

# The Multivariate Spline Method for Scattered Data Fitting and Numerical Solutions of Partial Differential Equations

Gerard Awanou, Ming-Jun Lai and Paul Wenston

Dedicated to Charles K. Chui on the Occasion of his 65th Birthday

**Abstract.** Multivariate spline functions are smooth piecewise polynomial functions over triangulations consisting of  $n$ -simplices in the Euclidean space  $\mathbb{R}^n$ . A straightforward method for using these spline functions to fit given scattered data and numerically solve elliptic partial differential equations is presented. This method does not require constructing macro-elements or locally supported basis functions nor computing the dimension of the finite element spaces or spline spaces. The method for splines in  $\mathbb{R}^2$  and  $\mathbb{R}^3$  has been implemented in MATLAB. Several numerical examples are shown to demonstrate the effectiveness and efficiency of the method.

## Table of Contents

- [0] Introduction
- [1] Definition of Various Multivariate Spline Spaces
- [2] B-form Representation of Multivariate Splines
- [3] The Constrained Multivariate Spline Method and a Matrix Iterative Algorithm
- [4] The Minimal Energy Method for Scattered Data Interpolation
- [5] The Discrete Least Squares Method for Scattered Data Fitting
- [6] The Penalized Least Squares Method for Scattered Data Fitting
- [7] Numerical Solution of Poisson and Other Second Order Elliptic Equations
- [8] Numerical Solution of Biharmonic Equations
- [9] Remarks
- References

## §0. Introduction

Multivariate spline functions are piecewise polynomial functions which have certain smoothness over given polygonal domains. They are very flexible for approximating known or unknown functions or any given data sets. For example, finite elements are typical spline functions. They are very popular for numerical solutions of partial differential equations (PDE's) and are very useful for surface designs. There are many finite elements and macro-elements available in the literature. See [58], [19], [51], [15], [14], and [10]. For general spline functions, many kinds of locally supported spline functions have been recently constructed based on special triangulations (cf. [37], [38], [33], [34], [35], [41], [43], [44], [45], [46], [1], and [2].) However, the implementation of finite elements and locally supported spline functions with high order smoothness is very difficult. For example,  $C^2$  finite elements or spline functions in the bivariate and trivariate settings are seldom implemented for applications. See [52] and [23] for preliminary testing of  $C^2$  spline functions. See, e.g., [20] and [55] for implementation of bivariate  $C^1$  quadratic and  $C^1$  quintic splines for applications. For another example, the construction of  $C^1$  trivariate splines using polynomials of degree 3 is already very complicated (cf. [57]) and its implementation is even more complicated based on our experience (cf. [48]).

To make multivariate spline functions available for applications, we propose a direct approach without constructing locally supported basis functions and finite elements. This approach may be summarized as follows. First multivariate spline functions are represented by using the B-form of splines. (cf. [24] and [11].) Each spline function  $s$  can be identified by its B-coefficient vector  $\mathbf{c}$ . Since  $s$  has certain smoothness, the smoothness conditions can be expressed by a linear system  $H\mathbf{c} = 0$ . Boundary conditions for elliptic PDE's or interpolation conditions for scattered data fitting provide additional constraints on the B-coefficient vector  $\mathbf{c}$ . These constraints can be given by a linear system  $B\mathbf{c} = \mathbf{g}$ . The constraints are enforced through Lagrange multipliers by minimizing suitable functionals, the classical variational forms for PDE's or the popular thin plate energy functional for scattered data problems. The resulting Euler-Lagrange equations are solved using a matrix iterative method, a variant of the augmented Lagrangian algorithm, which we introduce in §3. We shall prove that the matrix iterative method converges. Similarly, we can fit the given data using the discrete least squares method and the penalized least squares method. The most important feature about this approach is to use multivariate splines for applications without constructing locally supported splines or finite elements and without computing the dimension. Our approach is not completely new. In [25], a similar method of applying bivariate  $C^1$  cubic splines for scattered data fitting

was implemented. What is new about our approach is that we start with discontinuous piecewise polynomials, that we use an iterative method to solve the linear systems arising from the Lagrange multiplier method and that our implementation and experiments with this approach are done for bivariate and trivariate spline functions of arbitrary degrees  $d$  and arbitrary smoothness  $r$  with  $r < d$  over any triangulations or tetrahedral partitions.

The advantages of this approach are

- 1) We are able to use multivariate spline functions in an Euclidean space of arbitrary dimension.
- 2) We are able to use multivariate spline functions of variable degrees and variable smoothness across any given domain. These flexibilities make spline surfaces more user-friendly.
- 3) We are able to use piecewise quadratic or cubic polynomials to approximate the boundary of any given domain instead of piecewise linear polynomials.
- 4) The linear systems arising from our approach are more easily assembled than those from the finite (macro) elements or locally supported spline basis functions. The linear systems are sparser than that from any macro-FEM method. Also, the assembling can be done in parallel.
- 5) A special iterative method with excellent convergence rate is introduced to solve this linear system.

For bivariate and trivariate settings, we have already implemented this approach in MATLAB successfully. That is, we are able to use bivariate and trivariate spline functions of any degree  $d$  and any smoothness  $r$  with  $r < d$  over any polygonal domains. Many numerical experiments for scattered data interpolation/fitting and numerical solutions of partial differential equations will be reported in later sections. Our experience on personal computers shows that we can use bivariate polynomials of degree 10 or less and trivariate polynomials of degree 7 or less for applications for various smooth surfaces.

The disadvantage of this approach over the popular finite element method (FEM) is the necessity of solving linear systems of larger size. However, efficiency is gained in assembling linear systems and in using spline functions of high degrees and/or of high orders of smoothness. On the other hand, a matrix iterative algorithm which will be introduced in Section 3 works particularly well for the systems that arise from our approach. These systems are classical saddle point problems. The matrix iterative method converges very fast in general. Our experience tells us that two or three iterations will be enough for 4 significant digits.

The rest of the paper is organized as follows. We first present in detail the space of multivariate spline functions in §1. Then we present the B-form representation of multivariate splines, smoothness conditions and

degree reduction conditions in §2. In §3, we outline our general minimization algorithm and a matrix iterative algorithm. With these preparations, we shall deal with scattered data interpolation and fitting in §4–6. Then we shall discuss how to solve Poisson and biharmonic equations in §7 and §8. Multivariate spline solutions of nonlinear PDE's such as the Navier-Stokes equations are reported elsewhere (cf. [50] and [5].)

### §1. Definition of Various Multivariate Spline Spaces

Let  $\Delta$  be a collection of  $n$ -simplices in the Euclidean  $\mathbb{R}^n$ . We say that  $\Delta$  is a triangulation if it satisfies the following property: for any  $t, t' \in \Delta$ ,  $t \cap t'$  is either empty or a common  $k$ -simplex of  $t$  and  $t'$  for some  $k$ ,  $0 \leq k \leq n-1$ . Let  $\Omega = \cup_{t \in \Delta} t$  be a polygonal domain in  $\mathbb{R}^n$ . Given two integers  $d \geq 0$  and  $0 \leq r < d$ , let

$$S_d^r(\Delta) := \{s \in C^r(\Omega) : s|_t \in \mathbb{P}_d, \forall t \in \Delta\}$$

be the multivariate spline space of degree  $d$  and smoothness  $r$ , where  $\mathbb{P}_d$  denotes the space of all polynomials of degree  $\leq d$ . It is a standard spline space. Typically,  $S_d^0(\Delta)$  is the continuous spline space of degree  $d$  which is a very useful finite element space.

Next we introduce super spline subspaces. For simplicity, let us consider bivariate spline spaces first. Let  $\rho = \{\rho_v, v \in \mathcal{V}\}$  be a set of integers  $\rho_v \geq 0$  associated with vertices  $\mathcal{V}$  of  $\Delta$  and  $\mathbf{r} = \{r_e, e \in \mathcal{E}_I\}$  be a set of integers  $r_e \geq 0$  associated with interior edges  $\mathcal{E}_I$  of  $\Delta$ . Suppose that  $\rho_v \geq r \geq 0$  for all  $v \in \mathcal{V}$  and  $r_e \geq r \geq 0$  for all  $e \in \mathcal{E}_I$ . Let

$$S_d^{\mathbf{r}, \rho}(\Delta) = \{s \in S_d^r(\Delta), s \in C^{\rho_v} \text{ at } v \in \mathcal{V} \text{ and } s \in C^{r_e} \text{ across } e \in \mathcal{E}_I\}$$

be the spline subspace of super smoothness  $\rho$ , smoothness  $\mathbf{r}$  and degree  $d$ . Similarly, we can define such spline subspaces in  $\mathbb{R}^n$ . Let  $\Delta$  be an  $n$ -simplicial partition in  $\mathbb{R}^n$ . Let  $\rho$  be a set of integers  $\geq 0$  associated with all  $k$ -simplices with  $0 \leq k < n-1$  and  $\mathbf{r}$  be another set of integers  $\geq 0$  associated with all  $(n-1)$ -simplices of  $\Delta$ . Then  $S_d^{\mathbf{r}, \rho}(\Delta)$  is a super spline space on  $\mathbb{R}^n$  consisting of piecewise polynomials of total degree  $d$  which satisfy smoothness of order  $\rho$  around all  $k$ -simplices for  $0 \leq k < n-1$  and smoothness of order  $\mathbf{r}$  across all  $(n-1)$ -simplices of  $\Delta$ . This spline space of variable smoothness is useful for designing surfaces with variable smoothness across given domains. For numerical solutions of PDE's, it is reasonable to let a spline solution have more smoothness inside the domain and less smoothness near the boundary of the domain because of the regularity theory of PDE's (cf. [22]), that is, the weak solutions usually possess higher smoothness inside  $\Omega$  than near the boundary  $\partial\Omega$ . Although adding the extra smoothness conditions increases the approximation errors, it may be better for the visualization of the solution.

Finally, we introduce another degree of freedom, the degree of the spline functions. Let us again consider bivariate spline spaces first. Let  $\mathbf{d} = \{d_t, t \in \Delta\}$  be a set of integers  $d_t \geq 0$  associated with triangles of  $\Delta$ . Let  $\mathbf{r}$  and  $\rho$  as above. Define

$$S_{\mathbf{d}}^{\mathbf{r},\rho}(\Delta) = \{s \in S_{\mathbf{d}}^{\mathbf{r},\rho}(\Delta), s|_t \in \mathbb{P}_{d_t}, t \in \Delta\}$$

be the spline space of variable smoothness  $\rho$ ,  $\mathbf{r}$  and variable degree  $\mathbf{d}$  associated with the vertices, interior edges and triangles of  $\Delta$ , where  $\mathbb{P}_{d_t}$  stands for the space of polynomials of degree  $d_t$ . This is a user-friendly spline space allowing one to choose a spline function using polynomials of less degree in certain areas and higher degree in other areas. It is especially useful to trim off the oscillations of interpolatory surfaces.

## §2. B-form Representation of Multivariate Splines

Let  $t = \langle v^{(0)}, \dots, v^{(n)} \rangle \in \mathbb{R}^n$  be an  $n$ -simplex with  $n+1$  distinct points  $v^{(k)}, k = 0, 1, \dots, n$ . Suppose that the  $n$ -simplex  $t$  has nonzero volume. Then for any point  $x \in \mathbb{R}^n$ ,  $x - v^{(0)}$  can be uniquely expressed by a linear combination of  $v^{(i)} - v^{(0)}, i = 1, \dots, n$ . That is,

$$x = v^{(0)} + \sum_{i=1}^n \lambda_i (v^{(i)} - v^{(0)}).$$

Let  $\lambda_0 = 1 - \sum_{i=1}^n \lambda_i$ . Then the  $(n+1)$ -tuple  $(\lambda_0, \lambda_1, \dots, \lambda_n)$  is called the barycentric coordinate of  $x$  with respect to  $t$ . It is easy to see that each  $\lambda_i$  is a linear function of  $x$ . Next let  $\mathbf{Z}^{n+1}$  be the set of all multi-integers in  $\mathbb{R}^{n+1}$ . For a multi-integer  $\alpha = (\alpha_0, \dots, \alpha_n) \in \mathbf{Z}^{n+1}$  with  $|\alpha| = \alpha_0 + \dots + \alpha_n \geq 0$ , let

$$B_{\alpha}^t(x) := \frac{|\alpha|!}{\alpha!} \lambda^{\alpha},$$

where  $\alpha! = \alpha_0! \dots \alpha_n!$  and

$$\lambda^{\alpha} = \prod_{i=0}^n \lambda_i^{\alpha_i}.$$

Then it is clear that  $B_{\alpha}^t(x)$  is a polynomial of degree  $|\alpha|$  in  $x$ . It can be shown that  $\{B_{\alpha}^t(x), \alpha \in \mathbf{Z}^{n+1}, |\alpha| = d\}$  forms a basis for polynomials of degree  $\leq d$  (cf. [deBoor'87]). Thus, any polynomial  $p$  of total degree  $d$  may be written in terms of  $B_{\alpha}^t(x)$ 's as

$$p(x) = \sum_{|\alpha|=d} c_{\alpha}^t B_{\alpha}^t(x) \quad (1)$$

for some coefficients  $c_\alpha^t$ 's depending on  $t$ . Thus, any spline function  $s$  is given by

$$s(x) = \sum_{|\alpha|=d} c_\alpha^t B_\alpha^t(x), \quad x \in t \in \Delta \quad (2)$$

with B-coefficient vector  $\{c_\alpha^t, |\alpha| = d, t \in \Delta\}$  of length  $\hat{d}T$ , where  $T$  denotes the number of  $n$ -simplices in  $\Delta$  and

$$\hat{d} = \binom{d+n}{n}.$$

This representation of the spline function  $s$  is called the B-form of  $s$ . (Cf. [24] and [11].)

One simple property of the B-form of polynomials is:

**Lemma 1.** *Let  $t = \langle v^{(0)}, \dots, v^{(n)} \rangle$  be an  $n$ -simplex in  $\mathbb{R}^n$  and let  $p(x)$  be a polynomial of degree  $d$  given in B-form (1) with respect to  $t$ . Then*

$$p(v^{(k)}) = c_{d\mathbf{e}^k}^t, \quad \forall \quad 0 \leq k \leq n,$$

where  $\mathbf{e}^k = (0, \dots, 0, 1, 0, \dots, 0)$  with 1 appearing in the  $(k+1)^{th}$  place.

To evaluate  $p(x)$  in B-form (1), we use the so-called de Casteljau algorithm. The derivative of  $p(x)$  in B-form can be given in B-form again. The integration of a polynomial  $p$  in B-form is a sum of all coefficients of  $p$  with multiplication by an appropriate constant. See, e.g., [17] for all these properties. Another important property is the following Markov inequality:

**Lemma 2.** *Let  $1 \leq q \leq \infty$ . There exists a constant  $N$  depending only on  $d$  such that*

$$\frac{\|\{c_\alpha^t, |\alpha| = d\}\|_q}{N} \leq \|p\|_{q,t} \leq \|\{c_\alpha^t, |\alpha| = d\}\|_q.$$

for any polynomial  $p(x) = \sum_{|\alpha|=d} c_\alpha^t B_\alpha^t(x)$ , where  $\|p\|_{q,t}$  denotes the standard  $L_q$  norm over the  $n$ -simplex  $t$  and  $\|\{c_\alpha^t, |\alpha| = d\}\|_q$  denotes the  $\ell_q$  norm of the sequence  $\{c_\alpha^t, |\alpha| = d\}$ .

We refer the interested reader to [42] for a proof in the bivariate setting which can be generalized to the multivariate setting easily.

Next we look at the smoothness conditions. Let

$$t_1 = \langle v^{(0)}, \dots, v^{(k)}, v^{(k+1)}, \dots, v^{(n)} \rangle$$

and

$$t_2 = \langle v^{(0)}, \dots, v^{(k)}, u^{(k+1)}, \dots, u^{(n)} \rangle$$

be two  $n$ -simplices in  $\mathbb{R}^n$  and  $\tilde{t} = \langle v^{(0)}, \dots, v^{(k)} \rangle$  the  $k$ -simplex which is a common facet of  $t_1$  and  $t_2$ , with  $0 \leq k < n$ . Let  $F$  be a function defined on  $t_1 \cup t_2$  by

$$F(x) = \begin{cases} p_d(x) = \sum_{|\alpha|=n} a_\alpha B_\alpha^{t_1}(x), & \text{if } x \in t_1 \\ q_d(x) = \sum_{|\alpha|=n} b_\alpha B_\alpha^{t_2}(x), & \text{if } x \in t_2. \end{cases}$$

Let us assume that  $F$  is well defined on  $\tilde{t}$ . Writing  $u^{(j)} = \sum_{i=0}^n c_{ji} v^{(i)}$ ,  $j = k+1, \dots, n$ , we have the following:

**Theorem 3.** *Suppose that  $t_1$  and  $t_2$  are two  $n$ -simplices such that  $\tilde{t} = t_1 \cap t_2$  is a  $(n-1)$ -simplex in  $\mathbb{R}^n$ . Let  $F$  be the function defined above. Then  $F \in C^r(t_1 \cup t_2)$  if and only if the following conditions hold*

$$b_{(\alpha_0, \dots, \alpha_{n-1}, \ell)} = \sum_{|\gamma|=\ell} a_{(\alpha_0, \dots, \alpha_{n-1}, 0) + \gamma} B_\gamma^{t_1}(u^{(n)}). \quad (3)$$

for  $0 \leq \ell \leq r$ .

This is the well-known smoothness conditions (cf. [24] and [11]). For the geometric meaning of the smoothness conditions in the bivariate setting, see [39]. Next we look at the degree reduction conditions. These conditions allow us to constrain the spline function to be of variable degree over the simplices. Let

$$\Delta_{ij} c_\alpha = c_{\alpha+e_i} - c_{\alpha+e_j}$$

be a difference operator, where  $e_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbf{Z}^{n+1}$  with 1 in the  $i^{\text{th}}$  entry and similar for  $e_j$ . Inductively, let

$$\Delta_{ij}^k = \Delta_{ij}(\Delta_{ij}^{k-1})$$

for  $k \geq 2$ . For any multi-integer  $\beta = (\beta_1, \dots, \beta_n)$ , let

$$\Delta^\beta = \Delta_{10}^{\beta_1} \dots \Delta_{n0}^{\beta_n}$$

be a difference operator of order  $|\beta|$ . Our degree reduction conditions are:

**Theorem 4.** *Let  $p = \sum_{|\alpha|=d} c_\alpha B_\alpha^t$  be a polynomial of degree  $d$  in  $B$ -form with respect to  $t$ . Then  $p$  is a polynomial of degree  $d_t < d$  if*

$$\Delta^\beta c_\alpha = 0, \quad d_t < |\beta| \leq d, \quad |\alpha| = d - |\beta|, \quad (4)$$

where  $\Delta^\beta c_\alpha = \Delta^\beta c_{|\alpha}$ , that is, the difference operators are applied first before the evaluation at the index  $\alpha$ .

The conditions can be verified easily and are left to the interested reader. It is easy to see that both conditions (3) and (4) are linear relations among the B-coefficients of polynomials.

Let us summarize the discussions above as follows: For each spline function in

$$\mathcal{S} := S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta) \quad (5)$$

the spline space of smoothness  $\mathbf{r}$ , super smoothness  $\rho$  and degree  $\mathbf{d}$  for three fixed sequences  $\rho, \mathbf{r}$ , and  $\mathbf{d}$  associated with the  $k$ -simplices with  $0 \leq k < n - 1$ , interior  $n - 1$  simplices, and  $n$ -simplices of  $\Delta$ , we write

$$s = \sum_{t \in \Delta} \sum_{|\alpha|=d_t} c_\alpha^t \mathcal{B}_\alpha^t, \quad (6)$$

with  $\mathbf{c} = (c_\alpha^t, |\alpha| = d_t, t \in \Delta) \in \mathbb{R}^N$ ,  $N = \sum_{t \in \Delta} \hat{d}_t$  with  $\hat{d}_t = \binom{d_t+n}{n}$  and

$$\mathcal{B}_\alpha^t(x) = \begin{cases} B_\alpha^t(x), & \text{if } x \in t \\ 0, & \text{if } x \in \Delta \setminus \{t\}. \end{cases}$$

In addition,  $\mathbf{c}$  satisfies the constraints  $H\mathbf{c} = 0$  for the smoothness conditions that  $\mathcal{S}$  possesses and  $D\mathbf{c} = 0$  for the degree reduction conditions.

### §3. Multivariate Spline Method and a Matrix Iterative Method

Let  $\mathcal{S}$  be the spline space defined in (5) and let

$$J(u) = \frac{1}{2}a(u, u) - b(u)$$

define a functional on  $\mathcal{S}$  where  $a$  is a continuous bilinear form and  $b$  a continuous linear functional. We are interested in subsets of  $\mathcal{S}$  satisfying additional constraints  $L(u) = G$  such as boundary conditions and/or interpolation conditions.

We consider the following abstract problem: Find  $s_G \in \mathcal{S}$  such that

$$J(s_G) = \min\{J(s) : s \in \mathcal{S}, L(s) = G\}.$$

In this section, we explain how to solve this minimization problem.

We first denote by  $A = (a(\mathcal{B}_\alpha^t, \mathcal{B}_\beta^{t'}))_{\substack{|\alpha|=d, t \in \Delta \\ |\beta|=d, t' \in \Delta}}$  the matrix associated with the bilinear form  $a$  and  $F = (b(\mathcal{B}_\alpha^t))_{|\alpha|=d, t \in \Delta}$  the vector associated with linear functional  $b$ . Similarly, we may also express the side conditions  $L(u) = G$  in matrix form:

$$L\mathbf{c} = \mathbf{g},$$

where we have to approximate the constraints  $L(u) = G$  if  $B(u)$  is a nonlinear functional.

Using the B-form (6) of spline functions, the above abstract problem can be expressed in the following form:

$$\begin{aligned} \min J(\mathbf{c}) &= \frac{1}{2}\mathbf{c}^T A\mathbf{c} - \mathbf{c}^T F \\ \text{subject to } H\mathbf{c} &= 0, D\mathbf{c} = 0, L\mathbf{c} = \mathbf{g}. \end{aligned} \quad (7)$$

By the theory of Lagrange multipliers, letting

$$\mathcal{L}(\mathbf{c}, \lambda_1, \lambda_2, \lambda_3) = \frac{1}{2}\mathbf{c}^T A\mathbf{c} - \mathbf{c}^T F + \lambda_1^T H\mathbf{c} + \lambda_2 D\mathbf{c} + \lambda_3^T (L\mathbf{c} - \mathbf{g}),$$

there exist  $\lambda_1, \lambda_2, \lambda_3$  such that

$$\begin{aligned} A\mathbf{c} + H^T \lambda_1 + D^T \lambda_2 + L^T \lambda_3 &= F \\ H\mathbf{c} &= 0 \\ D\mathbf{c} &= 0 \\ L\mathbf{c} &= \mathbf{g}. \end{aligned} \quad (8)$$

In general, the linear system above is not invertible. In particular,  $A$  may be singular. Thus, we can not solve it directly. However, we can solve it using a least squares method. Indeed, assuming the existence and uniqueness of the solution of  $\mathbf{c}$ , any least squares solution will satisfy the above linear system exactly and the part of solution corresponding to  $\mathbf{c}$  is the solution for (7). For linear systems of small size, the least squares method works well according to our experiments. However, we do not recommend it when the size of this system is very large.

Next we present an iterative algorithm to solve the linear systems of large size. For simplicity, we write the Lagrange equations in the following matrix form:

$$\begin{bmatrix} B^T & A \\ 0 & B \end{bmatrix} \begin{bmatrix} \lambda \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} F \\ G \end{bmatrix} \quad (9)$$

with appropriate matrices  $B$  and  $G$ . This is a typical saddle point problem. There is a vast literature on the numerical solution to (9) (cf., e.g., [6], [53], [21], [8], [9]). All the methods discussed in the references above require  $A$  to be nonsingular. However, in our situation, the matrix  $A$  is singular. Nevertheless, we require  $A$  positive definite on the kernel of  $B$ . As mentioned in the introduction, the matrix iterative method described below is a variant of the augmented Lagrangian algorithm, c.f. [4]. Unlike in [Brenner and Scott'94], we do not require  $B^T$  to be injective, that is we allow non uniqueness of the multiplier  $\lambda$ . Consider the following sequence of problems:

$$\begin{bmatrix} B^T & A \\ -\epsilon I & B \end{bmatrix} \begin{bmatrix} \lambda^{(k+1)} \\ \mathbf{c}^{(k+1)} \end{bmatrix} = \begin{bmatrix} F \\ G - \epsilon\lambda^{(k)} \end{bmatrix} \quad (10)$$

for  $k = 0, 1, \dots$ , with an initial guess  $\lambda^{(0)}$ , e.g.,  $\lambda^{(0)} = 0$ , and  $I$  the identity matrix of  $\mathbb{R}^m$ , assuming the size of  $B$  is  $m \times n$ . Note that (10) reads

$$\begin{aligned} A\mathbf{c}^{(k+1)} + B^T\lambda^{(k+1)} &= F \\ B\mathbf{c}^{(k+1)} - \epsilon\lambda^{(k+1)} &= G - \epsilon\lambda^{(k)}. \end{aligned} \quad (11)$$

Multiplying on the left of the second equation in (11) by  $B^T$ , we get

$$B^T B\mathbf{c}^{(k+1)} - \epsilon B^T\lambda^{(k+1)} = B^T G - \epsilon B^T\lambda^{(k)}$$

or  $B^T\lambda^{(k+1)} = \frac{1}{\epsilon}B^T B\mathbf{c}^{(k+1)} - \frac{1}{\epsilon}B^T G + B^T\lambda^{(k)}$  and substitute it into the first equation in (11) to get

$$\left(A + \frac{1}{\epsilon}B^T B\right)\mathbf{c}^{(k+1)} = F + \frac{1}{\epsilon}B^T G - B^T\lambda^{(k)}. \quad (12)$$

It follows that

$$\left(A + \frac{1}{\epsilon}B^T B\right)\mathbf{c}^{(1)} = F + \frac{1}{\epsilon}B^T G - B^T\lambda^{(0)}.$$

Using the first equation in (11), i.e.,  $A\mathbf{c}^{(k)} = F - B^T\lambda^{(k)}$  to replace  $F$  in (12), we have

$$\left(A + \frac{1}{\epsilon}B^T B\right)\mathbf{c}^{(k+1)} = A\mathbf{c}^{(k)} + \frac{1}{\epsilon}B^T G$$

for all  $k \geq 1$ . This suggests the following:

**Algorithm 5.** Let  $I$  be the identity matrix of  $\mathbb{R}^m$ . Fix  $\epsilon > 0$ . Given an initial guess  $\lambda^{(0)} \in \text{Im}(B)$ , we first compute

$$\mathbf{c}^{(1)} = \left(A + \frac{1}{\epsilon}B^T B\right)^{-1} \left(F + \frac{1}{\epsilon}B^T G - B^T\lambda^{(0)}\right)$$

and iteratively compute

$$\mathbf{c}^{(k+1)} = \left(A + \frac{1}{\epsilon}B^T B\right)^{-1} \left(A\mathbf{c}^{(k)} + \frac{1}{\epsilon}B^T G\right) \quad (13)$$

for  $k = 1, 2, \dots$ , where  $\text{Im}(B)$  is the range of  $B$ .

We remark that this is not a brand new method. It was discussed in [Gunzburger'89] briefly. But no proof of the convergence was given there, but see [4] for a linear convergence rate proof. We now present another proof which shows that  $\mathbf{c}^{(k)}$  converges to the  $\mathbf{c}$ .

**Theorem 6.** *Let us assume that the linear system (9) has a solution  $(\lambda, \mathbf{c})$  with  $\mathbf{c}$  unique. Assume also that  $A$  is symmetric positive definite with respect to  $B$ , that is,  $\mathbf{x}^T A \mathbf{x} \geq 0$  and  $\mathbf{x}^T A \mathbf{x} = 0$  with  $B \mathbf{x} = 0$  imply that  $\mathbf{x} = 0$ . Then there exists a constant  $C_1(\epsilon)$  depending on  $\epsilon$  but independent of  $k$  such that*

$$\|\mathbf{c} - \mathbf{c}^{(k+1)}\| \leq C_1(\epsilon) \left( \frac{C_2 \epsilon}{1 + C_2 \epsilon} \right)^{k+1}$$

for  $k \geq 1$ , where  $C_2 = \|B^+\|^2 \|A\|$  and  $B^+$  stands for the pseudo inverse of  $B$ .

**Proof:** We first show that  $E = A + \frac{1}{\epsilon} B^T B$  is invertible. If  $E \mathbf{x} = 0$ , we have

$$\mathbf{x}^T E \mathbf{x} = \mathbf{x}^T A \mathbf{x} + \frac{1}{\epsilon} \|B \mathbf{x}\|^2 = 0$$

which implies that  $\mathbf{x}^T A \mathbf{x} = 0$  and  $B \mathbf{x} = 0$  since  $A$  is nonnegative definite. Since  $A$  is also positive definite with respect to  $B$ , we have  $\mathbf{x} = 0$ . Thus,  $E$  is invertible and hence the sequence  $\{\mathbf{c}^{(k)}\}$  is well-defined. Let  $C_1 = \|E^{-1}\|$  which is obviously dependent on  $\epsilon$ . But our numerical experiments show that  $C_1$  has an upper bound independent of  $\epsilon$ .

Next we show the convergence. We divide the proof into several steps.

Step 1. From (12), we have

$$E \mathbf{c}^{(k+1)} = F + \frac{1}{\epsilon} B^T G - B^T \lambda^{(k)}.$$

Similarly, we use (9) to get

$$E \mathbf{c} = F + \frac{1}{\epsilon} B^T G - B^T \lambda.$$

Therefore, we have

$$\mathbf{c}^{(k+1)} - \mathbf{c} = E^{-1} B^T (\lambda - \lambda^{(k)}). \quad (14)$$

Step 2. Using the second equation in (11) and (12), we have

$$-\epsilon(\lambda^{(k+1)} - \lambda) = -\epsilon(\lambda^{(k)} - \lambda) + G - B \mathbf{c}^{(k+1)}$$

and

$$\mathbf{c}^{(k+1)} = E^{-1} F + \frac{1}{\epsilon} E^{-1} B^T G - E^{-1} B^T \lambda^{(k)}.$$

It follows that

$$\begin{aligned} \epsilon(\lambda^{(k+1)} - \lambda) &= \epsilon(\lambda^{(k)} - \lambda) - G + B E^{-1} F \\ &\quad + \frac{1}{\epsilon} B E^{-1} B^T G - B E^{-1} B^T \lambda^{(k)} \\ &= \epsilon(\lambda^{(k)} - \lambda) - G + B E^{-1} (F + \frac{1}{\epsilon} B^T G - B^T \lambda^{(k)}). \end{aligned}$$

Because  $F = A\mathbf{c} + B^T\lambda$  and  $G = B\mathbf{c}$ , we have

$$\begin{aligned} F + \frac{1}{\epsilon}B^TG &= A\mathbf{c} + B^T\lambda + \frac{1}{\epsilon}B^TB\mathbf{c} \\ &= \left(A + \frac{1}{\epsilon}B^TB\right)\mathbf{c} + B^T\lambda \\ &= E\mathbf{c} + B^T\lambda. \end{aligned}$$

Consequently, we obtain

$$\begin{aligned} \epsilon(\lambda^{(k+1)} - \lambda) &= \epsilon(\lambda^{(k)} - \lambda) - G + BE^{-1}(E\mathbf{c} + B^T\lambda - B^T\lambda^{(k)}) \\ &= \epsilon(\lambda^{(k)} - \lambda) - D(\lambda^{(k)} - \lambda) \end{aligned}$$

where  $D = BE^{-1}B^T$ . That is, we have

$$\lambda^{(k+1)} - \lambda = \left(I - \frac{1}{\epsilon}D\right)(\lambda^{(k)} - \lambda). \quad (15)$$

Step 3. We also need the following

**Lemma 7.** *Let  $\text{Ker}(B)$  be the kernel of  $B$  and  $\text{Im}(B)$  the range of  $B$ . Then  $\mathbb{R}^m = \text{Ker}(B^T) \oplus \text{Im}(B)$ . Similarly,  $\mathbb{R}^n = \text{Ker}(B) \oplus \text{Im}(B^T)$ .*

**Proof:** We write  $\mathbb{R}^m = \text{Im}(B) \oplus \text{Im}(B)^\perp$ . For any  $v \in \text{Im}(B)^\perp$ , we have

$$0 = \langle Bu, v \rangle = \langle u, B^Tv \rangle$$

for any  $u \in \mathbb{R}^n$ . Thus,  $B^Tv = 0$  or  $v \in \text{Ker}(B^T)$ . Equivalently, if  $v \in \text{Ker}(B^T)$ , then  $v \in \text{Im}(B)^\perp$ . Therefore  $\mathbb{R}^m = \text{Ker}(B^T) \oplus \text{Im}(B)$ . Similar for the decomposition of  $\mathbb{R}^n$ .  $\square$

Using Lemma 7, we may assume that  $\lambda \in \text{Im}(B)$  since the component of  $\lambda$  in  $\text{Ker}(B^T)$  will not make any contribution in (9). It follows from the second equation in (11) that

$$B(\mathbf{c}^{(k+1)} - \mathbf{c}) = \epsilon(\lambda^{(k)} - \lambda^{(k+1)}).$$

That is,  $\lambda^{(k)} - \lambda^{(k+1)}$  is in the  $\text{Im}(B)$ . Since

$$\lambda^{(k)} - \lambda = \sum_{j=1}^k (\lambda^{(j)} - \lambda^{(j-1)}) + (\lambda^{(0)} - \lambda),$$

we have  $\lambda^{(k)} - \lambda \in \text{Im}(B)$  for each  $k$ . From (15), we need to estimate the norm of  $I - \frac{1}{\epsilon}D$  restricted to  $\text{Im}(B)$  in order to estimate the norm of

$\lambda^{(k+1)} - \lambda$ . We write  $\|I - \frac{1}{\epsilon}D\|$  for  $\left\| \left(I - \frac{1}{\epsilon}D\right) \Big|_{\text{Im}(B)} \right\|$  and we have:

$$\|\lambda^{(k+1)} - \lambda\| \leq \left\| \left(I - \frac{1}{\epsilon}D\right) \Big|_{\text{Im}(B)} \right\| \|\lambda^{(k)} - \lambda\|.$$

Step 4. We claim that

$$\|I - \frac{1}{\epsilon}D\| \leq \frac{C_2\epsilon}{1 + C_2\epsilon}, \quad (16)$$

for some constant  $C_2 > 0$ . Indeed, by the Rayleigh-Ritz quotient, we have

$$\|I - \frac{1}{\epsilon}D\| = \max_{0 \neq v \in \text{Im}(B)} \left(1 - \frac{1}{\epsilon} \frac{v^T D v}{v^T v}\right). \quad (17)$$

We now use a technique from [54] to prove that

$$R(v) = \frac{v^T D v}{\epsilon v^T v} > \frac{1}{1 + C_2\epsilon}, \quad \forall v \in \text{Im}(B), \quad v \neq 0.$$

We have

$$\begin{aligned} R(v) &= \frac{v^T B E^{-1} B^T v}{\epsilon v^T v} = \frac{v^T B E^{-1} E E^{-1} B^T v}{\epsilon v^T v} \\ &= \frac{(v^T B E^{-1}) E (E^{-1} B^T v)}{\epsilon v^T v} = \frac{\|E^{-1} B^T v\|_E^2}{\epsilon v^T v}, \end{aligned}$$

where we have used a norm  $\|\cdot\|_E$  associated with the positive definite matrix  $E$ ; That is,

$$\|u\|_E = b(u, u) \quad \text{with} \quad b(u, v) = v^T E u.$$

Next we have

$$\begin{aligned} \|E^{-1} B^T v\|_E &= \sup_{u \neq 0} \frac{b(E^{-1} B^T v, u)}{\|u\|_E} \\ &\geq \frac{u^T E E^{-1} B^T v}{\|u\|_E} = \frac{u^T B^T v}{\|u\|_E} = \frac{v^T B u}{\|u\|_E}, \quad \forall u \neq 0. \end{aligned}$$

It follows from Lemma 7 that since  $v \in \text{Im}(B)$ , there is  $u_v$  in  $\text{Im}(B^T)$  such that  $v = B u_v$ . Then

$$\begin{aligned} R(v) &\geq \frac{1}{\epsilon \|B u_v\|^2} \left( \frac{(B u_v)^T (B u_v)}{\|u_v\|_E} \right)^2 = \frac{1}{\epsilon} \frac{\|B u_v\|^2}{\|u_v\|_E^2} \\ &= \frac{\|B u_v\|^2}{\epsilon u_v^T A u_v + u_v^T B^T B u_v} \geq \frac{\|B u_v\|^2}{\epsilon \|u_v\|^2 \|A\| + \|B u_v\|^2} \geq \frac{1}{\epsilon \frac{\|u_v\|^2 \|A\|}{\|B u_v\|^2} + 1}. \end{aligned}$$

Since  $B$  maps  $\text{Im}(B^T)$  into  $\text{Im}(B)$ ,  $B^+$  is defined from  $\text{Im}(B)$  into  $\text{Im}(B^T)$  and

$$\frac{\|u_v\|}{\|B u_v\|} \leq \sup_{B u \neq 0} \frac{\|u\|}{\|B u\|} = \sup_{0 \neq v \in \text{Im}(B)} \frac{\|B^+ v\|}{\|v\|} = \|B^+\|.$$

So

$$R(v) \geq \frac{1}{\epsilon \|B^+\|^2 \|A\| + 1} \geq \frac{1}{C_2 \epsilon + 1},$$

if we let  $C_2 = \|B^+\|^2 \|A\|$ . It follows from (17) that

$$\|I - \frac{1}{\epsilon} D\| \leq 1 - \min_{0 \neq v \in \text{Im}(B)} R(v) \leq \frac{C_2 \epsilon}{1 + C_2 \epsilon}$$

which is (16).

Finally we get

$$\|\lambda^{(k+1)} - \lambda\| \leq \left( \frac{C_2 \epsilon}{C_2 \epsilon + 1} \right) \|\lambda^{(k)} - \lambda\|$$

from (15) and (16) and

$$\|\mathbf{c}^{(k+1)} - \mathbf{c}\| \leq \|E^{-1}\| \|B\| \left( \frac{C_2 \epsilon}{C_2 \epsilon + 1} \right)^{k+1} \|\lambda^{(0)} - \lambda\|$$

from (14), which is the desired result.  $\square$

In the following applications, we only need to verify that the matrix  $A$  is symmetric and positive definite with respect to the side condition matrix block  $B$ . It turns out that existence and uniqueness will guarantee it. Hence, the matrix iterative method will be well-defined in each subsequent application.

#### §4. Minimal Energy Method for Scattered Data Interpolation

Let  $\{(x^{(i)}, f_i), i = 1, \dots, V\}$  be a given set of scattered data in  $\mathbb{R}^n$ . Assume that the  $x^{(i)}$ 's are distinct. A general problem is to find a smooth surface  $s_f$  which interpolates the given data:

$$s_f(x^{(i)}) = f_i, \quad i = 1, \dots, V. \quad (18)$$

For many applications,  $s$  has to be a smooth surface such as a  $C^r$  surface for a fixed integer  $r \geq 1$ . Depending on particular applications,  $s$  sometimes has to be smoother in certain regions of the same domain. This requires spline functions of variable smoothness.

We shall discuss a popular method called the minimal energy method to find such an interpolating surface. Let  $\Delta$  be a triangulation of the given data locations  $\{x_i, i = 1, \dots, V\}$ . For example, we can use the well-known Delaunay method to find such a triangulation. Let

$$\mathcal{S} := S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$$

be a spline space of smoothness  $\mathbf{r}$ , super smoothness  $\rho$  and degree  $\mathbf{d}$  for three fixed sequences  $\rho$ ,  $\mathbf{r}$ , and  $\mathbf{d}$  associated with  $k$ -simplices,  $0 \leq k < n-1$ , interior  $n-1$  simplices, and  $n$ -simplices of  $\Delta$ . Let  $\Lambda(f)$  be the set of interpolating splines in  $\mathcal{S}$ , i.e.,

$$\Lambda(f) := \{s \in \mathcal{S}, s(x^{(i)}) = f_i, i = 1, \dots, V\}. \quad (19)$$

Assume that the global smoothness  $\min\{\mathbf{r}, \rho\}$  of the spline space is at least 1. Our minimal energy method is to find  $s_f \in \mathcal{S}$  satisfying the interpolation condition (18), i.e.  $s_f \in \Lambda(f)$  and minimizing the energy functional

$$E(s) := \sum_{t \in \Delta} \int_t \left( \sum_{|\beta|=2} (D^\beta s)^2 \right) dx. \quad (20)$$

Some extensions of the minimal energy method will be given later in the section. The convergence of the minimal energy method for bivariate spline interpolation was given in [27]. Our first result is about the existence and uniqueness of such a  $s_f$ . We prefer to give an elementary proof of the following:

**Theorem 8.** *Suppose that  $\Lambda(f)$  is not empty. Then there exists a unique  $s_f \in \Lambda(f)$  minimizing (20).*

**Proof:** We first prove this theorem in the bivariate setting. Any spline function  $s \in \mathcal{S}$  can be written

$$s(x, y)|_t = \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x, y), \quad (x, y) \in t \in \Delta.$$

Let  $\mathbf{c} = (c_{ijk}^t, i+j+k=d, t \in \Delta)$  be the coefficient vector associated with  $s$ . The length of the vector  $\mathbf{c}$  is  $\hat{d}T$  with  $T$  being the number of triangles in  $\Delta$  and  $\hat{d} = (d+1)(d+2)/2$ . Since  $s \in \mathcal{S}$ ,  $s$  satisfies the smoothness and super smoothness conditions which can be expressed by a linear system  $H\mathbf{c} = 0$  using the smoothness conditions (3). Also,  $s$  satisfies the degree reduction conditions which can also be given by another linear system  $D\mathbf{c} = 0$  by using (4). The energy functional  $E(s)$  can be written in terms of  $\mathbf{c}$  as

$$E(s) = \mathbf{c}^T K \mathbf{c},$$

where  $K = \text{diag}(K_t, t \in \Delta)$  is a block diagonal matrix with

$$K_t = \left[ \int_t (D_x^2 B_{ijk}^t D_x^2 B_{p,q,r}^t + 2D_x D_y B_{ijk}^t D_x D_y B_{p,q,r}^t + D_y^2 B_{ijk}^t D_y^2 B_{p,q,r}^t) dx dy \right]_{\substack{i+j+k=d \\ p+q+r=d}}.$$

Since  $\Lambda(f)$  is not empty, choose an  $s_0 \in \Lambda(f)$ . Then let

$$A = \{\mathbf{c} \in \mathbb{R}^{\hat{d}T} : E(s) = \mathbf{c}^T K \mathbf{c} \leq E(s_0)\}.$$

We first show that  $A$  is a bounded and closed set. Clearly, we have

$$\int_t |D_x^2 s(x, y)|^2 dx dy \leq E(s_0),$$

for any triangle  $t \in \Delta$ . Since  $D_x^2 s(x, y)|_t$  is in the space of polynomials of degree  $\leq d - 2$ , it follows that

$$\max_{(x, y) \in t} |D_x^2 s(x, y)| \leq CE(s_0)$$

for some  $C$  depending on the triangle  $t$ . Similar for  $D_x D_y s$  and  $D_y^2 s$ . Thus, we have  $|s|_{2, \infty, t} \leq CE(s_0)$  for another constant  $C > 0$  depending on  $t$ , where

$$|s|_{2, \infty, t} := \max_{(x, y) \in t} \max\{|D_x^2 s(x, y)|, |D_x D_y s(x, y)|, |D_y^2 s(x, y)|\}.$$

Next we need to show that  $\max\{|s(x, y)|, (x, y) \in t\}$  is bounded. To this end, we write  $v_i = (x_i, y_i)$  and  $f(v_i) = f_i$  for  $i = 1, \dots, V$ . For convenience, we consider  $t = \langle v_1, v_2, v_3 \rangle$  and  $v = (x, y) \in t$ . Using the Taylor expansion, we have

$$f(v_i) = s(v_i) = s(v) + \nabla s(v) \cdot (v_i - v) + O(|s|_{2, \infty, t} |t|^2) \quad (21)$$

for  $i = 1, 2, 3$ , where  $|t|$  denote the diameter of  $t$ . It follows that

$$\begin{aligned} f(v_2) - f(v_1) &= \nabla s(v) \cdot (v_2 - v_1) + O(|s|_{2, \infty, t} |t|^2) \\ f(v_3) - f(v_1) &= \nabla s(v) \cdot (v_3 - v_1) + O(|s|_{2, \infty, t} |t|^2). \end{aligned}$$

Solving this linear system for  $\nabla s(v) = (D_x s(v), D_y s(v))$  gives

$$\begin{aligned} D_x s(v) &= O(|t|^3 |s|_{2, \infty, t} / A_t) + |f(v_2) - f(v_1)| |t| / A_t \\ D_y s(v) &= O(|t|^3 |s|_{2, \infty, t} / A_t) + |f(v_3) - f(v_1)| |t| / A_t, \end{aligned}$$

where  $A_t$  stands for the area of the triangle  $t$ . Inserting these estimates for  $\nabla s(v)$  in (21), we immediately get

$$|s(v)| \leq C \left( (1 + |t|^2 / A_t) \|\mathbf{f}\|_\infty + |t|^4 |s|_{2, \infty, t} / A_t \right).$$

where  $\|\mathbf{f}\|_\infty = \max\{|f_\ell|, \ell = 1, \dots, V\}$ . Hence,

$$|s(x, y)| \leq C_t (\|\mathbf{f}\|_\infty + E(s_0))$$

for another constant  $C_t > 0$  depending on  $t$ . It follows that

$$\max_{i+j+k=d} |c_{ijk}^t| \leq N \max_{(x,y) \in t} |s(x,y)| \leq NC_t (\|\mathbf{f}\|_\infty + E(s_0))$$

by Lemma 2 and hence,

$$\|\mathbf{c}\|_\infty := \max_{\substack{i+j+k=d \\ t \in \Delta}} |c_{ijk}^t| \leq \max_{t \in \Delta} NC_t (\|\mathbf{f}\|_\infty + E(s_0)).$$

Thus,  $A$  is a bounded set. It is clear that  $A$  is a closed set since  $E(s) = \mathbf{c}^T K \mathbf{c}$  is a continuous function of  $\mathbf{c}$ . It follows that  $E(s)$  achieves its minimum in  $A$ . Let  $\mathbf{c}_f$  be a minimizer.

We now show that the minimizer is unique. It is clear that  $K$  is nonnegative definite and  $E(s) = \mathbf{c}^T K \mathbf{c} = 0$  if and only if  $s$  is a linear function. By the interpolation conditions, the values  $f_i, i = 1, \dots, V$  have to be obtained from a linear function. In this case,  $s$  is the unique solution. If the  $f_i$ 's are not the values of a linear function, then  $E(s) = \mathbf{c}^T K \mathbf{c} > 0$  for any  $\mathbf{c} \in \Lambda(f)$ . Hence,  $E(s)$  is a strictly convex function. Thus, the minimizer is unique. Hence, we have completed the proof.

We note that the above arguments can be generalized to the multivariate setting and the existence and uniqueness of the interpolatory spline of minimal energy follow similarly.  $\square$

To solve the minimal energy interpolation problem, we first note that it is equivalent to the following constrained minimization problem:

$$\begin{aligned} & \min \mathbf{c}^T K \mathbf{c} \\ & \text{subject to} \\ & H\mathbf{c} = 0, I\mathbf{c} = \mathbf{f}, D\mathbf{c} = 0 \end{aligned} \tag{22}$$

where  $\mathbf{f} = (f_1, \dots, f_V)$  is the data value vector and  $I\mathbf{c} = \mathbf{f}$  is a linear system associated with the interpolation condition (18) since the spline value at a vertex of an  $n$ -simplex in  $\Delta$  is the same as the corresponding B-coefficient value (cf. Lemma 1). Also,  $D\mathbf{c} = 0$  denotes the degree reduction conditions.

This is a typical example of our model problem discussed in §3. We shall use the Lagrange multiplier method. Let

$$\mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma) := \mathbf{c}^T K \mathbf{c} + \alpha^T H \mathbf{c} + \beta^T D \mathbf{c} + \gamma^T (I \mathbf{c} - \mathbf{f})$$

be a Lagrangian function. We need to find a minimizer of  $\mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma)$ . That is,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{c}} \mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma) &= 0, \quad \frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma) = 0, \\ \frac{\partial}{\partial \beta} \mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma) &= 0, \quad \frac{\partial}{\partial \gamma} \mathcal{L}(\mathbf{c}, \alpha, \beta, \gamma) = 0. \end{aligned}$$

It follows that

$$\begin{aligned} 2\mathbf{c}^T K + \alpha^T H + \beta^T D + \gamma^T I &= 0, \\ H\mathbf{c} = 0, D\mathbf{c} = 0, I\mathbf{c} &= \mathbf{f}. \end{aligned}$$

In other words, we have

$$\begin{bmatrix} H^T & D^T & I^T & 2K \\ 0 & 0 & 0 & H \\ 0 & 0 & 0 & D \\ 0 & 0 & 0 & I \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \mathbf{f} \end{bmatrix}. \quad (23)$$

We shall apply the matrix iterative method for solving the above linear system. To establish the matrix iterative method, we need to show that  $2K$  is symmetric and positive definite with respect to  $[H; D; I]$ . Clearly,  $2K$  is symmetric. For any  $\mathbf{c}$  such that  $\mathbf{c}^T K \mathbf{c} = 0$  and  $H\mathbf{c} = 0, D\mathbf{c} = 0, I\mathbf{c} = 0$ , it follows from Theorem 8 above that  $\mathbf{c} = 0$  is the unique solution in  $\mathcal{S}$  satisfying the interpolation condition.

Let us make a remark on the nonemptiness of  $\Lambda(f)$ . It is possible to use the smoothness, degree reduction, and interpolation matrices to check if  $\Lambda(f) = \emptyset$  or not. That is, letting  $\mathbf{c}$  be a least squares solution of

$$\begin{bmatrix} H \\ D \\ I \end{bmatrix} \mathbf{c} = \begin{bmatrix} 0 \\ 0 \\ \mathbf{f} \end{bmatrix},$$

if  $\mathbf{c}$  solves the above linear system exactly, then  $\Lambda(f) \neq \emptyset$ . Furthermore, when we solve the minimal energy interpolation problem, we use the maximum norm

$$\left\| \begin{bmatrix} H\mathbf{c} \\ D\mathbf{c} \\ I\mathbf{c} - \mathbf{f} \end{bmatrix} \right\|_{\infty} \quad (24)$$

to measure if the least squares solution does satisfy the interpolation condition, smoothness conditions, and degree reduction conditions. For convenience, we will call such a measurement the exactness of the solution.

We have implemented the above algorithms for 2D and 3D scattered data interpolation using bivariate and trivariate splines of any degrees and variable smoothness over triangulations. Let us present some numerical experiments.

**Example 9.** Consider a set of scattered data locations used in [26] and triangulate it using Delaunay triangulation method. See Figure 1 for the triangulation  $\Delta$ . We use the following function

$$f(x, y) = \sin(\pi(x^2 + 2y^2))$$

evaluated at the data locations to have a set of scattered data. We use spline spaces of  $S_3^1(\Delta)$ ,  $S_4^1(\Delta)$ ,  $S_5^1(\Delta)$ , etc. to find interpolation surfaces and then compute the maximum errors of the interpolation against this function. Our maximum errors are computed based on  $101 \times 101$  equally-spaced points over  $[-0.05, 1.045] \times [-0.031, 1.051]$ . We tabulate the maximum errors and cpu times for computing the interpolation surfaces. The cpu times are based on a PC with 450Mhz. The exactness (24) has been checked and is less than  $10^{-8}$  for all the cases listed in the Table 1.

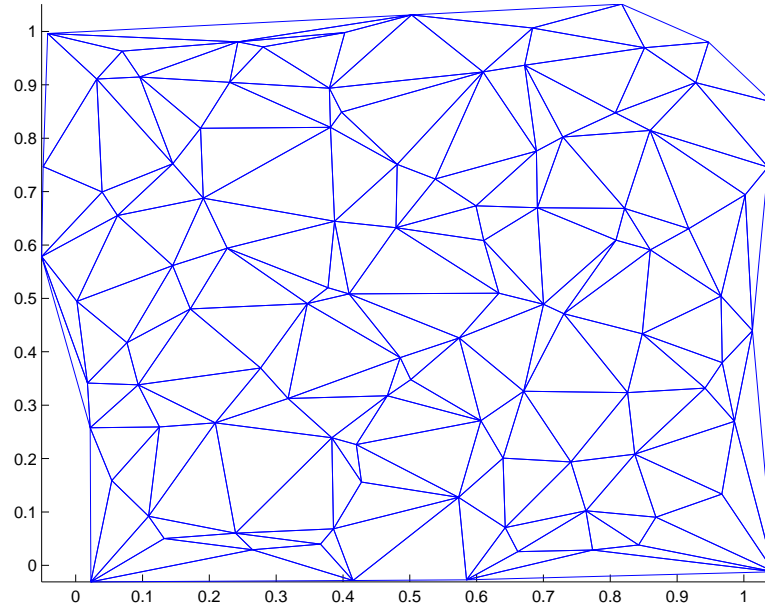


Fig. 1 A triangulation

Table 1. The Numerical Approximation using Minimal Energy Method

	Maximum Errors	CPU Times
$S_3^1(\Delta)$	0.22423758660928	16.6s
$S_4^1(\Delta)$	0.20516709129396	33.2s
$S_5^1(\Delta)$	0.19199158596566	63.5s
$S_5^2(\Delta)$	0.18141014563253	125.6s
$S_6^2(\Delta)$	0.19675675300293	212.4s

The minimal energy method for computing interpolation surface has many extensions and generalizations. We give three possible extensions here.

**Example 10.** Instead of the energy functional, we may minimize other positive functionals, e.g., triharmonic functional. Consider

$$E_3(s) = \sum_{t \in \Delta} \int_t \left[ \sum_{|\beta|=3} (D^\beta s)^2 \right] dx dy.$$

Then we formulate the following interpolation problem: find  $s_f \in \mathcal{S}$  with global smoothness  $\geq 2$  such that  $s_f$  satisfies the interpolation condition (18) and

$$E_3(s_f) = \min\{E_3(s) : s \in \mathcal{S}, s(x^{(i)}) = f_i, i = 1, \dots, V\}.$$

By using similar arguments, we can obtain results similar to Theorem 8. We leave the details to the interested reader. We have implemented this method for a scattered data interpolation problem.

**Example 11.** Instead of interpolating function values, we may consider the gradient values at the vertices of  $\Delta$ . For simplicity, let us consider a problem in the bivariate setting, i.e., in  $\mathbb{R}^2$ . That is, we need to find a smooth surface  $s$  such that

$$D_x s(x_i, y_i) = f_{x,i}, \quad D_y s(x_i, y_i) = f_{y,i}, i = 1, \dots, V. \quad (25)$$

By fixing

$$s(x_1, y_1) = 0, \quad (26)$$

we formulate the following problem: Find  $s_f \in \mathcal{S}$  such that  $s_f(x_1, y_1) = 0$ ,  $s_f$  satisfies the above interpolation condition (25), and

$$E(s_f) = \min\{E(s) : s \in \mathcal{S}, s \text{ satisfies (25) and (26)}\}.$$

Similar to the arguments of Theorem 8, we can show the existence and uniqueness of a solution  $s_f$  to this problem. We leave the details to the interested reader.

**Example 12.** Instead of interpolating given data values, we may consider the problem of filling polygonal holes. Assume that a given smooth surface  $S$  has a hole which can be projected onto a plane  $P$  and the projection is a polygonal domain  $\Omega$ . We write  $g$  to be the function value of the given surface  $S$  on  $\partial\Omega$  and  $h$  to be the normal derivative value of  $S$  on  $\partial\Omega$ . Then we can formulate the following problem: find  $s_f \in \mathcal{S}$  such that

$$E(s_f) = \min\{E(s) : s|_{\partial\Omega} = g \text{ and } \frac{\partial}{\partial \mathbf{n}} s|_{\partial\Omega} = h\}.$$

Similar to the arguments of Theorem 8, we can show the existence and uniqueness of a spline function which fills the polygonal hole. We again leave the details to the interested reader. Cf. [18] for some numerical examples using bivariate  $C^1$  cubic splines.

### §5. The Discrete Least Squares Method for Data Fitting

Let  $\{(x^{(i)}, f_i), i = 1, \dots, N\}$  be a given set of scattered data in  $\mathbb{R}^n$ , where  $N$  is a relatively large integer. For example, in the bivariate setting,  $N \geq 1000$  with  $x^{(i)} \in [0, 1] \times [0, 1]$ . Also, the  $x^{(i)}$ 's may not be distinct. A general problem is to find a smooth surface which resembles the given data. In this section, we discuss the discrete least squares method for finding such a surface.

To this end, let  $\Omega$  be the convex hull of the given data locations and  $\Delta$  a triangulation of  $\Omega$ . Consider a spline space  $\mathcal{S} = S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$  for fixed sequences  $\rho, \mathbf{r}$ , and  $\mathbf{d}$  associated with  $k$ -simplices,  $0 \leq k < n - 1$ , interior  $(n - 1)$ -simplices, and  $n$ -simplices of  $\Delta$ . The discrete least squares method is to find  $s_f \in \mathcal{S}$  such that

$$\sum_{i=1}^N |s_f(x^{(i)}) - f_i|^2 = \min\left\{\sum_{i=1}^N |s(x^{(i)}) - f_i|^2, s \in \mathcal{S}\right\}. \quad (27)$$

We first discuss the existence and uniqueness of the solution. Let us introduce the following

**Definition 13.** For a given spline space  $S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$ , let  $d = \max\{d_t, t \in \Delta\}$ . We say that the given data locations  $x^{(i)}, i = 1, \dots, N$  are evenly distributed over  $\Delta$  with respect to  $d$  if for each triangle  $t \in \Delta$ , the matrix

$$[B_{\alpha}^t(x^{(\ell)}), \alpha \in \mathbf{Z}^{n+1}, |\alpha| = d, x^{(\ell)} \in t]$$

is of full rank.

**Theorem 14.** Suppose that the data locations  $x^{(i)}, i = 1, \dots, N$  are evenly distributed over  $\Delta$  with respect to  $d$ . Then there exists a unique discrete least squares solution  $s_f$  satisfying (27).

**Proof:** For convenience, we consider the existence and uniqueness in the bivariate setting. Recall that any  $s \in \mathcal{S}$  can be given by

$$s(x, y) = \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x, y), \quad \text{if } (x, y) \in t \in \Delta.$$

Let  $\mathbf{c} = (c_{ijk}^t, i + j + k = d, t \in \Delta)$  be the coefficient vector of  $s$ . For any  $s \in \mathcal{S}$ , we have  $H\mathbf{c} = 0$  and  $D\mathbf{c} = 0$  for the smoothness conditions and degree reduction conditions. Let

$$\begin{aligned} L(\mathbf{c}) &= \sum_{i=1}^N |s(x^{(i)}) - f_i|^2 \\ &= \sum_{t \in \Delta} \sum_{(x_{\ell}, y_{\ell}) \in t} \left( \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x_{\ell}, y_{\ell}) - f_{\ell} \right)^2. \end{aligned}$$

Note that  $L(0) = \|\mathbf{f}\|_2^2$ , with  $\mathbf{f} = (f_\ell, \ell = 1, \dots, N)$  being a data value vector and  $\|\mathbf{f}\|_2 := \left( \sum_{i=1}^N |f_i|^2 \right)^{1/2}$  denoting the standard Euclidean norm of the vector  $\mathbf{f}$ . Consider

$$A = \{\mathbf{c}, L(\mathbf{c}) \leq \|\mathbf{f}\|_2^2\}.$$

We now show that  $A$  is a bounded and closed set. Fix any triangle  $t \in \Delta$ . For any  $\mathbf{c} \in A$ , we have

$$\left| \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x_\ell, y_\ell) - f_\ell \right| \leq \|\mathbf{f}\|_2, \quad \forall (x_\ell, y_\ell) \in t.$$

It follows that

$$\left| \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x_\ell, y_\ell) \right| \leq 2\|\mathbf{f}\|_2, \quad \forall (x_\ell, y_\ell) \in t.$$

Since the data locations are evenly distributed, the matrix

$$[B_{ijk}^t(x_\ell, y_\ell)]_{\substack{i+j+k=d \\ (x_\ell, y_\ell) \in t}}$$

is of full rank and hence, there exists an index set  $I_t \subset \{1, 2, \dots, N\}$  such that the square matrix

$$B_t := [B_{ijk}^t(x_\ell, y_\ell)]_{\substack{i+j+k=d \\ \ell \in I_t}}$$

is invertible. It follows that

$$\|(c_{ijk}^t, i+j+k=d)\|_2 \leq C_t,$$

where  $C_t$  is a positive constant depending only on  $\|\mathbf{f}\|_2$  and the  $\ell_2$  norm of the inverse matrix of  $B_t$ . Hence,  $\|\mathbf{c}\|_2$  is bounded above. Thus,  $A$  is bounded. It is easy to see that  $A$  is closed. Also, it is easy to see that  $A_s := \{\mathbf{c}, H\mathbf{c} = 0, D\mathbf{c} = 0\}$  is a closed set and hence the set  $A \cap A_s$  is compact.

It is clear that  $L(\mathbf{c})$  is a continuous function of  $\mathbf{c}$ . Therefore, there exists a  $\mathbf{c}_f \in A \cap A_s$  such that  $L(\mathbf{c})$  achieves its minimum at  $\mathbf{c}_f$ .

Next we show that  $\mathbf{c}_f$  is the unique solution. We first note that  $L(\mathbf{c})$  is a convex function. Suppose that there exist two minimizers  $\mathbf{c}_f$  and  $\hat{\mathbf{c}}_f$ . Then any convex combination of  $\mathbf{c}_f$  and  $\hat{\mathbf{c}}_f$  is a minimizer. That is,

$$L(\mathbf{c}_f + z(\hat{\mathbf{c}}_f - \mathbf{c}_f)) = L(\mathbf{c}_f)$$

for any  $z \in [0, 1]$ . Thus,

$$\begin{aligned}
& \frac{1}{2} \frac{d}{dz} L(\mathbf{c}_f + z(\hat{\mathbf{c}}_f - \mathbf{c}_f)) \\
&= \sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \left( \sum_{i+j+k} c_{ijk}^t + z(\hat{c}_{ijk}^t - c_{ijk}^t) B_{ijk}^t(x_\ell, y_\ell) - f_\ell \right) \times \\
&\quad (\hat{c}_{ijk}^t - c_{ijk}^t) B_{ijk}^t(x_\ell, y_\ell) \\
&= z \sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \sum_{i+j+k} (\hat{c}_{ijk}^t - c_{ijk}^t)^2 B_{ijk}^t(x_\ell, y_\ell)^2 \\
&\quad + \sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \sum_{i+j+k} c_{ijk}^t (\hat{c}_{ijk}^t - c_{ijk}^t) B_{ijk}^t(x_\ell, y_\ell)^2 \\
&\quad - \sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \sum_{i+j+k} f_\ell (\hat{c}_{ijk}^t - c_{ijk}^t) B_{ijk}^t(x_\ell, y_\ell) \\
&= 0
\end{aligned}$$

for any  $z \in (0, 1)$ . It follows that

$$\sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \sum_{i+j+k} (\hat{c}_{ijk}^t - c_{ijk}^t)^2 B_{ijk}^t(x_\ell, y_\ell)^2 = 0.$$

That is,  $\mathbf{c}_f = \hat{\mathbf{c}}_f$  because the data locations are evenly distributed over each triangle of  $\Delta$ . Hence, the minimizer is unique.

Clearly, the whole arguments above can be generalized to the multivariate setting  $\mathbb{R}^n$  with  $n > 2$  easily. We have thus completed the proof.  $\square$

To find the minimizer, we note that  $L(\mathbf{c})$  is a convex function. Thus, any local minimizer is the global minimizer. Hence, we only need to find a local minimizer. By the Lagrange multiplier method, we let

$$\mathcal{F}(\mathbf{c}, \alpha, \beta) := L(\mathbf{c}) + \alpha^T H \mathbf{c} + \beta^T D \mathbf{c}$$

and set

$$\frac{\partial}{\partial \mathbf{c}} \mathcal{F}(\mathbf{c}, \alpha, \beta) = 0, \quad \frac{\partial}{\partial \alpha} \mathcal{F}(\mathbf{c}, \alpha, \beta) = 0, \quad \frac{\partial}{\partial \beta} \mathcal{F}(\mathbf{c}, \alpha, \beta) = 0.$$

It follows that we need to solve the following linear system

$$\begin{bmatrix} H^T & D^T & 2B \\ 0 & 0 & H \\ 0 & 0 & D \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} 2\mathbf{b} \\ 0 \\ 0 \end{bmatrix}$$

where  $B = \text{diag} (B^t, t \in \Delta)$  with

$$B^t = \left( \sum_{x^{(\ell)} \in t} B_{\alpha}^t(x^{(\ell)}) B_{\alpha'}^t(x^{(\ell)}) \right)_{\substack{\alpha \in \mathbf{Z}^{n+1}, |\alpha|=d \\ \alpha' \in \mathbf{Z}^{n+1}, |\alpha'|=d}}$$

being a matrix of size  $\hat{d} \times \hat{d}$  with  $\hat{d} = \binom{d+n}{n}$  and  $\mathbf{b} = (b_{\alpha}^t, \alpha \in \mathbf{Z}^{n+1}, |\alpha| = d, t \in \Delta)$  with

$$b_{\alpha}^t = \sum_{x^{(\ell)} \in t} f_{\ell} B_{\alpha}^t(x_{\ell}, y_{\ell}).$$

We again use the matrix iterative method to solve the above linear system.

Note that when the data set is evenly distributed, the above Theorem 14 implies that  $B$  is positive definite with respect to  $[H; D]$ . Therefore, our matrix iterative method can be applied. If the data set is not evenly distributed, we may use any least squares method instead of the inverses in (11) to compute a fitting spline. See Example 16. When we solve the discrete least squares problem, we use the maximum norm

$$\left\| \begin{bmatrix} H\mathbf{c} \\ D\mathbf{c} \end{bmatrix} \right\|_{\infty} \tag{28}$$

to measure if the iterative solution does satisfy the smoothness conditions and degree reduction conditions. For convenience, we will call such a measurement the exactness of the solution.

**Example 15.** Consider 1000 random points  $(x_i, y_i)$ 's over  $[0, 1] \times [0, 1]$  as shown in Fig. 2. Let  $\{(x_i, y_i, f(x_i, y_i)), i = 1, \dots, 1000\}$  be a scattered data set, where

$$f(x, y) = \sin(\pi(x^2 + 2y^2)).$$

We use the bivariate spline spaces  $S_d^r(\Delta)$  to find the discrete least squares fitting splines and then compare the maximum errors against the exact function, where  $\Delta$  is the triangulation given in Fig. 2. The maximum errors are measured using  $101 \times 101$  equally-spaced points over  $[0, 1] \times [0, 1]$ . We have checked the exactness (28) of the solutions for all the cases listed in Table 2. The exactness is always less than  $10^{-8}$ .

Table 2. The approximation errors using discrete least squares splines

	Errors	CPU		Errors	CPU
$S_3^0(\Delta)$	0.273500	5s	$S_3^1(\Delta)$	0.462761	5s
$S_4^0(\Delta)$	0.076285	10s	$S_4^1(\Delta)$	0.197462	10.7s
$S_5^0(\Delta)$	0.014693	20s	$S_5^1(\Delta)$	0.052269	21.5s

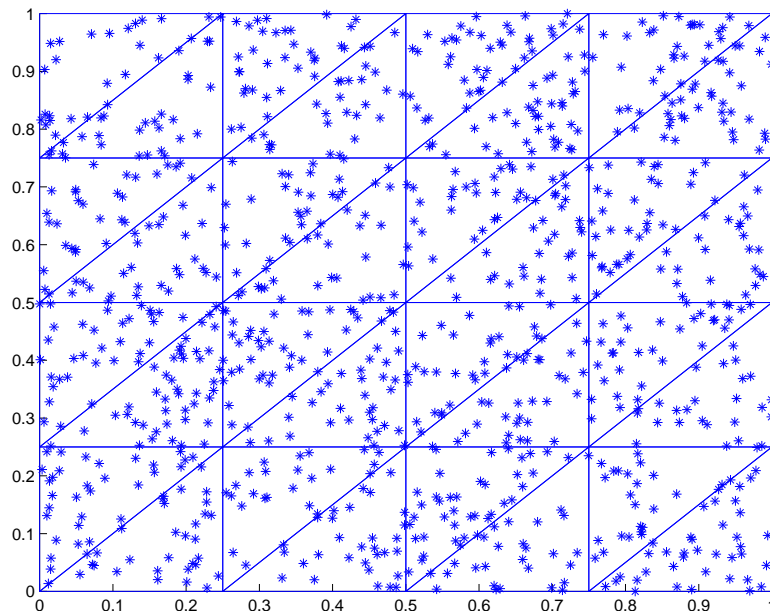


Fig. 2 A set of 1000 scattered data points

**Example 16.** We use the same data locations as in the previous example while refining the triangulation uniformly. See Fig. 3 below for the refined triangulation. Note that on the upper left corner, there are fewer points in two of the triangles. Thus, the data set is not evenly distributed over all triangles. We use  $f(x, y) = \sin(\pi(x^2 + 2y^2))$  to get the data values at these scattered data locations. Our method produces a  $C^1$  cubic fitting spline as shown in Fig. 4. (The exactness (28) is less than  $10^{-8}$  and thus, it is a  $C^1$  surface). The spline surface closely resembles the surface  $z = f(x, y)$  except near the upper left corner because there is not enough given data there.

**Example 17.** Consider 10,000 random points  $(x_i, y_i)$ 's over  $[0, 1] \times [0, 1]$ . Let  $\{(x_i, y_i, f(x_i, y_i)), i = 1, \dots, 10000\}$  be a scattered data set, where

$$f(x, y) = \sin(\pi(x^2 + 2y^2)).$$

We use the bivariate spline spaces  $S_d^r(\Delta)$  to find the discrete least squares fitting splines and then compare the maximum errors against the exact function, where  $\Delta$  is the triangulation given as in Example 15. The maximum errors are measured using  $101 \times 101$  equally-spaced points over  $[0, 1] \times [0, 1]$ . The exactness (28) is within  $10^{-8}$  for all the cases listed in the Table 3.

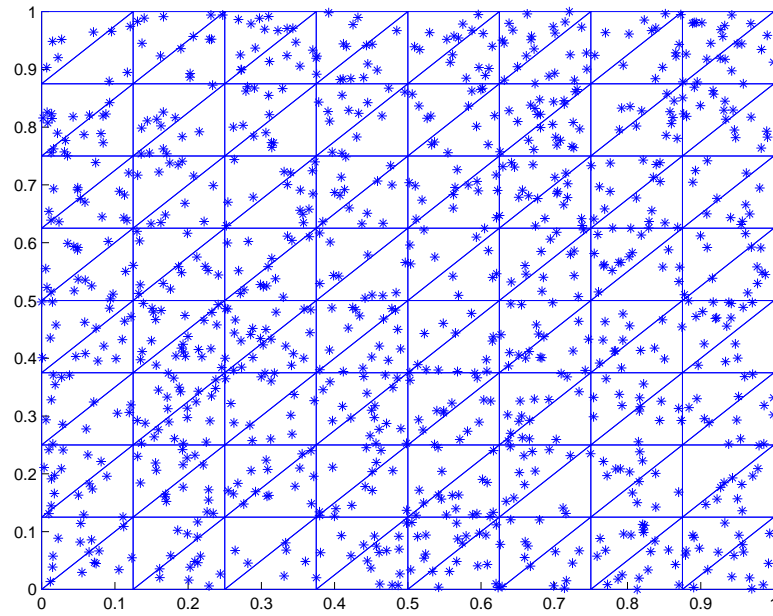


Fig. 3 Refined triangulation

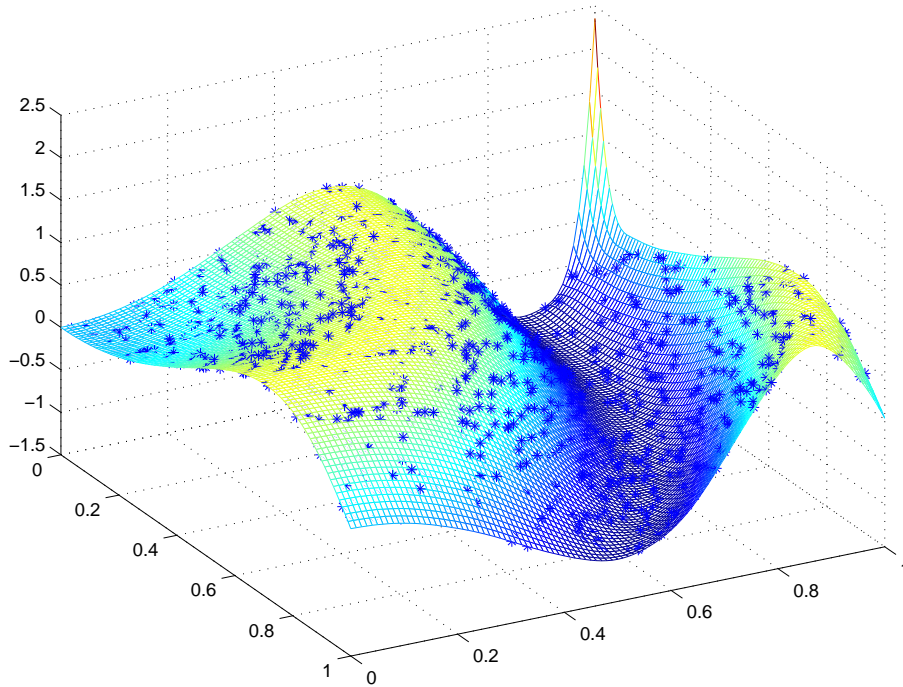


Fig. 4 Discrete Least Squares Fitting of  $f(x, y) = \sin(\pi(x^2 + 2y^2))$

Table 3. The approximation errors from discrete least squares splines

	Errors	CPU		Errors	CPU
$S_3^1(\Delta)$	0.092324	60s	$S_6^1(\Delta)$	0.000105	335s
$S_4^1(\Delta)$	0.007112	112s	$S_5^2(\Delta)$	0.002193	205s
$S_5^1(\Delta)$	0.001352	193s	$S_6^2(\Delta)$	0.000292	345s

### §6. The Penalized Least Squares Method for Data Fitting

Let  $\{(x^{(i)}, f_i), i = 1, \dots, N\}$  be a given set of scattered data in  $\mathbb{R}^n$ , where  $N$  is a relatively large integer. Also, the  $x^{(i)}$ 's may not be distinct. Without the assumption of even distribution of the data locations, we need to find a smooth surface which resembles the given data. In this situation, we use the penalized least squares method for finding such a surface.

To this end, let  $\Omega$  be the convex hull of the given data locations and  $\Delta$  be a triangulation of  $\Omega$ . Consider a spline space  $\mathcal{S} = S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$  with fixed smoothness  $\mathbf{r}$  and local smoothness sequence  $\rho$  associated with  $k$ -simplices,  $0 \leq k < n - 1$ , and degree sequence  $\mathbf{d}$  associated with  $n$ -simplices of  $\Delta$ . Assume that the global smoothness  $\min\{\mathbf{r}, \rho\}$  of the spline space  $\mathcal{S}$  is at least 1. The penalized least squares method is to find  $s_f \in \mathcal{S}$  such that

$$P_\lambda(s_f) = \min\{P_\lambda(s) : s \in \mathcal{S}\} \quad (29)$$

where  $\lambda > 0$  is a positive weight,

$$P_\lambda(s) := \sum_{i=1}^N |s(x^{(i)}) - f_i|^2 + \lambda E(s), \quad (30)$$

and  $E(s)$  denotes the usual energy functional defined in (20). It is clear to see that when  $\lambda \gg 1$ ,  $s_f$  is close to a linear plane fit of the given data and when  $\lambda \ll 1$ ,  $s_f$  is close to the discrete least squares fitting. A common method to choose an appropriate weight  $\lambda$  is the cross validation method (cf. [56]).

We first discuss the existence and uniqueness of the solution  $s_f \in \mathcal{S}$  satisfying (29) and a numerical method to compute  $s_f$ .

**Theorem 18.** *Fix a  $\lambda > 0$ . Suppose that all vertices of  $\Delta$  are part of the data locations. Then there exists a unique  $s_f \in \mathcal{S}$  satisfying (29).*

**Proof:** For convenience, we first prove the existence and uniqueness in the bivariate setting. We shall use  $x^{(i)} = (x_i, y_i)$ . Let us write any spline function  $s \in \mathcal{S}$  as

$$s(x, y)|_t = \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x, y), \quad (x, y) \in t \in \Delta,$$

where  $d = \max\{d_i, i = 1, \dots, T\}$ . Let  $\mathbf{c} = (c_{ijk}^t, i + j + k = d, t \in \Delta)$  be the coefficient vector associated with  $s$ . The length of vector  $\mathbf{c}$  is  $\hat{d}T$  with  $T$  being the number of triangles in  $\Delta$  and  $\hat{d} = (d+1)(d+2)/2$ . Certainly,  $s$  satisfies the smoothness and super smoothness conditions which can be expressed by  $H\mathbf{c} = 0$ . Also,  $s$  satisfies the degree reduction condition  $D\mathbf{c} = 0$ . The energy functional  $E(s)$  can be expressed in terms of  $\mathbf{c}$  as

$$E(s) = \mathbf{c}^T K \mathbf{c}$$

with  $K = \text{diag}(K_t, t \in \Delta)$  and

$$K_t = \left[ \int_t (D_x^2 B_{ijk}^t D_x^2 B_{p,q,r}^t + 2D_x D_y B_{ijk}^t D_x D_y B_{p,q,r}^t + D_y^2 B_{ijk}^t D_y^2 B_{p,q,r}^t) dx dy \right]_{\substack{i+j+k=d \\ p+q+r=d}}$$

We have

$$\begin{aligned} & \sum_{i=1}^N |s(x_i, y_i) - f_i|^2 \\ &= \sum_{t \in \Delta} \sum_{(x_\ell, y_\ell) \in t} \left( \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x_\ell, y_\ell) - f_\ell \right)^2 \\ &= \mathbf{c}^T B \mathbf{c} - 2\mathbf{b}^T \mathbf{c} + \|\mathbf{f}\|^2, \end{aligned}$$

where  $\mathbf{f} = (f_\ell, \ell = 1, \dots, N)$  is a data value vector,  $B$  and  $\mathbf{b}$  are as in the previous section. Thus, we have

$$P_\lambda(s) = \lambda \mathbf{c}^T K \mathbf{c} + \mathbf{c}^T B \mathbf{c} - 2\mathbf{b}^T \mathbf{c} + \|\mathbf{f}\|_2^2. \quad (31)$$

It is easy to see that  $P_\lambda(0) = \|\mathbf{f}\|_2^2$ . Consider

$$A = \{\mathbf{c}, P_\lambda(s) \leq \|\mathbf{f}\|_2^2\}.$$

We now show that  $A$  is a bounded and closed set. Clearly, for  $\mathbf{c} \in A$ , we have

$$\int_t |D_x^2 s(x, y)|^2 dx dy \leq \frac{1}{\lambda} \|\mathbf{f}\|_2^2$$

for any triangle  $t \in \Delta$ . Since  $D_x^2 s(x, y)|_t$  is in the space of polynomials of degree  $\leq d-2$ , it follows that

$$\max_{(x,y) \in t} |D_x^2 s(x, y)| \leq C \frac{1}{\sqrt{\lambda}} \|\mathbf{f}\|_2$$

for some  $C$  depending on the triangle  $t$ . Similar for  $D_x D_y s$  and  $D_y^2 s$ . Thus, we have  $|s|_{2,\infty,t} \leq C \|\mathbf{f}\|_2 / \sqrt{\lambda}$  for another constant  $C > 0$  depending on  $t$ , where

$$|s|_{2,\infty,t} := \max_{(x,y) \in t} \max\{|D_x^2 s(x,y)|, |D_x D_y s(x,y)|, |D_y^2 s(x,y)|\}.$$

Next we claim that  $\max\{|s(x,y)|, (x,y) \in t\}$  is bounded. Since all the vertices of  $\triangle$  are part of the data locations, writing  $t = \langle v_i, v_j, v_k \rangle$ , we have  $v_i, v_j, v_k \in \{(x_\ell, y_\ell), \ell = 1, \dots, N\}$ . For simplicity, let  $f_i$  be the data value associated with  $v_i$ . It follows that, for  $\mathbf{c} \in A$ ,

$$\begin{aligned} |s(v_i)| &\leq |s(v_i) - f_i| + |f_i| \leq \left( \sum_{(x_\ell, y_\ell) \in t} (s(x_\ell, y_\ell) - f_\ell)^2 \right)^{1/2} + \|\mathbf{f}\|_2 \\ &\leq (P_\lambda(s))^{1/2} + \|\mathbf{f}\|_2 \\ &\leq \left( \frac{1}{\sqrt{\lambda}} + 1 \right) \|\mathbf{f}\|_2. \end{aligned}$$

Similarly, we have the same estimate for  $|s(v_j)|$  and  $|s(v_k)|$ . For any point  $v = (x, y) \in t$ , we use the Taylor expansion to get

$$s(v_i) = s(v) + \nabla s(v) \cdot (v_i - v) + O(|s|_{2,\infty,t} |t|^2), \quad (32)$$

where  $|t|$  denotes the diameter of  $t$ . Similar for  $s(v_j)$  and  $s(v_k)$ . It follows that

$$\begin{aligned} s(v_j) - s(v_i) &= \nabla s(v) \cdot (v_j - v_i) + O(|s|_{2,\infty,t} |t|^2) \\ s(v_k) - s(v_i) &= \nabla s(v) \cdot (v_k - v_i) + O(|s|_{2,\infty,t} |t|^2). \end{aligned}$$

Solving this linear system for  $\nabla s(v) = (D_x s(v), D_y s(v))$  gives

$$\begin{aligned} D_x s(v) &= O(|t|^3 |s|_{2,\infty,t} / A_t) + |s(v_j)| + |s(v_i)| |t| / A_t \\ D_y s(v) &= O(|t|^3 |s|_{2,\infty,t} / A_t) + |s(v_k)| + |s(v_i)| |t| / A_t, \end{aligned}$$

where  $A_t$  stands for the area of the triangle  $t$ . Inserting these estimates for  $\nabla s(v)$  in (32), we immediately get

$$|s(v)| \leq C \left( (1 + |t|^2 / A_t) \|\mathbf{f}\|_\infty + |t|^4 |s|_{2,\infty,t} / A_t \right).$$

Hence, we have

$$|s(x, y)| \leq C_t \left( \frac{1}{\sqrt{\lambda}} + 1 \right) \|\mathbf{f}\|_2$$

for another constant  $C_t > 0$  depending on  $t$ . It follows from Lemma 2 that

$$\max_{i+j+k=d} |c_{ijk}^t| \leq N \max_{(x,y) \in t} |s(x, y)| \leq N C_t \left( \frac{1}{\sqrt{\lambda}} + 1 \right) \|\mathbf{f}\|_2.$$

Therefore  $\mathbf{c}$  is bounded and thus,  $A$  is a bounded set. Clearly,  $A$  is a closed set and hence,  $A$  is compact.

By (31), it is clear that  $P_\lambda$  is a continuous function of the variable  $\mathbf{c}$ . Hence,  $P_\lambda$  achieves its minimum over the compact set  $A$ . That is, there exists a spline  $s_f$  solving the minimization problem (29).

Next we show the uniqueness of the minimizer  $s_f$ . Suppose that we have two solutions  $s_f$  and  $\hat{s}_f$ . Let  $\mathbf{c}_f$  and  $\hat{\mathbf{c}}_f$  be the two coefficients associated with  $s_f$  and  $\hat{s}_f$ , respectively. Since  $P_\lambda$  is a convex functional, we have, for any  $z \in [0, 1]$ ,

$$P_\lambda(zs_f + (1-z)\hat{s}_f) \leq zP_\lambda(s_f) + (1-z)P_\lambda(\hat{s}_f) = P_\lambda(s_f).$$

That is,  $P_\lambda(\hat{s}_f + z(s_f - \hat{s}_f))$  is a constant function of  $z \in [0, 1]$ . It follows that  $\frac{\partial}{\partial z}P_\lambda(\hat{s}_f + z(s_f - \hat{s}_f)) = 0$  for all  $z \in (0, 1)$ . That is,

$$\begin{aligned} 0 &= \frac{\partial}{\partial z}P_\lambda(\hat{s}_f + z(s_f - \hat{s}_f)) \\ &= 2\lambda z(\mathbf{c}_f - \hat{\mathbf{c}}_f)^T K(\mathbf{c}_f - \hat{\mathbf{c}}_f) + 2z(\mathbf{c}_f - \hat{\mathbf{c}}_f)^T B(\mathbf{c}_f - \hat{\mathbf{c}}_f) \\ &\quad - 2\mathbf{b}^T(\mathbf{c}_f - \hat{\mathbf{c}}_f) \end{aligned}$$

for all  $z \in (0, 1)$ . Thus, we have

$$(\mathbf{c}_f - \hat{\mathbf{c}}_f)^T K(\mathbf{c}_f - \hat{\mathbf{c}}_f) = 0 \text{ and } (\mathbf{c}_f - \hat{\mathbf{c}}_f)^T B(\mathbf{c}_f - \hat{\mathbf{c}}_f) = 0.$$

because both  $K$  and  $B$  are nonnegative definite. The first equation is equivalent to  $E(s_f - \hat{s}_f) = 0$  which implies that  $s_f - \hat{s}_f$  is a linear polynomial. The second equation implies that  $s_f - \hat{s}_f$  is equal to zero at all vertices of  $\Delta$ . Thus,  $s_f - \hat{s}_f \equiv 0$ . Hence, the minimizer is unique. We have thus completed the proof for the bivariate setting.

We note that the above arguments can be easily generalized to the multivariate setting, and we leave the generalization to the interested reader.  $\square$

Next we look at how to compute the minimizer. Since  $P_\lambda(s)$  is a convex functional, any local minimizer is the global minimizer. To find a local minimizer, we use the Lagrange multiplier method by letting

$$\mathcal{F}(\mathbf{c}, \alpha, \beta) = P_\lambda(s) + \alpha^T H\mathbf{c} + \beta^T D\mathbf{c}$$

and compute

$$\frac{\partial}{\partial \mathbf{c}}\mathcal{F}(\mathbf{c}, \alpha, \beta) = 0, \frac{\partial}{\partial \alpha}\mathcal{F}(\mathbf{c}, \alpha, \beta) = 0, \frac{\partial}{\partial \beta}\mathcal{F}(\mathbf{c}, \alpha, \beta) = 0.$$

Using the expression (31), we have

$$\begin{bmatrix} H^T & D^T & 2(B + \lambda K) \\ 0 & 0 & H \\ 0 & 0 & D \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} 2\mathbf{b} \\ 0 \\ 0 \end{bmatrix}. \quad (33)$$

We shall use the matrix iterative method introduced in §3 to solve the above linear system. Since the uniqueness of the solution implies that  $B + \lambda K$  is positive definite with respect to  $[H; D]$ , Theorem 6 is applicable. This establishes our numerical method for penalized least squares spline method for scattered data fitting. Furthermore, we use the maximum norm

$$\left\| \begin{bmatrix} H\mathbf{c} \\ D\mathbf{c} \end{bmatrix} \right\|_{\infty} \quad (34)$$

to measure if the iterative solution does satisfy the smoothness conditions and degree reduction conditions. We say that the solution is exact if it's zero in the above norm.

**Example 19.** Consider 1000 random points  $(x_i, y_i)$  over  $[0, 1] \times [0, 1]$  as shown in Fig. 2. Let  $\{(x_i, y_i, f(x_i, y_i)), i = 1, \dots, 1000\}$  be a scattered data set, where

$$f(x, y) = \sin(\pi(x^2 + 2y^2)).$$

Let  $\Delta$  be the two triangles obtained from  $[0, 1] \times [0, 1]$  by adding one diagonal. We uniformly refine  $\Delta$  one, two, and three times to obtain new triangulations  $\Delta_1, \Delta_2$ , and  $\Delta_3$ . We use bivariate spline spaces  $S_5^1(\Delta_i), i = 1, 2, 3$  to find the penalized least squares fitting splines with different  $\lambda$  and then compare the maximum errors against the exact function. The maximum errors are measured using  $100 \times 100$  equally-spaced points over  $[0, 1] \times [0, 1]$  and given in Table 4. The CPU time of the computations is measured in seconds. The exactness (34) of the solution is checked and is less than  $10^{-8}$  for all the cases.

Table 4. Maximum Errors from Penalized Least Squares Method

Splines $\setminus \lambda$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	CPU
$S_5^1(\Delta_1)$	1.2345	0.6621	0.2826	0.6587		6s
$S_5^1(\Delta_2)$	1.2222	0.7908	0.2951	0.1217	0.0595	22s
$S_5^1(\Delta_3)$	1.2344	0.8033	0.3054	0.1484	0.1046	99s

Next we use 10,000 scattered data locations in  $[0, 1] \times [0, 1]$  and perform the same experiments as above. The maximum errors are listed in Table 5.

Table 5. Maximum Errors from Penalized Least Squares Method (cont.)

Splines \ $\lambda$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	CPU
$S_5^1(\Delta_1)$	0.9273	0.1979	0.7181				20s
$S_5^1(\Delta_2)$	1.0050	0.3685	0.1121	0.0256	0.0108	0.01233	55s
$S_5^1(\Delta_3)$	1.0127	0.3840	0.1085	0.0278	0.0090	0.00224	206s

**Example 20.** We also perform a similar numerical experiment using trivariate splines. That is, we use 10,000 scattered data located in  $[0, 1] \times [0, 1] \times [0, 1]$ . We partition the domain into 12 tetrahedra by connecting the midpoint  $(0.5, 0.5, 0.5)$  to the six square faces and dividing each of the six resulting pyramids into two tetrahedra. These 12 tetrahedra form a tetrahedral partition  $\Delta$  of the given domain  $[0, 1] \times [0, 1] \times [0, 1]$ . The test function is  $f(x, y, z) = \exp(x + y + z)$ . We use the trivariate spline space  $S_5^1(\Delta)$  to find the penalized least squares fitting splines with different  $\lambda$  and then compare the maximum errors against the exact function. The maximum errors are measured using  $20 \times 20 \times 20$  equally-spaced points over  $[0, 1] \times [0, 1] \times [0, 1]$  and given in Table 6. The exactness (34) of the solution is checked and is less than  $10^{-8}$  for all the cases.

Table 6. Maximum Errors from Penalized Least Squares Method (cont.)

Splines \ $\lambda$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	CPU
$S_5^1(\Delta)$	0.6904	0.1793	0.0320	0.0070	0.0041	84s

## §7. Numerical Solution of Poisson Equations and

### Other Second Order Elliptic Equations

In this section, we will show how to solve the Poisson equation and other second order elliptic equations by using multivariate splines of variable degree and variable smoothness. These spline functions will provide a versatile tool for numerical solution of PDE's because of the flexibility in choosing the degrees and smoothness when constructing numerical solutions. For example, it is known that the weak solution of the Poisson equation over a polygonal domain  $\Omega$  is at least  $H^2(V)$  for any open set  $V \subset \Omega$  (cf. [22]). We should choose spline functions which are  $C^1$  inside  $\Omega$  and  $C^0$  near the boundary of  $\Omega$  to find an approximate weak solution.

Let us begin with the Poisson equation:

$$\begin{cases} -\Delta u = f & \text{in } \Omega \\ u = g, & \text{on } \partial\Omega \end{cases}$$

where  $\Omega$  is a polygonal domain in  $\mathbb{R}^n$ ,  $f \in L_2(\Omega)$ , and  $g$  is continuous over the boundary  $\partial\Omega$  of  $\Omega$ . The weak formulation of the Poisson equation is to find  $u \in H^1(\Omega)$  which satisfy  $u = g$  on  $\partial\Omega$  and

$$a(u, v) = \langle f, v \rangle, \quad \forall v \in H_0^1(\Omega),$$

where  $a(u, v)$  is the bilinear form defined by

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v dx dy$$

and  $\langle f, v \rangle = \int_{\Omega} f v dx dy$  stands for the standard  $L_2$  inner product of  $f$  and  $v$ . Here  $H^1(\Omega)$  and  $H_0^1(\Omega)$  are standard Sobolev spaces. By the standard calculus of variations, the Poisson equation is the Euler-Lagrange equation of the energy functional

$$E(w) = \int_{\Omega} \left( \frac{1}{2} \nabla w \cdot \nabla w - w f \right) dx.$$

It is known that the weak solution of the Poisson equation is the minimizer of the energy functional  $E(w)$  among the class of admissible functions

$$\mathcal{A} = \{w \in H^1(\Omega), w = g \text{ on } \partial\Omega\}.$$

(Cf. [22], §8.2.3.) That is, the weak solution  $u$  satisfies

$$E(u) = \min_{w \in \mathcal{A}} E(w). \quad (35)$$

Also any minimizer satisfying (35) is the weak solution.

Next we discuss how to compute approximate weak solutions which are multivariate spline functions. For convenience, let us consider the Poisson equation in the bivariate setting first. Let  $\Delta$  be a triangulation of the domain  $\Omega \in \mathbb{R}^2$  and let

$$\mathcal{S} := S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$$

be a spline space with fixed (global) smoothness  $\mathbf{r}$ , local smoothness vector  $\rho$  and degree vector  $\mathbf{d}$  associated with vertices, interior edges, and triangles of  $\Delta$ . Let  $d$  be the largest integer in  $\mathbf{d}$ . Instead of piecewise linear boundary of  $\Omega' = \cup_{t \in \Delta} t$ , we may use piecewise quadratic polynomials to approximate the boundary of  $\Omega$ . That is, for each boundary triangle  $t \in \Delta$ , if the boundary edge  $e_t$  of  $t$  is not a part of boundary  $\partial\Omega$  of  $\Omega$ , we use a circular arc  $\tilde{e}_t$  which passes through two vertices of  $e_t$  and another point on  $\partial\Omega$  between the two vertices to replace  $e_t$ . Let  $\tilde{t}$  be the convex hull of the vertices of  $t$  and the circular arc. All the interior triangles and

new boundary triangles (with curved side) form a new domain  $\tilde{\Omega}$  which is a better approximation of  $\Omega$  than  $\Omega'$ . Since each spline function  $s \in \mathcal{S}$  can be extended naturally to  $\tilde{\Omega}$ , we may consider that  $\mathcal{S}$  are defined on  $\tilde{\Omega}$ .

We remark that when solving the Poisson equation with Dirichlet boundary condition, we require spline functions to have less smoothness near the boundary while having more smoothness inside the domain according to the regularity theory of the weak solution of the Poisson and general elliptic PDE's (cf. [22]). In general, there is no spline function in  $\mathcal{S}$  satisfying the boundary condition exactly. Let  $\tilde{\mathcal{A}}$  be the subset of  $\mathcal{S}$  satisfying the boundary condition approximately in the sense that  $s_u \in \mathcal{S}$  interpolates  $g$  at  $2d + 1$  distinct points at each curve edge and  $d + 1$  distinct points over each straight boundary edge. Here, we have assumed that the degrees of the spline functions in  $\mathcal{S}$  are  $d$  over each boundary triangle. Otherwise, we modify the interpolation conditions appropriately. We compute the approximation  $s_u \in \mathcal{S}$  satisfying

$$E(s_u) = \min_{w \in \tilde{\mathcal{A}}} E(w).$$

Following the same arguments in [22], the minimizer  $s_u$  is the approximate weak solution in  $\mathcal{S}$ .

We next give an algorithm to compute such an  $s_u$  with the assumption that  $s_u$  exists and is unique. The proof of the existence and uniqueness is well-known and will be mentioned briefly later.

Let us write any spline function  $s \in \mathcal{S}$  as in (2), where  $d = \max\{d_t, t \in \Delta\}$ . That is,  $s \in \mathcal{S}$  may be expressed by

$$s(x, y)|_t = \sum_{i+j+k=d} c_{ijk}^t B_{ijk}^t(x, y), \quad (x, y) \in t \in \Delta.$$

Let  $\mathbf{c} = (c_{ijk}^t, i + j + k = d, t \in \Delta)$  be the coefficient vector associated with  $s$ . The length of the vector  $\mathbf{c}$  is  $\hat{d}T$  with  $T$  being the number of triangles in  $\Delta$  and  $\hat{d} = (d + 1)(d + 2)/2$ . The smoothness and super smoothness conditions that  $s$  satisfies can be expressed by  $H\mathbf{c} = 0$ . Also,  $s$  satisfies the degree reduction conditions, i.e.,  $D\mathbf{c} = 0$ .

Then the bilinear form  $a(s, \hat{s})$  can be expressed in terms of  $\mathbf{c}$  and  $\hat{\mathbf{c}}$  by

$$a(s, \hat{s}) = \mathbf{c}^T K \hat{\mathbf{c}}$$

where  $K = \text{diag}(K_t, t \in \Delta)$  with

$$K_t = \left[ \int_t \nabla B_{ijk}^t \cdot \nabla B_{p,q,r}^t dx dy \right]_{\substack{i+j+k=d \\ p+q+r=d}}.$$

Note that the inner product  $\langle f, \hat{s} \rangle$  can be approximated by  $\langle s_f, \hat{s} \rangle$  where  $s_f \in S_d^{-1}(\Delta)$ , the space of piecewise polynomials of degree  $d$  on each triangle, interpolates  $f$  over the domain points of each triangle  $t$ . Thus,

$$\langle f, \hat{s} \rangle \approx \hat{\mathbf{c}}^T M \mathbf{c}_f,$$

where  $M = \text{diag}(M^t, t \in \Delta)$  is a block diagonal matrix with square blocks

$$M^t = \left[ \int_t B_{ijk}^t(x, y) B_{p,q,r}^t(x, y) dx dy \right]_{\substack{i+j+k=d \\ p+q+r=d}}$$

and  $\mathbf{c}_f$  encodes the coefficients of  $s_f$ . We need to solve the following minimization problem:

$$\begin{aligned} & \min \frac{1}{2} \mathbf{c}^T K \mathbf{c} - \mathbf{c}^T M \mathbf{c}_f \\ & \text{subject to} \\ & H \mathbf{c} = 0, D \mathbf{c} = 0, B \mathbf{c} = \mathbf{g}, \end{aligned}$$

where  $B \mathbf{c} = \mathbf{g}$  denotes a linear system associated with the boundary conditions. Indeed, based on the de Casteljau algorithm, the evaluation of  $s_u$  at any point on a curved edge is a linear equation in terms of the unknown coefficients of  $s_u$ . As we can show that there exists a unique approximate weak solution  $s_u \in \mathcal{S}$ , we know that the minimization problem has a unique solution. Since the energy functional is convex, any local minimum is the global minimum. Let us compute a local minimum by using the Lagrange multiplier method. Letting

$$\mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = \frac{1}{2} \mathbf{c}^T K \mathbf{c} - \mathbf{c}^T M \mathbf{c}_f + \theta^T H \mathbf{c} + \eta^T D \mathbf{c} + \nu^T (B \mathbf{c} - \mathbf{g}),$$

we compute

$$\begin{aligned} \frac{\partial}{\partial \mathbf{c}} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) &= 0, \quad \frac{\partial}{\partial \theta} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = 0, \\ \frac{\partial}{\partial \eta} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) &= 0, \quad \frac{\partial}{\partial \nu} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = 0. \end{aligned}$$

It follows that

$$\begin{aligned} K \mathbf{c} + H^T \theta + D^T \eta + B^T \nu &= M \mathbf{c}_f \\ H \mathbf{c} = 0, D \mathbf{c} = 0, B \mathbf{c} &= \mathbf{g}. \end{aligned}$$

In other words, we have

$$\begin{bmatrix} B^T & D^T & H^T & K \\ 0 & 0 & 0 & H \\ 0 & 0 & 0 & D \\ 0 & 0 & 0 & B \end{bmatrix} \begin{bmatrix} \theta \\ \eta \\ \nu \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} M \mathbf{c}_f \\ 0 \\ 0 \\ \mathbf{g} \end{bmatrix}. \quad (36)$$

We shall apply the matrix iterative method for solving the above linear system when it is of large size. The uniqueness of the weak solution implies that  $K$  is positive definite with respect to  $[B; H; D]$ . Therefore, the matrix iterative method is well-defined. We remark that assembling the matrices  $M$  and  $K$  is particularly easy and can be done without knowing the relations among the triangles in any given triangulation partition. This is also true in the multivariate setting.

This leads to a numerical method to compute approximate weak solutions for Poisson equations in  $\mathbb{R}^2$ . It is clear that the above arguments can be generalized to the multivariate setting.

We have implemented this method using bi- and tri-variate spline spaces of any degree and any smoothness over any triangulation of any polygonal domain to solve 2D and 3D Poisson equation. We will provide several numerical experiments near the end of this section.

Next let us briefly discuss the existence and uniqueness of the approximate weak solution  $s_u$ . The discussion is parallel to the one using finite elements. Mainly we use the well-known Lax-Milgram Theorem. Since  $\mathcal{S}$  is a finite dimensional space, we may find a basis  $\{\phi_i, i = 1, \dots, \dim(\mathcal{S})\}$  which may not be locally supported. For any spline function  $s \in \mathcal{S}$ , we write  $s = \sum_i s_i \phi_i$  for some coefficients  $s_i$ 's. Thus, for  $s \in \mathcal{S} \cap H_0^1(\Omega)$ , the bilinear form can be given by

$$a(s, \hat{s}) = \mathbf{s}K'\hat{\mathbf{s}}$$

with a new stiffness matrix  $K'$ . Because  $a(\cdot, \cdot)$  is coercive, it can be easily shown that  $K'$  is positive definite over  $\mathcal{S} \cap H_0^1(\Omega)$ . Thus the existence and uniqueness of the approximation weak solution  $s_u$  follows.

We remark that the Poisson equation with Neumann boundary condition

$$\begin{cases} -\Delta u = f, & \text{in } \Omega \\ \frac{\partial u}{\partial \mathbf{n}} = h, & \text{on } \partial\Omega \\ \int_{\Omega} u dx = 0 \end{cases} \quad (37)$$

can be numerically solved in the same fashion. We leave the details to the interested reader. We have implemented the algorithm for solving (37) numerically in the bivariate and trivariate setting.

We now turn our attention to general second order elliptic equations. Consider

$$\begin{cases} -\sum_{i,j=1}^n \frac{\partial}{\partial x_i} \left( a_{ij}(x) \frac{\partial}{\partial x_j} u \right) = f, & x \in \Omega \\ u(x) = g(x), & x \in \partial\Omega, \end{cases} \quad (38)$$

where  $a_{ij}(x) = a_{ji}(x) \in L_{\infty}(\Omega)$  for  $i, j = 1, \dots, n$  and satisfy

$$\sum_{i,j=1}^n a_{ij} \lambda_i \lambda_j \geq m \sum_{i=1}^n \lambda_i^2, \forall \lambda_i, i = 1, \dots, n$$

for a positive constant  $m > 0$ . Using the results in ([22], §8.2.3), we can show that the weak solution of (38) is the minimizer of

$$E(w) = \int_{\Omega} \left( \frac{1}{2} \sum_{i,j=1}^n a_{ij} \frac{\partial}{\partial x_i} w \frac{\partial}{\partial x_j} w - wf \right) dx$$

over the set  $\mathcal{A}$  of admissible functions. Thus, to find an approximate weak solution in  $\mathcal{S}$ , we need to solve the following minimization problem:

$$\begin{aligned} & \min \frac{1}{2} \mathbf{c}^T \mathcal{K} \mathbf{c} - \mathbf{c}^T M \mathbf{f} \\ & \text{subject to} \\ & H \mathbf{c} = 0, \quad D \mathbf{c} = 0, \quad B \mathbf{c} = \mathbf{g}, \end{aligned}$$

where  $\mathcal{K} = \text{diag} (\tilde{K}_t, t \in \Delta)$  is a block diagonal matrix with

$$\tilde{K}_t = \left[ \int_t \sum_{i,j=1}^n a_{ij} \frac{\partial}{\partial x_i} B_{\alpha}^t \frac{\partial}{\partial x_j} B_{\alpha}^t dx \right]_{\substack{\alpha \in \mathbf{z}^{n+1}, |\alpha|=d \\ \hat{\alpha} \in \mathbf{z}^{n+1}, |\hat{\alpha}|=d}}.$$

The Lagrange multiplier method implies that we need to solve the following linear system:

$$\begin{bmatrix} B^T & D^T & H^T & \mathcal{K} \\ 0 & 0 & 0 & H \\ 0 & 0 & 0 & D \\ 0 & 0 & 0 & B \end{bmatrix} \begin{bmatrix} \theta \\ \eta \\ \nu \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} M \mathbf{c}_f \\ 0 \\ 0 \\ \mathbf{g} \end{bmatrix}. \quad (39)$$

Again the uniqueness of the weak solution implies that the matrix iterative method is well defined. Furthermore, we use the maximum norm

$$\left\| \left[ \begin{array}{c} H \mathbf{c} \\ D \mathbf{c} \\ B \mathbf{c} - \mathbf{g} \end{array} \right] \right\| \quad (40)$$

to check if the iterative solution does satisfy the smoothness conditions, degree reduction conditions, and boundary conditions. The solution will be said exact if it is zero in such a norm.

Let us report on some numerical experiments for the 2D and 3D Poisson equations.

**Example 21.** Consider the Poisson equation with exact solution

$$u(x, y) = 10 \exp(-(x^2 + y^2))$$

over a square domain:

$$\begin{cases} -\Delta u = 40 \exp(-(x^2 + y^2))(1 - x^2 - y^2) & (x, y) \in [-2, 2] \times [-2, 2] \\ u(x, y) = 10 \exp(-(x^2 + y^2)), & (x, y) \in \partial[-2, 2] \times [-2, 2] \end{cases}$$

The solution is relatively large inside the domain as compared to its values on the boundary. Our spline solutions can approximate it very well. We use a triangulation similar to Fig. 2 with 25 vertices and 32 triangles. We test many spline spaces and list the maximum errors of approximate weak spline solutions against the exact solution in Table 7. The maximum errors are computed based on  $101 \times 101$  equally-spaced points over  $[-2, 2] \times [-2, 2]$ . The exactness (34) is checked for all the spline spaces list Table 7 and is less than  $10^{-8}$ .w

Table 7. Approximation Errors from Bivariate Spline Spaces

	Maximum Errors	CPU Times
$S_3^1(\Delta)$	0.732222	0.40s
$S_4^1(\Delta)$	0.063235	0.48s
$S_5^1(\Delta)$	0.010793	0.78s
$S_6^1(\Delta)$	0.001382	1.06s
$S_7^1(\Delta)$	0.000502	1.65s
$S_8^1(\Delta)$	0.000173	2.56s
$S_9^1(\Delta)$	0.000013	4.03s

**Example 22.** Consider the 3D Poisson equation with exact solution

$$u(x, y, z) = 10 \exp(-(x^2 + y^2 + z^2))$$

over an octahedron  $\Omega := \langle (1, 0, 0), (0, 1, 0), (-1, 0, 0), (0, -1, 0), (0, 0, 1), (0, 0, -1) \rangle$ . We split  $\Omega$  into 8 tetrahedra by three coordinate planes. Let  $\Delta$  denote the collection of all 8 tetrahedra. We find approximate weak solutions in the 3D spline spaces  $S_d^1(\Delta)$  for  $d = 3, \dots, 7$ . The maximum errors are computed based on  $20 \times 20 \times 20$  equally-spaced points over  $\Omega$  and listed in Table 8. The exactness (40) is checked and is less than  $10^{-8}$ .

Table 8. Approximation Errors from Trivariate Spline Spaces

	matrix size	Maximum Errors	CPU Times
$S_3^1(\Delta)$	$160 \times 160$	0.17127	0.07s
$S_4^1(\Delta)$	$280 \times 280$	0.02737	0.17s
$S_5^1(\Delta)$	$448 \times 448$	0.00749	0.625s
$S_6^1(\Delta)$	$672 \times 672$	0.000842	1.67s
$S_7^1(\Delta)$	$960 \times 960$	0.0004601	5.18s

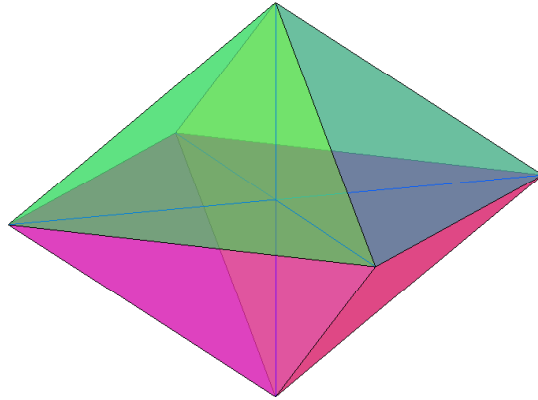


Fig. 5 A simple tetrahedron partition

Next we compare the errors by using refinements of underlying triangulations and by increasing the degrees of spline functions.

**Example 23.** We solve the Poisson equation with Dirichlet boundary condition over a star domain as shown in Fig. 6 with exact solution  $u = \exp(x + y)$  using  $C^1$  cubic splines over successively refined triangulations. We can only refine 3 times within the capacity of our PC and the results are listed in Table 9. In Table 10, the degrees of the spline spaces are varied.

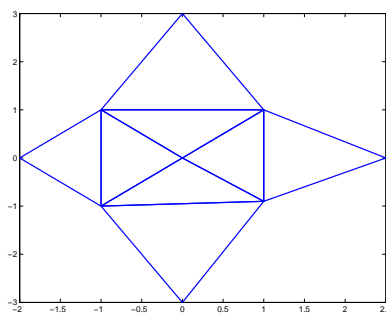


Fig. 6 An Initial Triangulation of a Star-shape Domain

Table 9. Approximation from Uniform Refinements (Dirichlet Problem)

Refinement Levels	Matrix Size	Maximum Errors
1	80 × 80	0.254346956
2	320 × 320	0.029554301
3	1280 × 1280	0.004515225
4	5120 × 5120	0.000535312

Table 10. Approximation from Degree Increase (Dirichlet Problem)

Polynomial Degrees	Matrix Size	Maximum Errors
3	80 × 80	0.25434695641
4	120 × 120	0.04251752024
5	168 × 168	0.00608204535
6	224 × 224	0.00080855135
7	288 × 288	0.00009770118
8	360 × 360	0.00001031184
9	440 × 440	0.00000096358
10	528 × 528	0.00000008441
11	624 × 624	0.00000000697

**Example 24.** Next we solve the Poisson equation with Neumann boundary condition over a star domain as shown in Figure 6 with exact solution

$$u = 10 \exp(-(x^2 + y^2))$$

using  $C^1$  quartic splines over succesively refined triangulation. The numerical results are given in Tables 11 and 12 with refinements and degree increases.

Table 11. Approximation from Uniform Refinements (Neumann Problem)

Refinement Levels	Matrix Size	Maximum Errors
1	120 × 120	$2.38 \times 10^{-2}$
2	480 × 480	$6.81 \times 10^{-4}$
3	1920 × 1920	$3.15 \times 10^{-5}$
4	7680 × 7680	$1.29 \times 10^{-6}$

**Example 25.** In this example, we show the spline approximation of a highly oscillatory solution of the Poisson equation:

$$\begin{cases} -\Delta u = f(x, y) & (x, y) \in [0, 1] \times [0, 1] \\ u(x, y) = 10 \sin(x^2 + y^2) + \sin(25(x^2 + y^2)), & (x, y) \in \partial[0, 1] \times [0, 1] \end{cases}$$

Table 12. Approximation from Degree Increase (Neumann Problem)

Polynomial Degrees	Matrix Size	Maximum Errors
4	$120 \times 120$	$2.38 \times 10^{-2}$
5	$168 \times 168$	$5.53 \times 10^{-3}$
6	$224 \times 224$	$1.67 \times 10^{-4}$
7	$288 \times 288$	$1.24 \times 10^{-5}$
8	$360 \times 360$	$7.87 \times 10^{-6}$
9	$440 \times 440$	$2.52 \times 10^{-6}$
10	$528 \times 528$	$3.00 \times 10^{-7}$
11	$624 \times 634$	$4.08 \times 10^{-8}$

where  $f(x, y) = 40(\cos(x^2 + y^2) - (x^2 + y^2)\sin(x^2 + y^2)) + 50\cos(25(x^2 + y^2)) - 2500\sin(25(x^2 + y^2))$ . The exact solution contains a high frequency part which is very hard to approximate with linear finite elements. Our spline approximation yields a good approximation of such solution. In Table 13, we give the maximum errors of the spline approximation using degrees 5, 6 and 7 over uniformly refined triangulations.

Table 13. Spline Approximation of Oscillatory Solution

Levels	No. of Triangles	Degree 5	Degree 6	Degree 7
1	32	6.738	10.08	4.194
2	128	1.616	0.845	0.212
3	512	0.0391	0.0086	0.0011

## §8. Numerical Solution of Biharmonic Equations

In this section, we show how to solve biharmonic equations using multivariate splines of variable degree and variable smoothness. The biharmonic equation is given as follows:

$$\begin{cases} \Delta^2 u = f, & \text{in } \Omega \\ u = g, & \text{on } \partial\Omega \\ \frac{\partial}{\partial \mathbf{n}} u = h, & \text{on } \partial\Omega, \end{cases} \quad (41)$$

where  $\Omega$  is a polygonal domain in  $\mathbb{R}^n$ ,  $f \in L_2(\Omega)$ ,  $g$  and  $h$  are in  $C(\partial\Omega)$ , and  $\mathbf{n}$  stands for the normal direction of the boundary  $\partial\Omega$ . A typical biharmonic equation is the 2D Stokes equations in the stream function formulation (cf., e.g., [47]). The weak formulation for biharmonic equation is to find  $u \in H^2(\Omega)$  such that  $u$  satisfies the boundary conditions in (41) and

$$a(u, v) = \langle f, v \rangle, \quad \forall v \in H_0^2(\Omega),$$

where  $a(u, v)$  is a bilinear form defined by

$$a(u, v) = \int_{\Omega} \Delta u \Delta v dx$$

and  $\langle f, v \rangle = \int_{\Omega} f v dx$  stands for the standard  $L_2$  inner product of  $f$  and  $v$ . Here  $H^2(\Omega)$  and  $H_0^2(\Omega)$  are standard Sobolev spaces. With the assumption that the boundary conditions are compatible, that is, there exists a  $u_b \in H^2(\Omega)$  satisfying both boundary conditions in (41), we can show that the weak solution exists and is unique by the well-known Lax-Milgram Theorem (cf. [50]). Let

$$E_2(w) = \int_{\Omega} \left( \frac{1}{2} (\Delta w)^2 - wf \right) dx$$

be an energy functional and

$$\mathcal{A}_2 = \left\{ w \in H^2(\Omega), w = g, \frac{\partial}{\partial \mathbf{n}} w = h, \text{ on } \partial\Omega \right\}$$

be the class of admissible functions. By the compatibility of the boundary conditions, we know that  $\mathcal{A}_2$  is not empty. As before, we shall consider the following minimization problem: Find  $u \in \mathcal{A}_2$  such that

$$E_2(u) = \min\{E_2(w) : w \in \mathcal{A}_2\}.$$

Based on the standard calculus of variations, it is easy to prove that any minimizer  $u$  is a weak solution. Since the weak solution is unique, so is the minimizer.

To find an approximation of the weak solution  $u$ , we use multivariate splines of variable degree and variable smoothness. Let  $\Delta$  be a triangulation of the domain  $\Omega \subset \mathbb{R}^n$  and let

$$\mathcal{S} := S_{\mathbf{d}}^{\mathbf{r}, \rho}(\Delta)$$

be the spline space of fixed smoothness  $\rho$ ,  $\mathbf{r}$  and degree  $\mathbf{d}$  associated with  $k$ -simplices,  $0 \leq k < n - 1$ ,  $(n - 1)$ -simplices, and  $n$ -simplices of  $\Delta$ . We assume that the global smoothness  $\min\{\mathbf{r}, \rho\}$  of  $\mathcal{S}$  is bigger or equal to 1 so that  $\mathcal{S} \subset H^2(\Omega)$ . Let  $d$  be the large integer in  $\mathbf{d}$ . As the same as in the previous section, we will extend  $\mathcal{S}$  to be defined over  $\tilde{\Omega}$ . We should point out that in general, the weak solution  $u$  is smoother inside the domain  $\Omega$  (cf. [31]). Thus, we should choose  $\mathcal{S}$  such that each spline function in  $\mathcal{S}$  is more smooth than near the boundary. Let  $\tilde{\mathcal{A}}_2$  be the class of spline functions  $s \in \mathcal{S}$  satisfying the boundary conditions approximately, i.e.,  $s \in \mathcal{S}$  interpolates  $g$  at  $2d + 1$  distinct points over each curved edge and

$d + 1$  distinct points over each straight edge and  $\frac{\partial}{\partial \mathbf{n}}s$  interpolates  $h$  at  $2d - 1$  distinct points over each curved edge and  $d$  distinct points over each straight edge. Our algorithm is to find  $s_u \in \tilde{\mathcal{A}}_2$  such that

$$E_2(s_u) = \min\{E_2(s) : s \in \tilde{\mathcal{A}}_2\}.$$

More precisely, let us write any spline function  $s \in \mathcal{S}$  as

$$s(x)|_t = \sum_{\substack{\alpha \in \mathbf{z}^{n+1} \\ |\alpha|=d}} c_\alpha^t B_\alpha^t(x), \quad x \in t \in \Delta,$$

where  $d = \max\{d_t, t \in \Delta\}$ . Let  $\mathbf{c} = (c_\alpha^t, \alpha \in \mathbf{Z}^{n+1}, |\alpha| = d, t \in \Delta)$  be the coefficient vector associated with  $s$ . The smoothness and super smoothness conditions that  $s$  satisfies can be expressed by  $H\mathbf{c} = 0$ . Also,  $s$  satisfies the degree reduction conditions  $D\mathbf{c} = 0$ .

Then the bilinear form  $a(s, \hat{s})$  can be expressed in terms of  $\mathbf{c}$  and  $\hat{\mathbf{c}}$  by

$$a(s, \hat{s}) = \mathbf{c}^T K \hat{\mathbf{c}}$$

where  $K = \text{diag}(K_t, t \in \Delta)$  with

$$K_t = \left[ \int_t \Delta B_\alpha^t(x) \Delta B_\gamma^t(x) dx \right]_{\substack{\alpha \in \mathbf{z}^{n+1}, |\alpha|=d \\ \gamma \in \mathbf{z}^{n+1}, |\gamma|=d}}.$$

Note that the inner product  $\langle f, \hat{s} \rangle$  can be approximated by  $\langle s_f, \hat{s} \rangle$  for a spline  $s_f$  which interpolates  $f$  over the domain points of each  $n$ -simplex  $t$ . Thus

$$\langle f, \hat{s} \rangle \approx \hat{\mathbf{c}}^T M \mathbf{c}_f$$

where  $M = \text{diag}(M^t, t \in \Delta)$  is a block diagonal matrix with square blocks

$$M^t = \left[ \int_t B_\alpha^t(x) B_\gamma^t(x) dx \right]_{\substack{|\alpha|=d \\ |\gamma|=d}}$$

and  $\mathbf{c}_f$  is the coefficient vector for  $s_f$ . We need to solve the following minimization problem:

$$\min \frac{1}{2} \mathbf{c}^T K \mathbf{c} - \mathbf{c}^T M \mathbf{c}_f$$

subject to

$$H\mathbf{c} = 0, \quad D\mathbf{c} = 0, \quad B\mathbf{c} = \mathbf{g},$$

where  $B\mathbf{c} = \mathbf{g}$  denotes the linear system associated with the approximate boundary conditions. Note that the minimization problem has a unique

solution. Since the energy functional is convex, any local minimum is the global minimum. Let us compute a local minimum by using the Lagrange multiplier method. Letting

$$\mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = \frac{1}{2} \mathbf{c}^T K \mathbf{c} - \mathbf{c}^T M \mathbf{f} + \theta^T H \mathbf{c} + \eta^T D \mathbf{c} + \nu^T (B \mathbf{c} - \mathbf{g}),$$

we compute

$$\begin{aligned} \frac{\partial}{\partial \mathbf{c}} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) &= 0, \quad \frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = 0, \\ \frac{\partial}{\partial \eta} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) &= 0, \quad \frac{\partial}{\partial \gamma} \mathcal{L}(\mathbf{c}, \theta, \eta, \nu) = 0. \end{aligned}$$

It follows that

$$\begin{aligned} K \mathbf{c} + H^T \theta + D^T \eta + B^T \nu &= M \mathbf{c}_f \\ H \mathbf{c} = 0, D \mathbf{c} = 0, B \mathbf{c} &= \mathbf{g}. \end{aligned}$$

In other words, we need to solve the following linear system

$$\begin{bmatrix} H^T & D^T & B^T & K \\ 0 & 0 & 0 & H \\ 0 & 0 & 0 & D \\ 0 & 0 & 0 & B \end{bmatrix} \begin{bmatrix} \theta \\ \eta \\ \nu \\ \mathbf{c} \end{bmatrix} = \begin{bmatrix} M \mathbf{c}_f \\ 0 \\ 0 \\ \mathbf{g} \end{bmatrix}. \quad (42)$$

As discussed in the previous sections, the uniqueness of the weak solution implies that the matrix  $K$  is positive definite with respect to  $[H; D; B]$ . Thus, the matrix iterative method is well-defined.

To make sure that the iterative solution is the weak solution of the biharmonic equation, we use the maximum norm

$$\left\| \left[ \begin{array}{c} H \mathbf{c} \\ D \mathbf{c} \\ B \mathbf{c} - \mathbf{g} \end{array} \right] \right\|_{\infty} \quad (43)$$

to check if it does satisfy the boundary conditions, smoothness conditions, and degree reduction conditions. The solution will be said exact if it's zero in the above norm.

We remark that the above algorithm also gives a numerical method to determine if the boundary conditions are compatible or not. That is, if the least squares solution in the norm (43) is not close to zero as the underlying triangulations are refined or degrees of the spline functions increase, then the boundary conditions are not compatible since  $S_d^1(\Delta)$  becomes dense in  $H^2(\Omega)$  if  $d$  increases to  $\infty$  and/or  $|\Delta|$  decreases to 0.

We have implemented this method for 2D and 3D biharmonic equations using bivariate and trivariate spline spaces of any degree and any smoothness. That is, we are able to numerically solve biharmonic equations over any polygonal domain in the bivariate or trivariate setting. Let us present several numerical examples below.

**Example 26.** Consider a 2D biharmonic equation with exact solution  $u(x, y) = \exp(x + y)$  over a unit square domain:

$$\begin{cases} \Delta^2 u = 4 \exp(x + y) & (x, y) \in [0, 1] \times [0, 1] \\ u(x, y) = \exp(x + y), & (x, y) \in \partial[0, 1] \times [0, 1] \\ \frac{\partial}{\partial x} u(x, y) = \exp(x + y), & (x, y) \in \partial[0, 1] \times [0, 1] \\ \frac{\partial}{\partial y} u(x, y) = \exp(x + y), & (x, y) \in \partial[0, 1] \times [0, 1]. \end{cases}$$

We used the triangulation as in Fig. 2 with 25 vertices and 32 triangles and tested many spline spaces. The maximum errors of approximate weak spline solutions against the exact solution are given below. The maximum errors are computed based on  $101 \times 101$  equally-spaced points over  $[0, 1] \times [0, 1]$ . The exactness (43) is checked and is less than  $10^{-8}$  for all the spline spaces listed in Table 14.

Table 14. Numerical Approximation of Biharmonic Equation over a Standard Square Domain

	Maximum Errors	CPU Times
$S_5^1(\Delta)$	$5.7959 \times 10^{-7}$	2.8s
$S_6^1(\Delta)$	$1.1001 \times 10^{-8}$	4.5s
$S_7^1(\Delta)$	$1.3208 \times 10^{-10}$	6.2s
$S_8^1(\Delta)$	$1.1465 \times 10^{-11}$	9.8s
$S_5^2(\Delta)$	$1.1187 \times 10^{-5}$	3.5s
$S_6^2(\Delta)$	$3.2605 \times 10^{-8}$	5.2s
$S_7^2(\Delta)$	$2.7998 \times 10^{-10}$	7.9s
$S_8^2(\Delta)$	$1.1982 \times 10^{-11}$	13.2s

**Example 27.** Consider a 2D biharmonic equation with exact solution  $u(x, y) = 10 \exp(-(x^2 + y^2))$  over a unit circular domain:

$$\begin{cases} \Delta^2 u = 160 \exp(-(x^2 + y^2)) \times \\ \quad \times (x^4 + y^4 + 2x^2y^2 + 2 - 4x^2 - 4y^2) & (x, y) \in \{(x, y), x^2 + y^2 < 1\} \\ u(x, y) = 10 \exp(-(x^2 + y^2)), & (x, y) \in \{(x, y), x^2 + y^2 = 1\} \\ \frac{\partial}{\partial x} u(x, y) = -20x \exp(-(x^2 + y^2)), & (x, y) \in \{(x, y), x^2 + y^2 = 1\} \\ \frac{\partial}{\partial y} u(x, y) = -20y \exp(-(x^2 + y^2)), & (x, y) \in \{(x, y), x^2 + y^2 = 1\}. \end{cases}$$

We use the following triangulation and test many spline spaces. The maximum errors of approximate weak spline solutions against the exact solution are given in Table 15. The maximum errors are computed based on  $101 \times 101$  equally-spaced points over  $[-1, 1] \times [-1, 1]$  within the circular domain.

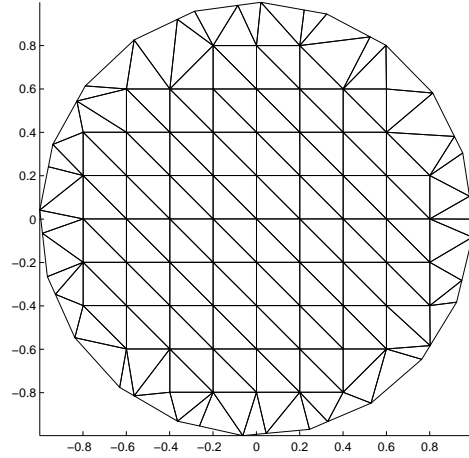


Fig. 7 A triangulation of the unit circular domain

Table 15. Numerical Approximation of Biharmonic Equation over a Circular Domain

Spline Spaces	Matrix Sizes	Maximum Errors
$S_3^1(\Delta)$	$1990 \times 1990$	$6.6819 \times 10^{-2}$
$S_4^1(\Delta)$	$2985 \times 2985$	$2.0199 \times 10^{-4}$
$S_5^1(\Delta)$	$4179 \times 4179$	$1.3653 \times 10^{-6}$
$S_6^1(\Delta)$	$5572 \times 5572$	$7.6779 \times 10^{-8}$
$S_7^1(\Delta)$	$7164 \times 7164$	$1.7841 \times 10^{-9}$
$S_8^1(\Delta)$	$8955 \times 8955$	$4.4959 \times 10^{-10}$

**Example 28.** Consider a 3D biharmonic equation with exact solution

$$u(x, y, z) = 10 \exp(-(x^2 + y^2 + z^2))$$

over an octahedron  $\Omega$  as in Example 22. We use the same tetrahedral partition as above. We find approximate weak solutions from 3D spline spaces  $S_d^1(\Delta)$  for  $d = 3, \dots, 7$  and Table 16 is a list of maximum errors against the exact solution evaluated at  $20 \times 20 \times 20$  points over  $\Omega$ .

### §9. Remarks

We have the following remarks in order:

**Remark 29.** In the theory of bivariate splines, we know that the spline spaces  $S_d^r(\Delta)$  will achieve the full approximation order as long as  $d \geq 3r+2$

Table 16. Approximation Errors from Trivariate Spline Spaces

Spline Spaces	Matrix Size	Maximum Errors
$S_3^1(\Delta)$	$160 \times 160$	0.248542
$S_4^1(\Delta)$	$280 \times 280$	0.048342
$S_5^1(\Delta)$	$448 \times 448$	0.014806
$S_6^1(\Delta)$	$672 \times 672$	0.001903
$S_7^1(\Delta)$	$960 \times 960$	0.000756

over any triangulation  $\Delta$ . For some triangulations  $\Delta$ , the approximation order will not be full when  $d < 3r + 2$ . (Cf. [12], [42] and [13]). Also, the approximation order of  $S_{3r+2}^r(\Delta)$  may be realized by super spline subspaces (cf. [16]). In general, we do not know what the approximation order of  $S_d^{r,p}(\Delta)$  is. Nevertheless many bivariate spline spaces over several special triangulations such as Clough-Tocher, Powell-Sabin, and Fraeijis de Veubek and Sander triangulation do possess the full approximation order for some  $d < 3r + 2$ . (Cf. [36], [37], [38], [41], [43], [44], [45], and [46]). If possible, one should use the spline spaces over these special triangulations to ensure the full approximation power of spline spaces. We know much less about the approximation power of trivariate spline spaces.

**Remark 30.** We have also generalized the method for numerical solution of nonlinear partial differential equations, e.g., 2D and 3D Navier-Stokes equations. We also extended the method for scattered data interpolation and fitting using nonlinear functionals, e.g.,  $L_1$  spline method. We report these numerical methods and results elsewhere. See [3], [4], [49] and [50].

**Remark 31.** In addition, we used spherical splines to treat spherical scattered data fitting problems based on the ideas in this paper. Numerical methods and results are reported in [7].

**Remark 32.** The convergence analysis of the minimal energy method, discrete least squares method, and penalized least squares method using multivariate splines were given in [27], [28], [29], and [30].

**Remark 33.** Even using the matrix iterative algorithm, the linear systems from our algorithms are still too large for problems in the trivariate setting using spline functions of higher degrees. Currently we are working on domain decomposition methods to further reduce the size of the systems.

## References

1. Alfeld, P. and L. L. Schumaker, Smooth finite elements based on Clough-Tocher triangular splits, Numer. Math. **90** (2002), 597–616.

2. Alfeld, P. and L. L. Schumaker, Smooth finite elements based on Powell-Sabin triangular splits, *Adv. Comp. Math.* **16** (2002), 29–46
3. Awanou, G., Energy Methods in 3D Spline Approximations of Navier-Stokes Equations, Ph.D. Dissertation, University of Georgia, Athens, Georgia, (2003).
4. Awanou, G. and M. J. Lai, Trivariate spline approximations of 3D Navier-Stokes equations, *Math. Comp.* **74** (2005), 585–601.
5. Awanou, G. and M.J. Lai, On convergence rate of the augmented Lagrangian algorithm for nonsymmetric saddle point problems, *Appl. Numer. Math.* **54** (2005), no. 2, 122–134.
6. Bank, R., B. D. Welfert and H. Yerentant, A class of iterative methods for solving saddle point problems, *Numer. Math.*, **56** (1990), 645–666.
7. Baramidze, V., M. J. Lai, and C. K. Shum, Spherical Splines for Data Interpolation and Fitting, to appear in *SIAM J. Scientific Computing*, (2005).
8. Bramble, J. H., J. E. Pasciak, and A. T. Vassilev, Analysis of the inexact Uzawa algorithm for saddle point problems, *SIAM J. Num. Anal.*, **34** (1997), 1072–1092.
9. Bramble, J. H., J. E. Pasciak, and A. T. Vassilev, Uzawa type algorithms for nonsymmetric saddle point problems, *Math. Comp.*, **69** (1999), 667–689.
10. Brenner, S. C. and L. R. Scott, *The mathematical theory of finite element methods*, Springer Verlag, New York, 1994.
11. deBoor, C., B-form basics, *Geometric Modeling: Algorithms and New Trends*, G. Farin (ed), SIAM Publication, 1987, 131–148.
12. de Boor, C. and K. Höllig, Approximation power of smooth bivariate pp functions, *Math. Z.*, **197** (1988), 343–363.
13. de Boor, C. and R. Q. Jia, A sharp upper bound on the approximation order of smooth bivariate pp functions, *J. Approx. Theory*, **72** (1993), 24–33.
14. Brezzi, F. and M. Fortin, *Mixed and Hybrid Finite Element Methods*, Springer Verlag, New York, 1991.
15. Chui, C. K. *Multivariate Splines*, SIAM Publication, Philadelphia, 1988.
16. Chui, C. K. and M. J. Lai, On bivariate super vertex splines, *Constr. Approx.*, **6** (1990), 399–419.
17. Chui, C. K. and M. J. Lai, Multivariate vertex splines and finite elements, *J. Approx. Theory*, **60** (1990), 245–343.
18. Chui, C. K. and M. J. Lai, Filling polygonal holes using  $C^1$  cubic triangular spline patches, *Comp. Aided Geom. Design.*, **17** (2000), 297–307.

19. Ciarlet, F., *The Finite Element Method for Elliptic Problems*, North-Holland, Amsterdam, New York, 1978.
20. Dierckx, P., *Curve and Surface Fitting with Splines*, Oxford University Press, Oxford, 1985.
21. Elman, H. C. and G. H. Golub, Inexact and preconditioned Uzawa algorithms for saddle point problems, *SIAM J. Numer. Anal.*, (1994), 1645–1661.
22. Evans, L., *Partial Differential Equations*, American Math. Soc., Providence, 1998.
23. Farmer, K. W. and M. J. Lai, Scattered data interpolation by  $C^2$  quintic wplines using energy minimization, in *Approximation Theory IX: Computational Aspects* Charles K. Chui and Larry L. Schumaker (eds.) Vanderbilt University Press (Nashville), 1998, 47–54.
24. Farin, G., Triangular BernsteinBézier patches, *Comput. Aided Geom. Design* **3** (1986), 83–127.
25. Fasshauer, G. and L. L. Schumaker, Multi-patch parametric surfaces with minimal energy, *Comp. Aided Geom. Design*, **13** (1996), 45–79.
26. Franke, R., 1982, Scattered data interpolation: tests of some methods, *Math. Comp.*, **38** (1982), 181–200.
27. von Golitschek, M., M. J. Lai, L. L. Schumaker, Bounds for minimal energy bivariate polynomial splines, *Numer. Math.* **93** (2002), 315–331
28. von Golitschek, M. and L. L. Schumaker, Penalized Least Square Fitting, *Algorithms for Approximation II*, J. Mason and M. G. Cox (eds.), Chapman and Hall, London (1990), 210–227.
29. von Golitschek, M. and L. L. Schumaker, Bounds on projections onto bivariate polynomial spline spaces with stable local bases, *Constr. Approx.* **18** (2002), 241–254.
30. von Golitschek, M. and L. L. Schumaker, Penalized least squares fitting, *Serdica* **18** (2002), 1001–1020.
31. Grisvard, P., *Elliptic Problems in Nonsmooth Domains*, Pitman, Boston, 1985.
32. Gunzburger, M. D., *Finite Element Methods for Viscous Incompressible Flows*, Academic Press, Boston, 1989.
33. Laghchim-Lahlou, M and P. Sablonnière, Composite quadrilateral finite elements of Class  $C^r$ , in *Mathematical Methods in Computer Aided Geometric Design*, edited by T. Lyche and L. L. Schumaker, Academic Press, New York, 1989, 413–418.
34. Laghchim-Lahlou, M. and P. Sablonnière, Eléments finis polynomiaux composés de classe  $C^r$ , *C. R. Acad. Sci., série I*, **136** (1993), 503–508.

35. Laghchim-Lahlou, M. and P. Sablonnière, Quadrilateral finite elements of FVS type and class  $C^r$ , *Numer. Math.*, **70** (1995), 229–243.
36. Lai, M. J., 1994, Approximation order from bivariate  $C^1$  cubics on a four-directional mesh is full, *Computer Aided Geometric Design*, **11** (1994), 215–223.
37. Lai, M. J., Scattered data interpolation and approximation by  $C^1$  piecewise cubic polynomials, *Computer Aided Geometric Design*, **13** (1996), 81–88.
38. Lai, M. J., On  $C^2$  quintic spline functions over triangulations of Powell-Sabin's type, *J. Computational and Applied Mathematics*, **73** (1996), 135–155.
39. Lai, M. J., Geometric interpretation of smoothness conditions of triangular polynomial patches, *Computer Aided Geometric Design*, **14** (1997), 191–199.
40. Lai, M. J., Convex preserving scattered data interpolation using bivariate  $C^1$  cubic splines, *J. Comput. Applied Math.*, **119** (2000), 249–258.
41. Lai, M. J. and L. L. Schumaker, Scattered data interpolation using piecewise polynomials of degree six, *SIAM Numer. Anal.*, **34** (1997), 905–921.
42. Lai, M. J. and L. L. Schumaker, Approximation power of bivariate splines, *Advances in Comput. Math.*, **9** (1998), 251–279.
43. Lai, M. J. and L. L. Schumaker, On the approximation power of splines on triangulated quadrangulations, *SIAM Numerical Analysis*, **36** (1999), 143–159.
44. Lai, M. J. and L. L. Schumaker, Macro-Elements and stable local bases for splines on Clough-Tocher triangulations, *Numer. Math.*, **88** (2001), 105–119.
45. Lai, M. J. and L. L. Schumaker, Quadrilateral Macro-Elements, *SIAM J. Math. Anal.* **33** (2002), 1107–1116.
46. Lai, M. J. and L. L. Schumaker, Macro-elements and stable local bases for splines on Powell-Sabin triangulations, *Math. Comp.* **72** (2003), 335–354
47. Lai, M. J. and P. Wenston, Bivariate spline method for numerical solution of Navier-Stokes equations over polygons in stream function formulation, *Numerical Methods for P.D.E.*, **16** (2000), 147–183.
48. Lai, M. J. and P. Wenston, Trivariate  $C^1$  cubic splines for numerical solution of biharmonic equation, in: *Trends in Approximation Theory*, K. Kopotun, T. Lyche, and M. Neamtu (eds.), Vanderbilt University Press, Nashville, 2001, 224–234.

49. Lai, M. J. and P. Wenston,  $L_1$  spline methods for scattered data interpolation and approximation, *Adv. Comp. Math*, **21** (2004), 293–315.
50. Lai, M. J. and P. Wenston, Bivariate splines for fluid flows, *Computer & Fluids* **33** (2004), 1047–1073.
51. Le Mehauté, A., Interpolation et approximation par des fonctions polynomiales the computed solutions from the domain decomposition method. par morceaux dans  $\mathbb{R}^n$ , Ph.D. Thesis, Univ. Rennes, France, 1984.
52. Rescorla, K. L., Cardinal interpolation: a bivariate polynomial example, *Comput. Aided Geom. Design*, **3** (1986), 313–321
53. Rusten, T. and R. Winther, A preconditioned iterative method for saddlepoint problems, *SIAM J. Matrix Anal. App.*, **13** (1992), 887–904.
54. Segal, A., On the numerical solution of the Stokes equations using the finite element method, *Computer Methods in Applied Mechanics and Engineering*, **19** (1979), 165–185.
55. Späth, H, Two Dimensional Spline Interpolation Algorithms, A. K. Peters, Wellesley, 1995.
56. Wahba, G., *Splines for Observation Data*, SIAM Pub., Philadelphia, 1990.
57. Worsey, A. J. and G. Farin, An  $n$ - dimensional Clough-Tocher interpolant, *Constr. Approx.*, **3** (1987), 99–110.
58. Ženiček, A., Polynomial approximation on tetrahedrons in the finite element method, *J. Approx. Theory*, **7** (1973), 34–351.

Gerard Awanou  
Department of Mathematical Sciences,  
DeKalb, IL, 60115,  
Northern Illinois University  
awanou@math.niu.edu

Ming-Jun Lai  
Department of Mathematics,  
The University of Georgia, Athens, GA 30602  
mjlai@math.uga.edu

Paul Wenston  
Department of Mathematics,  
The University of Georgia, Athens, GA 30602  
paul@math.uga.edu