

## Chapter 5

# Convergence

In Section 3.4 we presented three general problems, Problems A, B and C. In Section 4.1 we formulated the problems for the Galerkin approximations. These are Problems AG, BG and CG.

### 5.1 Equilibrium problem

In this section we consider the convergence of the solution of Problem BG to the solution of Problem B.

Assume that  $u^h \in S^h$  is the solution of

$$b(u^h, v) = (f, v) \text{ for all } v \in S^h \quad (5.1.1)$$

and  $u \in V$  is the solution of

$$b(u, v) = (f, v) \text{ for all } v \in V. \quad (5.1.2)$$

In the proof of the theorem, we will use the projection  $P$ , defined in Section 4.6.

#### Theorem 5.1.1

1.  $\|u^h - u\|_E \rightarrow 0$  if  $h \rightarrow 0$ .

2. If  $u \in H^4 \cap V$ , then

$$\|u^h - u\|_E \leq \widehat{C}h^{k^*-2}|u|_4,$$

and

$$\|u^h - u\|_0 \leq \widehat{C}h^{k^*}|u|_4.$$

**Proof** If (5.1.2) is subtracted from (5.1.1), it follows that

$$b(u - u^h, v) = 0 \text{ for all } v \in S^h.$$

This means that  $u^h = Pu$ . The first part of the theorem follows directly from Lemma 4.6.2. The estimates in the second part of the theorem follow from Lemma 4.6.1 and Lemma 4.6.3.  $\square$

## 5.2 Eigenvalue problem

Our approach is based on the ideas of [BDSW], [BF] and [SF] and we follow the presentation in [SF]. The theory in the book of Strang and Fix, [SF, Section 6.3], concerns eigenvalue problems for general symmetric elliptic operators. Most of the presentation is written in a style which encourage abstraction. In collaboration with others, [ZVGV2], we verified that the theory is valid for abstract eigenvalue problems such as Problem C. In this thesis we present this abstract version, and also offer a number of modest improvements.

The rate of convergence for eigenvalue problems also depends on the regularity of the eigenvectors. In the absence of such theory for interface problems, we pose the following assumption which we showed to be true in the one-dimensional case. (See Section 3.5.)

**Regularity Assumption** The eigenvectors of the eigenvalue problem, Problem C, are in  $H^k \cap V$  for  $k = 4$  or  $6$ , and there exists a constant  $C_b$ —depending on the bilinear forms  $b$  and  $(\cdot, \cdot)$ —such that for each eigenvector  $y$

$$\|y\|_k \leq C_b \lambda_i^\alpha \|y\|_0, \text{ where } \alpha = \frac{k}{4}.$$

### 5.2.1 The Rayleigh quotient and the Minmax principle

To analyse the convergence of eigenvalues and eigenvectors, some preparation is necessary. First we establish bounds for the approximate eigenvalues using the Rayleigh quotient and the minimax principle.

It is well-known that the eigenvalues are the stationary values of the Rayleigh quotient. However, the following result gives a more convenient characterization of the eigenvalues. See [SF, p 221].

**Lemma 5.2.1** *Minmax principle*

Let  $\mathcal{T}$  denote the class of subspaces of  $V$  having dimension  $j$ , then

$$\lambda_j = \min_{S \in \mathcal{T}} \max_{v \in S} R(v).$$

We may assume that the eigenvalues are ordered

$$\lambda_1 \leq \lambda_2 \leq \lambda_3 \dots$$

For some integer  $m$ , consider the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_m$  and corresponding eigenvectors  $y_1, y_2, \dots, y_m$ . Equal eigenvalues are possible but we assume that  $\lambda_j \neq \lambda_m$  for each  $j > m$ .

In the finite dimensional subspace  $S^h$  we have  $\lambda_1^h, \lambda_2^h, \dots, \lambda_m^h$  (also ordered) and corresponding eigenvectors  $y_1^h, y_2^h, \dots, y_m^h$ . (Equal eigenvalues do not matter. In the case of multiplicity,  $y_j^h$  is not uniquely determined, but it does not influence any proof.)

### 5.2.2 Bounds for the approximate eigenvalues

The minimax principle yields lower bounds for the approximate eigenvalues.

**Lemma 5.2.2**  $\lambda_j^h \geq \lambda_j$  for each  $j$ .

**Proof** The minimax principle is also true for the space  $S^h$ . Any subspace of  $S^h$  is a subspace of  $V$ .  $\square$

**Notation**  $y_i$  will be used to denote the normalised eigenvectors, i.e.  $\|y_i\| = 1$ .

For  $j = 1, 2, \dots, m$ , let  $E_j$  denote the subspace of  $V$  spanned by  $\{y_1, y_2, \dots, y_j\}$ .

Consider the subspaces  $S_j$  where

$$S_j = PE_j \text{ for } j = 1, 2, \dots, m.$$

$P = P_h$  is the projection defined in Section 4.6.

An upper bound for some approximate eigenvalue depends on the construction of  $S^h$ . Clearly the construction of  $S^h$  must be such that  $\dim S^h = N > m$ . However, it is still possible that  $\dim S_m < m$ .

**Assumption** The construction of  $S^h$  is such that  $\dim S_m = m$ .

We define a quantity  $\mu_m^h$  to measure the “distortion” of the projection of the unit ball  $B_m = \{y \in E_m : \|y\| = 1\}$ : Set

$$\mu_m^h = \inf\{\|Py\|^2 : y \in B_m\}.$$

**Proposition 5.2.1**  $\mu_m^h > 0$  if and only if  $\dim S_m = m$ .

**Proof** The function  $\|Py\|^2$  has a minimum on the compact set  $B_m$ . Hence  $\mu_m^h > 0$  if and only if  $Py \neq 0$  for each  $y \in B_m$ . But this is so if and only if the vectors  $Py_1, Py_2, \dots, Py_m$  form a linearly independent set.  $\square$

**Proposition 5.2.2**  $\lambda_m^h \leq \max\{R(Py) : y \in B_m\}$ .

**Proof**

Since  $\dim S_m = m$ , it follows from the minimax principle that

$$\lambda_m^h \leq \max\{R(v) : v \in S_m\}.$$

Consider any  $v \in S_m$ ,  $v \neq 0$ . There exists a vector  $y \in E_m$  such that  $Py = v$ .

Note that  $R(P(\alpha y)) = R(\alpha v) = R(v)$ . Choose  $\alpha$  such that  $\alpha y \in B_m$ . Consequently,

$$\max\{R(Pz) : z \in B_m\} = \max\{R(v) : v \in S_m\}.$$

$\square$

The following result is crucial.

**Lemma 5.2.3**  $\lambda_m^h \leq \frac{\lambda_m}{\mu_m^h}$ .

**Proof** If  $y = \sum_{i=1}^m c_i y_i$ , then  $b(y, y) = \sum_{i=1}^m c_i^2 \lambda_i$ , since  $\{y_1, y_2, \dots, y_m\}$  is an orthonormal set. Hence

$$b(y, y) \leq \lambda_m \sum_{i=1}^m c_i^2 = \lambda_m \|y\|^2 \text{ for each } y \in E_m.$$

Since  $P$  is a projection with respect to the inner product  $b$ ,

$$b(Py, Py) \leq \lambda_m \text{ for each } y \in B_m.$$

From the definition of  $\mu_m^h$ , we have

$$R(Py) = \frac{b(Py, Py)}{\|Py\|^2} \leq \frac{\lambda_m}{\mu_m^h}.$$

Now use Proposition 5.2.2.  $\square$

**Corollary 5.2.1**  $\mu_m^h \leq 1$ .

This is a direct consequence of Lemmas 5.2.2 and 5.2.3. It is convenient to formulate error estimates in terms of the quantity  $1 - \mu_m^h$ .

**Notation**  $\sigma_m^h = 1 - \mu_m^h$ .

**Corollary 5.2.2**  $0 \leq \sigma_m^h < 1$  and  $\lambda_m^h - \lambda_m \leq \lambda_m \sigma_m^h$ .

Since the eigenvalue error is bounded by  $\sigma_m^h$ , it is sufficient to estimate  $\sigma_m^h$  and prove that  $\sigma_m^h \rightarrow 0$ .

### 5.2.3 Estimates

**Proposition 5.2.3**  $\sigma_m^h = \max\{2(y, y - Py) - \|y - Py\|^2 : y \in B_m\}$ .

**Proof**

$$\begin{aligned} \|y - Py\|^2 &= \|Py\|^2 + \|y\|^2 - 2(y, Py) \\ &= \|Py\|^2 + 2(y, y) - 2(y, Py) - 1 \quad (\text{since } \|y\|^2 = 1). \end{aligned}$$

As a consequence  $1 - \|Py\|^2 = 2(y, y - Py) - \|y - Py\|^2$ . The result follows from the definition of  $\sigma_m^h$ .  $\square$

**Remark** In [SF, Section 6.3]  $\sigma_m^h$  is defined by

$$\sigma_m^h = \max\{|2(y, y - Py) - \|y - Py\|^2| : y \in B_m\}.$$

The absolute value is not necessary since

$$\max\{2(y, y - Py) - \|y - Py\|^2 : y \in B_m\} \geq 0.$$

The assumption is then made that  $\sigma_m^h < 1$ , and they prove that  $\dim S_m = m$ .

We proved the fact that  $\dim S_m = m$  is equivalent to  $\sigma_m^h < 1$  and we believe that it is important to take note of this equivalence.

**Notation** For any  $y \in E_m$ , let  $y^* = \sum_{i=1}^m c_i \lambda_i^{-1} y_i$  where  $y = \sum_{i=1}^m c_i y_i$ .

**Proposition 5.2.4** For any  $y \in E_m$ ,

$$(y, y - Py) = b(y^* - Py^*, y - Py).$$

**Proof**

$$b(y_i - Py_i, y - Py) = b(y_i, y - Py) \text{ since } b(y - Py, Py_i) = 0.$$

Hence,

$$\lambda_i(y_i, y - Py) = b(y_i - Py_i, y - Py).$$

Multiply by  $c_i \lambda_i^{-1}$  and sum over  $i$ . We have

$$\begin{aligned} (y, y - Py) &= \sum_{i=1}^m c_i \lambda_i^{-1} b(y_i - Py_i, y - Py) \\ &= b \left( \sum_{i=1}^m c_i \lambda_i^{-1} y_i - \sum_{i=1}^m c_i \lambda_i^{-1} Py_i, y - Py \right). \end{aligned}$$

□

The following result also differ from [SF].

**Lemma 5.2.4**  $\sigma_m^h \leq \max\{2\|y^* - Py^*\|_E \|y - Py\|_E : y \in B_m\}$ .

**Proof** Consider the result of Proposition 5.2.3. We have demonstrated that the quantity

$$2(y, y - Py) - \|y - Py\|^2$$

must have a non negative maximum (Corollary 5.2.2). Consequently

$$\sigma_m^h \leq \max\{2(y, y - Py) : y \in B_m\}.$$

Use Proposition 5.2.4 and the Schwartz inequality for the inner product  $b$ . □

**Proposition 5.2.5** For any  $\varepsilon > 0$  there exists a  $\delta > 0$  such that for  $h < \delta$ ,

$$\begin{aligned} \|y^* - Py^*\|_E &< \varepsilon \text{ for each } y \in B_m, \\ \|y - Py\|_E &< \varepsilon \text{ for each } y \in B_m. \end{aligned}$$

**Proof** From Lemma 4.6.2 there exist positive numbers  $\delta_1, \delta_2, \dots, \delta_n$  such that for each  $i$

$$\|y_i - Py_i\|_E \leq \varepsilon \text{ if } h < \delta_i.$$

Now, suppose  $h < \min_i \delta_i$ , then

$$\|y - Py\|_E \leq \sum_{i=1}^m |c_i| \|y_i - Py_i\|_E \leq \varepsilon m, \text{ since } |c_i| < 1.$$

The same arguments are valid for  $\|y^* - Py^*\|_E$ . □

**Lemma 5.2.5** *For any  $\varepsilon > 0$  there exists a  $\delta > 0$  such that*

$$\sigma_m^h < \varepsilon \text{ if } h < \delta.$$

**Proof**

For any  $\varepsilon > 0$  there exists a  $\delta > 0$  such that if  $h < \delta$ , then

$$\|u - Pu\|_E < \varepsilon \text{ for each } u \in B_m.$$

The result follows from Lemma 5.2.4 and Proposition 5.2.5. □

**Proposition 5.2.6** *If Problem C satisfies the regularity assumption, then for any  $y \in B_m$*

$$\|y^* - Py^*\|_E \leq \widehat{C} C_b \lambda_m^{\alpha-1} h^{k^*-2}$$

and

$$\|y - Py\|_E \leq \widehat{C} C_b \lambda_m^\alpha h^{k^*-2}.$$

**Proof** We may assume that  $c_i \geq 0$  for each  $i$ .

First estimate:

$$\begin{aligned} \|y^* - Py^*\|_E &\leq \sum_{i=1}^m c_i \lambda_i^{-1} \|y_i - Py_i\|_E \\ &\leq \widehat{C} \sum_{i=1}^m c_i \lambda_i^{-1} |y_i|_{k^*} h^{k^*-2} \\ &\leq C_b \widehat{C} \sum_{i=1}^m c_i \lambda_i^{\alpha-1} \|y_i\| h^{k^*-2} \\ &\leq \widehat{C} C_b \lambda_m^{\alpha-1} h^{k^*-2}. \end{aligned}$$



Second estimate:

$$\begin{aligned} \|y - Py\|_E &\leq \sum_{i=1}^m c_i \|y_i - Py_i\|_E \\ &\leq \widehat{C} C_b \lambda_i^\alpha h^{k^* - 2}, \end{aligned}$$

using the same arguments as for the first estimate.  $\square$

**Lemma 5.2.6** *If Problem C satisfies the regularity assumption, then*

$$\sigma_m^h \leq \widehat{C} C_b \lambda_m^{2\alpha - 1} h^{2(k^* - 2)}.$$

**Proof**

Use Lemma 5.2.4 and Proposition 5.2.6.  $\square$

## 5.2.4 Convergence of eigenvalues

We may now use the results of the previous subsection to establish the convergence of  $\lambda_m^h$  to  $\lambda_m$ .

**Lemma 5.2.7** *There exists a  $\delta > 0$  such that for  $h < \delta$ ,*

$$\lambda_m^h - \lambda_m \leq 2\lambda_m \sigma_m^h.$$

**Proof**

Choose  $\delta$  such that  $\sigma_m^h < \frac{1}{2}$ . Consequently  $\lambda_m^h < 2\lambda_m$ .  $\square$

**Theorem 5.2.1**

1.  $\lambda_m^h - \lambda_m \rightarrow 0$  as  $h \rightarrow 0$ .
2. *If Problem C satisfies the regularity assumption, then*

$$\lambda_m^h - \lambda_m \leq \widehat{C} C_b \lambda_m^{2\alpha} h^{2(k^* - 2)}.$$

**Proof**

1. This is a direct consequence of Lemmas 5.2.5 and 5.2.7.
2. Use Lemmas 5.2.6 and 5.2.7.  $\square$

### 5.2.5 Convergence of eigenvectors

To estimate the error  $\|y_m - y_m^h\|$ , we need to estimate the difference  $\|Py_m - y_m^h\|$ . It is necessary to consider the possibility that  $\lambda_m$  has multiplicity more than one. Suppose that the multiplicity of  $\lambda_m$  is  $r$  and let  $\Lambda = 1, 2, \dots, m - r, m + 1, \dots, N$ . From Theorem 5.2.1 it follows that there exist real numbers  $\rho > 0$  and  $\delta > 0$  such that if  $h < \delta$ , then

$$|\lambda_m - \lambda_j^h| > \rho \text{ for each } j \in \Lambda. \quad (5.2.1)$$

**Assumption** Assume that  $h$  is sufficiently small for (5.2.1) to hold.

Suppose  $\{y_{m-r+1}^h, y_{m-r+2}^h, \dots, y_m^h\}$  is an orthonormal set of eigenvectors corresponding to  $\lambda_{m-r+1}^h, \lambda_{m-r+2}^h, \dots, \lambda_m^h$ . The strategy now is to estimate the distance between  $y_{m-r+i}^h$  and some (uniquely defined) vector in  $E_{\lambda_m}$ , the eigenspace corresponding to  $\lambda_m$ .

We define a projection  $P_m$  with domain  $P(E_{\lambda_m})$ :

$$P_m w = \sum_{j=m-r+1}^m (w, y_j^h) y_j^h \text{ for each } w \in P(E_{\lambda_m}).$$

This projection enables us to deal with the case of a repeated eigenvalue.

Here we differ from [SF]. Although most of the computations are the same, we believe that our construction of the projection  $P_m$  is a worthwhile contribution. We will show that  $P_m P$  (and hence  $P_m$ ) is invertible for  $h$  sufficiently small.

**Proposition 5.2.7** For each  $j \in \Lambda$  and each  $y \in E_{\lambda_m}$ ,

$$(\lambda_j^h - \lambda_m)(Py, y_j^h) = \lambda_m(y - Py, y_j^h).$$

**Proof**

It is only necessary to prove that

$$\lambda_j^h(Py, y_j^h) = \lambda_m(y, y_j^h) \quad (5.2.2)$$

since the term  $-\lambda_m(Py, y_j^h)$  appears on both sides of the equation.

Since  $y_j^h$  and  $y$  are eigenvectors, it follows that

$$\lambda_j^h(Py, y_j^h) = b(Py, y_j^h) \text{ and } \lambda_m(y, y_j^h) = b(y, y_j^h).$$

But  $b(Py - y, y_j^h) = 0$  for each  $j$ , thus (5.2.2) follows.  $\square$

**Lemma 5.2.8**

$$\|Py - P_m Py\| \leq \lambda_m \rho^{-1} \|y - Py\| \text{ for each } y \in E_{\lambda_m}$$

**Proof** From the assumption we have the estimate

$$\frac{\lambda_m}{|\lambda_j^h - \lambda_m|} \leq \rho_m \text{ for each } j \in \Lambda, \quad (5.2.3)$$

where  $\rho_m = \lambda_m \rho^{-1}$ .

The set  $y_1^h, y_2^h, \dots, y_N^h$  form an orthonormal basis for  $S^h$ , hence

$$Py = \sum_{j=1}^N (Py, y_j^h) y_j^h.$$

Consequently,

$$Py - P_m Py = \sum_{j \in \Lambda} (Py, y_j^h) y_j^h.$$

If  $y \in E_{\lambda_m}$ , then

$$\|Py - P_m Py\|^2 = \sum_{j \in \Lambda} (Py, y_j^h)^2.$$

We now use Proposition 5.2.7.

$$\begin{aligned} \|Py - P_m Py\|^2 &= \sum_{j \in \Lambda} \left( \frac{\lambda_m}{|\lambda_j^h - \lambda_m|} \right)^2 (y - Py, y_j^h)^2 \\ &\leq \rho_m^2 \sum_{j \in \Lambda} (y - Py, y_j^h)^2 \quad (\text{Inequality (5.2.3)}) \\ &\leq \rho_m^2 \sum_{j=1}^N (y - Py, y_j^h)^2 \\ &= \rho_m^2 \|y - Py\|^2. \end{aligned}$$

□

**Lemma 5.2.9**

$$\|y - P_m P y\| \leq (1 + \lambda_m \rho^{-1}) \|y - P y\| \text{ for each } y \in E_{\lambda_m}.$$

**Proof**

$$\begin{aligned} \|y - P_m P y\| &\leq \|y - P y\| + \|P y - P_m P y\| \\ &\leq (1 + \lambda_m \rho^{-1}) \|y - P y\|. \end{aligned}$$

□

**Corollary 5.2.3**  $P_m P$  is invertible for  $h$  sufficiently small.

**Proof** Let  $y \in B_m \cap E_{\lambda_m}$ . Then

$$\|P y - P_m P y\| \leq \lambda_m \rho^{-1} \|y - P y\| < \frac{1}{2},$$

for  $h$  sufficiently small. Since

$$\|P y\|^2 = \|P_m P y\|^2 + \|P y - P_m P y\|^2,$$

it follows that  $\|P_m P y\| > \frac{1}{4}$ . Consequently

$$\|P_m P y\| > \frac{1}{4} \|y\| \text{ for each } y \in E_{\lambda_m}.$$

□

**Corollary 5.2.4** If  $h$  is sufficiently small, then for each  $j$ ,  $j = 1, 2, \dots, r$  there exists a unique  $x_j \in E_{\lambda_m}$  with  $\|x_j\| = 1$  such that

$$\|x_j - y_{m-r+j}^h\| \leq 4(1 + \rho^{-1} \lambda_m) \|x_j - P x_j\|.$$

**Proof** There exists a unique  $y \in E_{\lambda_m}$  such that  $y = (P_m P)^{-1} y_{m-r+j}$ . Hence,

$$\|y - y_{m-r+j}^h\| \leq (1 + \rho^{-1} \lambda_m) \|y - P y\|.$$

Let  $\beta$  be a real number such that  $|\beta| = \|y\|$  and let  $x_j = \beta^{-1} y$ . We can choose  $x_j$  such that  $\beta > 0$ . As a consequence  $\|y\| = \beta$ .

It follows that

$$\|x_j - y\| = |\beta - 1| = \left| \|y\| - \|y_{m-r+j}^h\| \right| \leq \|y - y_{m-r+j}^h\|.$$

Hence

$$\begin{aligned} \|x_j - y_{m-r+j}^h\| &\leq \|x_j - y\| + \|y - y_{m-r+j}^h\| \\ &\leq 2(1 + \rho^{-1}\lambda_m)\|y - Py\|. \end{aligned}$$

□

It is important to realise that one compute the approximation  $y_{m-r+j}^h$ . The result above guarantees the existence of an exact eigenvector, with norm one close to the approximate one.

The following result from [SF] shows that an error estimate in the energy norm depends on error estimates in the norm  $\|\cdot\|$  and eigenvalue errors. We modified it slightly to make it useful for the case of repeated eigenvalues.

**Lemma 5.2.10**

$$\|y_m - y_j^h\|_E^2 = \lambda_m \|y_m - y_j^h\|^2 + \lambda_j^h - \lambda_m, \quad j \notin \Lambda.$$

**Proof**

$$\begin{aligned} b(y_m - y_j^h, y_m - y_j^h) &= b(y_m, y_m) - 2b(y_m, y_j^h) + b(y_j^h, y_j^h) \\ &= \lambda_m \|y_m\|^2 - 2\lambda_m (y_m, y_j^h) + \lambda_j^h \|y_j^h\|^2 \\ &= \lambda_m - 2\lambda_m (y_m, y_j^h) + \lambda_j^h \\ &= \lambda_m [2 - 2(y_m, y_j^h)] + \lambda_j^h - \lambda_m \\ &= \lambda_m [\|y_m\|^2 - 2(y_m, y_j^h)] + \|y_j^h\|^2 + \lambda_j^h - \lambda_m \\ &= \lambda_m \|y_m - y_j^h\|^2 + \lambda_j^h - \lambda_m. \end{aligned}$$

□

**Theorem 5.2.2**

1. Let  $\varepsilon > 0$  be arbitrary. If  $h$  is sufficiently small, then for each  $j$ ,  $j = 1, 2, \dots, r$  there exists a unique  $x_j \in E_{\lambda_m}$  with  $\|x_j\| = 1$  such that

$$\|x_j - y_{m-r+j}^h\|_E \leq \varepsilon.$$

2. Suppose Problem C satisfies the regularity assumption. If  $h$  is sufficiently small, then for each  $j$ ,  $j = 1, 2, \dots, r$  there exists a unique  $x_j \in E_{\lambda_m}$  with  $\|x_j\| = 1$  such that

$$\|x_j - y_{m-r+j}^h\|_E \leq \widehat{C} C_b \lambda_m^\alpha h^{(k^*-2)}.$$

**Proof**

1. Use Corollary 5.2.4: There exists a unique  $x_j \in E_{\lambda_m}$  such that

$$\|x_j - y_{m-r+j}^h\| \leq 2(1 + \rho^{-1}\lambda_m)\|x_j - Px_j\|.$$

But

$$\|x_j - y_{m-r+j}^h\|_E^2 \leq \lambda_m \|x_j - y_{m-r+j}^h\|^2 + \lambda_{m-r+j}^h - \lambda_m$$

(Lemma 5.2.10). Hence

$$\|x_j - y_{m-r+j}^h\|_E^2 \leq 4\lambda_m(1 + \rho^{-1}\lambda_m)^2\|x_j - Px_j\|^2 + \lambda_{m-r+j}^h - \lambda_m \quad (5.2.4)$$

Now use Proposition 5.2.5 and Theorem 5.2.1.

2. Consider the Inequality (5.2.4). We have the estimates

$$\|x_j - Px_j\| \leq \widehat{C} C_b \lambda_m^\alpha h^{k^*-2} \quad (5.2.5)$$

from Proposition 5.2.6 and

$$\lambda_{m-r+j}^h - \lambda_m \leq \widehat{C} C_b \lambda_m^{2\alpha} h^{2(k^*-2)} \quad (5.2.6)$$

from Theorem 5.2.1. The result follows from (5.2.5) and (5.2.6).

□

## 5.3 Vibration problem

Our concern is the difference between the solution  $u$  of Problem  $A$  and the solution  $u_h$  of Problem  $AG$ . It is possible to estimate this error in terms of the projection error (Section 4.6) and errors for the initial conditions. See [SF, Section 7.3]. This is called a projection method and was first used for parabolic problems. For second order hyperbolic problems, it appear that credit is due to [D], [De] and [SF]. Research in this direction was also done by [Ba].

After this it appear that abstract methods became popular. See for example [Sh, Section 6.4]. In Section 3 of an invited paper, [BIt], very general results are given. (Incidentally they use results in [Sh].)

A general approximation theory, using functional analysis, is obviously important. However, we found that the basic error inequality mentioned before ([SF, Section 7.3] and [D, Lemma 1]) is valid for an abstract problem as general as Problem  $A$ . As a final remark we mention the paper [FXX] where the authors also use what they term a “partial projection” method to obtain  $L^2$ -error estimates.

### 5.3.1 Discretization error

In this section we show that the convergence proof sketched by Strang and Fix [SF, Section 7.3] can be applied to Problem  $A$  in Section 3.4 and Problem  $AG$  in Section 4.1. In this proof the projection operator  $P$  defined in Section 4.6 is used to find an estimate on the discretization error

$$\|u(t) - u_h(t)\|_E \text{ for } t \in [0, \infty).$$

We also use the symbol  $P$  to denote the “projection”  $Pu$  of the solution  $u$  of Problem  $A$ , i.e.  $(Pu)(t) = Pu(t)$  for each  $t \geq 0$ .

Let  $e(t) = Pu(t) - u_h(t)$  and  $e_p(t) = u(t) - Pu(t)$ . Then

$$\|u(t) - u_h(t)\|_E \leq \|e_p(t)\|_E + \|e(t)\|_E. \quad (5.3.1)$$

The following result is required for the main result of this section. Note that differentiability with respect to the energy norm is required to prove that the projection function  $Pu$  is differentiable. This regularity requirement is not stated by [SF].

**Lemma 5.3.1** *If  $u \in C^2([0, \infty), V)$ , then  $Pu \in C^2([0, \infty), V)$  with*

$$(Pu)'(t) = Pu'(t) \text{ and } (Pu)''(t) = Pu''(t).$$

**Proof** As the projection operator  $P$  is a bounded linear operator with norm less than one, it follows that

$$\|(\delta t)^{-1}(Pu(t + \delta t) - Pu(t)) - Pu'(t)\|_E \leq \|(\delta t)^{-1}(u(t + \delta t) - u(t)) - u'(t)\|_E.$$

This implies that  $Pu \in C^1([0, \infty), V)$  and  $(Pu)'(t) = Pu'(t)$ .

In exactly the same way we prove that  $(Pu)' \in C^1([0, \infty), V)$  and  $(Pu)''(t) = Pu''(t)$ .

□

Since we already have an estimate for the projection error  $e_p(t)$ , it is only necessary to estimate the other part of the error.

In the next proof the following “energy” expression will be convenient:

$$\begin{aligned} E(t) &= \frac{1}{2}(e'(t), e'(t)) + \frac{1}{2}b(e(t), e(t)) \\ &= \frac{1}{2}\|e'(t)\|^2 + \frac{1}{2}\|e(t)\|_E^2. \end{aligned} \quad (5.3.2)$$

**Lemma 5.3.2** *Assume that  $u \in C^2([0, \infty), V)$ . Then, for any  $t \geq 0$ ,*

$$\|e(t)\|_E \leq \|P\alpha - \alpha_h\|_E + \|P\beta - \beta_h\| + \int_0^t \|e_p''(s)\| + \frac{k}{C_I}\|e_p'(s)\|_0 ds.$$

**Proof** From Problem A and the Galerkin approximation (Problem AG) we deduce that for any  $v \in S^h$ ,

$$(u''(t) - u_h''(t), v) + a(u'(t) - u_h'(t), v) + b(u(t) - u_h(t), v) = 0. \quad (5.3.3)$$

Since  $P$  is a projection, we have

$$b(u(t) - Pu(t), v) = b(u'(t) - Pu'(t), v) = 0 \text{ for all } v \in S^h.$$

Using the fact that  $Pu''(t) = (Pu)''(t)$ , (5.3.3) can be written as

$$\begin{aligned} (e''(t), v) + b(e(t), v) &= -(e_p''(t), v) - k(e_p'(t), v)_0 - k(e'(t), v)_0 \\ &\quad - \mu b(e'(t), v) \text{ for all } v \in S^h. \end{aligned} \quad (5.3.4)$$



(Note that  $a(u, v) = \mu b(u, v) + k(u, v)_0$  where  $\mu$  or  $k$  or both can be zero.)

We will use the fact that

$$E'(t) = (e''(t), e'(t)) + b(e(t), e'(t)).$$

As  $e(t) \in S^h$  it follows that  $e'(t) \in S^h$ . Choose  $v = e'(t)$  in (5.3.4), then

$$\begin{aligned} E'(t) &= -(e''_p(t), e'(t)) - (ke'_p(t), e'(t))_0 - k(e'(t), e'(t))_0 - \mu b(e'(t), e'(t)) \\ &\leq \left( \|e''_p(t)\| + \frac{k}{C_I} \|e'_p(t)\|_0 \right) \|e'(t)\|. \end{aligned}$$

From (5.3.2),  $\|e'(t)\| \leq \sqrt{2E(t)}$ . Thus

$$E'(t) \leq \sqrt{2E(t)} \left( \|e''_p(t)\| + \frac{k}{C_I} \|e'_p(t)\|_0 \right)$$

and consequently

$$\frac{d}{dt} \sqrt{E(t)} \leq \frac{1}{\sqrt{2}} \left( \|e''_p(t)\| + \frac{k}{C_I} \|e'_p(t)\|_0 \right).$$

This yields that

$$\sqrt{E(t)} \leq \sqrt{E(0)} + \frac{1}{\sqrt{2}} \int_0^t \left( \|e''_p(s)\| + \frac{k}{C_I} \|e'_p(s)\|_0 \right) ds. \quad (5.3.5)$$

As

$$E(0) = \frac{1}{2} \|P\beta - \beta_h\|^2 + \frac{1}{2} \|P\alpha - \alpha_h\|_E^2$$

and  $\|e(t)\|_E \leq \sqrt{2E(t)}$ , again from (5.3.2), the result follows from (5.3.5).  $\square$

**Theorem 5.3.1** *Assume that  $u \in C^2([0, \infty), V)$ . Then, for any  $t \geq 0$ .*

$$\begin{aligned} \|u(t) - u_h(t)\|_E &\leq \|e_p(t)\|_E + \|P\alpha - \alpha_h\|_E + \|P\beta - \beta_h\| \\ &\quad + \int_0^t \left( \|e''_p(s)\| + \frac{k}{C_I} \|e'_p(s)\|_0 \right) ds. \end{aligned}$$

**Proof** Use Lemma 5.3.2 and Equation (5.3.1).  $\square$

To prove the convergence results, it is now necessary to consider the terms on the right side of the inequality in this theorem.

### 5.3.2 Convergence

The main factor that determines the rate of convergence of the solution  $u_h$  of Problem AG to the solution  $u$  of Problem A as  $h$  tends to zero, is the regularity of the weak solution  $u$ . The regularity of  $u$  depends on the regularity of the initial values  $\alpha$  and  $\beta$ , as we pointed out in Section 3.4. [Ra] gave an example to show that the regularity of the solution is necessary to obtain optimal order convergence.

The rate of convergence is also directly influenced by the choice of the initial values  $\alpha_h$  and  $\beta_h$  for the solution  $u_h$  of Problem AG. We will consider two cases, i.e.  $\alpha_h = \Pi\alpha$ ,  $\beta_h = \Pi\beta$  and  $\alpha_h = P\alpha$ ,  $\beta_h = P\beta$ . In the following result we show that the rate of convergence in the energy norm is of order  $h^2$  if certain regularity conditions are satisfied. The estimates are expressed in terms of the constants  $C_I$  and  $C_E$  defined in Section 3.4 as well as  $\widehat{C}$  defined in Section 4.5.

**Theorem 5.3.2** *Let  $\alpha_h = \Pi\alpha$  and  $\beta_h = \Pi\beta$ . Assume that  $u \in C^2([0, \infty), V)$  and that  $u(t)$ ,  $u'(t)$  and  $u''(t)$  are in  $H^4 \cap V$  for  $t \geq 0$ . Then,*

$$\|u(t) - u_h(t)\|_E \leq \widehat{C} \left( |\alpha|_4 + C_E^{-1} |\beta|_4 + |u(t)|_4 + k(C_E C_I)^{-1} t \max_{s \in [0, t]} |u'(s)|_4 + C_E^{-1} t \max_{s \in [0, t]} |u''(s)|_4 \right) h^2 \text{ for } t \in [0, \infty).$$

**Proof** From Theorem 5.3.1,

$$\|u(t) - u_h(t)\|_E \leq \|e_p(t)\|_E + \|P\alpha - \Pi\alpha\|_E + \|P\beta - \Pi\beta\|_E + \int_0^t \left( \|e_p''(s)\|_E + \frac{k}{C_I} \|e_p'(s)\|_0 \right) ds.$$

All that remains to be done is to apply the approximation results from Corollary 4.6.1 to each of the terms in this expression:

$$\|P\alpha - \Pi\alpha\|_E \leq \widehat{C} |\alpha|_4 h^2,$$

$$\|P\beta - \Pi\beta\|_E \leq C_E^{-1} \|P\beta - \Pi\beta\|_E \leq C_E^{-1} \widehat{C} |\beta|_4 h^2$$

and

$$\|e_p(t)\|_E = \|u(t) - Pu(t)\|_E \leq \widehat{C} |u(t)|_4 h^2.$$

From Lemma 5.3.1,  $(Pu)' = Pu'$  and hence  $e_p'(t) = u'(t) - Pu'(t)$ . This yields that

$$\|e_p'(s)\|_0 \leq (C_E C_I)^{-1} \|e_p'(s)\|_E \leq (C_E C_I)^{-1} \widehat{C} |u'(s)|_4 h^2$$

and

$$\int_0^t \|e_p'(s)\|_0 ds \leq (C_E C_I)^{-1} \widehat{C} t \max_{s \in [0, t]} |u'(s)|_4 h^2.$$

Similarly,

$$\int_0^t \|e_p''(s)\| ds \leq C_E^{-1} \widehat{C} t \max_{s \in [0, t]} |u''(s)|_4 h^2.$$

□

Under less strict regularity conditions we can still show that the solution  $u_h$  of Problem AG converges to the solution  $u$  of Problem A in the energy norm if  $h$  tends to zero.

**Theorem 5.3.3** *Let  $\alpha_h = \Pi\alpha$  and  $\beta_h = \Pi\beta$ . Assume that  $\alpha \in V$ ,  $\beta \in V$  and  $u \in C^2([0, \infty), V)$ , then*

$$\lim_{h \rightarrow 0} \|u(t) - u_h(t)\|_E = 0 \text{ for } t \in [0, \tau].$$

**Proof** From Theorem 5.3.1,

$$\begin{aligned} \|u(t) - u_h(t)\|_E &\leq \|e_p(t)\|_E + \|P\alpha - \Pi\alpha\|_E + \|P\beta - \Pi\beta\| \\ &\quad + \int_0^t \left( \|e_p''(s)\| + \frac{k}{C_I} \|e_p'(s)\|_0 \right) ds. \end{aligned}$$

From the approximation results we know that for any  $\varepsilon > 0$ , each term is less than  $\varepsilon$ , provided that  $h$  is sufficiently small. □

### 5.3.3 Inertia norm estimate

In a final result we show that the Aubin-Nitsche trick can also be applied to this problem to find inertia norm estimates for the discretization error.

**Theorem 5.3.4** Let  $\alpha_h = P\alpha$  and  $\beta_h = P\beta$ . Assume that  $u(t)$ ,  $u'(t)$  and  $u''(t)$  are all in  $V \cap H^4$  for all  $t \geq 0$ . Then,

$$\|u(t) - u_h(t)\| \leq \widehat{C} \left( |u(t)|_4 + kt(C_I^2 C_E)^{-1} \max_{s \in [0, t]} |u'(s)|_4 + tC_E^{-1} \max_{s \in [0, t]} |u''(s)|_4 \right) h^4$$

for  $t \in [0, \infty)$ .

**Proof** From Theorem 5.3.2,

$$\begin{aligned} \|u(t) - u_h(t)\| &\leq \|e_p(t)\| + \|e(t)\| \\ &\leq \|e_p(t)\| + C_E^{-1} \|e(t)\|_E \\ &\leq \|e_p(t)\| + C_E^{-1} \int_0^t \left( \|e_p''(s)\| + \frac{k}{C_I} \|e_p'(s)\|_0 \right) ds. \end{aligned}$$

For a fixed  $t \geq 0$ , we consider  $e_p(t) = u(t) - Pu(t)$ .

We conclude from Corollary 4.6.2 that

$$\|e_p(t)\| \leq \widehat{C} |u(t)|_4 h^4.$$

Similar arguments yield that

$$\|e_p'(t)\|_0 \leq C_I^{-1} \|e_p'(t)\| \leq C_I^{-1} \widehat{C} |u'(t)|_4 h^4 \text{ and } \|e_p''(t)\| \leq \widehat{C} |u''(t)|_4 h^4. \quad (5.3.6)$$

□

A useful result is also obtained if the Aubin-Nitsche trick is used only for the terms containing the integrals.

**Theorem 5.3.5** Let  $\alpha_h = \Pi\alpha$  and  $\beta_h = \Pi\beta$ . Assume that  $\alpha$ ,  $\beta$ ,  $u(t)$ ,  $u'(t)$  and  $u''(t)$  are all in  $V \cap H^4$  for all  $t$ . Then,

$$\|u(t) - u_h(t)\|_E \leq \widehat{C} (|\alpha|_4 + C_E^{-1} |\beta|_4 + |u(t)|_4) h^2 +$$

$$\widehat{C} \left( ktC_I^{-2} \max_{s \in [0, t]} |u'(s)|_4 + t \max_{s \in [0, t]} |u''(s)|_4 \right) h^4 \text{ for } t \in [0, \infty).$$

**Proof** The proof is exactly the same as the proof of Theorem 5.3.2. The estimates in (5.3.6) are used for the terms containing the integral. □

**Remark** We consider this result to be significant. It is advantageous to have an error estimate in the energy norm, while the terms containing  $t$  are “suppressed” by  $h^4$ .

## 5.4 Finite Differences

In this section we consider the system of ordinary differential equations, Problem AD in Section 4.1, and the finite difference method for approximating the solution. The objective is to prove that the solution of the discretized problem converges to the solution of the Galerkin approximation. This method has been extensively studied—even in the context of finite difference methods for second order hyperbolic partial differential equations. However, one must be careful when matching the estimates. Although all norms are equivalent in the finite dimensional space  $S^h$ , the “constants” may depend on the dimension of  $S^h$ . Presenting error estimates for semi-discrete and fully discrete systems in the same presentation is a line also followed by others. See for example [D], [Ba] and [FXX].

We consider Problem AG in Section 4.1 and the finite difference scheme proposed in Section 4.4. In the first subsection we estimate the local error and then proceed to establish stability results.

### 5.4.1 Local error

The first step is to derive finite difference formulas similar to the Newmark schemes [Zi]. Since we need error estimates in terms of the unknown function or its derivatives, it is necessary to derive the formulas.

We will use Taylor’s theorem in the following form:

$$g(t) = g(t_0) + (t - t_0)g'(t_0) + \dots + \frac{(t - t_0)^{n-1}}{(n-1)!}g^{(n-1)}(t_0) + R(t)$$

where  $R(t) = \frac{1}{(n-1)!} \int_{t_0}^t (t - \theta)^{n-1} g^{(n)}(\theta) d\theta$ . It is also true for  $t < t_0$ .

See [Cl, p 179] or [Ap, p 279].

The following notation is introduced for convenience.

**Notation**  $R_n^+(t) = \frac{1}{(n-1)!} \int_t^{t+\delta t} (t + \delta t - \theta)^{n-1} g^{(n)}(\theta) d\theta$  and

$$R_n^-(t) = \frac{1}{(n-1)!} \int_t^{t-\delta t} (t - \delta t - \theta)^{n-1} g^{(n)}(\theta) d\theta.$$

The first proposition contains well-known results and the proofs are trivial.

**Proposition 5.4.1**

1. If the real valued function  $g$  is in  $C^3[t - \delta t, t + \delta t]$ , then

$$g(t + \delta t) - g(t - \delta t) = 2\delta t g'(t) + R_3^+(t) - R_3^-(t). \quad (5.4.1)$$

2. If the real valued function  $g$  is in  $C^4[t - \delta t, t + \delta t]$ , then

$$g(t + \delta t) - 2g(t) + g(t - \delta t) = (\delta t)^2 g''(t) + R_4^+(t) + R_4^-(t). \quad (5.4.2)$$

**Proof**

1. Use Taylor's theorem to get:

$$g(t + \delta t) = g(t) + \delta t g'(t) + \frac{(\delta t)^2}{2} g''(t) + R_3^+(t)$$

and

$$g(t - \delta t) = g(t) - \delta t g'(t) + \frac{(\delta t)^2}{2} g''(t) + R_3^-(t).$$

Clearly

$$g(t + \delta t) - g(t - \delta t) = 2\delta t g'(t) + R_3^+(t) - R_3^-(t).$$

2. Approximate  $g$  by a polynomial of degree three and compute  $g(t + \delta t) + g(t - \delta t)$ .

□

We gave the proof of part one in detail because we use the result in the next proposition.

**Proposition 5.4.2** Let  $\rho_0$  and  $\rho_1$  be real numbers such that  $\rho_0 + 2\rho_1 = 1$ .

1. If the real valued function  $g$  is in  $C^4[t - \delta t, t + \delta t]$ , then

$$\begin{aligned} g(t + \delta t) - g(t - \delta t) \\ = 2\delta t (\rho_1 g'(t + \delta t) + \rho_0 g'(t) + \rho_1 g'(t - \delta t)) + \tilde{R}_4(t), \end{aligned} \quad (5.4.3)$$

where

$$\begin{aligned} \tilde{R}_4(t) = \omega_1 \{ R_3^+(t) - R_3^-(t) \} + \omega_2 \left\{ R_4^+(t) + R_4^-(t) - \right. \\ \left. \frac{\delta t}{6} \int_t^{t+\delta t} (t + \delta t - \theta)^2 g^{(4)}(\theta) d\theta - \frac{\delta t}{6} \int_t^{t-\delta t} (t - \delta t - \theta)^2 g^{(4)}(\theta) d\theta \right\}. \end{aligned}$$

2. Suppose the real valued function  $g$  is in  $C^5[t - \delta t, t + \delta t]$ , then

$$g(t + \delta t) - 2g(t) + g(t - \delta t) = (\delta t)^2 (\rho_1 g''(t + \delta t) + \rho_0 g''(t) + \rho_1 g''(t - \delta t)) + \tilde{R}_5(t), \quad (5.4.4)$$

where  $\tilde{R}_5(t) = \omega_1 \{R_4^+(t) + R_4^-(t)\} + \omega_2 \left\{R_5^+(t) + R_5^-(t) - \frac{(\delta t)^2}{24} \int_t^{t+\delta t} (t + \delta t - \theta)^2 g^{(5)}(\theta) d\theta - \frac{(\delta t)^2}{24} \int_t^{t-\delta t} (t - \delta t - \theta)^2 g^{(5)}(\theta) d\theta\right\}$ .

**Proof**

1. Use Taylor's theorem to get:

$$g(t + \delta t) = g(t) + \delta t g'(t) + \frac{(\delta t)^2}{2} g''(t) + \frac{(\delta t)^3}{6} g'''(t) + R_4^+(t)$$

and

$$g(t - \delta t) = g(t) - \delta t g'(t) + \frac{(\delta t)^2}{2} g''(t) - \frac{(\delta t)^3}{6} g'''(t) + R_4^-(t).$$

This yields

$$g(t + \delta t) - g(t - \delta t) = 2\delta t g'(t) + \frac{(\delta t)^3}{3} g'''(t) + R_4^+(t) - R_4^-(t). \quad (5.4.5)$$

Applying Taylor's theorem once more on  $g'$  we obtain

$$g'(t + \delta t) = g'(t) + \delta t g''(t) + \frac{(\delta t)^2}{2} g'''(t) + \frac{1}{2} \int_t^{t+\delta t} (t + \delta t - \theta)^2 g^{(4)}(\theta) d\theta$$

and

$$g'(t - \delta t) = g'(t) - \delta t g''(t) + \frac{(\delta t)^2}{2} g'''(t) + \frac{1}{2} \int_t^{t-\delta t} (t - \delta t - \theta)^2 g^{(4)}(\theta) d\theta.$$

The two equations yield

$$g'(t + \delta t) + g'(t - \delta t) = 2g'(t) + (\delta t)^2 g'''(t) + \frac{1}{2} \int_t^{t+\delta t} (t + \delta t - \theta)^2 g^{(4)}(\theta) d\theta + \frac{1}{2} \int_t^{t-\delta t} (t - \delta t - \theta)^2 g^{(4)}(\theta) d\theta.$$

From this we get an expression for  $(\delta t)^2 g'''(t)$  which can substituted into (5.4.5). The result is

$$\begin{aligned} g(t + \delta t) - g(t - \delta t) &= \frac{(\delta t)}{3} [g'(t + \delta t) + 4g'(t) + g'(t - \delta t)] + R_4^+(t) \\ &\quad + R_4^-(t) - \frac{\delta t}{6} \int_t^{t+\delta t} (t + \delta t - \theta)^2 g^{(4)}(\theta) d\theta \\ &\quad - \frac{\delta t}{6} \int_t^{t-\delta t} (t - \delta t - \theta)^2 g^{(4)}(\theta) d\theta. \end{aligned} \quad (5.4.6)$$

Finally we combine (5.4.1) and (5.4.6) with weights  $\omega_1$  and  $\omega_2$  to get the desired result.

2. This proof is similar to the proof in (1).  $\square$

**Remark** The results above will also be used in the case where the function  $g$  is not defined for  $t < 0$ . In this case we may extend  $g$  by using the polynomial approximation on  $[t - \delta t, 0)$ . This will only influence the result in so far as there will be fewer remainder terms.

The second step is to apply the difference formulas to Problem AG and to estimate the errors.

**Assumption** We assume that  $f \in C^3[0, \tau]$  so that the solution  $u_h$  of Problem AG is in  $C^5[0, \tau]$ .

**Notation**  $\|u_h\|_{5, \max}^E = \sum_{k=0}^5 \max_{t \in [0, \tau]} \|u_h^{(k)}(t)\|_E.$

**Notation** In the rest of this section  $C_b$  will denote a generic constant that depends on the bilinear forms, i.e.  $C_b$  is a combination of  $C_E$  and  $C_I$ .

**Notation**  $\|f\|_{3, \max} = \sum_{k=0}^3 \max_{t \in [0, \tau]} \|f^{(k)}(t)\|.$



**Proposition 5.4.3** Suppose  $u_i \in C^5[t - \delta t, t + \delta t]$  for  $i = 1, 2, \dots, n$  and  $\{\phi_1, \phi_2, \dots, \phi_n\}$  is the basis for  $S^h$  and let  $u_h(t) = \sum_{i=1}^n u_i(t)\phi_i$ . Suppose also that  $\rho_0$  and  $\rho_1$  are real numbers such that  $\rho_0 + 2\rho_1 = 1$ .

If  $u_h(t + \delta t) - 2u_h(t) + u_h(t - \delta t)$

$$= (\delta t)^2(\rho_1 u_h''(t + \delta t) + \rho_0 u_h''(t) + \rho_1 u_h''(t - \delta t)) + e_1^h, \quad (5.4.7)$$

$u_h(t + \delta t) - u_h(t - \delta t)$

$$= 2\delta t(\rho_1 u_h'(t + \delta t) + \rho_0 u_h'(t) + \rho_1 u_h'(t - \delta t)) + e_2^h \quad (5.4.8)$$

and

$$u_h(t + \delta t) - u_h(t - \delta t) = 2\delta t u_h'(t) + e_3^h, \quad (5.4.9)$$

then

$$\|e_i^h\| \leq K(\delta t)^3 \left\{ \max_{\theta \in [t, t+\delta t]} \|u_h^{(3)}(\theta)\| + \max_{\theta \in [t, t+\delta t]} \|u_h^{(4)}(\theta)\| + \max_{\theta \in [t, t+\delta t]} \|u_h^{(5)}(\theta)\| \right\}$$

and

$$\|e_i^h\|_E \leq K(\delta t)^3 \left\{ \max_{\theta \in [t, t+\delta t]} \|u_h^{(3)}(\theta)\|_E + \max_{\theta \in [t, t+\delta t]} \|u_h^{(4)}(\theta)\|_E + \max_{\theta \in [t, t+\delta t]} \|u_h^{(5)}(\theta)\|_E \right\}.$$

### Proof

Consider (5.4.7) as an example: Use (5.4.4) in Proposition 5.4.2 for  $u_i$  and denote the remainder by  $\tilde{R}_{5i}(t)$ . Now, each term in (5.4.7) can be written as a linear combination, for example,

$$u_h(t + \delta t) = \sum_{i=1}^n u_i(t + \delta t)\phi_i.$$

Consequently we have (5.4.7), if we set  $e_1^h(t) = \sum_{i=1}^n \tilde{R}_{5i}(t)\phi_i$ .

It remains to estimate the error term  $e_1^h(t)$ , which is actually the sum of six error terms. Consider one of the terms: For any  $v \in S^h$

$$\begin{aligned} & \left| \left( \omega_1 \sum_{i=1}^n \int_t^{t+\delta t} (t + \delta t - \theta)^3 u_i^{(4)}(\theta) \phi_i d\theta, v \right) \right| \\ &= \left| \int_t^{t+\delta t} \omega_1 (t + \delta t - \theta)^3 (u_h^{(4)}(\theta), v) d\theta \right| \\ &\leq \int_t^{t+\delta t} |\omega_1| (t + \delta t - \theta)^3 \|u_h^{(4)}(\theta)\| \|v\| d\theta \\ &\leq \frac{1}{4} (\delta t)^4 |\omega_1| \|v\| \max_{\theta \in [t, t+\delta t]} \|u_h^{(4)}(\theta)\|. \end{aligned}$$

Hence there exists a constant  $K$ , which depends only on the weights  $\omega_1$  and  $\omega_2$  such that  $|(e_i^h(t), v)| \leq K(\delta t)^3 \|v\| \left\{ \max_{\theta \in [t, t+\delta t]} \|u_h^{(3)}(\theta)\| + \max_{\theta \in [t, t+\delta t]} \|u_h^{(4)}(\theta)\| + \max_{\theta \in [t, t+\delta t]} \|u_h^{(5)}(\theta)\| \right\}$ .

Note that the worst of the errors are of order  $(\delta t)^3$ . Since  $v$  is arbitrary, we have the desired result. The same procedure yields estimates in the energy norm.

□

**Lemma 5.4.1** *Suppose  $u_h$  is the solution of Problem AG. Let  $\rho_0$  and  $\rho_1$  be real numbers such that  $\rho_0 + 2\rho_1 = 1$ . If  $u^*(t, \delta t)$  is defined by*

$$\begin{aligned} & (u^*(t, \delta t) - 2u_h(t) + u_h(t - \delta t), v) + \frac{(\delta t)}{2} a(u^*(t, \delta t) - u_h(t - \delta t), v) \\ & + (\delta t)^2 b(\rho_1 u^*(t, \delta t) + \rho_0 u_h(t) + \rho_1 u_h(t - \delta t), v) \\ & = (\delta t)^2 (\rho_1 f(t + \delta t) + \rho_0 f(t) + \rho_1 f(t - \delta t), v)_0 \quad \text{for each } v \in S^h, \end{aligned} \tag{5.4.10}$$

then  $\|u_h(t + \delta t) - u^*(t, \delta t)\| \leq C_b (\delta t)^3 \|u_h\|_{5, \max}^E$ .

**Proof** Using Proposition 5.4.3 we have

$$\begin{aligned} & (u_h(t + \delta t) - 2u_h(t) + u_h(t - \delta t), v) + \frac{(\delta t)}{2} a(u_h(t + \delta t) - u_h(t - \delta t), v) \\ & = (\delta t)^2 (\rho_1 u_h''(t + \delta t) + \rho_0 u_h''(t) + \rho_1 u_h''(t - \delta t), v) + (e_1^h, v) \\ & + (\delta t)^2 a(\rho_1 u_h'(t + \delta t) + \rho_0 u_h'(t) + \rho_1 u_h'(t - \delta t), v) + \frac{(\delta t)}{2} a(e_2^h, v). \end{aligned}$$

Now use the fact that  $u_h$  is the solution of Problem AG to prove that  $u_h$  satisfies (5.4.10) with  $u(t + \delta t)$  in stead of  $u^*(t, \delta t)$  provided that the error terms  $(e_1^h, v)$  and  $\frac{(\delta t)}{2}a(e_2^h, v)$  are included.

Consequently,  $(u(t + \delta t) - u^*(t, \delta t), v) = (e_1^h, v) + \frac{(\delta t)}{2}a(e_2^h, v)$  for each  $v \in S^h$ . Replace  $v$  by  $u^*(t, \delta t) - u_h(t + \delta t)$  to obtain the estimate.  $\square$

Reconsider the semi discrete system in Section 4.4.

$$M\bar{u}''(t) + L\bar{u}'(t) + K\bar{u}(t) = \bar{f}(t) \quad (5.4.11)$$

$$\bar{u}(0) = \bar{\alpha}, \quad \bar{u}'(0) = \bar{\beta}.$$

To estimate the local errors for a finite difference scheme, we consider a one-to-one correspondence between  $S^h$  and  $\mathbb{R}_n$ .

**Definition 5.4.1** For  $u^h \in S^h$ , the vector  $\bar{u} = Qu^h$  has components  $u_i$  where  $u^h = \sum_{i=1}^n u_i \phi_i$ .

If we use the norm  $\|\bar{u}\|_M = (M\bar{u} \cdot \bar{u})^{\frac{1}{2}}$  for  $\mathbb{R}_n$ , then  $\|Qu^h\|_M = \|u^h\|$ .

In our next result use the fact that  $\bar{u}$  is a solution of (5.4.11) if and only if  $u_h$  is a solution of Problem AG.

**Corollary 5.4.1** If  $\bar{u}$  is a solution of the system of differential equations (5.4.11) and  $\bar{u}^*(t, \delta t)$  is defined by

$$\begin{aligned} & M [\bar{u}^*(t, \delta t) - 2\bar{u}(t) + \bar{u}(t - \delta t)] + \frac{(\delta t)}{2} L [\bar{u}^*(t, \delta t) - \bar{u}(t - \delta t)] \\ & + (\delta t)^2 K [\rho_1 \bar{u}^*(t, \delta t) + \rho_0 \bar{u}(t) + \rho_1 \bar{u}(t - \delta t)] \\ & = (\delta t)^2 [\rho_1 \bar{f}(t + \delta t) + \rho_0 \bar{f}(t) + \rho_1 \bar{f}(t - \delta t)], \end{aligned} \quad (5.4.12)$$

then  $\|\bar{u}^*(t, \delta t) - \bar{u}(t + \delta t)\|_M \leq C_b(\delta t)^3 \|u_h\|_{5, \max}^E$ .

**Proof** Consider the terms in (5.4.10). If  $\bar{v} = Qv^h$ , then

$$(u_h(t + \delta t), v^h) = M\bar{u}(t + \delta t) \cdot \bar{v}.$$

In this way we can associate each term in (5.4.12) with a corresponding term in (5.4.10). The result follows from the fact that  $\bar{u}^*(t, \delta t) = Qu^*(t, \delta t)$  and  $\bar{u}(t + \delta t) = Qu_h(t + \delta t)$ .  $\square$

2. If we assume that  $f$  is merely continuous and hence  $u_h$  twice continuously differentiable, we could still estimate the local errors but not obtain the same order. The results would be of the form: *Given  $\epsilon > 0$ , there exists a real number  $\Delta > 0$  such that the error will be less than  $\epsilon \delta t K$  for  $\delta t < \Delta$ . ( $K$  a constant depending on  $u_h$  and  $f$ .)*

## 5.4.2 Transformation

Due to symmetry considerations, it will be more convenient to consider a transformed system for stability analysis. Since  $M$  is symmetric and positive definite, there exists a symmetric positive definite matrix  $N$  such that  $N^2 = M$ . Set  $\bar{v}(t) = N\bar{u}(t)$ , then  $\bar{v}$  is a solution of the problem

$$\bar{v}'' + N^{-1}LN^{-1}\bar{v}' + N^{-1}KN^{-1}\bar{v} = N^{-1}\bar{f}$$

or

$$\bar{v}'' + \tilde{L}\bar{v}' + \tilde{K}\bar{v} = \tilde{g}. \quad (5.4.14)$$

where  $\tilde{L} = N^{-1}LN^{-1}$ ,  $\tilde{K} = N^{-1}KN^{-1}$  and  $\tilde{g} = N^{-1}\bar{f}$ .

The advantage of the transformation is that the matrix  $\tilde{K}$  is symmetric, and hence has orthogonal eigenvectors.

Let  $\bar{y} = N\bar{x}$ , then

$$K\bar{x} = \lambda M\bar{x}$$

if and only if

$$\tilde{K}\bar{y} = N^{-1}KN\bar{y} = \lambda\bar{y}.$$

The eigenvalues of  $\tilde{K}$  are the eigenvalues of the eigenvalue problem CG (See Sections 4.1 and 5.2.)

We use the norm  $\|\bar{x}\|_2 = (\bar{x} \cdot \bar{x})^{\frac{1}{2}}$ , and in the remaining part of this section  $\|\cdot\|$  will refer to  $\|\cdot\|_2$  unless stated otherwise.

**Corollary 5.4.3** *If  $\bar{v}$  is a solution of the system of differential equations (5.4.14), and  $\bar{v}^*(t, \delta t)$  is defined by*

$$\begin{aligned} & [\bar{v}^*(t, \delta t) - 2\bar{v}(t) + \bar{v}(t - \delta t)] + \frac{(\delta t)}{2} \tilde{L} [\bar{v}^*(t, \delta t) - \bar{v}(t - \delta t)] \\ & + (\delta t)^2 \tilde{K} [\rho_1 \bar{v}^*(t, \delta t) + \rho_0 \bar{v}(t) + \rho_1 \bar{v}(t - \delta t)] \\ = & (\delta t)^2 [\rho_1 \tilde{g}(t + \delta t) + \rho_0 \tilde{g}(t) + \rho_1 \tilde{g}(t - \delta t)], \end{aligned}$$

then  $\|\bar{v}^*(t, \delta t) - \bar{v}(t + \delta t)\| \leq C_b(\delta t)^3 (\|u_h\|_{5, \max}^E + (\delta t)^3 \|f\|_{3, \max})$ .

**Proof** Direct from Corollary 5.4.1, since  $\bar{v}^*(t, \delta t) = N\bar{u}^*(t, \delta t)$ . □

**Corollary 5.4.4** *If  $\bar{v}$  is a solution of the system of differential equations (5.4.14), and  $\bar{v}^{**}(t, \delta t)$  is defined by*

$$\begin{aligned} & 2[\bar{v}^{**}(t, \delta t) - \bar{v}(t)] + \frac{(\delta t)^2}{2} \tilde{K} [\bar{v}^{**}(t, \delta t) + \bar{v}(t)] \\ = & \frac{(\delta t)^2}{2} [\bar{g}(t + \delta t) + \bar{g}(t)] + 2\delta t \bar{v}'(t) - (\delta t)^2 \tilde{L} \bar{v}'(t) + \frac{(\delta t)^3}{2} \tilde{K} \bar{v}'(t) \\ & - \frac{(\delta t)^3}{2} \bar{g}'(t), \end{aligned}$$

then  $\|\bar{v}^{**}(t, \delta t) - \bar{v}(t + \delta t)\| \leq C_b(\delta t)^3 (\|u_h\|_{5, \max}^E + (\delta t)^3 \|f\|_{3, \max})$ .

**Proof** See Corollary 5.4.3. □

**Remark** The result remains true for  $t = 0$ .

### 5.4.3 Global error

We approximate the solution of (5.4.14) on the interval  $[0, \tau]$ . Let  $\delta t$  indicate the time step length, i.e.  $\delta t = \tau/N$ , and let  $\bar{w}_k$  denote the approximation for  $\bar{v}(t_k)$ .

We use the difference scheme (which corresponds to (5.4.12))

$$\begin{aligned} & (\bar{w}_{k+1} - 2\bar{w}_k + \bar{w}_{k-1}) + \frac{(\delta t)}{2} \tilde{L}(\bar{w}_{k+1} - \bar{w}_{k-1}) \\ & + (\delta t)^2 \tilde{K}(\rho_1 \bar{w}_{k+1} + \rho_0 \bar{w}_k + \rho_1 \bar{w}_{k-1}) = (\delta t)^2 (\rho_1 \bar{g}_{k+1} + \rho_0 \bar{g}_k + \rho_1 \bar{g}_{k-1}). \end{aligned} \tag{5.4.15}$$

The initial conditions for the system of differential equations are

$$\bar{v}(0) = N\bar{\alpha} \text{ and } \bar{v}'(0) = N\bar{\beta},$$

and the initial conditions of the finite difference system are

$$w_0 = N\bar{\alpha} \text{ and } (2\delta t)^{-1}(\bar{w}_1 - \bar{w}_{-1}) = N\bar{\beta}.$$

To estimate local errors the following scheme will also be used:

$$\begin{aligned}
 & 2(\bar{w}_{k+1} - 2\bar{w}_k + \bar{w}_{k-1}) + \delta t \tilde{L}(\bar{w}_{k+1} - \bar{w}_{k-1}) + \frac{(\delta t)^2}{2} \tilde{K}(\bar{w}_{k+1} + \bar{w}_k) \\
 = & \frac{(\delta t)^2}{2} (\bar{g}_{k+1} + \bar{g}_k) + 2\delta t \bar{v}'(t_k) + (\delta t)^2 L\bar{v}'(t_k) + \frac{(\delta t)^3}{2} \tilde{K}\bar{v}'(t_k) - \frac{(\delta t)^3}{2} g'(t_k).
 \end{aligned} \tag{5.4.16}$$

To estimate the global error  $\bar{w}_N - \bar{v}(\tau)$ , we introduce artificial numerical solutions  $\bar{w}_k^{(i)}$ . For each  $i$ ,  $\bar{w}_k^{(i)}$  satisfies (5.4.14) with  $\bar{w}_i^{(i)} = \bar{v}(t_i)$  and  $\bar{w}_{i-1}^{(i)} = \bar{w}_{i+1}^{(i)} - 2\delta t \bar{v}'(t_i)$ . Note that  $\bar{w}_k = \bar{w}_k^{(0)}$ .

For the global error we have

$$\|\bar{v}(\tau) - \bar{w}_N\| \leq \|\bar{v}(\tau) - \bar{w}_N^{(N-1)}\| + \|\bar{w}_N^{(N-1)} - \bar{w}_N^{(N-2)}\| + \dots + \|\bar{w}_N^{(1)} - \bar{w}_N\|. \tag{5.4.17}$$

(Note that the global error for the original system can be derived from this error.)

It is clearly necessary to estimate  $\|\bar{w}_N^{(i)} - \bar{w}_N^{(i-1)}\|$ . The next two subsections will be devoted to the estimation of the differences between “neighbouring numerical solutions”.

#### 5.4.4 Consistency

In this subsection we consider the differences  $\|\bar{w}_{i+1}^{(i)} - \bar{v}(t_{i+1})\|$  and  $\|\bar{w}_{i+2}^{(i)} - \bar{w}_{i+2}^{(i+1)}\|$ . For simplicity we denote  $\bar{v}(t_i)$  by  $\bar{v}_i$ . The first lemma deals with the “starting” error.

##### Lemma 5.4.3

$$\|\bar{v}_{i+1} - \bar{w}_{i+1}^{(i)}\| \leq C_b (\delta t)^3 (\|u_h\|_{5,\max}^E + (\delta t)^3 \|f\|_{3,\max}).$$

**Proof** This is a direct consequence of Corollary 5.4.4. □

Next we have the error at the second step.

**Lemma 5.4.4**

$$\|\bar{w}_{i+2}^{(i)} - \bar{v}_{i+2}\| \leq C_b(\delta t)^3 (\|u_h\|_{5,\max}^E + (\delta t)^3 \|f\|_{3,\max}).$$

**Proof** Combine the results of Corollary 5.4.3 and Lemma 5.4.3.  $\square$

Lemma 5.4.3 provide an estimate for the difference  $\bar{w}_i^{(i)} - \bar{w}_{i+1}^{(i)}$ . The following result provide an estimate for the difference at the second step.

**Corollary 5.4.5**

$$\|w_{i+2}^{(i)} - w_{i+2}^{(i+1)}\| \leq C_b(\delta t)^3 (\|u_h\|_{5,\max}^E + \|f\|_{3,\max}).$$

**Proof** Use Lemmas 5.4.3 and 5.4.4.

$$\|w_{i+2}^{(i)} - w_{i+2}^{(i+1)}\| \leq \|\bar{w}_{i+2}^{(i)} - \bar{v}_{i+2}\| + \|\bar{v}_{i+1} - \bar{w}_{i+2}^{(i+1)}\|.$$

$\square$

### 5.4.5 Stability

For the stability analysis we introduce the following matrices:

$$\begin{aligned} A &= I + \frac{(\delta t)}{2} \tilde{L} + \rho_1(\delta t)^2 \tilde{K}, \\ B &= -2I + \rho_0(\delta t)^2 \tilde{K}, \\ C &= I - \frac{(\delta t)}{2} \tilde{L} + \rho_1(\delta t)^2 \tilde{K}. \end{aligned}$$

The system (5.4.15 ) is now

$$A\bar{w}_{k+1} + B\bar{w}_k + C\bar{w}_{k-1} = (\delta t)^2(\rho_1\bar{g}_{k+1} + \rho_0\bar{g}_k + \rho_1\bar{g}_{k-1}) \quad (5.4.18)$$

As mentioned at the end of Subsection 5.4.3, we need to estimate the difference  $\bar{w}_N^{(i)} - \bar{w}_N^{(i+1)}$  for each  $i$ . Since both  $\bar{w}_j^{(i)}$  and  $\bar{w}_j^{(i+1)}$  satisfy the system (5.4.18) it follows that the error  $\bar{e}_j = \bar{w}_j^{(i)} - \bar{w}_j^{(i+1)}$  must satisfy

$$A\bar{e}_{j+1} + B\bar{e}_j + C\bar{e}_{j-1} = \bar{0}, \quad (5.4.19)$$

with the starting values, the local errors  $\bar{e}_{i+1}$  and  $\bar{e}_{i+2}$ , already estimated.

For the case  $L = \mu K$ , we derive the eigenvalues of  $A$ ,  $B$  and  $C$ . If  $\tilde{K}\bar{y} = \lambda\bar{y}$ , then

$$\begin{aligned} A\bar{y} &= \bar{y} + \mu \frac{(\delta t)}{2} \tilde{K}\bar{y} + \rho_1(\delta t)^2 \tilde{K}\bar{y} = \left(1 + \frac{\delta t}{2}(\mu\lambda) + \rho_1(\delta t)^2\lambda\right) \bar{y}, \\ B\bar{y} &= -2\bar{y} + \rho_0(\delta t)^2 \tilde{K}\bar{y} = (-2 + \rho_0(\delta t)^2\lambda) \bar{y}, \\ C\bar{y} &= \left(1 - \frac{\delta t}{2}(\mu\lambda) + \rho_1(\delta t)^2\lambda\right) \bar{y}. \end{aligned}$$

It is now possible to solve (5.4.19). Let  $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n$  denote the normalized eigenvectors of  $\tilde{K}$  and suppose  $\bar{e}_{i+1} = \sum_{i=1}^n \eta_i \bar{y}_i$  and  $\bar{e}_{i+2} = \sum_{i=1}^n \xi_i \bar{y}_i$ .

Since the eigenvectors are orthogonal, it is sufficient to solve difference equations of the form

$$\alpha_i r_{k+1} + \beta_i r_k + \gamma_i r_{k-1} = 0 \text{ for } i = 1, 2, \dots, n,$$

where  $\alpha_i, \beta_i$  and  $\gamma_i$  denote the eigenvalues of the matrices  $A, B$  and  $C$  respectively.

The following result can be obtained by elementary calculations. Note that we do not use the subscripts for the coefficients  $\alpha, \beta$  and  $\gamma$ . We take  $r_1 = \xi$  and  $r_0 = \eta$ .

### Solution of the difference equation

**Case 1**  $\beta^2 < 4\alpha\gamma$ .

Note that in this case  $\gamma > 0$ . The solution is of the form

$$r_k = p^k (A \cos \omega k + B \sin \omega k),$$

where  $p = \sqrt{\gamma/\alpha}$ ,  $\cos \omega = -\beta/(2\sqrt{\alpha\gamma})$ ,  $A = \eta$  and  $B = (\xi - p\eta \cos \omega)/(\sin \omega)$ .

**Case 2**  $\beta^2 = 4\alpha\gamma$ .

The solution is of the form

$$r_k = \eta(1 - k)r^k + \xi k r^{k-1}, \text{ with } r = -\frac{\beta}{2\alpha}.$$

**Case 3**  $\beta^2 > 4\alpha\gamma$ .



The solution is of the form

$$r_k = Ar_1^k + Br_2^k,$$

with  $r_1$  and  $r_2$  the real roots of the equation  $\alpha r^2 + \beta r + \gamma = 0$ . The constants are  $A = (\xi - \eta r_2)/(r_1 - r_2)$  and  $B = (\xi - \eta r_1)/(r_2 - r_1)$ .

We now prove the stability result. Bear in mind that  $\lambda_k \rightarrow \infty$  as  $k \rightarrow \infty$ .

**Lemma 5.4.5 Stability**

If  $\rho_0 \leq 2\rho_1$ , then there exists a constant  $K$ —independent of the dimension of  $S^h$ —such that

$$\|w_N^{(i)} - w_N^{(i+1)}\| \leq K \left( \|v_{i+1} - w_{i+1}^{(i)}\| + \|w_{i+2}^{(i)} - w_{i+2}^{(i+1)}\| \right).$$

**Proof** For the eigenvalues

$$\begin{aligned} \alpha &= 1 + \rho_1(\delta t)^2 \lambda + \frac{\delta t}{2} \lambda \mu, \\ \beta &= -2 + \rho_0(\delta t)^2 \lambda, \\ \gamma &= 1 + \rho_1(\delta t)^2 \lambda - \frac{\delta t}{2} \lambda \mu \end{aligned}$$

of  $A$ ,  $B$  and  $C$ , we get

$$\begin{aligned} \beta^2 - 4\alpha\gamma &= -4\rho_0(\delta t)^2 \lambda + \rho_0^2(\delta t)^4 \lambda^2 - 8\rho_1(\delta t)^2 \lambda - 4\rho_1^2(\delta t)^4 \lambda^2 + (\delta t)^2 \mu^2 \lambda^2 \\ &= \lambda(\delta t)^2 [\mu^2 \lambda - 4 + (\rho_0^2 - 4\rho_1^2)(\delta t)^2 \lambda]. \end{aligned}$$

Consider the different cases:

**Case 1** If  $\beta^2 < 4\alpha\gamma$  then  $r_k$  is bounded if  $\gamma \leq \alpha$ .

This case,  $\beta^2 - 4\alpha\gamma < 0$ , is possible only for a finite number of small eigenvalues and only if  $\rho_0 \geq 2\rho_1$ . Since  $\gamma/\alpha < 1$ , the corresponding modes will not cause error growth.

**Case 2** If  $\beta^2 = 4\alpha\gamma$  then  $r_k$  is bounded if  $|\beta/\alpha| < 2$ .

If  $\rho_0(\delta t)^2 \lambda < 2$ , we have  $|\beta|/|\alpha| < \frac{2}{1}$ .

If  $\rho_0(\delta t)^2 \lambda > 2$ , we have  $|\beta|/|\alpha| < \frac{\rho_0(\delta t)^2 \lambda}{\rho_1(\delta t)^2 \lambda} = \frac{\rho_0}{\rho_1} \leq 2$  if  $\rho_0 \leq 2\rho_1$ .

**Case 3** If  $\beta^2 > 4\alpha\gamma$  then  $r_k$  is bounded if both roots of  $\alpha r^2 + \beta r + \gamma = 0$  are less than one in absolute value. Let

$$\Delta = \beta^2 - 4\alpha\gamma \leq \lambda^2(\delta t)^2 [d^2 + (\rho_0^2 - 4\rho_1^2)(\delta t)^2]$$

and let  $r_{\max}$  denote the absolute value of the root largest in absolute value.

$$\begin{aligned} r_{\max} &= \frac{\beta + \sqrt{\Delta}}{2\alpha} \\ &\leq \frac{\rho_0(\delta t)^2\lambda + \lambda\delta t\sqrt{\mu^2 + (\rho_0^2 - 4\rho_1^2)(\delta t)^2}}{2\rho_1(\delta t)^2\lambda + \delta t\lambda\mu} \\ &= \frac{\rho_0\delta t + \sqrt{\mu^2 + (\rho_0^2 - 4\rho_1^2)(\delta t)^2}}{2\rho_1\delta t + \mu} \\ &\leq \frac{\rho_0\delta t + \mu}{2\rho_1\delta t + \mu} \quad (\text{if } \rho_0 \leq 2\rho_1) \\ &\leq 1 \quad (\text{if } \rho_0 \leq 2\rho_1). \end{aligned}$$

□

**Remark** If damping is excluded, the difference system is considered to be unconditionally stable for  $\rho_1 = \frac{1}{4}$  and  $\rho_0 = \frac{1}{2}$ , see [RM] or [Zi]. However, the bound may depend on the eigenvalues.

$$\begin{aligned} \cos \omega &= \frac{-\beta}{2\sqrt{\alpha\gamma}} = \frac{-\rho_0(\delta t)^2\lambda + 2}{2(1 + \rho_1(\delta t)^2\lambda)} \\ &= \frac{-\rho_0(\delta t)^2 + 2\lambda^{-1}}{2\rho_1(\delta t)^2 + 2\lambda^{-1}} \rightarrow \frac{-\rho_0}{2\rho_1} \text{ as } \lambda \rightarrow \infty. \end{aligned}$$

Consequently,  $\sin \omega \rightarrow 0$  as  $\lambda \rightarrow \infty$  and  $\sin \omega$  is present in the numerator of a constant.

### Remarks

1. Exactly the same results hold if we assume that rotary inertia can be ignored and we have only viscous damping. In this case  $L = kM_0$  and  $M = M_0$ .
2. The eigenvalues of  $K\bar{x} = \lambda M_0\bar{x}$  are much larger than the eigenvalues of  $K\bar{x} = \lambda M\bar{x}$  (with rotary inertia).
3. Rotary inertia and Kelvin-Voigt damping both enhance stability.

### 5.4.6 Convergence

**Lemma 5.4.6** *Global error*

$$\|\bar{e}_N\| \leq C_b N (\delta t)^3 (\|u_h\|_{5,\max}^E + \|f\|_{3,\max})$$

(where  $N$  is the number of steps).

**Proof**

$$\begin{aligned} \|\bar{e}_N\| &\leq \|\bar{v}(\tau) - \bar{w}_N^{N-1}\| + \dots + \|w_N^{(1)} - \bar{w}_N\| \\ &\leq KN \max_i \{\|\bar{v}_{i+1} - \bar{w}_{i+1}^{(i)}\| + \|\bar{v}_{i+1} - \bar{w}_{i+2}^{(i)}\|\}. \end{aligned}$$

Now use Lemmas 5.4.3 and 5.4.4.

□

To the sequence of finite difference vectors  $\bar{v}_k$ , correspond a sequence of approximations for  $u^h$ :  $u_h^{(k)} = (QJ)^{-1} \bar{v}_k \in S^h$ .

**Theorem 5.4.1** *If  $u_h$  is the solution of Problem AG, then*

$$\|u_h(\tau) - u_N^h\| \leq C_b (\delta t)^2 (\|u_h\|_{5,\max}^E + \|f\|_{3,\max}),$$

where  $u_N^h = Q^{-1} \bar{u}_N$ .

**Proof** Since  $\bar{e}_N = QN(u_n(\tau) - u_n^N)$ , we have

$$\|u_h(\tau) - u_h^N\| = \|\bar{e}_N\|,$$

Now use Lemma 5.4.6.

□

**Remark** Error estimates for the fully discrete system is obtained by combining Theorem 5.4.1 with the results of Section 5.3. Note that the error estimates are with respect to the inertia norm.

## Chapter 6

# Application. Damaged beam

### 6.1 Introduction

We consider Problem 1 (from Section 2.3). This model for a damaged beam was proposed in [VV]. See also [JVRV].

The detection of damage in structures or materials is clearly of great importance. Ideally it should be possible to infer the location and extent of damage from indirect measurements or signals. To facilitate such deduction, a mathematical model of the object or structure is necessary. See [VV] for details and numerous other references.

Viljoen *et al.* [VV] use changes in the natural frequencies of the beam to locate and quantify the damage. The natural angular frequencies for the damaged beam are calculated from the characteristic equation obtained from the associated eigenvalue problem. As is well-known, only the first few natural angular frequencies and modes are usually calculated with this method, because of computational difficulties with the hyperbolic functions. Due to this limitation, the need arises for a numerical method to simulate the dynamical behaviour of the beam.

In a joint paper, [ZVV], we developed a finite element method (FEM) to approximate the solution of the model problem for arbitrary initial conditions. (Ironically we also found it possible to calculate eigenvalues and eigenfunctions more accurately with the FEM.) It was necessary to adapt standard procedures to deal with the discontinuity in the derivative that arises as a

result of the elastic joint. We made the assumption that damping would not influence the solution significantly on a small time scale. We now investigate the validity of this assumption, and also deem it prudent to include the effect of rotary inertia.

In the paper [ZVV], only Hermite piecewise cubics were used as basis functions. In this investigation we also demonstrate the effectiveness of Hermite piecewise quintics.

From Section 3.1 we have the variational formulation. The Galerkin approximations for the eigenvalue and initial value problems are given in Section 4.1. (We do not consider the equilibrium problem.) In Section 4.2 we showed how the standard basis functions are adapted to deal with the discontinuity in the derivative.

In Section 6.2 we compute the natural angular frequencies and modes of vibration from the characteristic equation for comparison purposes. The computation of the matrices is discussed in Section 6.3.

In Sections 6.4 and 6.5 numerical results are presented that demonstrate not only the effect of damage on the motion of a beam but also the effect of damping and rotary inertia. We also investigate the use of Hermite piecewise quintics as basis functions instead of Hermite piecewise cubics.

## 6.2 Natural frequencies and modes of vibration

One way to calculate the natural angular frequencies and modes of vibration for the damaged beam is to apply the method of separation of variables directly to Problem 1 (from Section 2.3).

For the case  $r = 0$  (without rotary inertia), we have the following eigenvalue problem:

$$\begin{aligned} w^{(4)} - \lambda w &= 0, \quad 0 < x < 1, \quad x \neq \alpha, \\ w(0) &= w'(0) = w''(1) = w'''(1) = 0, \\ w(\alpha^+) &= w(\alpha^-), \\ w''(\alpha^+) &= w''(\alpha^-), \\ w'''(\alpha^+) &= w'''(\alpha^-), \\ w''(\alpha) &= \frac{1}{\delta}(w'(\alpha^+) - w'(\alpha^-)). \end{aligned}$$

For this eigenvalue problem it is possible to find so called exact solutions. It is convenient to introduce the positive real number  $\nu$ , with  $\lambda = \nu^4$ . Consequently  $\nu^2 = \sqrt{\lambda}$  is a natural angular frequency. Analogous to the case of the undamaged beam, the corresponding mode is of the form

$$w(x) = \begin{cases} A \sin(\nu x) - A \sinh(\nu x) + B \cos(\nu x) - B \cosh(\nu x) & \text{for } 0 < x < \alpha, \\ (C + A) \sin(\nu x) + (D - A) \sinh(\nu x) \\ \quad + (E + B) \cos(\nu x) + (F - B) \cosh(\nu x) & \text{for } \alpha < x < 1. \end{cases}$$

Note that the boundary conditions at  $x = 0$  have already been taken into account.

From the continuity conditions and the jump condition at  $x = \alpha$ , the constants  $C$ ,  $D$ ,  $E$ , and  $F$  can be expressed in terms of  $A$  and  $B$ . Finally, from the two boundary conditions at  $x = 1$ , the characteristic equation for  $\nu$  can be constructed from

$$\det \begin{bmatrix} d_1 & d_2 \\ d_3 & d_4 \end{bmatrix} = 0 \tag{6.2.1}$$

where

$$\begin{aligned}
 d_1 &= -(\sin \nu + \sinh \nu) + \frac{\delta \nu}{2}(\sin \nu \alpha + \sinh \nu \alpha) \times \\
 &\quad (\sin \nu \cos \nu \alpha - \sinh \nu \cosh \nu \alpha - \cos \nu \sin \nu \alpha + \cosh \nu \sinh \nu \alpha), \\
 d_2 &= -(\cos \nu + \cosh \nu) + \frac{\delta \nu}{2}(\cos \nu \alpha + \cosh \nu \alpha) \times \\
 &\quad (\sin \nu \cos \nu \alpha - \sinh \nu \cosh \nu \alpha - \cos \nu \sin \nu \alpha + \cosh \nu \sinh \nu \alpha), \\
 d_3 &= -(\cos \nu + \cosh \nu) + \frac{\delta \nu}{2}(\sin \nu \alpha + \sinh \nu \alpha) \times \\
 &\quad (\cos \nu \cos \nu \alpha - \cosh \nu \cosh \nu \alpha + \sin \nu \sin \nu \alpha + \sinh \nu \sinh \nu \alpha), \\
 d_4 &= (\sin \nu - \sinh \nu) + \frac{\delta \nu}{2}(\cos \nu \alpha + \cosh \nu \alpha) \times \\
 &\quad (\cos \nu \cos \nu \alpha - \cosh \nu \cosh \nu \alpha + \sin \nu \sin \nu \alpha + \sinh \nu \sinh \nu \alpha).
 \end{aligned}$$

Solving equation (6.2.1) numerically using the Newton-Raphson method, yields the natural angular frequencies for the damaged beam. For each natural angular frequency a corresponding mode can then be obtained. As is expected, only the first few natural angular frequencies and modes could be calculated, as it is difficult to handle the hyperbolic functions numerically for large values of  $\nu$ .

Numerical results obtained using the finite element method—using cubics as well as quintics as basis functions— are given in Section 6.4.

### 6.3 Computation of Matrices

The matrices  $K$ ,  $L$  and  $M$  are defined in Section 4.1 in terms of the bilinear forms defined in Section 3.1. The computation of the matrices is complicated by the interface conditions which results in non-standard basis elements. In this section we give an indication of how we went about in computing these matrices. The first step is to reorder the basis elements constructed in Subsection 4.2.2.

Consider the matrix  $M_0$ :

$$[M_0]_{ij} = (\tilde{\phi}_i, \tilde{\phi}_j) = \int_0^\alpha \phi_{i1} \phi_{j1} + \int_\alpha^1 \phi_{i2} \phi_{j2}.$$

Note that  $\tilde{\phi}_i = \langle 0, \phi_i \rangle$  or  $\langle \phi_i, 0 \rangle$  except when we are dealing with a basis element associated with the node  $x_p = \alpha$ , the location of the damage. In general then, the entries will be those of the standard mass matrix for an undamaged beam. Now suppose one of the basis elements are associated with  $x_p$ :

If  $\tilde{\phi}_i = \tilde{\phi}_p^{(0)}$ , then

$$[M_0]_{ij} = \int_0^\alpha \phi_{i1} \phi_{j1} + \int_\alpha^1 \phi_{i2} \phi_{j2} = \int_0^1 \phi_i \phi_j.$$

Again the result will be the same as in the standard case. The same for  $\tilde{\phi}_p^{(2)}$ .

On the other hand, suppose  $\tilde{\phi}_i = \tilde{\phi}_{pL}^{(1)}$ , then

$$[M_0]_{ij} = \int_0^\alpha \phi_{j1} \phi_{i1} + 0,$$

which is not the same as for an undamaged beam. Similarly for  $\tilde{\phi}_i = \tilde{\phi}_{pR}^{(1)}$ . Thus the standard matrix has to be modified for the damaged beam.

We have the same situation for the matrix  $M_r$  where we define

$$[M_r]_{ij} = (\tilde{\phi}'_i, \tilde{\phi}'_j) = \int_0^\alpha \phi'_{i1} \phi'_{j1} + \int_\alpha^1 \phi'_{i2} \phi'_{j2}.$$

There is an additional complication for the  $K$ -matrix:

$$K_{ij} = b(\tilde{\phi}_i, \tilde{\phi}_j) = \int_0^\alpha \phi''_{i1} \phi''_{j1} + \int_\alpha^1 \phi''_{i2} \phi''_{j2} + \frac{1}{\delta} (\phi'_{j2}(\alpha) - \phi'_{j1}(\alpha)) (\phi'_{i2}(\alpha) - \phi'_{i1}(\alpha)).$$



Only four entries in the standard  $K$ -matrix will change due to the additional term,  $(u'_2(\alpha) - u'_1(\alpha))(v'_2(\alpha) - v'_1(\alpha))/\delta$ , in the bilinear form  $b$ .

For greater clarity we will explain the procedure in another way. In the discussion that follows, we refer to  $\tilde{\phi}_i^{(k)}$  as a Type  $k$  basis function.

In modifying the matrices for an undamaged beam to the matrices for a damaged beam, we have to keep in mind that the Type 1 basis function associated with  $x_p = \alpha$ , has changed. By replacing the row and column associated with the Type 1 basis function at  $x_p$ , in the matrix of the undamaged beam, by two rows and columns respectively, provision is made for the modified basis function. The values in the matrix in these two rows and columns have to be modified accordingly. For the  $K$ -matrix one must also keep the additional term in mind.

Having computed  $M_0$ ,  $M_r$  and  $K$  we are done since  $L = \mu K + kM_0$  and  $M = M_0 + M_r$ .

## 6.4 Numerical results. Eigenvalue problem

Cubics as basis functions, are usually sufficiently accurate in solving one-dimensional vibration problems with the finite element method. In a joint paper [ZVV] we discussed the use of cubics as basis functions for the damaged beam.

In this section of the thesis we also consider numerical convergence of the eigenvalues. The order of convergence that is suggested by the numerical results is also compared to the order obtained from the theory. Additionally, quintics are considered as basis functions. The main reason for this is that cubics are not compatible with reduced quintics in plate beam models. We also investigate the effect of rotary inertia.

### 6.4.1 Cubics

Natural angular frequencies and modes for the vibration problem are calculated by solving the eigenvalue problem with the FEM. We developed the code to construct the relevant matrices in Matlab and use standard Matlab subroutines to calculate the eigenvalues and eigenvectors of the generalised eigenvalue problem.

It is possible to compare only the first few FEM eigenvalues to the so called exact eigenvalues calculated from the characteristic equation. Thereafter the exact values can not be computed accurately and the FEM is used to calculate the eigenvalues.

In Table 6.1 we list values for the eigenvalues obtained from the characteristic equation (see Section 6.2) and values obtained by the FEM using cubics as basis functions with 20, 40, 80 and 160 subintervals respectively. This give approximations for respectively the first 40, 80, 160 and 320 eigenvalues.

$i$	$\lambda_i$	$\lambda_i^{(20)}$	$\lambda_i^{(40)}$	$\lambda_i^{(80)}$	$\lambda_i^{(160)}$
1	11.81469	11.81469	11.81469	11.81469	11.81469
2	406.01614	406.01757	406.01623	406.01615	406.01615
3	3806.05283	3806.17742	3806.06067	3806.05332	3806.05287
4	12544.12940	12545.47137	12544.21439	12544.13473	12544.12972
5	39943.82322	39957.35763	39944.68387	39943.87724	39943.82661

Table 6.1: Eigenvalues from the characteristic equation as well as FEM eigenvalues using cubics as basis functions with  $\delta = 0.1$  and  $\alpha = 0.5$ .

Throughout this section  $n$  denote the number of subintervals. (All of equal length.)

To investigate the convergence of the FEM eigenvalues, we calculate the relative difference between FEM approximations, that is  $(\lambda^{(2n)} - \lambda^{(n)})/\lambda^{(2n)}$ . These differences are calculated and listed in Table 6.2 for  $n = 20, 40, 80$  and 160 subintervals respectively.

$i$	$ \lambda_i^{(2n)} - \lambda_i^{(n)} /\lambda_i^{(2n)}$			
	$n = 20$	$n = 40$	$n = 80$	$n = 160$
6	$6.1 \times 10^{-4}$	$3.9 \times 10^{-5}$	$2.5 \times 10^{-6}$	$1.6 \times 10^{-7}$
12	$1.1 \times 10^{-2}$	$7.6 \times 10^{-4}$	$4.9 \times 10^{-5}$	$3.1 \times 10^{-6}$
24	$2.2 \times 10^{-1}$	$1.2 \times 10^{-2}$	$8.6 \times 10^{-4}$	$5.5 \times 10^{-5}$
48	—	$2.2 \times 10^{-1}$	$1.3 \times 10^{-2}$	$9.2 \times 10^{-4}$

Table 6.2: Relative differences for FEM eigenvalues using cubics as basis functions.

The tendency of the relative difference to decrease (by roughly a factor 10) each time that the number of subintervals is doubled, is empirical verification that there is convergence of the FEM eigenvalues. We found that the eigenvalues computed from the characteristic equation were less dependable.

It is necessary to determine a relationship between the number of FEM eigenvalues that is sufficiently accurate (criterion to be specified) and the number of subintervals used.

A relative difference strictly less than  $10^{-3}$  is considered sufficiently accurate for our purpose. Using this as criterion, we find that approximately a seventh of the  $2n$  eigenvalues calculated using  $n$  subintervals, yields a relative difference,  $(\lambda^{(2n)} - \lambda^{(n)})/\lambda^{(2n)}$ , strictly less than  $10^{-3}$ , see Table 6.2.

The relative difference between the FEM eigenvalues with 160 and 320 subintervals is an indication of the relative error between the exact eigenvalue and the FEM eigenvalue using 320 subintervals.

Since we use  $(\lambda^{(320)} - \lambda^{(160)})/\lambda^{(320)}$  as measure of the relative error,  $(\lambda - \lambda^{(320)})/\lambda$ , we conclude that the first 90 eigenvalues obtained using 320 subintervals yield a relative error that is sufficiently accurate.

An indication of the order of convergence of the FEM eigenvalues can be obtained from the ratio of two successive differences

$$|\lambda_i^{(2n)} - \lambda_i^{(n)}|/|\lambda_i^{(4n)} - \lambda_i^{(2n)}|.$$

Typical results are listed in Table 6.3.

$i$	$ \lambda_i^{(2n)} - \lambda_i^{(n)} / \lambda_i^{(4n)} - \lambda_i^{(2n)} $		
	$n = 20$	$n = 40$	$n = 80$
1	10.25	0.03	0.14
3	15.88	16.41	6.57
6	15.53	15.88	15.72
12	14.23	15.57	15.90
24	18.29	14.47	15.62

Table 6.3: Relationship between successive relative differences with cubics as basis functions.

These relative differences decrease by roughly a factor 16 if the number of subintervals is doubled. From this it would appear that the convergence is of order  $h^4$  which matches the theory, Section 5.2.

It is observed that those differences not yielding a factor 16 typically occur in the right top part as well as the left bottom part of Table 6.3. These deviations are illustrated by the first, third and 24th eigenvalues:

Firstly, the accuracy of an approximation can decrease if the number of subintervals is increased. This is due to an increase in the roundoff error and has significant effects in situations where the errors are already small. For example, FEM approximations for the first eigenvalue yield

$$(\lambda_1^{(40)} - \lambda_1^{(20)}) = -1.1 \times 10^{-13} \quad \text{while} \quad (\lambda_1^{(80)} - \lambda_1^{(40)}) = 3.8 \times 10^{-12}.$$

From the theory, Section 5.2, we know that the FEM approximations of an eigenvalue will decrease if the number of subintervals is increased. This can

be used to detect cases where the effect of the roundoff error is greater than the advantageous effect of an increase in the number of subintervals used.

Rounding error also explain the decrease in the ratios for the third eigenvalue from roughly a factor 16 to 6.5. This situation differ from the first eigenvalue in that the decrease (improvement) in the relative difference was just partially cancelled by the increase in the roundoff error.

Secondly, as we have showed previously, there is a relationship between the number of FEM eigenvalues that can be calculated sufficiently accurately and the number of subintervals used. The 24th eigenvalue is such an example. The effect of the poor approximation of  $\lambda_{24}^{(20)}$ , is seen in Table 6.3 in that  $18.29 > 14.47$ . This is expected as only the first six eigenvalues obtained, using 20 subintervals, yield relative errors less than  $10^{-3}$ .

### 6.4.2 Quintics

We now consider quintics as basis functions, and compare the results to the case where we used cubics.

In Table 6.4 we list values for the eigenvalues obtained from the characteristic equation and values obtained by the FEM using quintics as basis functions with 2, 4, 8 and 16 subintervals respectively. This gives approximations for respectively the first 6, 12, 24 and 48 eigenvalues.

$i$	$\lambda_i$	$\lambda_i^{(2)}$	$\lambda_i^{(4)}$	$\lambda_i^{(8)}$	$\lambda_i^{(16)}$
1	11.81469	11.81469	11.81469	11.81469	11.81469
2	406.01614	406.01954	406.01618	406.01614	406.01614
3	3806.05283	3822.51900	3806.09344	3806.05297	3806.05283
4	12544.12940	12844.88875	12544.53719	12544.13441	12544.12941
5	39943.82322	41569.18041	40042.72518	39944.02589	39943.82387

Table 6.4: *Eigenvalues from the characteristic equation as well as FEM eigenvalues using quintics as basis functions with  $\delta = 0.1$  and  $\alpha = 0.5$ .*

If these values are compared to those in Table 6.1, it seems as if the same accuracy can be obtained, using quintics as basis functions, with less subintervals, than in the case where cubics were used as basis functions. For example, the fifth FEM eigenvalue using quintics as basis functions with

16 subintervals, already yields a better approximation than using cubics with 80 subintervals.

As in the case with cubics as basis functions, we investigate the convergence of the FEM eigenvalues by considering relative differences,  $(\lambda^{(2n)} - \lambda^{(n)})/\lambda^{(2n)}$ . These values are listed in Table 6.5 for 2, 4, 8 and 16 subintervals respectively.

$ \lambda_i^{(2n)} - \lambda_i^{(n)} /\lambda_i^{(2n)}$				
$i$	$n = 2$	$n = 4$	$n = 8$	$n = 16$
1	$2.1 \times 10^{-8}$	$7.9 \times 10^{-11}$	$1.2 \times 10^{-11}$	$1.7 \times 10^{-11}$
2	$8.3 \times 10^{-6}$	$9.5 \times 10^{-8}$	$3.4 \times 10^{-10}$	$1.3 \times 10^{-11}$
4	$2.4 \times 10^{-2}$	$3.2 \times 10^{-5}$	$4.0 \times 10^{-7}$	$1.4 \times 10^{-9}$
8	—	$3.9 \times 10^{-2}$	$7.2 \times 10^{-5}$	$8.2 \times 10^{-7}$

Table 6.5: *Relative differences for FEM eigenvalues using quintics as basis functions.*

The numerical results suggest convergence of the FEM eigenvalues since the relative error decreases (by roughly a factor 100) each time that the number of subintervals is doubled, see Table 6.5.

For approximately a third of the  $3n$  eigenvalues computed, using  $n$  subintervals, the relative difference  $(\lambda^{(2n)} - \lambda^{(n)})/\lambda^{(2n)}$  is strictly less than  $10^{-3}$ .

As with the cubics, we now consider the ratio of two successive differences

$$|\lambda_i^{(2n)} - \lambda_i^{(n)}|/|\lambda_i^{(4n)} - \lambda_i^{(2n)}|$$

to get an idea of the order of convergence. Typical results are listed in Table 6.6.

As was the case in Table 6.3, the values in the top right of Table 6.6 exhibit effect of roundoff error and the values in the bottom left the result of eigenvalues not calculated sufficiently accurately. From this it would appear that the order of convergence is  $h^8$  which matches the theory, Section 5.2.

To compare the accuracy of the FEM eigenvalues using quintics as basis functions to the case using cubics as basis functions, we choose the number of subintervals in each of the cases such that the sizes of the matrices in the two cases are equal. For example, using 30 subintervals for cubics yield  $61 \times 61$  matrices and 20 subintervals for quintics  $62 \times 62$  matrices. We then

$i$	$ \lambda_i^{(2n)} - \lambda_i^{(n)}  /  \lambda_i^{(4n)} - \lambda_i^{(2n)} $		
	$n = 2$	$n = 4$	$n = 8$
1	265.60	6.75	0.69
2	87.21	278.22	25.45
3	405.90	278.67	319.70
4	745.69	80.59	291.89
5	15.43	488.57	311.65
6	16.24	330.71	307.27
7	304.98	166.95	286.96
8	302.37	540.94	87.70

Table 6.6: Relationship between successive relative differences with quintics basis functions.

compare the eigenvalues calculated in the two cases with the eigenvalues computed using cubics with 320 subintervals. (We use the first 90 FEM eigenvalues using cubics as basis functions with 320 subintervals as the FEM approximation to the first 90 exact eigenvalues.)

Our numerical experiments indicate that using quintics with  $n$  subintervals, yield at least double the number of eigenvalues to the prescribed accuracy (relative error strictly less than  $10^{-3}$ ) than when cubics are used with  $3n/2$  subintervals. In Table 6.7 we give an example of results obtained.

In Table 6.7 we use the following notation:

- Let  $\lambda_i$  denote the  $i$ th FEM eigenvalue that we use as approximation for the exact eigenvalue. (In this case those FEM eigenvalues obtained using cubics as basis functions with 320 subintervals.)
- To distinguish between the FEM eigenvalues computed using quintics and cubics as basis functions, we denote the  $i$ th FEM eigenvalue using cubics with 30 subintervals by  $\lambda_i^{(c)}$  and using quintics with 20 subintervals by  $\lambda_i^{(q)}$ .

Note that the FEM approximations for the first eigenvalue are identical in both cases.

From Table 6.7 we see that using cubics, the first 9 eigenvalues (approximately a seventh of the number of eigenvalues calculated,  $61/7 \approx 8.7$ ) have

$i$	$(\lambda_i - \lambda_i^{(c)})/\lambda_i$	$(\lambda_i - \lambda_i^{(q)})/\lambda_i$
1	$2.6 \times 10^{-6}$	$2.6 \times 10^{-6}$
5	$6.8 \times 10^{-5}$	$9.6 \times 10^{-9}$
9	$8.5 \times 10^{-4}$	$4.3 \times 10^{-7}$
10	$1.2 \times 10^{-3}$	$8.6 \times 10^{-7}$
15	$6.8 \times 10^{-3}$	$4.6 \times 10^{-5}$
20	$1.9 \times 10^{-2}$	$1.4 \times 10^{-4}$

Table 6.7: Comparing FEM eigenvalues using quintics with 20 subintervals to FEM eigenvalues using cubics with 30 subintervals.

relative difference less than  $10^{-3}$ . Using quintics, the first 20 eigenvalues, that is approximately a third of the number of eigenvalues calculated, have relative difference less than  $10^{-3}$ .

In conclusion, for the same computational effort (same size of the matrices), quintics yield twice as many eigenvalues sufficiently accurate than when cubics are used i.e. to obtain the first  $k$  FEM eigenvalues with relative difference less than  $10^{-3}$ ,  $7k/2$  subintervals must be used with cubics as basis functions and  $k$  subintervals with quintics.

### 6.4.3 The effect of rotary inertia

We now consider the effect of rotary inertia on the eigenvalues and use quintics as basis functions.

Note that this eigenvalue problem differs from the one excluding rotary inertia, Section 3.5. The parameter  $r$  is a measure of the effect of rotary inertia, Section 2.2.

We start by establishing convergence of the FEM eigenvalues for the case where rotary inertia is included, thereafter, we investigate the effect of rotary inertia on the eigenvalues.

As for the case without rotary inertia, the numerical results indicate convergence of the FEM eigenvalues. In Table 6.8 typical results for the relative differences,  $(\lambda^{(2n)} - \lambda^{(n)})/\lambda^{(2n)}$ , including rotary inertia, are listed for 2, 4, 8 and 16 subintervals respectively.



$i$	$ \lambda_i^{(2n)} - \lambda_i^{(n)} /\lambda_i^{(2n)}$			
	$n = 2$	$n = 4$	$n = 8$	$n = 16$
1	$2.1 \times 10^{-8}$	$7.8 \times 10^{-11}$	$2.9 \times 10^{-11}$	$7.6 \times 10^{-10}$
2	$7.8 \times 10^{-6}$	$8.9 \times 10^{-8}$	$3.2 \times 10^{-10}$	$1.4 \times 10^{-11}$
4	$1.6 \times 10^{-3}$	$4.0 \times 10^{-6}$	$2.0 \times 10^{-8}$	$7.9 \times 10^{-11}$
8	—	$1.6 \times 10^{-2}$	$9.5 \times 10^{-5}$	$3.3 \times 10^{-7}$

Table 6.8: Relative differences for FEM eigenvalues including rotary inertia with  $1/r = 4800$ .

The numerical results again suggests convergence of the FEM eigenvalues. The same pattern with respect to the order of convergence is observed as for the case without rotary inertia.

The presence of rotary inertia decreases the values of corresponding eigenvalues in comparison to the case without rotary inertia. Furthermore, the bigger the parameter  $r$ , the greater the change in the eigenvalues in comparison to the case without rotary inertia. In Table 6.9 we list eigenvalues for different values of  $r$  as well as for the case without rotary inertia ( $r = 0$ ). We use 32 subintervals for these approximations.

In Table 6.9  $\lambda_i$  denotes the  $i$ th FEM eigenvalue.

$i$	$\lambda_i$ with $r = 0$	$1/r = 19200$	$1/r = 4800$	$1/r = 1200$
1	11.81469	11.81138	11.80145	11.76186
2	406.01614	403.97925	397.71100	370.06960
4	12544.12940	11057.05461	5401.71657	3576.19925
8	273293.79309	169832.12061	158744.64019	125970.79578

Table 6.9: FEM eigenvalues for different effects of rotary inertia using 32 subintervals.

These results are for the dimensionless case. Where rotary inertia is included, two dimensionless constants,  $T$  and  $r$ , must be calculated if the results is to be connected to a specific beam, Section 2.2.

### Modes

In [ZVV] we showed that only up to the seventh so called exact mode can be computed before computational difficulties are encountered. Therefore we

consider the convergence of the FEM modes using quintics as basis functions and include rotary inertia.

Let  $\bar{w}_i^{(n)}$  denote the FEM approximation for the  $i$ th mode using  $n$  elements, normalised with respect to the infinity norm,  $\|\cdot\|_\infty$ .

The way in which we ordered our basis elements implies that the first  $n + 1$  components of  $\bar{w}_i^{(n)}$  are associated with the function values at the  $n + 1$  nodes. The next  $n + 2$  values represent the values of the first order derivatives at the nodes. Two values are associated with the point where the damage occurs. Quintics as basis functions also yield approximations for the values of the second order derivatives, and the last  $n + 1$  values of  $\bar{w}_i^{(n)}$  represent the values of the second order derivatives at the nodes.

In Table 6.10 the numerical convergence of the FEM modes are illustrated. We list the differences  $\|\bar{w}_i^{(2n)} - \bar{w}_i^{(n)}\|_\infty$ ,  $\|(\bar{w}_i^{(2n)})' - (\bar{w}_i^{(n)})'\|_\infty$  and  $\|(\bar{w}_i^{(2n)})'' - (\bar{w}_i^{(n)})''\|_\infty$  for different values of  $n$ .

$i$	$\ \bar{w}_i^{(2n)} - \bar{w}_i^{(n)}\ _\infty$ $\ (\bar{w}_i^{(2n)})' - (\bar{w}_i^{(n)})'\ _\infty$ $\ (\bar{w}_i^{(2n)})'' - (\bar{w}_i^{(n)})''\ _\infty$		
	$n = 4$	$n = 8$	$n = 16$
1	$6.83140 \times 10^{-6}$	$4.56592 \times 10^{-7}$	$2.578292 \times 10^{-8}$
	$9.65478 \times 10^{-6}$	$6.45299 \times 10^{-7}$	$3.64876 \times 10^{-8}$
	$2.48714 \times 10^{-6}$	$6.33818 \times 10^{-8}$	$8.43408 \times 10^{-9}$
2	$2.49236 \times 10^{-5}$	$2.09579 \times 10^{-6}$	$1.47615 \times 10^{-7}$
	$1.18092 \times 10^{-4}$	$9.93009 \times 10^{-6}$	$6.99423 \times 10^{-7}$
	$1.96527 \times 10^{-4}$	$6.46058 \times 10^{-6}$	$2.06282 \times 10^{-7}$
4	$1.46719 \times 10^{-5}$	$6.47026 \times 10^{-6}$	$5.44719 \times 10^{-7}$
	$1.61154 \times 10^{-4}$	$7.01558 \times 10^{-5}$	$5.90054 \times 10^{-6}$
	$8.16804 \times 10^{-3}$	$2.53189 \times 10^{-4}$	$8.28369 \times 10^{-6}$
8	$3.70776 \times 10^{-4}$	$6.16539 \times 10^{-6}$	$1.75076 \times 10^{-6}$
	$7.90080 \times 10^{-3}$	$1.30886 \times 10^{-4}$	$3.71925 \times 10^{-5}$
	$6.26608 \times 10^{-2}$	$1.00243 \times 10^{-3}$	$2.90161 \times 10^{-5}$
14	$3.20162 \times 10^{-4}$	$9.48592 \times 10^{-5}$	$9.65450 \times 10^{-7}$
	$1.68312 \times 10^{-2}$	$4.23303 \times 10^{-3}$	$4.29824 \times 10^{-5}$
	1.17602	$1.65933 \times 10^{-1}$	$1.00174 \times 10^{-2}$
16	—	$9.03027 \times 10^{-5}$	$4.81429 \times 10^{-6}$
	—	$4.78628 \times 10^{-3}$	$2.60210 \times 10^{-4}$
	—	$1.91004 \times 10^{-1}$	$2.44992 \times 10^{-2}$

Table 6.10: Convergence of FEM modes with  $\delta = 0.1$ ,  $\alpha = 0.5$  and  $1/r = 4800$ .

The rate at which convergence of the function values and the first order derivatives occur, differ from the convergence rate of the second order derivatives, which is much slower. Those modes that are associated with the first eigenvalues, starting with the smallest, converges faster than the modes associated with later eigenvalues. (Convergence in the energy norm implies that the second order derivative converges in the mean.)

## 6.5 Numerical results. Initial value problem

Consider the initial value problem. From Section 4.1 we have the following system of differential equations

$$M\bar{u}''(t) = -L\bar{u}'(t) - K\bar{u}(t).$$

For the numerical experimentations, we choose the following initial conditions:  $u'_h(0) = 0$  and  $u_h(0)$  a quintic “solitary wave”.

To approximate the solution of this problem we use the difference scheme in Section 5.4 with  $\rho_0 = 2\rho_1 = 1/2$ .

$$\begin{aligned} \left(\frac{M}{\delta t^2} + \frac{L}{2\delta t} + \frac{1}{4}K\right)\bar{u}_{k+1} + \left(-2\frac{M}{\delta t^2} + \frac{1}{2}K\right)\bar{u}_k + \\ \left(\frac{M}{\delta t^2} - \frac{L}{2\delta t} + \frac{1}{4}K\right)\bar{u}_{k-1} = 0. \end{aligned}$$

Since the initial velocity is zero, we have  $\bar{u}_1 = \bar{u}_{-1}$ .

The results obtained for the eigenvalue problem motivated us to use quintics as basis functions.

### Convergence

To verify convergence, we choose a fixed spacial discretization and a fixed final time  $\tau$ . Then, starting with 10 time intervals, we increased the number of intervals until the relative difference is strictly less than  $10^{-3}$ . This approximation is then considered as sufficiently accurate for the system of differential equations.

Decreasing the time step size, we found the first order derivatives needed approximately double the number of time steps to yield the same relative difference in  $\|\cdot\|_\infty$  than the function values do. It seems as if the second order derivatives do not converge point wise, if they do, the convergence is very slow. This is not altogether surprising (see Section 5.4).

To establish the number of elements needed for our approximation, we choose a fixed final time,  $\tau$ , and time step size,  $\delta t$ . Then the number of elements, starting with 10, is doubled until the relative difference satisfy our criterion.

### Simulation of the motion of beam

We are primarily concerned with the detection of damage. In this section we give an indication of the effect of respectively damage, damping and rotary inertia on the motion of a beam.

Our experiments indicate that measurable differences between the undamaged and damaged beams occur in displacements as well as gradients. (Table 6.11.) Viscous damping has no significant effect on the motion. Looking at the modal analysis this was expected, since it only effects the first few modes. Adding Kelvin-Voigt damping, the differences between the damaged and undamaged cases decrease, but is still clearly detectable. (Table 6.12.) The presence of rotary inertia can have a more significant effect on the difference between the motion of the damaged and undamaged beams. (Tables 6.13 and 6.14).

To illustrate the above effects we compare the motion of an undamaged beam to that of a damaged beam where the initial velocity is zero and the initial position a 'solitary wave'. For this simulation we choose  $\alpha = 0.4$ ,  $\delta = 0.1$  which is rather excessive, 80 elements,  $\tau = 0.02$  and 400 time subintervals.

Almost immediately after the first wave front pass through the damaged point, measurable differences in displacements as well as gradients between the two cases occur. (See Figure 6.1.) In Table 6.11 we compare the displacement of the damaged and undamaged beams on  $\tau = 0.02$  at  $x = 0.3$  and  $x = 0.7$ .

$x$	Undamaged beam	Damaged beam	% difference
0.3	$2.325 \times 10^{-1}$	$1.505 \times 10^{-1}$	8.2
0.7	$4.733 \times 10^{-1}$	$5.508 \times 10^{-1}$	7.8

Table 6.11: *Effect of damage during motion where  $\delta = 0.1$ ,  $\alpha = 0.4$ ,  $\tau = 0.02$ .*

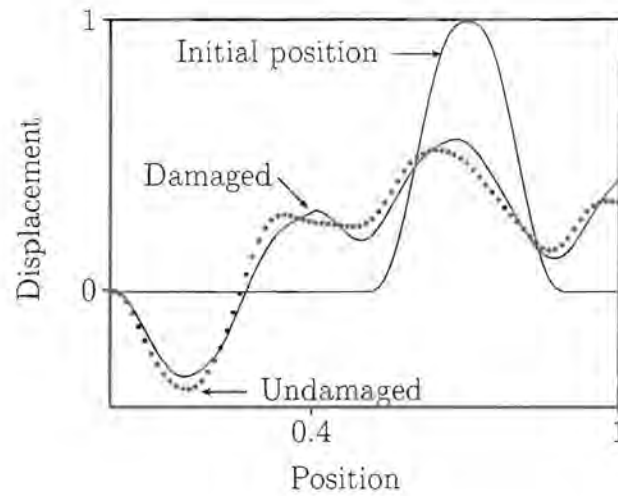


Figure 6.1: Comparing the motion of an undamaged beam to that of a damaged beam where  $\delta = 0.1$ ,  $\alpha = 0.4$  and  $\tau = 0.02$ .

We now add Kelvin-Voigt damping to the same situation as in the previous case. We use  $\mu = 3.469 \times 10^{-5}$ . This value for  $\mu$  was obtained from [JVRV].

$x$	Undamaged beam	Damaged beam	% difference
0.3	$2.099 \times 10^{-1}$	$1.378 \times 10^{-1}$	7.2
0.7	$4.716 \times 10^{-1}$	$5.436 \times 10^{-1}$	7.2

Table 6.12: Effect of Kelvin-Voigt damping on the damage during motion where  $\delta = 0.1$ ,  $\alpha = 0.4$ ,  $\tau = 0.02$  and  $\mu = 3.469 \times 10^{-5}$ .

The presence of rotary inertia can make the differences more difficult to detect. An example is given in Tables 6.13 and 6.14.

$x$	Undamaged beam	Damaged beam	% difference
0.2	$-2.525 \times 10^{-1}$	$-2.065 \times 10^{-1}$	4.6
0.6	$4.762 \times 10^{-1}$	$5.248 \times 10^{-1}$	4.9

Table 6.13: *Effect of Rotary inertia on the damage during motion where  $\delta = 0.1$ ,  $\alpha = 0.4$ ,  $\tau = 0.02$  and  $1/r = 19200$ .*

$x$	Undamaged beam	Damaged beam	% difference
0.2	$-2.236 \times 10^{-1}$	$-2.092 \times 10^{-1}$	1.4
0.6	$4.140 \times 10^{-1}$	$4.287 \times 10^{-1}$	1.5

Table 6.14: *Effect of Rotary inertia on the damage during motion where  $\delta = 0.1$ ,  $\alpha = 0.4$ ,  $\tau = 0.02$  and  $1/r = 4800$ .*

## Chapter 7

# Application. Plate beam model

### 7.1 Introduction

We consider Problem 3 (from Section 2.6). It is a mathematical model for a plate connected to two beams. Problems of this type are clearly of great practical importance. The plate can be rigidly connected to the beams or simply supported by the beams. The same model can be used for an I-shaped structural member (depending on the type of vibration). For simplicity we restrict our investigation to the case of a plate supported by beams. If the plate is rigidly connected to the beams, it may result in a problem with six unknown functions (excluding shear) due to dynamical effects. Even in our restricted case, one may easily encounter very large matrices.

In collaboration with others, [ZVGV1], we considered the equilibrium and eigenvalue problems of a rectangular plate supported by two beams at the boundary. In this thesis we extend the investigation and include the effect of rotary inertia.

The computation of the matrices is explained in Section 7.2. We use reduced quintics for the plate, which necessitates the use of quintics for the beams. We treat the equilibrium problem in Section 7.3 and the eigenvalue problem Section 7.4.



## 7.2 Computation of matrices

For the numerical experimentation we consider a square plate,  $\Omega$ , rigidly supported at two opposing sides and supported by identical beams at the remaining sides. The plate has thickness  $h$  and the beams are of square profile with thickness  $d$ . Furthermore, we assume the plate and beams are of the same material. (These restrictions are evidently not necessary.)

The reference configuration  $\Omega$  is the rectangle with  $0 < x_1 < 1$  and  $0 < x_2 < 1$ . Those parts of the boundary where  $x_1 = 0$  and  $x_1 = 1$  are denoted by  $\Sigma_0$  and  $\Sigma_1$  respectively and correspond to the rigidly supported parts of the boundary. Those parts where  $x_2 = 0$  and  $x_2 = 1$  are denoted by  $\Gamma_0$  and  $\Gamma_1$  respectively and correspond to the sections of the boundary supported by beams.

### 7.2.1 Basis elements

For the plate we use only reduced quintics as basis functions. These functions are in  $H^2(\Omega)$  or fully conforming, in finite element language. They are defined on a triangular mesh. The mesh for the rectangle  $\bar{\Omega}$  is generated in the following way: The interval  $[0, 1]$  is divided into  $n_1$  subintervals and the interval  $[0, 1]$  into  $n_2$  subintervals. This partition of the intervals yields  $n_1 \times n_2$  rectangles. The final triangular mesh is then obtained by dividing each of these rectangles into two triangles by connecting the lower left corner with the upper right corner. The rectangle  $\Omega$  is divided into  $2n_1 \times n_2$  triangles. Consequently we have  $2n_1 \times n_2$  elements  $\Omega_i$ .

Reduced quintics are defined in Section 4.3.1. The computation of the coefficients is not trivial and we describe it in Appendix C. The choice of reduced quintics “force” one to use quintics for the beams. Hermite piecewise quintics are defined in Section 4.2.1.

## 7.2.2 Standard Matrices

First we compute standard matrices for the two beams with quintics. The procedure is the same as with cubics.

$$\begin{aligned} [M_0^{\Gamma_0}]_{ij} &= \int_0^1 (\gamma_0 \phi_i)(\gamma_0 \phi_j), & [M_1^{\Gamma_0}]_{ij} &= \int_0^1 (\gamma_0 \phi_i)'(\gamma_0 \phi_j)', \\ [M_0^{\Gamma_1}]_{ij} &= \int_0^1 (\gamma_1 \phi_i)(\gamma_1 \phi_j), & [M_1^{\Gamma_1}]_{ij} &= \int_0^1 (\gamma_1 \phi_i)'(\gamma_1 \phi_j)', \end{aligned}$$

as well as

$$K_{ij}^{\Gamma_0} = \int_0^1 (\gamma_0 \phi_i)''(\gamma_0 \phi_j)'' \quad \text{and} \quad K_{ij}^{\Gamma_1} = \int_0^1 (\gamma_1 \phi_i)''(\gamma_1 \phi_j)''.$$

Next we compute standard matrices for the plate. These computations are quite involved and we provide some detail in Appendix C.

$$[M_0^\Omega]_{ij} = \int_\Omega \phi_i \phi_j, \quad [M_1^\Omega]_{ij} = \int_\Omega \text{grad } \phi_i \cdot \text{grad } \phi_j \quad \text{and} \quad K_{ij}^\Omega = b_\Omega(\phi_i, \phi_j).$$

The bilinear forms are given in Section 3.3. Each basis element is of the form  $\tilde{\phi}_i = \langle \phi_i, \gamma_0 \phi_i, \gamma_1 \phi_i \rangle$ . (The forced boundary conditions are satisfied by eliminating certain basis elements.) Now, consider for example

$$c(\tilde{\phi}_i, \tilde{\phi}_j) = c_\Omega(\phi_i, \phi_j) + \beta c_0(\phi_i, \phi_j) + \beta c_1(\phi_i, \phi_j).$$

$c_0(\phi_i, \phi_j)$  involves the restriction of basis functions  $\phi_i$  and  $\phi_j$  to the boundary  $\Gamma_0$ . These restrictions are non-zero only for some of the basis functions associated with nodes on  $\Gamma_0$ . (The restriction of a reduced quintic on  $\Gamma_0$  is an one-dimensional quintic.) The result is  $M_{ij} = M_{ij}^\Omega + \beta M_{ij}^{\Gamma_0} + \beta M_{ij}^{\Gamma_1}$ . Consequently  $M = M^\Omega + \beta M^{\Gamma_0} + \beta M^{\Gamma_1}$ ,

where

$$M^\Omega = M_0^\Omega + r M_1^\Omega, \quad M^{\Gamma_0} = M_0^{\Gamma_0} + r_b M_1^{\Gamma_0} \quad \text{and} \quad M^{\Gamma_1} = M_0^{\Gamma_1} + r_b M_1^{\Gamma_1}.$$

The computation of the  $K$ -matrix is the similar,

$$K = K^\Omega + \alpha K^{\Gamma_0} + \alpha K^{\Gamma_1}.$$

### 7.3 Equilibrium problem

To find the Galerkin approximation for the solutions of the equilibrium problem, we solve a system of linear equations.

#### Problem BD

$$K\bar{u} = \bar{F}, \text{ where } F_i = (f, \phi_i).$$

The parameter  $\alpha$  gives an indication of the stiffness of the beams in comparison to that of the plate. Increasing the value of  $\alpha$  implies an increase in the stiffness of the beams and  $\alpha = 0$  corresponds to the case where two sides are free.

For different values of  $\alpha$  we compare in Table 7.1 the FEM approximations for the maximum displacement, to the so called exact solution. See [TW]. (Interesting historical remarks are found in [TW].)

Note that the maximum displacement occurs at the centre of the plate as a result of symmetry.

We consider a square plate with the same number of equal intervals per side. We denote this number by  $n$ , and use it to distinguish between different meshes.

Denote the maximum displacement obtained from the so called exact solution by  $u_{\max}$  and the FEM approximation of the maximum displacement where  $n$  subintervals are used, by  $u_{\max}^{(n)}$ . Choose Poisson's ratio  $\nu = 0.3$ .

$\alpha$	Exact	$(u_{\max} - u_{\max}^{(n)})/u_{\max}$		
		$n = 2$	$n = 4$	$n = 8$
100	$4.09 \times 10^{-3}$	$2.0421 \times 10^{-3}$	$2.9308 \times 10^{-4}$	$2.3653 \times 10^{-4}$
30	$4.16 \times 10^{-3}$	$1.0507 \times 10^{-3}$	$6.5975 \times 10^{-4}$	$7.1510 \times 10^{-4}$
10	$4.34 \times 10^{-3}$	$1.7896 \times 10^{-3}$	$1.7460 \times 10^{-4}$	$1.2220 \times 10^{-4}$
6	$4.54 \times 10^{-3}$	$3.5724 \times 10^{-3}$	$5.0933 \times 10^{-3}$	$5.1428 \times 10^{-3}$
4	$4.72 \times 10^{-3}$	$2.9835 \times 10^{-3}$	$1.547 \times 10^{-3}$	$1.5006 \times 10^{-3}$
2	$5.29 \times 10^{-3}$	$3.2719 \times 10^{-3}$	$2.0580 \times 10^{-3}$	$2.0181 \times 10^{-3}$
1	$6.24 \times 10^{-3}$	$1.0228 \times 10^{-3}$	$9.3231 \times 10^{-5}$	$6.2112 \times 10^{-5}$
0.5	$7.56 \times 10^{-3}$	$2.2617 \times 10^{-3}$	$1.6195 \times 10^{-3}$	$1.5973 \times 10^{-3}$
0	$1.309 \times 10^{-2}$	$3.2828 \times 10^{-4}$	$2.8584 \times 10^{-4}$	$2.8129 \times 10^{-4}$

Table 7.1: Comparison of exact values with FEM approximations of maximum displacement.

The fact that the relative error originally improves if we double the number of intervals from 2 to 4 and then remains almost the same, suggests that the so called exact solution is not very accurate, as could be expected since only a few significant digits are given.

The relative difference between consecutive FEM approximations strengthens this observation as can be seen in Table 7.2.

$\alpha$	$(u_{\max}^{(4)} - u_{\max}^{(2)})/u_{\max}^{(4)}$	$(u_{\max}^{(8)} - u_{\max}^{(4)})/u_{\max}^{(8)}$
100	$1.748505 \times 10^{-3}$	$5.653982 \times 10^{-5}$
30	$1.711612 \times 10^{-3}$	$5.539486 \times 10^{-5}$
10	$1.614713 \times 10^{-3}$	$5.238781 \times 10^{-5}$
6	$1.528720 \times 10^{-3}$	$4.971930 \times 10^{-5}$
4	$1.433896 \times 10^{-3}$	$4.677682 \times 10^{-5}$
2	$1.211367 \times 10^{-3}$	$3.987090 \times 10^{-5}$
1	$9.294900 \times 10^{-4}$	$3.111738 \times 10^{-5}$
0.5	$6.412267 \times 10^{-4}$	$2.214781 \times 10^{-5}$
0	$4.242264 \times 10^{-5}$	$4.550660 \times 10^{-6}$

Table 7.2: Comparison of FEM approximations for the maximum displacement.

## 7.4 Eigenvalue problem

As mentioned before, Section 4.4, the occurrence of eigenvalues has a highly irregular pattern in the two-dimensional case. We have an elementary example to illustrate this, and also to show how difficult it can be to identify eigenvalues with multiplicity.

### 7.4.1 Multiplicity of eigenvalues

Consider the following eigenvalue problem,

$$-\nabla^2 u = \lambda u \text{ on the unit square with } u = 0 \text{ on the boundary.}$$

Clearly,

$u(x, y) = \sin(n\pi x) \sin(m\pi y)$  is an eigenfunction, for  $n$  and  $m$  integers. The corresponding eigenvalue is  $\lambda = n^2 + m^2$ .

The popular difference scheme is

$$-h^{-2} [u_{i,j+1} + u_{i+1,j} - 4u_{i,j} + u_{i,j-1} + u_{i-1,j}] = \lambda u_{i,j}$$

or

$$-h^{-2} [u_{i,j+1} - 2u_{i,j} + u_{i,j-1}] - h^{-2} [u_{i+1,j} - 2u_{i,j} + u_{i-1,j}] = \lambda u_{i,j},$$

where  $h$  is the length of a subinterval.

Let  $u_{i,j} = \sin(i\omega_k) \sin(j\omega_\ell)$ , then

$$u_{i+1,j} + u_{i-1,j} = 2 \cos(\omega_k) u_{i,j},$$

Hence

$$-h^{-2} [u_{i+1,j} - 2u_{i,j} + u_{i-1,j}] = \lambda_k u_{i,j}, \text{ where } \lambda_k = h^{-2}(2 - 2 \cos \omega_k).$$

$u_{i,j}$  satisfies the boundary conditions if  $\omega_k = (k\pi)/(n+1)$  and  $\omega_\ell = (\ell\pi)/(n+1)$ . It satisfies the difference equations if  $\lambda_\ell = h^{-2}(2 - 2 \cos \omega_\ell)$ . Hence  $u_{i,j}$  is an eigenvector and every eigenvalue is of the form  $\lambda_k + \lambda_\ell$ .

In Table 7.3 we list the exact eigenvalues for this problem as well as the numerical approximations obtained for different subinterval lengths. We give

$i$	Exact	$h = 0.2$	$h = 0.1$	$h = 0.05$	$h = 0.005$
1	19.7	19.3	19.6	19.7	19.7
2	49.3	45.6	48.2	49.0	49.3
3	49.3	45.6	48.2	49.0	49.3
4	79.0	72	76.8	78.4	78.8
5	98.7	81.6	93.3	97.2	98.3
6	98.7	81.6	93.3	97.2	98.3
7	128	108	122	127	128
8	128	108	122	127	128
9	168	118	151	163	167
10	168	118	151	163	167
11	178	144	167	175	177
12	197	144	180	192	196
13	197	144	180	192	196
14	247	144	217	240	245
15	247	144	217	240	245
16	257	170	225	245	254
17	257	170	225	245	254
18	286	180	246	275	283
19	286	180	246	275	283
20	316	206	283	307	313

Table 7.3: Comparison of finite difference eigenvalues to the exact eigenvalues.

only three significant digits as it is sufficient to illustrate difficulties of matching exact eigenvalues and approximate eigenvalues.

Let  $i$  denote the number of the eigenvalue and  $h$  the length of a subinterval.

Interpreting numerical results with respect to multiplicity of eigenvalues is difficult. Great care should be taken to establish whether approximate eigenvalues that are close together are an indication of multiplicity of exact eigenvalues, or not.

For example, for  $h = 0.2$ , the eleventh to fifteenth eigenvalues seem to be one eigenvalue with multiplicity five, while it actually approximates three different eigenvalues. Another example is the fifteenth and sixteenth eigenvalues for  $h = 0.05$ . These eigenvalues seem close and one might expect them to approximate the same eigenvalue with multiplicity more than one.

### 7.4.2 Plate beam

The generalized eigenvalue problem associated with the plate beam model is given by

**Problem CD**

$$K\bar{w} = \lambda M\bar{w}.$$

In a joint report [ZVGV1] we consider this eigenvalue problem for the plate beam model excluding rotary inertia. In this subsection we investigate the effect of rotary inertia if included in the model.

The ratio  $\alpha/\beta$  of the dimensionless constants,  $\alpha = (E_b I_b)/(aD)$  and  $\beta = (\rho_b A)(\rho a h)$ , defined in Section 2.6, with the plate of thickness  $h$  and the beams of square profile with thickness  $d$ , is a measure of the stiffness of the beams in comparison to that of the plate.

In the special case where both the beams and the plate are of the same material, we have

$$\frac{\alpha}{\beta} = \left(\frac{d}{h}\right)^2 (1 - \nu^2).$$

As the values of  $d/h$  increase, i.e. the stiffness of the beams is increased, the situation approaches the plate problem where all four sides of the plate are rigidly supported. For this problem the eigenvalues and eigenfunctions are known. The eigenvalues are of the form

$$((n\pi)^2 + (m\pi)^2)^2$$

with corresponding eigenfunctions

$$\sin(n\pi x) \sin(m\pi y).$$

Since the exact eigenvalues for the plate beam problem are not available, the FEM approximations for the eigenvalues for large values of  $d/h$  can be compared to the eigenvalues of this limiting case, see Table 7.4.

Denote the  $i$ th eigenvalue for the case where all four sides are rigidly supported by  $\lambda_i$ . The eigenvalues are ordered according to size. The FEM approximation of the  $i$ th eigenvalue is denoted by  $\lambda_i^{(n)}$  where  $n$  subintervals are used.

Throughout this subsection we use Poisson's ratio as  $\nu = 0.3$ .

$i$	$\lambda_i^{(8)}$ for different values of $d/h$				$\lambda_i$
	$d/h = 1$	$d/h = 10$	$d/h = 100$	$d/h = 200$	
1	92.6654	386.3556	389.6361	389.6364	389.6364
2	250.7783	2359.4575	2435.2366	2435.2398	2435.2273
3	1264.1968	2433.5697	2435.2500	2435.2525	2435.2273
4	1514.0745	6221.0210	6234.2338	6234.2346	6234.1818
5	2142.1461	7345.9342	9741.4914	9741.5319	9740.9091
8	7725.6133	11308.1573	16463.7086	16463.7127	16462.1364
10	11599.1799	16455.3398	28158.9627	28159.0563	28151.2273

Table 7.4: Comparison of FEM eigenvalues for different values of  $d/h$  with eigenvalues of a rigidly supported plate. Rotary inertia is excluded.

As  $d/h$  increases, the FEM approximation of the eigenvalues approaches the eigenvalues of the plate rigidly supported on all four sides.

For values of  $d/h$  that do not correspond to the limit case, the numerical convergence of the FEM eigenvalues are illustrated in Table 7.5.

$i$	$d/h = 10$		$d/h = 100$	
	$(\lambda_i^{(2n)} - \lambda_i^{(n)}) / \lambda_i^{(2n)}$		$(\lambda_i^{(2n)} - \lambda_i^{(n)}) / \lambda_i^{(2n)}$	
	$n = 2$	$n = 4$	$n = 2$	$n = 4$
1	$6.94630 \times 10^{-4}$	$8.2136 \times 10^{-6}$	$7.0601 \times 10^{-4}$	$8.3506 \times 10^{-6}$
2	$2.1888 \times 10^{-2}$	$4.0239 \times 10^{-4}$	$1.2764 \times 10^{-2}$	$2.9014 \times 10^{-4}$
3	$4.1506 \times 10^{-2}$	$4.2474 \times 10^{-4}$	$5.3144 \times 10^{-2}$	$5.6098 \times 10^{-4}$
4	$6.6013 \times 10^{-2}$	$7.3256 \times 10^{-4}$	$6.6438 \times 10^{-2}$	$7.3537 \times 10^{-4}$
5	$2.3179 \times 10^{-2}$	$9.5076 \times 10^{-4}$	$4.8512 \times 10^{-2}$	$3.1069 \times 10^{-4}$
10	$2.3327 \times 10^{-1}$	$5.0927 \times 10^{-3}$	$5.7506 \times 10^{-1}$	$1.3180 \times 10^{-2}$

Table 7.5: Numerical convergence of FEM eigenvalues for different values of  $d/h$ . Rotary inertia is excluded.

**Remark** Choosing  $n = 8$ , yields  $(486 \times 486)$  matrices which are already very time consuming to handle with our available computer hardware and software. Therefore we do not consider more than 8 subintervals.

### Including rotary inertia in the model

In addition to the joint report [ZVGV1] we now establish the effect of rotary



inertia on the eigenvalues of the plate beam problem.

From Sections 2.5 and 2.6 we have the dimensionless constants  $r_b = I_b/(a^2d^2)$  and  $r = I/(a^2h)$ . In the experimentation we work with a fixed plate, i.e.  $a$  and  $h$  are fixed, and modify the beams by changing  $d$ . Consequently  $r_b$  and  $r$  depend on the relationship  $d/h$  and indicate the effect of rotary inertia.

In Table 7.6 we illustrate numerical convergence of the FEM eigenvalues if rotary inertia is included. For  $d/h = 50$  we have  $r_b = 2.083 \times 10^{-2}$  and  $r = 8.333 \times 10^{-6}$ .

$i$	$(\lambda_i^{(4)} - \lambda_i^{(2)}) / \lambda_i^{(4)}$	$(\lambda_i^{(8)} - \lambda_i^{(4)}) / \lambda_i^{(8)}$
1	$7.0599 \times 10^{-4}$	$8.3271 \times 10^{-6}$
5	$4.8512 \times 10^{-2}$	$3.1350 \times 10^{-3}$
10	$5.7479 \times 10^{-1}$	$1.312 \times 10^{-2}$

Table 7.6: Numerical convergence of FEM eigenvalues for  $d/h = 50$ . Including rotary inertia.

For the plate supported on all four sides, repeated eigenvalues are expected—and indeed observed. For the plate beam problem the symmetry is partially lost, and the question arises if repeated eigenvalues will occur, and whether those FEM eigenvalues will be observed as repeated eigenvalues?

As with the case excluding rotary inertia, the exact solution is not available. Again, the exact eigenvalues for the plate rigidly supported on all four sides are used to give an indication of what can be expected of the FEM eigenvalues.

As is expected, the presence of rotary inertia decreases the eigenvalues in comparison to the case without rotary inertia. The effect of rotary inertia is illustrated in Table 7.7. For  $d/h = 50$  we have  $r_b = 2.083 \times 10^{-2}$  and  $r = 8.333 \times 10^{-6}$ .

$i$	Excluding rotary inertia $\lambda_i^{(8)}$	Including rotary inertia $\lambda_i^{(8)}$	Exact value $\lambda_i$
1	389.6313	389.5673	389.6364
2	2435.1545	2434.1536	2435.2273
3	2435.2439	2434.2429	2435.2273
4	6234.2146	6230.1154	6234.1818
5	9740.7884	9732.7844	9740.9091
6	9741.5344	9733.5289	9740.9091
7	16463.2389	16445.6552	16462.1364
8	16463.6608	16446.0765	16462.1364
9	28154.6617	28115.3573	28151.2273
10	28158.7179	28119.4013	28151.2273
11	31564.1696	31517.5097	31560.5455

Table 7.7: *Effect of rotary inertia on the eigenvalues for  $d/h = 10$ .*

In Table 7.7 the multiplicity of eigenvalues of the plate rigidly supported on all four sides are observed. These repetitions give reason to expect that the corresponding FEM eigenvalues for the plate beam problem may also be repeated eigenvalues.

The question arises whether the FEM approximation will yield repeated eigenvalues or will the eigenvalues only be close?