_{m+1} \mathbf{q}_{m+1} \hat{\mathbf{e}}_m^T = \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m - \mathbf{Q}_m \mathbf{T}_m \quad (25)$$

where  $\hat{\mathbf{e}}_m^T$  is the  $m$ th row of  $\mathbf{I}_m$ , and  $\mathbf{T}_m$  is a tri-diagonal matrix made up of the coefficients  $\alpha_j$  and  $\beta_j$  :

$$\mathbf{T}_m = \begin{bmatrix} \alpha_1 & \beta_2 & & & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & & & \\ & & \cdot & \cdot & \cdot & & & \\ & & & \beta_{m-1} & \alpha_{m-1} & \beta_m & & \\ & & & & \beta_m & \alpha_m & & \end{bmatrix} \quad (26)$$

The Lanczos vectors obtained in this way are  $\mathbf{M}$ -orthonormal, i.e., they satisfy

$$\mathbf{Q}_m^T \mathbf{M} \mathbf{Q}_m = \mathbf{I}_m \quad (27)$$

where  $\mathbf{I}_m$  is the identity matrix of dimension  $m$ . After pre-multiplying (25) by  $\mathbf{Q}_m^T \mathbf{M}$  and using (27) we can obtain

$$\mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m = \mathbf{T}_m \quad (28)$$

That is, the  $\mathbf{M}$ -orthogonal projection of  $\mathbf{K}_\sigma^{-1} \mathbf{M}$  onto the Krylov subspace with basis  $\mathbf{Q}_m$  is a tri-diagonal matrix.

To take advantage of the  $\mathbf{M}$ -orthonormality and tri-diagonality between the Lanczos vectors, we pre-multiply the equation of motion (after using  $\mathbf{K} = \mathbf{K}_\sigma + \sigma \mathbf{M}$ )

$$\mathbf{M} \ddot{\mathbf{u}}(t) + [ (a + b \sigma) \mathbf{M} + b \mathbf{K}_\sigma ] \dot{\mathbf{u}}(t) + ( \mathbf{K}_\sigma + \sigma \mathbf{M} ) \mathbf{u}(t) = \hat{\mathbf{f}} \epsilon(t) \quad (29)$$

by  $\mathbf{K}_\sigma^{-1}$  to obtain a new equation :

$$\begin{aligned} \mathbf{K}_\sigma^{-1} \mathbf{M} \ddot{\mathbf{u}}(t) + [ (a + b \sigma) \mathbf{K}_\sigma^{-1} \mathbf{M} + b \mathbf{I} ] \dot{\mathbf{u}}(t) \\ + ( \mathbf{I} + \sigma \mathbf{K}_\sigma^{-1} \mathbf{M} ) \mathbf{u}(t) = \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} \epsilon(t) \end{aligned} \quad (30)$$

Here the operator  $\mathbf{K}_\sigma^{-1} \mathbf{M}$  is unsymmetric but is symmetric with respect to the  $\mathbf{M}$ -weighted inner product since

$$(\mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{u})^T \mathbf{M} (\mathbf{v}) = \mathbf{u}^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{v} = (\mathbf{u})^T \mathbf{M} (\mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{v}) \quad (31)$$

Therefore we can apply the  $\mathbf{M}$ -orthogonal projection method to equation (30) to find an approximate solution. By substituting  $\mathbf{u}_m(t) = \mathbf{Q}_m \mathbf{x}_m(t)$  and its time derivatives into equation (30), and requiring that the resulting residual be  $\mathbf{M}$ -orthogonal to  $\mathbf{Q}_m$ , we obtain

$$\begin{aligned} \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m \ddot{\mathbf{x}}_m(t) + [ (a + b \sigma) \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m + b \mathbf{Q}_m^T \mathbf{M} \mathbf{Q}_m ] \dot{\mathbf{x}}_m(t) \\ + ( \mathbf{Q}_m^T \mathbf{M} \mathbf{Q}_m + \sigma \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m ) \mathbf{x}_m(t) = \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} \epsilon(t) \end{aligned} \quad (32)$$

Using equations (27) and (28), we may rewrite the above equation to

$$\begin{aligned} \mathbf{T}_m \ddot{\mathbf{x}}_m(t) + [ (a + b \sigma) \mathbf{T}_m + b \mathbf{I}_m ] \dot{\mathbf{x}}_m(t) \\ + ( \mathbf{I}_m + \sigma \mathbf{T}_m ) \mathbf{x}_m(t) = \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} \epsilon(t) \end{aligned} \quad (33)$$

This reduced system is only slightly coupled and can be solved either by step-by-step integration methods or by finding eigensolutions and using mode superposition.

Here it should be emphasized that since equation (33) is a coupled system, we have to distinguish the  $j$ th component  $x_{j,m}(t)$  of the  $\mathbf{x}_m(t)$  from the  $j$ th component  $x_{j,m-1}(t)$  of the  $\mathbf{x}_{m-1}(t)$  for  $j = 1, \dots, m-1$ , although the difference between them may be slight for some  $j$ . The fact that the Lanczos coordinates  $x_{j,m}(t)$  changes when we increase the dimension  $m$  of the approximating subspace should be contrasted with the fact that the eigenvector coordinates  $v_j(t)$  remains fixed when we increase  $m$  as shown in Section 3. We will elucidate this point further by transforming the reduced set of equations into a decoupled set in the next Section.

If we denote the  $j$ th component  $\mathbf{q}_j^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}$  of the right-hand side vector in (33) by  $\hat{p}_j$ , then it follows that

$$\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \sum_{j=1}^n \mathbf{q}_j \hat{p}_j = \mathbf{Q} \hat{\mathbf{p}} \quad (34a)$$

That is, the consistent way of expressing the force vector  $\hat{\mathbf{f}}$  in terms of the basis  $\mathbf{Q}$  is

$$\hat{\mathbf{f}} = \sum_{j=1}^n \mathbf{K}_\sigma \mathbf{q}_j \hat{p}_j = \mathbf{K}_\sigma \mathbf{Q} \hat{\mathbf{p}} \quad (34b)$$

This expression is different from (14) for the eigenvectors. Moreover, it is not appropriate to call the  $\hat{p}_j$  participation factor because equation (33) is not a decoupled set.

The starting vector of the Lanczos algorithm is arbitrary, so we choose it to maximize the efficiency. For this purpose, we can choose the starting vector  $\mathbf{b}$  to be the *dynamic deflection shape*  $\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}$  so that  $\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \mathbf{b} = \beta_1 \mathbf{q}_1$ . By choosing  $\mathbf{b}$  this way, the right-hand side of the reduced system (33) is simplified to

$$\mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \hat{\mathbf{p}}_m^T = (\beta_1, 0, \dots, 0) \quad (35)$$

This indicates that the first element  $\hat{p}_1$  is equal to  $\beta_1$  and the remaining elements of the  $\hat{\mathbf{p}}_m$  are all equal to 0, independent of  $m$ . Physically, this choice guarantees that the vector  $\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}$  is completely represented. Algebraically, this choice tends to make  $x_{j,m}(t)$  decrease to zero as  $j$  increases so that the superposition process will converge quickly.

To see how to choose the shift  $\sigma$ , we consider a special case where the system is undamped and the loading  $\epsilon(t)$  is a harmonic function, say  $\sin \bar{\omega}(t)$ . Under this circumstance, the reduced system is simply

$$\mathbf{T}_m \ddot{\mathbf{x}}_m(t) + (\mathbf{I}_m + \sigma \mathbf{T}_m) \mathbf{x}_m(t) = \hat{\mathbf{p}}_m \sin \bar{\omega}(t) \quad (36)$$

If  $\sigma$  is chosen to be  $\bar{\omega}^2$ , the solution to the equation (36) is  $\mathbf{x}_m(t) = \hat{\mathbf{p}}_m \sin \bar{\omega}(t)$ . In this case, only  $x_{1,m}(t)$  is needed because all other  $x_{j,m}(t)$  are equal to 0. Hence, by choosing the dynamic deflection shape,  $(\mathbf{K} - \bar{\omega}^2 \mathbf{M})^{-1} \hat{\mathbf{f}}$ , as the first Lanczos vector we can find the exact solution,  $(\mathbf{K} - \bar{\omega}^2 \mathbf{M})^{-1} \hat{\mathbf{f}} \sin \bar{\omega}(t)$ , by only one vector and no more other Lanczos vectors are required. This is not surprising because the original system can be solved directly when the loading function  $\epsilon(t)$  is harmonic. On the other hand, if  $\sigma$  is not chosen to be  $\bar{\omega}^2$ , then in general all the  $x_{j,m}(t)$  will not be equal to 0 and accordingly more Lanczos vectors are needed to obtain a better approximation. Guided by this special case, we can pick a suitable shift to accelerate the convergence of the superposition process by using the information on the frequency content of the loading function, such as the Fourier expansion of the  $\epsilon(t)$ , when dealing with general cases where the loading is not harmonic or the system is damped, or both.

The error vector  $\mathbf{e}_m(t)$  of the approximation is

$$\mathbf{e}_m(t) = \mathbf{u}(t) - \mathbf{u}_m(t) = \mathbf{Q}_n \mathbf{x}_n(t) - \mathbf{Q}_m \mathbf{x}_m(t)$$

$$= \sum_{j=1}^n \mathbf{q}_j x_{j,n}(t) - \sum_{j=1}^m \mathbf{q}_j x_{j,m}(t) \quad (37)$$

The  $\mathbf{M}$ -norm of this error vector can be computed as

$$\begin{aligned} \|\mathbf{e}_m(t)\|_{\mathbf{M}} &= [\mathbf{e}_m^{\mathbf{T}}(t) \mathbf{M} \mathbf{e}_m(t)]^{1/2} \\ &= [\mathbf{x}_n^{\mathbf{T}}(t) \mathbf{Q}_n^{\mathbf{T}} \mathbf{M} \mathbf{Q}_n \mathbf{x}_n(t) - 2 \mathbf{x}_n^{\mathbf{T}}(t) \mathbf{Q}_n^{\mathbf{T}} \mathbf{M} \mathbf{Q}_m \mathbf{x}_m(t) + \mathbf{x}_m^{\mathbf{T}}(t) \mathbf{Q}_m^{\mathbf{T}} \mathbf{M} \mathbf{Q}_m \mathbf{x}_m(t)]^{1/2} \\ &= \left[ \sum_{j=1}^n x_{j,n}^2(t) - 2 \sum_{j=1}^m x_{j,n}(t) x_{j,m}(t) + \sum_{j=1}^m x_{j,m}^2(t) \right]^{1/2} \end{aligned} \quad (38)$$

where the  $\mathbf{M}$ -orthonormality between the Lanczos vectors has been used in the manipulation. The relative error is then given by

$$\frac{\|\mathbf{e}_m(t)\|_{\mathbf{M}}}{\|\mathbf{u}(t)\|_{\mathbf{M}}} = \frac{\left[ \sum_{j=1}^n x_{j,n}^2(t) - 2 \sum_{j=1}^m x_{j,n}(t) x_{j,m}(t) + \sum_{j=1}^m x_{j,m}^2(t) \right]^{1/2}}{\left[ \sum_{j=1}^n x_{j,n}^2(t) \right]^{1/2}} \quad (39)$$

We note that these expressions for the absolute and relative errors are different from those given in Section 3. Besides, there is no simple connection between  $\hat{p}_j$  and  $x_{j,m}(t)$ .

## RITZ VECTORS

To decouple the system represented by equation (33), we solve the eigenproblem

$$\mathbf{T}_m \mathbf{S}_m = \mathbf{S}_m \mathbf{\Theta}_m^{-2} \quad (40)$$

to find the set of vectors  $\mathbf{S}_m = [\mathbf{s}_{1,m}, \dots, \mathbf{s}_{m,m}]$  which satisfies

$$\mathbf{S}_m^{\mathbf{T}} \mathbf{S}_m = \mathbf{I}_m \quad \mathbf{S}_m^{\mathbf{T}} \mathbf{T}_m \mathbf{S}_m = \mathbf{\Theta}_m^{-2} \quad (41)$$

where the  $\mathbf{\Theta}_m$  is a diagonal matrix and contains the Ritz values  $\theta_j$ . Introducing the transformation  $\mathbf{x}_m(t) = \mathbf{S}_m \mathbf{z}_m(t)$  and pre-multiplying both sides of equation (33) by  $\mathbf{S}_m^{\mathbf{T}}$ , we obtain the decoupled set of equations :

$$\begin{aligned} \mathbf{\Theta}_m^{-2} \dot{\mathbf{z}}_m(t) + [ (a + b \sigma) \mathbf{\Theta}_m^{-2} + b \mathbf{I}_m ] \mathbf{z}_m(t) \\ + ( \mathbf{I}_m + \sigma \mathbf{\Theta}_m^{-2} ) \mathbf{z}_m(t) = \mathbf{Y}_m^{\mathbf{T}} \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} \epsilon(t) \end{aligned} \quad (42)$$

where the Ritz vectors  $\mathbf{Y}_m$  is equal to  $\mathbf{Q}_m \mathbf{S}_m$  and has the following properties :

$$\mathbf{Y}_m^T \mathbf{M} \mathbf{Y}_m = \mathbf{S}_m^T \mathbf{Q}_m^T \mathbf{M} \mathbf{Q}_m \mathbf{S}_m = \mathbf{S}_m^T \mathbf{S}_m = \mathbf{I}_m \quad (43a)$$

$$\mathbf{Y}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Y}_m = \mathbf{S}_m^T \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M} \mathbf{Q}_m \mathbf{S}_m = \mathbf{S}_m^T \mathbf{T}_m \mathbf{S}_m = \Theta_m^{-2} \quad (43b)$$

That is, the set of Ritz vectors  $\mathbf{Y}_m$  satisfies the  $\mathbf{M}$ -orthonormality and the  $\mathbf{M} \mathbf{K}_\sigma^{-1} \mathbf{M}$ -orthogonality. With these properties, we can easily verify that (42) actually is the projected system resulted from using the Ritz vectors as a basis to express the approximate solution. This projected system is equivalent to the one given by (33) since  $\mathbf{Y}_m$  and  $\mathbf{Q}_m$  span the same subspace and the solution  $\mathbf{Y}_m \mathbf{z}_m(t)$  is the same as the solution  $\mathbf{Q}_m \mathbf{x}_m(t)$ . Therefore we can use (42) to analyze the approximate solution  $\mathbf{u}_m(t)$  obtained by solving equation (33).

As in previous discussion, if we denote the  $j$ th component  $y_j^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}$  of the right-hand side vector in (42) by  $\tilde{p}_j$ , then it follows that

$$\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \sum_{j=1}^n y_j \tilde{p}_j = \mathbf{Y} \tilde{\mathbf{p}} \quad (44a)$$

which is equivalent to

$$\hat{\mathbf{f}} = \sum_{j=1}^n \mathbf{K}_\sigma y_j \tilde{p}_j = \mathbf{K}_\sigma \mathbf{Y} \tilde{\mathbf{p}} \quad (44b)$$

Note that this expression is similar to (34) but different from (14). The  $\mathbf{M}$ -norm of the vector  $\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}$  is

$$\|\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}\|_{\mathbf{M}} = (\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}})^T \mathbf{M} (\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}) = \tilde{\mathbf{p}}^T \mathbf{Y}^T \mathbf{M} \mathbf{Y} \tilde{\mathbf{p}} = \tilde{\mathbf{p}}^T \tilde{\mathbf{p}} = \sum_{j=1}^n \tilde{p}_j^2 \quad (45a)$$

or, in terms of  $\hat{\mathbf{p}}$ ,

$$\|\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}\|_{\mathbf{M}} = (\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}})^T \mathbf{M} (\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}}) = \hat{\mathbf{p}}^T \mathbf{Q}^T \mathbf{M} \mathbf{Q} \hat{\mathbf{p}} = \hat{\mathbf{p}}^T \hat{\mathbf{p}} = \sum_{j=1}^n \hat{p}_j^2 \quad (45b)$$

The relation between the right-hand side vectors  $\tilde{\mathbf{p}}_m$  and  $\hat{\mathbf{p}}_m$  is given by

$$\tilde{\mathbf{p}}_m = \mathbf{Y}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \mathbf{S}_m^T \mathbf{Q}_m^T \mathbf{M} \mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \mathbf{S}_m^T \hat{\mathbf{p}}_m \quad (46)$$

Using this relation and  $\mathbf{S}_m \mathbf{S}_m^T = \mathbf{I}_m$ , we can verify that  $\tilde{\mathbf{p}}_m^T \tilde{\mathbf{p}}_m$  is identically equal to  $\hat{\mathbf{p}}_m^T \hat{\mathbf{p}}_m$ , as it should. Further, if the first Lanczos vector is chosen to be the dynamic deflection shape  $\mathbf{K}_\sigma^{-1} \hat{\mathbf{f}} = \beta_1 \mathbf{q}_1$ , we have

$$\tilde{p}_{j,m} = \mathbf{s}_{j,m}^T \hat{\mathbf{p}}_m = s_{1,j,m} \beta_1 \quad (47)$$

where  $s_{1,j,m}$  is the first element of  $\mathbf{s}_{j,m}$ , the  $j$ th vector of  $\mathbf{S}_m$ . Here we use an extra subscript  $m$  to label the components  $\tilde{p}_{j,m}$  of the right-hand side vector in (42) because the values of these components will change as the dimension  $m$  of the subspace is increased. Recall that the components  $p_j$  of the right-hand side vector in eqn (13) are independent of  $m$ . Therefore, we should distinguish the interpretation of the  $\tilde{p}_{j,m}$  from the interpretation of the  $p_j$  in Section 3. In particular, we need not compute the quantity  $\sum_{j=1}^m \tilde{p}_{j,m}^2$  in the present case since it is always equal to  $\beta_1^2$  due to the fact that  $\mathbf{S}_m \mathbf{S}_m^T = \mathbf{I}_m$  independent of  $m$ .

To relate the  $z_{j,m}(t)$  to  $\tilde{p}_{j,m}$ , we can solve the  $j$ th equation of the projected system (42)

$$\ddot{z}_{j,m}(t) + [a + b(\sigma + \theta_{j,m}^2)] \dot{z}_{j,m}(t) + (\theta_{j,m}^2 + \sigma) z_{j,m}(t) = \theta_{j,m}^2 \tilde{p}_{j,m} \epsilon(t) \quad (48)$$

Introducing the notation  $\tilde{\theta}_{j,m}^2$  for  $\theta_{j,m}^2 + \sigma$  and  $2\eta_{j,m} \tilde{\theta}_{j,m}$  for  $a + b\tilde{\theta}_{j,m}^2$ , we can rewrite this equation to

$$\ddot{z}_{j,m}(t) + 2\eta_{j,m} \tilde{\theta}_{j,m} \dot{z}_{j,m}(t) + \tilde{\theta}_{j,m}^2 z_{j,m}(t) = \tilde{p}_{j,m} \theta_{j,m}^2 \epsilon(t) \quad (49)$$

The solution to this equation is

$$z_{j,m}(t) = \tilde{p}_{j,m} \frac{\theta_{j,m}^2}{\tilde{\theta}_{j,m}^2} D_{j,m}(t) \quad (50)$$

where the *dynamic load factor*  $D_{j,m}(t)$  is given by (23) but with  $\omega_j$  and  $\xi_j$  replaced by  $\tilde{\theta}_{j,m}$  and  $\eta_{j,m}$  respectively. Here we note that in the expression of  $z_{j,m}(t)$  there is a factor of  $\theta_{j,m}^2 / \tilde{\theta}_{j,m}^2$ , which is equal to one when the shift  $\sigma$  is chosen to be zero. This indicates that the higher modes have the same weight as the lower modes, which is again different from the fact given by (22) where the eigenvectors is used. We can obtain the error expression in terms of the Ritz coordinates  $z_{j,m}(t)$  simply by replacing the  $x$ 's in (38) and (39) by the  $z$ 's. However, since  $\sum_{j=1}^m \tilde{p}_{j,m}^2$  is always equal to  $\beta_1^2$  and  $z_{j,m+1}(t)$  will be different from  $z_{j,m}(t)$ , we cannot judge when to terminate the superposition by the same criterion as stated in section 3.



## RECOMMENDATIONS

Although the Lanczos vectors or the Ritz vectors are superior to the eigenvectors for computing dynamic responses in a common class of problems, the fundamental mechanisms of the three sets of vectors are not the same. When the Lanczos vectors or the Ritz vectors are used, the conventional criterion for deciding how many vectors are enough is theoretical not sound. As shown in previous sections, the load vector may be completely represented by a set of Lanczos vectors or Ritz vectors due to the special choice of the starting vector, but it does not follow that the exact solution will be obtained by this set of Lanczos vectors or Ritz vectors. The error in the response, as defined in this paper, rather than the error in the representation of the load vector is actually what one is concerned with. More effort is needed to construct a convenient but justifiable error estimate so that one can know how many Lanczos vectors or Ritz vectors are needed for a good approximation.

## REFERENCES

- 1 Wilson, E. L., Yuan M. and Dickens, J. M. Dynamic analysis by direct superposition of Ritz vectors, *Earthquake Eng. Struct. Dyn.*, **10**, 813-821 (1982)
- 2 Bayo, E. P. and Wilson, E. L. Use of Ritz vectors in wave propagation and foundation response, *Earthquake Eng. Struct. Dyn.*, **12**, 499-505 (1984)
- 3 Bayo, E. P. and Wilson, E. L. Finite element and Ritz vector techniques for the solution to three-dimensional soil-structure interaction problems in the time domain, *Eng. Comput.*, **1**, 298-311 (1984)
- 4 Arnold, R. R., Citerley, R. L., Chargin, M. and Galant, D. Application of Ritz vectors for dynamics analysis of large structures, *Comput. Struct.*, **21**, 901-907 (1985)
- 5 Wilson, E. L. and Bayo, E. P. Use of special Ritz vectors in dynamic and substructure analysis', *J. Struct. Engng. Div., ASCE*, **112** (1986)
- 6 Nour-Omid, B. and Clough, R. W. Dynamic analysis of structures using Lanczos coordinates, *Earthquake Eng. Struct. Dyn.*, **12**, 565-577 (1984)
- 7 B. N. Parlett, *The symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ (1980)