





#### ABSTRACT

## BLENDING-FUNCTION TECHNIQUES WITH APPLICATIONS TO DISCRETE LEAST SQUARES

By

Dale Russel Doty

The theory of blending function spaces (bivariate interpolation) is developed in the general setting of interpolation spaces. In this setting it is shown that blending function spaces have the desirable quality of doubling the order of accuracy with less computation when compared to standard tensor product spaces.

The dimensionality of discretized blending function spaces is derived, and several bases are explicitly constructed. The special example of Hermite spline blended piecewise polynomials is developed, showing that these spaces have bases with small support which are easy to calculate. These spaces offer maximum order of convergence for a minimum number of basis elements. For example, linearly blended piecewise cubic polynomials offer a fourth order approxima tion scheme, and cubic Hermite spline blended piecewise polynomials offer an eighth order scheme. Next, using the exponential decay of the natural cubic cardinal splines and natural cubic spline blending, a derivative-free approximation scheme is developed, which is eighth order in the interior of the domain.

Algorithms with corresponding error estimates are given for solving the discrete least squares problem with unstructured data. For the univariate case, algorithms are developed using the space of cubic splines. The resulting error analysis indicates the necessary restrictions to be placed on the number and distribution of the data points to insure that the discrete least squares fit will be  $O(h^m)$  to a function  $f \in C^m[a, b]$  from which the data arises, where h is the mesh size and  $l \leq m \leq 4$ . An example is given to illustrate that the discrete least squares fit need not be close to f if these conditions are not realized. For the bivariate case, algorithms and error analyses are given for the spaces of bicubic splines and discretized blending function spaces. It is shown that the discrete least squares fit to a bivariate function f is of the same order accuracy as the corresponding interpolation accuracy.

Discrete least squares is considered on general domains which have curved boundaries and are possibly multiply connected. This general domain is subdivided into "standard" subdomains, and explicit mappings from the unit square to these standard subdomains are constructed which are one-one, onto, and have easily calculated inverses. Thus, discrete least squares over general domains reduces to the cases previously considered.

Finally, an extensive computational error analysis is given for a constrained least squares algorithm.

## BLENDING-FUNCTION TECHNIQUES WI TH APPLICATIONS TO DISCRETE LEAST SQUARES

By

Dale Russel Doty

## A DISSERTATION

Submitted to Michigan State University in partial fulfillment of the requirements for the degree of

### DOCTOR OF PHILOSOPHY

Department of Mathematics

1975

### ACKNOWLEDGMENTS

I wish to express my appreciation to Professor Gerald D. Taylor for his encouragement and guidance during the preparation of this thesis.

# TABLE OF CONTENTS

CHAPTER 1	COMPUTA FOR DISC	COMPUTATIONAL ERROR ANALYSIS FOR DISCRETE LEAST SQUARES													
	Section 1	Least Squares with Constraints 2													
		Section 1.1 Pivoting 5 Section 1.2 A Constrained Least Squares Algorithm 7													
	Section 2	Defining the Perturbation													
		Problem													
	Section 3	Matrix Norm Inequalities 12													
	Section 4	Condition Numbers of Non-													
		square Matrices17													
	Section 5	An Upper Bound for $  \delta Q   \dots \dots 20$													
	Section 6	Basic Relations in the Perturba-													
		tional Problem													
	Section 7	Effects of Perturbational Errors 37													
	Section 8	Rounding Errors in Computation 45													
		Section 8.1 Vector Addition 46 Section 8.2 Matrix Multiplication 46													
		angular Set of Equations . 47 Section 8.4 Orthogonal Transfor-													
		mation													
	Section 9	Rounding Error Analysis 51													
CHAPTER 2	BLENDIN	G FUNCTION THEORY 65													
	Section 1	Spline and Blending Function Interpolation													
	Section 2	Dimension of Discretized Blending													
	_	Function Spaces													
	Section 3	Natural Cubic Blending 102													
	Section 4	Exponential Decay of Natural													
	Cubic Cardinal Splines iii														

Page

CHAPTER 3	DISCRET	E LEAST SQUARES	21									
	Section 1	Uniform Error Estimates 12	22									
	Section 2	A Uniform Bound 12	29									
	Section 3	Univariate Discrete Least										
		Squares	40									
	Section 4	Bivariate Least Squares										
		with Data on Mesh Lines 1	55									
	Section 5	Bivariate Least Squares										
		with Unstructured Data 1	72									
		Section 5.1 Bicubic Splines	72									
		Section 5.2 Cubically Blended										
		Cubic Splines 1	78									
		Section 5.3 Hermite Blended										
		Piecewise Polynomials. 18	80									
		Section 5.4 Linear Blending 18	87									
	Section 6	Domain Transformations 19										
		Section 6.1 Type 1 Domain										
		Transformations 20	იი									
		Section 6.2 Type 2 Transfor-	50									
		mations	11									
		Section 6.3 General Domains 2	13									
		Section 6.4 Discrete Least										
		Squares over $\Gamma$ 22	15									
BIBLIOGRAP	HY		18									
	• •		~ ~									

# LIST OF FIGURES

																														Page
Figure	1	••	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	108
Figure	2	The	e ]	Re	egi	.01	n <b>S</b>	Ω	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	200
Figure	3	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	214
Figure	4	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	214
Figure	5	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	215

.

#### CHAPTER 1

### COMPUTATIONAL ERROR ANALYSIS FOR DISCRETE LEAST SQUARES

The following chapters will deal with solving the constrained least squares problems which arise from applications of Blending Function spaces.

Due to computational considerations, it is necessary to use methods of orthogonal factorization and pseudo-inverses to obtain solutions to these problems. It is therefore necessary to perform a perturbational and computational analysis of these methods to show how the numerical solutions are affected. We would hope that the condition numbers,  $\chi$ , associated with our least squares and constraint equations, appear only to the first power in the error analysis, be cause our solutions contain only inverses to the first power. This, however, has been shown not to be the case, quoting van der Sluis [31]:

> "It caused something of a shock, therefore, when in 1966 Golub and Wilkinson asserted that already the multiplications QA and QB may produce errors in the solution containing a factor  $\chi^2(A)$ ."

We will prove that the condition numbers associated with the <sup>cons</sup> trained least squares equations appear only linearly in the error <sup>analysis</sup> except for the coefficient of the residual. Section 1. Least Squares With Constraints

We are given the two matrices

(1.1) 
$$A_{\epsilon} \mathbb{R}^{m \times n}$$
, and

$$(1.2) \qquad C \in \mathbb{R}^{1-1},$$

where r < n and  $m \ge n$ .

rYn

Then we seek a solution  $x \in \mathbb{R}^n$  which minimizes the norm of R

(1.3) 
$$R = Ax - f$$
, where  $f \in \mathbb{R}^{11}$ ,

and satisfies

(1.4) 
$$Cx = g$$
, where  $g \in \mathbb{R}^{r}$ .

We can think of (1.3) as a least squares problem subject to the **constraint** equations given in (1.4).

We will develop here the method of Halliday and Hayes [21] for finding a solution to (1.3) and (1.4). But first we need to state some of the pertinent theories of factorization by orthogonal transformations. Householder [22] developed the theory of factorization into orthogonal transformations. Precisely, he has shown that if at step i the matrix  $C_i$  has the form

$$(1.5) C_i = \left[\frac{L_i \mid 0}{T_i}\right]$$

where  $L_i \in \mathbb{R}^{(i-1)x(i-1)}$  is a lower-triangular matrix,  $T_i \in \mathbb{R}^{(r-(i-1))xn}$ is a rectangular matrix, and  $C_1 = C$ , then the orthogonal matrix  $P_i \in \mathbb{R}^{nxn}$  which introduces zeros into row i, from i+1 to n of the matrix  $C_{i+1}$ , and leaves the first i-l columns of  $C_i$  unaltered is given by

(1.6) 
$$P_i = I - v_i v_i^T / H_i$$
,

where

(1.7) 
$$C_{i+1} = C_i P_i$$
.

If we represent the entries of  $C_i$  by  $c_{jk}^i$  for  $l \leq j \leq r$ ,  $l \leq k \leq n$ , then the vector  $v_i \in \mathbb{R}^n$  is given by

(1.8) 
$$v_i^T = (0, \dots, 0, c_{ii}^i + sgn(c_{ii}^i) S_i, c_{i,i+1}^i, \dots, c_{i,n}^i)$$
.

where  $sgn(\cdot)$  is a function of a real variable defined by

(1.9) 
$$\operatorname{sgn}(x) = \begin{cases} +1 & \text{if } x \ge 0 \\ -1 & \text{if } x < 0 \end{cases}$$

S<sub>i</sub> is defined by

$$(1.10) \qquad S_{i} = \begin{bmatrix} n \\ \Sigma \\ j=1 \end{bmatrix} (c_{ij}^{i})^{2}$$

Finally  $H_i$  is defined by

(1.11) 
$$H_i = S_i^2 + |c_{ii}| |S_i|$$

Also, he has shown that

$$(1.12)$$
  $|c_{i,i}^{i+1}| = S_i$ ,

or that the new diagonal element has the Euclidean length of the old row i, from i to n. This is a property of orthogonal transformations, that the Euclidean length of any transformed row remains invariant.

If for some i,  $S_i = 0$ , then row i of  $C_i$  is zero from i to n, and we take  $P_i = I$ . If this is not the case, then both  $S_i$  and  $H_i$  are strictly positive. In either case, the factorization proceeds until the completion of step r, and we have a matrix,  $C_{r+1}$ , which is lower triangular

(1.13)  $[L|0] = C_{r+1}$ , where  $L = L_{r+1}$ .

Define Q by

(1.14)  $Q = P_1 \cdot P_2 \cdot \cdots \cdot P_r$ .

Then Q is also an orthogonal matrix, i.e.,  $Q^{-1} = Q^{T}$ , and it follows from (1.7) that

- (1.15a) [L|0] = CQ, or
- (1.15b) L = CQ<sub>1</sub>,

where  $Q_1$  is the first r columns of  $Q_2$ ,

 $(1.15 c) \qquad Q_1 = Q\left[\frac{I_r}{0}\right].$ 

Then we may write (1.4) as

(1.16)  $[L|0] Q^{T} x = g.$ 

#### Section 1.1. Pivoting

Row pivoting can be included in the above algorithm for factoring C. This is done in the following manner. At step i of the factorization, row k, where  $i \le k \le r$ , is chosen as the next pivotal row by some pivotal strategy. Then we premultiply by a matrix  $E_i \in \mathbb{R}^{rxr}$ , which interchanges the rows i and k, see Deskins [7, p. 551]. We now factor  $P_i^{-1}$  out of the matrix  $E_i C_i$ and define

$$(1.17) C_{i+1} = (E_i C_i) P_i,$$

where  $P_i$  is obtained from row i of  $E_i C_i$ . We can conclude from (1.4), (1.14) and (1.17) that

 $(1.18) \qquad \begin{bmatrix} L \mid 0 \end{bmatrix} Q^{T} \mathbf{x} = \mathbf{E}_{r} \cdot \cdots \cdot \mathbf{E}_{l} \mathbf{g}$  $= \mathbf{E}_{r} \cdot \cdots \cdot \mathbf{E}_{l} \mathbf{C} \mathbf{x}.$ 

Therefore, for theoretical purposes, we will assume that the matrix <sup>C</sup> and vector g with which we are working, have had their rows or <sup>elements</sup> permuted beforehand. So that, when we apply our pivotal <sup>strategy</sup>, the pivotal rows will be chosen sequentially from 1 to r.

The type of pivoting that is usually used in practice is called "maximal row pivoting". At step i, this type of pivoting chooses the next pivot by the criterion that  $S_i$  is maximized. If two or more rows give the same value, then the first of these is chosen as the pivotal row. For this type of pivoting, Golub [11] has shown that (1.19)  $S_1 \ge S_2 \ge \cdots \ge S_r \ge 0$ .

If for some i,  $S_i = 0$ , then from (1.12), L will have a zero at position i on its diagonal, which means that the rank of L will be less than r. But from (1.15a) we have  $C = [L|0]Q^T$  which implies that the rank of C will also be less than r, see Deskins [7, p. 550]. Therefore, if we assume that C is of full rank r, then for  $1 \le i \le r$  we must have

(1.20)  $S_i > 0$ ,

and because the  $S_i$  for  $1 \leq i \leq r$  are the absolute value of the diagonal entries of L we have

$$(1.21)$$
 L<sup>-1</sup> exists.

It should be noted that if  $r \ge n$ , then  $C^{T}$  can be factored in the following way

$$C^{T} P_{1} \cdot \cdots \cdot P_{n} = [L|0],$$

where by taking transposes and using  $P_i^T = P_i \epsilon R^{r \times r}$ 

$$P_n \cdot \cdots \cdot P_l C = \left[\frac{W}{0}\right],$$

and  $W = L^{T}$  is an upper-triangular matrix. We will use both types of factorizations in what follows. Section 1.2. A Constrained Least Squares Algorithm

Next, we shall describe an algorithm for the solution of constrained least squares which was developed by Hayes and Halliday [21]. This will be presented in detail, because we need their formulas for later reference.

<u>Algorithm 1.1.</u> (Hayes and Halliday [21]). If we have a least squares problem with constraints, as defined in (1.1)-(1.4), where C is of full rank, then the solution x is obtained as follows.

<u>Step 1:</u> Let C be of full rank, then for  $1 \le i \le r$ , we have from (1. 20) that  $S_i > 0$  and L is nonsingular. Thus there exists a Q given by (1. 14) such that in (1. 16)

(1.16) [L] 0] 
$$Q^T x = g$$
.

Step 2: Define the vectors  $d_1 \in \mathbb{R}^r$  and  $d_2 \in \mathbb{R}^{n-r}$  such that (1.23)  $\begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = Q^T x$ .

Then from (1.16) and (1.23) we have

$$(1.24) d_1 = L^{-1} g.$$

Step 3: Solve for d<sub>2</sub>, which upon inverting (1.23) will yield a solution x

$$(1.25) x = Q \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}$$

Toward this end, using (1.3) and (1.23) we have

(1.26) A x = A Q Q<sup>T</sup> x  
= 
$$[B_1 | B_2] Q^T x$$
  
=  $B_1 d_1 + B_2 d_2$   
= f + R,

where  $B_1$  is the first r columns, and  $B_2$  the remaining n-r columns of the matrix A Q

(1.27) 
$$[B_1 | B_2] = AQ.$$

From (1.26) we have

(1.28) 
$$B_2 d_2 = f - B_1 d_1 + R$$
.

Let B<sub>2</sub> be of full rank, then the factorization of B<sub>2</sub>, using orthogonal transformations, proceeds in the following way. Define

$$(1.29)$$
  $B_{21} = B_{2}$ 

and for 14 i4n-r

(1.30)  $B_{2,i+1} = U_i B_{2,i}$ 

where at step i, the orthogonal transformation  $U_{i} \in \mathbb{R}^{m \times m}$  is defined

by (1.6), where  $B_{2,i}^{T}$  plays the role of  $C_{i}$  in (1.8)-(1.12), and we are factoring from the left. Define

(1.31) 
$$V = U_{n-r} \dots U_1$$

and

(1.32a) 
$$\left[\frac{W}{0}\right] = V B_2 = B_{2,m+1}$$

or

(1.32b) 
$$W = V_1 B_2$$
,

where from (1.20) and (1.21),  $W_{\xi} \mathbb{R}^{(n-r)x(n-r)}$  is an upper triangular nonsingular matrix. Then the solution of (1.28) which minimizes R, [see 21], is given by

(1.33) 
$$d_2 = W^{-1} V_1 \{ f - B_1 d_1 \}$$

where  $V_1$  is the first n-r rows of V

(1.34) 
$$V_1 = [I_{n-r} | 0] V$$
.

It is shown by Peters and Wilkinson [28] that  $d_2$  as defined (1.33) is unique. Then by using (1.24), (1.33) and (1.25) we have the solution

(1.35) 
$$\mathbf{x} = \mathbf{Q} \begin{bmatrix} \mathbf{L}^{-1} \mathbf{g} \\ \\ \\ \\ \mathbf{W}^{-1} \mathbf{V}_{1} (\mathbf{f} - \mathbf{B}_{1} \mathbf{L}^{-1} \mathbf{g}) \end{bmatrix}$$

#### Section 2. Defining the Perturbation Problem

In this section, we shall derive an upper bound on the computational errors of the solution defined by Algorithm 1.1. Toward this end, we will develop a perturbational analysis of the problem defined in (1.1) to (1.4). Due to calculation and rounding errors, the problem which is stored in the computer is not the exact matrices and vectors as defined in (1.3) and (1.4). Instead, the original quantities have been perturbed to give the new problem:

(2.1) We are given the two matrices  $\hat{A} \in \mathbb{R}^{m \times n}$  and  $\hat{C} \in \mathbb{R}^{m \times n}$  where

$$(2.2) r < n and n \leqslant m.$$

Then the solution  $\hat{\mathbf{x}}$  of the perturbed problem is the vector which minimizes the norm of

(2.3)  $\hat{\mathbf{R}} = \hat{\mathbf{A}}\hat{\mathbf{x}} - \hat{\mathbf{f}}, \text{ where } \hat{\mathbf{f}} \in \mathbf{R}^{\mathbf{m}},$ 

and we require that  $\hat{\mathbf{x}}$  satisfy the constraint

(2.4) 
$$\hat{C}\hat{x} = \hat{g}$$
, where  $\hat{g} \in \mathbb{R}^r$ .

We will denote with a "hat" all of the perturbed quantities, and their meaning will be the same as in Section 1. If  $\hat{C}$  and  $\hat{B}_2$  are of full rank, then Algorithm 1.1 gives the solution  $\hat{x}$ 

(2.5) 
$$\hat{\mathbf{x}} = \hat{\mathbf{Q}} \begin{bmatrix} \hat{\mathbf{L}}^{-1} \hat{\mathbf{g}} \\ \\ \\ \\ \hat{\mathbf{W}}^{-1} \hat{\mathbf{V}}_{1} \{ \hat{\mathbf{f}} - \hat{\mathbf{B}}_{1} \hat{\mathbf{L}}^{-1} \hat{\mathbf{g}} \} \end{bmatrix}$$

We now define the perturbational quantities  $\delta A$ ,  $\delta f$ ,  $\delta C$ ,  $\delta g$ ,  $\delta P_i$ ,  $\delta L$  etc. as the difference between a perturbed and unperturbed quantity. For example,  $\delta A = \hat{A} - A$ ,  $\delta f = \hat{f} - f$ , etc.

Before we can estimate  $||\hat{\mathbf{x}} - \mathbf{x}||$ , we need some well known facts from the theory of matrices. Due to the lack of references for the following norm relations for non-square matrices, we will include in the next section a somewhat detailed development for completeness.

#### Section 3. Matrix Norm Inequalities

Definition 3.1: Let  $A \in \mathbb{R}^{n \times n}$ , then the complex number  $\lambda_i(A)$  is an eigenvalue of the matrix A if and only if there exists a non-zero eigenvector  $\mathbf{x}_i \in \mathbb{C}^n$  such that  $A \times_i = \lambda_i(A) \times_i$ . The spectral radius  $\rho(\cdot)$  of A is defined to be  $\rho(A) = \max |\lambda_i(A)|$ , see Varga [32, p. 9], and let  $\mu(\cdot)$  of A be defined by  $\mu(A) = \min_i \{|\lambda_i(A)|| \lambda_i(A) \neq 0\}$ .

Definition 3.2: The Euclidean norm  $||\cdot||$  of a vector  $\mathbf{x} \in \mathbb{R}^{m}$  is defined to be  $||\mathbf{x}|| = (\sum_{i=1}^{m} (\mathbf{x}_{i})^{2})^{\frac{1}{2}}$ . The Euclidean (Schur) norm  $||\cdot||_{E}$  of a rectangular matrix  $A \in \mathbb{R}^{n \times m}$  is defined to be  $||A||_{E} = \begin{bmatrix} n & m \\ \sum & \sum \\ i=1 & j=1 \end{bmatrix}^{\frac{1}{2}}$ , see [36, p. 81].

Definition 3.3: The spectral matrix norm  $||\cdot||$  of a rectangular matrix  $A \in \mathbb{R}^{n \times m}$ , induced by the Euclidean vector norm, is defined to be  $||A|| = \sup ||A \times ||$ , see Varga [32, p. 9]. ||x||=1 $x \in \mathbb{R}^{m}$ 

<u>Remark</u>: For a vector  $x \in \mathbb{R}^m$ , the Euclidean vector norm, Euclidean matrix norm and spectral matrix norm are the same, see the follow-ing lemma.

Lemma 3.1. Given the rectangular matrix  $A \in \mathbb{R}^{n \times m}$  and vector  $x \in \mathbb{R}^{m}$  then

(3.1)  $||A|| = [\rho(A^T A)]^{\frac{1}{2}}$ .

(3.2a) 
$$\rho(A^T A) = \rho(AA^T)$$
.  
(3.2b)  $\mu(A^T A) = \mu(A A^T)$ .  
(3.2b)  $\mu(A^T A) = \mu(A A^T)$ .  
(3.3) If  $A A^T \in \mathbb{R}^{nxm}$  is nonsingular, then  
 $||(A A^T)^{-1}|| = 1/\mu(A A^T)$ .  
(3.4)  $||A||_E / \sqrt{\pi} \le ||A|| \le 1||A||_E$ .  
(3.5)  $||A^T||_E = ||A||_E$ .  
(3.6)  $||A^T||_E = ||A||_E$ .  
(3.7) If  $Q \in \mathbb{R}^{nxm}$  is orthogonal (i. e.  $Q^T = Q^{-1}$ ), then  $||Q|| = 1$ .  
(3.8) If  $Q \in \mathbb{R}^{nxm}$  is orthogonal, then  $||QA|| = ||A||$ .  
(3.9) If  $Q \in \mathbb{R}^{mxm}$  is orthogonal, then  $||AQ|| = ||A||$ .  
(3.10)  $||Ax|| \le ||A|| ||x||$ .  
(3.11) If  $B \in \mathbb{R}^{mxm}$ , then  $||AB|| \le ||A|| ||B||$ .  
(3.12) If  $B \in \mathbb{R}^{nxm}$ , then  $||A+B|| \le ||A|| + ||B||$ .  
(3.13) If  $A_1 \in \mathbb{R}^{nxm}$  is a matrix obtained by deleting the first  
or last r columns of A, then  $||A_1|| \le ||A||$ .  
(3.14) If  $A_1 \in \mathbb{R}^{(n-r)xm}$  is a matrix obtained by deleting the  
first or last r rows of A, then  $||A_1|| \le ||A||$ .  
(3.15) If  $Q \in \mathbb{R}^{nxm}$  is an orthogonal matrix and  $Q_1$  is the first  
r rows or columns of Q, then  $||Q_1|| \le 1$ .  
(3.16) If  $B \in \mathbb{R}^{nxm}$  and  $||B|| < 1$ , then the following inverses  
exist and  
 $||(I + B)^{-1}|| \le 1/(1 - ||B||)$ .

•

13

(3.17) If B, 
$$\hat{B}$$
,  $\delta B \in \mathbb{R}^{n \times n}$ , both  $B^{-1}$  and  $\hat{B}^{-1}$  exist and  $\delta B = \hat{B} - B$ , then  $\hat{B}^{-1} = B^{-1} - \hat{B}^{-1} \delta B B^{-1}$ .

(3.18) If the vector v is any row or column of A, then
$$||v|| \leq ||A||.$$

<u>Proof (3.1)</u>: From Varga [32, p. 11] we have that because  $A^{T} A$  is a non-negative definite (i.e.  $x^{T} A^{T} A x \ge 0$ ) symmetric square matrix, the eigenvalues are all real and non-negative. Since  $||A x||^{2}/||x||^{2} = x^{T} A^{T} A x/x^{T} x$ , we have  $||A x||^{2}/||x||^{2} \le \rho(A^{T} A)$ , where equality is taken on for an eigenvector corresponding to  $\rho(A^{T} A)$ .

(3.2a), (3.2b): We will show that the set of non-zero eigenvalues of  $A^T A_{\epsilon} \mathbb{R}^{m \times m}$  is the same as that of  $A A^T {}_{\epsilon} \mathbb{R}^{n \times m}$ . Let  $\lambda$  be a non-zero eigenvalue of  $A^T A$ , then  $\lambda$  is real and there exists  $x \neq 0$  such that  $A^T A x = \lambda x$ . Premultiply by A and we have  $(A A^T)A x = \lambda A x$ . The vector A x is non-zero, because if it were zero we would have  $\lambda x = A^T (A x) = 0$ , and because  $\lambda \neq 0$ , this would imply that x = 0, which is a contradiction. Therefore,  $\lambda$  is an eigenvalue of  $A A^T$ . The reverse inclusion is proved in an identical manner.

(3.3):  $A A^{T}$  is positive definite because  $x^{T}AA^{T}x = ||A^{T}x||^{2} \ge 0$ , and for  $x \ne 0$  we have, by using an argument similar to the one given in (3.2a) that  $||A^{T}x|| > 0$ . Because  $AA^{T}$  is a

square matrix, we can use property G of Wilkinson  $\begin{bmatrix} 34 & p. & 290 \end{bmatrix}$  and conclude our result.

(3.4): If 
$$||\mathbf{x}|| = 1$$
, then  $||\mathbf{A}\mathbf{x}||^2 = \sum_{i=1}^{n} \left[ \sum_{j=1}^{m} \mathbf{a}_{ij} \mathbf{x}_j \right]^2$   
 $\leq \sum_{i=1}^{n} \left[ \sum_{j=1}^{m} \mathbf{a}_{ij}^2 \cdot \sum_{j=1}^{m} \mathbf{x}_j^2 \right] = ||\mathbf{A}||_{\mathbf{E}},$ 

where we have used the Schwarz inequality on each term in brackets. Therefore, by using the definition of the spectral norm, we have  $||A|| \leq ||A||_{E}$ .

To prove the left hand side of the inequality we use (3.1) and (3.2) to obtain  $||A||^2 = \rho(AA^T)$ . Because  $AA^T \in \mathbb{R}^{n\times n}$ , we have from Marcus and Ming [26, p. 23] that  $\sum_{i=1}^{n} (AA^T)_{ii} = \sum_{i=1}^{n} \lambda_i (AA^T)$ . Because  $x^T AA^T x = ||A^T x||^2 \ge 0$ , and using [26, p. 69], we have that  $AA^T$  is non-negative definite and  $\lambda_i(AA^T) \ge 0$  for  $1 \le i \le n$ . Therefore  $n \cdot \rho(AA^T) \ge \sum_{i=1}^{n} \lambda_i(AA^T) = \sum_{i=1}^{n} (AA^T)_{ii} = \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} a_{ij}^2 = ||A||_E^2$ , and we have proved (3.4).

(3.5): Use (3.1) and (3.2a).

(3.6): Use the definition of  $||A||_{E}$ . (3.7): From (3.1) we have  $||Q||^{2} = \rho(Q^{T} Q) = \rho(I) = 1$ ,

because all eigenvalues of the identity are 1.

(3.8): 
$$||QA||^2 = \rho(QA)^TQA = \rho(A^TA) = ||A||^2$$
.  
(3.9): Because  $Q^T$  is orthogonal, use (3.5) and (3.8).  
(3.10): Follows from the definition of  $||A||$ .

(3.11): Using Varga [32, p. 11], the norm of AB is obtained for an unit eigenvector x corresponding to  $\rho((AB)^{T}(AB))$ . Then  $||AB|| = ||A(Bx)|| \leq ||A|| ||Bx|| \leq ||A|| ||B||$ , where we have used (3.10).

(3.12): Similar to (3.11).

(3.13): If  $A_1$  is obtained by deleting the last r columns of A, then  $A_1 = A\left[\frac{I}{0}\right]$ , where  $I \in \mathbb{R}^{(m-r)x(m-r)}$  is the identity matrix. From (3.11) and (3.1),  $||A_1||^2 \leq ||A||^2 \rho([I|0]\left[\frac{I}{0}\right]) = ||A||^2$ . The other case is done in the same way.

from which our result follows.

(3.18): We will assume that  $v \in \mathbb{R}^n$  is row l of A, where  $l \leq l \leq n$ . Define the vector  $e_l \in \mathbb{R}^n$  to have zeros in all entries except for a unit in entry l. It follows that  $v^T = (e_l)^T A$  and  $||e_l^T||^2 = ||e_l||^2 = \rho(e_l^T e_l) = \rho(1) = 1$ .

### Section 4. Condition Numbers of Non-square Matrices

<u>Definition 4.1</u>: If  $A_{\epsilon} \mathbb{R}^{m \times n}$  is any matrix, then the condition number  $\chi(\cdot)$  of A is defined to be

(4.1) 
$$\chi(A) = \left[\rho(A^{T} A) / \mu(A^{T} A)\right]^{\frac{1}{2}} \ge 1,$$

if A is of full rank, otherwise  $\chi(A)$  is undefined. It should be recalled from Definition 3.1 that  $\mu(A^T A)$  is the smallest <u>non-zero</u> eigenvalue of  $A^T A$  in absolute value.

There are many definitions for condition numbers, depending upon the type of matrix. We will list some of them here and show that they are equivalent to (4.1).

If  $A \in \mathbb{R}^{n \times n}$  is a nonsingular matrix, then the condition number  $\chi(\cdot)$  of A is defined to be

(4.2)  $\chi(A) = ||A^{-1}|| \cdot ||A|| \ge 1.$ 

If A is singular, then  $\chi(A)$  is undefined, (see [23, p. 81]).

We would like to show that (4.1) and (4.2) are equivalent for m = n. This is clear, because if A is not of full rank, then  $\chi(A)$ is undefined in both (4.1) and (4.2). If A is of full rank and square, then  $A^{-1}$  exists, and using property (I) of Wilkinson [34, p. 290] and (3.3) we have  $||A^{-1}||^2 = ||A^{-1}A^{-T}|| = ||(A^T A)^{-1}|| = 1/\mu(A^T A)$ . Thus we have that (4.1) is equivalent to (4.2).

If  $A \in \mathbb{R}^{m \times n}$  is a rectangular matrix, where  $m \ge n$ , then the condition number  $\chi(\cdot)$  of A is defined to be

(4.3) 
$$\chi(\mathbf{A}) = \sup_{\substack{||\mathbf{x}||=1\\\mathbf{x} \in \mathbb{R}^n}} ||\mathbf{A} \mathbf{x}|| / \inf_{\substack{||\mathbf{x}||=1\\\mathbf{x} \in \mathbb{R}^n}} ||\mathbf{A} \mathbf{x}|| \geqslant 1,$$

if A is of full rank n, otherwise  $\chi(A)$  is undefined, (see Bjorck [3]).

To see that (4.1) and (4.3) are equivalent for  $m \ge n$ , we have from (3.1) and Definition 3.3 that  $\sup_{\substack{||\mathbf{x}||=1}} ||\mathbf{A}\mathbf{x}|| = \rho^{\frac{1}{2}} (\mathbf{A}^T \mathbf{A})$ . Also,

we have that  $A^{T} A \in \mathbb{R}^{n \times n}$  is of full rank, because if it were not, there would exist a  $x \in \mathbb{R}^{n}$  such that  $x \neq 0$  and  $A^{T} A x = 0$ . But this implies that  $x^{T} A^{T} A x = ||A x||^{2} = 0$ , and in turn that A x = 0. This cannot happen if A is of full rank n and  $x \neq 0$ . Therefore, we have that the smallest eigenvalue of  $A^{T} A$  is non-zero and  $\min_{i} \lambda_{i}(A^{T} A) = \mu(A^{T} A)$ . Varga [32, p. 11], has shown that  $l \leq i \leq n$ on for the unit eigenvector x corresponding to  $\mu(A^{T} A)$ . Therefore,  $\inf_{i} ||A x|| = \mu^{\frac{1}{2}} (A^{T} A)$ , and we see that (4.3) is equivalent to (4.1). ||x||=1

<u>Remark</u>: The importance of defining  $\mu(A^T A)$  to be the smalles <u>non-zero</u> eigenvalue in absolute value comes from the case where m < n. Here we have that  $A^T A$  has a zero eigenvalue even if A is of full rank, however,  $A A^T$  has not. Using (3.2b), we have that  $\mu(A^T A) = \mu(A A^T) \neq 0$ , and our condition number is finite. **<u>Remark</u>**: The typical situation in error analysis is that we must bound the following

(4.4) 
$$||A x|| ||A^{-1}y|| \leq \chi(A) ||x|| ||y||,$$

where A is nonsingular and square. Note that the left hand side of (4.4) will be much smaller than the right hand side unless both the vectors x and y are eigenvectors corresponding to  $\rho(A^T A)$  and  $\mu(A^T A)$  respectively. Therefore, in a practical problem, if A is ill-conditioned, we would expect (4.4) to be a gross overestimate of  $||A x|| ||A^{-1} y||$ .

<u>**Remark:**</u> If Q is orthogonal, then  $\chi(Q) = 1$ .

# Section 5. An Upper Bound for $||\delta Q||$

We will develop here a bound for the norm of the perturbational quantity  $\delta Q = \hat{Q} - Q$ , where Q is defined by (1.14). It should be noted that Q and  $\hat{Q}$  are not the unique orthogonal matrices which reduce C and  $\hat{C}$  respectively if r < n. For example, if we consider Q, then recalling (1.15a) we have [L|0] = CQ. If we post multiply by the orthogonal matrix  $P = I - 2uu^{T} \epsilon \mathbb{R}^{n \times n}$ , where  $u \epsilon \mathbb{R}^{n}$ is any unit vector which has its first r entries set to zero, then C(QP) = [L|0]. This is true because P will leave unaltered the first r columns of CQ. Therefore, we cannot expect to obtain a bound for  $||\delta Q||$  by manipulating the formula (1.15a). Instead, because Q is uniquely defined by (1.14), we will use the definition of the  $P_i$  given in (1.6) to obtain our estimate.

In doing this, we shall use the following equality

(5.1) 
$$(1+K)^{s} = 1 + \sum_{j=1}^{s} K (1+K)^{j-1}$$
,

where K is a non-negative real number and s is any natural number.

Also, for the following lemma, we recall that  $c_{kk}^{k}$  and  $\hat{c}_{kk}^{k}$ are the k<sup>th</sup> diagonal elements of the matrices  $C_{k}$  and  $\hat{C}_{k}$  respectively. The matrix  $C_{k}$  and correspondingly  $\hat{C}_{k}$ , as defined in (1.7), represent the factorization at step k. Lemma 5.1. If we use the factorization method of Section 1 by orthogonal transformation with maximal row pivoting, assume after pivoting that both of the matrices C and  $\hat{C}$  have the same ordering of rows, also at step k of the factorization  $\operatorname{sgn}(c_{kk}^{k}) = \operatorname{sgn}(\hat{c}_{kk}^{k})$ for  $1 \leq k \leq r$ , and finally that there exists a real non-negative number  $\Delta$  such that

$$(5.2a) \qquad ||\delta C||/||C|| \leq \Delta / \left[ (9+4\sqrt{2}+8\Delta)^{r-1} \chi(C) \right] ,$$

then it follows that

(5.2b)  $||\delta Q|| < K(r) \chi(C) ||\delta C||/||C||,$ where  $K(r) = (1 + K)^{r} - 1$  and  $K = 4(2 + \sqrt{2} + 2\Delta).$ 

**Remark:** In (5.2a) we will use the fact that  $9 + 4\sqrt{2} + 8\Delta = 1 + K$ .

**Proof:** To estimate  $||\delta Q||$ , we express it as the following telescoping series

$$(5.3) \qquad ||\delta Q|| = ||\hat{P}_{1} \cdot \hat{P}_{2} \cdot \cdots \cdot \hat{P}_{r} - P_{1} \cdot P_{2} \cdot \cdots \cdot P_{r}||$$
$$= ||\delta P_{1} \cdot \hat{P}_{2} \cdot \cdots \cdot \hat{P}_{r} + P_{1} \cdot \delta P_{2} \cdot \hat{P}_{3} \cdot \cdots \cdot \hat{P}_{r}$$
$$+ \cdots + P_{1} \cdot P_{2} \cdot \cdots \cdot P_{r-1} \cdot \delta P_{r}||$$
$$\leqslant \sum_{k=1}^{r} ||\delta P_{k}||,$$

where we have used (3.8) and (3.9) to obtain the last inequality.

To estimate  $\delta P_k$  we recall the following definitions from sections one and two, which we tabulate here for convenience,  $l \leq k \leq r$ 

(5.4) 
$$C_{1} = C, C_{k+1} = C_{k} P_{k},$$
  $\hat{C}_{1} = \hat{C}, \hat{C}_{k+1} = \hat{C}_{k} \hat{P}_{k},$   
(5.5)  $[L|0] = C Q,$   $[\hat{L}|0] = \hat{C} \hat{Q},$   
(5.6)  $P_{k} = I - v_{k} v_{k}^{T} / H_{k},$   $\hat{P}_{k} = I - \hat{v}_{k} \hat{v}_{k}^{T} / \hat{H}_{k},$   
(5.7)  $S_{k} = \begin{bmatrix} n \\ \Sigma \\ \ell = k \end{bmatrix} (c_{k\ell}^{k})^{2} \frac{1}{2},$   $\hat{S}_{k} = \begin{bmatrix} n \\ \Sigma \\ \ell = k \end{bmatrix} (c_{k\ell}^{k})^{2} \frac{1}{2},$   
(5.8)  $H_{k} = S_{k}^{2} + |c_{kk}^{k}| S_{k},$   $\hat{H}_{k} = \hat{S}_{k}^{2} + |\hat{c}_{kk}^{k}| \hat{S}_{k}$ .

The 
$$v_k \hat{v}_k \in \mathbb{P}^n$$
 are defined as  
(5.9)  $v_k^T = (0, \dots, 0, c_{kk}^k + sgn(c_{kk}^k) S_k, c_{k,k+1}^k, \dots, c_{k,n}^k),$   
(5.10)  $\hat{v}_k^T = (0, \dots, 0, \hat{c}_{kk}^k + sgn(\hat{c}_{kk}^k) \hat{S}_k, \hat{c}_{k,k+1}^k, \dots, \hat{c}_{k,n}^k),$ 

where the first k-1 entries are zero. Finally, we have from (5.9) (5.10), (5.7) and (5.8) the relations

(5.11)  $||\mathbf{v}_{k}||^{2} = 2(S_{k}^{2} + |c_{kk}^{k}| S_{k}) = 2 H_{k}$ (5.12)  $||\mathbf{v}_{k}^{*}||^{2} = 2(S_{k}^{2} + |\mathbf{c}_{kk}^{k}| \mathbf{s}_{k}) = 2 \hat{H}_{k}$ .

We will now relate the perturbation  $\delta P_k$  to the perturbation at step k in  $\delta C_k = \hat{C}_k - C_k$ . From (5.6) we have

(5.13) 
$$\delta P_{k} = (I - \hat{v}_{k} \hat{v}_{k}^{T} / \hat{H}_{k}) - (I - v_{k} v_{k}^{T} / H_{k})$$
$$= (v_{k} v_{k}^{T} - \hat{v}_{k} \hat{v}_{k}^{T}) / H_{k} + \hat{v}_{k} \hat{v}_{k}^{T} (\hat{H}_{k} - H_{k}) / (\hat{H}_{k} \cdot H_{k}) .$$

We define  $\delta v_k \in \mathbb{R}^n$  by

(5.14) 
$$\delta \mathbf{v}_{k} = \hat{\mathbf{v}}_{k} - \mathbf{v}_{k}$$
  
=  $(0, \dots, 0, \delta c_{kk}^{k} + \text{sgn} (c_{kk}^{k})(\hat{\mathbf{S}}_{k} - \mathbf{S}_{k}), \delta c_{k, k+1}^{k}, \dots, \delta c_{k, n}^{k})^{T}$ ,

where we have made use of our assumption that  $\operatorname{sgn}(c_{kk}^k) = \operatorname{sgn}(c_{kk}^k)$ . This assumption means that the initial perturbation of  $\mathcal{S}C$  in C is not large enough to cause a sign change in  $c_{kk}^k$  at step k of the reduction. If the signs were different, then we can conclude that both

$$(5.15) \qquad | \stackrel{\diamond}{c}_{kk}^{k} |, \quad | \stackrel{\diamond}{c}_{kk}^{k} | \leqslant | \delta \stackrel{\diamond}{c}_{kk}^{k} |$$

Later in our proof we will obtain an estimate for  $\delta C_k$  which will enable us to test if this can happen. If the signs do differ, however, we cannot guarantee that the matrix  $\delta P_k$  will have a small norm, which implies that this type of error analysis does not apply.

If we make use of (5.14), and substitute for  $v_k$  in (5.13), we obtain the following expression for  $\delta P_k$ 

$$(5.16) \quad \delta \mathbf{P}_{k} = -(\mathbf{v}_{k} \cdot \delta \mathbf{v}_{k}^{T} + \delta \mathbf{v}_{k} \cdot \mathbf{v}_{k}^{T} + \delta \mathbf{v}_{k} \cdot \delta \mathbf{v}_{k}^{T})/\mathbf{H}_{k} + \hat{\mathbf{v}}_{k} \cdot \hat{\mathbf{v}}_{k}^{T} (\hat{\mathbf{H}}_{k} - \mathbf{H}_{k})/(\hat{\mathbf{H}}_{k} \cdot \mathbf{H}_{k}).$$

We will now introduce the following notation which will help in our making norm estimates. Define the column vectors  $C_i^k$ ,  $\hat{C}_i^k$ ,  $\delta C_i^k \in \mathbb{R}^n$  to be the transpose of row i of the matrices  $C_k$ ,  $\hat{C}_k$  and  $\delta C_k$  respectively. Also, define the vectors  $\tilde{C}_i^k$ ,  $\hat{C}_i^k$ ,  $\delta \tilde{C}_i^k \in \mathbb{R}^n$  to be obtained by setting the first k-1 coordinates equal to zero of the vectors  $C_i^k$ ,  $\hat{C}_i^k$  and  $\delta C_i^k$ , respectively. If we apply the above definitions to (5.7) we have

(5.17) 
$$||\hat{c}_k^k|| = S_k \text{ and } ||\hat{c}_k^k|| = \hat{S}_k.$$

We can now use (5.17) and the above definitions to obtain a convenient expression for the norm of (5.14)

(5.18) 
$$|| \delta v_k ||^2 = || \delta \hat{C}_k^k ||^2 + 2 \delta c_{kk}^k \operatorname{sgn}(c_{kk}^k)(||\hat{\tilde{C}}_k^k|| - || \hat{\tilde{C}}_k^k ||)$$
  
  $+ (||\hat{\tilde{C}}_k^k|| - ||\tilde{\tilde{C}}_k^k||)^2.$ 

By the triangle inequality we have

(5.19) 
$$\left| ||\hat{\tilde{c}}_{k}^{k}|| - ||\tilde{c}_{k}^{k}|| \right| \leq ||\hat{\tilde{c}}_{k}^{k} - \tilde{c}_{k}^{k}||$$
  
 $\leq ||\delta \tilde{c}_{k}^{k}||.$ 

Using (5.18) and (5.19) we find

(5.20) 
$$||\delta v_{k}||^{2} \leq 2||\delta \widetilde{C}_{k}^{k}||^{2} + 2|\delta c_{kk}^{k}|\cdot||\delta \widetilde{C}_{k}^{k}||$$
  
 $\leq 4||\delta \widetilde{C}_{k}^{k}||^{2},$ 

where in the last inequality we have made use of the fact that  $|\delta c_{kk}^{k}| \leq ||\delta C_{k}^{k}||$ . This gives a bound for  $||\delta v_{k}||$ (5.21)  $||\delta v_{k}|| \leq 2||\delta \widetilde{C}_{kk}^{k}||$ .

We use (5.8), (5.17) and (5.19) to obtain a bound for  $|\hat{H}_k - H_k|$ 

$$(5.22) \qquad |\hat{H}_{k} - H_{k}| \leq |\hat{S}_{k}^{2} + |\hat{c}_{kk}^{k}|\hat{S}_{k} - (S_{k}^{2} + |c_{kk}^{k}|S_{k})| \\ \leq (\hat{S}_{k} + S_{k}) \cdot |\hat{S}_{k} - S_{k}| + |\hat{c}_{kk}^{k}| \cdot |\hat{S}_{k} - S_{k}| + |\hat{c}_{kk}^{k}| - |c_{kk}^{k}| |\cdot S_{k}|$$

$$\leq 2(||\hat{\widetilde{C}}_{k}^{k}|| + ||\widetilde{C}_{k}^{k}||) \cdot ||_{\delta}\widetilde{C}_{k}^{k}|$$

$$\leq 2(2S_{k} + ||_{\delta}\widetilde{C}_{k}^{k}||) ||_{\delta}\widetilde{C}_{k}^{k}||.$$

If we substitute (5.21) and (5.22) into (5.16) and use (3.5) and (3.11) we obtain the following after taking norms

$$(5.23) \qquad || \delta P_{k} || \leq || \delta v_{k} || (2 || v_{k} || + || \delta v_{k} ||) / H_{k} + (|| \hat{v}_{k} ||^{2} / \hat{H}_{k}) \cdot |\hat{H}_{k} - H_{k} | / H_{k} \leq 4 || \delta \tilde{C}_{k}^{k} || [(|| v_{k} || + || \delta \tilde{C}_{k}^{k} ||) / H_{k} + (2S_{k} + || \delta \tilde{C}_{k}^{k} ||) / H_{k}],$$

where we have made use of (5.12) in the last inequality. From (5.8) we have

(5.24) 
$$H_k \ge S_k^2$$
,

and using this with (5.11) leads to the estimate

(5.25) 
$$(||v_k|| + ||\delta \widetilde{C}_k^k||)/H_k \leq \sqrt{2}/\sqrt{H_k} + ||\delta \widetilde{C}_k^k||/H_k$$
  
 $\leq (\sqrt{2} + ||\delta \widetilde{C}_k^k||/S_k)/S_k.$ 

If we combine (5.11), (5.23) and (5.25), we have the following bound for  $||\delta P_k||$ , using  $||\delta \tilde{C}_k^k|| \leq ||\delta C_k^k||$ (5.26)  $||\delta P_k|| \leq 4||\delta C_k^k|| \cdot [(2+\sqrt{2})+2||\delta C_k^k||/S_k]/S_k$ .

A condition will be given later, as to when we have a bound for  $||_{\delta}C_k^k||/s_k^k.$ 

We now relate the perturbation in step k, which is  $\delta C_k$ , to the initial perturbation  $\delta C$ . We will do this by tracing the error backwards step by step. Toward this end we calculate the following estimates. From (5.4) we have for  $1 \le i \le r$ 

(5.27) 
$$\delta C_{i+1} = \delta C_i \hat{P}_i + C_i \delta P_i .$$

For  $1 \leq \ell \leq r$ , row  $\ell$  of the matrix  $\delta C_{i+1}$  is given by

(5.28) 
$$(\delta C_{\ell}^{i+1})^{T} = (\delta C_{\ell}^{i})^{T} \dot{P}_{i} + (C_{\ell}^{i})^{T} \delta P_{i}$$
.

We want to replace the last term in (5.28) by something more convenient for our purposes. It will be shown that

(5.29) 
$$(C_{\ell}^{i})^{T} \delta P_{i} = (\widetilde{C}_{\ell}^{i})^{T} \delta P_{i}$$

To accomplish this, we recall (5.13), which gives

(5.30) 
$$\delta P_i = v_i v_i^T / H_i - \hat{v}_i \hat{v}_i^T / \hat{H}_i$$

From (5.9) and (5.10), the first i-1 entries of  $v_i$  and  $\hat{v}_i$  are zero, which, when we combine this with (5.30), implies that the first i-1 rows of  $\delta P_i$  are zero. This implies that irrespective of how the first i-1 entries of  $C_{\ell}^i$  are changed, the resulting product will be the same, as long as the remaining entries of  $C_{\ell}^i$  remain unchanged. Using the fact that  $\tilde{C}_{\ell}^i$  agrees with  $C_{\ell}^i$  in all but the first i-1 entries, we have

(5.31) 
$$(\delta C_{\boldsymbol{\ell}}^{i+1})^{\mathrm{T}} = (\delta C_{\boldsymbol{\ell}}^{i})^{\mathrm{T}} \hat{P}_{i} + (\tilde{C}_{\boldsymbol{\ell}}^{i})^{\mathrm{T}} \delta P_{i}$$
(5.32) 
$$\mu_{\ell i} = ||\tilde{c}_{\ell}^{i}|| / ||\tilde{c}_{\ell}^{i}||$$
$$= ||\tilde{c}_{\ell}^{i}|| / s_{i},$$

which, because of (1.20), are all finite.

Taking norms in (5.31) and using (5.32) and (3.9) we have

(5.33) 
$$||\delta C_{\ell}^{i+1}|| \leq ||\delta C_{\ell}^{i}|| + \mu_{\ell i} S_{i} ||\delta P_{i}||.$$

We will now prove by induction on i, for  $1 \le i \le r$  and  $1 \le l \le r$ , that the following is true

(5.34) 
$$||\delta C_{\ell}^{i+1}|| \leq \sum_{j=1}^{1} \mu_{\ell j} S_{j}^{i} ||\delta P_{j}|| + ||\delta C_{\ell}^{1}||.$$

For i = 1, where  $1 \le l \le r$ , inequality (5.34) is just inequality (5.33). Assume that (5.34) is true for some i, such that  $1 \le i \le r-1$ , then we will prove that (5.34) is also true for i+1. Using (5.33) for i+1, and our induction hypothesis, we obtain for  $1 \le l \le r$ 

$$(5.35) \qquad || \delta C_{\ell}^{i+2} || \leq || \delta C_{\ell}^{i+1} || + \mu_{\ell, i+1} S_{i+1} || \delta P_{i+1} || \\ \leq \sum_{j=1}^{i} \mu_{\ell j} S_{j} || \delta P_{j} || + || \delta C_{\ell}^{1} || \\ + \mu_{\ell, i+1} S_{i+1} || \delta P_{i+1} || \\ \leq \sum_{j=1}^{i+1} \mu_{\ell j} S_{j} || \delta P_{j} || + || \delta C_{\ell}^{1} || .$$

At this point, we will make use of our assumption of "maximal row pivoting". This assumption implies that all of the  $\mu_{\ell j} \leq 1$  for  $\ell \geq j$ , because of the manner in which the next pivotal row is chosen (see Section 1.1). We have retained the  $\mu_{\ell j}$  until this time, because we will refer to (5.34) later when we discuss pivotal strategies. Using this assumption, (5.34) now becomes for  $\ell = i+1$ 

(5.36) 
$$|| \delta C_{i+1}^{i+1} || \leq \sum_{j=1}^{i} S_{j} || \delta P_{j} || + || \delta C_{i+1}^{1} ||.$$

At this point in our proof, we would like to show that (5.36)and our assumption (5.2a) leads to

(5.37) 
$$||_{\delta}C_{k}^{k}||/S_{k} \leq \Delta$$
, for  $l \leq k \leq r$ ,

and

(5.38) 
$$||\delta C_k^k|| \leq (K+1)^{k-1} ||\delta C|| \text{ for } l \leq k \leq r.$$

Where (5.37) would yield the following simplification of (5.26)

(5.39) 
$$||\delta P_k|| \leq K ||\delta C_k^k|| / S_k$$
,

where  $K = 4(2 + \sqrt{2} + 2\triangle)$ .

Toward this end we will prove the following for  $1 \leqslant k \leqslant r$ 

(5.40) 
$$||c||/s_k \leq \chi(c)$$
.

This is true because of our assumptions of "maximal row pivoting", and that C is of full rank, and because we have from (1.19) and (1.20) that

$$(5.41)$$
  $||C||/S_k \leq ||C||/S_r$ ,

where  $S_r > 0$ . Because  $S_r$  is the absolute value of element r on the diagonal of L, which is a lower triangular matrix, we have  $|\lambda_r(L)| = S_r$ , where  $\lambda_r(L)$  is an eigenvalue of L. Therefore, there is a unit vector  $x_{\epsilon} \mathbb{R}^r$ , such that  $Lx = \lambda_r(L)x$ . From (1.21)  $L^{-1}$  exists, so we have  $L^{-1}x = (1/\lambda_r(L))x$ . Taking norms, we find

(5.42) 
$$1/S_{r} = ||(1/\lambda_{r}(L))x||$$
  
=  $||L^{-1}x||$   
 $\leq ||L^{-1}||$ .

If we use (3.1) we have  $||L^{-1}||^2 = \rho(L^{-T}L^{-1}) = \rho((LL^{T})^{-1})$ . Using Wilkinson [34, p. 290] property (F), because  $(LL^{T})^{-1}$  is a symmetric matrix we have  $\rho((LL^{T})^{-1}) = ||(LL^{T})^{-1}||$ . Using (3.3), we have  $||(LL^{T})^{-1}|| = 1/\mu(LL^{T})$ , and from (5.5) and (3.2b) we have  $\mu(LL^{T}) = \mu(CC^{T}) = \mu(C^{T}C)$ . This proves that (5.40) is true, using Definition 4.1 for the condition number  $\chi(C)$ .

We shall prove (5.37), (5.38) and (5.39) simultaneously by induction on k, for  $1 \le k \le r$ , by using (5.2a) and (5.36). For k=1, using (5.40), we have that (5.37) is valid and hence also (5.39) because

(5.43) 
$$||\delta C_{1}^{1}|| / S_{1} \leq (||\delta C|| / ||C||) (||C|| / S_{1})$$
$$\leq \left[ \bigtriangleup / (\chi(C)(1+K)^{r-1}) \right] \chi(C)$$
$$\leq \bigtriangleup / (1+K)^{r-1}$$
$$\leq \bigtriangleup .$$

Also, (5.38) is trivial for k=1. Assume that (5.37), (5.38) and (5.39) are valid for all l such that  $1 \leq l \leq k$  where  $1 \leq k < r$ , then we will prove that they are also valid for k+1. Using (5.36) and our induction hypothesis, (5.38) is valid because

$$(5.44) \qquad ||\delta C_{k+1}^{k+1}|| \leq \sum_{j=1}^{k} S_{j} ||\delta P_{j}|| + ||\delta C_{k+1}^{1}||$$
$$\leq \sum_{j=1}^{k} K ||\delta C_{j}^{j}|| + ||\delta C_{k+1}^{1}||$$
$$\leq \sum_{j=1}^{k} K (K+1)^{j-1} ||\delta C|| + ||\delta C||$$
$$\leq (K+1)^{k} ||\delta C|| ,$$

where we have used (5.1). To see that (5.37) is valid for k+1

(5.45) 
$$||\delta C_{k+1}^{k+1}||/S_{k+1} \leq [(K+1)^{k} ||\delta C||/||C||] \cdot [||C||/S_{k+1}]$$
  
 $\leq \triangle (K+1)^{k} / (K+1)^{r-1}$   
 $\leq \triangle$ ,

where we have used (5.44), (5.2a) and (5.40). Finally (5.39) follows from (5.37) and (5.26), and our induction is complete.

What we have accomplished is to express the perturbation in row k at step k of our factorization in terms of the initial perturbation in C. We have also shown in (5.39) that the perturbation in the elementary orthogonal transformation at step k is bounded in terms of the perturbation in row k of  $C_k$ . Therefore, if we combine (5.38), (5.39) and (5.40), we obtain the following estimate for  $||\delta P_k||$ (5.46)  $||\delta P_k|| \leq K(1+K)^{k-1} ||\delta C||/S_k$  $\leq K(1+K)^{k-1} \chi(C) ||\delta C||/||C||.$ 

We are now in a position to obtain a bound for  $||\delta Q||$ . Using (5.3), (5.46) and (5.1), it is clear that

(5.47) 
$$|| \delta Q || \leq \begin{bmatrix} \mathbf{r} \\ \Sigma \\ \mathbf{k} = 1 \end{bmatrix} \chi(\mathbf{C}) || \delta C || / || \mathbf{C} ||$$
  
 $\leq K(\mathbf{r}) \chi(\mathbf{C}) || \delta C || / || \mathbf{C} ||,$ 

where

(5.48) 
$$K(r) = (1+K)^{r} - 1$$

<u>Remark</u>: In reference to the assumption that  $sgn(c_{kk}^{k}) = sgn(c_{kk}^{k})$ : if it is not the case that we have equality for some k, then recalling (5.15)

.

$$(5.15) \qquad |\hat{c}_{kk}^{k}|, |c_{kk}^{k}| \leq |\delta c_{kk}^{k}|,$$

let k be the initial natural number such that we do not have equality, then (5.15) is valid, and we have from (5.38)

(5.49) 
$$|\hat{c}_{kk}^{k}|, |c_{kk}^{k}| \leq (1+K)^{k-1} ||\delta C||.$$

This is a necessary condition for the signs to be different. The typical situation is as follows. We have a numerical bound for  $||\delta C||$ , and we are performing the factorization on the perturbed

matrix  $\hat{C}$ . We check at each step, k, of the factorization to see if (5.49) cannot hold for any k, then we can conclude that the factorization of the matrix C would proceed in the same way, and  $sgn(\hat{c}_{kk}^{k}) =$  $sgn(c_{kk}^{k})$  for  $1 \le k \le r$ .

<u>Remark</u>: On pivoting. The usual pivotal strategy used is "maximal row pivoting", see Golub [11]. Other strategies of various types have been tried with various degrees of computational success, see Jennings and Osborn [25]. Usually, no theoretical justification is given as to how the computational errors are affected by the choice of pivotal strategy.

We are able to give here a justification of sorts as to the desirability of "maximal row pivoting". Assume (5.39) holds, then we have from (5.34)

(5.50) 
$$||\delta C_{i+1}^{i+1}|| \leq \sum_{j=1}^{i} (K \mu_{i+1,j}) ||\delta C_{j}^{j}|| + ||\delta C_{i+j}^{1}||$$

Intuitively, it is seen by examining (5.50) that, by repeated back substitution, we obtain an expression which is the sum of products, each of which is made up of repeated factors of the type  $(K\mu_{st})$ . Because "maximal row pivoting" gives  $\mu_{st} \leq 1$ , it helps reduce the effect of the power of k in the error analysing giving, in general, a better bound then a strategy which allows the  $\mu_{st}$  to be larger than unity. Section 6. Basic Relations in the Perturbational Problem

Using the notation of Sections 1 through 5, we will establish the following relations, which will be used repeatedly.

<u>Lemma 6.1</u> . (6.1) $CC^{T} = LL^{T}$ , $\hat{C}\hat{C}^{T} = \hat{L}\hat{L}^{T}$ , $B_{2}^{T}B_{2} = W^{T}W$ and $\hat{B}^{T}\hat{B} = \hat{W}^{T}\hat{W}$
(6.2) $L^{-}C = Q_1 \text{ and } B_2 W^{-} = V_1$ .
(6.3) $\chi$ (L) = $\chi$ (C).
(6.4) $\chi(W) = \chi(B_2)$ .
If $  L^{-1} \delta L   \leq \Delta \leq 1$ , then
(6.5a) $  L   \leq e_1   L  $ , where $e_1 = (1 + \Delta)$ ,
(6.5b) $  _{L}^{\Lambda-1}   \leq e_{2}   _{L}^{-1}  $ , where $e_{2} = 1/(1 - \Delta)$ ,
(6.5c) $\chi(\hat{L}) \leq e_3 \chi(C)$ , where $e_3 = (1 + \Delta)/(1 - \Delta)$ .
If $   \delta W W^{-1}    \leq \Delta < 1$ , then
(6.6a) $  \hat{W}   \leq e_1   W  $ , where $e_1 = (1 + \Delta)$ ,
(6.6b) $  \hat{W}^{-1}   \leq e_2   W^{-1}  $ , where $e_2 = 1/(1 - \Delta)$ ,
(6.6c) $\chi(\hat{W}) \leq e_3 \chi(B_2)$ , where $e_3 = (1 + \Delta)/(1 - \Delta)$ .
If $\sqrt{r/2}   L^{-1}\delta L   \leq \Delta < 1$ ,
then

- (6.7a)  $||L^{-1}\delta L|| \leq K_1(r) \chi(C) ||\xi C|| / ||C||$ ,
- (6.7b)  $|| \delta Q_1 || \leq K_2(\mathbf{r}) \chi(C) || \delta C || /|| C ||,$

(6.7c) 
$$||\delta L|| \leq K_{3}(r) \mathcal{N}(C) ||\delta C||$$
,  
where  $K_{1}(r) = \sqrt{r} (2 + \Delta)/(\sqrt{2}(1 - \Delta))$ ,  
 $K_{2}(r) = 1 + K_{1}(r)$ , and  $K_{3}(r) = K_{2}(r) + 1$ .  
If  $\sqrt{(n-r)/2} ||\delta W W^{-1}|| \leq \Delta < 1$ ,

then

(6.8a) 
$$|| \delta W W^{-1} || \leq K_1(n-r) \chi(B_2) || \delta B_2 || / || B_2 ||,$$

(6.8b) 
$$|| \delta V_1 || \leq K_2(n-r) \chi(B_2) || \delta B_2 || /| |B_2||$$
,

(6.8c) 
$$|| \delta W || \leq K_3^{(n-r)} \chi(B_2) || \delta B_2^{(n-r)} || \delta B_2^{(n-r)}$$

where  $K_1(n-r)$ ,  $K_2(n-r)$  and  $K_3(n-r)$  are defined in (6.7). Under the hypothesis of Lemma 5.1 we have

(6.9a) 
$$||\delta B_{1}||, ||\delta B_{2}||, ||[\delta B_{1}|\delta B_{2}]|| \leq (||\delta A||/||A|| + K(r) \chi(C) ||\delta C||/||C||)||A||,$$

$$(6.9b) \qquad ||\delta B_{1}|| \leq (||\delta A||/||A|| + ||\delta Q_{1})||A||.$$
If  $\mathbf{x} = Q \cdot \begin{bmatrix} d_{1} \\ d_{2} \end{bmatrix}$ , where  $\mathbf{x} \in \mathbb{R}^{n}$ ,  $Q \in \mathbb{R}^{n \times n}$  is orthogonal,  $d_{1} \in \mathbb{R}^{n1}$ ,  
 $d_{2} \in \mathbb{R}^{n2}$  and  $n_{1} + n_{2} = n$ , then  

$$(6.10) \qquad ||d_{1}||, ||d_{2}|| \leq ||\mathbf{x}||.$$

<u>Proof</u> (6.1): From (1.15a) we have  $CC^{T} = [L|0]Q^{t}Q[\frac{L^{T}}{0}] = LL^{T}$ ,

and from (1.32a) we have

$$\mathbf{B}_{2}^{\mathrm{T}}\mathbf{B}_{2} = \left[\mathbf{W}^{\mathrm{T}} \mid \mathbf{0}\right] \mathbf{V} \mathbf{V}^{\mathrm{T}} \left[\frac{\mathbf{W}}{\mathbf{0}}\right] = \mathbf{W}^{\mathrm{T}} \mathbf{W}$$

(6.2): From (1.15a) we have  $L^{-1}C = L^{-1}[L|0]Q^{T} = Q_{1}^{T}$ , where the proof for W is similar.

(6.3): If we use (6.1), (3.2a) and (3.2b), then (6.3) follows from Definition 4.1.

(6.4): Same as (6.3).

(6.5a): Because  $\hat{L} = L (I + L^{-1} \delta L)$  is valid, the result follows by taking norms.

(6.5b): Because  $\hat{L}^{-1} = (I + L^{-1}\delta L)^{-1} L^{-1}$  is valid, the result follows by taking norms and using (3.16).

(6.5c): This follows from (6.5a), (6.5b) and (6.3).

(6.6a)-(6.6c): These are proved in a similar way to (6.5a) through (6.5c).

(6.7a): This is a result of Jennings and Osborne [25, p. 327] inequality (2.6) of their paper.

(6.7b): From (1.15a) we have  $\delta C = \hat{L} \hat{Q}_1 - L Q_1^T = \delta L \hat{Q}_1^T + L \delta Q_1^T$  which gives  $L \delta Q_1^T = \delta C - \delta L \hat{Q}_1^T$ . Premultiplying by  $L^{-1}$  we have  $\delta Q_1^T = L^{-1} \delta C - L^{-1} \delta L \hat{Q}_1^T$ . Taking norms and using (3.5), (3.11), (3.13) and (3.7) we obtain  $||\delta Q_1|| \leq ||L^{-1}|| ||\delta C|| + ||L^{-1}\delta L||$ . After using (6.1), (6.3) and (6.7a) our result follows.

(6.7c): Using (1.15b) and taking norms we have  $||\delta L|| \leq ||\delta C|| + ||C|| ||\delta Q_1||$ , where we have used (3.11), (3.13) and (3.7). If we use (6.7b) and the fact that  $\chi(C) \geq 1$  we have our result.

(6.8a)-(6.8c): These are proved in almost an identical way as (6.7a) through (6.7c).

i w

101

(6.9a): From (1.27) we have  $\left[\delta B_1 \mid \delta B_2\right] = \delta A \hat{Q} + A \delta Q$ .

36

If we apply (3.13), (3.9) and Lemma 5.1 we have the result.

(6.9b):  $\delta B_1 = \delta A \hat{Q}_1 + A \delta Q_1$ .

(6.10): Use (3.8) and the definition of the Euclidean vector norm on  $\begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = Q^T x$ .

#### Section 7. Effects of Perturbational Errors

In this section we will show the effect of perturbational errors on a solution given by the method of Algorithm 1.1. In order to obtain an expression which is compact, it will be necessary to merge all terms of higher powers into those of the first power. These terms represent relative errors, which, if the perturbations and condition numbers are reasonable, should be very much smaller than unity. Therefore, it is not unreasonable to make the assumption that they are bounded away from unity. This is the reason for making assumptions (7.1a) through (7.1e). It should be clear, that even terms (7.1a) and (7.1b) represent relative errors, because  $||L^{-1}\delta L|| \leq$  $\chi(L) || \delta L || / || L ||$ . Also, (7.1c) should be reasonable, because in the unconstrained case it is always true that  $||\mathbf{R}||/||\mathbf{f}|| \leq 1$ , and we would hope that in a practical problem, f would be very near the space spanned by the columns of A, which would yield a very small relative error for the residual.

<u>Theorem 7.1</u>. The notation of Sections 1 through 6 is used. If there exists a real non-negative  $\Delta < 1$ , such that

- (7.1a)  $h_1 || L^{-1} \delta L|| \leq \Delta$ , where  $h_1 = \max \{1, \sqrt{r/2}\}$ ,
- (7.1b)  $h_2 || \delta W W^{-1} || \leq \Delta$ , where  $h_2 = \max \{1, \sqrt{(n-r)/2}\}$ ,
- (7.1c)  $||\mathbf{R}||/||\mathbf{f}|| \leq \Delta, \mathbf{f} \neq 0$
- (7.1d)  $|| \delta \Omega_1 || \leq K_2(\mathbf{r}) \chi(\mathbf{C}) || \delta \mathbf{C} || /|| \mathbf{C} || \leq \Delta$ ,
- $(7.1c) \qquad ||\delta A||/||A|| \leq \Delta,$

37

where C,  $\hat{C}$ ,  $\hat{B}_2$ ,  $\hat{B}_2$  are of full rank and the hypothesis of Lemma 5.1 is satisfied then

$$(7.2) \qquad ||\hat{\mathbf{x}} - \mathbf{x}|| / ||\mathbf{x}|| \leq (\mathbf{e}_{2} || \delta \mathbf{g} || / || \mathbf{g} || + \mathbf{K}_{5}(\mathbf{r}) || \delta \mathbf{C} || / || \mathbf{C} ||) \chi(\mathbf{C}) + (\mathbf{e}_{2}^{2} || \delta \mathbf{f} || / || \mathbf{f} || + 2\mathbf{e}_{2} || \delta \mathbf{A} || / || \mathbf{A} ||) \cdot (|| \mathbf{A} || / || \mathbf{B}_{2} ||) \chi(\mathbf{B}_{2}) + (\mathbf{e}_{2}^{2} \mathbf{e}_{4} || \delta \mathbf{g} || / || \mathbf{g} || + \mathbf{e}_{2} \mathbf{K}_{4}(\mathbf{r}) || \delta \mathbf{C} || / || \mathbf{C} ||) (|| \mathbf{A} || / || \mathbf{B}_{2} ||) \chi(\mathbf{C}) \chi(\mathbf{B}_{2}) + \mathbf{e}_{2}^{2} \mathbf{K}_{2}(\mathbf{n} - \mathbf{r}) (|| \mathbf{A} || / || \mathbf{B}_{2} ||)^{2} (|| \delta \mathbf{A} || / || \mathbf{A} || + \mathbf{K}(\mathbf{r}) \chi(\mathbf{C}) || \delta \mathbf{C} || / || \mathbf{C} ||) \chi^{2} (\mathbf{B}_{2}) || \mathbf{R} || / || \mathbf{f} || ,$$

where the constants  $e_2$ ,  $e_4$ , K(r),  $K_2(n-r)$ ,  $K_4(r)$  and  $K_5(r)$  are given by (6.5b), (7.21), (5.48), (6.7), (7.23) and (7.26) respectively.

<u>Proof</u>: We have from Algorithm 1.1 that solutions x and  $\hat{x}$  exist to both problems, where x and  $\hat{x}$  are given by

(7.3a) 
$$\mathbf{x} = \mathbf{Q} \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \end{bmatrix} = \mathbf{Q} \begin{bmatrix} \mathbf{L}^{-1} \mathbf{g} \\ \mathbf{W}^{-1} \mathbf{V}_1 \{ \mathbf{f} - \mathbf{B}_1 \mathbf{L}^{-1} \mathbf{g} \} \end{bmatrix}$$

(7.3b) 
$$\hat{\mathbf{x}} = \hat{\mathbf{Q}} \begin{bmatrix} \hat{\mathbf{d}}_1 \\ \hat{\mathbf{d}}_2 \end{bmatrix} = \hat{\mathbf{Q}} \begin{bmatrix} \hat{\mathbf{L}}^{-1} \hat{\mathbf{g}} \\ \\ \hat{\mathbf{w}}^{-1} \hat{\mathbf{v}}_1 \{ \hat{\mathbf{f}} - \hat{\mathbf{B}}_1 \hat{\mathbf{L}}^{-1} \hat{\mathbf{g}} \} \end{bmatrix}$$

respectively. Subtracting x from  $\hat{\mathbf{x}}$  we obtain

(7.4) 
$$\hat{\mathbf{x}} - \mathbf{x} = \hat{\mathbf{Q}} \begin{bmatrix} \delta \mathbf{d}_1 \\ \delta \mathbf{d}_2 \end{bmatrix} + \delta \mathbf{Q} \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \end{bmatrix}$$

where  $\delta d_1 = \hat{d}_1 - d_1$  and  $\delta d_2 = \hat{d}_2 - d_2$ . Using (6.10), (7.3a) and (7.4), we find after taking norms

(7.5) 
$$||\hat{\mathbf{x}}-\mathbf{x}|| \leq ||\delta d_1|| + ||\delta d_2|| + ||\delta Q|| ||\mathbf{x}||.$$

Lemma 5.1 gives an upper bound for  $||\delta Q||$ 

(7.6) 
$$|| \delta Q || \leq K(r) \chi(C) || \delta C || / || C ||$$
.

Therefore, all we have to find are upper bounds for  $||\delta d_1||$  and  $||\delta d_2||$ . Using (7.3a) the expression for  $\delta d_1$  is

(7.7) 
$$\delta d_1 = \hat{L}^{-1} \hat{g} - L^{-1} g$$
  
=  $\hat{L}^{-1} (\delta g - \delta L d_1),$ 

where we have made use of (3.17) and the definition of  $d_1$  to obtain our simplification. Because  $\hat{L}^{-1} = (I + L^{-1} \delta L)^{-1} L^{-1}$ , we obtain the following upon taking norms

(7.8) 
$$||\delta d_1|| \leq e_2 [\chi(C)||\delta g||/(||C||||x||)+||L^{-1}\delta L||] ||x||,$$

where we have made use of (3.16), (6.3) and (6.10). From (1.4) we have that  $||C|| ||x|| \ge ||g||$ , and substituting this and (6.7a) into (7.8) we have

(7.9) 
$$||\delta d_1|| \leq e_2 [||\delta g||/||g|| + K_1(r)||\delta C||/||C||] \chi(C)||x||.$$

,

Before estimating  $||\delta d_2||$ , define the vector  $\xi \in \mathbb{R}^m$  by

(7.10) 
$$\xi = \{ \hat{f} - \hat{B}_1 \hat{d}_1 \} - \{ f - B_1 d_1 \}$$
$$= \delta f - \hat{B}_1 \delta d_1 - \delta B_1 d_1 .$$

If we make use of the equalities (1.33), (3.17) and  $\hat{V}_1 = V_1 + \delta V_1$ , we obtain the following simplified expression for  $\delta d_2$ 

$$(7.11) \qquad \delta d_{2} = \hat{W}^{-1} \hat{V}_{1} \{\hat{f} - \hat{B}_{1} \hat{d}_{1}\} - W^{-1} V_{1} \{f - B_{1} d_{1}\} \\ = \left[\hat{W}^{-1} \hat{V}_{1} - W^{-1} V_{1}\right] \{f - B_{1} d_{1}\} + \hat{W}^{-1} \hat{V}_{1} \xi \\ = \hat{W}^{-1} \left[\delta V_{1} \{f - B_{1} d_{1}\} - \delta W d_{2} + \hat{V}_{1} \xi\right].$$

In the last expression we have made a simplification by using the definition of d<sub>2</sub> given in (1.33). From (1.28) we have  $B_2 d_2 - R = f - B_1 d_1$ . Substituting this into (7.11) we find (7.12)  $\delta d_2 = \hat{W}^{-1} \left[ \delta V_1 B_2 d_2 - \delta V_1 R - \delta W d_2 + \hat{V}_1 \xi \right]$  $= \hat{W}^{-1} \left[ (\delta V_1 B_2 - \delta W) d_2 - \delta V_1 R + \hat{V}_1 \xi \right].$ 

Using (1.32b) we see that

(7.13) 
$$\delta W = \hat{V}_1 \hat{B}_2 - V_1 \hat{B}_2$$
$$= \hat{V}_1 \delta \hat{B}_2 + \delta \hat{V}_1 \hat{B}_2$$

If (7.13) is substituted into (7.12), we find after the cancellation of  $\delta V_1 B_2$ 

(7.14) 
$$\delta d_2 = \hat{W}^{-1} \left[ -\hat{V}_1 \delta B_2 d_2 - \delta V_1 R + \hat{V}_1 \xi \right]$$

If substitution of (7.10) into (7.14) is made, we have after regrouping

(7.15) 
$$\delta d_{2} = \hat{W}^{-1} \hat{V}_{1} \left[ \delta f - \hat{B}_{1} \delta d_{1} - \delta B_{1} d_{1} - \delta B_{2} d_{2} \right] - \hat{W}^{-1} \delta V_{1} R.$$

Taking norms in (7.15), and using (6.6b) and (6.10) we have

$$(7.16) \qquad ||\delta d_2|| \leq e_2 \Big[ ||\delta f||/(||A|| ||x||) + (||B_1||/||A|| + ||\delta B_1||/||A|| + ||\delta B_1||/||A|| + ||\delta B_2||/||A|| \Big] ||\delta d_1||/||x|| + ||\delta B_1||/||A|| + ||\delta B_2||/||A|| \Big] ||A|| ||w^{-1}|| ||x|| + e_2 ||w^{-1}|| ||\delta V_1|| ||R||.$$

From (1.27), by using (3.9) and (3.13), we have that

$$(7.17) \qquad ||B_1||, ||B_2|| \leq ||A||.$$

From (1.3) and our hypothesis (7.1c) it follows

(7.18)  $||A|| ||x|| \ge ||f+R||$  $\ge ||f|| - ||R||$  $\ge ||f||/e_2 > 0$ ,

where  $e_2 = 1/(1 - \Delta)$ .

Using inequalities (7.17) and (7.18) to simplify (7.16) we obtain

$$(7.19) \qquad || \delta d_2 || \leq e_2 \Big[ e_2 || \delta f || / || f || + (1 + || \delta B_1 || / || A ||) \\ \cdot || \delta d_1 || / || x || + || \delta B_1 || / || A || \\ + || \delta B_2 || / || A || \Big] (|| A || / || B_2 ||) \chi(B_2) || x || \\ + e_2 || W^{-1} || || \delta V_1 || || R ||,$$

where (6.4) has been used.

From (6.9b), and our hypotheses (7.1d) and (7.1e) it follows

$$(7.20) \qquad ||_{\delta B}|| \leq 2 \Delta ||A||,$$

thereby simplifying the coefficient of  $||\delta d_1||/||x||$  to be

(7.21) 
$$e_4 = 1 + 2\Delta$$
.

Using (6.9a) as an upper bound for  $||\delta B_1||$  and  $||\delta B_2||$ , and substituting (7.9) for  $||\delta d_1||$ , we obtain after regrouping (7.22) $||\delta d_2|| \le e_2 [e_2 ||\delta f||/||f|| + 2||\delta A||/||A||$  $+(e_2 e_4 ||\delta g||/||g|| + K_4(r) ||\delta C||/||C||)\chi(C)]$  $\cdot (||A||/||B_2||\chi(B_2)||x|| + e_2 ||W^{-1}|| ||\delta V_1|| ||R||,$ 

where

(7.23) 
$$K_4(r) = e_2 e_4 K_1(r) + 2 K(r)$$
,

and  $K_1(r) = \sqrt{r} (2 + \Delta) / (\sqrt{2}(1 - \Delta))$ .

We will now obtain an upper bound for the coefficient of  $||\mathbf{R}||$ , by applying the estimates given in (6.4), (6.8b), (6.9a) and (7.18)

$$(7.24) \quad e_{2} ||W^{-1}|| ||\delta V_{1}|| ||R|| \\ \leq e_{2} K_{2}(n-r) \chi^{2}(B_{2})||R|| ||\delta B_{2}||/||B_{2}||^{2} \\ \leq e_{2}^{2} K_{2}(n-r) [||\delta A||/||A|| + K(r) \chi(C)||\delta C||/||C||] \\ \cdot \chi^{2}(B_{2})(||A||/||B_{2}||)^{2}(||R||/||f||) ||x||.$$

If (7.6), (7.9), (7.22) and (7.24) are substituted into (7.5) and

regrouped according to condition numbers we obtain

$$(7.25) \qquad ||\hat{\mathbf{x}} - \mathbf{x}||/||\mathbf{x}|| \leq (\mathbf{e}_{2} || \delta \mathbf{g} ||/||\mathbf{g} || \\ + \mathbf{K}_{5}(\mathbf{r})|| \delta \mathbf{C} ||/||\mathbf{C}||) \chi(\mathbf{C}) + (\mathbf{e}_{2}^{2} || \delta \mathbf{f} ||/||\mathbf{f} || \\ + 2\mathbf{e}_{2} || \delta \mathbf{A} ||/||\mathbf{A} ||)(||\mathbf{A} ||/||\mathbf{B}_{2} ||) \chi(\mathbf{B}_{2}) \\ + (\mathbf{e}_{2}^{2} \mathbf{e}_{4} || \delta \mathbf{g} ||/||\mathbf{g} || \\ + \mathbf{e}_{2} \mathbf{K}_{4}(\mathbf{r}) || \delta \mathbf{C} ||/||\mathbf{C} ||)(||\mathbf{A} ||/||\mathbf{B}_{2} ||) \chi(\mathbf{C}) \chi(\mathbf{B}_{2}) \\ + \mathbf{e}_{2}^{2} \mathbf{K}_{4}(\mathbf{r}) || \delta \mathbf{C} ||/||\mathbf{C} ||)(||\mathbf{A} ||/||\mathbf{B}_{2} ||) \chi(\mathbf{C}) \chi(\mathbf{B}_{2}) \\ + \mathbf{e}_{2}^{2} \mathbf{K}_{2}(\mathbf{n} - \mathbf{r})(||\mathbf{A} ||/||\mathbf{B}_{2} ||)^{2}(|| \delta \mathbf{A} ||/||\mathbf{A} || \\ + \mathbf{K}(\mathbf{r}) \chi(\mathbf{C}) || \mathbf{S} \mathbf{C} ||/||\mathbf{C} ||) \chi^{2}(\mathbf{B}_{2}) ||\mathbf{R} ||/||\mathbf{f} ||,$$

where

(7.26) 
$$K_5(r) = e_2 K_1(r) + K(r)$$
.

We have shown that the condition numbers appear only linearly in (7.25), except for the coefficient of  $||\mathbf{R}||/||\mathbf{f}||$ . The question could be asked, if whether the term  $\chi^2(\mathbf{B}_2)$  is reasonable, or if it might be the fault of the error analysis. Van der Sluis [31] has shown, using a geometrical argument, that squaring of the condition number in the unconstrained case can indeed be realized. If  $\chi(\mathbf{B}_2) ||\mathbf{R}||/||\mathbf{f}||$  is not too large, however, the effect of  $\chi^2(\mathbf{B}_2)$ can be minimized, and we would expect the condition numbers to appear only linearly in the error.

Also, the full effect of a condition number is seldom or possibly never realized, see Wilkinson [34].

Finally, because A and  $B_2$  are calculated during our solution of (1.1) to (1.4), we have the bound

۰.

.

$$||A||/||B_2|| \leq \sqrt{n-r} ||A||_E/||B_2||_E$$
,

which is easy to calculate.

#### Section 8. Rounding Errors in Computation

We will state here for later reference some known results on computational errors as developed by Wilkinson [36], [35]. For our purposes, we will limit ourselves strictly to floating point computation.

To develop our results, we introduce the very general notation  $fl(\cdot op \cdot)$ , which is used by Wilkinson [36] and others. The "fl" notation means, that if we have two quantities A and B, which could be numbers, vectors, matrices, etc., and a corresponding operator "op" such as scalar addition, inner product, matrix multiplication, etc., then fl (A op B) is that quantity which would result in the computer by performing that operation by some computer program using floating point arithmetic.

A floating point binary number x consists of two parts, the binary exponent b and the binary mantissa a, where  $-\frac{1}{2} \ge a \ge -1$  or  $\frac{1}{2} \le a \le 1$ . Then x is expressed as  $x = a \cdot 2^b$ . The computer memory allocated for x is limited to a certain number of digits, which are divided up in some way between a and b. We will denote by t the number of binary digits allocated to the mantissa. Also, we will assume that the number of digits allotted to the exponent is large enough to accommodate all of our calculations, and will not concern ourselves here with the problem of exponent under- or over-flow. Finally, we will assume that the computer with which we are working has a double-precision accumulator. This means when single precision floating point numbers are retrieved from the memory to perform certain arithmetic operations, they are allocated a double precision mantissa before the operation is performed. The operation is carried out in double-precision, and the result is not rounded to single-precision until it is sent back to the memory.

Under these assumptions Wilkinson [36] has shown the following results.

## Section 8.1. Vector Addition

Given two real numbers a and b, there exists a real  $\mathcal{E}$ , such that

(8.1)  $f\ell (a + b) = (a + b) (1 + \xi)$ ,

where

 $(8.2) \qquad |\mathcal{E}| \leq 2^{-t}.$ 

Therefore, if we have two vectors x,  $y \in \mathbb{R}^n$ ,

(8.3) 
$$||fl(x + y) - (x + y)|| \leq 2^{-t} ||x + y||$$

### Section 8.2. Matrix Multiplication

Given a matrix  $A \in \mathbb{R}^{m \times n}$  and a vector  $b \in \mathbb{R}^{n}$ , let E denote the computational error in calculating  $A \cdot b$ 

$$(8.4) E = f \ell (A \cdot b) - A \cdot b,$$

where the entries of E are represented by  $e_i$  for  $1 \le i \le m$ .

Wilkinson [36], p. 83 has shown that under the assumption  $n 2^{-t} < .1$ 

(8.5) 
$$|\mathbf{e}_{i}| \leq 1.06 \cdot 2^{-t} \Big[ n |\mathbf{a}_{i1}| |\mathbf{b}_{1}| + \sum_{k=2}^{n} (n - (k - 2)) |\mathbf{a}_{ik}| |\mathbf{b}_{k}| \Big]$$
  
 $\leq 1.06 \cdot 2^{-t} n \sum_{k=1}^{n} |\mathbf{a}_{ik}| |\mathbf{b}_{k}|$   
 $\leq 1.06 \cdot 2^{-t} n (\sum_{k=1}^{n} a_{ik}^{2})^{\frac{1}{2}} (\sum_{k=1}^{n} b_{k}^{2})^{\frac{1}{2}} ,$ 

where the last inequality is obtained by using Schwarz's inequality. Therefore, it follows from the definition of  $||\cdot||_{E}$  and (3.4) that (8.6)  $||E|| \leq 1.06 \cdot 2^{-t} n ||A||_{E} ||b||$  $\leq 2^{-t} \overline{K_{1}}(n) ||A|| ||b||$ ,

where

(8.7)  $\overline{K_1}(n) = 1.06 n^{3/2}$ 

# Section 8.3. Solution of a Triangular Set of Equations

Let  $L \in \mathbb{R}^{n \times n}$  be an upper- or lower-triangular non-singular matrix. We now wish to know the computational error introduced by solving the linear equation

(8.8) 
$$Lx = b$$
, where  $x, b \in \mathbb{R}^n$ ,

for the vector x. Without loss of generality we will assume here that L is lower triangular. Wilkinson [36] has shown that the entries  $x_i$  of x are calculated sequentially from i=1 to i=n by using

(8.9) 
$$\mathbf{x}_{\mathbf{r}} = f\ell((-\ell_{\mathbf{r}1} \mathbf{x}_{1} - \ell_{\mathbf{r}2} \mathbf{x}_{2} - \cdots - \ell_{\mathbf{r}, \mathbf{r}-1} \mathbf{x}_{\mathbf{r}-1} + \mathbf{b}_{\mathbf{r}})/\ell_{\mathbf{r}\mathbf{r}}),$$

for  $1 \leq r \leq n$ . Note that the vector x which we have constructed in (8.9) is in general not the solution of (8.8), but rather, it is an exact solution of the perturbed problem

(8.10) 
$$(L + \delta L) x = b$$
,

where the lower triangular matrix  $\delta L$  is bounded by [36, p. 103]

(8.11) 
$$||\delta L|| \leq ||\delta L||_{E}$$
  
 $\leq 1.06 \ 2^{-t}(1 + 1.06 \cdot 2^{-t} \cdot 3(n+2)/2) ||L||_{E}$   
 $\leq 2^{-t} \overline{K_{2}}(n) ||L||,$ 

wi th

(8.12) 
$$\overline{K_2}(n) = 1.06 \sqrt{n} (1 + 1.06 2^{-t} 3(n + 2)/2),$$

and the value of n restricted sufficiently to make the term 1.06  $2^{-t}3(n+2)/2 \ll 1$ . Also, if  $L^{-1}$  exists, then Wilkinson [36] has shown that  $(L + \delta L)^{-1}$  also exists.

## Section 8.4. Orthogonal Transformation.

An elementary orthogonal transformation  $P \in \mathbb{R}^{n \times n}$  has the form  $P = I - 2ww^{T}$ , where  $w \in \mathbb{R}^{n}$  and ||w|| = 1. Given the vector w, we will denote the matrix "P", calculated using floating point arithmetic, by the symbol  $\overleftarrow{P}$ , (which may no longer be orthogonal), i.e.

$$(8.13) \qquad \overline{\mathbf{P}} = \mathbf{f} \boldsymbol{\ell}(\mathbf{P}) ,$$

where it is assumed that  $n2^{-t} < .1$ . Wilkinson [35] has shown that if  $A \in \mathbb{R}^{n\times m}$  is any matrix, premultiplication by P using floating point arithmetic gives

(8.14) 
$$||f\ell(\overline{P}A) - PA|| \leq 2^{-t} \beta \sqrt{m} ||A||,$$

where

$$(8.15) \quad \beta = 12.36.$$

Also if A is premultiplied sequentially by orthogonal transformations  $P_1, P_2, \dots, P_m$ , and for  $1 \le i \le m$  we define

$$(8.16a) \qquad \overline{P_i} = fl (P_i)$$

(8.16b) 
$$\overline{A}_{i+1} = f\ell (\overline{P}_i \cdot \overline{A}_i), \text{ where } \overline{A}_1 = A.$$

Then there exists a  $\delta A \in \mathbb{R}^{n \times m}$  such that

(8.17) 
$$\overline{\mathbf{A}}_{m+1} = \mathbf{P}_m \cdot \mathbf{P}_{m-1} \cdot \cdots \cdot \mathbf{P}_1 (\mathbf{A} + \boldsymbol{\delta}\mathbf{A}),$$

wi th

(8.18) 
$$||_{\delta A}|| \leq ||_{\delta A}||_{E}$$
  
 $\leq m 2^{-t} \beta (1 + 2^{-t} \beta)^{m-1} ||A||_{E}$   
 $\leq 2^{-t} \overline{K_{3}}(m) ||A||,$ 

where  $\beta$  is given in (8.15) and  $\overline{K_3}(m)$  is defined by

•

(8.19) 
$$\overline{K}_{3}(m) = \beta m^{3/2} (1 + 2^{-t}\beta)^{m-1}$$

The results of (8.14) through (8.19) hold also for post multi-

plication by orthogonal transformations. This is true because

(8.20) B P = (P B<sup>T</sup>)<sup>T</sup>,

where  $B \in \mathbb{R}^{m \times n}$  and  $P^{T} = P$ , and our norms are invariant under transposition.

## Section 9. Rounding Error Analysis

We now wish to analyze the effect of computational errors on our least squares solution with constraints. We will show that most of the error can be accounted for by perturbational induced errors in the initial problem.

Algorithm 1.1 gives the exact solution (1.1) to (1.4) to be

(9.1) 
$$\mathbf{x} = \mathbf{Q} \begin{bmatrix} \mathbf{L}^{-1} \mathbf{g} \\ \mathbf{w}^{-1} \mathbf{v}_{1} \{\mathbf{f} - \mathbf{B}_{1} \mathbf{L}^{-1} \mathbf{g} \} \end{bmatrix}$$

We will now define stepwise the order of computation defined in (9.1). This will determine the effect of the errors.

<u>Step 1</u>: <u>Reduce C by orthogonal transformations to obtain L.</u> Define for  $1 \leq i \leq r$ 

(9.2) 
$$\overline{C}_{i+1} = fl(\overline{C}_i \overline{P}_i)$$
, where  $\overline{C}_1 = C$  and  $[\overline{L}|0] = \overline{C}_{r+1}$ ,

and

(9.3)  $\overline{\mathbf{P}}_{\mathbf{i}} = \mathbf{fl} (\hat{\mathbf{P}}_{\mathbf{i}}),$ 

where  $\overrightarrow{P}_i \in \mathbb{R}^{n \times n}$  is that elementary orthogonal matrix which exactly reduces row i of  $\overrightarrow{C}_i$ . Define the orthogonal matrix  $\hat{Q}$  by

(9.4) 
$$\hat{Q} = \hat{P}_1 \cdot \hat{P}_2 \cdot \cdots \cdot \hat{P}_r$$

From (8.17) and (8.18) there exists a perturbation SC in C such that

(9.5) 
$$\left[\overline{\mathbf{L}}\right] = (\mathbf{C} + \delta \mathbf{C}) \hat{\mathbf{Q}},$$

where

(9.6a) 
$$|\delta C| \leq 2^{-t} \overline{K}_{3}(r) |C|$$
, or

(9.6b) 
$$||\delta C|| = O(2^{-t} ||C||).$$

Step 2: Calculate  $\overline{L}^{-1}g$ . From (8.10) and (8.11) there exists a perturbation  $\delta L^*$  in  $\overline{L}$  such that

(9.7) 
$$f\ell (\overline{L}^{-1} g) = (\overline{L} + \delta L^*)^{-1} g_{\ell}$$

where  $(\overline{L} + \delta L^*)^{-1}$  exists because  $\overline{L}^{-1}$  does and

- (9.8)  $||\delta L^*|| \leq 2^{-t} \overline{K}_2(r) ||\overline{L}||.$
- <u>Step 3:</u> <u>Calculate  $\overline{B}_1$  and  $\overline{B}_2$ . Define for  $1 \le i \le r$ </u>

(9.9) 
$$\overline{A}_{i+1} = f\ell (\overline{A}_i \overline{P}_i) \text{ where } \overline{A}_1 = A$$
,

(9.10) 
$$\left[\overline{B}_{1} | \overline{B}_{2}\right] = \overline{A}_{r+1}$$
,

where  $\overline{P_i}$  is defined by (9.3). From (8.17) and (8.18) there exists a matrix  $\delta A^* \in \mathbb{R}^{m \times n}$  such that

(9.11) 
$$\left[\overline{B}_{1}|\overline{B}_{2}\right] = (A + \delta A^{*})\hat{Q},$$

where

(9.12a) 
$$||\delta A^*|| \leq 2^{-t} \overline{K}_3(n) ||A||, \text{ or}$$

(9.12b) 
$$||\delta A^*|| = O(2^{-t} ||A||).$$

Step 4: Calculate 
$$\overline{B}_{1}(\overline{L} + \delta L^{*})^{-1} g$$
.  
Define  $\mathcal{E}_{1} \in \mathbb{R}^{m}$  by  
(9.13)  $fl(\overline{B}_{1} \cdot (\overline{L} + \delta L^{*})^{-1} g) = \overline{B}_{1}(\overline{L} + \delta L^{*})^{-1} g + \mathcal{E}_{1}$ .

A bound for  $\xi_1$  is given in (8.6)

(9.14) 
$$||\xi_1|| \leq 2^{-t} \overline{K}_1(r) ||\overline{B}_1|| ||(\overline{L} + \delta L^*)^{-1} g||.$$

Step 5: Calculate 
$$f\ell(f - [\overline{B}_1(\overline{L} + \delta L^*)^{-1} g + \xi_1])$$
.  
Define  $\xi_2 \in \mathbb{R}^m$  by  
(9.15)  $f\ell(f - [\overline{B}_1(\overline{L} + \delta L^*)^{-1} g + \xi_1])$   
 $= f - (\overline{B}_1(\overline{L} + \delta L^*)^{-1} g + \xi_1) + \xi_2$ .

A bound for  $\mathcal{E}_2$  is given by (8.3)

(9.16) 
$$||\mathcal{E}_{2}|| \leq 2^{-t} ||f - (\overline{B}_{1}(\overline{L} + \delta L^{*})^{-1}g + \mathcal{E}_{1})||$$

<u>Step 6</u>: <u>Calculate  $\overline{W}$  by reducing  $\overline{B}_2$  with elementary orthogonal</u> <u>transformations</u>.

Define for  $1 \leq i \leq n-r$ 

(9.17) 
$$\overline{B}_{2,i+r} = f\ell(\overline{U}_i \overline{B}_{2,i}) \text{ and } \overline{B}_{2,1} = B_2$$
,

where

(9.18)  $\overline{U}_i = f\ell(\hat{U}_i)$ .

 $\hat{U}_i \in \mathbb{R}^{m \times m}$  is that elementary orthogonal transformation which exactly reduces  $\overline{B}_{2,i}$ . Define the orthogonal matrix  $\hat{V}$  by

$$(9.19) \qquad \hat{\mathbf{v}} = \hat{\mathbf{U}}_{\mathbf{n}-\mathbf{r}} \cdot \cdots \cdot \hat{\mathbf{U}}_{\mathbf{l}}$$

From (8.17) and (8.18) there exists a matrix  $\delta B_2^* \in \mathbb{R}^{mx(n-r)}$  such that

(9.20) 
$$\left[\frac{\overline{W}}{0}\right] = \overline{B}_{2,((n-r)+1)} = \hat{V}(\overline{B}_{2} + \delta B_{2}^{*})$$

and

$$(9.21) \qquad || \delta B_2^* || \leq 2^{-t} \overline{K}_3(n-r) || \overline{B}_2 ||,$$

where  $\overline{W} \in \mathbb{R}^{(n-r)x(n-r)}$  is an upper triangular matrix.

# Step 7: Premultiplication by $\hat{V}$ .

Define

(9.22) 
$$\overline{z}_1 = \{ f - (\overline{B}_1(\overline{L} + \delta L^*)^{-1} g + \xi_1) + \xi_2 \}$$

and for  $1 \leq i \leq n-r$ 

(9.23)  $\overline{z}_{i+1} = f\ell(\overline{U}_i \overline{z}_i)$ .

From (8.17) and (8.18) there exists a vector  $\delta z \in \mathbb{R}^m$  such that

(9.24) 
$$\overline{z}_{(n-r)+1} = \hat{V}(\overline{z}_1 + \delta z)$$
,

where

(9.25) 
$$||\delta z|| \leq 2^{-t} \overline{K}_{3}(n-r) ||\overline{z}_{1}||.$$

We now define  $y \in \mathbb{R}^{n-r}$  by

(9.26) 
$$\mathbf{y} = \begin{bmatrix} \mathbf{I}_{\mathbf{n}-\mathbf{r}} \mid \mathbf{0} \end{bmatrix} \overline{\mathbf{z}}_{\mathbf{r}+1} = \mathbf{\hat{V}}_{1} (\overline{\mathbf{z}}_{1} + \mathbf{\delta} \mathbf{z}),$$

where

(9.27) 
$$\hat{\mathbf{V}}_{1} = \begin{bmatrix} \mathbf{I}_{n-\mathbf{r}} \mid \mathbf{0} \end{bmatrix} \hat{\mathbf{V}}.$$

It is clear that  $\hat{V}_l$  is the first n-r rows of  $\hat{V}$ .

Step 8: Calculate  $\overline{W}^{-1} y$ .

We have from (8.10) and (8.11), that there exists a perturbation  $\delta W^*$  such that

(9.28) 
$$fl(\overline{W}^{-1} y) = (\overline{W} + \delta W^*)^{-1} y$$
,

where  $(\overline{W} + \delta W^*)^{-1}$  exists because  $\overline{W}^{-1}$  exists and

(9.29)  $||\delta W^*|| \leq 2^{-t} \overline{K}_{2}(n-r)||\overline{W}||.$ 

# Step 9: Premultiply by Q.

Define  $\overline{l}_1 \in \mathbb{R}^n$  from (9.7), (9.28), (9.26) and (9.22) by

(9.30) 
$$\overline{\ell}_{1} = \begin{bmatrix} (\overline{L} + \delta L^{*})^{-1} g \\ (\overline{W} + \delta W^{*})^{-1} \widehat{V}_{1} \{\delta z + f - \overline{B}_{1} (\overline{L} + \delta L^{*})^{-1} g - \xi_{1} + \xi_{2} \} \end{bmatrix}$$

and for  $1 \leq i \leq r$ 

(9.31) 
$$\overline{l}_{i+1} = fl(\overline{P}_{(r+1-i)}, \overline{l}_i)$$
.

Therefore our calculated solution is given by

$$(9.32) \qquad \overline{\mathbf{x}} = \overline{\boldsymbol{\ell}}_{r+1} ,$$

and from (8.17) and (8.18) there exists a vector  $\delta \ell^* \epsilon \mathbb{R}^n$  such that

(9.33) 
$$\overline{\mathbf{x}} = \hat{\mathbf{Q}} (\overline{\boldsymbol{l}}_1 + \delta \boldsymbol{l}^*) ,$$

where

$$(9.34) \qquad ||_{\delta} \ell^* || \leq 2^{-t} \overline{K}_{3}(1) || \overline{\ell}_{1} ||.$$

From (9.33) and (9.30) we have an <u>exact</u> representation of our <u>calcu-lated</u> least squares vector  $\overline{x}$  using floating point arithmetic. We now proceed to estimate how far  $\overline{x}$  is from the true solution x given in (9.1). This will be accomplished in two parts. First we will estimate how far  $\overline{x}$  is from an intermediate vector  $\hat{x}$  which is the exact solution of a perturbed problem, and then estimate how far  $\hat{x}$  is from the exact solution x of the original problem.

In order to obtain a compact expression, we will make the assumption that the following quantities are moderately small. For a real  $\triangle$  satisfying  $0 < \triangle < 1$  assume that

- (9.35a)  $h_1 ||L^{-1}L|| \leq \Delta$  where  $h_1 = \max\{1, \sqrt{r/2}\}$ ,
- (9.35b)  $2^{-t}\overline{K}_2(r) \chi(C) \leq (\Delta \Delta^2)/(1 + \Delta)$ ,

(9.35c) 
$$h_2 || \delta W W^{-1} || \leq \Delta$$
 where  $h_2 = \max \{1, \sqrt{(n-r)/2} \}$ ,

(9.35d) 
$$2^{-t} \overline{K}_2(n-r) \chi(B_2) \leq (\Delta - \Delta^2)/(1 + \Delta)$$
,

- (9.35e)  $||\hat{R}||/||f|| \leq \Delta$ ,
- $(9.35f) \qquad ||\mathbf{\hat{x}} \mathbf{x}||/||\mathbf{x}|| \leq \Delta,$

where  $\hat{x}$  and  $\hat{R}$  will be defined in (9.37) and (9.38). Also, assume that the conditions needed to satisfy the hypotheses of Lemma 5.1 and Theorem 7.1 are met.

It should be clear from (9.8), (6.5c), (9.35a) and (9.35b) that

(9.35g) 
$$||\overline{L}^{-1}|| ||\delta L^*|| \leq \Delta$$
,

where  $\hat{L} = \overline{L} = \delta L + L$ . Likewise from (9.29), (6.6c), (9.35c) and (9.35d) we have

(9.35h) 
$$||\delta W^*|| ||\overline{W}^{-1}|| \leq \Delta$$
,

where  $\hat{W} = \overline{W} = \delta W + W$ .

We will now account for most of our errors by interpreting them as errors induced by a perturbation of the initial problem (9.1), and then apply Theorem 7.1 to bound this portion of the error. Consider the following constrained least squares problem

(9.36a) 
$$\hat{C} = C + \delta C, \ \hat{g} = g, \ i.e. \ \delta g = 0$$
,

where  $\delta C$  is defined in (9.5),

(9.36b) 
$$\hat{A} = A + \delta A$$
,  $\hat{f} = f$ , i.e.  $\delta f = 0$ ,

where

(9.36c) 
$$\delta A = \left[0 \mid \delta B_2^*\right] \hat{Q}^T + \delta A^*$$
,

 $\delta B_2^*$  and  $\delta A^*$  are defined in (9.20) and (9.11) respectively, with solution  $\hat{x}$  which satisfies

$$(9.37) \qquad \overset{\wedge}{\mathbf{Cx}} = \mathbf{g},$$

minimizing the norm of

(9.38) 
$$\hat{R} = A x - f$$

From (9.5) it is clear that the orthogonal matrix  $\hat{Q}$  reduces  $\hat{C}$  exactly to  $[\bar{L}|0]$ , where (9.6b) gives a bound for  $||\delta C||$ . Also, from (9.36b) and (9.11) we have

(9.39)  $\widehat{A} \widehat{Q} = \left[\overline{B}_1 | \overline{B}_2 + \delta B_2^*\right],$ 

and from (9.20) the orthogonal matrix  $\hat{V}$  reduces  $\overline{B}_2 + \delta B_2^*$  to  $\left[\frac{\overline{W}}{0}\right]$ . This is the reason we have  $\hat{L} = \overline{L}$  and  $\hat{W} = \overline{W}$ .

Therefore, using Algorithm 1.1 we have

where

(9.41) 
$$(\overline{B}_2 + \delta B_2^*) \dot{d}_2 = f - \overline{B}_1 \dot{d}_1 + \hat{R}.$$

In order to obtain a bound for  $||\overline{B}_1||$ ,  $||\overline{B}_2||$  we use (9.11) and apply (3.13), (3.9) and (9.12b) to obtain

(9.42) 
$$||\overline{B}_1||, ||\overline{B}_2|| = O(||A||).$$

Using (9.42) on (9.21) we have

(9.43) 
$$||_{\delta B_2}^*|| = O(2^{-t} ||A||).$$

Therefore, from (9.36c), (9.43), (3.13) and (9.12b) it follows that

(9.44) 
$$||_{\delta}A|| = O(2^{-t} ||A||),$$
  
where  $\delta A = \delta B_2^*[0|I] Q^T + \delta A^*.$ 

If we use Theorem 7.1, (9.44) and (9.6b) it follows that

$$(9.45) \qquad ||_{\mathbf{X}-\mathbf{X}}^{A}||/||_{\mathbf{X}}|| = O(2^{-t} \chi(C)) + O(2^{-t} \chi(B_{2}) ||A||/||B_{2}||) + O(2^{-t} \chi(C) \chi(B_{2}) ||A||/||B_{2}||) + O(2^{-t} \chi(C) \chi^{2}(B_{2})(||\mathbf{R}||/||f||)(||A||^{2}/ ||B_{2}||^{2})),$$

where in the last expression we have merged terms by using the fact that  $\chi(C) \ge 1$  .

We now want to estimate the error which we cannot account for by perturbation. Toward this end use (9.33), (9.30) and (9.40) to define  $\Delta_1 \in \mathbb{R}^r$  and  $\Delta_2 \in \mathbb{R}^{n-r}$ (9.46)  $\overline{x} - \hat{x} = \hat{Q} \delta \ell^* + \hat{Q} \begin{bmatrix} \Delta_1 \\ \Delta_2 \end{bmatrix}$ .

Taking norms, it follows from (9.46) and (6.10) that

(9.47) 
$$||\bar{\mathbf{x}}-\hat{\mathbf{x}}|| \leq ||\delta l^*|| + ||\Delta_1|| + ||\Delta_2||.$$

We will consider first the vector  $\Delta_1$ , and use (9.30), (9.40) and (3.17) to obtain

(9.48) 
$$\Delta_{1} = (\overline{L} + \delta L^{*})^{-1} g - \overline{L}^{-1} g$$
$$= - (I + \overline{L}^{-1} \delta L^{*})^{-1} \overline{L}^{-1} \delta L^{*} \hat{d}_{1} .$$

Taking norms and using (9.35g) and (3.16) bounds the first inverse, using (9.8), (9.35a), (6.5c) and (6.10) yields

(9.49) 
$$||\Delta_1|| = O(2^{-t} \chi(C) ||\hat{x}||).$$

Before evaluating  $||\Delta_2||$ , we will bound the vectors  $\mathcal{E}_1$  and  $\mathcal{E}_2$ . To accomplish this, we show the following result from (9.39) by using (3.9) and (3.13)

$$(9.50) \qquad ||\overline{B}_{1}||, ||\overline{B}_{2} + \delta B_{2}^{*}|| \leq ||\hat{A}||.$$

Also, it follows from the fact that  $\overline{L}^{-1} g = \hat{d}_{1}$ 

(9.51) 
$$(\overline{L} + \delta L^*)^{-1} g = (I + \overline{L}^{-1} \delta L^*)^{-1} \hat{d}_1.$$

Therefore, from (9.14), (9.50), and (9.51) we have using (9.35g), (3.16) and (6.10) that

(9.52) 
$$||\xi_1|| = O(2^{-t} ||\hat{A}|| ||\hat{x}||).$$

From (9.16), a bound for  $||\xi_2||$  follows from the simplification of the expression

$$(9.53) fint{f} - B_{1}(\overline{L} + \delta L^{*})^{-1} g - \mathcal{E}_{1} = f - \overline{B}_{1} \overline{L}^{-1} g + \overline{B}_{1} (\overline{L} + \delta L^{*})^{-1} \delta L^{*} \overline{L}^{-1} g - \mathcal{E}_{1} = (\overline{B}_{2} + \delta B_{2}^{*}) \hat{d}_{2} - \hat{R} + \overline{B}_{1}(I + \overline{L}^{-1} \delta L^{*})^{-1} \overline{L}^{-1} \delta L^{*} \hat{d}_{1} - \mathcal{E}_{1},$$

where we have used (3.17), (9.40) and (9.41). In (9.35g) it is assumed that  $||\overline{L}^{-1}|| ||\delta L^*|| \leq \Delta \leq 1$ , upon applying this <u>twice</u> and using (9.50), (3.16), (6.10) and (9.52) we have that the norm of

(9.53) is 
$$O(||\hat{A}|| ||\hat{x}||) + O(||\hat{R}||)$$
, hence from (9.16)

$$(9.54) \qquad ||\mathcal{E}_{2}|| = O(2^{-2} ||\hat{A}|| ||\hat{x}||) + O(2^{-t} \cdot [||\hat{R}||/(||\hat{A}|| ||\hat{x}||)]||\hat{A}|| ||\hat{x}||).$$

Using (9.38) we find  $||\hat{A}|| ||\hat{x}|| \ge ||f|| - ||\hat{R}||$ , which, when combined with (9.54) and the assumption (9.35e) that  $||\hat{R}||/||f|| \le \triangle \le 1$ , yields

(9.55) 
$$||\mathcal{E}_{2}|| = O(2^{-t} ||\hat{A}|| ||\hat{x}||).$$

We now define the vector  $\xi_3 \in \mathbb{R}^m$  by

(9.56) 
$$\mathcal{E}_{3} = \{ \mathbf{f} - \overline{\mathbf{B}}_{1} (\overline{\mathbf{L}} + \delta \mathbf{L}^{*})^{-1} \mathbf{g} - \mathcal{E}_{1} + \mathcal{E}_{2} \} - \{ \mathbf{f} - \overline{\mathbf{B}}_{1} \overline{\mathbf{L}}^{-1} \mathbf{g} \}$$
$$= \overline{\mathbf{B}}_{1} (\mathbf{I} + \overline{\mathbf{L}}^{-1} \delta \mathbf{L}^{*})^{-1} \overline{\mathbf{L}}^{-1} \delta \mathbf{L}^{*} \hat{\mathbf{d}}_{1} - \mathcal{E}_{1} + \mathcal{E}_{2} .$$

A bound for  $||\mathcal{E}_3||$  is obtained in similar fashion as that employed for  $||\mathcal{E}_2||$ , therefore we have from (9.50), (9.35g), (3.16), (9.8), (9.35a), (6.5c), (9.52) and (9.55) that

(9.57) 
$$||\xi_3|| = O(2^{-t} \chi(C) ||\hat{A}|| ||\hat{x}||).$$

Finally, using (9.25) and (9.22) we have

(9.58) 
$$|\delta z| \leq 2^{-t} \overline{K}_{3}(n-r) ||f - \overline{B}_{1}(\overline{L} + \delta L^{*})^{-1} g - \xi_{1} + \xi_{2}||.$$

Upon examination of (9.16) and (9.58), we see that they differ only by  $\mathcal{E}_2$  and the factor  $\overline{K}_3(n-r)$ , therefore using (9.55) it follows that (9.59)  $||\delta z|| = O(2^{-t} ||\hat{A}|| ||\hat{x}||)$ .

The definition of  $\Delta_2$  comes from (9.46), (9.40), (9.33) and (9.30),
which yields the following simplification using (9.56), (3.17) and (9.40)

$$(9.60) \qquad \bigtriangleup_{2} = (\overline{W} + \delta W^{*})^{-1} \widehat{V}_{1} \left[ \delta z + \{f - \overline{B}_{1} (\overline{L} + \delta L^{*})^{-1} g - \xi_{1} + \xi_{2} \} \right] - \overline{W}^{-1} \widehat{V}_{1} \left[ f - \overline{B}_{1} \overline{L}^{-1} g \right]$$
$$= \overline{W}^{-1} (I + \delta W^{*} \overline{W}^{-1})^{-1} \left[ -\delta W^{*} \widehat{d}_{2} + \widehat{V}_{1} (\xi_{3} + \delta z) \right].$$

Using (9.35h), (3.16), (9.29), (6.10), (3.14), (3.7), (9.57), (9.59), (6.6b) and the fact that  $\hat{W} = \overline{W}$ , it follows upon taking norms

(9.61) 
$$||\Delta_2|| = O(2^{-t} \chi(C) ||w^{-1}|| ||\hat{A}|| ||\hat{x}||),$$

where we have used the fact that  $||\overline{W}|| = ||\overline{B}_2 + \delta B_2^*|| \leq ||\widehat{A}||$ , and  $\chi(C) \geq 1$ .

Finally, from (9.44) it follows that ||A|| = O(||A||), hence

(9.62) 
$$||\Delta_2|| = O(2^{-t} \chi(C) \chi(B_2) ||\mathbf{x}|| ||\mathbf{A}||/||\mathbf{B}_2||),$$

where we have used the fact that  $||W|| = ||B_2||$  and (6.4).

We are now in a position to bound  $||\delta l^*||$ . From (9.34) we have

(9.63) 
$$||_{\delta} \boldsymbol{\ell}^{*}|| = O(2^{-t} ||\overline{\boldsymbol{\ell}}_{1}||),$$

and substitution of (9.33) into (9.46) yields after rearrangement

(9.64) 
$$\overline{l}_1 = \begin{bmatrix} \Delta_1 \\ \Delta_2 \end{bmatrix} + \hat{Q}^T \hat{x} .$$

Therefore, a bound for  $||\delta l^*||$  follows from (9.63), (9.64), (9.49) and (9.62)

$$(9.65) \qquad ||\delta \ell^*|| = O(2^{-2t} \chi(C) ||\hat{x}||) + O(2^{-2t} \chi(C) \chi(B_2) ||\hat{x}|| ||A||/||B_2||) + O(2^{-t} ||\hat{x}||).$$

We now have our bound for  $||\bar{x}\cdot\hat{x}||$  from (9.47), (9.49), (9.62), (9.65) and the fact that  $\chi(C) \ge 1$ 

(9.66) 
$$||\bar{\mathbf{x}}-\hat{\mathbf{x}}|| = O(2^{-t} \chi(C) ||\hat{\mathbf{x}}||)$$
  
+  $O(2^{-t} \chi(C) \chi(B_2) ||\hat{\mathbf{x}}|| ||A||/||B_2||).$ 

From (9.45), we have a bound for  $||\hat{\mathbf{x}}-\mathbf{x}||/||\mathbf{x}||$ , which can be tested to see if our assumption that  $||\hat{\mathbf{x}}-\mathbf{x}||/||\mathbf{x}|| \leq \Delta$  is reasonable. If it is, then from (9.66) we have

(9.67) 
$$||\bar{\mathbf{x}}-\hat{\mathbf{x}}|| = O(2^{-t} \chi(C) ||\mathbf{x}||)$$
  
+  $O(2^{-t} \chi(C) \chi(B_2) ||\mathbf{x}|| ||\mathbf{A}||/||B_2||).$ 

From (9.45) and (9.67) follows our final result

$$(9.68) \qquad ||\overline{\mathbf{x}} - \mathbf{x}|| / ||\mathbf{x}|| = O(2^{-t} \chi(C)) + O(2^{-t} \chi(B_2) ||\mathbf{A}|| / ||\mathbf{B}_2||) + O(2^{-t} \chi(C) \chi(B_2) ||\mathbf{A}|| / ||\mathbf{B}_2||) + O(2^{-t} \chi(C) \chi^2(B_2)(||\mathbf{R}|| / ||\mathbf{f}||)(||\mathbf{A}|| / ||\mathbf{B}_2||)^2).$$

(9.68) shows the errors introduced by solving (9.1) using floating point arithmetic. Here we have suppressed all quantities which are constant (depend only on  $\Delta$ , r, n, and m) to illustrate the dependence on the machine accuracy,  $2^{-t}$ , and condition of the matrices C and B<sub>2</sub>. As hoped for, the influence of the condition of the constraint equations C appears only as a linear factor in (9.68). However, the term of major influence is the last one, in which the condition of B<sub>2</sub> is squared. The effect of this term will be minimized if the quantity  $\chi(B_2) ||R||/||f||$  is not too large, i.e. the residual R is small compared to f.

### CHAPTER 2

#### BLENDING FUNCTION THEORY

In Section 1, an introduction to blending function theory as developed by Gordon and Hall [12, 13, 14, 15, 16, 17] is given. The presentation will be a generalization of their results to the more general setting of interpolation spaces. Most of the proofs of Section 1 will follow from their own. However, the more general setting of Section 1 will enable the application of blending function techniques in the following chapter to be accomplished with greater ease.

In Section 2 the dimension of discretized blending function spaces will be shown and several bases will be explicitly developed in terms of the cardinal bases for the corresponding interpolation spaces.

In Sections 3 and 4, natural cubic spline blending will be developed with error bounds given for "approximate" natural cubic spline blending.

Finally, for  $f \in C[a, b]$  and  $g \in C([a, b]x[c, d])$  we define  $||f|| = \sup_{a \leq x \leq b} |f(x)|$  and  $||g|| = \sup_{a \leq x \leq b} |g(x, y)|$ .  $a \leq x \leq b$  $c \leq y \leq d$ 

65

Section 1. Spline and Blending Function Interpolants

Given a mesh  $\pi_x$ :  $a = x_1 < x_2 < \cdots < x_M = b$ , the space of cubic splines on  $\pi_x$  is defined to be

(1.1) 
$$S^{2}(\pi_{x}) = \{ s \in C^{2}[a, b] | s(x) \text{ is a cubic polynomial on}$$
  
 $\begin{bmatrix} x_{i}, x_{i+1} \end{bmatrix} \text{ for } 1 \leq i \leq M-1 \},$ 

see Schultz 29 ].

Let  $f \in C^{(1)}[a, b]$ , then the type 1 cubic spline interpolant,  $s_{f}$ , associated with f is defined to be the unique spline which satisfies

- (1.2a)  $s_f(x_i) = f(x_i), \quad 1 \leq i \leq M$
- (1.2b)  $s'_{f}(x_{i}) = f'(x_{i}), i = 1, M.$

The following theorem of Carlson and Hall [5] gives error bounds for type 1 spline interpolation.

<u>Theorem 1.1</u> (Carlson and Hall  $\begin{bmatrix} 5 \end{bmatrix}$ ). Let  $f \in C^{(m)}[a, b]$  and  $\pi_x$  be a mesh on [a, b]. Then for  $1 \le m \le 4$ ,  $0 \le r \le \min\{m, 3\}$ 

(1.3)  $||(s_f^{-f})^{(r)}|| \leq \mathcal{E}_{mr} ||f^{(m)}|| h_x^{m-r},$ 

where the mesh size  $h_x = \max_{\substack{1 \leq i \leq M-1}} \Delta x_i; \Delta x_i = x_{i+1} - x_i$ , the mesh

ratio 
$$M_{\pi} = \max_{i} \Delta x_{i} / \min_{i} \Delta x_{i}$$
 and  $\xi_{mr}$  is given in  
 $l \leq i \leq M-1$   $l \leq i \leq M-1$ 

Table 1.

Emr	<b>r</b> = 0	r = 1	r = 2	<b>r</b> = 3
m = 1	15/4	14		
m = 2	9/8	4	10	
m = 3	71/216	31/27	5	$(63 + 8 M_{\pi}^2)/9$
m = 4	5/384	(9+√3)/216	5	$(2 + M_{\pi}^2)/4$

TABLE 1

## **Bivariate Functions:**

We will present here an introduction to blending functions as developed by Gordon and Hall [12], [13], [14], [16]. To accomplish this, it will be desirable to introduce the following general notation, which will allow us to develop at one time many of their results. Although, it should be pointed out that the methods of proof used here are just slight generalizations of their own.

In what follows, we use the notation

(1.4) 
$$f^{(k,\ell)} = D^{(k,\ell)} [f] = \partial^{k+\ell} f / \partial x^k \partial y^\ell$$

<u>Definition 1.1</u>: Let  $C^{(m, n)}([a, b] \times [c, d]) = C^{(m, n)}$  be the space of real valued functions with domain  $[a, b] \times [c, d]$  such that if  $f \in C^{(m, n)}$  then  $f^{(k, \ell)}$  exists and is continuous for  $0 \le k \le m$  and  $0 \le \ell \le n$ . Definition 1.2: A projection operator (projector) P is an idempotent linear operator from a function space onto a subspace of the function space.

We will consider here only projectors which separate the x and y variables, where the dependence on one of the variables can be represented by an element from a finite dimensional vector space. This space is to satisfy an interpolation property for certain values of the function and its derivatives at specified points. Therefore, we introduce the following notation:

We are given the mesh  $\pi_x :a \leq x_1 \leq x_2 \leq \cdots \leq x_M \leq b$  where the points are not necessarily distinct. Also, the non-negative integer m, the interpolation function  $\alpha : I_M \longrightarrow J_m$ , where  $I_M =$ {i is an integer  $|1 \leq i \leq M$ } and  $J_m = \{j \text{ is an integer } |0 \leq j \leq m\}$ , satisfying the following restraint: if there exists distinct natural numbers s and t such that  $x_s = x_t$  then  $\alpha(s) \neq \alpha(t)$ .

Definition 1.3:  $V(\pi_x, M, m) \subseteq C^{(m)}[a, b]$  is a finite dimensional interpolation vector space with respect to  $\alpha$  on the mesh  $\pi_x$  if and only if there exists an algebraic basis  $\{\varphi_i\}_{i=1}^M$  satisfying the following cardinality conditions for  $1 \leq i, j \leq M$ .

(1.5) 
$$(D^{(\alpha(i))}(\phi_i))(x_j) = \delta_{ij} = \begin{cases} 1 \text{ if } i=j\\ 0 \text{ if } i\neq j \end{cases}$$

Such a basis for  $V(\pi_x, M, m)$  is called a cardinal basis.

In terms of the cardinal basis, if  $f_{\varepsilon} C^{(m)}[a,b]$ , the the function

(1.6) 
$$p(x) = \sum_{i=1}^{M} f^{(\alpha(i))}(x_i) \phi_i(x)$$

satisfies the interpolation conditions

(1.7) 
$$(D^{(\alpha(i))}p)(x_i) = f^{(\alpha(i))}(x_i), \text{ for } 1 \le i \le M.$$

Example 1.1: Let  $V(\pi_x, M, 0)$  be the space of Lagrange interpolation polynomials of degree M-1 on [a, b] interpolating on the mesh  $\pi_x:a \leq x_1 \leq x_2 \leq \cdots \leq x_M \leq b$ , then for  $1 \leq i \leq M$ ,  $\alpha(i) = 0$  and the cardinal basis is  $\{\ell_i\}_{i=1}^M$  where

(1.8) 
$$\ell_{i}(x) = \frac{M}{\prod_{j=1}^{m} (x-x_{j})} / \frac{M}{\prod_{j=1}^{m} (x_{i}-x_{j})}, \\j \neq i \qquad j \neq i$$

Given  $f_{\epsilon} C^{(0)}[a, b]$ , then the interpolation polynomial is given by

(1.9) 
$$p(x) = \sum_{i=1}^{M} f^{(\alpha(i))}(x_i) \ell_i(x),$$

where  $f^{(\alpha(i))}(x_i) = f(x_i)$ .

Example 1.2: Type 1 cubic splines on  $M-2 \ge 2$  knots. Then  $V(\pi_x, M, 1) = S^2(\pi_x)$  on the mesh  $\pi_x:a = x_1 = x_2 < x_3 < \cdots$   $< x_{M-2} < x_{M-1} = x_M = b$ , where  $\alpha(1) = \alpha(M) = 1$  and for  $2 \le i \le M-1$ ,  $\alpha(i) = 0$ . Also, the interpolation spline satisfies conditions (1.2a) and (1.2b), and the cardinal basis  $\{C_i(x)\}_{i=1}^M$  satisfies for  $1 \le j \le M$ 

(1.10a) 
$$C_i(x_j) = \delta_{ij}, C_i'(x_1) = C_i(x_M) = 0 \text{ for } 2 \le i \le M-1$$

(1.10b) 
$$C'_{1}(x_{1}) = C'_{M}(x_{M}) = 1, C'_{i}(x_{j}) = 0$$
 for  $i = 1, M$ .

# Example 1.3: Cubic Hermite splines on M knots.

Then  $V(\pi_x, 2M, 1) = H^1(\pi_x)$  on the mesh  $\pi_x: a = x_1 = x_2 < x_3 = x_4$  $< \cdots < x_{2M-1} = x_{2M} = b$ , where  $\alpha(i) = (i-1) \mod 2$  for  $1 \le i \le 2M$ . Given  $f_{\epsilon} C^{(1)}[a, b]$ , the Hermite cubic spline interpolant p(x) satisfies  $p^{(\alpha(i))}(x_i) = f^{(\alpha(i))}(x_i)$  for  $1 \le i \le 2M$ . The cardinal basis  $\{H_i\}_{i=1}^{2M}$  satisfies for  $1 \le j \le M$  and  $1 \le i \le M$ 

(1.11a) 
$$H_{2i-1}(x_{2j-1}) = \delta_{ij}, H_{2i-1}(x_{2j}) = 0$$
,

(1.11b) 
$$H_{2i}(x_{2j-1}) = 0, \quad H'_{2i}(x_{2j}) = \delta_{ij}$$
.

<u>Remark</u>: Even though our notation  $V(\pi_x, M, m)$  does not include a reference to  $\alpha$  and  $\{\phi_i\}_{i=1}^{M}$ , it is understood that these two quantities are always associated with  $V(\pi_y, M, m)$ .

We are now in a position to define the projection operator  $P_x$ of a bivariate function. Given the interpolation vector space  $V(\pi_x, M, m)$ , the corresponding interpolation function  $\alpha$ , basis  $\{\phi_i\}_{i=1}^M$  and  $f \in C^{(m,0)}$ , we define  $P_x[f]$  pointwise to be  $(1.12a) \qquad (P_x[f])(x, y) = \sum_{i=1}^M f^{(\alpha(i), 0)}(x_i, y) \phi_i(x)$ ,

where it is clear that

(1.12b) 
$$(D^{(\alpha(i), 0)} P_{\mathbf{x}}[f]) (\mathbf{x}_{i}, y) = f^{(\alpha(i), 0)} (\mathbf{x}_{i}, y)$$

for  $1 \le i \le M$ . Examination of (1.12a) and (1.12b) shows that  $P_x$  is indeed a projection operator on  $C^{(m,0)}$ .

# We define the univariate interpolation vector space $V(\pi_y, N, n) \leq C^{(n)}[c, d]$ in the variable y with interpolation function $\beta:I_N \rightarrow J_N$ and cardinal basis $\{\psi_j\}_{j=1}^N$ in an identical manner as above. The corresponding projection operator $P_y[f]$ for $f \in C^{(0, n)}$ is defined by

(1.13) 
$$(P_{y}[f])(x, y) = \sum_{j=1}^{N} f^{(0, \beta(j))}(x, y_{j}) \psi_{j}(y),$$

where for  $1 \leq j \leq N$ 

(1.14) 
$$(D^{(0,\beta(j))} P_{y}[f]) (x, y_{j}) = f^{(0,\beta(j))} (x, y_{j}).$$

In the usual manner, we define the sum or difference of two operators to be the pointwise sum or difference, and the product as composition. Direct calculation shows that the following statements are valid on  $C^{(m_1, n_1)}$  where  $m_1 \ge m$  and  $n_1 \ge n$ , and  $P_x$  and  $P_y$  are defined as above.

(1.15a)  $P_{x} P_{y} = P_{y} P_{x}$ ,

(1.15b) 
$$D^{(0,\ell)} P_x = P_x D^{(0,\ell)}$$
 for  $0 \le \ell \le n_1$ , and

(1.15c)  $D^{(k,0)} P_y = P_y D^{(k,0)}$  for  $0 \le k \le m_1$ .

We define the projection operators

- (1.16a)  $R_x = I P_x$ ,
- (1.16b)  $R_y = I P_y$ ,

where I is the identity operator. The following relations are a direct consequence of (1.15a) to (1.16b). On C  $\binom{(m_1, n_1)}{m_1 \ge m}$  where  $m_1 \ge m$  and  $n_1 \ge n$ 

- $(1.17a) \qquad \begin{array}{c} R_{x} R_{y} = R_{y} R_{x} \\ y = x \end{array}$
- (1.17b)  $D^{(0,l)} R_{x} = R_{x} D^{(0,l)}$  for  $0 \le l \le n_{1}$ , and
- (1.17c)  $D^{(k,0)} R_y = R_y D^{(k,0)}$  for  $0 \le k \le m_1$ .

The first observation that we make is that given  $f_{\varepsilon} C^{(m, n)}$ , then  $P_x P_y[f]$  is the tensor product interpolant to f, where

(1.18) 
$$(P_{\mathbf{x}} P_{\mathbf{y}}[f])(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{M} \sum_{j=1}^{N} f^{(\alpha(i), \beta(j))}(\mathbf{x}_{i}, \mathbf{y}_{j}) \phi_{i}(\mathbf{x}) \psi_{j}(\mathbf{y}),$$

see Gordon and Hall [14]. It follows from (1.18) that for  $1 \le i \le M$ and  $1 \le j \le N$ 

(1.19) 
$$(D^{(\alpha(i), \beta(j))} P_x P_y[f]) (x_i, y_j) = f^{(\alpha(i), \beta(j))} (x_i, y_j).$$

The error of the tensor product interpolant to f is given by

(1.20) 
$$(I - P_x P_y)[f] = R_x[f] + R_y[f] - R_x R_y[f].$$

Gordon and Hall [17] (see Theorem 1.2 below) have shown that there is an interpolation scheme using the projection operators  $P_x$  and  $P_y$  which will eliminate the terms  $R_x[f]$  and  $R_y[f]$  in (1.20). Definition: The operator on  $C^{(m,n)}$  defined by

(1.21) 
$$P_x \bigoplus P_y = P_x + P_y - P_x P_y$$

is the blending function operator and  $P_x \bigoplus P_y[f]$  is the blending function interpolant to the function  $f_{\varepsilon} C^{(m, n)}$ .

<u>Theorem 1.2.</u> (Gordon and Hall [17]). If  $f \in C^{(m, n)}$  then

- (1.22)  $(I P_x \bigoplus P_y)[f] = R_x R_y[f].$
- $\frac{\text{Proof:}}{\text{Proof:}} I (P_x + P_y P_x P_y) = I P_x + (I P_x) P_y$  $= (I P_x) (I P_y)$  $= R_x R_y.$

The set of functional values and derivatives on which  $P_x \bigoplus P_y$ interpolates f is given in the following theorem.

<u>Theorem 1.3.</u> (Gordon and Hall [17]). If  $f \in C^{(m, n)}$ , then for  $1 \leq i \leq M$  and  $y \in [c, d]$ (1.23a)  $(D^{(\alpha(i), 0)} P_x \bigoplus P_y[f])(x_i, y) = f^{(\alpha(i), 0)}(x_i, y)$ , also, for  $l \leq j \leq N$  and  $x \in [a, b]$ 

(1.23b) 
$$(D^{(0,\beta(j))} P_x \oplus P_y[f]) (x, y_j) = f^{(0,\beta(j))} (x, y_j).$$

**Proof:** By direct calculation.

What we have accomplished is the approximation of a bivariate function with a sum of univariate functions. Therefore, it should be clear that the image of  $C^{(m,n)}$  under the operator  $P_x \bigoplus P_y$  is not a finite dimensional space. In comparison, from (1.18), it is clear that the image of  $C^{(m, n)}$  under the tensor product operator  $P_{\mathbf{x}} P_{\mathbf{y}}$ is a finite dimensional space. Therefore, to implement the blending function method, it is usually necessary to make second level approximations to  $f^{(\alpha(i), 0)}(x_i, \cdot)$  and  $f^{(0, \beta(j))}(\cdot, y_i)$  which have an accuracy compatible with  $P_x \bigoplus P_v[f]$ . This is the sacrifice which must be made in order to achieve this gain in accuracy. Next, we want to show how much more accuracy is gained by blending function techniques, to see if it will justify the extra work of implementation. Toward this end, we note that the accuracy of blending function methods depends upon the interpolation accuracy of the two univariate spaces  $V(\pi_x, M, m)$  and  $V(\pi_v, N, n)$ . Therefore, we introduce the following general notation for an univariate error bound.

Definition 1.4: Let  $V(\pi_x, M, m_1)$  be an interpolation vector space, with interpolation function  $\alpha$ , and cardinal basis  $\{\phi_i\}_{i=1}^M$ . We say  $V(\pi_x, M, m_1)$  has an error bound for the k<sup>th</sup> derivative if and only if there exists a function  $g_x(k) = g_x(h_x, m_1, m_2, k)$  such that if  $\hat{x} \in [a, b]$  is any point where  $\phi_i^{(k)}(\hat{x})$  exists for  $1 \le i \le M$ , and  $f \in C^{(m_2)}[a, b]$  where  $m_2 \ge \max\{m_1, k\}$  then

(1.24) 
$$|f^{(k)}(\hat{x}) - \sum_{i=1}^{M} f^{(\alpha(i))}(x_i) \phi_i^{(k)}(\hat{x})| \leq g_{x}(k) || f^{(m_2)}||,$$

where  $h_x$  is a parameter of the mesh  $\pi_x$ . It is clear, if  $V(\pi_x, M, m_1)$  has an error bound for the  $k^{th}$  derivative,  $(m_2, n)_x$   $g \in C$ ,  $y \in [c, d]$ ,  $P_x$  is the corresponding projection operator and x is any point satisfying Definition 1.4, then

- (1.25a)  $(D^{(k, 0)} P_x[f])(\hat{x}, y)$  exists, and
- (1.25b)  $(D^{(k,0)} P_{\mathbf{x}}[f]) (\hat{\mathbf{x}}, \cdot) \in C^{(n)} ([c,d]).$

Therefore, it follows from (1.15b) that for  $0 \leq l \leq n$ 

(1.26a) 
$$(D^{(k, l)} P_{\mathbf{x}}[f])(\hat{\mathbf{x}}, \mathbf{y}) = (D^{(k, 0)} P_{\mathbf{x}}[f^{(0, l)}])(\hat{\mathbf{x}}, \mathbf{y}),$$

correspondingly

(1.26b) 
$$(D^{(k,\ell)} R_x[f])(\hat{x}, y) = (D^{(k,0)} R_x[f^{(0,\ell)}])(\hat{x}, y)$$

and finally

(1.26c) 
$$|(D^{(k,\ell)} R_{\mathbf{x}}[f])(\hat{\mathbf{x}},\mathbf{y})| \leq g_{\mathbf{x}}(k) \sup_{\mathbf{x} \in \mathbf{t} \leq \mathbf{b}} |f^{(m_2,\ell)}(t,\mathbf{y})|.$$

For the interpolation space  $V(\pi_y, N, n_1)$ , interpolation function  $\beta:I_N \longrightarrow J_n$  and cardinal basis  $\{\psi_j\}_{j=1}^N$ , we use an analogous definition for an error bound for the  $\ell^{\text{th}}$  derivative, and we have conclusions similar to those given by (1.25a) to (1.26c).

We are now in a position to prove the following theorem, which is a generalization of the error analysis given by Gordon and Hall [17].

Theorem 1.4. Let  $V(\pi_x, M, m_1)$  be an interpolation vector space which has an error bound for the k<sup>th</sup> derivative for  $0 \le k \le m_3$ , where  $m_3$  is an integer such that  $0 \le m_3 \le m_2$ . Also, let  $V(\pi_y, N, n_1)$  be an interpolation vector space which has an error bound for the  $\ell^{\text{th}}$  derivative for  $0 \le \ell \le n_3$ , where  $n_3$  is an integer such that  $0 \le n_3 \le n_2$ . If  $f \in C$ , then  $(1.27) \qquad |((I - P_x \bigoplus P_y)[f])^{(k,\ell)}(\hat{x}, \hat{y})| \le g_x(k) g_y(\ell)| |f^{(m_2, n_2)}||,$ 

(1.28) 
$$|((I - P_{x} P_{y})[f])^{(k, \ell)}(\hat{x}, \hat{y})| \leq g_{x}^{(k)}||f^{(m_{2}, \ell)}||$$
  
$$+ g_{y}^{(\ell)}||f^{(k, n_{2})}|| + g_{x}^{(k)}g_{y}^{(\ell)}||f^{(m_{2}, n_{2})}||,$$

where  $\hat{x}$ ,  $\hat{y}$ ,  $m_2$ ,  $n_2$ ,  $g_x(k)$  and  $g_y(l)$  are given in Definition 1.4. <u>Proof</u>: For  $0 \le q \le n_2$  and for each  $y \in [c, d]$  it follows from (1.26c)

(1.29) 
$$|(D^{(k,q)} \operatorname{R}_{\mathbf{x}}[f])(\hat{\mathbf{x}}, \mathbf{y})| \leq g_{\mathbf{x}}(k) \sup_{\substack{a \leq \mathbf{x} \leq \mathbf{b}}} |f^{(m_2,q)}(\cdot, \mathbf{y})| .$$

Because  $D^{(k,0)} R_{x}[f](\hat{x}, \cdot) \in C^{(n_{2})}[c, d]$  and  $\psi_{j}^{(\ell)}(\hat{y})$  exists for  $1 \leq j \leq N$  we have that  $(D^{(k,\ell)} R_{y} R_{x}[f])(\hat{x}, \hat{y})$  exists and from (1.29) and Theorem 1.2

(1.30) 
$$|(D^{(k,\ell)} R_{y} R_{x}[f]) (\hat{x}, \hat{y})| = |(D^{(0,\ell)} R_{y} D^{(k,0)} R_{x}[f]) (\hat{x}, \hat{y})|$$

$$\leq g_{y}(\ell) \sup_{c \leq y \leq d} |D^{(k,n_{2})} R_{x}[f] (\hat{x}, \cdot)|$$

$$\leq g_{x}(k) g_{y}(\ell) ||f^{(m_{2},n_{2})}||.$$

The proof of (1.28) follows in a similar way from (1.20), (1.29) and (1.30)

<u>Corollary 1.1</u>. (Gordon and Hall [17]). If  $f \in C^{(m, n)}([o, h] \times [o, h])$ , h < 1, and the spaces  $V(\pi_x, m, 0)$  and  $V(\pi_y, n, 0)$  are polynomial spaces as defined in Example 1.1, then

(1.31) 
$$||((I - P_x \oplus P_y)[f])^{(k, \ell)}|| \leq \overline{\mathcal{E}}_{mk} \overline{\mathcal{E}}_{n\ell} ||f^{(m, n)}||h^{m+n-(k+\ell)}|$$

 $\mathtt{and}$ 

(1.32) 
$$||((I - P_x P_y)[f])^{(k, \ell)}|| \leq \overline{\mathcal{E}}_{mk} ||f^{(m, \ell)}||h^{m-k}$$

+ 
$$\overline{\mathcal{E}}_{n\ell}$$
 ||  $f^{(k,n)}$  ||  $h^{n-\ell}$  +  $\overline{\mathcal{E}}_{mk}$   $\overline{\mathcal{E}}_{n\ell}$  ||  $f^{(m,n)}$  ||  $h^{m+n-(k+\ell)}$ ,

where  $0 \leq k \leq m$ ,  $0 \leq l \leq n$ ,  $\overline{\xi}_{mk} = 1/(m-k)!$  and  $\overline{\xi}_{nl} = 1/(n-l)!$ .

<u>Proof</u>: From [24, p. 289], if  $f \in C^{(m)}[0, h]$ , then

(1.33) 
$$|f^{(k)}(x) - \sum_{i=1}^{m} f^{(\alpha(i))}(x_i) \ell_i^{(k)}(x)| \leq \overline{\mathcal{E}}_{mk} ||f^{(m)}||h^{m-k},$$

hence  $g_{\mathbf{x}}(h, 0, m, k) = \overline{\xi}_{mk} h^{m-k}$ . An application of Theorem 1.4 completes the proof.

<u>Corollary 1.2</u>. (Carlson and Hall  $\begin{bmatrix} 5 \end{bmatrix}$ ). If  $f \in C^{(m,n)}([a, b]x[c, d])$ ,  $l \leq m, n \leq 4$  and the spaces  $V(\pi_x, M, 1)$  and  $V(\pi_y, N, 1)$  are cubic spline spaces given in Example 1.2, then

(1.34) 
$$||((I - P_x \oplus P_y)[f])^{(k,\ell)}|| \leq \mathcal{E}_{\substack{m,k \ n\ell}} \mathcal{E}_{\substack{m,n \ n\ell}} ||f^{(m,n)}||_{x} \mathcal{E}_{\substack{m-k \ n-\ell}}^{(m,n)},$$

and

(1.35) 
$$||((I - P_x P_y)[f])^{(k,\ell)}|| \leq \mathcal{E}_{mk} ||f^{(m,\ell)}||h_x^{m-k} + \mathcal{E}_{n\ell} ||f^{(k,n)}||h_y^{n-\ell} + \mathcal{E}_{mk} \mathcal{E}_{n\ell} ||f^{(m,n)}||h_x^{m-k}h_y^{n-\ell}$$

where  $0 \le k \le \min\{m, 3\}$ ,  $0 \le l \le \min\{n, 3\}$ , also  $\mathcal{E}_{mk}$  and  $\mathcal{E}_{nl}$ are given in Table 1,  $h = \max_{\substack{x \\ 2 \le i \le M-2}} (x_{i+1} - x_i)$  and

 $h_{y} = \max_{2 \le j \le N-2} (y_{j+1} - y_{j}).$ 

<u>Proof</u>: We have an error bound for the derivatives given by Theorem 1.1, where  $g_x(h_x, 1, m, k) = \xi_{mk} h_x^{m-k}$ .

Carlson and Hall in the same paper  $\begin{bmatrix} 5 \end{bmatrix}$  have given other error bounds for  $P_x P_y[f]$  under weaker continuity requirements for f.

<u>Corollary 1.3</u>. If  $f \in C^{(m,n)}([a,b] \times [c,d])$ ,  $l \leq m,n \leq 4$  and the spaces  $V(\pi_x, 2M, 1)$  and  $V(\pi_y, 2N, 1)$  are Hermite cubic spline spaces defined in Example 1.3 then

(1.36) 
$$||((I - P_x \oplus P_y)[f])^{(h, l)}|| \leq \hat{\mathcal{E}}_{mk} \hat{\mathcal{E}}_{nl} ||f^{(m, n)}||h_x^{m-k}h_y^{n-l}|$$

and

(1.37) 
$$||((I - P_x P_y)[f])^{(k, \ell)}|| \leq \hat{\mathcal{E}}_{mk} ||f^{(m, \ell)}||h_x^{m-k} + \hat{\mathcal{E}}_{n\ell} ||f^{(k, n)}||h_y^{n-\ell} + \hat{\mathcal{E}}_{mk} \hat{\mathcal{E}}_{n\ell} ||f^{(m, n)}||h_x^{m-k}h_y^{n-\ell}$$

where  $0 \le k \le \min\{m, 3\}$ ,  $0 \le l \le \min\{n, 3\}$  also  $\hat{\mathcal{E}}_{mk}$  and  $\hat{\mathcal{E}}_{nl}$ are given in Table 2,  $h_x = \max_{1 \le i \le M-1} (x_{2i+1} - x_{2i})$ ,

 $h_{y} = \max_{1 \le j \le N-1} (y_{2j+1} - y_{2j}).$ 

TABLE 2

Ê <sub>mk</sub>	k = 0	k = 1	k = 2	k = 3
m = 1	5/4	4		
m = 2	3/4	5/2	12	
m = 3	7/24	1	5	7
m = 4	1/384	<del>√3/</del> 216	1/12	1/2

Proof: From Carlson and Hall  $\begin{bmatrix} 5 \end{bmatrix}$ , if  $f \in C^{(m)}[a, b]$ , then (1.38)  $||f^{(k)} - \sum_{i=1}^{2M} f^{(\alpha(i))}(x_i) H_i|| \leq \hat{\mathcal{E}}_{mk} ||f^{(m)}|| h_x^{m-k}$ ,

implying that  $g_x(h_x, l, m, k) = \hat{\xi}_{mk} h_x^{m-k}$ .

We are now in a position to describe the second level decomposition of the blending function interpolant, see [17]. The purpose of this is to create a finite dimensional interpolation scheme while preserving the accuracy of our blending techniques.

Define the

Meshes:  $\overline{\pi}_{\mathbf{x}}$ :  $\mathbf{a} \leqslant \overline{\mathbf{x}}_{1} \leqslant \overline{\mathbf{x}}_{2} \leqslant \cdots \leqslant \overline{\mathbf{x}}_{\overline{\mathbf{M}}} \leqslant \mathbf{b}$  $\overline{\pi}_{\mathbf{y}}$ :  $\mathbf{c} \leqslant \overline{\mathbf{y}}_{1} \leqslant \overline{\mathbf{y}}_{2} \leqslant \cdots \leqslant \overline{\mathbf{y}}_{\overline{\mathbf{N}}} \leqslant \mathbf{d}$ 

Interpolation functions:  $\overline{\alpha}$ :  $I_{\overline{M}} \longrightarrow J_{\overline{m}}$ 

$$\overline{\beta}: \ I_{\overline{N}} \longrightarrow J_{\overline{n}}$$

Cardinal bases:  $\{\overline{\phi}_i\}_{i=1}^{\overline{M}}$ 

$$\{\overline{\psi}_j\}_{j=1}^{\overline{N}}$$

Projection operators:  $\overline{P}_{x}$  and  $\overline{R}_{x} = I - \overline{P}_{x}$ 

$$\overline{P}_{y}$$
 and  $\overline{R}_{y} = I - \overline{P}_{y}$ 

From this define the interpolation spaces  $\overline{V}(\overline{n}_x, \overline{M}, \overline{m}) \subseteq C^{(\overline{m})}([a, b])$ and  $\overline{V}(\overline{n}_y, \overline{N}, \overline{n}) \subseteq C^{(\overline{n})}([c, d]).$  For  $f \in C^{(m^*, n^*)}$ ,  $m^* = \max\{m, \overline{m}\}$  and  $n^* = \max\{n, \overline{n}\}$  define the discrete blending approximation to  $P_x \bigoplus P_y[f]$  as

(1.39) 
$$\overline{\mathbf{P}_{\mathbf{x}} \oplus \mathbf{P}_{\mathbf{y}}}[f] = \overline{\mathbf{P}}_{\mathbf{x}} \mathbf{P}_{\mathbf{y}}[f] + \overline{\mathbf{P}}_{\mathbf{y}} \mathbf{P}_{\mathbf{x}}[f] - \mathbf{P}_{\mathbf{x}} \mathbf{P}_{\mathbf{y}}[f],$$

where, for example

(1.40) 
$$\widehat{P}_{\mathbf{x}} P_{\mathbf{y}}[f](\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{M} \sum_{j=1}^{N} f^{(\overline{\alpha}(i), \beta(j))}(\overline{\mathbf{x}}_{i}, \mathbf{y}_{j}) \overline{\phi}_{i}(\mathbf{x}) \psi_{j}(\mathbf{y}) .$$

In general, the discrete blending approximation  $\overline{P_x \oplus P_y}$  [f] does not interpolate values of f and its derivatives. However, with the following restriction on the interpolation spaces, we will prove that  $\overline{P_x \oplus P_y}$  [f] does indeed interpolate.

<u>Definition 1.5</u>: Let  $V(\pi_x, M, m)$  and  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$  be interpolation vector spaces with interpolation functions  $\alpha$  and  $\overline{\alpha}$  respectively, then  $V(\pi_x, M, m)$  is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$  if and only if  $m \leq \overline{m}$  and for each i where  $1 \leq i \leq M$  there exists an  $\overline{i}$  such that  $1 \leq \overline{i} \leq \overline{M}, x_i = \overline{x_i}$  and  $\alpha(i) = \overline{\alpha(i)}$ .

We have the following generalization of a theorem due to Gordon and Hall  $\begin{bmatrix} 17 \end{bmatrix}$ .

<u>Theorem 1.5</u>. Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$  where  $V(\pi_x, M, m)$  is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$  and  $V(\pi_y, N, n)$  is subordinate to  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$ , if  $f \in C^{(m, n)}$  then for  $1 \le i \le M$ ,  $1 \le j \le N$ ,  $1 \le \overline{i} \le \overline{M}$ 

and 
$$1 \leq \overline{j} \leq \overline{N}$$
  
(1.41)  $(D^{(\overline{\alpha}(\overline{i}),\beta(\overline{j}))}\overline{P_{x} \oplus P_{y}}[f])(\overline{x}_{\overline{i}},y_{j}) = f^{(\overline{\alpha}(\overline{i}),\beta(\overline{j}))}(\overline{x}_{\overline{i}},y_{j})$ 

and

(1.42) 
$$(D^{(\alpha(i),\overline{\beta}(\overline{j}))}\overline{P_{x} \oplus P_{y}}[f])(x_{i},\overline{y_{j}}) = f^{(\alpha(i),\overline{\beta}(\overline{j}))}(x_{i},\overline{y_{j}}) .$$

**Proof:** By direct calculation.

Even if our blending spaces are not subordinate to those used in the second level decomposition,  $\overline{P_x \oplus P_y}[f]$  is still an approximation to  $P_x \oplus P_y[f]$  and hence to f also. This is the conclusion of the following theorem.

<u>Theorem 1.6.</u> (Gordon and Hall [17]). Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n})$ , if  $f \in C^{(m^*, n^*)}$  where  $m^* = \max\{m, \overline{m}\}$  and  $n^* = \max\{n, \overline{n}\}$  then

(1.43) 
$$(I - \overline{P_{\mathbf{x}} \oplus P_{\mathbf{y}}})[f] = \overline{R_{\mathbf{x}}}[f] + \overline{R_{\mathbf{y}}}[f] + R_{\mathbf{x}} R_{\mathbf{y}}[f]$$

$$- \overline{R}_{x} R_{y} [f] - \overline{R}_{y} R_{x} [f] .$$

$$\underline{Proof}: I - \overline{P_{x} \oplus P_{y}} = I - \overline{P_{x}} P_{y} - \overline{P_{y}} P_{x} + P_{x} P_{y}$$

$$= I - P_{x} \oplus P_{y} + \overline{R}_{x} P_{y} + \overline{R}_{y} P_{x}$$

$$= R_{x} R_{y} + \overline{R}_{x} - \overline{R}_{x} R_{y} + \overline{R}_{y} - \overline{R}_{y} R_{x} .$$

An examination of Theorem 1.6 and Theorem 1.4 yields the following generalization of a theorem by Gordon and Hall  $\begin{bmatrix} 17 \end{bmatrix}$ .

 $\begin{array}{ll} \underline{\text{Theorem 1.7}}. & \text{Given the nonnegative integers } \mathbf{m}_{3} \leq \min\{\mathbf{m}_{2}, \overline{\mathbf{m}}_{2}\}\\\\ \text{and } \mathbf{n}_{3} \leq \min\{\mathbf{n}_{2}, \overline{\mathbf{n}}_{2}\}, \text{ the interpolation spaces } V(\pi_{\mathbf{x}}, \mathbf{M}, \mathbf{m}_{1}) \text{ and } \\\\ \overline{V}(\overline{\pi}_{\mathbf{x}}, \overline{\mathbf{M}}, \overline{\mathbf{m}}_{1}) \text{ which have error bounds for } \mathbf{k}, \text{ where } 0 \leq \mathbf{k} \leq \mathbf{m}_{3},\\\\ \text{also the interpolation spaces } V(\pi_{\mathbf{y}}, \mathbf{N}, \mathbf{n}_{1}) \text{ and } \overline{V}(\overline{\pi}_{\mathbf{y}}, \overline{\mathbf{N}}, \overline{\mathbf{n}}_{1}) \text{ which have error bounds for } \mathbf{k}, \text{ where } 0 \leq \mathbf{k} \leq \mathbf{m}_{3},\\\\ \text{and the error bounds for } \mathbf{\ell} \text{ where } 0 \leq \mathbf{\ell} \leq \mathbf{n}_{3}, \text{ and if } \mathbf{f} \in \mathbf{C}^{(\mathbf{m}_{2}^{*}, \mathbf{n}_{2}^{*})}\\\\ \text{where } \mathbf{m}_{2}^{*} = \max\{\mathbf{m}_{2}, \overline{\mathbf{m}}_{2}\} \text{ and } \mathbf{n}_{2}^{*} = \max\{\mathbf{n}_{2}, \overline{\mathbf{n}}_{2}\} \text{ then }\\\\ (1.44) \quad \left|((\mathbf{I} - \overline{\mathbf{P}_{\mathbf{x}} \oplus \mathbf{P}_{\mathbf{y}}})[\mathbf{f}])^{(\mathbf{k}, \mathbf{\ell})}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})\right| \leq \overline{\mathbf{g}}_{\mathbf{x}}(\mathbf{k}) \left||\mathbf{f}^{(\overline{\mathbf{m}}_{2}, \mathbf{n}_{2})}\right||\\\\ \quad + \overline{\mathbf{g}}_{\mathbf{y}}(\mathbf{\ell})|\left|\mathbf{f}^{(\mathbf{k}, \overline{\mathbf{n}}_{2})}\right|| + \mathbf{g}_{\mathbf{x}}(\mathbf{k}) \mathbf{g}_{\mathbf{y}}(\mathbf{\ell})|\left|\mathbf{f}^{(\mathbf{m}_{2}, \mathbf{n}_{2})}\right||\\\\ \quad + \overline{\mathbf{g}}_{\mathbf{x}}(\mathbf{k}) \mathbf{g}_{\mathbf{y}}(\mathbf{\ell})|\left|\mathbf{f}^{(\overline{\mathbf{m}}_{2}, \mathbf{n}_{2}\right|| + \mathbf{g}_{\mathbf{x}}(\mathbf{k}) \overline{\mathbf{g}}_{\mathbf{y}}(\mathbf{\ell})|\left|\mathbf{f}^{(\mathbf{m}_{2}, \overline{\mathbf{n}}_{2})}\right||, \end{aligned}$ 

where  $\hat{x}$  and  $\hat{y}$  satisfy the conditions of Definition 1.4.

Proof: Similar to the proof of Theorem 1.4.

If it is possible, by some procedure, to increase the accuracy of the spaces  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m}_1)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n}_1)$ , we see from Theorem 1.7 that the total increase in accuracy for  $\overline{P_x \oplus P_y}$  is limited by the term  $g_x(k) g_y(l) || f^{(m_2, n_2)} ||$ . Therefore, it is not possible to increase the accuracy beyond that of the original blending approximation. Example 1.4: If we take for each of our interpolation spaces the space of cubic splines defined in Example 1.2, then  $V(\pi_x, M, 1) = S^2(\pi_x)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, 1) = S^2(\overline{\pi}_x)$ ,  $V(\pi_y, N, 1) = S^2(\pi_y)$ ,  $\overline{V}(\overline{\pi}_y, \overline{M}, 1) = S^2(\overline{\pi}_y)$ ,  $m_2 = \overline{m}_2$ ,  $n_2 = \overline{n}_2$ ,  $m_3 = \min\{m_2, 3\}$ ,  $n_3 = \min\{n_2, 3\}$   $1 \le m_2$ ,  $n_2 \le 4$  and if  $f \in C^{(m_2, n_2)}$  then (1.45)  $||((I - \overline{P_x \oplus P_y}([f])^{(k, \ell)})|| \le \mathcal{E}_{m_2, k}||f^{(m_2 \ell)}||\overline{h_x}^{m_2 - k}$   $+ \mathcal{E}_{n_2, \ell}||f^{(k, n_2)}||\overline{h_y}^{n_2 - \ell}$   $+ \mathcal{E}_{m_2, k} \mathcal{E}_{n_2, \ell}||f^{(m_2, n_2)}||h_x^{m_2 - k}h_y^{n_2 - \ell}$   $+ \mathcal{E}_{m_2, k} \mathcal{E}_{n_2, \ell}||f^{(m_2, n_2)}||\overline{h_x}^{m_2 - k}h_y^{n_2 - \ell}$  $+ \mathcal{E}_{m_2, k} \mathcal{E}_{n_2, \ell}||f^{(m_2, n_2)}||\overline{h_x}^{m_2 - k}h_y^{n_2 - \ell}$ .

For ease of comparison, we will let  $h = \max \{h_x, h_y\}$ ,  $\overline{h} = \max \{\overline{h}_x, \overline{h}_y\}$ ,  $m_2 = n_2 = 4$  and k = 1 = 0. Then the accuracy in (1.45) is limited to  $O(h^8)$ , and therefore we take  $\overline{h} = h^2$  to preserve the overall accuracy of the scheme. This tells us how much the meshes  $\overline{\pi}_x$  and  $\overline{\pi}_y$  must be refined in (1.45) to obtain a scheme which is  $O(h^8)$ .

In comparison, for bicubic splines (tensor product (1.35)) the accuracy is  $O(h^4)$ . In the next section we will compute the dimension of discretized blending function spaces, and therefore will be able to show that the increase in accuracy is worth the extra labor of implementation. <u>Remark</u>: In our notation  $V(\pi_x, M, m)$ , m can be any integer which satisfies  $m^* \leq m \leq m^{**}$ , where  $m^*$  is the smallest integer which allows enough continuity to perform our interpolation and  $m^{**}$ is the largest integer (if it exists) for which  $V(\pi_x, M, m) \subseteq C^{(m^{**})}[a, b]$ . In Example 1.2, for type 1 cubic spline interpolation, m can be equal to 1 or 2. In Example 1.3, for cubic Hermite splines, m = 1 is the only choice. Note that by increasing m, we restrict the choice of f which will satisfy our theorems. However, sometimes it is necessary to use a  $m > m^*$ , if, for example, we wish to satisfy Definition 1.5 or Theorem 1.5.

#### Section 2. Dimension of Discretized Blending Function Spaces

In applications other than interpolation, such as discrete least squares, Ritz-Galerkin methods and collocation, it is necessary to have a finite dimensional space on which to do the computation. Therefore, we introduce the following definition.

Definition 2.1: Let  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n})$  be interpolation spaces (see Definition 1.3), then define the discretized blending function space DBF( $V(\pi_x, M, m)$ ,  $V(\pi_y, N, n)$ ;  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$ ,  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n})$ ) to be the image of  $C^{(m^*, n^*)}([a, b]x[c, d])$  under the linear discretized blending function operator  $\overline{P}_x P_y + \overline{P}_y P_x - P_x P_y$ , (denoted by  $\overline{P_x \oplus P_y}$ ), where  $m^* = max \{m, \overline{m}\}$  and  $n^* = max \{n, \overline{n}\}$ .

From this notation, it is understood that the blending spaces are  $V(\pi_x, M, m)$  and  $V(\pi_y, N, n)$ , and the spaces which give the second level approximations are  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n})$ .

When the interpolation spaces are understood, we will drop them from our notation and denote the space of discretized blending functions as DBF.

It is clear that DBF is indeed a vector space.

Finally, for this section,  $\alpha$ ,  $\overline{\alpha}$ ,  $\beta$  and  $\overline{\beta}$  will be the interpolation functions and  $\{\phi_i\}_{i=1}^{\overline{M}}, \{\overline{\phi}_i\}_{i=1}^{\overline{M}}, \{\psi_j\}_{j=1}^{\overline{N}}$  and  $\{\overline{\psi}_j\}_{j=1}^{\overline{N}}$  will be

the cardinal bases for  $V(\pi_x, M, m)$ ,  $\overline{V(\pi_x}, \overline{M}, \overline{m})$ ,  $V(\pi_y, N, n)$  and  $\overline{V(\pi_y}, \overline{N}, \overline{n})$ , respectively.

We will need the following information about product space. If  $V(\pi_x, M, m)$  and  $V(\pi_y, N, n)$  are interpolation spaces, then their product is defined to be

(2.1) 
$$V(\pi_{x}, M, m) \bigotimes V(\pi_{y}, N, n) = \operatorname{span} \{ fg | f \in V(\pi_{x}, M, m) \text{ and} g \in V(\pi_{y}, N, n) \}.$$

An element  $f \bigotimes g_{\epsilon} V(\pi_x, M, m) \bigotimes V(\pi_y, N, n)$  will be abbreviated by  $fg = f \bigotimes g$ , where  $f_{\epsilon} V(\pi_x, M, m)$ ,  $g_{\epsilon} V(\pi_y, N, n)$  and fg(x, y) = f(x) g(y) for  $(x, y)_{\epsilon} [a, b] x [c, d]$ .

If  $\{u_i\}_{i=1}^{M}$  is any basis for  $V(\pi_x, M, m)$  and  $\{w_j\}_{j=1}^{N}$  is any basis for  $V(\pi_y, N, n)$  then  $\{u_i, w_j\}_{\substack{i \\ j \\ l \le j \le N}}$  is a basis for  $V(\pi_x, M, m)$  (x)

 $V(\pi_v, N, n)$ . Therefore,

(2.2) Dimension 
$$(V(\pi_x, M, m) \otimes V(\pi_y, N, n))$$
  
= Dim  $(V(\pi_x, M, m)) \cdot$  Dim  $(V(\pi_y, N, n))$   
=  $M \cdot N$ .

Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n})$ , and, for example, if  $f \in V(\pi_x, M, m) \otimes V(\pi_y, N, n)$  then (2.3)  $f = \sum_{i=1}^{M} \sum_{j=1}^{N} f^{(\alpha(i), \beta(j))}(x_i, y_j) \phi_i \psi_j$ ,

by the uniqueness of representation. Also

(2.4) 
$$(V(\pi_x, M, m) \cap \overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})) \otimes V(\pi_y, N, n)$$
  
=  $(V(\pi_x, M, m) \otimes V(\pi_y, N, n))$   
 $\cap (\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m}) \otimes V(\pi_y, N, n)),$ 

and

(2.5) 
$$V(\pi_x, M, m) \otimes (V(\pi_y, N, n) \cap \overline{V}(\overline{\pi}_y, \overline{N}, \overline{n}))$$
  
=  $(V(\pi_x, M, m) \otimes V(\pi_y, N, n))$   
 $\cap (V(\pi_x, M, m) \otimes \overline{V}(\overline{\pi}_y, \overline{N}, \overline{n}))$ .

We see that (2.4) is true, because the left hand side is clearly contained in the right hand side. Also, if f is an element of the right hand side of (2.4), then

(2.6) 
$$f = \sum_{j=1}^{N} \begin{cases} M \\ \Sigma \\ i=1 \end{cases} f^{(\alpha(i),\beta(j))}(\mathbf{x}_{i},\mathbf{y}_{j}) \phi_{i} \} \psi_{j}$$
$$= \sum_{j=1}^{N} \begin{cases} \overline{M} \\ \Sigma \\ i=1 \end{cases} f^{(\overline{\alpha}(i),\beta(j))}(\overline{\mathbf{x}}_{i},\mathbf{y}_{j}) \phi_{i} \} \psi_{j}$$

We will have proved our result if the bracketed expressions in (2.6) are contained in  $V(\pi_x, M, m) \cap \overline{V}(\overline{\pi_x}, \overline{M}, \overline{m})$  for  $1 \leq j \leq N$ . Using the cardinality conditions, we have for each j satisfying  $1 \leq j \leq N$ 

•

(2.7) 
$$f^{(0,\beta(j))}(\cdot, \mathbf{y}_{j}) = \sum_{i=1}^{M} f^{(\alpha(i),\beta(j))}(\mathbf{x}_{i}, \mathbf{y}_{j}) \phi_{i} \in V(\pi_{\mathbf{x}}, \mathbf{M}, \mathbf{m})$$
$$= \sum_{i=1}^{\overline{M}} f^{(\overline{\alpha}(i),\beta(j))}(\overline{\mathbf{x}}_{i}, \mathbf{y}_{j}) \overline{\phi}_{i} \in \overline{V}(\overline{\pi}_{\mathbf{x}}, \overline{\mathbf{M}}, \overline{\mathbf{m}}),$$

which proves (2.4). We note that (2.5) follows in a similar manner.

We now show that the following statement is valid for  $m\leqslant\overline{m}$  or  $n\leqslant\overline{n}$ 

(2.8) 
$$(\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n)) \cap (V(\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n})) \subseteq V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n).$$

To see that (2.8) is true, let f be an element of the left hand side. Then

(2.9) 
$$f = \sum_{j=1}^{N} \begin{cases} \overline{M} \\ \Sigma \\ i=1 \end{cases} f^{(\overline{\alpha}(i), \beta(j))}(\overline{x}_{i}, y_{j}) \overline{\phi}_{i} \} \psi_{j}$$
$$= \sum_{\substack{N \\ i=1 \end{cases}} \sum_{j=1}^{N} f^{(\alpha(i), \overline{\beta}(j))}(x_{i}, \overline{y}_{j}) \phi_{i} \overline{\psi}_{j} .$$

We will be done if it can be shown that the bracketed expression in (2.9) is in  $V(\pi_x, M, m)$  for  $l \leq j \leq N$ . Using  $n \leq \overline{n}$  and the cardinality conditions for  $l \leq \widehat{j} \leq N$ , then

$$(2.10) \qquad f^{(0,\beta(\hat{j}))}(\cdot,y_{\hat{j}}) = \sum_{i=1}^{\overline{M}} f^{(\overline{\alpha}(i),\beta(\hat{j}))}(\overline{x}_{i},y_{\hat{j}}) \overline{\phi}_{i}$$
$$= \sum_{i=1}^{M} \left[ \sum_{j=1}^{\overline{N}} f^{(\alpha(i),\overline{\beta}(j))}(x_{i},\overline{y}_{j}) \overline{\psi}_{j}^{(\beta(\hat{j}))}(y_{\hat{j}}) \right] \phi_{i}$$
$$\in V(\pi_{x}, M, m) .$$

Finally, from (1.39) and (1.40) it is clear that

(2.11) DBF 
$$\leq$$
 span { $(\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n))$   
 $\cup (V\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n}))$   
 $\cup (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n))$ }.

Usually, however, DBF is a proper subset. The above yields the following theorem.

<u>Proof:</u> From [7, p. 468], we have for finite dimensional vector spaces  $V_1$  and  $V_2$ 

(2.13)  $\operatorname{Dim}(\operatorname{span}\{V_1 \cup V_2\}) = \operatorname{Dim}(V_1) + \operatorname{Dim}(V_2) - \operatorname{Dim}(V_1 \cap V_2).$ Therefore, from (2.11) we have

$$(2.14) \quad \text{Dim} (\text{DBF}) \leq \text{Dim} (\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n)) \\ + \text{Dim} (V(\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n})) \\ + \text{Dim} (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n)) \\ - \text{Dim}((\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\overline{\pi}_{y}, N, n)) \cap (V(\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n}))) \\ - \text{Dim}((\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n)) \cap (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n))) \\ - \text{Dim}((\overline{V}(\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n})) \cap (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n))) \\ - \text{Dim}((\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n)) \cap (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n))) \\ + \text{Dim}((\overline{V}(\overline{\pi}_{x}, \overline{M}, \overline{m}) \otimes V(\pi_{y}, N, n)) \cap (V(\pi_{x}, M, m) \otimes \overline{V}(\overline{\pi}_{y}, \overline{N}, \overline{n})) \\ \cap (V(\pi_{x}, M, m) \otimes V(\pi_{y}, N, n))) .$$

The fourth and last terms together are non-positive, hence, can be dropped (if  $m \leq \overline{m}$  or  $n \leq \overline{n}$  then 2.8 would show that they are equal and the bound is best possible). Using (2.2), (2.4) and (2.5) we have our result.

For the proof of the next theorem, we introduce the following notation.

If  $V(\pi_x, M, m)$  is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, \overline{m})$ , then from Definition 1.5, for each i satisfying  $1 \le i \le M$ , there corresponds a unique  $\overline{i}$  such that  $1 \le \overline{i} \le \overline{M}$ ,  $\overline{x_{\overline{i}}} = x_{\overline{i}}$  and  $\alpha(i) = \overline{\alpha}(\overline{i})$ . Thus, we define the index set IM to be the set which contains each of the  $\overline{i}$ satisfying the above conditions. Finally, we define the index set  $\overline{IM}$ to be the set of the remaining  $\overline{M}$ -M integers.

Correspondingly, if  $V(\pi_y, N, n)$  is subordinate to  $(\overline{V}(\overline{\pi}_y, \overline{N}, \overline{n}), \overline{N}, \overline{n})$ , then we define the index sets JN and  $\overline{JN}$  in an analogous way. We will now prove a lemma which gives a lower bound on Dim (DBF), and eventually a basis for DBF.

Lemma 2.1. Given the interpolation spaces  $V(\pi_x, M, m)$   $V(\pi_y, N, n), \overline{V(\pi_x, M, m)}$  and  $\overline{V(\pi_y, N, n)}$  such that  $V(\pi_x, M, m)$ is subordinate to  $\overline{V(\pi_x, M, m)}$  and  $V(\pi_y, N, n)$  is subordinate to  $\overline{V(\pi_y, N, n)}$ , and if

$$\begin{aligned} &\Omega_{1} = \{\overline{\varphi_{i}} \ \psi_{j}\}_{i \in \overline{IM}, i \leq j \leq N}, \quad \Omega_{2} = \{\varphi_{i} \ \overline{\psi_{j}}\}_{j \in \overline{JN}, i \leq i \leq M}, \\ &\Omega_{3} = \{\overline{\varphi_{i}} \ \psi_{j} + P_{x} \ \overline{P}_{y} \left[\overline{\varphi_{i}} \ \psi_{j}\right] - P_{x} \left[\overline{\varphi_{i}} \ \psi_{j}\right]\}_{j \in IM, i \leq j \leq N}, \\ &\Omega_{4} = \{\varphi_{i} \ \overline{\psi}_{j} + \overline{P}_{x} \ P_{y} \left[\varphi_{i} \overline{\psi}_{j}\right] - P_{y} \left[\varphi_{i} \ \overline{\psi}_{j}\right]\}_{j \in JN, i < i \leq M}, \\ &\Omega_{5} = \{P_{y} \left[\overline{\varphi_{i}} \ \overline{\psi}_{j}\right] + P_{x} \left[\overline{\varphi_{i}} \ \overline{\psi}_{j}\right] - P_{x} \ P_{y} \left[\overline{\varphi_{i}} \ \overline{\psi}_{j}\right]\}_{i \in IM, j \in JN}, \\ &\Omega_{6} = \{\overline{P}_{x} \left[\varphi_{i} \ \psi_{j}\right] + \overline{P}_{y} \left[\varphi_{i} \ \psi_{j}\right] - \varphi_{i} \ \psi_{j}\}_{i \leq i \leq M}, i \leq j \leq N, \\ &T_{1} = \Omega_{1} \bigcup \Omega_{2} \bigcup \Omega_{3}, \quad T_{2} = \Omega_{1} \bigcup \Omega_{2} \bigcup \Omega_{4}, \\ &T_{3} = \Omega_{1} \bigcup \Omega_{2} \bigcup \Omega_{5}, \quad T_{4} = \Omega_{1} \bigcup \Omega_{2} \bigcup \Omega_{6}, \end{aligned}$$

then each of the sets  $T_1$ ,  $T_2$ ,  $T_3$  and  $T_4$  is a linearly independent subset of DBF.

<u>Proof</u>: The sets  $Q_3$  through  $Q_6$  are clearly contained in DBF, thus we need only show that  $Q_1$  and  $Q_2$  are subsets of DBF. Because  $\overline{\phi}_i \psi_j \in C^{(m,n)}$ , then for  $\overline{i} \in \overline{IM}$  and  $1 \leq j \leq N$  we have (2.15)  $\overline{P_x \oplus P_y} \left[ \overline{\phi}_{\overline{i}} \psi_j \right] = \overline{\phi}_{\overline{i}} \psi_j + \overline{P}_y P_x \left[ \overline{\phi}_{\overline{i}} \psi_j \right] - P_x \left[ \overline{\phi}_{\overline{i}} \psi_j \right].$ 

We make the observation that for  $1 \leq i \leq M$  and  $i \in IM$ 

(2.16) 
$$\overline{\phi_{i}}^{(\alpha(i))}(x_{i}) = 0$$
,

which implies that  $P_x[\overline{\phi_i} \psi_j] = 0$ . Hence  $Q_1$ , and in a similar manner also  $Q_2$  are contained in DBF. It then follows that  $T_1$  through  $T_4$  are subsets of DBF.

We now show the linear independence of the sets  $T_1$  through  $T_4$ . For notational convenience and also to save space we will only prove the linear independence of  $T_3$ . If  $T_3$  is not linearly independent, then there exists real numbers  $a_{ij}$  for  $i \in \overline{IM}$ ,  $1 \le j \le N$ . also  $b_{ij}$  for  $1 \le i \le M$ ,  $j \in \overline{JN}$  and  $c_{ij}$  for  $i \in IM$  and  $j \in JM$  not all zero such that

(2.17) 
$$\sum_{i \in \overline{IM}} \sum_{1 \leq j \leq N} a_{ij} \overline{\phi}_{i} \psi_{j} + \sum_{l \leq i \leq M} \sum_{j \in \overline{JN}} b_{ij} \phi_{i} \overline{\psi}_{j}$$

$$+ \sum_{i \in IM} \sum_{j \in JN} c_{ij} \left( P_{y} \left[ \overline{\phi_{i}} \overline{\psi_{j}} \right] + P_{x} \left[ \overline{\phi_{i}} \overline{\psi_{j}} \right] - P_{x} P_{y} \left[ \overline{\phi_{i}} \overline{\psi_{j}} \right] \right) = 0.$$

We will show that all  $c_{ij}$  are equal to zero. To accomplish this, apply the cardinality conditions for each  $\hat{i} \in IM$  and  $\hat{j} \in JN$  to (2.17). This means the application of the differential operator  $D^{(\overline{\alpha}(\hat{i}), \overline{\beta}(\hat{j}))}$  to (2.17) and evaluating at the point  $(\overline{x}_{\hat{i}}, \overline{y}_{\hat{j}})$ . Because of the definition of IM and JN, we have for  $i \in IM$  and  $j \in JN$ (2.18)  $\overline{\phi}_{i}^{(\overline{\alpha}(\hat{i}))}(\overline{x}_{\hat{i}}) = 0$  and  $\overline{\psi}_{i}^{(\overline{\beta}(\hat{j}))}(\overline{y}_{\hat{i}}) = 0$ ,

$$(2.19) \qquad (D^{(\overline{\alpha}(\hat{i}),\overline{\beta}(\hat{j}))}(\sum_{i \in IM} \sum_{j \in JN} c_{ij} P_{y}[\overline{\phi}_{i}\overline{\psi}_{j}]))(\overline{x}_{i},\overline{y}_{j}) = c_{ij}^{\wedge \wedge},$$

(2.20) 
$$(D^{(\overline{\alpha}(\hat{i}), \overline{\beta}(\hat{j}))}(\sum_{i \in IM \ j \in JN} \sum_{ij} P_{\mathbf{x}}[\overline{\phi}_{i}\overline{\psi}_{j}])(\overline{\mathbf{x}}_{i}, \overline{\mathbf{y}}_{j}) = c_{ij}^{\wedge \wedge}$$

and

(2.21) 
$$(D^{(\overline{\alpha}(\hat{i}), \overline{\beta}(\hat{j}))}(\sum_{i \in IM} \sum_{j \in JN} c_{ij} P_{x} P_{y}[\overline{\phi_{i}} \overline{\psi_{j}}]))(\overline{x_{j}}, \overline{y_{j}}) = c_{ij}^{A},$$

from which it follows, using (2.17), that  $c_{ij} = 0$  for  $i \in IM$  and  $\hat{j} \in JN$ . To conserve space, we will only show that (2.21) is valid, and merely note that (2.19) and (2.20) are similar.

Because our blending spaces are subordinate, there corresponds an i where  $1 \leq i_0 \leq M$ , and  $j_0$  where  $1 \leq j_0 \leq N$  such that  $\overline{\alpha}(\hat{i}) = \alpha(i_0)$ ,  $\overline{x}_{\hat{1}} = x_{\hat{i}_0}$ ,  $\overline{\beta}(\hat{j}) = \beta(j_0)$  and  $\overline{y}_{\hat{j}} = y_{\hat{j}_0}$ . Then (2.22)  $(D^{(\overline{\alpha}(\hat{i}), \overline{\beta}(\hat{j}))}(\sum_{i \in IM} \sum_{j \in JN} c_{ij} P_x P_y[\overline{\phi}_i \overline{\psi}_j]))(\overline{x}_{\hat{i}}, \overline{y}_{\hat{j}})$   $i \in IM \ j \in JN} \sum_{i \in IM} \overline{\phi}_i^{(\alpha(s))}(x_s) \overline{\psi}_j^{(\beta(t))}(y_t) \phi_s^{(\overline{\alpha}(\hat{i}))}(\overline{x}_{\hat{i}}) \psi_t^{(\overline{\beta}(\hat{j}))}(\overline{y}_{\hat{j}})$  $= c_{\hat{i}\hat{j}},$ 

where we have used the cardinal interpolation conditions and the correspondence given above to obtain our result.

With all  $c_{ij} = 0$ , it is now a simple matter, using the appropriate cardinality conditions, to prove that all  $a_{ij} = 0$  for  $i \in \overline{IM}$ ,  $l \leq j \leq N$ , and all  $b_{ij} = 0$  for  $l \leq i \leq M$  and  $j \in \overline{JN}$ . Therefore,  $T_3$  cannot be linearly dependent, and we have completed our proof.

For discretized blending function spaces which have subordinate interpolation spaces, we now have a lower bound for Dim (DBF).

<u>Theorem 2.2.</u> Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$  such that  $V(\pi_x, M, m)$ is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$  and  $V(\pi_y, N, n)$  is subordinate to  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$ , then (2.23a) Dim (DBF)  $\geq \overline{M} \cdot N + M \cdot \overline{N} - M \cdot N$ (2.23b) Dim (DBF)  $\leq \overline{M} N + M \overline{N} + M \cdot N$   $- N \cdot \text{Dim} (\overline{V}(\overline{\pi}_x, \overline{M}, m) \cap V(\pi_x, M, m))$  $- M \cdot \text{Dim} (\overline{V}(\overline{\pi}_y, \overline{N}, n) \cap V(\pi_y, N, n))$ .

<u>Proof</u>: The upper bound follows from Theorem 2.1 and the lower bound from Lemma 2.1.

We would like to know the exact dimension of DBF, and have a basis for the space. Examination of Theorem 2.2 shows the assumption needed to obtain our result.

<u>Theorem 2.3.</u> Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V(\pi_x, M, m)}$ ,  $V(\pi_y, N, n)$  and  $\overline{V(\pi_y, N, n)}$  such that  $V(\pi_x, M, m)$ is subordinate to  $\overline{V(\pi_x, M, m)}$ ,  $V(\pi_y, N, n)$  is subordinate to  $\overline{V(\pi_y, N, n)}$ ,  $V(\pi_x, M, m) \subseteq \overline{V(\pi_x, M, m)}$ ,  $V(\pi_y, N, n) \subseteq$   $\overline{V(\pi_y, N, n)}$ , and if  $S_1 = \{\overline{\phi_i} \ \psi_j\}_{i \in \overline{IM}}, \ l \leq j \leq N'$ ,  $S_2 = \{\phi_i \ \overline{\psi_j}\}_{j \in \overline{JN}}, \ l \leq i \leq M$ ,

$$\begin{split} \mathbf{S}_{3} &= \{\overline{\boldsymbol{\varphi}_{i}} \ \boldsymbol{\psi}_{j}\}_{i \in \mathrm{IM}, \ 1 \leq j \leq \mathrm{N}}, \ \mathbf{S}_{4} &= \{\boldsymbol{\varphi}_{i} \ \overline{\boldsymbol{\psi}_{j}}\}_{j \in \mathrm{JN}, \ 1 \leq i \leq \mathrm{M}}, \\ \mathbf{S}_{5} &= \{\mathbf{P}_{\mathbf{y}}\left[\overline{\boldsymbol{\varphi}_{i}} \ \overline{\boldsymbol{\psi}_{j}}\right] + \mathbf{P}_{\mathbf{x}}\left[\overline{\boldsymbol{\varphi}_{i}} \ \overline{\boldsymbol{\psi}_{j}}\right] - \mathbf{P}_{\mathbf{x}} \ \mathbf{P}_{\mathbf{y}}\left[\boldsymbol{\varphi}_{i} \ \boldsymbol{\psi}_{j}\right]\}_{i \in \mathrm{IM}, \ j \in \mathrm{JN}}, \\ \mathbf{S}_{6} &= \{\boldsymbol{\varphi}_{i} \ \boldsymbol{\psi}_{j}\}_{1 \leq i \leq \mathrm{M}, \ 1 \leq j \leq \mathrm{N}}, \\ \mathbf{T}_{1} &= \mathbf{S}_{1} \bigcup \mathbf{S}_{2} \bigcup \mathbf{S}_{3}, \qquad \mathbf{T}_{2} = \mathbf{S}_{1} \bigcup \mathbf{S}_{2} \bigcup \mathbf{S}_{4}, \\ \mathbf{T}_{3} &= \mathbf{S}_{1} \bigcup \mathbf{S}_{2} \bigcup \mathbf{S}_{5}, \qquad \mathbf{T}_{4} = \mathbf{S}_{1} \bigcup \mathbf{S}_{2} \bigcup \mathbf{S}_{6}, \end{split}$$

then

$$(2.24) \qquad \text{Dim} (\text{DBF}) = \overline{M} \cdot N + M \cdot \overline{N} - M \cdot N$$

and each of the sets  $T_1$ ,  $T_2$ ,  $T_3$  and  $T_4$  forms a basis for DBF. <u>Proof</u>: From Theorem 2.2, we have that (2.24) is true because of the containment of our blending spaces.

Also, because of the uniqueness of representation, we make the observation that

(2.25a)  $\overline{P}_{y}\left[\overline{\phi}_{i}\psi_{j}\right] = \overline{\phi}_{i}\psi_{j} \text{ for } l \leq i \leq \overline{M}, \ l \leq j \leq N,$ 

(2.25b) 
$$\overline{P}_{\mathbf{x}}\left[\phi_{i} \overline{\psi}_{j}\right] = \phi_{i} \overline{\psi}_{j} \text{ for } 1 \leq i \leq M, \ 1 \leq j \leq \overline{N},$$

and

(2.25c) 
$$\overline{P}_{\mathbf{x}}\left[\phi_{i} \psi_{j}\right] = \overline{P}_{\mathbf{y}}\left[\phi_{i} \psi_{j}\right] = \phi_{i} \psi_{j} \text{ for } 1 \leq i \leq M, \ 1 \leq j \leq N.$$

Using (2.25a), (2.25b) and (2.25c), then in Lemma 2.1,  $Q_3$  through  $Q_6$  reduce to  $S_3$  through  $S_6$  respectively. Thus from Lemma 2.1, it follows that  $T_1$  through  $T_4$  each forms a basis for DBF, and we have proved our result.

<u>Remark</u>:  $S_5$  is not as complicated as it appears. Because our blending interpolation spaces are subordinate, for each  $i \in IM$  and  $j \in JN$ there exists an  $i_0$  where  $1 \le i_0 \le M$  and  $j_0$  where  $1 \le j_0 \le N$  such that

(2.26) 
$$S_5 = \{\overline{\phi}_i \psi_j + \phi_i \overline{\psi}_j - \phi_i \psi_j\}.$$

Often in practice, it is not convenient to work with the cardinal basis for DBF. If one is willing to work with a spanning set rather than a basis, the following theorem is useful.

<u>Theorem 2.4.</u> Given the interpolation spaces  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$  such that  $V(\pi_x, M, m)$ is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  is subordinate to  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$ ,  $V(\pi_x, M, m) \subseteq \overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n) \subseteq$   $\overline{V}(\overline{\pi}_y, \overline{N}, n)$  and if  $\{A_i\}_{i=1}^M$ ,  $\{\overline{A}_i\}_{i=1}^{\overline{M}}$ ,  $\{B_j\}_{j=1}^N$  and  $\{\overline{B}\}_{j=1}^{\overline{N}}$  are any bases for  $V(\pi_x, M, m)$ ,  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  and  $\overline{V}(\overline{\pi}_y, \overline{N}, n)$ , respectively, then (2.27) T<sub>5</sub> =  $\{A_i, \overline{B}_j\}_{1 \le i \le M}$ ,  $1 \le j \le \overline{N} \bigcup \{\overline{A}_i, B_j\}_{1 \le i \le \overline{M}}$ ,  $1 \le j \le N$ 

is a spanning set for DBF.

<u>Proof</u>: We first show that  $T_5 \subseteq DBF$ . For  $1 \le i \le M$  and  $1 \le j \le \overline{N}$  we have  $A_i \xrightarrow{B_j} \in C^{(m,n)}$  and
$$(2.28) \qquad \overrightarrow{P_{x} \oplus P_{y}} \begin{bmatrix} A_{i} \ \overrightarrow{B}_{j} \end{bmatrix} = (\overrightarrow{P_{x}} \ P_{y} + \overrightarrow{P_{y}} \ P_{x} - P_{x} \ P_{y}) \begin{bmatrix} A_{i} \ \overrightarrow{B}_{j} \end{bmatrix}$$
$$= \overrightarrow{P_{x}} \ P_{y} \begin{bmatrix} A_{i} \ \overrightarrow{B}_{j} \end{bmatrix} + A_{i} \ \overrightarrow{B}_{j} - P_{y} \begin{bmatrix} A_{i} \ \overrightarrow{B}_{j} \end{bmatrix}$$
$$= A_{i} \ \overrightarrow{B}_{j},$$

since  $\overline{P}_{\mathbf{x}} \begin{bmatrix} A_i \ \overline{B}_j \end{bmatrix} = A_i \ \overline{B}_j$ , which implies  $A_i \ \overline{B}_j \in \text{DBF}$ . In a similar way, for  $1 \le i \le \overline{M}$  and  $1 \le j \le N$  we have  $\overline{A}_i \ B_j \in \text{DBF}$ , which implies that  $T_5 \subseteq \text{DBF}$ .

Because  $V(\pi_x, M, m) \subseteq \overline{V}(\overline{\pi}_x, \overline{M}, m)$ , it is clear from (2.1), which is the definition of a product space, that

(2.29) 
$$V(\pi_x, M, m) \otimes V(\pi_y, N, n) \subseteq \overline{V}(\overline{\pi}_x, \overline{M}, m) \otimes V(\pi_y, N, n)$$
.

It follows from (2.11), (2.29) and the statement following (2.1) that  $T_5$  is indeed a spanning set for DBF.

.

<u>Remark</u>: For certain spaces, if  $V(\pi_x, M, m) \subseteq \overline{V}(\overline{\pi}_x, \overline{M}, m)$ , then this implies that  $V(\pi_x, M, m)$  is subordinate to  $\overline{V}(\overline{\pi}_x, \overline{M}, m)$ . For example, consider the cubic spline spaces of Example 1.2, if  $S^2(\pi_x) \subseteq S^2(\overline{\pi}_x)$ , then  $S^2(\pi_x)$  is subordinate to  $S^2(\overline{\pi}_x)$ . This is clear, because if  $x_i \neq a$  and  $\neq b$  is a knot of the mesh  $\pi_x$ , then there exists a cubic spline in  $S^2(\pi_x)$  which has a jump discontinuity in the third derivative at  $x_i$ . However, all cubic splines in  $S^2(\overline{\pi}_x)$  are cubic polynomials between the knots of  $\overline{\pi}_x$ , from which it follows that  $x_i$  must be a knot of  $\overline{\pi}_x$ .

If we consider polynomial spaces of Example 1.1, then even if  $V(\pi_x, M, 0) \subseteq \overline{V}(\overline{\pi}_x, \overline{M}, 0)$ , it does not necessarily follow that  $V(\pi_x, M, 0)$  is subordinate to  $\overline{V(\pi_x, M, 0)}$ , because  $\pi_x$  and  $\overline{\pi_x}$  could be different.

Finally, let  $V(\pi_x, M, 0)$  be the space of piecewise linear functions and  $\overline{V}(\pi_x, M, 0)$  be the space of polynomials of degree M-1, then even though  $V(\pi_x, M, 0)$  is subordinate to  $\overline{V}(\pi_x, M, 0)$ , it is not the case that  $V(\pi_x, M, 0) \subseteq \overline{V}(\pi_x, M, 0)$ .

Example: In Theorem 2.1 and 2.2, even if the intersection of our interpolation spaces contains only the zero vector, we will show by the following example that the bounds on the dimension are still sharp. Let

(2.30a) 
$$\pi_{\mathbf{x}} = \mathbf{x}_{1} \in [1, 2], \ \overline{\pi}_{\mathbf{x}} = \overline{\mathbf{x}}_{1} \in [1, 2],$$
  
(2.30b)  $V(\pi_{\mathbf{x}}, 1, 0) = \text{span}\{1\}$  with cardinal basis  $\{1\}$ ,

(2.30c)  $\overline{V}(\overline{\pi}_{x}, 1, 0) = \operatorname{span} \{x\}$  with cardinal basis  $\{x/\overline{x}_{1}\}$ ,

(2.30d) 
$$\pi_{\mathbf{x}} = \mathbf{y}_{1} \epsilon [1, 2], \ \overline{\pi}_{\mathbf{y}} = \overline{\mathbf{y}}_{1} \epsilon [1, 2],$$

(2.30e) 
$$V(\pi_v, 1, 0) = \text{span} \{1\}$$
 with cardinal basis  $\{1\}$ ,

(2.30f)  $\overline{V}(\overline{\pi}_{y}, 1, 0) = \operatorname{span} \{y\}$  with cardinal basis  $\{y/\overline{y}_{i}\}$ , and if  $x_{1} \neq \overline{x}_{1}, y_{1} \neq \overline{y}_{1}$  and  $f \in C([1, 2] \times [1, 2])$ 

then

(2.31) 
$$\overline{\mathbf{P}_{\mathbf{x}} \oplus \mathbf{P}_{\mathbf{y}}} [\mathbf{f}](\mathbf{x}, \mathbf{y}) = \mathbf{f}(\overline{\mathbf{x}}_{1}, \mathbf{y}_{1}) \cdot (\mathbf{x}/\overline{\mathbf{x}}_{1}) \cdot \mathbf{1} + \mathbf{f}(\mathbf{x}_{1}, \overline{\mathbf{y}}_{1}) \cdot \mathbf{1} \cdot (\mathbf{y}/\overline{\mathbf{y}}_{1}) - \mathbf{f}(\mathbf{x}_{1}, \mathbf{y}_{1}) \cdot \mathbf{1} \cdot \mathbf{1}$$

and

therefore

(2.33) Dim (DBF) = 3,

which is the upper bound of Theorem 2.1 and Theorem 2.2.

Using (2.30a) through (2.30f), and if  $x_1 \neq x_1$  but  $y_1 = y_1$ , then from (2.31) we have

(2.34) DBF = 
$$\{a \cdot x + b \cdot (y/y_1 - 1) | a \text{ and } b \text{ are real } \}$$
,

therefore

(2.35) Dim (DBF) = 2.

Finally, if (2.30a) through (2.30f) hold, where  $V(\pi_x, 1, 0)$ and  $V(\pi_y, 1, 0)$  are subordinate to  $\overline{V}(\overline{\pi}_x, 1, 0)$  and  $\overline{V}(\overline{\pi}_y, 1, 0)$ , respectively, (i.e.,  $\overline{x}_1 = x_1$  and  $y_1 = \overline{y}_1$ ), then

(2.36) DBF = 
$$\{a(x/x_1 + y/y_1 - 1) | a \text{ is real} \},$$

and

(2.37) Dim (DBF) = 1,

which is the lower bound of Theorem 2.2.

<u>Remark</u>: If  $V(\pi_x, M, m)$ ,  $\overline{V(\pi_x}, \overline{M}, m)$ ,  $V(\pi_y, N, n)$  and  $\overline{V(\pi_y}, \overline{N}, n)$  are the spaces of cubic splines  $S^2(\pi_x)$ ,  $S^2(\overline{\pi_x})$ ,  $S^2(\pi_y)$ and  $S^2(\overline{\pi_y})$ , respectively, then a convenient basis with which to work is the basis of "B-splines" (see [29, p. 73]), which have support on at most four consecutive intervals and are non-negative. The B-spline basis is also relatively easy to construct when compared to the cardinal basis. In this case it is often advantageous to work with a spanning set of bicubic B-splines defined in Theorem 2.4, and use methods for overdetermined systems to solve the resulting equations.

### Section 3. Natural Cubic Blending

Often, information about the normal derivatives around the boundary does not exist, or it is not known to sufficient accuracy to be compatible with an eighth order method. In the case where some deterioration of accuracy is acceptable near the boundary of our domain, we can use natural cubic blending to obtain a method which is  $O(h^8)$  in the interior.

The natural cubic spline basis  $\{A_i\}_{i=0}^{M+1}$  for  $S^2(\pi_x) \subseteq C^{(2)}[a,b]$ on the mesh  $\pi_x:a = x_1 < x_2 < \cdots < x_M = b$  satisfies for  $l \leq i, j \leq M$  the conditions

(3.1a) 
$$A_i(x_j) = \delta_{ij}, A''_i(x_1) = A''_i(x_M) = 0$$
,

(3.1b) 
$$A_0(x_j) = 0$$
,  $A_0''(x_1) = 1$ ,  $A_0''(x_M) = 0$ 

(3.1c) 
$$A_{M+1}(x_j) = 0, A_{M+1}'(x_1) = 0, A_{M+1}''(x_M) = 1$$

The natural cubic spline interpolant to a function  $f \in C[a, b]$  is defined to be

(3.2) 
$$s(x) = \sum_{i=1}^{M} f(x_i) A_i(x) ,$$

where it is clear that  $s''(x_1) = s''(x_M) = 0$ .

The following theorem, which is crucial to the proof of Theorem 3.2, shows the behavior of the natural cubic spline interpolant of one variable, see Hall [19]. <u>Theorem 3.1.</u> (Hall [19]). Let s be the natural cubic spline interpolant to  $f \in C^{m}[a, b]$ , m = 2, 3 or 4, then for  $x \in [x_{i}, x_{i+1}]$ ,  $1 \leq i \leq M-1$ 

(3.3) 
$$|(\mathbf{s}-\mathbf{f})^{(\mathbf{k})}(\mathbf{x})| \leq ||\mathbf{f}^{(\mathbf{m})}|| \{\mathcal{E}_{\mathbf{mk}} \mathbf{h}_{\mathbf{x}}^{\mathbf{m}-\mathbf{k}} + \mathbf{K}_{\mathbf{m}} \mathbf{h}_{\mathbf{x}}^{\mathbf{m}-1} \alpha_{\mathbf{k}} \Delta_{\mathbf{i}} \}$$
  
+  $1/2 \mathbf{R} \mathbf{h}_{\mathbf{x}} \alpha_{\mathbf{k}} \Delta_{\mathbf{i}}$ ,

(3.4)  $R = \max \{ |f'(a)|, |f'(b)| \}$ 

for  $0 \leq k \leq 2$ , where the mesh size  $h_x$  is defined to be

(3.5) 
$$h_{\mathbf{x}} = \max_{1 \leq i \leq M-1} \Delta \mathbf{x}_{i}, \Delta \mathbf{x}_{i} = \mathbf{x}_{i+1} - \mathbf{x}_{i}, \Delta_{i} = \{2^{1-i} + 1\}$$

 $2^{1-M+i}$ , and the constants  $\mathcal{E}_{m,k}$ ,  $K_{m}$  and  $\alpha_{k}$  are given in the following tables.

Table 3

$\mathcal{E}_{mk}$	<b>k</b> = 0	k = 1	k = 2 .
m = 2	9/8	4	10
m = 3	71/216	31/27	5
m = 4	5/384	(9 + <del>√3</del> )/216	5

Table 4

α <sub>k</sub>	k = 0	k = 1	k = 2	
	$\Delta x_i/4$	1	6/∆× <sub>i</sub>	

Km	m = 4	m = 3	m = 2			
	7/24	1	5/2			
Table 6						
β <sub>l</sub>	$\ell = 0$	<b>f</b> = 1	<i>l</i> = 2			
	$\Delta y_{i}/4$	1	6/∆y <sub>j</sub>			

We now generalize this theorem to two variables. For a mesh on the y - axis  $\pi_y:c = y_1 < y_2 < \cdots < y_N = d$ , let  $S^2(\pi_y) \subseteq C^{(2)}[c, d]$ be the space of cubic splines in the variable y. The natural cubic spline basis  $\{B_j\}_{j=0}^{N+1}$  for  $S^2(\pi_y)$  satisfies conditions similar to (3.1a) - (3.1c). Define the projection operators  $P_x$  and  $P_y$  by (3.6)  $(P_x[f])(x, y) = \sum_{i=1}^{M} f(x_i, y) A_i(x)$ 

and

(3.7) 
$$(P_{y}[f])(x, y) = \sum_{j=1}^{N} f(x, y_{j}) B_{j}(y),$$

for each  $f \in C([a, b] \times [c, d])$ . The natural cubic spline blended interpolant to f is defined to be

(3.8) 
$$NB[f] = P_x[f] + P_y[f] - P_x P_y[f].$$

Table 5

Define 
$$h = \max_{y_{j} \in N-1} \Delta y_{j}$$
, where  $\Delta y_{j} = y_{j+1} - y_{j}$ 

We have the following theorem.

Theorem 3.2: Let NB be the natural cubic spline blended interpolant to  $f \in C^{(m,n)}([a,b]x[c,d])$ . If  $x \in [x_i, x_{i+1}]$ ,  $y \in [y_j, y_{j+1}]$ ,  $1 \leq i \leq M-1$ ,  $1 \leq j \leq N-1$ ,  $2 \leq m, n \leq 4$  and  $0 \leq k, l \leq 2$ , then  $|(((I-NB)[f])^{(k, \ell)})(x, y)|$ (3.9) $\leq ||\mathbf{f}^{(\mathbf{m},\mathbf{n})}|| \{ \mathcal{E}_{\mathbf{m}\mathbf{k}} \mathbf{h}_{\mathbf{x}}^{\mathbf{m}-\mathbf{k}} + \mathbf{K}_{\mathbf{m}} \mathbf{h}_{\mathbf{x}}^{\mathbf{m}-1} \alpha_{\mathbf{k}} \Delta_{\mathbf{i}} \} \{ \mathcal{E}_{\mathbf{n}\boldsymbol{\ell}} \mathbf{h}_{\mathbf{y}}^{\mathbf{n}-\boldsymbol{\ell}} \}$ +  $K_{n} h_{v}^{n-1} \beta_{l} \overline{\Delta}_{j}$  + 1/2 ||  $f^{(m, 2)}$  || { $\mathcal{E}_{mk} h_{x}^{m-k}$ +  $K_{\mathbf{n}} h_{\mathbf{x}}^{\mathbf{m}-1} \alpha_{\mathbf{k}} \Delta_{\mathbf{i}} \{ \overline{\Delta}_{\mathbf{i}} \beta_{\boldsymbol{\ell}} h_{\mathbf{v}} + 1/2 | | \mathbf{f}^{(2,\mathbf{n})} | | \{ \mathcal{E}_{\mathbf{n}\boldsymbol{\ell}} h_{\mathbf{v}}^{\mathbf{n}-\boldsymbol{\ell}} \}$ +  $K_n h_y^{n-1} \beta_{\ell} \overline{\Delta}_j$   $\Delta_i \alpha_k h_x + 1/4 || f^{(2,2)} || \Delta_i \overline{\Delta}_j \alpha_k \beta_{\ell} h_x h_y$ , where  $\Delta_i = \{2^{1-i} + 2^{1-M+i}\}, \ \overline{\Delta}_i = \{2^{1-j} + 2^{1-N+j}\}$ ,  $\mathcal{E}_{mk}$  and  $\mathcal{E}_{nl}$  are given in Table 3, K and K are given in Table 5,  $\alpha_k$  is given in Table 4,  $\beta_l$  is given in Table 6 and  $h_x$ and h are the mesh sizes of  $\pi_x$  and  $\pi_y$ , respectively.

<u>Proof</u>: From Section 1 we have  $I - NB = R_x R_y$ , where  $I - P_x = R_x$ and  $I - P_y = R_y$ . Also from (1.17b) and (1.17c) we have  $D^{(0,l)}R_x = R_x D^{(0,l)}$  and  $D^{(k,0)}R_y = R_y D^{(k,0)}$ . Define  $g = D^{(0,l)}R_y[f]$ . Then, for fixed x and y,

$$(3.10) \qquad \left| (((I - N B) [f])^{(k, \ell)})(x, y) \right| = \left| ((R_x R_y [f])^{(k, \ell)})(x, y) \right| \\ = \left| (D^{(k, 0)} R_x D^{(0, \ell)} R_y [f])(x, y) \right| \\ = \left| (D^{(k, 0)} R_x g) (x, y) \right| \\ = \left| (E(x, y) \right| .$$

For each fixed y, we consider  $P_{\mathbf{x}}[g](\cdot, \mathbf{y})$ , which is the natural cubic spline interpolant to the univariable function  $g(\cdot, \mathbf{y})$ , where  $g(\cdot, \mathbf{y}) \in C^{\mathbf{m}}[a, b]$ . From Theorem 3.1 we have

(3.11) 
$$|E(\mathbf{x}, \mathbf{y})| \leq ||g^{(m, 0)}(\cdot, \mathbf{y})|| \{\mathcal{E}_{mk} h_{\mathbf{x}}^{m-k} + K_{m} h_{\mathbf{x}}^{m-1} \alpha_{k} \Delta_{i} \}$$
  
+  $1/2 R_{1}(\mathbf{y}) \Delta_{i} \alpha_{k} h_{\mathbf{x}},$ 

where

(3.12) 
$$R_{1}(y) = \max \{ |g^{(2,0)}(a,y)|, |g^{(2,0)}(b,y)| \}$$

From the definition of g, we have for  $0 \leq t \leq m$ ,

(3.13) 
$$g^{(t,0)} = D^{(t,0)}(D^{(0,\ell)}R_{y}[f])$$
  
=  $D^{(0,\ell)}R_{y}[f^{(t,0)}].$ 

Therefore, applying Theorem 3.1 again, for any fixed  $\xi$  such that  $a \leq \xi \leq b$ , we have

$$(3.14) |g^{(t,0)}(\xi,y)| \leq ||f^{(t,n)}(\xi,\cdot)|| \{ \mathcal{E}_{n\ell} h_y^{n-\ell} + K_n h_y^{n-1} \beta_{\ell} \overline{\Delta}_j \} + 1/2 R_2(\xi) \overline{\Delta}_j \beta_{\ell} h_y,$$

where

(3.15) 
$$R_{2}(\xi) = \max \{ |f^{(t,2)}(\xi,c)|, |f^{(t,2)}(\xi,d)| \}$$
  
 $\leq ||f^{(t,2)}||.$ 

Taking the supremum over all  $\xi$  on the right of inequality (3.14) we have

(3.16) 
$$|g^{(t,0)}(\xi,y)| \leq ||f^{(t,n)}|| \{\mathcal{E}_{n\ell} h_{y}^{n-\ell} + K_{n} h_{y}^{n-1} \beta_{\ell} \overline{\Delta}_{j} + 1/2 ||f^{(t,2)}|| \overline{\Delta}_{j} \beta_{\ell} h_{y}.$$

Substituting (3.16), for t = m and t = 2, into (3.11) and (3.12) respectively, our conclusion follows after regrouping.

We will now make a series of remarks about Theorem 3.2.

Remark 3.1: The error still goes to zero even if only one of our meshes,  $\pi_x$  or  $\pi_y$ , is refined.

<u>Remark 3.2</u>: Because of the exponential decay of the terms  $\Delta_i$  and  $\overline{\Delta}_j$ , if  $a < a_1 < b_1 < b$  and  $c < c_1 < d_1 < d$ , then there exists a mesh fine enough so that the convergence on the subrectangle  $[a_1, b_1] \times [c_1, d_1]$  is  $O(h^{m+n-\ell-k})$ , where  $h = \max\{h_x, h_y\}$ .

<u>Remark 3.3</u>: Moreover, because the exponential decay of the terms  $\Delta_i$  and  $\overline{\Delta}_j$  tend toward zero faster than the term  $1/h^2 = (N-1)^2$ 

tends toward infinity, we have the area of higher order convergence increasing as we refine our meshes  $\pi_x$  and  $\pi_y$ .

Remark 3.4: Explicitly, the rate of increase is given by

$$a_1 - a = b - b_1 = c_1 - c = d - d_1 = (K \log (N-1))/(N-1) \longrightarrow 0$$

as  $N \longrightarrow \infty$ , where K is chosen with consideration of the bound on the derivatives of f.

<u>Remark 3.5</u>: For purposes of illustration, we will consider the special case of k = l = 0, a = c = 0, b = d = 1, M = N and m = n = 4. Therefore, both  $\alpha_0$  and  $\beta_0$  have a factor of h, and because of the exponential decay of the terms  $\Delta_i$  and  $\overline{\Delta}_j$  as we move away from the boundary, we have the following situation for our error, illustrated in Figure 1.



### Section 4. Exponential Decay of Natural Cubic Cardinal Splines

Given a mesh  $\pi_x$ :  $a = x_1 < x_2 < \cdots < x_M = b$  of M knots, the natural cardinal splines  $A_i(x) \in S^2(\pi_x)$  are uniquely determined by the M + 2 conditions given in (3.1a) to (3.1c).

We will now prove the exponential decay of the natural cubic cardinal splines by a series of lemmas and theorems. Much of what follows parallels the results of Birkhoff and De Boor [2], for cardinal splines with first derivative end conditions, and we will quote theorems from their paper which are also valid for natural splines. Finally, we will use the abbreviation M.V.T. for the Mean Value Theorem and I.V.T. for the Intermediate Value Theorem.

<u>Lemma 4.1</u>. (Birkhoff and De Boor [2]). If p(x) is a cubic polynomial which vanishes at 0 and  $h \neq 0$  then

- (4.1a) p'(h) = -2p'(0) h p''(0)/2
- (4.1b) p''(h)/2 = -3 p'(0)/h p''(0).

Corollary 4.1. For  $1 \le i \le M$ , the natural cubic cardinal splines satisfy

(4.2a) 
$$A'_{i}(x_{j+1}) = -2A'_{i}(x_{j}) - \Delta x_{j}A''_{i}(x_{j})/2$$

(4.2b) 
$$A_{i}^{\prime\prime}(x_{j+1})/2 = -3A_{i}^{\prime}(x_{j})/\Delta x_{j} - A_{i}^{\prime\prime}(x_{j}),$$

where  $l \leq j < i-1$ , and

(4.3a) 
$$A'_{i}(x_{j-1}) = -2 A'_{i}(x_{j}) + \Delta x_{j-1} A''_{i}(x_{j})/2$$

(4.3b) 
$$A_{i}''(x_{j-1})/2 = 3 A_{i}'(x_{j})/\Delta x_{j-1} - A_{i}''(x_{j})$$

where  $i+1 < j \leq M$ .

Define the function  $sgn(\cdot)$  of a real variable by

(4.4) 
$$sgn(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}$$

Lemma 4.2. (Birkhoff and De Boor  $\begin{bmatrix} 2 \end{bmatrix}$ ). Let s(x) be any cubic spline function with knots at the  $x_j$ , which satisfy for some i,  $s(x_{i-1}) = s(x_{i+1}) = 0$ ,  $s(x_i) > 0$ ,  $s'(x_{i-1}) s''(x_{i-1}) \ge 0$ ,  $s'(x_{i+1})$  $s''(x_{i+1}) \le 0$ . Then  $s'(x_{i-1}) \ge 0$ ,  $s''(x_{i-1}) \ge 0$ ,  $s'(x_{i+1}) \le 0$ ,  $s''(x_{i+1}) \ge 0$ ,  $s''(x_i) \le 0$  and  $s(x) \ge 0$  on  $[x_{i-1}, x_{i+1}]$ .

Lemma 4.3. For  $2 \le i \le M-1$ ,  $A_i(x)$  satisfies Lemma 4.2 on the interval  $[x_{i-1}, x_{i+1}]$ .

**Proof:** From Corollary 4.1 and the fact that  $A_i''(x_1) = A_i''(x_M) = 0$ .

Lemma 4.4. For  $1 \le i \le M$  the natural cubic cardinal splines satisfy

$$(4.5a) \qquad \operatorname{sgn} (A_i'(x_j)) = \operatorname{sgn} (A_i''(x_j)) \neq 0$$

(4.5b) 
$$\operatorname{sgn} (A_i'(x_j)) = -\operatorname{sgn} (A_i'(x_{j-1})) \neq 0$$
,

where  $2 \leq j \leq i-1$ , and for  $i+1 \leq j \leq M-1$ 

(4.6a) 
$$\operatorname{sgn} (A_{i}'(x_{j})) = -\operatorname{sgn} (A_{i}''(x_{j})) \neq 0$$

(4.6b) 
$$\operatorname{sgn}(A_i'(x_j)) = -\operatorname{sgn}(A_i'(x_{j+1})) \neq 0$$

# **Proof:** Case 1: $2 \leq i \leq M-1$ .

We will first prove that  $A'_i(x_1)$  and  $A'_i(x_N)$  are both non-zero. If this is not the case, say  $A'_i(x_1) = 0$ , then inductively  $A_i(x) \equiv 0$ on  $[x_1, x_{i-1}]$ . Therefore, on  $[x_{i-1}, x_i]$ ,  $A(x_{i-1}) = A'(x_{i-1}) =$   $A''(x_{i-1}) = 0$ ,  $A_i(x_i) = 1$  and from Lemmas 4.2 and 4.3 we also have  $A_i''(x_i) < 0$ . Several applications of the M.V.T. gives  $\xi \in (x_{i-1}, x_i)$ such that  $A''_i(\xi) > 0$  and the I.V.T. gives another distinct zero for  $A''_i$ . This implies that  $A_i$  is linear on  $[x_{i-1}, x_i]$  because  $A''_i \equiv 0$ on  $[x_{i-1}, x_i]$ . But this is a contradiction to the fact that  $A_i$  must satisfy the conditions  $A_i(x_{i-1}) = A'_i(x_{i-1}) = 0$  and  $A_i(x_i) = 1$ . An identical argument shows that  $A'_i(x_M) \neq 0$ . The result now follows inductively from Corollary 4.1.

Case 2: i = 1, M.

For i = M, if  $A'_{M}(x_{1}) = 0$  then  $A'_{M}(x) \equiv 0$  on  $[x_{1}, x_{M-1}]$ . We then know that  $A''_{M}$  has zeros at  $x_{M-1}$  and  $x_{M}$ , showing that  $A_{M}$  cannot satisfy all of the required conditions, therefore  $A'_{M}(x_{1}) \neq 0$ . An identical argument shows that  $A'_{1}(x_{M}) \neq 0$ . The result now follows inductively from Corollary 4.1.

Lemma 4.5. If 
$$1 \leq j \leq i-2$$
 and  $x \in [x_j, x_{j+1}]$ , then  $A'_i(x_{j+1}) A_i(x) \leq 0$  and if  $i+1 \leq j \leq M-1$  for  $x \in [x_j, x_{j+1}]$ , then  $A'_i(x_j) A_i(x) \geq 0$ .

<u>Proof</u>: From Lemma 4.4,  $\operatorname{sgn}(A'_i(x_j)) = -\operatorname{sgn}(A'_i(x_{j+1})) \neq 0$ . It is now clear that if the conclusion of the lemma is not true, then  $A_i$ will have four distinct zeros in  $[x_j, x_{j+1}]$ , which is false. The other case is proved in a similar way.

Lemma 4.6. For  $1 \leq i \leq M$  we have  $A'_i(x_{i-1}) > 0$  and  $A'_i(x_{i+1}) < 0$ .

<u>Proof</u>: Case 1:  $2 \leq i \leq M-1$ .

The result follows from Lemmas 4.2 through 4.4.

Case 2: i = 1, M.

For i = M,  $A_M(x_M) = 1$ ,  $A_M''(x_M) = 0$ . If  $A_M'(x_{M-1}) < 0$ , we have from Lemma 4.4 that  $A_M''(x_{M-1}) < 0$ . Two applications of the M. V. T. and one of the I. V. T. yields another zero for  $A_M''$  which cannot happen if  $A_M$  is to satisfy the required conditions at  $x_{M-1}$  and  $x_M$ . The case where i = 1 follows in the same manner.

8

.

<u>Lemma 4.7</u>. On  $[x_1, x_2]$ , max  $|A_1(x)| = 1$ , and on  $[x_{M-1}, x_M]$ max  $|A_M(x)| = 1$ .

<u>Proof</u>: On  $[x_{M-1}, x_M]$ , direct calculation yields a real number a such that  $0 > a > -1/(2\Delta x_{M-1}^3)$  and

(4.7) 
$$A_{M}(x) = a(x-x_{M})^{3} + (1-a(\Delta x_{M-1})^{3})(x-x_{M})/\Delta x_{M-1} + 1$$
.

Examination of (4.7) shows that  $A_M$  has no interior relative maximums or minimums on  $[x_{M-1}, x_M]$ . The bound for  $A_1$  follows in the same way.

We will now give a proof similar to that of Birkhoff and De Boor [2] for natural cardinal cubic splines.

Lemma 4.8. For 
$$1 \leq j \leq i-2$$
, if  $x_j \leq x \leq x_{j+1}$ , then  $|A_i(x)| \leq \Delta x_j |A_i'(x_{j+1})|$  and if  $i+1 \leq j \leq M-1$  for  $x_j \leq x \leq x_{j+1}$ , then  $|A_i(x)| \leq \Delta x_j |A_i'(x_j)|$ .

Proof: Case 1:  $l \leq j \leq i-2$ .

Without loss of generality assume  $A_i'(x_{j+1}) < 0$ . From Lemma 4.5  $A_i(x) \ge 0$  on  $[x_j, x_{j+1}]$ , and from Lemma 4.4  $A_i''(x_{j+1}) < 0$ ,  $A_i'(x_j) \ge 0$ ,  $A_i''(x_j) \ge 0$ . Our proof is by contradiction. Let  $\xi_1$  be the point where  $A_i$  obtains its absolute maximum on  $[x_j, x_{j+1}]$ , then  $A_i'(\xi_1) = 0$ . Assume  $A_i(\xi_1) \ge -A_i'(x_{j+1})\Delta x_j \ge 0$ , then because  $A_i(x_{j+1}) = 0$ , we have from the M. V. T. a  $\xi_2$ satisfying  $\xi_1 < \xi_2 < x_{j+1}$  such that  $A_i'(\xi_2) = -A_i(\xi_1)/\Delta x_j < A_i'(x_{j+1}) < 0$ . Applying the M. V. T. again we have the existence of  $\xi_3$  satisfying  $\xi_2 < \xi_3 < x_{j+1}$  such that  $A_i''(\xi_3) > 0$ . Because  $A_i''(x_{j+1}) < 0$ , the I. V. T. gives a  $\xi_4$  satisfying  $\xi_3 < \xi_4 < x_{j+1}$  such that  $A_i''(\xi_1) = 0$  and  $A_i''(x_j) \ge 0$  ( $A_i''(x_1) = 0$ ), we have a second zero of  $A_i''$ , which implies that  $A_i$  is linear on  $[x_j, x_{j+1}]$ . This contradicts the fact that  $A_i$  must satisfy  $A_i(x_j) = A_i(x_{j+1}) = 0$  and  $A_i'(x_j) > 0$ .

The proof of the other case is identical.

Paralleling the results of Birkhoff and De Boor [2], we have the following two lemmas for natural cubic cardinal splines.

Lemma 4.9: For  $1 \leq j \leq i-2$   $|A_i'(\mathbf{x}_j)| \leq 1/2 |A_i'(\mathbf{x}_{j+1})|$ , and for  $i+1 \leq j \leq M-1$  $|A_i'(\mathbf{x}_{j+1})| \leq 1/2 |A_i'(\mathbf{x}_j)|$ .

Proof: From Corollary 4.1.

Lemma 4.10. For  $1 \leq j \leq i-2$ , on  $[x_j, x_{j+1}]$  we have  $|A_i(x)| \leq 2^{j-i+2} |A_i'(x_{i-1})| \Delta x_j$ , and for  $i+2 \leq j \leq M$ , on  $[x_{j-1}, x_j]$  we have

$$|A_{i}(\mathbf{x})| \leq 2^{i+2-j} |A_{i}'(\mathbf{x}_{i+1})| \Delta \mathbf{x}_{j-1}$$

Proof: Follows inductively from Lemmas 4.8 and 4.9.

We have now shown through a sequence of Lemmas that the natural cubic cardinal splines behave in the same manner as the splines which have zero derivatives as their boundary conditions. Therefore, the proof of the following Lemma is now identical to that given by Birkhoff and De Boor [2], and we will only state their conclusion.

Lemma 4.11. (Birkhoff and De Boor  $\begin{bmatrix} 2 \end{bmatrix}$ ). For  $2 \le i \le M-1$ on  $[\mathbf{x}_{i-1}, \mathbf{x}_{i+1}]$  we have  $0 \le A_i(\mathbf{x}) \le L$  and  $|A_i'(\mathbf{x}_{i-1})| \le L/\Delta \mathbf{x}_{i-1}$ ,  $|A'(\mathbf{x}_{i+1})| \le L/\Delta \mathbf{x}_i$ , where  $L = 3 M_{\pi_X} (M_{\pi_X} + 1)^2 / (3 + 4M_{\pi_X})$ , the mesh ratio is  $M_{\pi_X} = \max \Delta \mathbf{x}_i / \min \Delta \mathbf{x}_i$ .

We will now show similar bounds for  $|A'_{1}(x_{2})|$  and  $|A'_{N}(x_{M-1})|$ by constructing a majorant,  $U_{i}(x)$ , for the end splines. Consider first the spline  $A_{M}$ , then define the unique cubic spline  $U_{M}$  on  $[x_{M-2}, x_{M}]$  satisfying  $U_{M}(x_{M-2}) = U(x_{M-1}) = U''_{M}(x_{M-2}) =$  
$$\begin{split} & U_{M}^{''}(\mathbf{x}_{M}) = 0 \quad \text{and} \quad U_{M}^{'}(\mathbf{x}_{M}) = 1. \quad \text{Define the cubic spline T on} \\ & \left[\mathbf{x}_{M-2}, \mathbf{x}_{M}\right] \quad \text{by } \mathbf{T}(\mathbf{x}) = (\mathbf{U}_{M} - \mathbf{A}_{M})(\mathbf{x}), \text{ then we will prove } \mathbf{T}'(\mathbf{x}_{M-1}) \\ & \geqslant 0. \quad \text{If } M-2 = 1, \text{ then } \mathbf{T}''(\mathbf{x}_{M-2}) = 0 \quad \text{and we are done. If } M-2 > 1, \\ \text{then } \mathbf{T}''(\mathbf{x}_{M-2}) \neq 0, \text{ and we will show that } \mathbf{T}'(\mathbf{x}_{M-1}) > 0. \quad \text{To do this,} \\ \text{we use the fact that } \mathbf{T}''(\mathbf{x}_{M-2}) = 0 - \mathbf{A}''_{M}(\mathbf{x}_{M-2}) > 0, \text{ which follows} \\ \text{from Lemmas } 4.4 \text{ and } 4.6, \text{ and that } \mathbf{T} \text{ satisfies } \mathbf{T}(\mathbf{x}_{M}) = \mathbf{T}''(\mathbf{x}_{M}) = \\ \mathbf{T}(\mathbf{x}_{M-1}) = \mathbf{T}(\mathbf{x}_{M-2}) = 0. \quad \text{Also, } \mathbf{T} \text{ satisfies the conditions of} \\ \text{Lemma } 4.1 \text{ and therefore } (4.3a) \text{ and } (4.3b), \text{ for } \mathbf{x}_{M-2}, \mathbf{x}_{M-1}, \mathbf{x}_{M} \\ \text{Because } \mathbf{T}''(\mathbf{x}_{M}) = 0 \quad \text{and } \mathbf{T}''(\mathbf{x}_{M-2}) > 0, \text{ we have that } \mathbf{T}'(\mathbf{x}_{M}) \neq 0. \\ \text{From } (4.3a) \text{ and } (4.3b) \text{ we have that } \text{sgn}(\mathbf{T}'(\mathbf{x}_{j})) = -\text{sgn}(\mathbf{T}'(\mathbf{x}_{j})) \\ \neq 0 \quad \text{for } \mathbf{j} = \mathbf{M}, \quad \mathbf{M} - 1 \quad \text{and } \text{sgn}(\mathbf{T}''(\mathbf{x}_{j})) = -\text{sgn}(\mathbf{T}'(\mathbf{x}_{j})) \neq 0 \quad \text{for } \mathbf{j} = \mathbf{M} - 2, \\ \mathbf{M} - 1. \quad \text{Therefore, because } \mathbf{T}''(\mathbf{x}_{M-1}) > 0. \quad \text{Using Lemma } 4.6, \text{ and our above} \\ \text{results, it follows that } U'_{M}(\mathbf{x}_{M-1}) \geqslant \mathbf{A}'_{M}(\mathbf{x}_{M-1}) > 0. \end{split}$$

 $U_{M}(x)$  is given explicitly by  $p_{1}(x)$  on  $\begin{bmatrix} x_{M-2}, x_{M-1} \end{bmatrix}$  and  $p_{2}(x)$  on  $\begin{bmatrix} x_{M-1}, x_{M} \end{bmatrix}$ , where

$$p_{1}(\mathbf{x}) = b(\mathbf{x} - \mathbf{x}_{M-1})(\mathbf{x} - \mathbf{x}_{M-2})(\mathbf{x} - \mathbf{x}_{M-1} + 2\Delta \mathbf{x}_{M-2})$$

$$p_{2}(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_{M-1})(-b\Delta \mathbf{x}_{M-2} [(\mathbf{x} - \mathbf{x}_{M})^{2} - \Delta \mathbf{x}_{M-1}(\mathbf{x} - \mathbf{x}_{M})/\Delta \mathbf{x}_{M-1} + 1/\Delta \mathbf{x}_{M-1})$$

where  $b = 1/[2\Delta x_{M-2}\Delta x_{M-1}(\Delta x_{M-1} + \Delta x_{M-2})]$ 

and  $U'_{M}(\mathbf{x}_{M-1}) = \Delta \mathbf{x}_{M-2} / (\Delta \mathbf{x}_{M-1} \left[ \Delta \mathbf{x}_{M-1} + \Delta \mathbf{x}_{M-2} \right])$ .

Therefore,

(4.8) 
$$0 < A'_{M}(x_{M-1}) \leq M_{\pi_{x}} / (\Delta x_{M-1} [1 + M_{\pi_{x}}]).$$

A bound for  $A'_1(x_2)$  follows in an identical way. The above yields the following theorem.

Theorem 4.1. If 
$$1 \le i \le M$$
,  $1 \le j \le i-2$  and  $x_j \le x \le x_{j+1}$ , then  
 $|A_i(x)| \le 2^{2+j-i} L M_{\pi_x}$ .

If  $i+2 \leq j \leq M$  and  $x_{j-1} \leq x \leq x_j$ , then  $|A_i(x)| \leq 2^{2+i-j} LM_{\pi_x}$ .

For 
$$2 \leq i \leq M-1$$
,  $|A_i(x)| \leq L$  on  $[x_{i-1}, x_{i+1}]$ ,  $|A_1(x)| \leq L$   
L on  $[x_1, x_2]$  and  $|A_M(x)| \leq L$  on  $[x_{M-1}, x_M]$ , where  
 $L = 3M_{\pi_x} (M_{\pi_x} + 1)^2 / (3 + 4M_{\pi_x})$  and  $M_{\pi_x}$  is the mesh ratio.

<u>Proof</u>: If  $2 \le i \le M-1$ , the result follows from Lemmas 4.10 and 4.11. For i = 1, M our result follows from Lemmas 4.7 and 4.10 and inequality (4.8), where we have used the liberal estimate that  $M_{\pi_{X}}/(1 + M_{\pi_{X}}) \le L$  and  $1 \le L$ .

We are now in a position to develop the error estimate for approximate natural spline blending. Usually, it is not possible to use all of the values of f along the mesh lines. Therefore, it is reasonable to replace the  $f(x_i, y)$  and  $f(x, y_j)$  with appropriate approximations  $p_i(y)$  and  $q_j(x)$ . The type of approximation is arbitrary as long as the overall accuracy of the scheme is preserved. For example, typical choices of approximation could be univariate spline or polynomial interpolation. We then blend these functions with natural splines to obtain an approximation to f. We would expect this approximation to be close to the original function f. Exactly how close is shown in the following theorem.

We are given the meshes  $\pi_x : a = x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y : c = y_1 < y_2 < \cdots < y_N = d$  and the natural cubic cardinal splines  $\{A_i\}_{i=1}^M \subseteq S^2(\pi_x)$  and  $\{B_j\}_{j=1}^N \subseteq S^2(\pi_y)$ .

<u>Theorem 4.2.</u> Let  $f \in C^{(m,n)}([a,b]x[c,d])$ , where  $2 \leq m, n \leq 4$ , and if the approximations to the mesh functions satisfy  $||f(x_i, \cdot) - p_i(\cdot)|| \leq \mathcal{E}_x$  and  $||f(\cdot, y_j) - q_j(\cdot)|| \leq \mathcal{E}_y$  for  $1 \leq i \leq M$  and  $1 \leq j \leq N$ , then the natural blending function approximation defined by

(4.9) 
$$\overline{\mathbf{s}}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{M} p_i(\mathbf{y}) A_i(\mathbf{x}) + \sum_{j=1}^{N} q_j(\mathbf{x}) B_i(\mathbf{y})$$
  
 $- \sum_{i=1}^{M} \sum_{j=1}^{N} (1/2 p_i(\mathbf{y}_j) + 1/2 q_j(\mathbf{x}_i)) A_i(\mathbf{x}) B_j(\mathbf{y})$   
 $i=1 \ j=1$ 

satisfies

(4.10) 
$$|(f-\bar{s})(x,y)| \leq |(I-NB)[f](x,y)| + \xi_{x}K(M_{\pi})$$

+
$$\mathcal{E}_{\mathbf{y}} \operatorname{K}(\mathbf{M}_{\mathbf{y}})$$
 +  $1/2(\mathcal{E}_{\mathbf{x}} + \mathcal{E}_{\mathbf{y}}) \operatorname{K}(\mathbf{M}_{\mathbf{y}}) \operatorname{K}(\mathbf{M}_{\mathbf{y}})$ ,

where a bound for the first term on the right of inequality (4.10) is given by Theorem 3.2, M and M are the mesh ratios, and  $\frac{\pi}{x}$  y  $K(\xi) = 6\xi(2\xi+1)(\xi+1)^2/(3+4\xi)$ .

Proof: Estimate the term

$$|(\mathrm{NB}[\mathbf{f}] - \overline{\mathbf{s}})(\mathbf{x}, \mathbf{y})| \leq \mathcal{E}_{\mathbf{x}} \sum_{i=1}^{M} |\mathbf{A}_{i}(\mathbf{x})| + \mathcal{E}_{\mathbf{y}} \sum_{j=1}^{N} |\mathbf{B}_{j}(\mathbf{y})|$$
$$+ 1/2(\mathcal{E}_{\mathbf{x}} + \mathcal{E}_{\mathbf{y}})(\sum_{i=1}^{M} |\mathbf{A}_{i}(\mathbf{x})|)(\sum_{j=1}^{N} |\mathbf{B}_{j}(\mathbf{y})|)$$
$$\leq \mathcal{E}_{\mathbf{x}} \operatorname{K}(\mathbf{M}_{\pi}) + \mathcal{E}_{\mathbf{y}} \operatorname{K}(\mathbf{M}_{\pi}) + 1/2(\mathcal{E}_{\mathbf{x}} + \mathcal{E}_{\mathbf{y}})$$
$$\cdot \operatorname{K}(\mathbf{M}_{\pi}) \operatorname{K}(\mathbf{M}_{\pi}) .$$

The proof is completed by applying the triangle inequality, and carefully observing the bound for each term in each interval given by Theorem 4.1 and using the fact that  $\sum_{i=0}^{\infty} 2^{i} = 2$ .

<u>Remark 4.</u> 1: It should be noted that the result of Theorem 3.2 is independent of the mesh ratios  $M_{x}$  and  $M_{y}$ , while the exponential decay of the basis functions is dependent on both  $M_{x}$  and  $M_{y}$ . <u>Remark 4.2</u>: To see the size of  $K(M_x)$ , consider the special case of equal spacing, where  $M_x = 1$ , then K(1) = 72/7.

<u>Remark 4.3</u>: If m = n = 4 and  $\mathcal{E} = O(h^8)$ , then the error satisfies Figure 1 from Remark 3.5 of Theorem 3.2.

#### CHAPTER 3

### DISCRETE LEAST SQUARES

This chapter is devoted to the development of discrete least squares algorithms on unstructured data sets, and to showing the accuracy that can be expected depending upon the distribution of the data points and the smoothness of the function f from which the data originates.

Sections 1 and 2 develop preliminary estimates for use in the following section. In Section 3, an example is given to show the necessity of having a sufficient number of data points reasonably distributed to guarantee that the discrete least squares fit will be close the function f. The remaining portion of this section gives error estimates for several univariate discrete least squares algorithms.

Sections 4 and 5 extend the error estimates of Section 3 to bivariate functions, and algorithms are developed using discretized blending function spaces.

Finally, in Section 6, general domains with curved boundarys are considered, and this case is reduced to the methods of Section 5.

121

# Section 1. Uniform Error Estimates

We will develop here an error analysis in the uniform norm for univariate cubic splines in terms of interpolation errors at the knots.

Let  $s \in S^2(\pi_x)$  be a cubic spline on the mesh  $\pi_x:a = x_1 < x_2 < \cdots < x_M = b$ , and set

(1.1) 
$$s_j = s(x_j), s'_j = s^{(1)}(x_j) \text{ and } s''_j = s^{(2)}(s_j) \text{ for } 1 \leq j \leq M.$$

The mesh  $\pi_x$  is uniform if  $x_{j+1} - x_j = h = (b-a)/(M-1)$  for  $l \leq j \leq M-1$ . Also, throughout this chapter, if  $f \in C[a,b]$ , then

$$(1.2) \qquad ||f|| = \sup_{\substack{a \leq x \leq b}} |f(x)|,$$

and if  $g \in C([a, b]x[c, d])$  then

(1.3) 
$$||g|| = \sup |g(x, y)|$$
.  
 $a \leq x \leq b$   
 $c \leq y \leq d$ 

Finally, for  $x_{\epsilon} \mathbb{R}^{N}$  where  $x = (x_{1}, x_{2}, \dots, x_{N})^{T}$  define the vector norm  $||x||_{\infty} = \max_{\substack{k \neq N \\ k \neq N}} |x_{i}|$ . The matrix norm  $||\cdot||_{\infty}$  of a matrix  $A = [a_{ij}]_{\epsilon} \mathbb{R}^{N \times N}$  subordinate to this vector norm is defined to be

$$\||\mathbf{A}\|_{\infty} = \max \sum_{\substack{i \in \mathbb{N} \\ 1 \leq i \leq N \\ j=1}}^{N} |\mathbf{a}_{ij}|, (\text{see} [9, p. 108]).$$

Lemma 1.1. Let  $s \in S^2(\pi_x)$  be a cubic spline on the uniform mesh  $\pi_x$  such that  $|s_j| \leq \xi$  and  $|s_j'| \leq \eta$  for  $1 \leq j \leq M$ , then (1.4)  $||s|| \leq \xi + h\eta/4$ . <u>Proof</u>: From [1, p. 12], if s is a cubic polynomial on  $\begin{bmatrix} x_{j-1}, x_i \end{bmatrix}$  for  $1 < j \le M$  then (1.5)  $s(x) = s'_{j-1} (x_j - x_j)^2 (x - x_{j-1})/h^2$   $- s'_j (x - x_{j-1})^2 (x_j - x)/h^2$   $+ s_{j-1} (x_j - x)^2 [2(x - x_{j-1}) + h]/h^3$  $+ s_j (x - x_{j-1})^2 [2(x_j - x) + h]/h^3$ ,

for  $x \in [x_{j-1}, x_j]$ . Taking absolute values in (1.5) and observing that all of the polynomial factors of  $s'_{j-1}$ ,  $s'_j$ ,  $s_{j-1}$  and  $s_j$  are positive, we have after a little algebra,

(1.6) 
$$|\mathbf{s}(\mathbf{x})| \leq \eta \left[ (\mathbf{x}_{j} - \mathbf{x})^{2} (\mathbf{x} - \mathbf{x}_{j-1}) + (\mathbf{x} - \mathbf{x}_{j-1})^{2} (\mathbf{x}_{j} - \mathbf{x}) \right] / h^{2} + \xi$$

The maximum of the right hand side of (1.6) occurs for  $x = x_{j-1} + h/2$ , and our result follows.

Lemma 1.2. Let  $s \in S(\pi_x)$  be a cubic spline on the uniform mesh  $\pi_x$  such that  $|s_j| \leq \xi$  for  $1 \leq j \leq M$  and  $|s_k'| \leq \eta$  for k = 1, M, then

(1.7) 
$$|s'_{j}| \leq 3\xi/h + \eta/2 \text{ for } 2 \leq j \leq M-1$$
.

<u>Proof</u>: From  $\begin{bmatrix} 1 & p \\ p & 12 \end{bmatrix}$ , if s is a cubic spline then for  $2 \leq j \leq M-1$ 

(1.8a) 
$$s'_{j-1} + 4s'_j + s'_{j+1} = 3(s_{j+1} - s_{j-1})/h$$
,

(1.8b) 
$$4s'_2 + s'_3 = 3(s_3 - s_1)/h - s'_1$$

and

(1.8c) 
$$s'_{M-2} + 4s'_{M-1} = 3(s_M - s_{M-2})/h - s'_M$$

Writing (1.8a) - (1.8c) in matrix notation we obtain



or AX = B, where A, X, and B correspond to the quantities in (1.9).

Premultiply the tridiagonal matrix  $A \in \mathbb{R}^{(M-2)x(M-2)}$  by the matrix D = (1/4)I, where I is the identity matrix. We obtain DA = I + C where



124

Examination of (1.10) yields  $||C||_{\infty} = 1/2$ , and using (3.16) of Chapter 1 we have the existence of  $(DA)^{-1}$  and

(1.11) 
$$||(DA)^{-1}||_{\infty} \leq 1/(1 - ||C||_{\infty})$$
  
 $\leq 2$ .

From (1.11) we can obtain a bound on the derivatives of s at the knots

(1.12) 
$$||X||_{\infty} = ||(DA)^{-1} DB||_{\infty}$$
  
 $\leq ||(DA)^{-1}||_{\infty}||D||_{\infty}||B||_{\infty}$   
 $\leq 1/2 ||B||_{\infty}.$ 

To obtain a bound on  $||B||_{\infty}$  we use (1.9) which yields

(1.13) 
$$||B||_{\infty} = \max \{\max_{\substack{3 \le j \le M-2}} |3(s_{j+1} - s_{j-1})/h|,$$
  
 $|3(s_{3} - s_{1})/h - s_{1}'|, |3(s_{M} - s_{M-2})/h - s_{M}'|\},$ 

where

(1.14) 
$$\max_{\substack{3 \le j \le M-2}} |3(s_{j+1} - s_{j-1})/h| \le 6\xi/h,$$

(1.15) 
$$|3(s_3-s_1)/h - s_1'| \leq 6\xi/h + \eta$$
,

and

(1.16) 
$$|3(s_M^{-s}M_{-2})/h - s_M^{\dagger}| \leq 6\xi/h + \eta$$
.

Therefore  $||B||_{\infty} \leq 6\xi/h + \eta$ , completing our proof of (1.7).

Combining Lemmas 1.1 and 1.2 gives a uniform norm estimate for s in terms of its values at the knots.

<u>Lemma 1.3</u>. Let  $s \in S^2(\pi_x)$  be a cubic spline on the uniform mesh  $\pi_x$  such that  $|s_j| \leq \xi$  for  $1 \leq j \leq M$  and  $|s_k'| \leq \eta$  for k = 1, M, then

(1.17) 
$$||s|| \leq (7/4)\xi + (1/4)h\eta$$
.

<u>Proof:</u> Using Lemma 1.2 and our hypothesis we have for  $1 \leq j \leq M$ 

(1.18) 
$$|s'_j| \leq \max \{\eta, 3\xi/h + \eta/2\}$$
  
 $\leq 3\xi/h + \eta$ .

Combining (1.18) and Lemma 1.1 completes our proof.

Lemma 1.3 is really a stability result for errors in interpolation. To see this, let  $s_1$  be the cubic spline which interpolates the following values of  $f \in C^{(m)}[a,b]$ ,  $l \leq m \leq 4$ 

$$(1.19) \qquad s_{1}(x_{i}) = f(x_{i}), \text{ for } l \leq i \leq M,$$

and

(1.20) 
$$s_1'(x_k) = f'(x_k)$$
, for  $k = 1$ , M.

Also, let s<sub>2</sub> be the cubic spline which interpolates the following functional values of f with errors

(1.21) 
$$s_2(x_i) = f(x_i) + \xi_i$$
, for  $1 \le i \le M$ ,

and

(1.22) 
$$s'_{2}(x_{k}) = f'(x_{k}) + \eta_{k}$$
, for  $k = 1, M$ ,

where  $|\xi_i| \leq \xi$  and  $|\eta_k| \leq \eta$ . Using Theorem 1.1 of Chapter 2 and Lemma 1.3 we have an error bound for the approximate interpolation spline  $s_2$ 

(1.23) 
$$||\mathbf{f} - \mathbf{s}_{2}|| \leq ||\mathbf{f} - \mathbf{s}_{1}|| + ||\mathbf{s}_{1} - \mathbf{s}_{2}||$$
  
 $\leq \xi_{m,0} ||\mathbf{f}^{(m)}|| \mathbf{h}^{m} + (7/4)\xi + (1/4) \mathbf{h} \cdot \eta,$ 

where  $\mathcal{E}_{m,0}$  are given in Table 1 of Chapter 2.

<u>Remark 1.1</u>: Often in performing spline interpolation, an approximate method must be employed to estimate the derivatives at the end points. It is clear from (1.23) that the values  $f'(x_1)$  and  $f'(x_M)$  must be approximated to at least  $O(h^{m-1})$  to preserve the accuracy of the interpolation scheme. Specifically, Lagrange interpolation polynomials can be used to estimate the derivatives. Let  $p_1$  and  $p_M$  be the Lagrange interpolation polynomials of degree m-1 on m points in the intervals  $[x_1, x_2]$  and  $[x_{M-1}, x_M]$  respectively. From [24, p. 289] we have for k = 1, M

(1.24) 
$$|f'(x_k) - p'_k(x_k)| \leq h^{m-1} ||f^{(m)}||/(m-1)!$$

Therefore, if we define  $s_3$  to be that cubic spline which interpolates  $s_3(x_i) = f(x_i) + \xi_i$  for  $1 \le i \le M$ ,  $s'_3(x_1) = p'_1(x_1)$  and  $s'_3(x_M) = p'_M(x_M)$ , then

(1.25) 
$$||\mathbf{f}-\mathbf{s}_{3}|| < (\mathcal{E}_{m,0} + 1/[4(m-1)!])||\mathbf{f}^{(m)}||\mathbf{h}^{m} + (7/4) \xi.$$

Examination of (1.25) indicates that the preservation of the order of the method depends upon limiting the interpolation errors,  $\xi_i$ , to be  $O(h^m)$ . We will use (1.24) later to obtain an error bound for univariate least squares.

#### Section 2. A Uniform Bound

For this section we assume that we are given the set  $X \subseteq [a, b]$ of  $\widehat{M} \ge 3(M-1)$  data points and the uniform mesh  $\pi_x : a = x_1 < x_2 < \cdots < x_M = b$ , where h = (b-a)/(M-1). For each j where  $2 \le j \le M-1$  we assume that there exists six fixed data points which are designated by  $\{\overline{x}_j^i\}_{i=1}^3 \bigcup \{x_j^i\}_{i=1}^3 \subseteq X$  such that

(2.1a) 
$$x_{j-1} \leq \overline{x}_j^3 < \overline{x}_j^2 < \overline{x}_j^1 < x_j < x_j^1 < x_j^2 < x_j^3 \leq x_{j+1}$$

and also the existence of  $\{x_l^i\}_{i=1}^3 \subseteq X$  and  $\{\overline{x}_M^i\} \subseteq X$  such that

(2.1b) 
$$x_1 < x_1^1 < x_1^2 < x_1^3 \leq x_2$$
 and  $x_{M-1} \leq \overline{x}_M^3 < \overline{x}_M^2 < \overline{x}_M^1 < x_M^3$ .

<u>Remark</u>: It need not be the case that  $x_{j-1}^{l} = \overline{x}_{j}^{3}$  or  $x_{j-1}^{2} = \overline{x}_{j}^{2}$  or  $x_{j-1}^{3} = \overline{x}_{j}^{l}$ , although equality in any of the above is acceptable as long as they satisfy the conditions of Lemma 2.1.

Also, for notational convenience, we define  $x_j = \overline{x}_j^0 = x_j^0$ , and powers of  $x_j^i$  will be written as  $(x_j^i)^n$  to avoid confusion.

We introduce the following notation for the k<sup>th</sup> divided difference of the cubic spline  $s \in \vec{S}(\pi_v)$ 

(2.2a) 
$$\Delta_j^k = s \left[ x_j^0, x_j^1, \cdots, x_j^k \right] = \sum_{i=0}^k s(x_j^i) \mu(k, j, i)$$
, and

(2.2b) 
$$\overline{\Delta}_{j}^{k} = s \left[\overline{x}_{j}^{0}, \overline{x}_{j}^{1}, \cdots, \overline{x}_{j}^{k}\right] = \sum_{i=0}^{k} s(\overline{x}_{j}^{i}) \mu(k, j, i),$$

where, for  $0 \leq i \leq k$ 

(2.3a) 
$$\mu(k, j, i) = 1 / \prod_{\substack{\ell=0 \\ \ell \neq i}}^{k} (x_{j}^{i} - x_{j}^{\ell}), \text{ and}$$
  
(2.3b)  $\overline{\mu}(k, j, i) = 1 / \prod_{\substack{\ell=0 \\ \ell \neq i}}^{k} (\overline{x}_{j}^{i} - \overline{x}_{j}^{\ell}),$ 

<u>Lemma 2.1</u>. Let  $s \in S^2(\pi_x)$  be a cubic spline such that  $|s'(x_1)|$ ,  $|s'(x_M)| \leq \eta$  and  $|s(\overline{x}_j^i)|$ ,  $|s(x_j^i)| \leq \xi$  where the data points  $\overline{x}_j^i$ and  $x_j^i$  satisfy (2.1a) and (2.1b). Also, define the real numbers  $\gamma$ ,  $\alpha$  and  $\hat{\alpha}$  by

(2.4) 
$$\begin{cases} \text{ min } \{ -h - \sum_{i=1}^{3} (x_{j} - x_{j}^{i}) \}, \\ \min \{ -h + \sum_{i=1}^{3} (x_{j} - x_{j}^{i}) \} \}, \\ 2 \leq j \leq M \end{cases}$$
(2.5a) 
$$\alpha h = \min \{ \min \min \min_{\substack{1 \leq j \leq M-1 \\ 1 \leq j \leq M-1 \\ 0 \leq i, l \leq 3 \\ i \neq l} } \{ |x_{j}^{i} - x_{j}^{l}| \}, \\ 1 \leq j \leq M-1 \\ 0 \leq i, l \leq 3 \\ i \neq l \end{cases}$$
(2.5b) 
$$\hat{\alpha} h^{3} = \min \{ \min \min_{\substack{1 \leq i \leq 3 \\ 1 \leq i \leq 3 \\ 0 \leq i, l \leq 3 \\ i \neq l \\ 1 \leq i \leq 3 \\ 0 \leq i, l \leq 3 \\ 0 \leq i, l \leq 3 \\ i \neq l \\ 1 \leq i \leq 3 \\ \end{cases}$$

If  $\gamma > 0$ , then

(2.6) 
$$||\mathbf{s}|| \leq 7(1-\alpha)(1-2\alpha) \left[ (3-6\alpha) \xi / \hat{\alpha} + h\eta \right] / (4\delta) + h\eta/4.$$

<u>Remark</u>: From (2.5a) and (2.5b) it follows that  $1/3 \ge \hat{\alpha} > 0$  and  $1 > \alpha > 0$ .

<u>Proof</u> From  $\begin{bmatrix} 1 & p & 10 \end{bmatrix}$ ,  $s \in S(\pi_x)$  is a cubic spline if and only if for  $2 \leq j \leq M-1$ 

$$(2.6a) \qquad (h^2/6)(s_{j-1}'' + 4s_j'' + s_{j+1}'') = s_{j-1} - 2s_j + s_{j+1}$$

(2.6b) 
$$(h^2/6)(2s_1'' + s_2'') = s_2 - s_1 - hs_1',$$

and

(2.6c) 
$$(h^2/6)(s_{M-1}'' + 2s_M'') = hs_M' - s_M + s_{M-1}$$
.

Because  $s \in C^{(2)}[a, b]$  is a cubic polynomial on  $[x_j, x_{j+1}]$  for  $1 \leq j \leq M-1$  we have  $s_{j+1}^{"'} = s_j^{"'} + s_j^{"''}(x_j)$ , where  $s_j^{"''}(x_j)$  is evaluated from the right. Also,  $s_j^{"''}(x) = \text{constant} = 6 \Delta_j^3$  on  $[x_j, x_{j+1}]$ , (see [24, p. 249]. Therefore,

(2.7) 
$$s''_{j+1} = s''_j + 6 h \Delta_j^3$$
, for  $l \leq j \leq M-1$ ,

and correspondingly

(2.8) 
$$s_{j-1}^{\prime\prime} = s_j^{\prime\prime} - 6h\overline{\Delta}_j^3$$
, for  $2 \leq j \leq M$ .

Because s is cubic on  $\begin{bmatrix} x_j, & x_{j+1} \end{bmatrix}$  for  $1 \le j \le M-1$ , we represent s on that interval as

(2.9) 
$$\mathbf{s}_{j} = \Delta_{j}^{0} + (\mathbf{x} - \mathbf{x}_{j}^{0}) \Delta_{j}^{1} + (\mathbf{x} - \mathbf{x}_{j}^{0})(\mathbf{x} - \mathbf{x}_{j}^{1}) \Delta_{j}^{2}$$

+ 
$$(x-x_{j}^{0})(x-x_{j}^{1})(x-x_{j}^{2}) \Delta_{j}^{3}$$
,

(see  $\begin{bmatrix} 24, p. 248 \end{bmatrix}$ ). Differentiating (2.9) twice and evaluating at  $x_j$  we obtain

(2.10) 
$$\mathbf{s}_{j}^{\prime\prime} = 2 \Delta_{j}^{2} + 2(2 \mathbf{x}_{j} - \mathbf{x}_{j}^{1} - \mathbf{s}_{j}^{2}) \Delta_{j}^{3}$$
.

Correspondingly, on  $[x_{j-1}, x_j]$  for  $2 \leq j \leq M$  we have

(2.11) 
$$\mathbf{s}_{j}^{"} = 2 \,\overline{\Delta}_{j}^{2} + 2(2\mathbf{x}_{j} - \overline{\mathbf{x}}_{j}^{1} - \overline{\mathbf{x}}_{j}^{2}) \,\overline{\Delta}_{j}^{3}.$$

Substituting (2.7), (2.8), (2.10) and (2.11) into (2.6a) through (2.6c) it follows that for  $2 \le j \le M-1$ 

(2.12a) 
$$h^{2} [(2x_{j} - \overline{x}_{j}^{1} - \overline{x}_{j}^{2} - h) \overline{\Delta}_{j}^{3} + \overline{\Delta}_{j}^{2} + \Delta_{j}^{2} + (2x_{j} - x_{j}^{1} - x_{j}^{2} + h) \Delta_{j}^{3}] = s_{j-1} - 2s_{j} + s_{j+1},$$
  
(2.12b)  $h^{2} [(2x_{1} - x_{1}^{1} - x_{1}^{2} + h) \Delta_{1}^{3} + \Delta_{1}^{2}] = s_{2} - s_{1} - hs',$ 

and

(2.12c) 
$$h^2 \left[ (2\mathbf{x}_M - \overline{\mathbf{x}}_M^1 - \overline{\mathbf{x}}_M^2 - \mathbf{h}) \overline{\Delta}_M^3 + \overline{\Delta}_M^2 \right] = h \mathbf{s}'_M - \mathbf{s}_M + \mathbf{s}_{M-1}$$

Substituting (2.2a) and (2.2b) for our divided differences into (2.12a) through (2.12c), we obtain the following after regrouping

(2.13a) 
$$h^{2} \Big[ (2x_{j} - \overline{x}_{j}^{1} - \overline{x}_{j}^{2} - h) \overline{\mu} (3, j, 0) + \overline{\mu} (2, j, 0) + \mu (2, j, 0) + (2x_{j} - x_{j}^{1} - x_{j}^{2} + h) \mu (3, j, 0) \Big] s_{j} + 2s_{j} - s_{j-1} - s_{j+1}$$

$$= -h^{2} \Big[ (2x_{j} - \overline{x}_{j}^{1} - \overline{x}_{j}^{2} - h) \frac{3}{\sum_{i=1}^{\infty}} s(\overline{x}_{j}^{i}) \overline{\mu} (3, j, i) \\ + \frac{2}{\sum_{i=1}^{\infty}} s(\overline{x}_{j}^{i}) \overline{\mu} (2, j, i) + \frac{2}{\sum_{i=1}^{\infty}} s(x_{j}^{i}) \mu (2, j, i) \\ + (2x_{j} - x_{j}^{1} - x_{j}^{2} + h) \frac{3}{\sum_{i=1}^{\infty}} s(x_{j}^{i}) \mu (3, j, i) \Big] ,$$

$$(2.13b) h^{2} \Big[ (2x_{1} - x_{1}^{1} - x_{1}^{2} + h) \mu (3, 1, 0) + \mu (2, 1, 0) \Big] s_{1} + s_{1} - s_{2} \\ = -h^{2} \Big[ (2x_{1} - x_{1}^{1} - x_{1}^{2} + h) \frac{3}{\sum_{i=1}^{\infty}} s(x_{1}^{i}) \mu (3, 1, i) \\ + \frac{2}{\sum_{i=1}^{\infty}} s(x_{1}^{i}) \mu (2, 1, i) \Big] - h s_{1}^{i} ,$$

and

(2.13c) 
$$h^2 \left[ (2x_M - \overline{x}_M^1 - \overline{x}_M^2 - h) \overline{\mu} (3, M, 0) + \overline{\mu} (2, M, 0) \right] s_M$$

$$+ {}^{s}{}_{M} - {}^{s}{}_{M-1}$$

$$= -h^{2} \left[ (2s_{M} - \overline{x}_{M}^{1} - \overline{x}_{M}^{2} - h) \sum_{i=1}^{3} s(\overline{x}_{M}^{i}) \overline{\mu} (3, M, 0) + \sum_{i=1}^{2} s(\overline{x}_{M}^{i}) \overline{\mu} (2, M, 0) \right] + h s'_{M}.$$

Using (2.3a) and (2.3b), we simplify (2.13a) through (2.13b) to obtain

(2.14a) 
$$-\mathbf{s}_{j-1} + (\overline{\xi}_j + 2 + \xi_j) \mathbf{s}_j - \mathbf{s}_{j+1} = \overline{\delta}_j + \delta_j$$
,

(2.14b) 
$$(1 + \xi_1) \mathbf{s}_1 - \mathbf{s}_2 = \delta_1 - \mathbf{hs}_1'$$
,

 $\mathtt{and}$ 

(2.14c) 
$$(\overline{\mathcal{E}}_{M} + 1) \mathbf{s}_{M} - \mathbf{s}_{M-1} = \overline{\delta}_{M} + \mathbf{hs}'_{M},$$
where for  $1 \leq j \leq M-1$ 

(2.15a) 
$$\mathcal{E}_{j} = h^{2} \left[ h + \sum_{i=1}^{3} (x_{j} - x_{j}^{i}) \right] \mu(3, j, 0) > 0$$
,

(2.15b) 
$$\delta_{j} = -h^{2} \sum_{i=1}^{3} \left[ h + \sum_{\ell=1}^{3} (x_{j} - x_{j}^{\ell}) \right] s(x_{j}^{i}) \mu (3, j, i) ,$$

and for  $2 \leqslant j \leqslant M$ 

(2.15c) 
$$\overline{\mathcal{E}}_{j} = h^{2} \left[ -h + \sum_{i=1}^{3} (x_{j} - \overline{x}_{j}^{i}) \right] \overline{\mu} (3, j, 0) > 0,$$

(2.15d) 
$$\overline{\delta_j} = -h^2 \sum_{i=1}^3 \left[ -h + \sum_{\substack{\ell=1 \\ \ell \neq i}}^3 (\mathbf{x}_j - \mathbf{\bar{x}}_j^\ell) \right] \mathbf{s}(\mathbf{\bar{x}}_j^i) \mathbf{\bar{\mu}} (3, j, i) .$$

In (2.15a) and (2.15c) we have used (2.4) and the fact that 3 > 0,  $\mu(3, j, 0) < 0$  and  $\overline{\mu}(3, j, 0) > 0$  to show that they are positive. In matrix form, equations (2.14a) through (2.14c) become



$$= \begin{bmatrix} \overline{\delta}_{1} - h s'_{1} \\ \overline{\delta}_{2} + \delta_{2} \\ \vdots \\ \vdots \\ \overline{\delta}_{M-1} + \delta_{M-1} \\ \overline{\delta}_{M} + h s'_{M} \end{bmatrix},$$

or AS = E, where A, S and E correspond to the quantities in (2.16). In order to estimate the norm of the vector S we define the diagonal matrix D



From (2.15a) and (2.15c) we see that all of the diagonal entries of the matrix D are positive, which implies that  $D^{-1}$  exists. Using this fact, and premultiplying A by  $D^{-1}$  we obtain

$$(2.18) D^{-1} A = I + B ,$$

where the tridiagonal matrix B is given by



If  $||B||_{\infty} < 1$ , then, from (3.16) of Chapter 1, we have that  $(D^{-1}A)^{-1}$  exists and

(2.20) 
$$||S||_{\infty} \leq ||(D^{-1}A)^{-1}||_{\infty}||D^{-1}E||_{\infty}$$
  
 $\leq [1/(1-||B||_{\infty})]||D^{-1}E||_{\infty}.$ 

We now proceed to calculate bounds for  $||B||_{\infty}$  and  $||D^{-1}E||_{\infty}$ .

Examination of (2.19), (2.15a) and (2.15c) shows that

(2.21) 
$$||B||_{\infty} = \max \{ \max \{ 2/(\overline{\xi}_{j}+2+\xi_{j}), 1/(1+\xi_{1}), 2 \le j \le M-1 \\ 1/(\overline{\xi}_{M}+1) \},$$

which implies that  $\overline{\mathcal{E}}_{j}$  and  $\mathcal{E}_{j}$  should be bounded from below to bound  $||B||_{\infty}$  from above.

Using (2.15a), (2.4) and (2.5a) it is clear that for  $1 \leq j \leq M-1$ 

$$(2.22) \qquad \begin{split} \mathcal{E}_{j} & \geqslant h^{3} \, \mathcal{V} \left| \mu(3, j, 0) \right| \\ & \geqslant h^{3} \, \mathcal{V} \left| (\mathbf{x}_{j} - \mathbf{x}_{j}^{1})(\mathbf{x}_{j} - \mathbf{x}_{j}^{2})(\mathbf{x}_{j} - \mathbf{x}_{j}^{3}) \right| \\ & \geqslant h^{3} \, \mathcal{V} \left[ h \cdot (1 - 2\alpha) \cdot h(1 - \alpha) \cdot h \right] \\ & \geqslant \mathcal{V} \left[ (1 - \alpha)(1 - 2\alpha) \right] \,. \end{split}$$

Correspondingly, from (2.15c), (2.4) and (2.5a) we have for

(2.23) 
$$\widetilde{\mathcal{E}}_{i} \geq \delta / [(1-\alpha)(1-2\alpha)].$$

 $2 \leq j \leq M$  that

Combining (2.21), (2.22) and (2.23) yields

(2.24) 
$$||B||_{\infty} \leq (1-\alpha)(1-2\alpha) / [(1-\alpha)(1-2\alpha) + \delta]$$
  
< 1.

In order to obtain a bound on  $||D^{-1}E||_{\infty}$  we use (2.16) and (2.17) which gives

$$(2.25) \quad D^{-1}E = \begin{bmatrix} \left(\delta_{1} - h s_{1}^{\prime}\right) / \left(1 + \mathcal{E}_{1}\right) \\ \left(\overline{\delta}_{2} + \delta_{2}\right) / \left(\overline{\mathcal{E}}_{2} + 2 + \mathcal{E}_{2}\right) \\ \vdots \\ \left(\overline{\delta}_{M-1} + \delta_{M-1}\right) / \left(\overline{\mathcal{E}}_{M-1} + 2 + \mathcal{E}_{M-1}\right) \\ \left(\overline{\delta}_{M} + h s_{M}^{\prime}\right) / \left(\overline{\mathcal{E}}_{M} + 1\right) \end{bmatrix}$$

Which implies using (2.22) and (2.23) that

(2.26) 
$$||D^{-1}E||_{\infty} \leq \left[(1-\alpha)(1-2\alpha)/((1-\alpha)(1-2\alpha)+\delta)\right] ||E||_{\infty}$$

where

$$(2.27) \stackrel{\wedge}{\mathbf{E}} = \begin{bmatrix} \delta_1 - \mathbf{h} \mathbf{s}'_1 \\ (\overline{\delta}_2 + \delta_2)/2 \\ \vdots \\ \vdots \\ (\overline{\delta}_{M-1} + \delta_{M-1})/2 \\ \overline{\delta}_M + \mathbf{h} \mathbf{s}'_M \end{bmatrix}.$$

An upper bound on  $||\stackrel{\wedge}{E}||_{\infty}$  is obtained by bounding  $\overline{\delta}_{j}$  and  $\delta_{j}$  from above. Using (2.5a) and carefully observing internal cancellations we have for  $1 \leq i \leq 3$ 

(2.28) 
$$|\mathbf{h} + \sum_{\substack{\ell=1\\ \ell\neq i}}^{3} (\mathbf{x}_{j} - \mathbf{x}_{j}^{\ell})| \leq \mathbf{h} [1 - \mathbf{i} \cdot \alpha]$$

and

(2.29) 
$$|-h + \sum_{\substack{\ell=1 \\ \ell \neq i}}^{3} (x_j - \overline{x}_j^{\ell})| \leq h [1 - i \cdot \alpha].$$

Combining (2.15b), (2.15d), (2.5b), (2.28) and (2.29) it follows that for  $1 \leq j \leq M-1$ 

(2.30) 
$$|\delta_{j}| \leq (3 - 6\alpha) \xi / \hat{\alpha},$$

and for  $2 \leq j \leq M$ 

(2.31) 
$$|\overline{\delta}_{j}| \leq (3 - 6\alpha) \xi / \hat{\alpha}$$

Therefore, from (2.27), (2.30) and (2.31) follows an upper bound for  $||\hat{E}||_{\infty}$ 

(2.32) 
$$|| \stackrel{\wedge}{\mathbf{E}} ||_{\infty} = \max \{ \max_{\substack{2 \leq j \leq M-1}} \{ |\overline{\delta}_{j} + \delta_{j} |/2 \}, \\ 2 \leq j \leq M-1$$
$$|\delta_{1} - h_{1} s_{1}' |, |\overline{\delta}_{M} + h s_{M}' | \}$$

$$\leq (3 - 6\alpha) \xi / \hat{\alpha} + h \eta$$
.

Combining (2.20), (2.24), (2.26) and (2.32) yields

(2.33) 
$$||S||_{\infty} \leq (1-\alpha)(1-2\alpha) \left[ (3-6\alpha)\xi / \alpha + h\eta \right] / \chi$$
.

Our result follows from (2.33) and Lemma 1.3.

Section 3. Univariate Discrete Least Squares

Assume that we are given an unstructured set  $X = \{\hat{x}_i\}_{i=1}^{\hat{M}} \subseteq [a, b]$  of  $\hat{M}$  data points, and the mesh  $\pi_x: a = x_1 < x_2 < \cdots < x_M = b$ , where  $h_x = \max_{\substack{1 \leq j \leq M-1}} (x_{j+1} - x_j)$ . If  $f \in C^{(m)}[a, b]$ , where  $l \leq m \leq 4$ ,

then a cubic spline  $s_{LS} \in S^2(\pi)$  which minimizes the Euclidean norm of the residual vector  $R \in \mathbb{R}^{\widehat{M}}$ , where component i of R is

(3.1) 
$$R_i = f(\hat{x}_i) - s_{LS}(\hat{x}_i)$$
,

i. e.

(3.2) 
$$||\mathbf{R}|| = (\sum_{i=1}^{\hat{M}} (f(\hat{\mathbf{x}}_{i}) - s_{LS}(\hat{\mathbf{x}}_{i}))^{2})^{1/2}$$
  
 $= \min_{\mathbf{s} \in S^{2}(\pi_{x})} (\sum_{i=1}^{\hat{M}} (f(\hat{\mathbf{x}}_{i}) - s(\hat{\mathbf{x}}_{i}))^{2})^{1/2}$ 

is a discrete least squares approximation to f on the unstructured data set X. It would be desirable to have an estimate as to how close  $s_{LS}$  is to f in the uniform norm similar to the estimate we have for the cubic interpolation spline  $s_f \in S^2(\pi_x)$ , (Theorem 1.1 of Chapter 2). We would hope that if f is smooth and the residual vector R is small, then  $s_{LS}$  is close to f. This, however, is not the case, even if the cubic interpolation spline  $s_f$  is close to f. Therefore, as a prelude to the next theorem, we will give the following example which illustrates the importance of having sufficient data points which are reasonably distributed on [a, b] with respect to the mesh  $\pi_y$ . Example 3.1: We will construct the discrete least squares approximation  $s_{LS} \in S^2(\pi_x)$  on the uniform mesh  $\pi_x : a = x_1 < x_2 < \cdots < x_M = M$ , where  $x_i = i$  for  $1 \le i \le M$ ,  $h_x = 1$  and M > 3 to the function  $f \in C^{(4)}[1, M]$ 

(3.3) 
$$f(x) = \begin{cases} \mathcal{E}(4(x-5/4))^5 & \text{if } 1 \leq x \leq 5/4 \\ \\ 0 & \text{if } 5/4 < x \leq M \end{cases}$$

Direct calculation using (3.3) yields  $||f|| = \mathcal{E}$  and  $||f^{(4)}|| =$ 30720  $\mathcal{E}$ . If  $s_f \in S^2(\pi_x)$  is the cubic interpolation spline which interpolates f on the mesh  $\pi_x$ , then from Theorem 1.1 of Chapter 2 (3.4)  $||f-s_f|| \leq 400 \mathcal{E}$ .

We will explicitly construct the discrete least squares spline approximation  $s_{LS}$  to f on the following M+2 data points

(3.5) 
$$X = \{1 + i/4\}_{i=0}^{3} \cup \{j + \mathcal{E}\}_{j=2}^{M-1} \subseteq [1, M]$$

where  $f(1) = -\mathcal{E}$ , f(1 + i/4) = 0 for  $1 \le i \le 3$ ,  $f(j + \mathcal{E}) = 0$  for  $2 \le j \le M-1$ , and we assume  $0 < \mathcal{E} \le 1$  2. Our construction will show that the residual vector  $\mathbb{R} \in \mathbb{R}^{M+2}$  is zero, and that  $s_{LS}$  is the unique discrete least squares solution to f on the data set X.

If R is zero, then  $s_{LS}$  interpolates f at the points {1 + i/4} $_{i=0}^{3}$ , which implies that on [1,2],  $s_{LS}$  is uniquely given by (3.6)  $s_{LS}^{(x)} = \mathcal{E} \cdot 32 \cdot (x - 5/4)(x - 3/2)(x - 7/4)/3$ ,

where

(3.7a) 
$$s_{1,S}(2) = \mathcal{E} > 0$$
,

(3.7b)  $s'_{LS}(2) > 0$ ,

and

(3.7c) 
$$s''_{LS}(2) > 0$$

Because  $s_{LS} \in C^{(2)}[a, b]$  and  $s_{LS}(j+\xi) = f(j+\xi) = 0$ , it follows that on the interval [j, j+1],  $s_{LS}$  has the unique representation

(3.8) 
$$s_{LS}(x) = -(s_{LS}'(j)/(2\xi) + s_{LS}'(j)/\xi^{2} + s_{LS}'(j)/\xi^{3})(x-j)^{3} + (s_{LS}'(j)/2)(x-j)^{2} + s_{LS}'(j)(x-j) + s_{LS}'(j) .$$

If  $s_{LS}$  has been uniquely constructed on the interval [1, j], then (3.8) gives a unique extension to [1, j+1], hence inductively to [1, M]. From (3.8), if  $s_{LS}(j)$ ,  $s'_{LS}(j)$  and  $s''_{LS}(j)$  are all nonzero and of the same sign, then  $s_{LS}(j+1)$ ,  $s'_{LS}(j+1)$  and  $s''_{LS}(j+1)$ are all non-zero with the opposite sign and

(3.9) 
$$|\mathbf{s}_{LS}(j+1)| > |\mathbf{s}_{LS}(j)|/\xi^2$$
,

where we have used our assumption that  $0 \le \le 1/2$  for the conservative lower bound in (3.9). Combining (3.7a), (3.7b), (3.7c) and (3.9) we conclude that for  $3 \le j \le M$ 

- (3.10)  $|s_{LS}(j)| > (1/\varepsilon)^{2j-5}$ ,
- $(3.11) \qquad ||f-s_{\rm LS}|| > (1/\mathcal{E})^{2M-5} ,$

and

(3.12) 
$$||s_{f} - s_{LS}|| > (1/\xi)^{2M-5}$$

where we have used the fact that  $f(M) = s_f(M) = 0$ . Therefore, by decreasing  $\mathcal{E}$ , we can cause f and its first four derivatives to be as small in norm as desired. However, the norm of the error in (3.11) can be made as large as desired, even though the residual vector R is zero.

From this example, it is clear that an acceptable upper bound on  $||f-s_{LS}||$  can be obtained only if we place restrictions on the number and distribution of the data points with respect to the mesh  $\pi_x$ . Finally, observe that even if  $s_f$  is close to f,  $s_{LS}$  need not be close to either  $s_f$  or f.

<u>Remark 3.1</u>: The observation should be made that the mesh size h = 1 was chosen only for convenience, and its size is not crucial to the above example.

The above example motivates the hypothesis of the following theorem.

<u>Theorem 3.1.</u> Let  $s_{LS} \in S^2(\pi_x)$  be a discrete least squares approximation to  $f \in C^{(m)}[a, b]$ ,  $1 \le m \le 4$ , on the set X of  $\hat{M}$  unstructured data points, where  $\pi_x$  is a uniform mesh. Assume for each j, where  $2 \le j \le M-1$ , that there exists six data points which we designate by  $\{\bar{x}_j^i\}_{i=1}^3 \cup \{x_j^i\}_{i=1}^3 \subseteq X$  satisfying (2.1a), and also

 $\{x_{j}^{i}\}_{i=1}^{3} \subseteq X \text{ and } \{\overline{x}_{M}^{i}\}_{i=1}^{3} \subseteq X \text{ satisfying (2.1b). Let the real numbers } \lambda, \alpha \text{ and } \hat{\alpha} \text{ be defined as in (2.4), (2.5a) and (2.5b), }$  respectively. If  $\lambda > 0$ ,  $|f'(x_{k}) - s'_{LS}(s_{k})| \leq \eta$ , for k = 1, M, and at all of the data points  $x_{j}^{i}$  and  $\overline{x}_{j}^{i}$  defined above  $|f(x_{j}^{i}) - s_{LS}(x_{j}^{i})| \leq \xi$  and  $|f(\overline{x}_{j}^{i}) - s_{LS}(\overline{x}_{j}^{i})| \leq \xi$ , then

(3.13) 
$$||f-s|_{LX}|| \leq \mathcal{E}_{m,0}(21(1-\alpha)(1-2\alpha)^2/(4\alpha)) + 1)||f^{(m)}||h_x^m$$

+ 
$$21(1-\alpha)(1-2\alpha)^2 \xi/(4\alpha \chi) + (7(1-\alpha)(1-2\alpha)/(4\chi) + 1/4) h_x \eta$$
,

where  $\mathcal{E}_{m,0}$  is given in Table 1 of Chapter 2 and  $h_x = (b-a)/(M-1)$ .

<u>Proof</u>: From Theorem 1.1 of Chapter 2, if  $s_f \in S^2(\pi_x)$  is the cubic interpolation spline to f, then

$$(3.14) \qquad ||f - s_{LS}|| \leq ||f - s_{f}|| + ||s_{f} - s_{LS}|| \\ \leq \mathcal{E}_{m,0} ||f^{(m)}|| h_{x}^{m} + ||s_{f} - s_{LS}||.$$

We now obtain a bound for  $||s_f - s_{LS}||$  from Lemma 2.1. First observe that  $s_f - s_{LS} \in S^2(\pi_x)$  and

$$(3.15) \qquad \left| \begin{array}{l} \mathbf{s}_{f}'(\mathbf{x}_{k}) - \mathbf{s}_{LS}'(\mathbf{x}_{k}) \right| = \left| f'(\mathbf{x}_{k}) - \mathbf{s}_{LS}'(\mathbf{x}_{k}) \right| \\ \leqslant \eta ,$$

for k = 1, M. Furthermore, for our specified data points we have

$$(3.16) \qquad |\mathbf{s}_{f}(\mathbf{x}_{j}^{i}) - \mathbf{s}_{LS}(\mathbf{x}_{j}^{i})| \leq |\mathbf{s}_{f}(\mathbf{x}_{j}^{i}) - \mathbf{f}(\mathbf{x}_{j}^{i})| + |\mathbf{f}(\mathbf{x}_{j}^{i}) - \mathbf{s}_{LS}(\mathbf{x}_{j}^{i})| \\ \leq ||\mathbf{s}_{f} - \mathbf{f}|| + \xi \\ \leq \mathcal{E}_{m,0} ||\mathbf{f}^{(m)}||\mathbf{h}_{\mathbf{x}}^{m} + \xi ,$$

and correspondingly

$$(3.17) \qquad |\mathbf{s}_{f}(\overline{\mathbf{x}}_{j}^{i}) - \mathbf{s}_{LS}(\overline{\mathbf{x}}_{j}^{i})| \leq \mathcal{E}_{m,0} ||\mathbf{f}^{(m)}||\mathbf{h}_{\mathbf{x}}^{m} + \boldsymbol{\xi} .$$

Therefore, from Lemma 2.1

(3.18) 
$$||s_{f} - s_{LS}|| \leq 7(1-\alpha)(1-2\alpha)(3-6\alpha)(\mathcal{E}_{m,0}||f^{(m)}||h_{x}^{m} + \xi)/4\hat{\alpha}\gamma + (7(1-\alpha)(1-2\alpha)/(4\delta) + 1/4)h_{x}\eta$$
.

Combining (3.14) and (3.18) yields (3.13) and the proof is complete.

۲

<u>Remark 3.2</u>: In Theorem 3.1 it is permissible to use either of the following estimates for  $\xi$ ,

(3.19)  $\xi = ||\mathbf{R}||_{\infty}$ ,

or

 $(3.20) \quad \xi = ||R||,$ 

where  $||R||_{\infty} \leq ||R||$ . In practice, however, ||R|| is usually a poor choice when compared to  $||R||_{\infty}$ , especially for large  $\stackrel{\Lambda}{M}$ .

<u>Remark 3.3</u>: It should be noted that during the numerical solution of a discrete least squares problem, the residual vector R is calculated, (see [4, 11, 21, 28]). Therefore, R is available for use in the estimates (3.19) and (3.20).

<u>Remark 3.4</u>: Examination of (2.3a) through (2.5b) yields the following lower bound for  $\checkmark$  and  $\hat{\alpha}$ 

 $(3.21) \qquad \forall \ge 6\alpha - 1$ 

and

 $(3.22) \qquad \stackrel{\wedge}{\alpha} \geqslant 2\alpha^3 \quad .$ 

However,  $6\alpha - 1$  and  $2\alpha^3$  are usually much smaller than  $\forall$  and  $\hat{\alpha}$ , respectively, and their use in Theorem 3.1 deteriorates the bound (3.13).

Example: If we assume that the data points X are distributed in such a way as to give  $\alpha \ge 1/4$ , then from (3.21) and (3.22)  $Y \ge 1/2$  and  $\hat{\alpha} \ge 1/32$ . This gives the following bound from (3.13)

(3.23) 
$$||f-s_{LS}|| \leq \mathcal{E}_{m,0}^{127} ||f^{(m)}||h_x^m + 126 \xi + (85/4) h_x^{3} \eta$$

Examination of Theorem 3.1 shows the need of insuring the smallness of  $\eta$ . We want to approximate  $f'(\mathbf{x}_k)$  for k = 1, M in some manner which is compatible with the  $h^m$ . One possibility is as follows. If each of the intervals  $[\mathbf{x}_1, \mathbf{x}_2]$  and  $[\mathbf{x}_{M-1}, \mathbf{x}_M]$  contains at least m data points, respectively (the hypotheses of Theorem 3.1 guarantees at least 3), then using Lagrange interpolation polynomials of degree m-1, we can approximate  $f'(\mathbf{x}_k)$  for k = 1, M with the derivative of these polynomials at the end points of our interval  $x_1$  and  $x_M$  (see Remark 1.1). If a constrained least squares algorithm is used to insure that the least squares spline  $s_{LS}$  has these values for its end derivatives, then inequality (1.24) gives a bound for M. Therefore, the following Corollary is an immediate consequence of Theorem 3.1 and the above construction.

<u>Corollary 3.1</u>. If the hypotheses of Theorem 3.1 are satisfied and there exists m data points in each of the intervals  $[x_1, x_2]$  and  $[x_{M-1}, x_M]$  on which the Lagrange interpolation polynomials of degree m-1 are constructed to approximate  $f'(x_k)$ , k = 1, M, and  $s'_{LS}(x_k)$  is equal to those approximations, then

(3.24) 
$$||f - s_{LS}|| \leq \left[ \frac{\mathcal{E}_{m,0}(21(1-\alpha)(1-2\alpha)^2/(4\alpha^2) + 1) + (7(1-\alpha)(1-2\alpha)/4\gamma + 1/4)/(m-1)! \right] ||f^{(m)}||h_x^m + 21(1-\alpha)(1-2\alpha)^2 \xi / (4\alpha^2) .$$

It would also be desirable to have an a priori bound on  $||f-s_{LS}||$ , knowing only that  $f \in C^{(m)}[a,b]$  and the distribution of the data points X. Toward this end we prove the following Corollary.

<u>Corollary 3.2</u>. If the hypotheses of both Theorem 3.1 and Corollary 3.1 are satisfied, then

$$(3.25) \qquad ||\mathbf{f} - \mathbf{s}_{\mathrm{LS}}|| \leq \left[ \mathcal{E}_{\mathbf{m}, 0}^{(21(1-\alpha)(1-2\alpha)^{2}(1+\widehat{M}^{1/2})/(4\widehat{\alpha}\chi) + 1)} + (7(1-\alpha)(1-2\alpha)/(4\chi) + 1/4)/(\mathbf{m}-1)! \right] ||\mathbf{f}^{(\mathbf{m})}|| \mathbf{h}_{\mathbf{x}}^{\mathbf{m}} ,$$

where  $\mathcal{E}_{m,0}$  is given in Table 1 of Chapter 2.

**Proof:** Examination of (3.24) shows that all that remains is to bound  $\xi$ . This will be accomplished by obtaining a bound for  $||\mathbf{R}||$ and using (3.20). The cubic interpolation spline  $\mathbf{s}_{f}$  is a candidate for the discrete least squares approximation to f, hence the norm of its residual vector must be greater than or equal to  $||\mathbf{R}||$ , i.e.,

$$(3.26) \qquad ||\mathbf{R}|| \leq \left(\sum_{i=1}^{M} (f(\hat{\mathbf{x}}_{i}) - \mathbf{s}_{f}(\hat{\mathbf{x}}_{i}))^{2}\right)^{1/2}$$
$$\leq \left(\sum_{i=1}^{\hat{M}} ||f - \mathbf{s}_{f}||^{2}\right)^{1/2}$$
$$\leq \mathcal{E}_{m,0} ||f^{(m)}|| h_{x}^{m} \hat{M}^{1/2} ,$$

where an application of Theorem 1.1 of Chapter 2 has been made. Combining (3.24), (3.20) and (3.26) completes the proof.

Because  $\hat{M}^{1/2} > (3(M-1))^{1/2} > h_x^{-1/2}$ , the power of  $h_x$  is no longer m, as (3.26) might indicate, but rather no larger than m-1/2, depending upon the number of data points. In practice, of course, we would expect (3.26) to be a rather crude estimate and would hope that  $\xi$  is much closer to  $\mathcal{E}_{m,0} ||f^{(m)}|| h_x^m$  instead of the bound given in (3.26).

We now give a more stringent condition on the data points, X, which will insure that the discrete least squares spline  $s_{LS}$  will be uniformly close to f. In preparation for the following theorem, the real non-negative number  $||\cdot||_X$  of a function f is defined to be

$$(3.27) \qquad ||f||_{X} = \max_{\substack{1 \leq i \leq \hat{M} \\ \hat{x}_{i} \in X}} |f(\hat{x}_{i})|,$$

which is used in the following Lemma (see  $\begin{bmatrix} 6 \end{bmatrix}$ ).

<u>Lemma 3.1</u>.  $\begin{bmatrix} 6 & p. & 91 \end{bmatrix}$ . Let P be an algebraic polynomial of degree  $\leq n$  on the interval [a, b], then (3.28)  $||P||(1 - n^2 \beta_1(X)) \leq ||P||_X$ ,

where

(3.29) 
$$\beta_1(X) = \max \min_{\substack{x \in [a, b]}} \frac{2|x - \hat{x}_i|}{b - a}$$

 $\mathtt{and}$ 

(3.30) 
$$||P||(1 - (n\hat{\beta}_{1}(X))^{2}/2) \leq ||P||_{X}$$

where

(3.31) 
$$\hat{\beta}_{1}(X) = \max_{x \in [a, b]} \min_{1 \le i \le M} |\cos^{-1}((2x - (a+b))/(b-a))|$$
  
 $-\cos^{-1}((2x - (a+b))/(b-a))|$ .

The above Lemma yields the following theorem.

<u>Theorem 3.2.</u> Let  $X = \{\hat{x}_i\}_{i=1}^{\hat{M}} \subseteq [a, b]$  be a set of  $\hat{M}$  unstructured data points and  $\pi_x: a = x_1 < x_2 < \cdots < x_M = b$  be a mesh on [a, b]. Let  $s_{LS} \in S^2(\pi_x)$  be a discrete least squares approximation to  $f \in C^{(m)}[a, b], 1 \leq m \leq 4$ , on the data set X with residual vector  $R \in \mathbb{R}^{\hat{M}}$ . Define  $\beta(X)$  and  $\hat{\beta}(X)$  by

$$(3.32) \qquad \beta(X) = \max \max \min_{\substack{1 \leq j \leq M-1 \\ i \leq j \leq M-1 \\ i \in [x_j, x_{j+1}]}} \max \prod_{\substack{1 \leq i \leq \widehat{M} \\ \widehat{x}_i \in [x_j, x_{j+1}]}}$$

$$\{2|\mathbf{x}-\hat{\mathbf{x}}_{i}|/(\mathbf{x}_{j+1}-\mathbf{x}_{j})\}$$
 ,

$$(3.33) \quad \widehat{\beta}(X) = \max \max \max \{ \underset{i \leq j \leq M-1}{\max} \max \{ \underset{j, j+1}{\max} \}$$

$$\{ |\cos^{-1}((2x - (x_j + x_{j+1}))/(x_{j+1} - x_j)) \\ - \cos^{-1}((2x_i - (x_j + x_{j+1}))/(x_{j+1} - x_j)) | \} .$$

If  $\beta(X) < 1/9$ , then

(3.34) 
$$||\mathbf{f} - \mathbf{s}_{LS}|| \leq [(2 - 9\beta(\mathbf{X}))\mathcal{E}_{m, 0}||\mathbf{f}^{(m)}||\mathbf{h}_{\mathbf{X}}^{m} + ||\mathbf{R}||_{\infty}]/$$

$$(1 - 9\beta(\mathbf{X})),$$

or if  $\hat{\beta}(X) < \sqrt{2}/3$ , then

(3.35) 
$$||f - s_{LS}|| \leq [(2 - (3\hat{\beta}(X))^2/2) \mathcal{E}_{m,0} ||f^{(m)}||h_x^m + ||R||_{\infty}]/(1 - (3\hat{\beta}(X))^2/2),$$

where  $\mathcal{E}_{m,0}$  are given in Table 1 of Chapter 2 and  $h_{\mathbf{x}} = \max_{\substack{1 \leq j \leq M-1}} (x_{j+1} - x_j) \cdot 1_{j \leq M-1}$ 

<u>Proof</u>: Let  $s_f \in S^2(\pi_x)$  be the interpolation spline to  $f \in C^{(m)}[a, b]$ . It follows from Theorem 1.1 of Chapter 2 that

(3.36) 
$$||\mathbf{f} \cdot \mathbf{s}_{LS}|| \leq ||\mathbf{f} \cdot \mathbf{s}_{f}|| + ||\mathbf{s}_{f} \cdot \mathbf{s}_{LS}||$$
  
 $\leq \mathcal{E}_{m,0}||\mathbf{f}^{(m)}||\mathbf{h}_{x}^{m} + ||\mathbf{s}_{f} \cdot \mathbf{s}_{LS}||.$ 

To prove (3.34), we first note that  $s_f - s_{LS}$  is a cubic polynomial on  $[x_j, x_{j+1}]$  for  $1 \le j \le M-1$ . Because  $\beta(X) \le 1/9$ , it follows by applying Lemma 3.1 to each interval  $[x_j, x_{j+1}]$  that

(3.37) 
$$||s_{f} - s_{LS}|| \leq ||s_{f} - s_{LS}||_{X}/(1 - 9\beta(X)).$$

Also, using (3.27) and the triangle inequality

$$(3.38) \qquad ||s_{f} - s_{LS}||_{X} \leq \max_{1 \leq i \leq \widehat{M}} \{|(f - s_{f})(\widehat{x}_{i})| + |(f - s_{LS})(\widehat{x}_{i})|\} \\ \leq ||f - s_{f}|| + ||R||_{\infty} \\ \leq \mathcal{E}_{m,0} ||f^{(m)}||h_{x}^{m} + ||R||_{\infty} .$$

Combining (3.36), (3.37) and (3.38) yields (3.34). The proof of (3.35) is identical to that of (3.34) and is omitted, thus completing our proof.

<u>Remark 3.5</u>: Because of the high density required of the data points, X, it is not necessary to specify an estimate for  $f'(\mathbf{x}_k)$ , k = 1, M, and therefore also unnecessary to employ a constrained least squares algorithm to solve for  $s_{LS}$ .

Remark 3.6: As was done in (3.26) we have

(3.39)  $||\mathbf{R}||_{\infty} \leq ||\mathbf{R}||$ 

$$\leq \sqrt{\hat{M}} \ \mathcal{E}_{m,0} ||f^{(m)}|| h_x^m$$

The bound (3. 39) can be used in Theorem 3.2 to obtain an a priori bound for  $||f-s_{LS}||$ , where we only need to know X and  $f \in C^{(m)}[a, b]$ . However, because  $\sqrt{\hat{M}} > \sqrt{M-1} \gg h_x^{-1/2}$ , the best a priori bound obtainable using this method can not have the exponent of h exceed m-1/2.

<u>Remark 3.7</u>: The observation should be made that a uniform mesh is required for an application of Theorem 3.1. However, the conclusion of Theorem 3.2 is independent of the mesh spacing.

This section is concluded by recording a few observations on the distribution of the data points required by each of the above The hypotheses of Theorem 3.1 are satisfied if there theorems. exists at least three data points in each interval  $[x_i, x_{i+1}]$  such that the spacing satisfies (2.1a), (2.1b) and in (2.4) we have  $\gamma > 0$ . The observation should be made that this is a reasonably weak condition to be placed on X. For example, in each interval, all of the three data points could lie in just half the interval, say  $[x_i, x_i+h_x/2]$ and  $\gamma > 0$  could still be realized. The data points  $x_i^l = h_x/3 + x_i$ ,  $x_i^2 = 5h_x/12 + x_i$  and  $x_i^3 = h_x/2 + x_i$ , all of which lie in the half interval, are surely acceptable, and give the value i = 1/4 > 0. If X contains an excess of data points beyond that required to fulfill the requirements (2.1a) and (2.1b), then a judicious selection of the three data points required for each interval can maximize the quantities  $\gamma$ ,  $\alpha$  and  $\hat{\alpha}$  and therefore minimize the upper bound given in (3.13) for  $||\mathbf{f}-\mathbf{s}_{\mathsf{T},\mathbf{S}}||$ .

For  $\beta(X) < 1/9$  in Theorem 3.2, the number of data points in each interval  $[x_j, x_{j+1}]$  must exceed 9 and be distributed in a reasonably uniform manner throughout each interval. Usually for unstructured data, this will require the number of data points in each interval to be far in excess of 9. In Theorem 3.2, the number of data points in each interval must exceed  $3\pi/(2\sqrt{2}) \cong 3.3$  to insure that  $\hat{\beta}(X) < \sqrt{2}/3$ . However, uniform distribution of the data points is not what is required. Examination of (3.33) discloses that more data points are required near the end points of the interval rather than near the center. Section 4. Bivariate Least Squares With Data on Mesh Lines

Let  $\pi_x : a = x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y : c = y_1 < y_2 < \cdots < y_N = d$  be univariate meshes on [a, b] and [c, d] respectively. On  $[a, b] \times [c, d]$  define the bivariate mesh  $\pi_x \oplus \pi_y$  which consists of the M vertical mesh lines  $x = x_1$  and N horizontal mesh lines  $y = y_1$ .

If unstructured data is given only on the mesh lines,  $\pi_y \oplus \pi_y$ , then it is possible to avoid solving a matrix problem of high dimension by solving a discrete least squares problem on each vertical and horizontal mesh line. The M + N discrete least squares solutions are then blended with natural splines to obtain a bivariate approximation. Toward this end, we note that sufficient theory has been developed to yield several algorithms and corresponding error bounds.

For the first algorithm, the M + N univariate discrete least squares cubic splines are constructed on each of the mesh lines of  $\pi_{\mathbf{x}} \bigoplus \pi_{\mathbf{y}}$  from the cubic spline spaces  $S^2(\pi_{\mathbf{x}})$  and  $S^2(\pi_{\mathbf{y}})$ . If we have sufficient data points on each of the M + N mesh lines, then Theorem 3.1 or Theorem 3.2 would yield an error estimate for each of these splines. If  $f \in C^{(m,n)}([a,b]\mathbf{x}[c,d]), 1 \le m, n \le 4$ , and the residual vector for each of the discrete least squares splines is small enough, then each would also be an approximation of order m or n, where  $1 \le m, n \le 4$ . Hence, we could blend these splines with linear functions, which would give a bivariate approximation to f which is of order min  $\{m, 2\} + \min \{n, 2\}$ .

Therefore, we introduce linear blending, see Gordon  $\begin{bmatrix} 12 \end{bmatrix}$ . In the notation of Chapter 2, define  $V(\pi_x, M, 0)$  to be the interpolation vector space of piecewise linear continuous functions, such that the basis functions  $\{\phi_i\}_{i=1}^M$  are defined by

(4.1) 
$$\phi_{i}(x) = \begin{cases} (x - x_{i-1})/(x_{i} - x_{i-1}) & \text{if } x_{i-1} \leq x \leq x_{i} \\ (x_{i+1} - x)/(x_{i+1} - x_{i}) & \text{if } x_{i} \leq x \leq x_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

and the interpolation function  $\alpha$  is defined by  $\alpha(i) = 0$  for

 $l \leq i \leq M$ . Define the corresponding space  $V(\pi_y, N, 0)$  with basis  $\{\psi_j\}_{j=1}^N$  and interpolation function  $\beta$ , of piecewise linear continuous functions on the mesh  $\pi_y$  in an identical fashion.  $V(\pi_x, M, 0)$  (and correspondingly  $V(\pi_y, N, 0)$ ) has the following error analysis: if  $g \in C^{(m)}[a, b]$  then

(4.2) 
$$||g - \sum_{i=1}^{M} g(x_i) \phi_i|| \leq K_m^* ||f^{(m)}||h_x^m$$

where  $1 \le m \le 2$ ,  $K_1^* = 2$  and  $K_2^* = 1/8$ . When m = 1, the proof of (4.2) follows from [24, pp. 248-249], where, on the interval  $\begin{bmatrix} x_i, x_{i+1} \end{bmatrix}$ 

(4.3) 
$$|g(\mathbf{x}) - \sum_{i=1}^{M} g(\mathbf{x}_{i}) \phi_{i}(\mathbf{x})| = |(\mathbf{x} - \mathbf{x}_{i})(\mathbf{x} - \mathbf{x}_{i+1})g[\mathbf{x}_{i}, \mathbf{x}_{i+1}, \mathbf{x}]$$
$$= |(\mathbf{x} - \mathbf{x}_{i+1}) \{g[\mathbf{x}_{i+1}, \mathbf{x}] - g[\mathbf{x}_{i}, \mathbf{x}_{i+1}]\}|$$
$$\leq 2h_{\mathbf{x}} ||g^{(1)}|| .$$

The proof of the case where m = 2 follows from [24, p. 248]. The error in linear blending interpolation to a function  $f \in C^{(m, n)}$ 

 $([a, b]x[c, d]), 1 \leq m, n \leq 2$  is given by Theorem 1.4 of Chapter 2 as

(4.4) 
$$||f - P_{\mathbf{x}} \bigoplus P_{\mathbf{y}} f|| \leq K_{\mathbf{m}}^* K_{\mathbf{n}}^* ||f^{(\mathbf{m},\mathbf{n})}||h_{\mathbf{x}}^{\mathbf{m}} h_{\mathbf{y}}^{\mathbf{n}}$$

where  $K_1^* = 2$ ,  $K_2^* = 1/8$ ,  $h_x = \max_{1 \le i \le M-1} (x_{i+1} - x_i)$  and

$$\underset{\mathbf{y}}{\overset{\mathbf{h}}{\underset{1 \leq j \leq N-1}{\max}}} (y_{j+1} - y_j) .$$

If  $h = \max(h_x, h_y)$  and m = n = 2, then in terms of h, this gives a fourth order approximation to f.

We now define stepwise an algorithm, which was described above. This algorithm will be designated as Algorithm 4.1.

## Algorithm 4.1: Linear Blended Discrete Least Squares with Approximation to Boundary Derivaties.

Let  $\pi_x :a = x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y :c = y_1 < y_2 < \cdots < y_N = d$  be uniform univariate meshes, and let  $\pi_x \oplus \pi_y$  be the bivariate mesh of mesh lines defined above. For  $1 \le j \le N$ , let

 $X_{j} = \{\hat{x}_{ji}\}_{i=1}^{\hat{M}_{j}} \leq [a, b] \text{ be an unstructured discrete data set of } \hat{M}_{j}$ points, and for  $1 \leq i \leq M$  let  $Y_{i} = \{\hat{y}_{ij}\}_{j=1}^{\hat{N}_{i}} \leq [c, d]$  be an unstructured discrete data set of  $\hat{N}_{i}$  points.

<u>Step 1</u>: Let  $f \in C^{(m,n)}([a,b]x[c,d])$  where  $1 \leq m,n \leq 4$ . For  $1 \leq j \leq N$  construct  $s_{LS,j}^{\mathbf{x}} \in S^{2}(\pi_{\mathbf{x}})$ , which is an univariate discrete least squares cubic spline which minimizes the Euclidean norm of the residual vector  $R_{j}^{\mathbf{x}} \in \mathbf{R}$ , where component i of  $R_{j}^{\mathbf{x}}$  is given by

(4.5) 
$$(R_{j}^{x})_{i} = f(\hat{x}_{ji}, y_{j}) - s_{LS, j}(\hat{x}_{ji})$$

and  $\hat{x}_{ji} \in X_j$ . Also for each k, k=1, M,  $(s_{LS,j}^x)'(x_k)$  is constrained to equal the approximation to  $f^{(1,0)}(x_k, y_j)$  given by the Lagrange interpolation polynomial of degree m-1 constructed on the m data points  $\{\hat{x}_{ji}\}_{i=1}^m \subseteq X_j$ , where  $\{\hat{x}_{ji}\}_{i=1}^m \subseteq [x_1, x_2]$  for k=1 or  $\subseteq [x_{M-1}, x_M]$  for k=2, which interpolates the values of  $f(\hat{x}_{ji}, y_j)$  (see Remark 1.1 for procedure). Correspondingly, for  $1 \leq i \leq M$ , construct  $s_{LS,i}^y \in S^2(\pi_y)$ , which is a univariate discrete least squares cubic spline which minimizes the Euclidean norm of the residual vector  $R_i^y \in \mathbb{R}^n$ , where component j of  $R_i^y$  is given by

(4.6) 
$$(R_{i}^{y})_{j} = f(x_{i}, \hat{y}_{ij}) - s_{LS, i}^{y}(\hat{y}_{ij}),$$

and  $\hat{y}_{ij} \in Y_i$ . Also, for each k, k=1, N,  $(s_{LS,i}^y)'(y_k)$  is constrained to equal the approximation to  $f^{(0,1)}(x_i, y_k)$  given by the Lagrange

interpolation polynomial of degree n-1 constructed on the n data  
points 
$$\{\hat{y}_{ij}\}_{j=1}^{n} \subseteq Y_{i}$$
, where  $\{\hat{y}_{ij}\}_{j=1}^{n} \subseteq [y_{1}, y_{2}]$  for k=1, or  
 $\subseteq [y_{N-1}, y_{N}]$  for k=N, which interpolates the values  $f(x_{i}, \hat{y}_{ij})$ .

Step 2: Linearly blend the M + N discrete least squares cubic splines to obtain the following approximation to f

(4.7) 
$$s(x, y) = \sum_{i=1}^{M} s_{LS, i}^{y}(y) \phi_{i}(x) + \sum_{j=1}^{N} s_{LS, j}^{x}(x) \psi_{j}(y) - \sum_{i=1}^{M} \sum_{j=1}^{N} (1/2)(s_{LS, i}^{y}(y) + s_{LS, j}^{x}(x_{i})) \phi_{i}(x) \psi_{j}(y) .$$

In order to obtain an error estimate, define the following parameters of the data sets  $X_j$  and  $Y_i$ . For each j,  $1 \le j \le N$ , assume that there are m distinct data points of  $X_j$  in each of the intervals  $[x_1, x_2]$  and  $[x_{M-1}, x_M]$ , and for 1 < i < M, assume that there exists fixed data points in  $X_j$  satisfying (2.1a) and also fixed data points satisfying (2.1b) with respect to the mesh  $\pi_x$ . For these fixed data points in  $X_j$ , define the real numbers  $\int_j^x$ ,  $\alpha_j^x$  and  $\hat{\alpha}_j^x$  with respect to the mesh  $\pi_x$  by (2.4), (2.5a) and (2.5b), respectively (recall the notation that  $x_i = \overline{x}_i^0 = x_i$  is a knot of the mesh  $\pi_x$ ). Finally, define the real numbers

$$(4.8) \qquad \qquad \gamma^{\mathbf{x}} = \min \gamma^{\mathbf{x}}_{j}, \\ \mathbf{1} \leq \mathbf{j} \leq \mathbf{N}$$

(4.9) 
$$\alpha^{\mathbf{X}} = \min_{\substack{\alpha \\ 1 \leq j \leq N}} \alpha^{\mathbf{X}}_{j} > 0 ,$$

and

(4.10) 
$$\hat{\alpha}^{\mathbf{X}} = \min_{\substack{\alpha \\ l \leq j \leq N}} \hat{\alpha}^{\mathbf{X}}_{j} > 0.$$

Correspondingly, for each i,  $1 \le i \le M$ , identical assumptions on the data sets  $Y_i$ , with respect to the mesh  $\pi_y$ , are made. Also, the real numbers  $\delta^y$ ,  $\alpha^y > 0$  and  $\hat{\alpha}^y > 0$  are defined in a manner analogous to (4.8), (4.9) and (4.10).

For Algorithm 4.1, the following error estimate is valid.

<u>Theorem 4.1</u>. Let s, defined in (4.7), be constructed by Algorithm 4.1. If  $f \in C^{(m,n)}([a,b]x[c,d]), \quad \chi^x > 0$  and  $\chi^y > 0$ , then

$$(4.11) \qquad ||f-s|| \leq K_{m^{*}}^{*} K_{n^{*}}^{*} ||f^{(m^{*}, n^{*})}||h_{x}^{m^{*}} h_{y}^{n^{*}} + (3/2) \Big[ \mathcal{E}_{m, 0}^{(21(1-\alpha^{x})(1-2\alpha^{x})^{2}/(4\alpha^{x}y^{x}) + 1)} + (7(1-\alpha^{x})(1-2\alpha^{x})/(4y^{x}) + 1/4)/(m-1) !\Big] ||f^{(m, 0)}||h_{x}^{m} + (63/2)(1-\alpha^{x})(1-2\alpha^{x})^{2} \max_{l \leq j \leq N} ||R_{j}^{x}||_{\infty} / (4\alpha^{x}y^{x}) + (3/2) \Big[ \mathcal{E}_{n, 0}^{(21(1-\alpha^{y})(1-2\alpha^{y})^{2}/(4\alpha^{y}y^{y}) + 1)} + (7(1-\alpha^{y})(1-2\alpha^{y})/(4y^{y}) + 1/4)/(n-1) !\Big] ||f^{(0, n)}||h_{y}^{n} + (63/2)(1-\alpha^{y})(1-2\alpha^{y})^{2} \max_{l \leq i \leq M} ||R_{i}^{y}||_{\infty} / (4\alpha^{y}y^{y}),$$

where  $m^* = \min \{m, 2\}$ ,  $n^* = \min \{n, 2\}$ ,  $\mathcal{E}_{m, 0}$  and  $\mathcal{E}_{n, 0}$  are

given in Table 1 of Chapter 2,  $K_1^* = 2$ ,  $K_2^* = 1/8$ ,  $h_x = (b-a)/(M-1)$ and  $h_y = (d-c)/(N-1)$ .

Proof: From (4.4), we have

$$(4.12) \qquad ||\mathbf{f}-\mathbf{s}|| \leq ||\mathbf{f}-\mathbf{P}_{\mathbf{x}} \bigoplus \mathbf{P}_{\mathbf{y}}[\mathbf{f}]|| + ||\mathbf{P}_{\mathbf{x}} \bigoplus \mathbf{P}_{\mathbf{y}}[\mathbf{f}]_{-\mathbf{s}}||$$
$$\leq K_{\mathbf{m}^{*}}^{*} K_{\mathbf{n}^{*}}^{*} ||\mathbf{f}^{(\mathbf{m}^{*},\mathbf{n}^{*})}||\mathbf{h}_{\mathbf{x}}^{\mathbf{m}^{*}} \mathbf{h}_{\mathbf{y}}^{\mathbf{n}} + ||\mathbf{P}_{\mathbf{x}} \bigoplus \mathbf{P}_{\mathbf{y}}[\mathbf{f}]_{-\mathbf{s}}||.$$

Writing out  $P_x \bigoplus P_y[f] - s$  and using (4.7), gives

(4.13) 
$$P_{\mathbf{x}} \bigoplus P_{\mathbf{y}}[f] - \mathbf{s} = \sum_{i=1}^{M} (f(\mathbf{x}_{i}, \cdot) - \mathbf{s}_{\mathrm{LS}, i}^{\mathbf{y}})\phi_{i}$$
$$+ \sum_{j=1}^{N} (f(\cdot, \mathbf{y}_{j}) - \mathbf{s}_{\mathrm{LS}, j}^{\mathbf{x}})\psi_{j}$$
$$M = N$$

$$\sum_{\substack{i=1 \ j=1}}^{\infty} \sum_{j=1}^{(1/2)(f(\mathbf{x}_i, \mathbf{y}_j) - \mathbf{s}_{LS, i}^{\mathbf{y}}(\mathbf{y}_j) + f(\mathbf{x}_i, \mathbf{y}_j) - \mathbf{s}_{LS, j}^{\mathbf{x}}(\mathbf{x}_i))\phi_i \psi_j .$$

Taking absolute values, observing that both  $\varphi_i$  and  $\psi_j$  are non- i

negative, 
$$\begin{array}{ccc} M & N\\ \Sigma & \phi_i = 1 \ \text{and} \ \ \Sigma & \psi_j = 1, \ \text{then (4.13) reduces to}\\ i=1 & j=1 \end{array}$$

(4.14) 
$$|(P_x \oplus P_y[f]-s)(x, y)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| \leq (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i, \cdot)| < (3/2) \max_{\substack{i \leq M \\ i \leq M}} ||f(x_i,$$

$$- s_{\text{LS,i}}^{y} || + (3/2) \max_{1 \leq j \leq N} || f(\cdot, y_j) - s_{\text{LS,j}}^{x} || .$$

Corollary 3.1 yields

(4.15) 
$$\max_{l \leq i \leq M} ||f(x_{i}, \cdot) - s_{LS, i}^{y}||$$

$$\leq \left[ \mathcal{E}_{n, 0}^{(21(1 - \alpha^{y})(1 - 2\alpha^{y})^{2}/(4\alpha^{y}y^{y}) + 1)} + (7(1 - \alpha^{y})(1 - 2\alpha^{y})/(4y^{y}) + 1/4)/(n - 1)! \right] ||f^{(0, n)}|| h_{y}^{n}$$

$$+ 21(1 - \alpha^{y})(1 - 2\alpha^{y})^{2} \max_{l \leq i \leq M} ||R_{i}^{y}||_{\infty}/(4\alpha^{y}y^{y}),$$

with a corresponding bound for  $\max_{1 \leq j \leq N} ||f(\cdot, y_j) - s_{LS, j}^{x}||$ . Com-

bining (4.12), (4.14) and (4.15) completes the proof of the Theorem.

We first observe that if  $T_1$  is the total number of data points, then  $T_1$  must satisfy

(4.16) 
$$T_1 = \sum_{i=1}^{M} \hat{N}_i + \sum_{j=1}^{N} \hat{M}_j \ge 6 M N - (M+N),$$

for sufficient data points to be available to apply Theorem 4.1. Next, because only univariate least squares is being performed, the size (total number of entries) of the largest matrix that must be stored in the computer at one time is given by

(4.17) 
$$T_{2} = \max \{ (M+2) \max \bigwedge_{\substack{i \leq j \leq N}}^{n} (N+2) \max \bigotimes_{\substack{i \leq i \leq M}}^{n} \}.$$

162

Instead, if a bivariate least squares solution was obtained using the discretized blending function space of linearly blended cubic splines, the matrix problem to be solved would be significantly larger than  $T_2$ , (usually of the order of  $(T_2)^2$ ).

Example 4.1: Consider the special case when m = n = 4 and  $h = max (h_x, h_y)$ . Hence  $m^* = n^* = 2$  and (4.11) gives a bound in which the power of h is four. Also, an a priori bound can be obtained using Theorem 4.1 and the reasoning of (3.26) on each residual vector  $R_j^x$  and  $R_i^y$ . However, the highest realizable power of h is 3 1/2.

If each of the data sets  $X_j$  and  $Y_i$  has sufficient data points, the following algorithm can be applied, where it is not necessary to approximate the normal derivatives or to use uniform meshes.

Algorithm 4.2: Let  $\pi_x := x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y := y_1 < y_2 < \cdots < y_N = d$  be univariate meshes and let  $\pi_x \oplus \pi_y$  be the corresponding bivariate mesh. For  $1 \le j \le N$ , let  $X_j = \{\hat{x}_{ji}\}_{i=1}^{\hat{M}_j} \subseteq [a, b]$  be an unstructured discrete data set of  $\hat{M}_j$  points and for  $1 \le i \le M$  let  $Y_i = \{\hat{y}_{ij}\}_{j=1}^{\hat{N}_i} \subseteq [c, d]$  be discrete unstructured data sets of  $\hat{N}_j$  points.

<u>Step 1</u>: Let  $f \in C^{(m,n)}([a,b]x[c,d])$ , where  $l \leq m, n \leq 4$ . For  $l \leq j \leq N$ , construct  $s_{LS,j}^{\mathbf{x}} \in S^{2}(\pi_{\mathbf{x}})$  which is a univariate discrete

least squares cubic spline which minimizes the Euclidean norm of the residual vector  $R_y^{\mathbf{x}} \in \mathbb{R}^{\hat{M}_j}$  given by (4.5). Correspondingly, for  $l \leq i \leq M$ , construct  $s_{LS,i}^{\mathbf{y}} \in S^2(\pi_y)$  which minimizes the norm of the residual vector  $R_i^{\mathbf{y}} \in \mathbb{R}^{\hat{N}_i}$  defined in (4.6).

<u>Step 2</u>: Linearly blend the discrete least squares splines to obtain the bivariate approximation s defined in (4.7).

In order to obtain an error bound, calculate for each j,  $1 \leq j \leq N$ , the real numbers  $\beta(X_j)$  and  $\hat{\beta}(X_j)$  with respect to the mesh  $\pi_x$  from (3.32) and (3.33), respectively. Define the real numbers  $\beta_x$  and  $\hat{\beta}_x$  by

(4.19) 
$$\hat{\beta}_{x} = \max \hat{\beta}(X_{j})$$
  
 $l \leq j \leq N$ 

Correspondingly, for each i,  $1 \le i \le M$  calculate the numbers  $\beta(Y_i)$  and  $\hat{\beta}(Y_i)$  with respect to the mesh  $\pi_y$  from (3.32) and (3.33), respectively, and define the real numbers  $\beta_y$  and  $\hat{\beta}_y$  by

(4.20) 
$$\beta = \max_{\substack{y \\ l \leq i \leq M}} \beta(Y_i)$$

(4.21) 
$$\hat{\beta}_{y} = \max_{\substack{y \\ l \leq i \leq M}} \hat{\beta}(Y_{i}).$$

The following error estimate is a consequence of Theorem 3.2.

$$\begin{aligned} \frac{\text{Theorem 4.2.}}{\text{4.2 and } f \in C^{(m,n)}([a,b]_{x}[c,d]). & \text{If } \beta_{x} < 1/9 \text{ and } \beta_{y} < 1/9, \text{ then}}{4.2 \text{ and } f \in C^{(m,n)}([a,b]_{x}[c,d]). & \text{If } \beta_{x} < 1/9 \text{ and } \beta_{y} < 1/9, \text{ then}} \\ (4.22) & ||f-s|| \leqslant K_{m^{*}}^{*} K_{n^{*}}^{*} ||f^{(m^{*},n^{*})}||h_{x}^{m^{*}} h_{y}^{n^{*}} \\ & + (3/2) \left[ (2-9 \beta_{x}) \mathcal{E}_{m,0} ||f^{(0,n)}||h_{x}^{m} + \max_{1 \leqslant j \leqslant N} ||R_{j}^{x}||_{\infty} \right] / (1-9\beta_{x}) \\ & + (3/2) \left[ (2-9\beta_{y}) \mathcal{E}_{n,0} ||f^{(0,n)}||h_{y}^{n} + \max_{1 \leqslant j \leqslant M} ||R_{1}^{y}||_{\infty} \right] / (1-9\beta_{y}), \end{aligned}$$
or if  $\beta_{x} < \sqrt{2}/3$  and  $\beta_{y} < \sqrt{2}/3$ , then
$$(4.23) \qquad ||f-s|| \leqslant K_{m^{*}} K_{n^{*}} ||f^{(m^{*},n^{*})}||h_{x}^{m^{*}} h_{y}^{n^{*}} \\ & + (3/2) \left[ (2-(3\beta_{x})^{2}/2) \mathcal{E}_{m,0} ||f^{(m,0)}||h_{x}^{m} + \max_{1 \leqslant j \leqslant N} ||R_{j}^{x}||_{\infty} \right] \\ & / (1-(3\beta_{x})^{2}/2) \\ & + (3/2) \left[ (2-(3\beta_{y})^{2}/2) \mathcal{E}_{n,0} ||f^{(0,n)}||h_{y}^{n} \\ & + \max_{1 \leqslant j \leqslant N} ||R_{1}^{y}||_{\infty} \right] / (1-(3\beta_{y})^{2}/2), \end{aligned}$$
where  $m^{*} = \min\{m, 2\}, n^{*} = \min\{n, 2\}, \mathcal{E}_{m,0}$  and  $\mathcal{E}_{n,0}$  are given in Table 1 of Chapter 2,  $K_{1}^{*} = 2, K_{2}^{*} = 1/8, \end{aligned}$ 

$$h_{\mathbf{x}} = \max_{\substack{1 \leq i \leq M-1}} (x_{i+1} - x_i) \text{ and } h'_{\mathbf{y}} = \max_{\substack{1 \leq j \leq N-1}} (y_{j+1} - y_j).$$

<u>Proof</u>: The proof of Theorem 4.2 is essentially identical to the proof of Theorem 4.1. We merely note that combining (4.12), (4.14) and the bounds of Theorem 3.2 yields our result.

Algorithms 4.1 and 4.2 are hybrid algorithms, which combine piecewise cubic and piecewise linear splines. If it is desired that cubic splines be used throughout, then the following algorithm, outlined below, may be useful.

We blend our discrete least squares cubic splines with natural cubic splines (see Section 3 of Chapter 2). If  $f \in C^{(m, n)}([a, b]_x[c, d])$ ,  $2 \leq m, n \leq 4$ , then this method has a potential accuracy of order m + n in the "interior of the region", (see Sections 3 and 4 of Chapter 2). Examination of Theorem 4.2 of Chapter 2 indicates that preservation of this accuracy requires the discrete least squares cubic spline approximations to be of order m + n, when their residual vectors are small. To accomplish this, on each mesh line of  $\pi_x \oplus \pi_y$ , we replace the meshes  $\pi_x$  and  $\pi_y$  by uniform univariate meshes  $\overline{\pi_x}$ :  $a = \overline{x_1} < \overline{x_2} < \cdots < \overline{x_M} = b$  and  $\overline{\pi_y}$ :  $c = \overline{y_1} < \overline{y_2} < \cdots < \overline{y_N} = d$ , respectively. From Theorem 1.1 of Chapter 2, the interpolation accuracy of the univariate cubic spline spaces  $S^2(\overline{\pi_x})$  and  $S^2(\overline{\pi_y})$  is  $O(\overline{h_x}^m)$  and  $O(\overline{h_y}^n)$  respectively, where  $\overline{h_x} = (b-a)/(\overline{M}-1)$  and  $\overline{h_y} = (d-c)/(\overline{N}-1)$ . Therefore, for the preservation of our

accuracy, it is necessary to refine the meshes  $\overline{\pi}_x$  and  $\overline{\pi}_y$ sufficiently to have  $\overline{h}_x^m \leqslant h_x^m h_y^n$  and  $\overline{h}_y^n \leqslant h_y^m h_y^n$ .

<u>Remark</u>: For the special case where m = n = 4 and  $h_x = h_y = h$ , this reduces to  $\overline{h_x} \le h^2$  and  $\overline{h_y} \le h^2$ .

The above observations lead to the following algorithm.

<u>Algorithm 4.3</u>: Let  $\pi_x :a = x_1 < x_2 < \cdots < x_M = b$ ,  $\pi_y :c = y_1 < y_2 < \cdots < y_N = d$ ,  $\overline{\pi}_x :a = \overline{x}_1 < \overline{x}_2 < \cdots < \overline{x}_M = b$  and  $\overline{\pi}_y :c = \overline{y}_1 < \overline{y}_2 < \cdots < \overline{y}_N = d$  be univariate meshes such that  $\overline{\pi}_x$  and  $\overline{\pi}_y$  are uniform,  $\overline{h}_x^m \leq h_x^m h_y^n$  and  $\overline{h}_y^n \leq h_x^m h_y^n$  (these last two conditions are not necessary for Theorem 4.3 to remain valid, however, to preserve the desirability of Algorithm 4.3, attempts should be made to insure that these two conditions hold, or at least nearly hold). Let  $\pi_x \bigoplus \pi_y$  be the bivariate mesh of M+N mesh lines. The data points will only be specified on the mesh  $\pi_x \bigoplus \pi_y$ . For  $1 \leq j \leq N$ , let  $X_j = \{\hat{x}_{ji}\}_{i=1}^{M_j} \subseteq [a, b]$  be a discrete unstructured data set of  $\widehat{M}_j$  points, and for  $1 \leq i \leq M$ , let  $Y_i = \{\hat{y}_{ij}\}_{j=1}^{j-1} \subseteq [c, d]$  be a discrete unstructured data set of  $\widehat{N}_i$  points.

<u>Step 1</u>: Let  $f \in C^{(m, n)}([a, b] \times [c, d])$  where  $2 \leq m, n \leq 4$ . For  $l \leq j \leq N$  construct  $s_{LS, j}^{\mathbf{x}} \in S^{2}(\overline{\pi}_{\mathbf{x}})$ , which is a univariate discrete least squares cubic spline which minimizes the Euclidean norm of

the residual vector  $\mathbf{R}_{i}^{\mathbf{x}} \in \mathbf{R}^{\mathbf{M}_{j}}$  defined by (4.5). Also, for  $k = 1, \overline{M}$ ,  $(s_{LS,i}^{x})'(x_{k})$  is constrained to equal the approximation to  $f^{(1,0)}(\bar{x}_k, y_i)$  given by the Lagrange interpolation polynomial of degree m-l constructed on the m data points  $\{ \hat{x}_{ii} \}_{i=1}^{m} \subseteq X_{i}$ , where  $\{\hat{x}_{ij}\}_{i=1}^{m} \subseteq [\overline{x}_1, \overline{x}_2]$  for k = 1 or  $\subseteq [\overline{x}_{M-1}, \overline{x}_M]$  for  $k = \overline{M}$ , which interpolates the values  $f(\mathbf{x}_{ii}, \mathbf{y}_{i})$ . Correspondingly, for  $l \leq i \leq M$ , construct  $s_{LS,i}^{y} \in S^{2}(\overline{\pi}_{v})$ , which is a univariate discrete least squares cubic spline which minimizes the Euclidean norm of the residual vector  $R_{i}^{y} \in \mathbb{R}^{N_{i}}$  defined by (4.6). Also, for k=1,  $\overline{N}$ ,  $(s_{LS,i}^{y})'(\overline{y}_{k})$  is constrained to equal the approximation to  $f^{(0,1)}(x_i, \overline{y}_k)$  given by the Lagrange interpolation polynomial of degree n-l constructed on the n data points  $\{\hat{y}_i\}_{i=1}^n \subseteq Y_i$ , where  $\{ \hat{\bar{y}}_{ii} \}_{i=1}^{n} \subseteq \left[ \bar{\bar{y}}_{1}, \bar{\bar{y}}_{2} \right] \text{ for } k=1 \text{ or } \subseteq \left[ \bar{\bar{y}}_{\overline{N}-1}, \bar{\bar{y}}_{\overline{N}} \right] \text{ for } k=\overline{N} \text{ which}$ interpolates the values  $f(x_i, \dot{y}_i)$ .

<u>Step 2</u>: Blend the above univariate discrete least squares cubic splines with natural cubic splines (see Section 3 of Chapter 2) to obtain the following bivariate approximation to f, where  $\{A_i\}_{i=1}^M \subseteq S^2(\pi_x)$  and  $\{B_j\}_{j=1}^N \subseteq S^2(\pi_y)$  are natural cardinal basis functions,

(4.24) 
$$\mathbf{s} = \sum_{i=1}^{M} \mathbf{s}_{LS,i}^{\mathbf{y}} \mathbf{A}_{i} + \sum_{j=1}^{N} \mathbf{s}_{LS,j}^{\mathbf{x}} \mathbf{B}_{j}$$
  
-  $\sum_{i=1}^{M} \sum_{j=1}^{N} (1/2)(\mathbf{s}_{LS,i}^{\mathbf{y}}(\mathbf{y}_{j}) + \mathbf{s}_{LS,j}^{\mathbf{x}}(\mathbf{x}_{i})) \mathbf{A}_{i} \mathbf{B}_{j}$ 

The following assumptions are made on the distribution of the data points so that an error analysis can be performed. For each  $j, l \leq j \leq N$ , assume that there exists m distinct data points of  $X_j$  in each interval  $[\bar{x}, \bar{x}_2]$  and  $[\bar{x}_{\bar{M}-1}, \bar{x}_{\bar{M}}]$ , and for  $l \leq i \leq \bar{M}$ , assume that there exists fixed data points in  $X_j$  satisfying (2.1a) and also fixed data points in  $X_j$  satisfying (2.1b) with respect to the mesh  $\bar{\pi}_x$ . For these fixed data points in  $X_j$ , define the real numbers  $\bar{J}_j^x, \bar{\alpha}_j^x$  and  $\frac{\Delta x}{\alpha_j}$  with respect to the mesh  $\bar{\pi}_x$  by (2.4), (2.5a) and (2.5b) respectively. Define the real numbers  $\bar{y}^x, \bar{\alpha}^x$  and  $\frac{\Delta x}{\alpha}$  to be the minimum of  $\bar{y}_j^x, \bar{\alpha}_j^x$  and  $\frac{\Delta x}{\alpha_j}$ , respectively, for  $l \leq j \leq N$ . Correspondingly, for each i,  $l \leq i \leq M$ , identical assumptions on the data sets  $Y_j$ , with respect to the mesh  $\bar{\pi}_y$ , are made. The real numbers  $\bar{\lambda}^y, \bar{\alpha}^y$  and  $\frac{\Delta y}{\alpha}^y$  are defined in an analogous manner.

Algorithm 4.3, along with the assumptions on the distribution of the data points, gives the following theorem.

<u>Theorem 4.3.</u> Let s, defined in (4.24), be constructed by Algorithm 4.3. If  $f \in C^{(m,n)}([a,b]x[c,d])$ ,  $\overline{\gamma}^{x} > 0$  and  $\overline{\delta}^{y} > 0$ , then for  $x \in [x_{i}, x_{i+1}]$  and  $y \in [y_{j}, y_{j+1}]$ (4.25)  $|f(x, y) - s(x, y)| \leq ||f^{(m,n)}|| \{\mathcal{E}_{m,0} h_{x}^{m} + K_{m}(h_{x}^{m-1} \Delta x_{i} \Delta_{i}/4) \{\mathcal{E}_{n,0} h_{y}^{n} + K_{n}(h_{y}^{n-1} \Delta y_{j}) \overline{\Delta}_{j}/4\}$
$$+ (1/8) ||f^{(m,2)}|| \{ \mathcal{E}_{m,0} h_{x}^{m} + K_{m}(h_{x}^{m-1} \Delta x_{i} \Delta_{i}/4) (h_{y} \Delta_{y}) \overline{\Delta}_{j}$$

$$+ (1/8) ||f^{(2,n)}|| \{ \mathcal{E}_{n,0} h_{y}^{n} + K_{n}(h_{y}^{n-1} \Delta_{y}) \overline{\Delta}_{j}/4 \} (h_{x} \Delta x_{i}) \Delta_{i}$$

$$+ (1/64) ||f^{(2,2)}|| \{ (h_{x} \Delta x_{i}) (h_{y} \Delta y_{j}) \} \Delta_{i} \overline{\Delta}_{j}$$

$$+ K(M_{\pi})(1 + K(M_{\pi})/2) \{ [\mathcal{E}_{m,0}(21(1-\overline{\alpha}^{x})(1-2\overline{\alpha}^{x})^{2}/(4\overline{\alpha}^{x}\overline{\delta}^{x}) + 1)$$

$$+ (7(1-\overline{\alpha}^{x})(1-2\overline{\alpha}^{x})/(4\overline{\delta}^{x}) + 1/4)/(m-1) !] ||f^{(m,0)}|| \overline{h}_{x}^{m}$$

$$+ 21(1-\overline{\alpha}^{x})(1-2\overline{\alpha}^{x})^{2} \max_{1 \leq j \leq N} ||R_{j}^{x}||_{\infty}/(4\overline{\alpha}^{x}\overline{\delta}^{x}) \}$$

$$+ K(M_{\pi})(1 + K(M_{\pi})/2) \{ [\mathcal{E}_{n,0}(21(1-\overline{\alpha}^{y})(1-2\overline{\alpha}^{y})^{2}/(4\overline{\alpha}^{y}\overline{\delta}^{y}) + 1)$$

$$+ (7(1-\overline{\alpha}^{y})(1-2\overline{\alpha}^{y})/(4\overline{\delta}^{y}) + 1/4)/(n-1) !] ||f^{(0,n)}|| \overline{h}_{y}^{m}$$

$$+ 21(1-\overline{\alpha}^{y})(1-2\overline{\alpha}^{y})^{2} \max_{1 \leq i \leq M} ||R_{i}^{y}||_{\infty}/(4\overline{\alpha}^{y}\overline{\delta}^{y}) \} ,$$

where  $\mathcal{E}_{m,0}$  and  $\mathcal{E}_{n,0}$  are given in Table 3 of Chapter 2,  $K_{m}$  and  $K_{n}$  are given in Table 5 of Chapter 2,  $\Delta x_{i} = x_{i+1} - x_{i}, \Delta y_{j} = y_{j+1} - y_{j}, h_{x} = \max_{1 \leq i \leq M-1} \Delta x_{i}, h_{y} = \max_{1 \leq j \leq N-1} \Delta y_{j}, \Delta_{i} = \{2^{1-i} + 2^{1-M+i}\}, \overline{\Delta}_{j} = \{2^{1-j} + 2^{1-N+j}\}, M_{\pi} = \max_{x} \Delta x_{i} / \min_{1 \leq i \leq M-1} \Delta x_{i}, M_{\pi} = \max_{x} 1 \leq i \leq M-1$  i  $\lambda x_{i}, M_{\pi} = \max_{x} 1 \leq i \leq M-1$  i  $\lambda y_{j} / \min_{1 \leq j \leq N-1} \Delta y_{j}, K(\xi) = 6\xi(2\xi+1)(\xi+1)^{2}/(3+4\xi), \overline{h}_{x} = (b-a)/(\overline{M}-1)$  and  $\overline{h}_{y} = (d-c)/(\overline{N}-1).$ 

Proof: This is a direct consequence of Corollary 3.1 and Theorem 4.2 of Chapter 2.

<u>Remark</u>: The observation should be made that the order of accuracy depends upon which rectangle,  $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$ , the point (x, y) belongs to, because of the exponential decay of the terms  $\Delta_i$ and  $\overline{\Delta}_j$  toward the center of the region. If the norm of the residual vectors are small enough, m = n = 4,  $h = \max(h_x, h_y)$ ,  $\overline{h_x} \leq h^2$ and  $\overline{h_y} \leq h^2$ , then the error satisfies Figure 1 of Chapter 2, i.e., eighth order near the center and less near the boundary (see Section 3 of Chapter 2).

<u>Remark</u>: The continuity required by Theorem 4.1 and Theorem 4.2 is  $C^{(m,0)} \cap C^{(0,n)} \cap C^{(m^*,n^*)} \ge C^{(m,n)}$ , where  $m^* = \min \{m, 2\}$ and  $n^* = \min \{n, 2\}$ . However, Theorem 4.3 requires continuity of  $C^{(m,0)} \cap C^{(0,n)} \cap C^{(m,n)} = C^{(m,n)}$ , hence, using cubic splines throughout the algorithm makes more efficient use of the available continuity.

172

Section 5. Bivariate Least Squares with Unstructured Data Let  $X = \{(\hat{x}_i, \hat{y}_i)\}_{i=1}^{\hat{M}} \subseteq [a, b] \times [c, d]$  be a set of  $\hat{M}$  unstructured data points. We will describe here several bivariate finite dimensional vector spaces from which discrete least squares fits can be calculated on the data set X. The first two spaces, bicubic splines and cubically blended cubic splines are known (see 29], p. 49 and  $\begin{bmatrix} 13 \end{bmatrix}$ , respectively). The Hermite blended piecewise polynomials are new and have some interesting properties. For this chapter, define the real number  $||\cdot||_{X}$  of a bivariate function  $f \in C^{(0, 0)}([a, b]x[c, d])$  to be

(5.1) 
$$||f||_{X} = \max_{\substack{i \leq i \leq M}} |f(\hat{x}_{i}, \hat{y}_{i})|.$$

# Section 5.1. Bicubic Splines

Let  $\pi_x := x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y := y_1 < y_2 < \cdots < y_N < y_N < \cdots < y_N < \cdots < y_N < y_N < \cdots <$  $y_{N} = d$  be univariate meshes, then the space of bicubic splines is defined to be  $S^{2}(\pi_{x}) \otimes S^{2}(\pi_{v})$ , where  $S^{2}(\pi_{x})$  and  $S^{2}(\pi_{v})$  are cubic spline spaces of Example 1.2 of Chapter 2 (see [29, p. 49]). From the remark following (2.1) of Chapter 2, if  $\{u_i\}_{i=1}^{M+2}$  and  $\{w_j\}_{j=1}^{N+2}$ are any bases for  $S^{2}(\pi_{x})$  and  $S^{2}(\pi_{y})$ , respectively, then  $\{u, w_j\}_{1 \leq i \leq M+2, \ l \leq j \leq N+2}$  is a basis for  $S^2(\pi_x) \bigotimes S^2(\pi_y)$ , and the dimension of this space is (M+2)(N+2).

Let  $f \in C^{(m,n)}([a,b]x[c,d]), 1 \leq m, n \leq 4$ , then  $s \in S^{2}(\pi_{-})$  (x)  $S^{2}(\pi_{v})$  is a discrete least squares cubic spline fit (for f on X) if it minimizes the Euclidean norm of the residual vector  $\mathbf{R} \in \mathbb{R}^{\widehat{M}}$ , where component i is given by

(5.2) 
$$R_i = f(\hat{x}_i, \hat{y}_i) - s(\hat{x}_i, \hat{y}_i),$$

for 
$$(\hat{x}_i, \hat{y}_i) \in X$$
 and  $l \leq i \leq \hat{M}$ .

Before deriving an error analysis, the bivariate generalization of Lemma 3.1 is given (see  $\begin{bmatrix} 6 \\ p \\ 91 \end{bmatrix}$ ).

Lemma 5.1. Let Q be a bivariate polynomial of degree  $\leq n$  in both variables, then

(5.3) 
$$||Q|| (1-n^2 \overline{\beta}(X)) \leq ||Q||_X$$
, where

(5.4) 
$$\overline{\beta}(X) = \max \min \{2 | x - \hat{x}_i | / (b-a)$$
  
 $a \leq x \leq b \\ c \leq y \leq d \quad (\hat{x}_i, \hat{y}_i) \in X$ 

$$+ 2|y-\hat{y}_i|/(c-d) \}$$
,

and

(5.5) 
$$||Q|| (1 - \frac{1}{2}(n \beta(X))^2) \leq ||Q||_X$$
, where

(5.6) 
$$\hat{\beta}(X) = \max_{\substack{a \leq x \leq b \\ c \leq y \leq d}} \min_{\substack{a \leq x \leq b \\ c \leq y \leq d}} \{ |\cos^{-1}((2x - (a+b))/(b-a)) | \\ -\cos^{-1}((2x - (a+b))/(b-a)) | \\ + |\cos^{-1}((2y - (c+d))/(d-c)) - \cos^{-1}((2y - (c+d))/(c-d)) | \}$$

<u>Proof</u>: The proof of (5.3) is along similar lines as the proof of (5.5)and is omitted. The proof of the univariate case of (5.5) is given in [6, pp. 91-92].

Let  $(\bar{x}, \bar{y}) \in [a, b] \times [c, d]$  be such that  $|Q(\bar{x}, \bar{y})| = ||Q||$ , then from (5.6), there exists  $(\hat{x}, \hat{y}) \in X$  such that if  $\bar{\alpha} = \cos^{-1}((2\bar{x} - (a+b))/(b-a))$ ,  $\bar{\theta} = \cos^{-1}((2\bar{y} - (c+d))/(d-c))$ ,  $\hat{\alpha} = \cos^{-1}((2\hat{x} - (a+b))/(b-a))$ and  $\hat{\theta} = \cos^{-1}((2\hat{y} - (c+d))/(d-c))$ , then

(5.7)  $\frac{\hat{\beta}}{\hat{\beta}}(X) \ge |\overline{\alpha} - \hat{\alpha}| + |\overline{\theta} - \hat{\theta}|$ .

Define the bivariate trigonometric polynomial of degree n to be  

$$R(\alpha, \theta) = Q(((b-a)\cos\alpha+a+b)/2, ((d-c)\cos\theta+c+d)/2). \text{ Because}$$

$$||Q|| = |R(\overline{\alpha}, \overline{\theta})| = ||R||_{[-\pi, 0]x[-\pi, 0]} = ||R||_{[-\infty, \infty]x[-\infty, \infty]} \text{ then}$$
(5.8) 
$$R^{(1, 0)}(\overline{\alpha}, \overline{\theta}) = R^{(0, 1)}(\overline{\alpha}, \overline{\theta}) = 0,$$

and we have

(5.9) 
$$R(\overline{\alpha}, \overline{\theta}) = \int_{C} R^{(1, 0)}(\alpha, \theta) d\alpha + R^{(0, 1)}(\alpha, \theta) d\theta + R(\widehat{\alpha}, \widehat{\theta}),$$

where C is a straight line from  $(\hat{\alpha}, \hat{\theta})$  to  $(\overline{\alpha}, \overline{\theta})$ . If C is parameterized by  $t \in [0, 1]$ , then

(5.10) 
$$R(\alpha, \theta) = \int_{0}^{1} \left[ \frac{R^{(1,0)}(t(\overline{\alpha} - \widehat{\alpha}) + \widehat{\alpha}, t(\overline{\theta} - \widehat{\theta}) + \widehat{\theta}) - R^{(1,0)}(\overline{\alpha}, \overline{\theta})}{(t-1)} \right]$$
(t-1)

+ 
$$\frac{\mathbf{R}^{(0,1)}(\mathbf{t}(\overline{\alpha}-\hat{\alpha})+\hat{\alpha},\mathbf{t}(\overline{\theta}-\hat{\theta})+\hat{\theta})-\mathbf{R}^{(0,1)}(\overline{\alpha},\overline{\theta})}{(t-1)} \left[ dt \right]$$

+  $R(\hat{\alpha}, \hat{\theta})$ .

By the Mean Value Theorem for each t  $\varepsilon$  [0, 1), there exists a  $\widetilde{t}\,\varepsilon$  (t, 1), such that

(5.11) 
$$(R^{(1, 0)}(t(\overline{\alpha} - \hat{\alpha}) + \hat{\alpha}, t(\overline{\theta} - \hat{\theta}) + \hat{\theta}) - R^{(1, 0)}(\overline{\alpha}, \overline{\theta}))/(t-1)$$
$$= R^{(2, 0)}(\tilde{t}(\overline{\alpha} - \hat{\alpha}) + \hat{\alpha}, \tilde{t}(\overline{\theta} - \hat{\theta}) + \hat{\theta}) (\overline{\alpha} - \hat{\alpha})$$
$$+ R^{(1, 1)}(\tilde{t}(\overline{\alpha} - \hat{\alpha}) + \hat{\alpha}, \tilde{t}(\overline{\theta} - \hat{\theta}) + \hat{\theta}) (\overline{\theta} - \hat{\theta}),$$

with a similar expression for the second term in brackets on the right hand side of (5.10). Taking absolute values and substituting (5.11) into (5.10) we have, after two applications of Bernstein's inequality (see  $\begin{bmatrix} 6 \\ 9 \\ 91 \end{bmatrix}$ ),

$$(5.12) \qquad ||Q|| = |\mathbf{R}(\overline{\alpha}, \overline{\theta})|$$

$$\leq ||Q||\mathbf{n}^{2}(|\overline{\alpha}-\widehat{\alpha}|+|\overline{\theta}-\widehat{\theta}|)^{2} \int_{0}^{1} |t-1|dt+|\mathbf{R}(\widehat{\alpha}, \widehat{\theta})|$$

$$\leq ||Q||\mathbf{n}^{2}(|\overline{\alpha}-\widehat{\alpha}|+|\overline{\theta}-\widehat{\theta}|)^{2}/2 + ||Q||_{\mathbf{X}}$$

$$\leq ||Q||(\mathbf{n}\overline{\beta}(\mathbf{X}))^{2}/2 + ||Q||_{\mathbf{X}},$$

where (5.7) has been used and the fact that  $|R(\hat{\alpha}, \hat{\theta})| = |Q(\hat{x}, \hat{y})| \leq ||Q||_X$ .

Define the real numbers  $\overline{\beta}(X, \pi_x, \pi_y)$  and  $\frac{\overline{\beta}}{\beta}(X, \pi_x, \pi_y)$  with respect to the univariate meshes  $\pi_x$  and  $\pi_y$  to be

$$(5.13) \qquad \beta(X, \pi_{x}, \pi_{y}) = \max_{\substack{1 \leq i \leq M-1 \\ 1 \leq j \leq N-1 \\ y \in [y_{j}, y_{j+1}]}} \max_{\substack{(\hat{x}_{k}, \hat{y}_{k}) \in X \\ \hat{y}_{k} \in [x_{i}, x_{i+1}]}} \min_{\substack{(\hat{x}_{k}, \hat{y}_{k}) \in X \\ \hat{y}_{k} \in [y_{j}, y_{j+1}]}}$$

$$\{2|\hat{x}_{k}-x|/(x_{i+1}-x_{i})+2|\hat{y}_{k}-y|/(y_{j+1}-y_{j})\},\$$

and

$$(5.14) \qquad \widehat{\widehat{\beta}}(X, \pi_{x}, \pi_{y}) = \max_{\substack{1 \leq i \leq M-1 \\ 1 \leq j \leq N-1 \\ j \leq N-1 \\ y \in [y_{j}, y_{j+1}] \\ y \in [y_{j}, y_{j+1}] \\ \widehat{x}_{k} \in [x_{i}, x_{i+1}] \\ \widehat{y}_{k} \in [y_{j}, y_{j+1}] \\ \left\{ |\cos^{-1}((2\hat{x}_{k}^{-}(x_{i} + x_{i+1}))/(x_{i+1}^{-} - x_{i})) - \cos^{-1}((2x - (x_{i}^{+} + x_{i+1}))/(x_{i+1}^{-} - x_{i})) + |\cos^{-1}((2x - (x_{i}^{+} + x_{i+1}))/(x_{i+1}^{-} - x_{i}))| + |\cos^{-1}((2\hat{y}_{k}^{-}(y_{j+1}^{-} + y_{j}))/(y_{j+1}^{-} - y_{j})) - \cos^{-1}((2y - (y_{j+1}^{-} - x_{i}^{-})))| + |\cos^{-1}((2\hat{y}_{k}^{-} - (y_{j+1}^{-} + y_{j}))/(y_{j+1}^{-} - y_{j}))| \} .$$

<u>Theorem 5.1</u>. Let  $s \in S^2(\pi_x) \bigotimes S^2(\pi_y)$  be a discrete least squares bicubic spline on the data set X to  $f \in C^{(m, n)}([a, b]x[c, d])$  for  $l \leq m, n \leq 4$ . If  $\overline{\beta}(x) < 1/9$ , then

$$(5.15) \qquad ||f-s|| \leq \{ \mathcal{E}_{m,0} ||f^{(m,0)}||h_x^m + \mathcal{E}_{n,0} ||f^{(0,n)}||h_y^n + \mathcal{E}_{m,0} \mathcal{E}_{n,0} ||f^{(m,n)}||h_x^m h_y^n \} \{ 1 + 1/(1 - 9 \overline{\beta} (X, \pi_x, \pi_y)) \} + ||R||_{\infty} / (1 - 9 \overline{\beta} (X, \pi_x, \pi_y)) .$$

If 
$$\frac{\hat{\beta}}{\hat{\beta}}(X, \pi_{x}, \pi_{y}) < \sqrt{2}/3$$
, then  
(5.16)  $||f-s|| \leq \{\mathcal{E}_{m,0}||f^{(m,0)}||h_{x}^{m} + \mathcal{E}_{n,0}||f^{(0,n)}||h_{y}^{n}$   
 $+ \mathcal{E}_{m,0} \mathcal{E}_{n,0}||f^{(m,n)}||h_{x}^{m} h_{y}^{n}\} \{1 + 1/(1 - (3\beta(X, \pi_{x}, \pi_{y}))^{2}/2)\}$   
 $+ ||R||_{\infty}/(1 - (3\beta(X, \pi_{x}, \pi_{y}))^{2}/2),$ 

where  $\mathcal{E}_{m,0}$  and  $\mathcal{E}_{n,0}$  are given in Table 1 of Chapter 2 and  $h_x$ and  $h_y$  are the mesh sizes of  $\pi_x$  and  $\pi_y$ , respectively.

Proof: From Corollary 1.2 of Chapter 2

$$(5.17) ||f-s|| \leq ||f-P_{x}P_{y}[f]|| + ||P_{x}P_{y}[f]-s||$$

$$\leq \xi_{m,0} ||f^{(m,0)}||h_{x}^{m} + \xi_{n,0}||f^{(0,n)}||h_{y}^{m}$$

$$+ \xi_{m,0}\xi_{n,0} ||f^{(m,n)}||h_{x}^{m}h_{y}^{n} + ||P_{x}P_{y}[f]-s||,$$

where  $P_x P_y[f]$  is the bicubic interpolant to f defined in (1.18) of Chapter 2. Note that  $P_x P_y[f]$ -s is a bivariate cubic polynomial on each rectangle  $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$  for  $1 \le i \le M-1$  and  $1 \le j \le N-1$ , hence from Lemma 5.1

(5.18) 
$$|(P_x P_y[f]-s)(x, y)| \leq ||P_x P_y[f]-s||_X/(1-9\bar{\beta}(X, \pi_x, \pi_y)),$$

because  $\overline{\beta}(X, \pi_x, \pi_y) < 1/9$ . It follows that

$$(5.19) \qquad || \mathbf{P}_{\mathbf{x}} \mathbf{P}_{\mathbf{y}}[\mathbf{f}] - \mathbf{s} ||_{\mathbf{X}} \leq \max_{\substack{\mathbf{i} \leq \mathbf{i} \leq \mathbf{M}}} |(\mathbf{P}_{\mathbf{x}} \mathbf{P}_{\mathbf{y}}[\mathbf{f}] - \mathbf{f})(\hat{\mathbf{x}}_{\mathbf{i}}, \hat{\mathbf{y}}_{\mathbf{i}})| + \max_{\substack{\mathbf{i} \leq \mathbf{i} \leq \mathbf{M}}} |(\mathbf{f} - \mathbf{s})(\hat{\mathbf{x}}_{\mathbf{i}}, \hat{\mathbf{y}}_{\mathbf{i}})|$$

$$\left| \left| \mathbf{P}_{\mathbf{X}} \mathbf{P}_{\mathbf{y}} \left[ \mathbf{f} \right] - \mathbf{f} \right| \right| + \left| \left| \mathbf{R} \right| \right|_{\infty}$$

Using Corollary 1.2 of Chapter 2 and combining (5.17), (5.18) and (5.19) completes the proof.

<u>Remark</u>: Note that we have the choice of using any basis of  $S^2(\pi_x)$ and  $S^2(\pi_y)$ . Usually, bicubic B-splines are chosen because they have the desirable property that they have support on at most sixteen adjacent rectangles, are non-negative, and easy to calculate, see [21]. This property minimizes the number of evaluations of the basis functions which must be made to construct the least squares matrix.

### Section 5.2 Cubically Blended Cubic Splines.

Let  $\pi_x: a = x_1 < x_2 < \cdots < x_M = b$ ,  $\overline{\pi}_x: a = \overline{x}_1 < \overline{x}_2 < \cdots < \overline{x}_M = b$ ,  $\overline{\pi}_y: c = \overline{y}_1 < \overline{y}_2 < \cdots < y_N = d$  and  $\overline{\pi}_y: c = \overline{y}_1 < \overline{y}_2 < \cdots < \overline{y}_N = d$  be univariate meshes with the corresponding cubic spline spaces  $S^2(\pi_x)$ ,  $S^2(\overline{\pi}_x)$ ,  $S^2(\pi_y)$  and  $S^2(\overline{\pi}_y)$ .

From Definition 2.1 of Chapter 2, the discretized blending function space of cubically blended cubic splines,  $DBF_1$ , is the image of  $C^{(1,1)}$  under the map  $\overline{P_x \oplus P_y} = \overline{P_x}P_y + P_x\overline{P_y} - P_xP_y$ .

 $\overline{\pi}_{x}$  is said to be a refinement of the mesh  $\pi_{x}$  if and only if the knots of the mesh  $\pi_{x}$  are knots of the mesh  $\overline{\pi}_{x}$ . For cubic spline spaces, if  $\overline{\pi}_{x}$  is a refinement of  $\pi_{x}$ , then  $S^{2}(\pi_{x})$  is subordinate to  $S^2(\overline{\pi}_x)$  and  $S^2(\pi_x) \subseteq S^2(\overline{\pi}_x)$  (see remark following (2.29) of Chapter 2).

If the meshes  $\overline{\pi}_x$  and  $\overline{\pi}_y$  are refinements of  $\pi_x$  and  $\pi_y$ , respectively, then from Theorem 2.3 of Chapter 2

(5.20) 
$$\operatorname{Dim}(\operatorname{DBF}_{1}) = (M+2)(\overline{N}+2) + (\overline{M}+2)(N+2) - (M+2)(N+2),$$

and  $T_1$  through  $T_4$ , given in that theorem, each form a basis for  $DBF_1$ .

Finally, a bound for  $||f - \overline{P_x + P_y}[f]||$  is given in Example 1.4 of Chapter 2.

Let  $f \in C^{(m, n)}([a, b] \times [c, d])$ ,  $l \leq m, n \leq 4$ , then  $s \in DBF_1$  is a discrete least squares solution on the data set X if the Euclidean norm of the residual vector R defined in (5.2) is minimized. Combining the above observations and Lemma 5.1 yields the following theorem.

Theorem 5.2. Let  $s \in DBF_1$  be a discrete least squares solution on the data set X to  $f \in C^{(m,n)}([a,b]x[c,d])$  for  $1 \leq m, n \leq 4$ . If  $\overline{\beta}(X, \overline{\pi}_x, \overline{\pi}_y) < 1/9$ , then (5.21)  $||f-s|| \leq \{\mathcal{E}_{m,0}||f^{(m,0)}||\overline{h}_x + \mathcal{E}_{n,0}||f^{(0,n)}||\overline{h}_y$  $+ \mathcal{E}_{m,0} \mathcal{E}_{n,0}||f^{(m,n)}||h_x^m h_y^n + \mathcal{E}_{m,0} \mathcal{E}_{n,0}||f^{(m,n)}||\overline{h}_x^m h_y^n$ 

$$\begin{split} &+ \mathcal{E}_{m,0} \mathcal{E}_{n,0} || f^{(m,n)} || h_{x}^{m} \overline{h}_{y}^{n} \} (1 + 1 / (1 - 9 \overline{\beta} (X, \overline{\pi}_{x}, \overline{\pi}_{y}))) \\ &+ || R ||_{\infty} / (1 - 9 \overline{\beta} (X, \overline{\pi}_{x}, \overline{\pi}_{y})) . \end{split}$$
If  $\frac{\Delta}{\beta} (X, \overline{\pi}_{x}, \overline{\pi}_{y}) < \sqrt{2}/3$ , then
$$(5.22) \qquad || f - s || \leq \{ \mathcal{E}_{m,0} || f^{(m,0)} || \overline{h}_{x} + \mathcal{E}_{n,0} || f^{(0,n)} || \overline{h}_{y} \\ &+ \mathcal{E}_{m,0} \mathcal{E}_{n,0} || f^{(m,n)} || h_{x}^{m} h_{y}^{n} + \mathcal{E}_{m,0} \mathcal{E}_{n,0} || f^{(m,n)} || \overline{h}_{x}^{m} h_{y}^{n} \\ &+ \mathcal{E}_{m,0} \mathcal{E}_{n,0} || f^{(m,n)} || h_{x}^{m} \overline{h}_{y}^{n} \} (1 + 1 / (1 - (3 \overline{\beta} (X, \overline{\pi}_{x}, \overline{\pi}_{y}, \overline{\pi}_{y}))^{2} / 2)) \\ &+ || R ||_{\infty} / (1 - (3 \overline{\beta} (X, \overline{\pi}_{x}, \overline{\pi}_{y}, \overline{\pi}_{y}))^{2} / 2) , \end{split}$$

where  $\mathcal{E}_{m,0}$  and  $\mathcal{E}_{n,0}$  are given in Table 1 of Chapter 2,  $h_x$ ,  $\overline{h}_x$ ,  $h_y$  and  $\overline{h}_y$  are the mesh sizes of  $\pi_x$ ,  $\overline{\pi}_x$ ,  $\pi_y$  and  $\overline{\pi}_y$ , respectively.

Proof: Identical to Theorem 5.1, hence omitted.

### Section 5.3. Hermite Blended Piecewise Polynomials.

Let  $\pi_x:a = x_1 = x_2 < x_3 = x_4 < \cdots < x_{2M-1} = x_{2M} = b$  and  $\pi_y:c = y_1 = y_2 < y_3 = y_4 < \cdots < y_{2N-1} = y_{2N} = d$  be univariate meshes and  $V(\pi_x, 2M, 1)$  and  $V(\pi_y, 2N, 1)$  be the interpolation spaces of cubic Hermite splines of Example 1.3 of Chapter 2. Let  $\{\phi_i\}_{i=1}^{2M}$  and  $\{\psi_j\}_{j=1}^{2N}$  be the cardinal bases (see [29, pp. 25-27] for explicit representation),  $\alpha$  and  $\beta$  be the interpolation functions

180

for the interpolation spaces  $V(\pi_x, 2M, 1)$ , and  $V(\pi_y, 2N, 1)$ , respectively, given in Example 1.3 of Chapter 2.

If 
$$g \in C^{(4)}[a, b]$$
, then

(5.23) 
$$||g - \sum_{i=1}^{2M} g^{(\alpha(i))}(x_i) \phi_i|| \leq (1/384) ||g^{(4)}|| h_x^4$$

with a corresponding error bound for the interpolation space  $V(\pi_v, 2N, 1)$ , see Carlson and Hall [5].

Refine the mesh  $\pi_{\mathbf{x}}$  by adding k-2 additional points between each of the knots  $\mathbf{x}_{2i}$  and  $\mathbf{x}_{2i+1}$ , for  $1 \leq i \leq M-1$ , to obtain  $\overline{\pi}_{\mathbf{x}}:\mathbf{a} = \overline{\mathbf{x}}_1 = \overline{\mathbf{x}}_2 < \overline{\mathbf{x}}_3 < \cdots < \overline{\mathbf{x}}_k < \overline{\mathbf{x}}_{k+1} = \overline{\mathbf{x}}_{k+2} < \overline{\mathbf{x}}_{k+3} < \cdots < \overline{\mathbf{x}}_{(M-1)k}$  $< \overline{\mathbf{x}}_{(M-1)k+1} = \overline{\mathbf{x}}_{(M-1)k+2} = b$ , where we have the correspondence for  $1 \leq i \leq M$ 

(5.24) 
$$x_{2i-1} = \overline{x}_{(i-1)k+1}$$
 and  $x_{2i} = \overline{x}_{(i-1)k+2}$ 

Define the interpolation function

(5.25) 
$$\overline{\alpha}(i) = \begin{cases} 1 & \text{if } (i-2) \mod k = 0 \\ 0 & \text{otherwise} \end{cases}$$

 $\overline{V(\pi_x, (M-1)k+2, 1)}$  is the interpolation space of continuously differentiable functions such that on the interval  $[x_{ik+1}, x_{(i+1)k+2}]$  each element is a polynomial of degree k+1 for  $0 \le i \le M-2$ . The cardinal basis for this space is represented by  $\{\overline{\phi_i}\}_{i=1}^{(M-1)k+2}$ .

By the manner in which the new points have been added we have that  $V(\pi_x, 2M, 1)$  is subordinate to  $\overline{V}(\overline{\pi}_x, (M-1)k+2, 1)$  and  $V(\pi_x, 2M, 1) \subseteq \overline{V}(\overline{\pi}_x, (M-1)k+2, 1).$  Correspondingly, refine the mesh  $\pi_y$  by adding  $\ell$ -2 points between each of the points  $y_{2j}$  and  $y_{2j+1}$  of  $\pi_y$  for  $1 \le j \le N-1$ to obtain the mesh  $\overline{\pi_y}: c=\overline{y_1} = \overline{y_2} < \overline{y_3} < \cdots < \overline{y_\ell} < \overline{y_{\ell+1}} = \overline{y_{\ell+2}} < \cdots$  $< \overline{y_{(N-1)\ell}} < \overline{y_{(N-1)\ell+1}} = \overline{y_{(N-1)\ell+2}} = d$ . Construct the interpolation space of piecewise polynomials of degree  $\ell + 1$ ,  $\overline{V}(\overline{\pi_y}, (N-1)\ell+2, 1)$ , with interpolation function  $\overline{\beta}$  and cardinal basis  $\{\overline{\psi}\}_{j=1}^{(N-1)\ell+2}$  in an analogous manner. Then  $V(\pi_y, 2N, 1)$  is subordinate to  $\overline{V}(\overline{\pi_y}, (N-1)\ell+2, 1)$  and  $V(\pi_y, 2N, 1) \subseteq \overline{V}(\overline{\pi_y}, (N-1)\ell+2, 1)$ .

Construct the discretized blending function space  $DBF_2$ ; then, from Theorem 2.3 of Chapter 2,

(5.26) 
$$Dim(DBF_2) = 2M[(N-1)\ell+2] + [2N (M-1)k+2] - 4MN,$$

and each of the sets  $T_1$  through  $T_4$  forms a basis for DBF<sub>2</sub>. We will examine the basis  $T_4$  in detail. If  $IM = \{i | i = (s-1)k+1 \text{ or} i = (s-1)k+2, 1 \leq s \leq M\}$ ,  $\overline{IM} = \{i | 1 \leq i \leq (M-1)k+2 \text{ and } i \notin IM\}$ ,  $JN = \{j | j = (s-1)\ell+1 \text{ or } j = (s-1)\ell+2, 1 \leq s \leq N\}$  and  $\overline{JN} = \{j | 1 \leq j \leq (N-1)\ell+2 \text{ and } j \notin JN\}$ , then  $T_4 = S_1 \cup S_2 \cup S_6$ , where  $S_1 = \{\overline{\varphi_i}, \psi_j\}_{i \in \overline{IM}}, 1 \leq j \leq 2N'$ ,  $S_2 = \{\varphi_i, \overline{\psi_j}\}_1 \leq i \leq 2M$ ,  $j \in \overline{JN}$  and  $S_6 = \{\varphi_i, \psi_j\}_{1 \leq i \leq 2M}, 1 \leq j \leq 2N'$ . For  $\overline{\varphi_i}\psi_j \in S_1$ , let  $\hat{i} = (2((i-1)-(i-1) \mod k)/k)+2$  and  $\hat{j} = j - (j+1) \mod 2$ , then  $\overline{\varphi_i}\psi_j$  has support on the rectangle  $[x_i, x_{i+1}]$   $x[y_{\hat{j}-1}, y_{\hat{j}+2}]$ , where we restrict ourselves to the domain [a, b]x[c, d](recall that  $y_j = y_{j+1}^{*}$ ). Correspondingly for  $\varphi_i \overline{\psi_j} \in S_2$ , if  $\hat{j} = (2((j-1)-(j-1) \mod \ell)/\ell) + 2 \text{ and } \hat{i} = i-(i+1) \mod 2, \text{ then } \phi_i \overline{\psi_j}$  has support on  $[x_{\hat{i}-1}, x_{\hat{i}+2}] \times [y_j^{\hat{i}}, y_{\hat{j}+1}], \text{ where } x_{\hat{i}}^{\hat{i}} = x_{\hat{i}+1}^{\hat{i}} \text{ and we}$  restrict ourselves to  $[a, b] \times [c, d]$ . Finally, for  $\phi_i \psi_i \in S_6$ , if  $\hat{i} = i-(i+1) \mod 2$  and  $\hat{j} = j-(j+1) \mod 2$ , then  $\phi_i \psi_j$  has support on  $[x_{\hat{i}-1}, x_{\hat{i}+2}] \times [y_{\hat{j}-1}, y_{\hat{j}+2}]$ . Paraphrasing the above,  $\overline{\phi_i} \psi_j$  and  $\phi_i \overline{\psi_j}$  have support on at most two adjacent rectangles for the meshes  $\pi_x$  and  $\pi_y$ , while the support of  $\phi_i \psi_j$  is over at most four adjacent rectangles.

The cubic Hermite bases  $\{\phi_i\}_{i=1}^{2M}$  and  $\{\psi_j\}_{j=1}^{2N}$  are easy to calculate and can be stored in the computer in polynomial form. Actually, only two basis functions need be stored, as the others can be derived from these two (for explicit representation see [29, pp. 25-27]). Correspondingly, the basis of piecewise polynomials  $\{\overline{\phi}_i\}_{i=1}^{(M-1)k+2}$  of degree k+1, and the basis of piecewise polynomials  $\{\overline{\psi}_j\}_{j=1}^{(N-1)\ell+2}$  of degree  $\ell$ +1 are easy to calculate and can be stored in the computer in polynomial form. Again, for similar reasons, only k+2 or  $\ell$ +2 polynomials need be stored as the others can be easily generated from these.

Note that the above basis circumvents many of the difficulties associated with the implementation of discretized blending function spaces, hence the main reason for its introduction.

It is a simple matter to calculate the interpolation accuracy of the interpolation spaces  $\overline{V}(\overline{\pi}_{x}, (M-1)k+2, 1)$  (correspondingly, for

$$\overline{V(\overline{\pi}_{y}, (N-1)\ell+2, 1)}. \quad \text{For } \mathbf{x} \in [\overline{\mathbf{x}_{sk+1}}, \overline{\mathbf{x}_{(s+1)k+2}}] \text{ and } g \in C^{(k+2)}[a, b]$$

$$(5.27) \quad g(\mathbf{x}) = \sum_{i=1}^{(M-1)k+2} g^{(\overline{\alpha}(i))}(\overline{\mathbf{x}_{i}}) \quad \overline{\phi_{i}}(\mathbf{x})$$

$$i=1$$

$$= \prod_{i=1}^{k+2} (x - \overline{x}_{sk+1}) g^{(k+2)} \xi(x) / (k+2)!,$$

where  $\xi(x) \in \left[\overline{x}_{sk+1}, \overline{x}_{(s+1)k+2}\right]$ , see [33, pp. 1-5].

Hence

(5.28) 
$$||g - \sum_{i=1}^{(M-1)k+2} g^{(\overline{\alpha}(i))}(\overline{x}_i) \overline{\phi}_i||$$

$$\leq \prod_{i=1}^{k+2} |\mathbf{x} - \overline{\mathbf{x}}_{sk+i}| ||g^{(k+2)}||/(k+2)!$$

$$\leq ||g^{(k+2)}||h_{x}^{k+2}/(k+2)|$$

Remark: Note that 
$$\max_{0 \le s \le M-2} (\overline{x}_{(s+1)k+2} - \overline{x}_{sk+1})$$

 $= \max_{\substack{0 \leq s \leq M-2}} (x_{2(s+1)+1} - x_{2s+1}) = h, \text{ the mesh size of } \pi_x.$ 

<u>Remark</u>: If k is fixed and the spacing of the points  $x_{sk+1}$  for  $3 \le i \le k$  is also fixed, then a more refined estimate can be given for this term. For example, if k = 6, and  $\overline{x}_{sk+i}$  are equally spaced for  $3 \le i \le k$ , then  $\prod_{i=1}^{8} |x - \overline{x}_{6s+i}| < h_x^8 24/5^6$ . Combining the above and using Theorem 1.7 of Chapter 2, we have the following error estimate for  $f \in C^{(k+2, l+2)}([a, b]x[c, d])$  (note that k and  $l \ge 2$ ),

$$(5.29) ||f - \overline{P_{x}} \oplus \overline{P_{y}}[f]|| \leq (1/(k+2)!) ||f^{(k+2,0)}||h_{x}^{k+2} + (1/(\ell+2)!)||f^{(0,\ell+2)}||h_{y}^{\ell+2} + (1/384)^{2}||f^{(4,4)}||h_{x}^{4}h_{y}^{4} + (1/(384\cdot(k+2)!))||f^{(k+2,4)}||h_{x}^{k+2}h_{y}^{4} + (1/(384(\ell+2)!)) + (1/(4\ell+2)!) + (1/($$

The limiting accuracy of the above interpolation scheme is  $h_x^4 h_y^4$ , thus we choose  $\ell$  and k sufficiently large to insure that  $h_x^{k+2}$ ,  $h_y^{\ell+2} \leq h_x^4 h_y^4$ . For example, if  $h = \max(h_x, h_y)$  then this would imply that in terms of h, the choice of  $\ell = k = 6$  would suffice.

<u>Remark</u>: If  $h = max (h_x, h_y)$ , and we assume sufficient continuity of f, then both of the discretized blending function spaces, cubically blended cubic splines = DBF<sub>1</sub>, and cubic Hermite blended piecewise polynomials = DBF<sub>2</sub>, give eighth order approximations to f. For large M and N, examination of (5.20) and (5.26) shows that Dim (DBF<sub>1</sub>) is much larger than Dim (DBF<sub>2</sub>). Thus, there will be a corresponding savings of computer memory needed to store the least squares matrix using DBF<sub>2</sub>. However, in order to obtain this savings, the continuity required to implement DBF<sub>2</sub> must be larger than DBF<sub>1</sub>. Let  $X = \{(\hat{x}_i, \hat{y}_i)\}_{i=1}^{\hat{M}} \subseteq [a, b]x[c, d]$  be a set of  $\hat{M}$  unstructured data points. Let  $f \in C^{(k+2, \ell+2)}([a, b]x[c, d])$ , then a discrete least squares solution  $s \in DBF_2$  is a function which minimizes the Euclidean norm of the residual vector  $R \in \mathbb{R}^{\hat{M}}$  defined in (5.2).

Theorem 5.3. Let  $s \in DBF_2$  be a discrete least squares solution on the data set X to  $f \in C^{(k+2, \ell+2)}([a, b] \times [c, d])$ . If  $\overline{\beta}(X, \pi_x, \pi_y) < 1/t^2$ , then

$$(5.30) ||f-s|| \leq \{(1/(k+2)!)||f^{(k+2,0)}||h_{x}^{k+2} + (1/(k+2)!)||f^{(0,\ell+2)}||h_{y}^{\ell+2} + (1/(k+2)!)||f^{(4,\ell+2)}||h_{x}^{\ell+2} + (1/(k+2)!)||f^{(k+2,\ell+2)}||h_{x}^{k+2} h_{y}^{4} + (1/(k+2)!))||f^{(k+2,\ell+2)}||h_{x}^{k+2} h_{y}^{4} + (1/(k+2)!))||f^{(4,\ell+2)}||h_{x}^{k} h_{y}^{\ell+2}\} \{1-1/(1-t^{2}\overline{\beta}(X, \pi_{x}, \pi_{y}))\} + ||R||_{\infty}/(1-t^{2}\overline{\beta}(X, \pi_{x}, \pi_{y})).$$

If 
$$\hat{\beta}(X, \pi_{x}, \pi_{y}) < \sqrt{2}/t$$
, then  
(5.31)  $||f-s|| \leq \{ (1/(k+2)!) ||f^{(k+2,0)}||h_{x}^{k+2} + (1/(384!)^{2})||f^{(4,4)}||h_{x}^{4}h_{y}^{4} + (1/(384!(k+2)!)) ||f^{(k+2,4)}||h_{x}^{k+2}h_{y}^{4} + (1/(384!(k+2)!)) ||f^{(k+2,4)}||h_{x}^{k+2}h_{y}^{4} + (1/(384!(\ell+2)!)) ||f^{(4,\ell+2)}||h_{x}^{4}h_{y}^{\ell+2} \} \{ 1-1/(1-(t\hat{\beta}(X, \pi_{x}, \pi_{y}))^{2}/2) \} + ||R||_{\infty} / (1-(t\hat{\beta}(X, \pi_{x}, \pi_{y}))^{2}/2) ,$   
where k,  $\ell \ge 2$  and t = max {k+1,  $\ell+1$ }.

Proof: Similar to Theorem 5.1, hence omitted.

# Section 5.4. Linear Blending

As a final example, we shall consider the discretized blending function space of linearly blended piecewise cubic polynomials.

Let  $\pi_x: a = x_1 < x_2 < \cdots < x_M = b$  and  $\pi_y: c = y_1 < y_2 < \cdots < y_N = d$  be univariate meshes. The interpolation spaces of piecewise linear continuous functions  $V(\pi_x, M, 0)$  and  $V(\pi_y, N, 0)$  were defined in Section 4, where the basis functions  $\{\phi_i\}_{i=1}^M$  and  $\{\psi_j\}_{j=1}^N$  were given by (4.1) and the error estimate by (4.2).

Refine the meshes  $\pi_x$  and  $\pi_y$  by adding two knots between each of the knots of  $\pi_x$  and  $\pi_y$  to obtain  $\overline{\pi_x}$ :a =  $\overline{x_1} < \overline{x_2} < \overline{x_3} < \cdots < \overline{x_{3M-2}} = b$  and  $\overline{\pi_y}$ :c =  $\overline{y_1} < \overline{y_2} < \overline{y_3} < \cdots < \overline{y_{3N-2}} = d$ , respectively, where

(5.32) 
$$x_i = \overline{x}_{3(i-1)+1}$$
 for  $1 \le i \le M$ , and  $y_j = y_{3(j-1)+1}$  for  $1 \le j \le N$ .

Define the interpolation space of piecewise continuous cubic polynomials  $\overline{V}(\overline{\pi}_x, 3M-2, 0)$  (respectively,  $\overline{V}(\overline{\pi}_y, 3N-2, 0)$ ) with interpolation function  $\overline{\alpha}(i) = 0$  for  $1 \le i \le 3M-2$  ( $\overline{\beta}(j) = 0$  for  $1 \le j \le 3N-2$ ) and cardinal basis of piecewise continuous cubic polynomials  $\{\overline{\varphi}_i\}_{i=1}^{3M-2}$  $(\{\overline{\psi}_j\}_{j=1}^{3N-2})$ . For a fixed i, such that  $1 \le i \le 3M-2$ , let  $s = ((i-1)-(i-1) \mod 3)/3$  and k = i - 3s, if  $s \le M-2$ , then  $\overline{\varphi}_i$  is the cubic Lagrange interpolation polynomial on the four points  $\{\overline{x}_{3s+t}\}_{t=1}^{4}$ such that  $\overline{\phi}_i(\overline{x}_{3s+t}) = \delta_{kt}$ . If k = 2, 3, then  $\overline{\phi}_i$  is identically zero off the interval  $[\overline{x}_{3s+1}, \overline{x}_{3(s+1)+1}] = [x_{s+1}, x_{s+2}]$ . If k = 1 and i > 1, or s = M-1, then  $\overline{\phi}_i$  is also the cubic Lagrange interpolation polynomial on the four points  $\{\overline{x}_{3(s-1)+t}\}_{t=1}^{4}$  such that  $\overline{\phi}_i(\overline{x}_{3(s-1)+t})$  $= \delta_{t,4}$  and identically zero off the interval  $[x_s, x_{s+2}]$ . Thus, the support of  $\{\phi_i\}_{i=1}^{M}$  and  $\{\overline{\phi}\}_{i=1}^{3M-2}$   $(\{\psi_j\}_{j=1}^{N}$  and  $\{\overline{\psi}_j\}_{j=1}^{3N-2}$  consists of at most two adjacent intervals of the mesh  $\pi_x$   $(\pi_y)$ .

Note that  $V(\pi_x, M, 0)$  is subordinate to  $\overline{V}(\overline{\pi}_x, 3M-2, 0)$  and  $V(\pi_x, M, 0) \subseteq \overline{V}(\overline{\pi}_x, 3M-2, 0)$  with corresponding relations holding for  $V(\pi_y, N, 0)$  and  $\overline{V}(\overline{\pi}_y, 3N-2, 0)$ . Theorem 2.3 of Chapter 2 gives the dimension of the discretized blending function space of linearly blended piecewise cubic polynomials = DBF<sub>3</sub> as

(5.33) 
$$\text{Dim}(\text{DBF}_3) = 5 \text{ MN} - 2(\text{M}+\text{N})$$
,

where  $T_1$  through  $T_4$  each forms a basis for DBF<sub>3</sub>. Note that the cardinal basis elements  $\phi_i$ ,  $\overline{\phi_i}$ ,  $\psi_j$  and  $\overline{\psi_j}$  are particularly simple, easy to store and compute, with support over at most two adjacent intervals of the meshes  $\pi_x$  or  $\pi_y$ . Hence, the basis elements for DBF<sub>3</sub> will share these desirable properties of being very simple, easy to store and compute, with support over at most four adjacent rectangles.

On the interval  $[x_i, x_{i+1}]$ , the interpolation error of  $\overline{V}(\overline{\pi}_x, 3M-2, 0)$  is just the difference of the cubic Lagrange

interpolation polynomial and a function  $g \in C^{(4)}[a, b]$ . Thus, on this interval, we have from [24, p. 249]

(5.34) 
$$|g(x) - \sum_{k=1}^{3M-2} g(\overline{x}_{k}) \overline{\phi}_{k}(x)| \leq \prod_{s=1}^{4} |x - \overline{x}_{3(i-1)+s}| ||g^{(4)}||/4!$$
  
 $\leq K h_{x}^{4} ||g^{(4)}||/4!$ ,

where K is a constant less than one.

<u>Remark</u>: If structure is given to the four data points  $\{x_{3(i-1)+s}\}_{s=1}^{4}$ , then various estimates for K can be given. For example, if the points are equally spaced in each interval  $[x_i, x_{i+1}]$ , then direct calculation shows that K = 1/81. If the points are the zeros of the cubic Chebyshev polynomial (see [24, pp. 228] for definition, explanation and bound), then K = 1/128, which is the best possible for any distribution of the four points in each interval.

From Theorem 1.7 of Chapter 2, for  $f \in C^{(4,4)}([a,b]x[c,d])$ (5.35)  $||f - \overline{P_x} \oplus \overline{P_y}[f]|| \leq (K/24) ||f^{(4,0)}|| h_x^4$  $+ (K/24) ||f^{(0,4)}||h_y^4 + (1/64)||f^{(2,2)}|| h_x^2 h_y^2$  $+ (K/192) ||f^{(4,2)}||h_x^4 h_y^2 + (K/192)||f^{(2,4)}||h_x^2 h_y^4$ ,

where K < 1 is some constant, which is defined above.

Let  $f \in C^{(4,4)}([a,b]x[c,d])$ , then a discrete least squares solution  $s \in DBF_3$  is a function which minimizes the Euclidean norm of the residual vector  $R \in \mathbb{R}^{\hat{M}}$  defined in (5.2) on the data set X.

Theorem 5.4. Let  $s \in DBF_3$  be a discrete least squares solution to  $f \in C^{(4, 4)}([a, b]x[c, d])$  on the data set X. If  $\overline{\beta}(X, \pi_x, \pi_y) < 1/9$ , then

$$(5.36) ||f-s|| \leq \{ (K/24) ||f^{(4,0)}||h_x^4 + (K/24)||f^{(0,4)}||h_y^4 + (1/64)||f^{(2,2)}||h_x^2 h_y^2 + (K/192)||f^{(4,2)}||h_x^4 h_y^2 + (K/192)||f^{(2,4)}||h_x^2 h_y^4 \} \{ 1-1/(1-9\overline{\beta}(X, \pi_x, \pi_y)) \}$$

+  $||R||_{\infty}/(1-9\overline{\beta}(X, \pi_{y}, \pi_{y}))$ .

If 
$$\hat{\beta}(X, \pi_{x}, \pi_{y}) < \sqrt{2}/3$$
, then  
(5.37)  $||f-s|| \leq \{(K/24)||f^{(4,0)}||h_{x}^{4} + (K/24)||f^{(0,4)}||h_{y}^{4} + (1/64)||f^{(2,2)}||h_{x}^{2}h_{y}^{2} + (K/192)||f^{(4,2)}||h_{x}^{4}h_{y}^{2} + (K/192)||f^{(2,4)}||h_{x}^{2}h_{y}^{4} + (K/192)||f^{(2,4)}||h_{x}^{2}h_{y}^{4}\}\{1-1/(1-(3\hat{\beta}(X, \pi_{x}, \pi_{y}))^{2}/2)\} + ||R||_{\infty}/(1-(3\hat{\beta}(X, \pi_{x}, \pi_{y}))^{2}/2),$ 

where K < 1 is some constant, depending upon the meshes  $\overline{\pi}_x$  and  $\overline{\pi}_y$ , and  $h_x$  and  $h_y$  are the mesh sizes of  $\pi_x$  and  $\pi_y$ .

**Proof:** Similar to the proof of Theorem 5.1.

A procedure is now presented which minimizes the computer storage needed to compute a discrete least squares solution. Note

that for the discretized blending function spaces considered, even for relatively small M and N, a large matrix must be stored in the computer, with a correspondingly large number of computational operations needed to obtain the numerical solution. The procedure given here is to solve a separate discrete least squares problem on each rectangle  $[x_i, x_{i+1}] \times [y_i, y_{i+1}]$  for  $1 \le i \le M-1$  and  $l \leqslant j \leqslant$  N-1. Restricted to this rectangle, the discretized blending function space of linearly blended piecewise continuous cubic polynomials will be denoted by DBF<sub>ii</sub>, where it is clear from the previous notation what the meshes, interpolation spaces, and their cardinal basis functions will be. From (5.33), it follows that Dim (DBF<sub>1</sub>) = 12, with a basis given by  $T_1$  through  $T_4$  of Theorem 2.3 of Chapter 2. If  $X_{ij} = X \cap ([x_i, x_{i+1}] \times [y_i, y_{i+1}])$ for  $1 \leq i \leq M-1$  and  $1 \leq j \leq N-1$  represents the  $\bigwedge_{ij}^{\wedge}$  data points which are in each of the rectangular regions (note that a data point could be in more than one  $X_{ij}$  if it lies on a mesh line of the bivariate mesh  $\pi_x \bigoplus \pi_v$ ), then the matrix problem to be stored at any one time has only  $12\dot{M}_{ii}$  elements. This is significantly less storage than that required by any of the previous methods. If  $\dot{M}_{ii}$ is not too large, then most computers, even those having a very small memory capacity, can store the full least squares matrix in central memory.

Let 
$$f \in C^{(4, 4)}([x_i, x_{i+1}] \times [y_j, y_{j+1}])$$
 and let  $s_{ij} \in DBF_{ij}$  be  
a discrete least squares solution which minimizes the residual  
vector  $R_{ij} \in \mathbb{R}^{\hat{M}_{ij}}$ , where each of the  $\hat{M}_{ij}$  elements of  $R_{ij}$  is the  
difference of f and  $s_{ij}$  evaluated at one of the  $\hat{M}_{ij}$  data points of  
 $X_{ij}$ . Theorem 5.4 gives an error estimate for  $||f - s_{ij}||$  on  
 $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$ . The set  $\{s_{ij}\}_{1 \leq i \leq M-1, 1 \leq j \leq N-1}$  forms a  
''patch network'' of discrete least squares functions over the full  
domain  $[a, b] \times [c, d]$ . However, this ''patch network'' is not neces-  
sarily continuous across the mesh lines of  $\pi_x \bigoplus \pi_y$ . In order to  
remedy this situation and obtain an approximation to f which has  
global continuity on  $[a, b] \times [c, d]$ , the following scheme is intro-  
duced which produces an approximation  $\tilde{s} \in DBF_3$  to f.

If  $T_4$  of Theorem 2.3 of Chapter 2 is chosen as a basis for DBF<sub>ij</sub> on  $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$ , then

$$(5.38) \quad \mathbf{T}_{4} = \{ \phi_{\mathbf{s}}^{1\mathbf{j}} \psi_{\mathbf{t}}^{1\mathbf{j}} \}_{2 \leq \mathbf{s} \leq 3, 1 \leq \mathbf{t} \leq 2} \bigcup \{ \phi_{\mathbf{s}}^{1\mathbf{j}} \overline{\psi}_{\mathbf{t}}^{1\mathbf{j}} \}_{\mathbf{t} \leq \mathbf{s} \leq 2, 2 \leq \mathbf{t} \leq 3}$$

$$\bigcup \{\phi_{\mathbf{s}}^{\mathbf{ij}} \psi_{\mathbf{t}}^{\mathbf{ij}}\}_{\mathbf{l}\leqslant \mathbf{s}, \mathbf{t}\leqslant \mathbf{2}},$$

and

•

$$(5.39) \quad \mathbf{s}_{ij} = \begin{array}{c} 3 & 2 \\ \Sigma & \Sigma \\ \mathbf{s}=2 \ t=1 \end{array} \mathbf{s}_{t}^{ij} \quad \mathbf{\phi}_{s}^{ij} \quad \mathbf{\psi}_{t}^{ij} + \begin{array}{c} 2 & 3 \\ \Sigma & \Sigma \\ \mathbf{s}=1 \ t=2 \end{array} \mathbf{s}_{s}^{ij} \quad \mathbf{\phi}_{s}^{ij} \quad \mathbf{\psi}_{t}^{ij} \\ \mathbf{s}=1 \ t=2 \end{array} \mathbf{s}_{s}^{ij} \quad \mathbf{\phi}_{s}^{ij} \quad \mathbf{\psi}_{t}^{ij} \\ + \begin{array}{c} 2 & 2 \\ \Sigma & \Sigma \\ \mathbf{s}=1 \ t=1 \end{array} \mathbf{c}_{st}^{ij} \quad \mathbf{\phi}_{s}^{ij} \quad \mathbf{\psi}_{t}^{ij} \quad \mathbf{s}_{s}^{ij} \quad \mathbf{\phi}_{s}^{ij} \quad \mathbf{\psi}_{t}^{ij} \end{array}$$

For each fixed i and j, let  $0 \le u, v \le l$ , and define the average of the "a", "b" and "c" coefficients to be

(5.40) 
$$\tilde{a}_{2+u, 1+v}^{ij} = \sum_{t=0}^{l} a_{2+u, 1+t}^{i, j+v-t} / DA(i, j, u, v),$$

(5.41) 
$$\tilde{b}_{1+u, 2+v}^{ij} = \frac{1}{\sum_{s=0}^{\infty} b_{1+s, 2+v}^{i+u-s, j}} / DB(i, j, u, v),$$

and

(5.42) 
$$\widetilde{c}_{1+u,1+v}^{ij} = \sum_{s=0}^{l} \sum_{t=0}^{l} c_{1+s,1+t}^{i+u-s,j+v-t} / DC(i, j, u, v)$$
.

where the indices are restricted to insure that  $1 \le i, i+u-s \le M-1$  and  $1 \le j, j+v-t \le N-1$ . The DA(i, j, u, v) = 1 or 2, DB(i, j, u, v) = 1 or 2 and DC(i, j, u, v) = 1, 2 or 4 are just the total number of terms which have been summed in (5.40), (5.41) and (5.42), respectively. For example, DC(1, 1, 0, 0) = 1, and if M, N > 2, then DC(1, 1, 1, 0) = 2 and DC(1, 1, 1, 1) = 4, etc.

Define

$$(5.43) \quad \widetilde{s}_{ij} = \begin{array}{c} 3 & 2 \\ \Sigma & \Sigma & s_{st} \\ s=2 & t=1 \end{array} \qquad \widetilde{s}_{st} \begin{array}{c} \widetilde{\phi}_{ij} & \widetilde{\phi}_{ij} \\ s_{st} & \psi_{t}^{ij} \\ s=1 \\ t=2 \end{array} \qquad \widetilde{s}_{st} \begin{array}{c} 2 & 3 \\ s_{st} & \widetilde{\phi}_{st}^{ij} \\ s=1 \\ t=2 \end{array} \qquad \widetilde{s}_{st} \begin{array}{c} \widetilde{\phi}_{s}^{ij} \\ \widetilde{\phi}_{s}^{ij} \\ \widetilde{\phi}_{t}^{ij} \\ \widetilde{\phi}_{s}^{ij} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{st} \\ \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} 2 & 3 \\ \widetilde{t}_{s} \\ s=1 \\ t=2 \end{array} \qquad \widetilde{t}_{st} \begin{array}{c} \widetilde{t}_{s}^{ij} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s}^{ij} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s}^{ij} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \\ \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \begin{array}{c} \widetilde{t}_{s} \end{array} \qquad \widetilde{t}_{s} \end{array}$$

to be that element in  $DBF_{ij}$  which has as its coefficients the average "a's", "b's" and "c's". Correspondingly, define  $\tilde{s} \in DBF_3$  by

$$(5.44) \quad \widetilde{s} = \frac{M-1}{\Sigma} \frac{3}{\Sigma} \frac{N}{\Sigma} \frac{3}{s_{1}} \frac{1}{\phi_{3(i-1)+s}} \psi_{j}$$

$$+ \frac{M}{\Sigma} \frac{N-1}{\Sigma} \frac{3}{b_{1t}} \frac{1}{\phi_{i}} \frac{1}{\psi_{3(j-1)+t}} + \frac{M}{\Sigma} \frac{N}{\Sigma} \frac{N}{c_{1,1}} \phi_{i} \psi_{j},$$

$$= \frac{M}{i=1} \frac{1}{j=1} \frac{1}{t=2} \frac{1}{t=2} \frac{1}{t=2} \frac{1}{t=1} \frac{1}{t=1}$$

where we define  $\tilde{a}_{s,1}^{i,N} = \tilde{a}_{s,2}^{i,N-1}$ ,  $\tilde{b}_{l,t}^{M,j} = \tilde{b}_{2,t}^{M-1,j}$ ,

 $\tilde{c}_{1,1}^{i,N} = \tilde{c}_{1,2}^{i,N-1}$ ,  $\tilde{c}_{1,1}^{M,j} = \tilde{c}_{2,1}^{M-1,j}$  and  $\tilde{c}_{1,1}^{M,N} = \tilde{c}_{2,2}^{M-1,N-1}$ .

Let the real number  $\mathcal{E} > 0$  be such that

(5.45) 
$$\max_{\substack{1 \leq i \leq M-1 \\ l \leq j \leq N-1}} ||f^{-s}_{ij}|| [x_i, x_{i+1}] x [y_j, y_{j+1}] \leq \mathcal{E} ,$$

then we have the following theorem.

<u>Theorem 5.5</u>. Let  $\tilde{s} \in DBF_3$ , defined in (5.44), be constructed as above, then

(5.46) 
$$||f-\tilde{s}|| \leq \{5/2 + 2\{\max_{\substack{1 \leq i \leq M-1 \\ 1 \leq i \leq M-1 \\ s = 2}} 3 \\ ||\bar{\phi}_{3(i-1)+s}|| \\ + \max_{\substack{1 \leq j \leq N-1 \\ 1 \leq j \leq N-1 \\ t = 2}} 3 \\ ||\bar{\psi}_{3(j-1)+t}|| \} \mathcal{E},$$

where  $\mathcal{E}$  satisfies (5.45).

<u>Proof</u>: It will be shown that for  $1 \le i \le M-1$ ,  $1 \le j \le N-1$ , and  $1 \le s, t \le 2$  that

(5.47) 
$$|\tilde{c}_{st}^{ij} - c_{st}^{ij}| \leq (3/2) \mathcal{E}$$
.

Because of the cardinality conditions, (5.39) and (5.45), we have for  $l \leq i \leq M-1$ ,  $l \leq j \leq N-1$ ,  $l \leq s, t \leq 2$ 

$$(5.48) \quad |c_{st}^{ij} - f(x_{i+s-1}, y_{j+t-1})| \leq \mathcal{E}.$$

Using (5.48), (5.42), the triangle inequality, and observing the internal cancellation of one of the terms yields (5.47) (notice that if DC(i, j, s, t) = k for k = 1, 2 then  $|\tilde{c}_{st}^{ij} - c_{st}^{ij}| \leq (k-1) \mathcal{E}$ ). Also, for  $2 \leq s \leq 3$  and  $1 \leq t \leq 2$  we have

(5.49) 
$$|\tilde{a}_{st}^{ij} - a_{st}^{ij}| \leq 2\mathcal{E} \text{ and } |\tilde{b}_{ts}^{ij} - b_{ts}^{ij}| \leq 2\mathcal{E}.$$

The proof is given only for the first inequality, as the other is nearly identical. From the cardinality conditions, (5.39) and (5.45), we have

$$(5.50) \quad \left|a_{st}^{ij} + \sum_{k=1}^{2} c_{kt}^{ij} \phi_{k}^{ij} (\overline{x}_{3(i-1)+s}) - f(\overline{x}_{3(i-1)+s}, y_{j+t-1})\right| \leq \mathcal{E},$$

for all  $1 \leq i \leq M-1$ ,  $1 \leq j \leq N-1$ ,  $2 \leq s \leq 3$  and  $1 \leq t \leq 2$ .

Consider the case where DA(i, j, u, v) = 2 (as the case where DA(i, j, u, v) = 1 is trivial), and without loss of generality let t = 1. Then from (5.40), (5.49) and the triangle inequality we have the following after some algebra and cancellation

$$(5.51) \quad |\tilde{a}_{s1}^{ij} - a_{s1}^{ij}| = \frac{1}{2} |a_{s,2}^{i,j-1} - a_{s,1}^{i,j}|$$
$$\leq \frac{1}{2} \left[ \mathcal{E}_{+} \mathcal{E}_{+} \sum_{k=1}^{2} |c_{k,2}^{i,j-1} - c_{k,1}^{i,j}| \phi_{k}^{i,j} (\overline{x}_{3(i-1)+s}) \right]$$

$$\leq \mathcal{E} + \frac{1}{2} \max_{k=1,2} |c_{k,2}^{i,j-1} - c_{k,1}^{i,j}|$$

$$\leq 2 \mathcal{E},$$

where the second expression was obtained using the fact that  $\phi_k^{ij} = \phi_k^{i,j-1} \ge 0$ , the third from the fact that  $\sum_{k=1}^{2} \phi_k^{ij} \equiv 1$ , and the last from the triangle inequality and (5.48).

From (5.39), (5.43), (5.45), (5.49), (5.47) and 
$$(x, y) \in [x_i, x_{i+1}] \times [y_j, y_{j+1}]$$

$$(5.52) \quad \left| (\mathbf{f} - \widetilde{\mathbf{s}}_{ij})(\mathbf{x}, \mathbf{y}) \right| \leq \left| (\mathbf{f} - \mathbf{s}_{ij})(\mathbf{x}, \mathbf{y}) \right| + \left| (\mathbf{s}_{ij} - \widetilde{\mathbf{s}}_{ij})(\mathbf{x}, \mathbf{y}) \right|$$

$$\leq \mathcal{E} + \sum_{\substack{s=2 \ t=1}}^{3} 2\mathcal{E}[\overline{\phi}_{s}^{ij}(\mathbf{x}) | \psi_{t}^{ij}(\mathbf{y}) \\ + \sum_{\substack{s=2 \ t=1}}^{2} 2\mathcal{E}[\overline{\psi}_{t}^{ij}(\mathbf{y}) | \phi_{s}^{ij}(\mathbf{x}) + \sum_{\substack{s=1 \ t=1}}^{2} \sum_{\substack{s=1 \ t=1}}^{2} (3\mathcal{E}/2) \phi_{s}^{ij}(\mathbf{x}) \psi_{t}^{ij}(\mathbf{y}) \\ \leq (5/2)\mathcal{E} + 2\mathcal{E} \sum_{\substack{s=2 \ s=2}}^{3} (|\overline{\phi}_{s}^{ij}(\mathbf{x})| + |\overline{\psi}_{s}^{ij}(\mathbf{y})|) \\ \leq (5/2)\mathcal{E} + 2\mathcal{E} \{ \max_{\substack{l \leq i \leq M-1}} \left[ \sum_{\substack{s=2 \ s=2}}^{3} ||\overline{\phi}_{3(i-1)+s}|| \right] \right]$$

+ max 
$$\begin{bmatrix} 3\\ \Sigma\\ t=2 \end{bmatrix} |\overline{\psi}_{3(j-1)+t}||$$
 },

where we have used the fact that  $\phi_s^{ij} \ge 0$ ,  $\psi_t^{ij} \ge 0$ ,  $\sum_{s=1}^{2} \phi_s^{ij} \ge 1$ ,

$$\sum_{s=1}^{2} \psi_{s}^{ij} \equiv 1, \quad \overline{\phi}_{s}^{ij} = \overline{\phi}_{3(i-1)+s} \quad \text{on} \quad [x_{i}, x_{i+1}] \quad \text{and} \quad \overline{\psi}_{s}^{ij} = \overline{\psi}_{3(j-1)+s}$$
  
on  $[y_{j}, y_{j+1}]$ .

We now show that  $\tilde{s}_{ij} = \tilde{s}$  on  $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$ . This is accomplished by making the following observations. For  $1 \le s, t \le 2$ ,  $1 \le i \le M-1$ ,  $1 \le j \le N-1$  we have  $\bar{\phi}_s^{ij} = \bar{\phi}_{3(i-1)+s}$ ,  $\bar{\psi}_t^{ij} = \bar{\psi}_{3(j-1)+t}$ ,  $\phi_s^{ij} = \phi_{i+s-1}$  and  $\psi_t^{ij} = \psi_{j+t-1}$ . The following three sets are sets of equal elements

$$\{\widetilde{a}_{2+u,1+t}^{i,j+v-t}\}_{0\leqslant t\leqslant 1}, \{\widetilde{b}_{1+s,2+v}^{i+u-s,j}\} \text{ and } \{\widetilde{c}_{1+s,1+t}^{i+u-s,j+v-t}\}_{0\leqslant s,t\leqslant 1}\}$$

for  $0 \le u, v \le 1$ , which follows from (5.40), (5.41) and (5.42), respectively. Direct substitution of the above sets and relations into (5.44) shows that  $\tilde{s}_{ij} = \tilde{s}$ .

<u>Remark</u>: Once the meshes  $\overline{\pi}_x$  and  $\overline{\pi}_y$  have been specified, then direct calculation gives a bound for

$$\max_{\substack{l \leq i \leq M-1}} \left[ \begin{array}{c} 3 \\ \Sigma \\ s=2 \end{array} \middle| \left| \overline{\phi}_{3(i-1)+s} \right| \left| \right] + \max_{\substack{l \leq j \leq N-1}} \left[ \begin{array}{c} 3 \\ \Sigma \\ t=2 \end{array} \middle| \left| \overline{\psi}_{3(j-1)+t} \right| \right| \right].$$

For example, if on each of the intervals  $[x_i, x_{i+1}]$   $([y_j, y_{j+1}])$ we assume that the four points  $\{\overline{x}_{3(i-1)+s}\}_{s=1}^4$   $(\{\overline{y}_{3(j-1)+t}\}_{t=1}^4)$ are equally spaced, then direct calculation shows that for Theorem

$$(5.53) \quad ||f-\widetilde{s}|| \leq (5/2 + 4(10 + 7\sqrt{7})/27) \mathcal{E}.$$

<u>Remark</u>: If Theorem 5.4 is applied to each of the spaces  $DBF_{ij}$ to obtain an upper bound for  $||f-s_{ij}||[x_i, x_{i+1}]x[y_j, y_{j+1}]$ , then the maximum of these can be used for  $\mathcal{E}$  in Theorem 5.5.

### Section 6. Domain Transformations

In the previous sections, the domain on which the data was given was always the rectangle [a, b]x[c, d]. In this section, a method will be presented which will remove this restriction and allow singly and multiply connected domains with curved boundarys. This will be accomplished by using vector valued blending techniques (see Gordon and Hall [16], [17]) for domain transformations. Their procedure is to divide the boundary of a bounded region  $\Omega < \mathbb{R}^2$  into four parameterized boundary curves. Blending these four curves yields a mapping  $\vec{U}: [0, 1] \times [0, 1] \longrightarrow \mathbb{R}^2$  which maps the boundary of  $[0, 1] \times [0, 1]$  onto the boundary of  $\Omega$ . Of major concern is, under what conditions is the map  $\vec{U}$  univalent? Some conditions have been given in [16] to insure univalency, however, the major responsibility of insuring univalency usually rests with the ability of the person implementing the scheme. For application to least squares, it is imperative to know the point  $\overrightarrow{U}^{-1}(x, y)$ , where  $(x, y) \in X \subset \Omega$  is a data point. Even though an inverse for  $\vec{U}$  may exist, it is not explicitly known. Thus, the following procedure is introduced which allows the inverse of  $\overrightarrow{U}$  to be calculated at any point  $(x, y) \in \Omega$ . This will be accomplished by considering several special domains and then subdividing the domain  $\Omega$  into a finite collection of these special domains.

## Section 6.1 Type 1 Domain Transformations

Consider the closed region  $\Omega < R^2$  which has the following form



Figure 2. The Region  $\Omega$ .

where the distinct points  $\vec{A} = (a_1, a_2)$ ,  $\vec{B} = (b_1, b_2)$ ,  $\vec{C} = (c_1, c_2)$  and  $\vec{D} = (d_1, d_2)$  are ordered as in Figure 2 to form a quadrilateral. We desire a map  $\vec{U}: I^2 \longrightarrow \Omega$ , where  $I^2 = [0, 1] \times [0, 1]$ , whose range is exactly  $\Omega$  and is univalent.

Let  $F \in C^{(1,1)}(\Omega^*)$ , where  $\Omega^*$  is some open region containing the curve F(x,y) = 0 from  $\overrightarrow{A}$  to  $\overrightarrow{B}$  where

(6.1) 
$$(\mathbf{F}^{(1,0)}(\mathbf{x},\mathbf{y}))^2 + (\mathbf{F}^{(0,1)}(\mathbf{x},\mathbf{y}))^2 > 0$$

for all points (x, y) on the curve F(x, y) = 0 from  $\vec{A}$  to  $\vec{B}$ .

For each  $s \in [0, 1]$ , construct the straight line  $\overrightarrow{L}(r;s)$  that contains the two points  $(1-s)\overrightarrow{A} + s\overrightarrow{B}$  and  $(1-s)\overrightarrow{D} + s\overrightarrow{C}$  (which are on the straight lines  $\overrightarrow{AB}$  and  $\overrightarrow{DC}$ , respectively)

(6.2) 
$$\overrightarrow{L}(r;s) = r\{(1-s)\overrightarrow{A} + s\overrightarrow{B} - ((1-s)\overrightarrow{D} + s\overrightarrow{C})\} + (1-s)\overrightarrow{D} + s\overrightarrow{C},$$

where  $r \in \mathbb{R}$  and s is a parameter. The x and y components of

 $\vec{L}(r;s)$  are expressed as  $L^{x}(r;s)$  and  $L^{y}(r;s)$ , respectively. The line  $\vec{L}(r;s)$  can also be expressed in the form

(6.3) 
$$y = \frac{(1-s)(a_2-d_2) + s(b_2-c_2)}{(1-s)(a_1-d_1) + s(b_1-c_1)} (x-((1-s)a_1+sb_1))+(1-s)a_2+sb_2$$

or

(6.4) 
$$y = \eta(s) x + \xi(s)$$
.

<u>Remark 6.1</u>: If for some  $s \in [0, 1]$ , the denominator of  $\eta(s)$  is zero, i.e.,  $(1-s)(a_1-d_1) + s(b_1-c_1) = 0$ , or if  $|\eta(s)| >> 1$ , then express the line  $\overrightarrow{L}(r; s)$  as  $x = \widetilde{\eta}(s) y + \widetilde{\xi}(s)$  where  $\widetilde{\eta}(s) = 1/\eta(s)$ and, in the following discussion, interchange the rolls of the x and y variables. Without loss of generality, throughout this section we will assume that  $\eta(s)$  is finite for the value of s considered.

The following assumptions about the lines  $\overrightarrow{L}(\mathbf{r}; \mathbf{s})$   $(\mathbf{y} = \eta(\mathbf{s}) \mathbf{x} + \xi(\mathbf{s}))$  and the curve  $F(\mathbf{x}, \mathbf{y}) = 0$  are made. For each  $\mathbf{s} \in [0, 1]$ , the straight line  $\overrightarrow{L}(\mathbf{r}; \mathbf{s})$   $(\mathbf{y} = \eta(\mathbf{s}) \mathbf{x} + \xi(\mathbf{s}))$  intersects the curve  $F(\mathbf{x}, \mathbf{y}) = 0$  once and only once. Also, each point of  $F(\mathbf{x}, \mathbf{y}) = 0$  from  $\overrightarrow{A}$  to  $\overrightarrow{B}$  lies on a line  $\overrightarrow{L}(\mathbf{r}; \mathbf{s})$   $(\mathbf{y} = \eta(\mathbf{s}) \mathbf{x} + \xi(\mathbf{s}))$  for some  $\mathbf{s} \in [0, 1]$ (these conditions are usually not too restrictive, since a domain often can be subdivided until it holds). The line  $\overrightarrow{L}(\mathbf{r}; \mathbf{s})$   $(\mathbf{y} = \eta(\mathbf{s}) \mathbf{x} + \xi(\mathbf{s}))$  does not intersect the curve  $F(\mathbf{x}, \mathbf{y}) = 0$  tangentially. Also, for  $\mathbf{s} \in [0, 1]$ , assume that for distinct values of  $\mathbf{s}$  the lines  $\overrightarrow{L}(\mathbf{r}; \mathbf{s})$   $(\mathbf{y} = \eta(\mathbf{s}) \mathbf{x} + \xi(\mathbf{s}))$  do not have a point of intersection in the region  $\Omega$  and that the curve F(x, y) = 0 does not intersect the interior of the line  $\overrightarrow{DC}$ . Finally, for  $s \in [0, 1]$ , if  $\overrightarrow{P}$  is that point of intersection of the straight line  $\overrightarrow{L}(r; s)$  and curve F(x, y) = 0, then  $\{t(\overrightarrow{P} - ((1-s)\overrightarrow{D} + s\overrightarrow{C})) + (1-s)\overrightarrow{D} + s\overrightarrow{C}| \quad 0 \leq t \leq 1 \} \leq \Omega$ .

Remark 6.2: This final assumption can be proved from the others, but its proof leads us away from the desired results of this section.

# Parameterization of F(x, y) = 0

We wish to locate the point designated by  $(x(s), y(s)) = \overrightarrow{F}(s) \in \Omega^*$  where the straight line  $\overrightarrow{L}(r; s)$   $(y = \eta(s) x + \xi(x))$  intersects the curve F(x, y) = 0. This can be accomplished by either one of the following two procedures which we develop simultaneously.

For a fixed  $s \in [0, 1]$ , let

(6.5a) 
$$z(r) = F(L(r; s))$$

and

(6.5b) 
$$w(x) = F(x, \eta(s) x + \xi(s))$$
.

The root  $\hat{\mathbf{r}}$  where  $\mathbf{z}(\hat{\mathbf{r}}) = 0$  or  $\mathbf{x}(\mathbf{s})$  where  $\mathbf{w}(\mathbf{x}(\mathbf{s})) = 0$  is the desired solution. One of the many root finding techniques for univariate functions can be employed to locate the root in (6.5a) and (6.5b); for example, Newton's method or the method of false position (see [24, pp. 97-100]). If Newton's method is chosen, then the derivatives  $\mathbf{z}'(\mathbf{r})$  and  $\mathbf{w}'(\mathbf{x})$  are calculated as follows.

If  $\overrightarrow{L}(r; s)$  and  $(x, \eta(s) x + \xi(s)) \in \Omega^*$ , then because  $F \in C^{(1,1)}(\Omega^*)$ , we have

$$(6.6a) z'(r) = F^{(1,0)}(\overrightarrow{L}(r;s))(L^{x}(r;s))' + F^{(0,1)}(\overrightarrow{L}(r;s))(L^{y}(r;s))'$$

$$= F^{(1,0)}(\overrightarrow{L}(r;s))\{(1-s)(a_{1}-d_{1}) + s(b_{1}-c_{1})\}$$

$$+ F^{(0,1)}(\overrightarrow{L}(r;s))\{(1-s)(a_{2}-d_{2}) + s(b_{2}-c_{2})\}$$

$$(1,0) (0,1)$$

(6.6b)  $\mathbf{w}'(\mathbf{x}) = \mathbf{F}^{(1,0)}(\mathbf{x},\eta(\mathbf{s})\mathbf{x} + \xi(\mathbf{s})) + \mathbf{F}^{(0,1)}(\mathbf{x},\eta(\mathbf{s})\mathbf{x} + \xi(\mathbf{s}))\eta(\mathbf{s}).$ 

Let  $\hat{r}$  and x(s) correspond to the roots of (6.6a) and (6.6b), respectively. It will be shown that both  $z'(\hat{r})$  and w'(x(s)) are non-zero.

If  $z'(\hat{r}) = 0$ , then without loss of generality, from (6.1) we assume that  $F^{(0,1)}(\vec{L}(\hat{r};s)) \neq 0$ . Thus from (6.6a)

(6.7) 
$$(L^{\mathbf{x}}(\hat{\mathbf{r}};\mathbf{s}))' \{-F^{(1,0)}(\overset{\rightarrow}{\mathbf{L}}(\hat{\mathbf{r}};\mathbf{s})) / F^{(0,1)}(\overset{\rightarrow}{\mathbf{L}}(\hat{\mathbf{r}};\mathbf{s}))\} = (L^{\mathbf{y}}(\hat{\mathbf{r}};\mathbf{s}))'.$$

If  $(L^{x}(\hat{r};s))' = 0$ , then  $(L^{y}(\hat{r};s))' = 0$ . Using (6.2) this implies that the straight lines  $\overline{AB}$  and  $\overline{DC}$  intersect, which cannot happen, implying that  $(L^{x}(\hat{r};s)) \neq 0$ . The Implicit Function Theorem (see [10, p. 256]) implies that the slope of the curve F(x, y) = 0 at  $\vec{L}(\hat{r};s)$  is given by the term in brackets of (6.7) which would equal  $(L^{y}(\hat{r};s))'/(L^{x}(\hat{r};s))' = \eta(s)$  which is the slope of the straight line  $(\vec{L}(r;s)$ . Our nontangential intersection assumption guarantees that this cannot happen, implying that  $z'(\hat{r}) \neq 0$ . The argument that  $w'(x(s)) \neq 0$ , is similar and is omitted.

For sufficiently close initial estimates  $r_1$  and  $x_1$ , for  $i = 1, 2, \cdots$ 

(6.8a) 
$$r_{i+1} = r_i - z(r_i)/z'(r_i)$$

and

(6.8b) 
$$x_{i+1} = x_i - w(x_i) / w'(x_i)$$
,

where the  $r_i$  converge to  $\hat{r}$  and the  $x_i$  converge to x(s).

<u>Remark 6.3</u>: If  $F \in C^{(2,2)}(\Omega^*)$ , then z(r) and w(x) are twice differentiable in a neighborhood of  $\hat{r}$  and x(s), respectively. Under these conditions, the convergence of Newton's method is quadratic (see [24, p. 98]).

<u>Remark 6.4:</u> Finding an initial estimate for  $\hat{r}$  is often easier than for x(s), because r = 1 corresponds to a point on the line  $\overline{AB}$  and the curve F(x, y) = 0 is usually "somewhat near"  $\overline{AB}$ .

Having calculated  $\hat{r}$  and x(s), then

$$(6.9a) \qquad \overrightarrow{L}(\mathbf{r}; \mathbf{s})$$

or

(6.9b) 
$$(x(s), \eta(s) x(s) + \xi(s)) = (x(s), y(s))$$

is the point of intersection of F(x, y) = 0 and  $\overrightarrow{L}(r; s)$   $(y = \eta(s) x + \xi(s))$ . This completes the parameterization of the curve F(x, y) = 0.
<u>Remark 6.5:</u> If the curve is given as y = f(x) or x = g(y), then the extension of the above to this case is obvious.

# Calculation of $\overrightarrow{U}(s,t)$

The straight lines  $\overline{DA}$ ,  $\overline{CB}$  and  $\overline{DC}$  are parameterized linearly as follows for  $s \in [0, 1]$  and  $t \in [0, 1]$ 

(6.10a)  $t(\overrightarrow{A} - \overrightarrow{D}) + \overrightarrow{D} = (1 - t)\overrightarrow{D} + t\overrightarrow{A}$ 

(6.10b) 
$$t(\overrightarrow{B} - \overrightarrow{C}) + \overrightarrow{C} = (1 - t) \overrightarrow{C} + t\overrightarrow{B}$$

and

(6.10c) 
$$\mathbf{s}(\overrightarrow{\mathbf{C}} - \overrightarrow{\mathbf{D}}) + \overrightarrow{\mathbf{D}} = (1 - \mathbf{s})\overrightarrow{\mathbf{D}} + \mathbf{s}\overrightarrow{\mathbf{C}}$$
.

For notational convenience we define the vector  $\vec{F}(s) = (x(s), y(s))$ , where there is no confusion between the vector  $\vec{F}(s)$  and the curve F(x, y) = 0 of which (x(s), y(s)) is a point.

Linearly blend the four curves  $\overrightarrow{DA}$ ,  $\overrightarrow{CB}$ ,  $\overrightarrow{DC}$  and F(x, y) = 0(see [16] for procedure) to obtain the vector valued map  $\overrightarrow{U}:I^2 \longrightarrow \mathbb{R}^2$ (6.11)  $\overrightarrow{U}(s, t) = (1-s) \{(1-t)\overrightarrow{D} + t\overrightarrow{A}\} + s\{(1-t) \overrightarrow{C} + t\overrightarrow{B}\}$   $+ (1-t) \{(1-s) \overrightarrow{D} + s\overrightarrow{C}\} + t\overrightarrow{F}(s)$   $- (1-s) (1-t) \overrightarrow{D} - (1-s) t\overrightarrow{A} - s(1-t)\overrightarrow{C} - st\overrightarrow{B}$  $= (1-t) \{(1-s) \overrightarrow{D} + s\overrightarrow{C}\} + t\overrightarrow{F}(s)$ . <u>Remark 6.6</u>: Observe that  $\overrightarrow{U}(s, 0) = (1-s)\overrightarrow{D} + s\overrightarrow{C}$ ,  $\overrightarrow{U}(s, 1) = \overrightarrow{F}(s)$ ,  $\overrightarrow{U}(0, t) = (1-t)\overrightarrow{D} + t\overrightarrow{F}(0) = (1-t)\overrightarrow{D} + t\overrightarrow{A}$  and  $\overrightarrow{U}(1, t) = (1-t)\overrightarrow{C} + t\overrightarrow{F}(1) = (1-t)\overrightarrow{C} + t\overrightarrow{B}$ . Thus the mapping  $\overrightarrow{U}$  carries the boundary of  $I^2$  onto the boundary of  $\Omega$ .

 $\overrightarrow{\underline{U}} \text{ is univalent. Let } (s_1, t_1), (s_2, t_2) \in I^2 \text{ be two points such that}$  $\overrightarrow{\overline{U}}(s_1, t_1) = \overrightarrow{\overline{U}}(s_2, t_2). \text{ Observe that the points } \overrightarrow{\overline{F}}(s_1) \text{ and } (1-s_1) \overrightarrow{\overline{D}} + s_1 \overrightarrow{\overline{C}} \text{ are on the line } \overrightarrow{\overline{L}}(r; s_1) \text{ by construction. From this and (6.11)}$ it is clear that  $\overrightarrow{\overline{U}}(s_1, t_1)$  is on the same line  $\overrightarrow{\overline{L}}(r; s_1)$  because  $(6.12) \qquad \overrightarrow{\overline{U}}(s_1, t) = t\{\overrightarrow{\overline{F}}(s_1) - [(1-s_1)\overrightarrow{\overline{D}} + s_1\overrightarrow{\overline{C}}]\} + (1-s_1)\overrightarrow{\overline{D}} + s_1\overrightarrow{\overline{C}}$ 

is a straight line passing through these two points. Making the corresponding observation that  $\overrightarrow{U}(s_2, t_2)$  is on the line  $\overrightarrow{L}(r; s_2)$ , then we have  $s_1 = s_2$  from the assumption that the lines  $\overrightarrow{L}(r; s)$  have no points of intersection in  $\Omega$  for distinct values of s and for  $t \in [0, 1]$ .

From (6.12) and our assumption that  $\overrightarrow{F}(s)$  does not intersect  $\overrightarrow{DC}$ , it is clear that  $t_1 = t_2$  because

(6.13) 
$$\overrightarrow{U}(s_1, t_2) - \overrightarrow{U}(s_1, t_1) = (t_2 - t_1)(\overrightarrow{F}(s_1) - [(1 - s_1)\overrightarrow{D} + s_1\overrightarrow{C}]),$$

hence, U is a univalent mapping.

<u>Remark 6.7</u>: For linear blending, the map  $\vec{U}$  will not always be univalent. However, with the special parameterization given, this problem is circumvented. <u>Continuity of  $\vec{U}$ </u>. If  $F \in C^{(1,1)}(\Omega^*)$ , then we will prove that  $\vec{U} \in C^{(1,1)}(I^2)$ . Examination of (6.11) shows that the continuity of  $\vec{U}$  is limited only by the continuity of  $\vec{F}(s) = (x(s), y(s))$ , hence it will be shown that x(s),  $y(s) \in C^{(1)}[0, 1]$ .

For a fixed  $\hat{s} \in [0, 1]$ , the denominator of  $\eta(\hat{s})$  is assumed non-zero, i.e.  $|(1-\hat{s})(a_1-d_1) + \hat{s}(b_1-c_1)| = \xi > 0$ . Hence, for  $\mathbf{s} \in (\hat{\mathbf{s}} - \delta, \hat{\mathbf{s}} + \delta) \cap [0, 1]$ , it follows that  $|(1 - \mathbf{s})(\mathbf{a}_1 - \mathbf{d}_1) + \mathbf{s}(\mathbf{b}_1 - \mathbf{c}_1)| \ge 1$  $\mathcal{E}/2 > 0$  where  $\delta = \mathcal{E}/2$  if  $|-(a_1-d_1)+(b_1-c_1)| \leq 1$  or  $\delta = \delta$  $\xi/(2|-(a_1-d_1)+(b_1-c_1)|)$  otherwise. Direct calculation shows that both  $\eta(s)$  and  $\xi(s)$  are continuously differentiable in the interval  $(\hat{s} - \delta, \hat{s} + \delta) \cap [0, 1]$ . Also observe that  $F(x, \eta(s) x + \xi(s))$  is continuously differentiable as a bivariate function of x and s for  $s \in (\hat{s} - \delta, \ \hat{s} + \delta) \cap [0, 1] \text{ and } x \text{ such that } (x, \ \eta(s) \ x + \xi(s)) \in \Omega^*$ because  $F \in C^{(1,1)}(\Omega^*)$ . The implicit function theorem (see [10, p. 257) implies the existence of x(s) which is a unique continuously differentiable function in some neighborhood N of  $\overset{\Lambda}{s}$  if  $\mathbf{F}^{(1,0)}(\mathbf{x}\ (\hat{s}),\ \eta(\hat{s})\ \mathbf{x}(\hat{s}) + \xi(\hat{s})) + \mathbf{F}^{(0,1)}(\mathbf{x}(\hat{s}),\ \eta(\hat{s})\ \mathbf{x}(\hat{s}) + \xi(\hat{s}))\ \eta(\hat{s})$  $\neq$  0 (note that our assumptions guarantee that x(s) and y(s) are both univalent, and we only need prove continuity). It has already been shown that this cannot be zero (see the discussion following (6.6b)). Thus

(6.14) 
$$\mathbf{x}'(\mathbf{s}) = -(\eta'(\mathbf{s}) \mathbf{x}(\mathbf{s}) + \xi'(\mathbf{s})) \mathbf{F}^{(0,1)} / (\mathbf{F}^{(1,0)} + \eta(\mathbf{s}) \mathbf{F}^{(0,1)}),$$

where  $F^{(0,1)}$  and  $F^{(1,0)}$  are evaluated at the point  $(x(s), \eta(s) x(s) + \xi(s))$  for  $s \in \mathbb{N}$ .

From (6.4), the first derivative of y(s) also exists for  $s \in N$ .

<u>Remark 6.8</u>: If higher continuity of F is assumed, then higher derivatives of x(s) and y(s) follow by repeated differentiation of (6.14) and (6.4) (observe that the denominator of (6.14) cannot be zero).

 $\vec{U}([0,1] \times [0,1]) = \Omega$ . Because  $\vec{U}$  is continuous, univalent, and maps the boundary of  $I^2$  onto the boundary of  $\Omega$  (see Remark 6.6), we have from Theorem 13.1 of [27, p. 121] that  $\vec{U}([0,1] \times [0,1]) =$  $\Omega$ , i.e., the range of  $\vec{U}$  is precisely  $\Omega$ .

<u>Remark 6.9</u>: Linear blending will not always yield a map whose range is  $\Omega$ . It is possible for the mapping to "spill over" the boundary of  $\Omega$  into the complement of  $\Omega$ , (see [16]).

Summarizing, to construct a univalent and onto bivariate map  $\vec{U}$ , parameterize F(x, y) = 0 with respect to s by the above method and use linear blending. The observation should be made, given the point  $(s,t) \in I^2$ , that F(x, y) = 0 need only be parameterized for that single value of s.

# Calculating $U^{-1}(x, y)$

A procedure is now given to calculate  $\vec{U}^{-1}(\hat{x}, \hat{y})$  for  $(\hat{x}, \hat{y}) \in \Omega$ . It was shown above that  $\vec{U}$  is a univalent onto map; thus a unique point  $(\hat{s}, \hat{t}) \in I^2$  exists such that  $\vec{U}(\hat{s}, \hat{t}) = (\hat{x}, \hat{y})$ .

Coordinate  $\hat{s}$  is calculated first. From (6.11), observe that  $(\hat{x}, \hat{y})$  lies on the straight line  $\vec{L}(r; \hat{s})$   $(y = \eta(\hat{s}) x + \xi(\hat{s}))$ . The two points  $(1-\hat{s}) \vec{A} + \hat{s} \vec{B}$  and  $(1-\hat{s}) \vec{D} + \hat{s} \vec{C}$  are also on the same line. Thus, we desire the value  $\hat{s}$  which causes the vector  $(v_1, v_2) =$   $(1-\hat{s}) \vec{A} + \hat{s} \vec{B} - [(1-\hat{s}) \vec{D} + \hat{s} \vec{C}]$  to be parallel to the vector  $(w_1, w_2) =$   $(\hat{x}, \hat{y}) - [(1-\hat{s}) \vec{D} + \hat{s} \vec{C}]$ , i.e.  $(w_1, w_2)$  to be a scalar multiple of  $(\dot{v}_1, v_2)$ . If  $(v_1, v_2) \neq (0, 0)$ , which is the case because  $\vec{A} \neq \vec{D}$  and  $\vec{B} \neq \vec{C}$ , then the above is equivalent to

(6.15) 
$$\det \begin{pmatrix} \mathbf{w}_1 & \mathbf{v}_1 \\ & & \\ \mathbf{w}_2 & \mathbf{v}_2 \end{pmatrix} = \mathbf{0} \, .$$

Direct calculation shows that s must satisfy

(6.16) 
$$0 = K_1 \hat{s}^2 + K_2 \hat{s} + K_3,$$

where

(6.17) 
$$K_1 = (d_1 - c_1) (b_2 - a_2) - (d_2 - c_2) (b_1 - a_1)$$

(6.18) 
$$K_2 = (d_1 - c_1) (a_2 - d_2) + (\hat{x} - d_1) (b_2 - c_2 + d_2 - a_2)$$
  
 $-(d_2 - c_2) (a_1 - d_1) - (\hat{y} - d_2) (b_1 - c_1 + d_1 - a_1)$ 

and

(6.19) 
$$K_3 = (\hat{x} - d_1)(a_2 - d_2) - (\hat{y} - d_2)(a_1 - d_1)$$
.

If  $K_1 = 0$ , then this is equivalent to

(6.20) 
$$\det(\overrightarrow{B} - \overrightarrow{A}, \overrightarrow{C} - \overrightarrow{D}) = 0$$

or that  $\overline{AB}$  be parallel to  $\overline{DC}$ . Hence, equation (6.16) reduces to a linear equation. In either case, because of our construction,  $\hat{s}$  is unique in the interval [0, 1].

The coordinate  $\hat{t}$  is computed as follows. First, calculate the point of intersection  $\vec{F}(\hat{s})$ , of the line  $\vec{L}(r; \hat{s})$   $(y = \eta(\hat{s}) x + \xi(\hat{s}))$ and the curve F(x, y) = 0 (the procedure is identical to that given above).  $\hat{t}$  is given by

(6.21) 
$$\hat{\mathbf{t}} = \left| \left| (\hat{\mathbf{x}}, \hat{\mathbf{y}}) - \left\{ (1 - \hat{\mathbf{s}}) \overrightarrow{\mathbf{D}} + \hat{\mathbf{s}} \overrightarrow{\mathbf{C}} \right\} \right| \left| \right| \left| \overrightarrow{\mathbf{F}}(\hat{\mathbf{s}}) - \left\{ (1 - \hat{\mathbf{s}}) \overrightarrow{\mathbf{D}} + \hat{\mathbf{s}} \overrightarrow{\mathbf{C}} \right\} \right| \right|.$$

The denominator of (6.21) is non-zero because it is assumed that the curve F(x, y) = 0 does not intersect the line  $\overline{DC}$ .

Using the fact that  $(1-\hat{s})\vec{A} + \hat{s}\vec{B}$ ,  $(1-\hat{s})\vec{D} + \hat{s}\vec{C}$ ,  $\vec{F}(\hat{s})$  and  $(\hat{x}, \hat{y})$ are all on the line  $\vec{L}(r; \hat{s})$   $(y = \eta(\hat{s}) x + \xi(\hat{s}))$  by construction, then direct calculation shows that  $\vec{U}(\hat{s}, \hat{t}) = (\hat{x}, \hat{y})$ .

To reiterate,  $\vec{U}^{-1}(\hat{x}, \hat{y})$  is calculated by first solving for  $\hat{s}$  from the quadratic (6.16), parameterizing F(x, y) = 0 to obtain  $\vec{F}(\hat{s})$ , and then calculating  $\hat{t}$  from (6.21).

### Section 6.2. Type 2 Transformations

The above procedure is modified by letting two of the points in Figure 2 coincide, either  $\overrightarrow{A} = \overrightarrow{D}$ ,  $\overrightarrow{B} = \overrightarrow{C}$ ,  $\overrightarrow{A} = \overrightarrow{B}$  or  $\overrightarrow{D} = \overrightarrow{C}$  and the resulting procedure will be denoted as a type 2 domain transformation.

Type 2a Transformations. The case where  $\vec{A} = \vec{D}$  will be examined first, observing that the case where  $\overrightarrow{B} = \overrightarrow{C}$  is nearly identical. Nothing is changed in parameterizing F(x, y) = 0 to obtain  $\overrightarrow{F}(s)$ (using the fact that if s = 0, then  $\overrightarrow{F}(0) = \overrightarrow{A} = \overrightarrow{D}$ ), and  $\overrightarrow{U}(s, t)$  is given by equation (6.11).  $\vec{U}$  is no longer univalent, because the points (0, t) for  $t \in [0, 1]$  are all mapped to the single point  $\vec{A} = \vec{D}$  (see Remark 6.6). However, the argument given above proves that  $\vec{U}$ is univalent on  $(0, 1] \times [0, 1]$ . Nothing is changed in the argument on the continuity of  $\vec{U}$ , taking  $\eta(s) = (b_2 - c_2)/(b_1 - c_1) = \eta$  for  $s \in [0,1]$ . However, the proof that  $\overrightarrow{U}(I^2) = \Omega$  is no longer valid for this case because  $\overrightarrow{U}$  is no longer univalent. To circumvent this problem, we first show that  $\stackrel{\rightarrow}{U}$  cannot map to a point outside of  $\Omega$ . Let  $(\hat{s}, \hat{t}) \in (0, 1] \times [0, 1]$ , then from the above construction,  $\vec{F}(\hat{s})$  and  $\vec{U}(\hat{s}, \hat{t})$  are both on the straight line  $\vec{L}(r; \hat{s})$  which can be expressed as  $\overrightarrow{U}(\hat{s}, t) = t\{\overrightarrow{F}(\hat{s}) - [(1-\hat{s})\overrightarrow{D} + \hat{s}\overrightarrow{C}]\} + (1-\hat{s})\overrightarrow{D} + \hat{s}\overrightarrow{C}$ . Because  $\hat{t} \in [0, 1]$ , our assumptions guarantee that  $\overrightarrow{U}(\hat{s}, \hat{t}) = (\hat{x}, \hat{y}) \in \Omega$ , proving that  $\overrightarrow{U}$ maps into  $\Omega$ . To prove that  $\overrightarrow{U}$  maps onto  $\Omega$ , let  $(\overset{\wedge}{x}, \overset{\wedge}{y}) \in \Omega$ . If  $(\hat{x}, \hat{y}) = \overline{A} = \overline{D}$  or  $(\hat{x}, \hat{y})$  is a point on the line  $\overline{BC}$ , then take  $\hat{s} = 0$  or

= 1, respectively. If  $(\hat{x}, \hat{y})$  is a point on the curve F(x, y) = 0 or the straight line  $\overrightarrow{AC} = \overrightarrow{DC}$ , then take  $\hat{s}$  such that  $(\hat{x}, \hat{y}) = \overrightarrow{F}(\hat{s})$  or = $(1-\hat{s})\vec{D}+\hat{s}\vec{C}$ , respectively. If  $(\hat{x},\hat{y})$  is not on the boundary of  $\Omega$ , then  $(\hat{x}, \hat{y})$  is a point on the straight line  $y = \eta(x - \hat{x}) + \hat{y}$ , which is parallel to the line  $\overline{BC}$  (recall  $\eta(s) = \eta = \text{constant for } s \in [0, 1]$ ). This line must intersect the boundary of  $\Omega$  at least twice and our assumptions imply that this line must also be one of the lines  $\vec{L}(r; \hat{s})$  for some  $\hat{s} \in (0, 1)$ . To see this, note that  $y = \eta(x - \hat{x}) + \hat{y}$ cannot intersect the point  $\overrightarrow{A} = \overrightarrow{D}$  or the line  $\overrightarrow{BC}$  which is parallel to it. If it intersects the curve F(x, y) = 0 at some point, say  $\vec{F}(\vec{s})$ , then the lines  $y = \eta(x \cdot \hat{x}) + \hat{y}$  and  $\vec{L}(r; \hat{s})$  are identical because they have the same slope and a point in common. A similar argument holds if  $y = \eta(x - \hat{x}) + \hat{y}$  intersects the curve  $\overline{DC}$  at some point  $(1-\hat{s})\overrightarrow{D}+\hat{s}\overrightarrow{C}$ , in either case take  $s = \hat{s}$ . Also note that this value of s is unique because  $y = \eta(x \cdot \hat{x}) + \hat{y}$  is a parameter line  $\vec{L}(r; \hat{s})$ .

If  $\hat{s} \neq 0$ , then compute  $\hat{t}$  from (6.21), and direct calculation shows that  $\overrightarrow{U}(\hat{s}, \hat{t}) = (\hat{x}, \hat{y})$  (if  $\hat{s} = 0$ , then any  $\hat{t} \in [0, 1]$  is acceptable).

<u>Type 2b Transformations</u>. Consider the case where  $\overrightarrow{D} = \overrightarrow{C}$  (note that the case where  $\overrightarrow{A} = \overrightarrow{B}$  is similar and less complex). All of the lines  $\overrightarrow{L}(r; s)$  ( $y = \eta(s) x + \xi(s)$ ) intersect at the point  $\overrightarrow{D} = \overrightarrow{C}$ , but no other point of  $\Omega$ . F(x, y) = 0 is parameterized as above and the map  $\overrightarrow{U}$  reduces to

(6.22) 
$$\overrightarrow{U}(s, t) = (1-t)\overrightarrow{D} + t\overrightarrow{F}(s)$$
.

 $\vec{U}$  is not univalent because the points (s, 0) for  $s \in [0, 1]$  are mapped to the single point  $\vec{D} = \vec{C}$ . However, the above proof for the univalency of  $\vec{U}$  holds for  $(s, t) \in [0, 1] \times (0, 1]$ . The differentiability of  $\vec{U}$  follows from the above argument. Finally, the proof that  $\vec{U}(I^2) = \Omega$  is similar to that of type 2a transformations where we work with the straight line  $t\{(x, y) - \vec{D}\} + \vec{D}$ .

Section 6.3. General Domains. It should be understood that for the region  $\Omega$ , the side on which the curve F(x, y) = 0 is located was specified as it was above for purposes of illustration. If it is desired, then the curve F(x, y) = 0 can be any one of the four sides of the quadrilateral region, where minor modifications of the above presentation will yield the appropriate procedure and formulas for the construction of  $\overrightarrow{U}$ .

We now consider closed regions  $\Gamma$  which can be subdivided into N subregions  $\Omega_i$  such that

$$(6.23) \qquad \Gamma = \bigcup_{i=1}^{\infty} \Omega_i,$$

and for  $1 \leq i, j \leq N, i \neq j$ 

$$(6.24) \quad \operatorname{Int}(\Omega_{i}) \cap \operatorname{Int}(\Omega_{i}) = \emptyset,$$

where each  $\Omega_{j}$  is a type 1 or type 2 region satisfying the assumptions of Section 6.1 or 6.2, respectively. For each subregion

construct the map  $\overrightarrow{U}_i: I^2 \longrightarrow \Omega_i$ . Several examples will be given to illustrate the procedure (where the numbers in the following illus-trations have the obvious correspondence).



### Figure 3

In Figure 3,  $\Omega_1$ ,  $\Omega_3$ ,  $\Omega_4$  and  $\Omega_6$  are type 2b regions (type 2a could be used with minor modifications) and  $\Omega_2$  and  $\Omega_5$  are type 1 regions. Note that the left hand side of  $I_1^2$  and  $I_4^2$  ( $I_3^2$  and  $I_6^2$ ) are identified in the obvious manner and the adjacent side of  $I_1^2$  and  $I_4^2$  ( $I_3^2$  and  $I_6^2$ ) are mapped to the single point  $\vec{P}_1$  ( $\vec{P}_2$ ).



In Figure 4  $\Omega_1$ ,  $\Omega_4$ ,  $\Omega_5$ ,  $\Omega_8$ ,  $\Omega_{11}$  and  $\Omega_{12}$  are type 2a regions and the remaining regions are type 1, where the left hand sides of  $I_1^2$  and  $I_5^2$  are mapped to the point  $\vec{P}_1$ , etc.



Figure 5

In Figure 5, the regions  $\Omega_i$  for  $1 \le i \le 10$  are all type 1, where the left hand side of  $I_1^2$  and  $I_6^2$  are identified with the right hand side of  $I_5^2$  and  $I_{10}^2$ , respectively.

The procedure of subdividing  $\Gamma$  should be lucid from the above examples.

Section 6.4. Discrete Least Squares Over  $\Gamma$ . Let  $f \in C^{(0,0)}(\Gamma)$ , where  $\Gamma$  satisfies (6.23) and (6.24) and let  $X = \{(\hat{x}_k, \hat{y}_k)\}_{k=1}^{\hat{M}} \subseteq \Gamma$ be a set of  $\hat{M}$  unstructured data points. A discrete least squares fit  $\tilde{g}_i:\Omega_i \longrightarrow \mathbb{R}$  to f is constructed over each region  $\Omega_i$  in the following manner. Choose a finite dimensional bivariate real valued function space  $S(I^2)$  over  $I^2$  (for example, any of those of Section 5) and construct the map  $\overrightarrow{U}_i: I^2 \longrightarrow \Omega_i$ . If  $\Omega_i$  is a type 2 domain, then constrain  $S(I^2)$  to insure that each element of  $S(I^2)$  is constant on that single boundary line of  $I^2$  which is mapped to the single point  $\vec{P}$  (see Section 6.2) (for the spaces given in Section 5, this can easily be accomplished by modifying the basis elements). With this restriction there is no ambiguity for the element  $s \cdot \vec{U}_i^{-1}$ , where  $s \in S(I^2)$ , because s is single valued on the inverse image (under  $\vec{U}_i^{-1}$ ) of each point of the type 2 region  $\Omega_i$ . Thus we define  $\vec{U}_i^{-1}(\vec{P})$  to be (for example) the midpoint of the boundary line of  $I^2$  whose image under  $\vec{U}_i$  is the point  $\vec{P}$ .

Define  $X_i = \{(\hat{x}_{ij}, \hat{y}_{ij})\}_{j=1}^{\hat{M}_i} = X \cap \Omega_i$  to be the set of  $\hat{M}_i$  data points which are in  $\Omega_i$  and define  $ST_i = \{\vec{U}_i^{-1}(\hat{x}_{ij}, \hat{y}_{ij})\}_{j=1}^{\hat{M}_i}$ , where for  $1 \leq j \leq \hat{M}_i$ 

(6.25) 
$$(\hat{s}_{ij}, \hat{t}_{ij}) = \overrightarrow{U}_i^{-1} (\hat{x}_{ij}, \hat{y}_{ij}).$$

Finally, define the finite dimensional function space  $V(\Omega_i) = \{s \circ \overrightarrow{U}_i^{-1} \mid s \in S(I^2)\}.$ 

We say  $\tilde{g}_i \in V(\Omega_i)$  ( $\tilde{s}_i \in S(I^2)$ ) is a discrete least squares fit to f (f  $\cdot \vec{U}_i$ ) on the data set  $X \leq \Omega_i$  (ST $_i \leq I^2$ ) if the Euclidean norm of the residual vector  $R_i \in \mathbb{R}^{\hat{M}_i}$  ( $\tilde{R}_i \in \mathbb{R}^{\hat{M}_i}$ ) is minimized, where component j of  $R_i$  is given by

(6.26) 
$$(\mathbf{R}_{i})_{j} = f(\hat{\mathbf{x}}_{ij}, \hat{\mathbf{y}}_{ij}) - \tilde{\mathbf{g}}_{i}(\hat{\mathbf{x}}_{ij}, \hat{\mathbf{y}}_{ij})$$

and component j of  $\widetilde{R}_{i}$  is given by

(6.27) 
$$(\widetilde{R}_{i})_{j} = (f \cdot \overrightarrow{U}_{i}) (\widehat{s}_{ij}, \widehat{t}_{ij}) - \widetilde{s}_{i}(\widehat{s}_{ij}, \widehat{t}_{ij}) .$$

It is clear from (6.25), (6.26) and (6.27) that if  $\widetilde{g}_{i}$  ( $\widetilde{s}_{i}$ ) is a discrete least squares fit, then there exists  $\mathbf{s}^{*} \in S(\mathbf{I}^{2})$  ( $\mathbf{g}^{*} = \widetilde{s}_{i} \circ \overrightarrow{U}_{i}^{-1} \in V(\Omega_{i})$ ) such that  $\mathbf{s}^{*} \circ \overrightarrow{U}_{i}^{-1} = \widetilde{g}_{i}$  ( $\mathbf{g}^{*} \circ \overrightarrow{U}_{i} = \widetilde{s}_{i}$ ) where (6.28)  $||\mathbf{R}_{i}|| = (\sum_{j=1}^{\widehat{M}_{i}} [(\mathbf{f} \circ \overrightarrow{U}_{i})(\widehat{s}_{ij}, \widehat{t}_{ij}) - \widetilde{\mathbf{s}}^{*}(\widehat{s}_{ij}, \widehat{t}_{ij})]^{2})^{1/2}$ ,

and

(6.29) 
$$||\widetilde{R}_{i}|| = (\sum_{j=1}^{\hat{M}_{i}} [f(\widehat{x}_{ij}, \widehat{y}_{ij}) - g^{*}(\widehat{x}_{ij}, \widehat{y}_{ij})]^{2})^{1/2}$$

showing that  $||R_i|| = ||\widetilde{R_i}||$  because  $||R_i|| \ge ||\widetilde{R_i}||$  and  $||R_i|| \le ||\widetilde{R_i}||$ . Therefore, a discrete least squares fit  $\widetilde{g}_i \in V(\Omega_i)$ can be calculated from a discrete least squares fit  $\widetilde{s}_i \in S(I^2)$  by taking

(6.30) 
$$\widetilde{g}_i = \widetilde{s}_j \circ \widetilde{U}_i^{-1}$$
.

Observing that

(6.31) 
$$||\mathbf{f} - \widetilde{\mathbf{g}}_i||_{\Omega_i} = ||\mathbf{f} \cdot \overrightarrow{\mathbf{U}}_i - \widetilde{\mathbf{s}}_i||_{I^2}$$

an error analysis for  $f - \tilde{g}_i$  over  $\Omega_i$  can be given in terms of  $f \cdot \vec{U}_i - \tilde{s}_i$  over  $I^2$ .

Thus, the case of least squares over a general domain  $\Gamma$  is reduced to least squares over rectangular domains and the error analysis reduces to bounding  $||f \circ \overrightarrow{U}_i - \widetilde{s}_i||$  where, for example,  $f \circ \overrightarrow{U}_i$ replaces f in the formulas (5.15), (5.16), (5.21), (5.22), (5.30), (5.31), (5.36) and (5.37).

## BIBLIOGRAPHY

,

#### BIBLIOGRAPHY

- 1. Ahlberg, Nilson, and Walsh, <u>The Theory of Splines and Their</u> Applications, Academic Press: New York, 1967.
- Birkhoff, G., and C. DeBoor, "Error Bounds for Spline Interpolation," Journal of Mathematics and Mechanics, Vol. 13, No. 5 (1964) pp. 827-835.
- 3. Björck, Å., "Solving Linear Least Squares Problems by Gram-Schmidt Orthogonalization," <u>Nordisk Tidskrift for Informations</u> -Behandling, Vol. 7 (1967).
- Businger, P., and G. H. Golub, "Linear Least Squares Solutions by Householder Transformations," <u>Numerische Mathematik</u>, Vol. 7 (1965) pp. 269-276.
- Carlson, R. E., and C. A. Hall, "Error Bounds for Bicubic Spline Interpolation," Journal of Approximation Theory, Vol. 7, No. 1, Jan. 1973, pp. 41-47.
- 6. Cheney, E. W., Introduction to Approximation Theory, McGraw-Hall: New York, 1966.
- 7. Deskins, W. E., Abstract Algebra, MacMillan Co.: New York, 1964.
- 8. Dugundji, J., Topology, Allyn and Bacon, Inc.: Boston, 1966.
- Faddeev, D. K., and V. N. Faddeeva, <u>Computational Methods of</u> <u>Linear Algebra</u>, W. H. Freedman and Company: San Francisco, 1963.
- Fulks, W., <u>Advanced Calculus An Introduction to Analysis</u>, John Wiley and Sons, Inc.: New York, 1961.
- 11. Golub, G., "Numerical Methods for Solving Linear Least Squares Problems," Numerische Mathematik, Vol. 7 (1965) pp. 206-216.

- Gordon, W. J. "Blending-Function' Methods of Bivariate and Multivariate Interpolation and Approximation," General Motors Research Publication, GMR-834, Warren, Michigan, October, 1968.
- Gordon, W. J., "Spline-Blended Surface Interpolation Through Curve Networks," Journal of Mathematics and Mechanics, Vol. 18, No. 10, April, 1969, pp. 931-952.
- Gordon, W. J., and C. A. Hall, "Discretization Error Bounds for Transfinite Elements," General Motors Research Publication, GMR-1196, Warren, Michigan, May, 1972.
- 15. Gordon, W. J., and C. A. Hall, "Geometric Aspects of the Finite Element Method," in <u>The Mathematical Foundations of the Finite</u> <u>Element Method with Applications to Partial Differential</u> <u>Equations</u>, Academic Press, Inc.: New York, 1972.
- 16. Gordon, W. J., and C. A. Hall, "Geometric Aspects of the Finite Element Method: Construction of Curvilinear Coordinate Systems and Their Application to Mesh Generation," General Motors Research Publication, GMR-1286, Warren, Michigan.
- Gordon, W. J., and C. A. Hall, "Transfinite Element Methods: Blending-Function Interpolation over Arbitrary Curved Element Domains," <u>Numerische Mathematik</u>, Vol. 21, 1973, pp. 109-129.
- Gordon, W. J., and J. A. Wixom, "Pseudo-Harmonic Interpolation on Convex Domains," General Motors Research Publication, GMR-1248, Warren, Michigan, August 1972.
- Hall, C. A., "Natural Cubic and Bicubic Spline Interpolation," <u>SIAM J. Numerical Analysis</u>, Vol. 10, No. 6, Dec, 1973, pp. 1055-1060.
- 20. Hall, C. A., "On Error Bounds for Spline Interpolation," Journal of Approximation Theory, Vol. 1, 1968, pp. 209-218.
- Halliday, J., and J. G. Hayes, "The Least-Squares Fitting of Cubic Spline Surfaces to General Data Sets," National Physical Laboratory, NPL Report NAC 22, December, 1972.
- 22. Householder, A. S., "Unitary Triangularization of a Nonsymmetric Matrix," Assoc. Comput. Mach., Vol. 5, 1958, pp. 339-342.

- 23. Householder, A. S., <u>The Theory of Matrices in Numerical</u> Analysis, Blardell Publishing Co.: 1964.
- Isaacson, E., and H. B. Keller, <u>Analysis of Numerical Methods</u>, Wiley and Sons, Inc.: New York, 1966.
- Jennings, L. S., and M. R. Osborne, "A Direct Error Analysis for Least Squares," <u>Numerische Mathematik</u>, Vol. 22, 1974, pp. 325-332.
- 26. Marcus, M., and H. Ming, <u>A Survey of Matrix Theory and</u> Matrix Inequalities, Allyn and Bacon Inc.: Boston, 1964.
- 27. Newman, M. H. A., <u>Topology of Plane Sets</u>, Cambridge University Press: 1964.
- Peters, G. and J. H Wilkinson, "The Least Squares Problem and Pseudo-Inverses," <u>The Computer Journal</u>, Vol. 13, No. 3, August, 1970, pp. 309-316.
- 29. Schultz, M H., <u>Spline Analysis</u>, Prentice-Hall, Inc.: Englewood Cliffs, N. J., 1973.
- 30. Soble, A. B., "Majorants of Polynomial Derivatives," <u>Ameri-</u> can Math Monthly, Vol. 64, 1957, pp. 639-643.
- van der Sluis, A., "Stability of the Solution of Linear Least Squares Problems," <u>Numerische Mathematik</u>, Vol. 23, 1975, pp. 241-254.
- 32. Varga, R. S. <u>Matrix Iterative Analysis</u>, Prentice-Hall, Inc.: Englewood Cliffs, N. J., 1962.
- 33. Wendroff, B., <u>Theoretical Numerical Analysis</u>, Academic Press: New York, 1966.
- Wilkinson, J. H., "Error Analysis of Direct Methods of Matrix Inversion," <u>Association for Computing Machinery</u>, Vol. 8, No. 1, Jan, 1961.
- 35. Wilkinson, J. H., "Error Analysis of Transformations Based on the Use of Matrices of the Form I-2ww<sup>H</sup>," in Error in <u>Digital Computation</u>, Vol. 2, Edited by Louis B. Rall, Wiley and Sons, Inc.: 1965.

36. Wilkinson, J. H., <u>Rounding Errors in Algebraic Processes</u>, Prentice-Hall, Inc.: Englewood Cliffs, N.J., 1962.

Second - Censurication of This FAGE (When Date Entered)			
REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM	
1. REPORT NUMBER	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER	
4. TITLE (and Subtitie)		5. TYPE OF REPORT & PERIOD COVERED	
BLENDING-FUNCTION TECHNIQUES WITH		Interim	
APPLICATIONS TO DISCRETE LEAST SQUARES		6. PERFORMING ORG. REPORT NUMBER	
7. AUTHOR(a)		8. CONTRACT OR GRANT NUMBER(=)	
Dale Russel Doty		AFOSR-72-2271	
9. PERFORMING ORGANIZATION NAME AND ADDRESS		10. PROGRAM ELEMENT, PROJECT, TASK	
Michigan State University			
Department of Mathematics		61102F 9749-03	
Last Lansing, NI 40024 11. Controlling office name and address		12. REPORT DATE	
Air Force Office of Scientific Research (NM)		August, 1975	
1400 Wilson Blvd.		13. NUMBER OF PAGES	
Arlington, VA 22209	t from Controlline Office)	L SU 15. SECURITY CLASS (of this report)	
		UNCLASSIFIED	
	:		
		15. DECLASSIFICATION/DOWNGRADING SCHEDULE	
16. DISTRIBUTION STATEMENT (of this Report)			
Approved for public release; distribution unlimited.			
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)			
18. SUPPLEMENTARY NOTES			
19. KEY WORDS (Continue on reverse eide if necessary and identify by block number)			
blending-functions, least squares, splines, error analysis, interpolation,			
natural splines, domain transformations			
20 ABCTRACT (Continue on control of the second of the seco			
The theory of blending function and can (hive viste internals tion) is developed			
in the general setting of interpolation spaces (Divariate Interpolation) is developed			
blending function spaces have the desirable quality of doubling the order of			
accuracy with less computation when compared to standard tensor product			
spaces.	F		
- The dimensionality of discretiv	red blanding from	tion ana cas is derived and	
The unnensionality of discretized	ine dimensionality of discretized blending function spaces is derived, and		

SECURITY CLASSIFICATION OF THIS PAGE (When Date Entered)

#### SECURITY CLASSIFICATION OF THIS PAGE(When Date Entered)

several bases are explicitly constructed. The special example of Hermite spline blended piecewise polynomials is developed, showing that these spaces have bases with small support which are easy to calculate. These spaces offer maximum order of convergence for a minimum number of basis elements. For example, linearly blended piecewise cubic polynomials offer a fourth order approximation scheme, and cubic Hermite spline blended piecewise polynomials offer an eighth order scheme.

Next, using the exponential decay of the natural cubic cardinal splines and natural cubic spline blending, a derivative-free approximation scheme is developed, which is eighth order in the interior of the domain.

Algorithms with corresponding error estimates are given for solving the discrete least squares problem with unstructured data. For the univariate case, algorithms are developed using the space of cubic splines. The resulting error analysis indicates the necessary restrictions to be placed on the number and distribution of the data points to insure that the discrete least squares fit will be  $O(h^{m})$  to a function  $f \in C^{m}[a, b]$  from which the data arises, where h is the mesh size and  $l \leq m \leq 4$ . An example is given to illustrate that the discrete least squares fit need not be close to f if these conditions are not realized. For the bivariate case, algorithms and error analyses are given for the spaces of bicubic splines and discretized blending function spaces. It is shown that the discrete least squares fit to a bivariate function f is of the same order accuracy as the corresponding interpolation accuracy.

Discrete least squares is considered on general domains which have curved boundaries and are possibly multiply connected. This general domain is subdivided into "standard" subdomains, and explicit mappings from the unit square to these standard subdomains are constructed which are one-one, onto, and have easily calculated inverses. Thus, discrete least squares over general domains reduces to the cases previously considered.

Finally, an extensive computational error analysis is given for a constrained least squares algorithm.

