

MATHEMATICAL EQUATION MODELS

Wen Hsiang Wei^{†1}

Department of Statistics, Tung Hai University, Taiwan

ABSTRACT

A class of models involving mathematical equations for fitting the data is proposed. The class of models consists of some commonly used statistical models, such as linear regression models, nonparametric regression models, linear mixed-effects models, and measurement error models. The equation of interest can be also a partial differential equation. Nonlinear programming methods can be used to estimate the underlying equation. Theoretical results for the methods of estimation are established. A simulation study and a modified example in thermodynamics are used to illustrate the proposed models and associated methods of estimation.

Key words and phrases: Mathematical equation models, Nonlinear programming, Partial differential equations, Penalty function methods, Reproducing kernel Hilbert space.
JEL classification: C61

[†]Correspondence to: Wen Hsiang Wei
E-mail: wenwei@thu.edu.tw

1. Introduction

Mathematical equations play a pivotal role in scientific research. Two classes of mathematical equations, nonlinear equations and differential equations, are commonly used. In the following examples, the data with means satisfying certain mathematical equations are analyzed and the corresponding statistical estimation problems are discussed.

In standard nonparametric regression setting, the mapping between the means of the covariates and responses is point to point. However, the means of the response and covariate variables might satisfy a nonlinear equation such that the point to point mapping is no longer true. A simple example of such nonlinear equations is the conic equation. Let the data $\mathbf{y}_{ij} = (y_{ij1}, y_{ij2})^t$, $i = 1, \dots, 629$, $j = 1, \dots, n_i$, generated from a bivariate normal distribution with mean vectors $\boldsymbol{\mu}_{y_i} = (\mu_{y_{i1}}, \mu_{y_{i2}})^t$ and identity variance-covariance matrix \mathbf{I} , where n_i is the number of repeated observations at site i and equal to 1 in the example. The data are shown in Figure 1. The mean vectors $\boldsymbol{\mu}_{y_i}$ satisfy the conic equation $F(\boldsymbol{\mu}_{y_i}) = 28\mu_{y_{i1}}^2 + 28\mu_{y_{i2}}^2 + 52\mu_{y_{i1}}\mu_{y_{i2}} - 162 = 0$, where F is a function defined on \mathbf{R}^2 . Note that y_{ij1} and y_{ij2} can be considered as the observations corresponding to the covariate and response variables, respectively. The fitted regression line along with the fits by other commonly used statistical methods, including polynomial regression, kernel smoother, and smoothing spline, are shown in Figure 1. Since the model assumption for these methods is $\mu_{y_{i2}} = f(\mu_{y_{i1}})$, these methods fail to discover the underlying equation, where f is a function defined on R . On the other hand, the blue dots are the fitted values based on the proposed models and associated method of estimation introduced in next section, which approximate the true values (the blue solid line) generated by the underlying equation well. In addition to the above equation, the pressure equation for helium at $273.15^\circ K$ corresponding to the equation of state for a gas in thermodynamics (see Britt and Luecke, 1973, Section 10) has the form

$$c_1\mu_{y_{i1}}^2\mu_{y_{i2}} + c_2\mu_{y_{i1}}\mu_{y_{i2}}^2 + c_3\mu_{y_{i1}}\mu_{y_{i2}} + c_4\mu_{y_{i1}} + c_5\mu_{y_{i2}} = 0,$$

where $\mu_{y_{i1}}$ and $\mu_{y_{i2}}$ are the mean pressures of the $(i - 1)$ th and the i th expansions and c_1, \dots, c_5 are some constants. The exact solutions of the equation (the blue solid line) based on the results given in Table 4 of Britt and Luecke (1973) using complete algorithm along with the data (the black points) generated by normal random variables with means satisfying the equation and coefficients of variation equal to 10%, and the fitted values (the blue dots) by the proposed models and associated method of estimation are given in the left part of Figure 2. The fitted values approximate the true values reasonably well. Furthermore, the encouraging results will be also obtained in Section 4.2 as considering the pressure equation for methane at $131.93^\circ K$ and both the data with larger variations and the real data given in Blancett, Hall and Canfield (1970) and Hoover (1965).

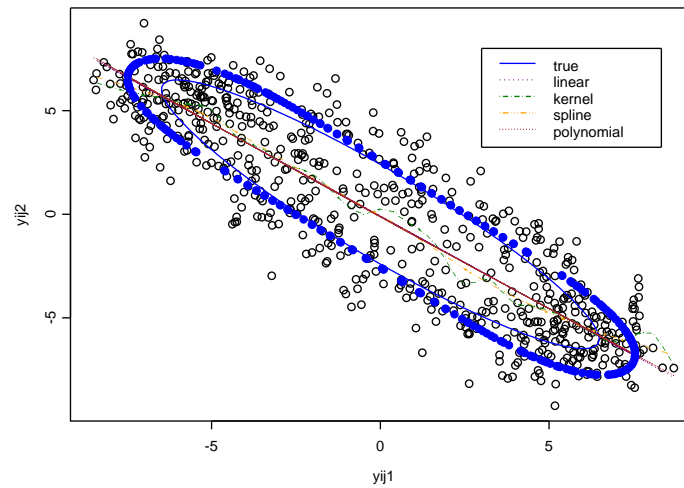


Fig. 1. The data with the means satisfying the conic equation.

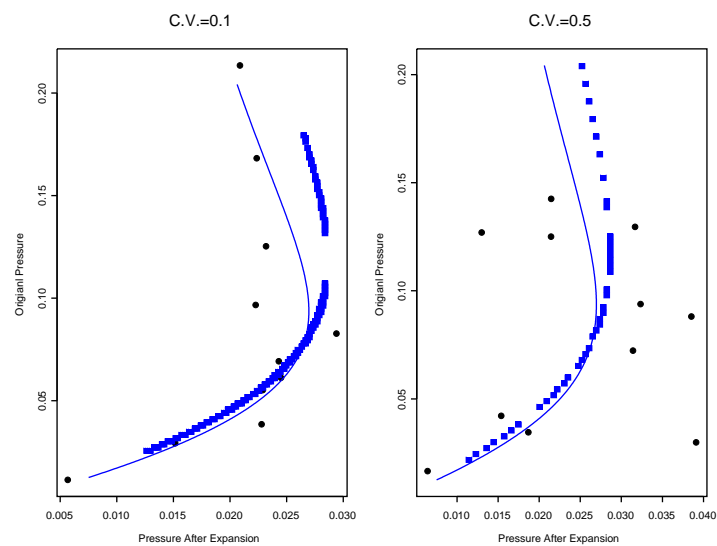


Fig. 2. The pressure data for helium at 273.15°K: Observed data (black ●); Fitted values (blue dot); True equation (blue line).

Partial differential equations (PDEs) are one of the intensively studied areas in mathematics. Therefore, the other example is the wave equation, which is a partial differential equation. The data $\mathbf{y}_{ij} = (y_{ij1}, y_{ij2})^t$, $i = 1, \dots, 63$, $j = 1, \dots, n_i$, are generated from a bivariate normal distribution with mean vectors $\boldsymbol{\mu}_{y_i} = (\mu_{y_{i1}}, \mu_{y_{i2}})^t$ and variance-covariance matrix $0.2^2 \mathbf{I}$, where n_i is also equal to 1 in this example. The data $y_{ij(3)}$ are generated from a normal distribution with means $F(\mu_{y_{i1}}, \mu_{y_{k2}})$ and variances equal to 0.2^2 , and $F(\mu_{y_{i1}}, \mu_{y_{k2}}) = 7.5 \cos(\mu_{y_{i1}} - 2\mu_{y_{k2}})$ is the underlying wave function, where $k = 1, \dots, 63$. The upper left part of Figure 3 gives the wave function. The mean vectors also satisfy the partial differential equation

$$\partial^2 F(\mu_{y_{i1}}, \mu_{y_{i2}}) / \partial \mu_{y_{i2}}^2 = 4 \partial^2 F(\mu_{y_{i1}}, \mu_{y_{i2}}) / \partial \mu_{y_{i1}}^2.$$

The plots in Figure 3 give the fitted functions based on the observed data $y_{ij(3)}$ with means equal to $F(\mu_{y_{i1}}, \mu_{y_{i2}})$, one incorporating with the partial differential equation and the other not. The proposed method of estimation incorporating with the partial differential equation could provide a sensible fit for the data. On the other hand, the fit without using the partial differential equation might not be sensible. In addition, if the variance-covariance matrix of \mathbf{y}_{ij} is equal to $2^2 \mathbf{I}$, the variances of the data $y_{ij(3)}$ displayed in the lower right part of Figure 3 are equal to 2^2 , and the initial conditions for the partial differential equation are known, the proposed method provides a very accurate fit even for the data with larger deviations, as will be illustrated in Section 4.1.

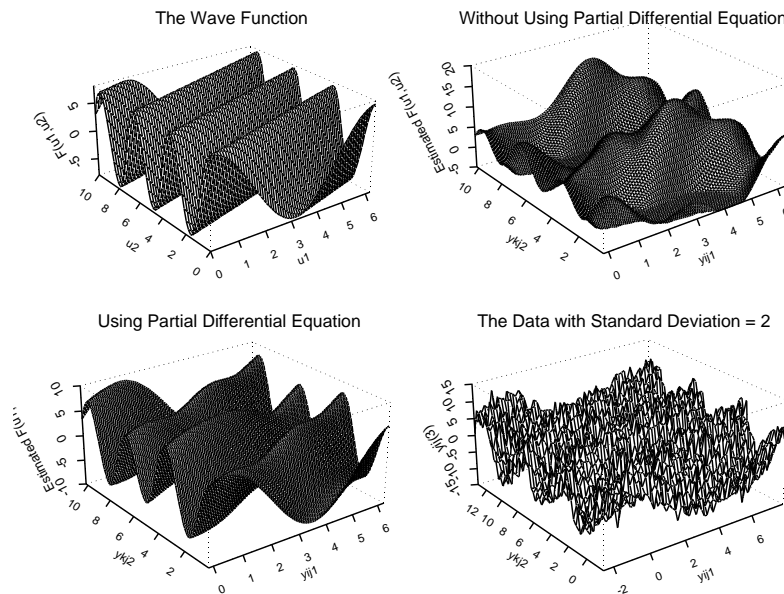


Fig. 3. The data with the means satisfying the wave equation.

The problem of estimating the unknown function in the first example is related to the one of estimating the parameters of nonlinear implicit functional models (see Britt and Luecke, 1973) and the parameters involved in the known nonlinear implicit functional relationship were of interest. The nonlinear implicit models could be employed to fit the data in chemical industry. However, relatively little has been done for the estimation of the unknown nonlinear implicit functional relationship itself (i.e., nonparametric implicit functional models), which is one of the goals of this article. On the other hand, the statistical inference related to the partial differential equations has not attracted much attention. Cavalier (2011) related to the estimation of the function given in the initial condition of the heat equation under the framework of statistical inverse problems. Nevertheless, the range of applications of the partial differential equations is enormous, for examples, astronomy, dynamics, elasticity, heat transfer, electromagnetic theory, quantum mechanics, and so on. Since the observed data might be subject to random errors, it might be reasonable to estimate either the solutions or the PDEs based on the statistical modeling. Therefore, another goal of this article is to model the PDEs of interest out of data and then estimate the corresponding solutions. If the PDEs depend on some unknown parameters or functions, the goal is then to estimate both these parameters or functions and the solutions. Above all, this article is to propose two classes of statistical models, one defined by a nonlinear equation and the other involving the PDEs, and to establish theoretical results for the methods of estimation. It turns out that the algorithms for nonlinear programming problems can be used to estimate the underlying equation. Nonlinear programming methods have been widely used in statistics (see Thisted, 1988, Chapter 4). In next section, two class of models incorporating the nonlinear equations and the differential equations with random errors are proposed. The associated estimators based on the nonlinear programming methods are also given in the section. The convergence results for the proposed methods of estimation are presented in Section 3. In Section 4, a simulation study is conducted to evaluate the proposed models and methods. Besides, a modified example in thermodynamics is given in the section. A concluding discussion is given in Section 5. Finally, some routine derivations used in Section 2.2, the proofs of the theoretical results in Section 2 and Section 3, and additional applications, including the models and methods of estimation for the correlated data, general constraints, system of equations, equation selection, and convergence of algorithms characterized by different maps, are delegated to the supplementary materials, which can be found at

<http://web.thu.edu.tw/wenwei/www/papers/jcsaSupplement.pdf/> .

Hereafter, the notation $\|\cdot\|_V$ is denoted as the norm of the normed space V . As V is a Hilbert space, the norm induced by the inner product is $\|\cdot\|_V = (\langle \cdot, \cdot \rangle_V)^{1/2}$. In addition, the Euclidean norm is used for \mathbf{R}^q , where q is a positive integer.

2. Mathematical Equation Models

The statistical models involving the nonlinear functional relationship and the differential equations have been explored in the literature, as indicated in the previous section. Therefore, two classes of statistical models involving mathematical equations for fitting the data given in the examples of Section 1 are proposed. The one involving a nonlinear equation is referred to as ordinary equation models, while the other involving a differential equation, possibly a partial differential equation, is referred to as differential equation models. The two classes of models fall in a broad class of models, referred to as mathematical equation models. Nevertheless, the general concept of the mathematical equation models is the models involving two essential ingredients: one is the operation or relation and the other is the statistical quantity related to the observed data. Therefore, other statistical models involving different types of operations or quantities of interest can be proposed, for some examples, the functions F being the solutions of integral equations or the parameters of interest being the variances of the observed data. These models can be also included in the family of the mathematical equation models.

2.1 Ordinary Equation Models

Intuitively, an ordinary equation model is a way of expressing the implicit relation of the mean vectors of some random vectors. Further, the implicit relation is the main interest. That is, given the random vectors $\mathbf{Y} = (Y_1, \dots, Y_p)^t$ with mean vectors $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^t$, the unknown function F satisfying $F(\boldsymbol{\mu}) = 0$ is of interest. A more formal statement of the ordinary equation models is given below.

Definition 2.1. Let V_Y be a collection of random vectors $\mathbf{Y} = (Y_1, \dots, Y_p)^t$ with mean vectors $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^t$ and $M_0 = \{\boldsymbol{\mu} : E(\mathbf{Y}) = \boldsymbol{\mu}, \mathbf{Y} \in V_Y\}$ be a subset of M , where M is a subset of \mathbf{R}^p . Let V_f be a subset of some normed space of real-valued functions defined on M and V_f^0 , the subset of V_f , be the set of functions F satisfying $F(\boldsymbol{\mu}) = 0$, $\boldsymbol{\mu} \in M_0$, i.e., V_f^0 is the set of "null" functions with respect to M_0 . An ordinary equation model is denoted by (V_Y, V_f, V_f^0) provided that V_f^0 is nonempty. If there exists a unique function F of which normed value equal to one and V_f^0 is the nonempty subset of the space spanned by F , the model is referred to as the unique ordinary equation model with respect to M_0 and M .

The ordinary equation models include some commonly used statistical models, such as linear regression models, nonparametric regression models, linear mixed-effects models, and measurement error models. In standard (conditional) linear and nonparametric regression models, $F(\boldsymbol{\mu}) = \mu_p - \beta_0 - \beta_1\mu_1 - \dots - \beta_{p-1}\mu_{p-1}$ and $F(\boldsymbol{\mu}) = \mu_p - f(\mu_1, \dots, \mu_{p-1})$, respectively, where $\beta_0, \dots, \beta_{p-1}$ are the parameters, Y_p is the response variable, and Y_1, \dots, Y_{p-1} are degenerated random variables. In linear and nonparamet-

ric measurement error models (or unconditional linear and nonparametric models), the functions F are equal to the ones in linear and nonparametric regression models but the random variables Y_1, \dots, Y_{p-1} are not degenerated. In linear mixed-effects models, the function associated with the i th observation is $F(\boldsymbol{\mu}) = \mu_p - \mu_1 - x_{i1}\mu_2 - \dots - x_{i(p-2)}\mu_{p-1}$, where $x_{i1}, \dots, x_{i(p-2)}$ are the observed values of covariates.

The uniqueness of the ordinary equation models relies on the choices of V_f . For a simple example, let $M_0 = \{-1, 1\}$ and V_f be the vector space consisting of all polynomials defined on $M = R$. Then, by the fundamental theorem of algebra, V_f^0 consists of the polynomials of the form $F(\mu) = c(\mu)(\mu - 1)^a(\mu + 1)^b$, where a, b are positive integers and $c(\mu)$ is any polynomial function. Thus, (V_Y, V_f, V_f^0) is not unique. However, if V_f is the vector space consisting of all polynomials with degrees less or equal to 2, (V_Y, V_f, V_f^0) is unique. In addition, the domain M also plays a crucial role in determining the uniqueness of the ordinary equation models. In the above example, if $M = M_0$, i.e., V_f is the vector space consisting of all polynomials defined on M_0 , then (V_Y, V_f, V_f^0) is also unique since all polynomials of the form $F(\mu) = c(\mu)(\mu - 1)^a(\mu + 1)^b$ are equal.

In this article, consider that V_f is a real Hilbert space with a norm $\|\cdot\|_{V_f}$ induced by the inner product on V_f . If F is considered as the minimizer of a specified objective functional, the following theorem indicates that the minimizer exists and falls in a finite dimensional subspace of V_f .

Theorem 2.1. *Let the objective functional be*

$$S(F) = \frac{1}{m} \sum_{i=1}^m (\langle F, \eta_{\mu_{y_i}} \rangle_{V_f})^2 + c \|P_{H^\perp}(F)\|_{V_f}^2,$$

where H is a finite dimensional subspace of V_f , $\eta_{\mu_{y_i}} \in V_f$ are the representers associated with $\mu_{y_1}, \dots, \mu_{y_m}$, the means of some random vectors, c is a positive constant, and P_{H^\perp} is a projection operator of V_f onto the orthogonal complement of H . Then, the minimizer

$$\hat{F} = \arg \min_{F \in V_f, \|F\|_{V_f}=1} S(F),$$

exists and has the form of $\sum_{l=1}^q \beta_l^* \psi_l$, where β_l^* are the coefficients and ψ_l are the basis functions of some finite dimensional subspace of V_f .

Remark 2.1. *In the above theorem, $\langle F, \eta_{\mu_{y_i}} \rangle_{V_f}$ related to $F(\mu_{y_i})$ provides the "quantitative" information about the fidelity of the function F to the data if the value of F evaluated at μ_{y_i} or its approximation exists. On the other hand, the term $\|P_{H^\perp}(F)\|_{V_f}^2$ could be associated with the smoothness of the functions of interest, as in nonparametric curve estimation using spline functions (see Berlinet and Thomas-Agnan, 2004, Chapter 3). In addition, if the space H is the space spanned by $\eta_{\mu_{y_i}}$, this term can rule out the "information" provided by the functions orthogonal to $\eta_{\mu_{y_i}}$. Intuitively, it means that only the "information" associated with the observations will be adopted.*

Therefore, ideally, the minimizer \hat{F} of the objective functional $S(F)$ results in both small values of $(\langle F, \eta_{\mu_{y_i}} \rangle_{V_f})^2$ and $\|P_{H^\perp}(F)\|_{V_f}^2$, respectively, i.e., the fidelity of F to the data reflected by the small value of $(\langle F, \eta_{\mu_{y_i}} \rangle_{V_f})^2$ and accurate approximation by the function in the finite dimensional space H reflected by the small value of $\|P_{H^\perp}(F)\|_{V_f}^2 = \|F - P_H(F)\|_{V_f}^2$. If $H = V_f$, i.e., $\|P_{H^\perp}(F)\|_{V_f}^2 = 0$, the minimizers are any elements in the subspace of V_f orthogonal to the proper subspace of V_f spanned by $\eta_{\mu_{y_i}}$.

If V_f has a reproducing kernel defined on $M \times M$ (see Aronszajn, 1950; Berlinet and Thomas-Agnan, 2004), then the pointwise value of F at μ exists and the results given in Theorem 2.1 hold, as indicated by the following corollary.

Corollary 2.1. *Let V_f be a Hilbert space with a reproducing kernel $K(\cdot, \cdot)$ defined on $\times M$. Then, as*

$$S(F) = \frac{1}{m} \sum_{i=1}^m F^2(\mu_{y_i}) + c \|P_{H^\perp}(F)\|_{V_f}^2,$$

the results given in Theorem 2.1 hold.

If $|F(\mu)| \leq c_\mu \|F\|_{V_f}$ for all F in V_f and all μ in M , V_f has a reproducing kernel (see Aubin, 2000, Theorem 5.9.1), where $c_\mu \geq 0$ depends on μ . An example of V_f is the completion of the tensor product of Paley-Wiener spaces (see Berlinet and Thomas-Agnan, 2004, p. 31, p. 304).

By the above theorem and for the purpose of computations, assume that the underlying function $F(\mu) = \sum_{l=1}^q \beta_l^* \psi_l(\mu)$, where F is either the minimizer given in Theorem 2.1 with $H = V_f$ and $\langle F, \eta_{\mu_{y_i}} \rangle_{V_f} = F(\mu_{y_i})$, or the function satisfying $F(\mu) = 0$ for $\mu \in M_0$, i.e., $F \in V_f^\circ$, and where $\psi_l \in V_f$ are some known basis functions defined on M . In the latter case, the coefficients β_l^* are referred to as the true coefficients. Thus, the finite dimensional optimization methods can be employed to estimate F . These basis functions can be used to approximate or generate the functions in V_f . For example, if V_f consists of all square-integrable functions, the orthogonal polynomials or the tensor product of the orthogonal polynomials can be used as the basis functions. For ease of exposition, let ψ_l be orthonormal and $\|F\|_{V_f} = 1$. Thus, the coefficient vector $\beta^* = (\beta_1^*, \dots, \beta_q^*)^t$ satisfies $\sum_{l=1}^q (\beta_l^*)^2 = 1$. Note that the imposed constraint mainly depends on the choices of the basis functions and the norm of F . In addition, assume that ψ_l are continuous functions defined on M .

Suppose that the independently observed data

$$\mathbf{y}_{ij} = (y_{ij1}, \dots, y_{ijp})^t, i = 1, \dots, m, j = 1, \dots, n_i,$$

come from the distributions with mean vectors μ_{y_i} , where m is the number of sites and $\sum_{i=1}^m n_i = n$. To estimate the function F , the vector $\hat{\beta}(\mathbf{T}_n)$ minimizing the mean sum

of squares,

$$\begin{aligned} S(\beta | \mathbf{T}_n) &= \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{n_i} \sum_{j=1}^{n_i} F^2[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\} \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{n_i} \left\{ \sum_{l=1}^q \beta_l n_i^{-1/2} \psi_l[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\}^2, \end{aligned}$$

subject to the constraint $\sum_{l=1}^q \beta_l^2 = 1$ needs to be obtained, where

$$\mathbf{T}_n = [\mathbf{T}_{11}^t(\mathbf{y}_1), \dots, \mathbf{T}_{1n_1}^t(\mathbf{y}_1), \dots, \mathbf{T}_{m1}^t(\mathbf{y}_m), \dots, \mathbf{T}_{mn_m}^t(\mathbf{y}_m)]^t,$$

$\beta = (\beta_1, \dots, \beta_q)^t$, $\mathbf{y}_i = (\mathbf{y}_{i1}^t, \dots, \mathbf{y}_{in_i}^t)^t$, and $\mathbf{T}_{ij}(\mathbf{y}_i)$ are sensible estimates of μ_{y_i} , for examples, the mean vectors $\mathbf{T}_{ij}(\mathbf{y}_i) = \sum_{j=1}^{n_i} \mathbf{y}_{ij}/n_i$ or $\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{y}_{ij}$. Both the trivial solution $\beta = \mathbf{0}$ and the other functions with normed values not equal to 1 in V_f^0 can be excluded by the constraint $\sum_{l=1}^q \beta_l^2 = 1$. The objective function $S(\beta | \mathbf{T}_n)$ can be considered as the estimator of the following objective function for the underlying mean vectors,

$$S^*(\beta | \mu_y) = \frac{1}{m} \sum_{i=1}^m F^2(\mu_{y_i}) = \beta^t [\Psi^*(\mu_y)]^t \Psi^*(\mu_y) \beta,$$

where $\Psi^*(\mu_y) = [m^{-1/2} \psi_l(\mu_{y_i})]$ is an $m \times q$ matrix of which rows are

$$[m^{-1/2} \psi_1(\mu_{y_i}), \dots, m^{-1/2} \psi_q(\mu_{y_i})],$$

and where $\mu_y = (\mu_{y_1}^t, \dots, \mu_{y_m}^t)^t$.

The search of the minimizer $\hat{\beta}(\mathbf{T}_n)$ can be considered as a nonlinear programming problem with a constraint, i.e., the search for the minimizer of a given function and a constraint imposed on the candidate solutions. To use the nonlinear programming algorithms available for unconstrained problems, the original constrained problem needs to be transformed into an unconstrained problem. The penalty function methods (see Bazaraa and Shetty, 1979, Chapter 9; Nocedal and Wright, 1999, Chapter 17) are commonly employed for the transformation. The penalty function multiplied by a positive penalty parameter λ can be added to the original objective function. Given some regularity conditions, the solutions of the transformed unconstrained problem could converge to the one of the original constrained problem as λ tends to infinity. The unconstrained objective function corresponding to $S(\beta | \mathbf{T}_n)$ is

$$\begin{aligned} S_\lambda(\beta | \mathbf{T}_n) &= S(\beta | \mathbf{T}_n) + \lambda \left(\sum_{l=1}^q \beta_l^2 - 1 \right)^2 \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{n_i} \left\{ \sum_{l=1}^q \beta_l n_i^{-1/2} \psi_l[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\}^2 + \lambda \left(\sum_{l=1}^q \beta_l^2 - 1 \right)^2, \end{aligned}$$

where $(\sum_{l=1}^q \beta_l^2 - 1)^2$ is the penalty function (see Bazaraa and Shetty, 1979, pp. 332-333). Note that other penalty functions can be also employed. The objective function in matrix form is

$$S_\lambda(\beta \mid \mathbf{T}_n) = \beta^t [\Psi(\mathbf{T}_n)]^t \Psi(\mathbf{T}_n) \beta + \lambda (\beta^t \beta - 1)^2, \quad (1)$$

where $\Psi(\mathbf{T}_n) = \{m^{-1/2}n_i^{-1/2}\psi_l[\mathbf{T}_{ij}(\mathbf{y}_i)]\}$ is an $n \times q$ matrix of which rows are

$$\left\{ m^{-1/2}n_i^{-1/2}\psi_1[\mathbf{T}_{ij}(\mathbf{y}_i)], \dots, m^{-1/2}n_i^{-1/2}\psi_q[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\}.$$

The minimizer of the objective function given in expression (1) is denoted as $\hat{\beta}_\lambda(\mathbf{T}_n)$. Several theoretical results concerning the consistency and asymptotic normality of the estimator $\hat{\beta}_\lambda(\mathbf{T}_n)$ are given in next section. The objective function $S_\lambda(\beta \mid \mathbf{T}_n)$ can be considered as the estimator of the following objective function for the underlying mean vectors,

$$\begin{aligned} S_\lambda^*(\beta \mid \mu_y) &= S^*(\beta \mid \mu_y) + \lambda (\beta^t \beta - 1)^2 \\ &= \beta^t [\Psi^*(\mu_y)]^t \Psi^*(\mu_y) \beta + \lambda (\beta^t \beta - 1)^2. \end{aligned} \quad (2)$$

The minimizers $\hat{\beta}^*(\mu_y)$ and $\hat{\beta}_\lambda^*(\mu_y)$ of the objective functions $S^*(\beta \mid \mu_y)$ subject to the constraint $\beta^t \beta = 1$ and the objective function $S_\lambda^*(\beta \mid \mu_y)$ given in expression (2), respectively, are equal to β^* , the vector of the true coefficients, provided that $M_0 = \{\mu_{y_i} : i = 1, \dots, m\}$, $F(\mu) = 0$ for $\mu \in M_0$, and the ordinary equation model is unique. If there are multiple minimizers of the above constrained objective function and $F \in V_f^\circ$, one of them is equal to the true coefficient vector for the ordinary equation model.

2.2 Differential Equation Models

If the function of interest F is a solution of a known differential equation or a differential equation depending on some unknown parameters or functions, the associated differential equation model is described below.

Definition 2.2. Let V_Y be a collection of $p \times 1$ random vectors \mathbf{Y} with mean vectors μ and $M = \{\mu : E(\mathbf{Y}) = \mu, \mathbf{Y} \in V_Y\}$ be an open subset of \mathbf{R}^p . A differential equation model is a model with two ingredients: one is the set M and the other is an equation involving the derivatives of an unknown real or complex function F on M . If $p \geq 2$, the differential equation model is referred to as the partial differential equation model.

Let $V_{\partial Y}$ be a collection of $p \times 1$ random vectors \mathbf{Y} with mean vectors $\tilde{\mu}$ and $\partial M = \{\tilde{\mu} : E(\mathbf{Y}) = \tilde{\mu}, \mathbf{Y} \in V_{\partial Y}\}$ be a set of points corresponding to initial and boundary conditions of the differential equations of interest. The issues of existence and uniqueness of the solutions of the partial differential equations are not completely

settled in mathematics. It is very crucial to prove the existence of the solutions of the partial differential equations of interest. Several methods can be employed to prove the existence of the solutions (see Jost, 2002). In addition, very few results existed for imposing the boundary conditions to determine a unique solution in a general setting (see Chester, 1970, Chapter 6-11). In this article, the equations with existed solutions or a unique solution are of interest.

There are several criteria for classifying PDEs (see Jost, 2002, pp. 4-6). One of the criteria is the order of the highest-occurring derivatives. For example, a second order PDE with $p = 2$ is

$$D \left[\boldsymbol{\mu}, F(\boldsymbol{\mu}), \frac{\partial F(\boldsymbol{\mu})}{\partial \mu_1}, \frac{\partial F(\boldsymbol{\mu})}{\partial \mu_2}, \frac{\partial^2 F(\boldsymbol{\mu})}{\partial \mu_1^2}, \frac{\partial^2 F(\boldsymbol{\mu})}{\partial \mu_1 \partial \mu_2}, \frac{\partial^2 F(\boldsymbol{\mu})}{\partial \mu_2 \partial \mu_1}, \frac{\partial^2 F(\boldsymbol{\mu})}{\partial \mu_2^2} \right] = 0,$$

where D is a real or complex function defined on the subset of $M \times \mathbf{K}^7$, where K is the scalar field, either the field of real numbers or the field of complex numbers. Hereafter, consider that the field of real numbers is of interest. Suppose that the underlying equation of interest is a d th order partial differential equation and the solution F falls in a real Hilbert space V_f . Further, suppose that the underlying d th order partial differential equation can be expressed as

$$\begin{aligned} & D \left[\boldsymbol{\mu}, F(\boldsymbol{\mu}), \frac{\partial F(\boldsymbol{\mu})}{\partial \mu_1}, \dots, \frac{\partial^{|\boldsymbol{\alpha}|} F(\boldsymbol{\mu})}{\partial \mu_{j_1}^{\alpha_1} \dots \partial \mu_{j_r}^{\alpha_r}}, \dots, \frac{\partial^d F(\boldsymbol{\mu})}{\partial \mu_p^d} \right] \\ &= [\mathcal{D}(F)](\boldsymbol{\mu}) \\ &= 0, \end{aligned}$$

where $\mathcal{D} : \text{Dom}(\mathcal{D}) \rightarrow V_f$ is an operator, $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_r)$, $|\boldsymbol{\alpha}| = \sum_{k=1}^r \alpha_k \leq d$, α_k are non-negative integers, $\{j_1, \dots, j_r\} \subset \{1, \dots, p\}$, and where $\text{Dom}(\mathcal{D})$ is the domain of the operator \mathcal{D} . As indicated by Evans (1998, pp. 239-240), the great advantage of interpreting PDE problem in the above form is that the general and elegant results of functional analysis can be used to study the solvability of various equations involving the operator \mathcal{D} . Frequently, $\mathcal{D}(F)$ is linear. To be continuous, the norm used in the domain of $\mathcal{D}(F)$ might be different from the one in V_f . For example, as considering the functions with compact support in an open set of \mathbf{R}^p , $\text{Dom}(\mathcal{D})$ can be the completion of the space of the functions infinitely differentiable with some inner product and V_f can be the space of the square integrable functions with another inner product. The following results analogous to Theorem 2.1 can be obtained based on the existence theorem given in Ekeland and T  mam (1999). The results can provide theoretical support for the approximation of the solution of the PDE of interest by the function falling in a finite dimensional subspace of V_f , i.e., the function having a finite basis function representation.

Theorem 2.2. *Let $f_j, j = 1, \dots, k_1$ and $f_j^*, j = 1, \dots, k_2$, corresponding to the initial conditions and boundary conditions, respectively, be the functions defined on subsets of*

\mathbf{R}^p which contain ∂M . Suppose that V_f^* , the closed subspace of V_f , is non-empty and H_0 is a finite dimensional subspace of V_f^* . Let the objective functional defined on V_f^* be

$$S(F) = \frac{1}{m_1} \sum_{i=1}^{m_1} d_i^2(F) + \lambda_1 S_1(F) + \lambda_2 S_2(F) + c \|P_{H_0^\perp}(F)\|_{V_f}^2,$$

where $P_{H_0^\perp}$ is a projection operator of V_f^* onto the orthogonal complement of H_0 ,

$$S_1(F) = \sum_{j=1}^{k_1} \left\{ \frac{1}{m_2} \sum_{i=1}^{m_2} [f_j(\tilde{\mu}_{2i}) - \langle F, \eta_{3ji} \rangle_{V_f}]^2 \right\} \\ + \sum_{j=1}^{k_2} \left\{ \frac{1}{m_3} \sum_{i=1}^{m_3} [f_j^*(\tilde{\mu}_{3i}) - \langle F, \eta_{4ji} \rangle_{V_f}]^2 \right\},$$

$$S_2(F) = \frac{1}{m_1} \sum_{i=1}^{m_1} \left[\frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij(p+1)} - \langle F, \eta_{2i} \rangle_{V_f})^2 \right],$$

and where $\lambda_1, \lambda_2 \geq 0$ and $c > 0$, $y_{ij(p+1)}$ are the observed values of some random variables with means $\langle F, \eta_{2i} \rangle_{V_f}$, $d_i(F) = a_0(\mu_{1i}) + \langle F, \eta_{1i} \rangle_{V_f}$, a_0 is a real-valued function defined on M , and where $\mu_{1i} \in M$, $\tilde{\mu}_{2i} \in \partial M$, and $\tilde{\mu}_{3i} \in \partial M$ are the means of some observable random vectors, η_{1i} , η_{2i} are representers associated with μ_{1i} , η_{3ji} are representers associated with $\tilde{\mu}_{2i}$, η_{4ji} are representers associated with $\tilde{\mu}_{3i}$, and $\eta_{1i}, \eta_{2i}, \eta_{3ji}, \eta_{4ji} \in V_f^*$. Then, the minimizer $\hat{F} = \arg \min_{F \in V_f^*} S(F)$ exists and falls in a

finite dimensional subspace of V_f^* .

Remark 2.2. Note that F in the above theorem can be the weak solutions falling in the real Hilbert space and F might take values on some normed space rather than \mathbf{R} then. The weak solutions of several well-known PDEs subject to prescribed boundary conditions fall in the Hilbert space, including the ones of second order elliptic PDEs, second order parabolic PDEs, second order hyperbolic PDEs, and Euler-Lagrange equation (see Adams and Fournier, 2003, Chapter 1 and Chapter 3; Evans, 1998, Chapter 6.2, Chapter 7.1, Chapter 7.2, and Chapter 8.2). Among these equations, the Euler-Lagrange equation is not a linear PDE. Note that the weak derivative (see Aubin, 2000, Chapter 9), i.e., the derivative in the sense of distributions, might not exist in the classical sense. However, those functions in Sobolev spaces with weak derivatives can be accurately approximated by smooth functions (see Adams and Fournier, 2003, 1.62; Meyers and Serrin, 1964). In addition, by Sobolev inequalities (see Adams and Fournier, 2003, Chapter 4; Evans, 1998, Chapter 5.6), that the equivalence class of some weak solution contains an element belonging to the space of smooth functions can be proved, i.e., the weak solution being imbedded into the space of functions having derivatives in the classical sense. Therefore, as the pointwise values of $F(\mu_{1i})$,

$[\mathcal{D}(F)](\boldsymbol{\mu}_{1i})$, $[\mathcal{D}_j(F)](\tilde{\boldsymbol{\mu}}_{2i})$, and $[\mathcal{D}_j^*(F)](\tilde{\boldsymbol{\mu}}_{3i})$ exist, it is natural to link these values with the values of $\langle F, \eta_{2i} \rangle_{V_f}$, $a_0(\boldsymbol{\mu}_{1i}) + \langle F, \eta_{1i} \rangle_{V_f}$, $\langle F, \eta_{3ji} \rangle_{V_f}$, and $\langle F, \eta_{4ji} \rangle_{V_f}$, respectively, where \mathcal{D}_j and \mathcal{D}_j^* are operators corresponding to the initial conditions and boundary conditions. To employ these pointwise values, it might be sensible to consider the Hilbert space with a reproducing kernel and assume that the operator \mathcal{D} depends on some continuous linear operator \mathcal{D}_0 , e.g. the one corresponding to a d th order linear nonhomogeneous partial differential equation, and the operators \mathcal{D}_j and \mathcal{D}_j^* are continuous linear operators. Thus, $F(\boldsymbol{\mu}_{1i}) = \langle F, \eta_{2i} \rangle_{V_f}$, $[\mathcal{D}_0(F)](\boldsymbol{\mu}_{1i}) = \langle F, \eta_{1i} \rangle_{V_f}$, $[\mathcal{D}_j(F)](\tilde{\boldsymbol{\mu}}_{2i}) = \langle F, \eta_{3ji} \rangle_{V_f}$, and $[\mathcal{D}_j^*(F)](\tilde{\boldsymbol{\mu}}_{3i}) = \langle F, \eta_{4ji} \rangle_{V_f}$ and the results given in the above theorem also hold, as indicated by the following corollary.

Corollary 2.2. Let $\mathcal{D}_j : \text{Dom}(\mathcal{D}_j) \rightarrow V_f$, $j = 1, \dots, k_1$, and $\mathcal{D}_j^* : \text{Dom}(\mathcal{D}_j^*) \rightarrow V_f$, $j = 1, \dots, k_2$, be continuous linear operators corresponding to the initial conditions and boundary conditions, where $\text{Dom}(\mathcal{D}_j)$ and $\text{Dom}(\mathcal{D}_j^*)$ are closed subspaces of V_f . Suppose that the underlying d th order partial differential equation can be expressed as

$$[\mathcal{D}(F)](\boldsymbol{\mu}) = [\mathcal{D}_0(F)](\boldsymbol{\mu}) + a_0(\boldsymbol{\mu}) = 0,$$

where $\mathcal{D}_0 : \text{Dom}(\mathcal{D}) \rightarrow V_f$ is a continuous linear operator, $\text{Dom}(\mathcal{D})$ is a closed subspace of V_f and $a_0 \in V_f$. Let $V_f^* = \text{Dom}(\mathcal{D}) \cap [\cap_{j=1}^{k_1} \text{Dom}(\mathcal{D}_j)] \cap [\cap_{j=1}^{k_2} \text{Dom}(\mathcal{D}_j^*)]$ and V_f be a Hilbert space with a reproducing kernel $K(\cdot, \cdot)$ defined on $\overline{M} \times \overline{M}$, where \overline{M} is a subset of \mathbf{R}^p , $M \subset \overline{M}$, and $\partial M \subset \overline{M}$. Let $d_i(F)$, $S_1(F)$, and $S_2(F)$ in the objective functional $S(F)$ be

$$d_i(F) = [\mathcal{D}_0(F)](\boldsymbol{\mu}_{1i}) + a_0(\boldsymbol{\mu}_{1i}),$$

$$\begin{aligned} S_1(F) = & \sum_{j=1}^{k_1} \left\{ \frac{1}{m_2} \sum_{i=1}^{m_2} \{f_j(\tilde{\boldsymbol{\mu}}_{2i}) - [\mathcal{D}_j(F)](\tilde{\boldsymbol{\mu}}_{2i})\}^2 \right\} \\ & + \sum_{j=1}^{k_2} \left\{ \frac{1}{m_3} \sum_{i=1}^{m_3} \{f_j^*(\tilde{\boldsymbol{\mu}}_{3i}) - [\mathcal{D}_j^*(F)](\tilde{\boldsymbol{\mu}}_{3i})\}^2 \right\}, \end{aligned}$$

and

$$S_2(F) = \frac{1}{m_1} \sum_{i=1}^{m_1} \left\{ \frac{1}{n_i} \sum_{j=1}^{n_i} [y_{ij(p+1)} - F(\boldsymbol{\mu}_{1i})]^2 \right\}.$$

Then, the minimizer $\hat{F} = \arg \min_{F \in V_f^*} S(F)$ exists and falls in a finite dimensional subspace of V_f^* .

The explicit expressions of the solutions of the PDEs are very few and numerical methods are commonly used (see Tveito and Winther, 1998). Therefore, consider $F(\boldsymbol{\mu}) = \sum_{l=1}^q \beta_l^* \psi_l(\boldsymbol{\mu})$ and assume that these basis functions $\psi_l \in V_f$ are smooth. Note that F is usually an approximation of the solution of the PDE of interest. Thus, some

errors might occur and an error analysis is given in Section 3.2. Since the normed values of the solutions of the partial differential equations might not be equal to one and the solutions are usually subject to some additional conditions, the constraint $\sum_{l=1}^q (\beta_l^*)^2 = 1$ can be taken away. Further, for ease of exposition, consider $p = 2$ and the second order linear nonhomogeneous partial differential equations, i.e.,

$$D \left[\boldsymbol{\mu}, \mathbf{F}^{(2)}(\boldsymbol{\mu}) \right] = a_0(\boldsymbol{\mu}) + \mathbf{a}(\boldsymbol{\mu}) \mathbf{F}^{(2)}(\boldsymbol{\mu}) = 0, \quad (3)$$

where $a_0(\boldsymbol{\mu})$ and $\mathbf{a}(\boldsymbol{\mu}) = [a_1(\boldsymbol{\mu}), \dots, a_6(\boldsymbol{\mu})]$ are constant or non-constant functions defined on M and

$$\mathbf{f}^{(2)} = \left[f(\boldsymbol{\mu}), \frac{\partial f(\boldsymbol{\mu})}{\partial \mu_1}, \frac{\partial f(\boldsymbol{\mu})}{\partial \mu_2}, \frac{\partial^2 f(\boldsymbol{\mu})}{\partial \mu_1^2}, \frac{\partial^2 f(\boldsymbol{\mu})}{\partial \mu_1 \partial \mu_2}, \frac{\partial^2 f(\boldsymbol{\mu})}{\partial \mu_2^2} \right]^t,$$

and where f is any twice differentiable function with continuous second and mixed second derivatives defined on M . For examples, the Poisson's equation based on the Laplace operator $\nabla^2 = \partial^2/\partial \mu_1^2 + \partial^2/\partial \mu_2^2$ is $\nabla^2 F(\boldsymbol{\mu}) - g(\boldsymbol{\mu}) = 0$, the heat equation is $c \partial^2 F(\boldsymbol{\mu})/\partial \mu_1^2 = \partial F(\boldsymbol{\mu})/\partial \mu_2$, and the wave equation is $c^2 \partial^2 F(\boldsymbol{\mu})/\partial \mu_1^2 = \partial^2 F(\boldsymbol{\mu})/\partial \mu_2^2$, where g is a real-valued function defined on M and c is a positive constant. The differential equations of interest depend on both the vectors $\boldsymbol{\mu}$ and $\boldsymbol{\beta}^*$ since $\partial F(\boldsymbol{\mu})/\partial \mu_{j_1} = \sum_{l=1}^q \beta_l^* \partial \psi_l(\boldsymbol{\mu})/\partial \mu_{j_1}$ and $\partial^2 F(\boldsymbol{\mu})/\partial \mu_{j_1} \partial \mu_{j_2} = \sum_{l=1}^q \beta_l^* \partial^2 \psi_l(\boldsymbol{\mu})/\partial \mu_{j_1} \partial \mu_{j_2}$, where $j_1 \in \{1, 2\}$, and $j_2 \in \{1, 2\}$. Let

$$S_0(\boldsymbol{\beta} \mid \mathbf{T}_n) = \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{n_i} \sum_{j=1}^{n_i} d^2(\boldsymbol{\beta} \mid \mathbf{T}_{ij}) \right],$$

where

$$\begin{aligned} d(\boldsymbol{\beta} \mid \mathbf{T}_{ij}) &= D \left\{ \mathbf{T}_{ij}(\mathbf{y}_i), \mathbf{F}^{(2)}[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\} \\ &= a_0[\mathbf{T}_{ij}(\mathbf{y}_i)] + \sum_{l=1}^q \beta_l \left\{ \mathbf{a}[\mathbf{T}_{ij}(\mathbf{y}_i)] \psi_l^{(2)}[\mathbf{T}_{ij}(\mathbf{y}_i)] \right\}. \end{aligned}$$

The unconstrained objective function corresponding to the partial differential equations given in expression (3) is

$$S(\boldsymbol{\beta} \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n) = S_t(\boldsymbol{\beta} \mid \tilde{\mathbf{T}}_n) + S_0(\boldsymbol{\beta} \mid \mathbf{T}_n), \quad t = 1, 2, \quad (4)$$

where $S_t(\mathbf{b} \mid \tilde{\mathbf{T}}_n)$ is the mean sum of squares corresponding to the initial conditions (IC) and boundary conditions (BC) or to $y_{ij(p+1)}$, the values of $Y_{ij(p+1)}$, and where $Y_{ij(p+1)}$ are the random variables with means $F(\boldsymbol{\mu}_{y_i})$ and $\tilde{\mathbf{T}}_n$ are the estimates corresponding to the parameters in ∂M . For a simple example, suppose that the initial conditions for the one-dimensional wave function with $M = (0, 1) \times (0, L)$ are $F(\mu_1, 0) = f_1(\mu_1, 0)$

and $\partial F(\mu_1, 0)/\partial \mu_2 = f_2(\mu_1, 0)$, where L is a positive constant and f_1 and f_2 are given functions mainly depending on μ_1 . In this case,

$$S_1(\beta | \tilde{\mathbf{T}}_n) = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{n_i} \sum_{j=1}^{n_i} \left\{ \left\{ f_1[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)] - F[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)] \right\}^2 + \left\{ f_2[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)] - \frac{\partial F[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)]}{\partial \mu_2} \right\}^2 \right\} \right\}, \quad (5)$$

where $\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)$ are sensible estimates of $\tilde{\boldsymbol{\mu}}_{y_i} = (\mu_{y_{i1}}, 0)^t$ and

$$\tilde{\mathbf{T}}_n = [\tilde{\mathbf{T}}_{11}^t(\mathbf{y}_1), \dots, \tilde{\mathbf{T}}_{1n_1}^t(\mathbf{y}_1), \dots, \tilde{\mathbf{T}}_{m1}^t(\mathbf{y}_m), \dots, \tilde{\mathbf{T}}_{mn_m}^t(\mathbf{y}_m)]^t.$$

In practice, the differences between $f_1[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)]$ and $F[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)]$ and the ones between $f_2[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)]$ and $\partial F[\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)]/\partial \mu_2$ might be significant. Therefore, $S_1(\beta | \tilde{\mathbf{T}}_n)$ rather than $\lambda S_1(\beta | \tilde{\mathbf{T}}_n)$ with large values of λ as given in the previous section is employed. However, $\lambda S_1(\beta | \tilde{\mathbf{T}}_n)$ with large values of λ is still a sensible alternative provided that $\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i)$ are accurate estimates of $\tilde{\boldsymbol{\mu}}_{y_i}$. If the boundary conditions for the one-dimensional wave function are required and are $F(0, \mu_2) = 0$ and $F(1, \mu_2) = 0$, the analogue mean sum of squares can be obtained and added to the objective function.

On the other hand, the mean sum of squares based on $y_{ij(p+1)}$ is

$$S_2(\beta | \tilde{\mathbf{T}}_n) = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{n_i} \sum_{j=1}^{n_i} \{y_{ij(p+1)} - F[\mathbf{T}_{ij}(\mathbf{y}_i)]\}^2 \right\}, \quad (6)$$

where $\tilde{\mathbf{T}}_n = \mathbf{T}_n$. The objective function $S(\beta | \mathbf{T}_n, \tilde{\mathbf{T}}_n)$ given in expression (4) can be considered as the estimator of the following objective function,

$$S^*(\beta | \boldsymbol{\mu}_y, \tilde{\boldsymbol{\mu}}_y) = S_t(\beta | \tilde{\boldsymbol{\mu}}_y) + \frac{1}{m} \sum_{i=1}^m d^2(\beta | \boldsymbol{\mu}_{y_i}), \quad (7)$$

where $\tilde{\boldsymbol{\mu}}_y = (\tilde{\boldsymbol{\mu}}_{y_1}^t, \dots, \tilde{\boldsymbol{\mu}}_{y_m}^t)^t$ are the parameters corresponding to the initial and boundary conditions as $t = 1$ or $\tilde{\boldsymbol{\mu}}_y = \boldsymbol{\mu}_y$ as $t = 2$, $S_t(\beta | \tilde{\boldsymbol{\mu}}_y)$ is the function of β by replacing $\tilde{\mathbf{T}}_{ij}$ in the function $S_t(\beta | \tilde{\mathbf{T}}_n)$ with their counterparts $\tilde{\boldsymbol{\mu}}_{y_i}$, and where $d(\beta | \boldsymbol{\mu}_{y_i}) = D[\boldsymbol{\mu}_{y_i}, \mathbf{F}^{(2)}(\boldsymbol{\mu}_{y_i})]$.

The objective function given in expression (4) incorporating with $S_2(\beta | \tilde{\mathbf{T}}_n)$ has the form

$$S(\beta | \mathbf{T}_n) = \beta^t \mathbf{A}_0(\mathbf{T}_n) \beta + [\mathbf{v}_0(\mathbf{T}_n)]^t \beta + c_0(\mathbf{T}_n), \quad (8)$$

where $c_0(\mathbf{T}_n)$ is a scalar, $\mathbf{v}_0(\mathbf{T}_n)$ is a $q \times 1$ vector, and $\mathbf{A}_0(\mathbf{T}_n)$ is a $q \times q$ matrix. On the other hand, if $S_1(\beta | \tilde{\mathbf{T}}_n)$ is a second-order polynomial function in β , the objective

function given in expression (4) incorporating with $S_1(\beta \mid \tilde{\mathbf{T}}_n)$ is analogous to the one incorporating with $S_2(\beta \mid \tilde{\mathbf{T}}_n)$, i.e.,

$$S(\beta \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n) = \beta^t \mathbf{A}_0(\mathbf{T}_n, \tilde{\mathbf{T}}_n) \beta + \left[\mathbf{v}_0(\mathbf{T}_n, \tilde{\mathbf{T}}_n) \right]^t \beta + c_0(\mathbf{T}_n, \tilde{\mathbf{T}}_n). \quad (9)$$

The expressions for c_0 , \mathbf{v}_0 , and \mathbf{A}_0 along with the ones corresponding to the objective function incorporating with $S_1(\beta \mid \tilde{\mathbf{T}}_n)$ for the wave function example are given in the supplementary materials. If the matrix \mathbf{A}_0 is positive definite, the minimizers of the objective functions $S(\beta \mid \mathbf{T}_n)$ and $S(\beta \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n)$ given in expressions (8) and (9) are $\hat{\beta} = (-1/2)\mathbf{A}_0^{-1}\mathbf{v}_0$.

The objective functions $S(\beta \mid \mathbf{T}_n)$ and $S(\beta \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n)$ given in expressions (8) and (9) can be considered as the estimators of the objective functions

$$S^*(\beta \mid \mu_y) = \beta^t \mathbf{A}_0^*(\mu_y) \beta + [\mathbf{v}_0^*(\mu_y)]^t \beta + c_0^*(\mu_y),$$

and

$$S^*(\beta \mid \mu_y, \tilde{\mu}_y) = \beta^t \mathbf{A}_0^*(\mu_y, \tilde{\mu}_y) \beta + [\mathbf{v}_0^*(\mu_y, \tilde{\mu}_y)]^t \beta + c_0^*(\mu_y, \tilde{\mu}_y),$$

respectively, where \mathbf{A}_0^* , \mathbf{v}_0^* , and c_0^* can be obtained by replacing \mathbf{T}_{ij} and $\tilde{\mathbf{T}}_{ij}$ in the functions \mathbf{A}_0 , \mathbf{v}_0 , and c_0 with their counterparts μ_{y_i} and $\tilde{\mu}_{y_i}$, respectively.

The above approaches can be extended to higher order linear partial differential equations or nonlinear partial differential equations. For nonlinear differential equations, the explicit form of $\hat{\beta}$ might not be available, but the nonlinear programming methods can be employed to find the minimizers. If the coefficient functions a_i or the functions involved in the initial and boundary conditions are unknown (see Cavalier, 2011, p. 20) and fall in some normed spaces, these functions could be expressed or approximated as a finite basis representation or considered as parameters. In the wave function example, c can be considered as a parameter and the unknown functions f_1 and f_2 could be expressed or approximated as $f_1(\tilde{\mu}) = \sum_{l=1}^{q_1} \beta_{f_1 l}^* \psi_{f_1 l}(\tilde{\mu})$ and $f_2(\tilde{\mu}) = \sum_{l=1}^{q_2} \beta_{f_2 l}^* \psi_{f_2 l}(\tilde{\mu})$, where $\psi_{f_1 l}$ and $\psi_{f_2 l}$ are some basis functions. Thus, the associated objective function is $S(\beta, \beta_{f_1}, \beta_{f_2}, c \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n)$ and the nonlinear programming methods can be employed to find the minimizers, where $\beta_{f_1} = (\beta_{f_1 1}, \dots, \beta_{f_1 q_1})^t$ and $\beta_{f_2} = (\beta_{f_2 1}, \dots, \beta_{f_2 q_2})^t$ are vectors of coefficients corresponding to $\beta_{f_1}^* = (\beta_{f_1 1}^*, \dots, \beta_{f_1 q_1}^*)^t$ and $\beta_{f_2}^* = (\beta_{f_2 1}^*, \dots, \beta_{f_2 q_2}^*)^t$, respectively.

If \mathbf{Y} is degenerated and both the values of μ_{y_i} and $Y_{ij(p+1)}$ are available, the objective function with large values of λ ,

$$S_2(\beta \mid \mu_y) + \lambda \left[\frac{1}{m} \sum_{i=1}^m d^2(\beta \mid \mu_{y_i}) + S_1(\beta \mid \tilde{\mu}_y) \right], \quad (10)$$

depending on the observations $y_{ij(p+1)}$, can be used to find the minimizers.

3. Asymptotic Aspects of Methods of Estimation

3.1 Nonlinear Programming Methods and Their Convergence

To find the minimizers of the objective functions given in expressions (1) and (4), the algorithms for the nonlinear programming problems can be employed. There are several techniques for the unconstrained optimization problems (see Bazaraa and Shetty, 1979; Nocedal and Wright, 1999; Rheinboldt, 1998). Basically, these methods can be classified according to the use of the derivatives of the objective function. The Fibonacci search procedure (see Kiefer, 1953), the method of Rosenbrock using line search (see Bazaraa and Shetty, 1979, Chapter 8.4; Rosenbrock, 1960), and the Nelder-Mead algorithm (Nelder and Mead, 1965) are derivative-free methods. Intuitively, the search direction and the distance along the direction are main quantities of several derivative-free algorithms. For example, as the objective function of several variables is $S_\lambda(\beta \mid \mathbf{T}_n)$, the method of Rosenbrock first determines a set of linearly independent orthogonal search directions $\{\mathbf{l}_j : j = 1, \dots, q\}$ based on Gram-Schmidt procedure and then finds the minimizer \hat{s}_j of the function $S_\lambda(\hat{\beta}_{k-1,j} + s\mathbf{l}_j \mid \mathbf{T}_n)$, where $\hat{\beta}_{k-1,j+1} = \hat{\beta}_{k-1,j} + \hat{s}_j\mathbf{l}_j$ is the estimated coefficient at the k th iteration with respect to the search direction \mathbf{l}_j and $\hat{\beta}_{0,1}$ is equal to the initial estimated coefficient. The new estimate at the $(k+1)$ th iteration is $\hat{\beta}_{k,1} = \hat{\beta}_{k-1,q} + \hat{s}_q\mathbf{l}_q$. The procedure stops as the distance of the current iterated point $\hat{\beta}_{k,1}$ and the previous iterated point $\hat{\beta}_{k-1,1}$ is smaller than the pre-specified termination scalar. On the other hand, the method of Newton, the methods using conjugate directions (see Bazaraa and Shetty, 1979, Chapter 8.6) which include the method of Davidon-Fletcher-Powell and its generalizations (see Broyden, 1967; Broyden, 1970; Davidon, 1959; Fletcher and Powell, 1963; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970), and the conjugate gradient method proposed by Fletcher and Reeves (1964) are the methods using derivatives. Since the direction of movement in the method of Davidon-Fletcher-Powell depends on the product of a positive definite matrix approximating the inverse of the Hessian matrix and the gradient vector, the method is also a quasi-Newton method.

In standard nonlinear programming problems, the objective functions involve some variables and known coefficients. For example, as the objective function is the Rosenbrock function, the nonlinear programming problem is

$$\min_{\beta_1, \beta_2 \in R} 100(\beta_2 - \beta_1^2)^2 + (1 - \beta_1)^2,$$

where β_1 and β_2 are variables (parameters in statistics) and the constants such as 100 and 1 are known coefficients. On the other hand, the objective functions for the mathematical equation models involve the non-random variables (parameters) β_1, \dots, β_q , random coefficients corresponding to the random vectors \mathbf{Y} , and known coefficients. Intuitively, it is similar to replace the constant coefficient 100 in the above function by a random variable Y with $E(Y) = \mu_y = 100$ and then to obtain the minimizer based

on the data generated from the random variable Y .

Convergence of the nonlinear programming algorithms is crucial. The convergence results for many iterative methods in the standard nonlinear programming problems have been well established (see Bazaraa and Shetty, 1979, Part 3). However, relatively little for the nonlinear programming problems involving the random coefficients, especially corresponding to the data in statistics, has been done. Therefore, developing convergence results by taking the randomness of some coefficients into account is required. The results concerning both the convergence of the estimators and the convergence of the optimization algorithms are established in this section. The first four theorems and associated corollary concern the convergence of the estimators and the sufficient and necessary conditions for the existence of the optimal estimators, while the last two theorems concern the convergence of the iterative sequence of estimators generated by the algorithms to the coefficients of interest. In a nutshell, the first two theorems concern the convergence of the minimizers based on the unconstrained objective functions S_λ to the vector of coefficients β^* , i.e., the consistency and asymptotical normality (weak convergence) of the estimated coefficients and associated estimator based on the objective function S_λ used in the ordinary equation models. The third and fourth theorems are the random versions of well-known Karush-Kuhn-Tucker (KKT) conditions (Karush, 1939; Kuhn and Tucker, 1951). To search for the minimizer of the unconstrained objective function of interest, such as S_λ , the nonlinear programming methods involve generating a sequence of vectors iteratively. The convergence of the algorithms (the iterative process) to the vector of coefficients β^* is crucial for the successes of the nonlinear programming methods. Therefore, the last two theorems concern the convergence of commonly used algorithms used in the mathematical equation models, such as the Newton's method. Assume that $M = \mathbf{R}^p$ in this sub-section.

The consistency of the statistics \mathbf{T}_{ij} with the mean vectors μ_{y_i} is crucial for the convergence of the nonlinear programming methods. In the following theorem, the convergence of a subsequence of minimizers $\{\hat{\beta}_\lambda(\mathbf{T}_n) : \lambda = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m\}$ to the coefficient vector β^* is established. Let

$$\mathbf{T}_{n_j} = \left[\mathbf{T}_{11}^t(\mathbf{y}_1), \dots, \mathbf{T}_{1n_{1j}}^t(\mathbf{y}_1), \dots, \mathbf{T}_{m1}^t(\mathbf{y}_m), \dots, \mathbf{T}_{mn_{mj}}^t(\mathbf{y}_m) \right]^t,$$

where $\{\mathbf{T}_{in_{ij}} : j = 1, 2, \dots\}$ is the subsequence of the sequence $\{\mathbf{T}_{in_i} : n_i = 1, 2, \dots\}$ for every i . Denote the notation \xrightarrow{p} as the convergence in probability. Note that the notations $\mu_y, \mathbf{u}, \tilde{\mu}_y, \tilde{\mathbf{u}} \in \mathbf{R}^{mp}$ and $\mu, \tilde{\mu}, \mu_{y_i} \in \mathbf{R}^p$ are used in the following.

The existence of the minimizer $\hat{\beta}_\lambda(\mathbf{T}_n)$ in Theorem 3.1 can be guaranteed by the following lemma.

Lemma 3.1. $\hat{\beta}_\lambda^*(\mathbf{u})$, the minimizer of $S_\lambda^*(\beta \mid \mathbf{u})$, exists for any $\mathbf{u} \in \mathbf{R}^{mp}$ and any $\lambda > 0$.

Note that there might exist multiple minimizers of $S_\lambda(\beta \mid \mathbf{T}_n)$. In such case, $\hat{\beta}_\lambda(\mathbf{T}_n)$ is equal to one of these minimizers.

Theorem 3.1. *There exists a subsequence $\{\hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) : j = 1, 2, \dots\}$ of the sequence of random vectors*

$$\left\{ \hat{\beta}_{\lambda}(\mathbf{T}_n) : \lambda = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m \right\},$$

converging to β^ in probability as $j \rightarrow \infty$, i.e.,*

$$\hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) \xrightarrow[p]{j \rightarrow \infty} \beta^*,$$

if the following conditions hold:

(i) *There exists a neighborhood of μ_y such that $S_{\lambda}^*(\beta | \mathbf{u})$ has a unique minimizer $\hat{\beta}_{\lambda}^*(\mathbf{u})$ for \mathbf{u} in the neighborhood and for large λ , i.e., $S_{\lambda}^*(\beta | \mathbf{u})$ having a unique minimizer provided that λ is greater than some positive integer.*

(ii)

$$\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i) \xrightarrow[p]{n_i \rightarrow \infty} \mu_{y_i}$$

for every i .

Note that the measurability of $\hat{\beta}_{\lambda}(\mathbf{T}_n)$ is assumed implicitly in the above theorem. If $\hat{\beta}_{\lambda}(\mathbf{T}_n)$ is not measurable, the outer probability measures could be employed to prove the convergence in probability in such case. As the V_f -valued estimator (see Da Prato and Zabczyk, 1992, Chapter 1) $\hat{F}_{\lambda_j} = \sum_{l=1}^q \hat{\beta}_{\lambda_j l}(\mathbf{T}_{n_j}) \psi_l(\mu)$ is of interest, the following result concerning the consistency of the estimator can be obtained by employing the above theorem, where $\hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) = [\hat{\beta}_{\lambda_j 1}(\mathbf{T}_{n_j}), \dots, \hat{\beta}_{\lambda_j q}(\mathbf{T}_{n_j})]^t$.

Corollary 3.1. *Let $\hat{F}_{\lambda_j} = \sum_{l=1}^q \hat{\beta}_{\lambda_j l}(\mathbf{T}_{n_j}) \psi_l(\mu)$ and $\psi_l \in V_f$. Then,*

$$\|\hat{F}_{\lambda_j} - F\|_{V_f} \xrightarrow[p]{j \rightarrow \infty} 0$$

if the conditions given in Theorem 3.1 hold.

As $\mathbf{T}_i(\mathbf{y}_i) = \sum_{j=1}^{n_i} \mathbf{y}_{ij}/n_i$, $y_{ijk}, j = 1, \dots, k = 1, \dots, p$, and $E|y_{i1k}| < \infty$, the condition for the consistency of the statistics $\mathbf{T}_i(\mathbf{y}_i)$ in Theorem 3.1 and Corollary 3.1 holds by the weak law of large numbers.

The following theorem concerns asymptotic normality of the V_f -valued estimator based on the estimated coefficients. For simplification, let $\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i)$, $n_i = N$, $i = 1, \dots, m$, i.e., equal number of repeated observations at each site, and $\mathbf{T}_N^* = [\mathbf{T}_1^t(\mathbf{y}_1), \dots, \mathbf{T}_m^t(\mathbf{y}_m)]^t$ be an $mp \times 1$ vector. The objective function $S_{\lambda}(\beta | \mathbf{T}_n)$ given in expression (1) depending on \mathbf{T}_N^* , i.e.,

$$S_{\lambda}(\beta | \mathbf{T}_n) = \beta^t [\Psi^*(\mathbf{T}_N^*)]^t \Psi^*(\mathbf{T}_N^*) \beta + \lambda (\beta^t \beta - 1)^2,$$

is also denoted as $S_{\lambda}(\beta | \mathbf{T}_N^*)$ and hence the corresponding minimizer, a q -dimensional random vector, is denoted as

$$\hat{\beta}_{\lambda}(\mathbf{T}_N^*) = \left[\hat{\beta}_{\lambda 1}(\mathbf{T}_N^*), \dots, \hat{\beta}_{\lambda q}(\mathbf{T}_N^*) \right]^t.$$

Furthermore, $S_\lambda(\beta | \mathbf{u}) = S_\lambda^*(\beta | \mathbf{u})$ at any $\mathbf{u} \in \mathbf{R}^{mp}$ in such case and hence $\hat{\beta}_\lambda(\mathbf{u}) = \hat{\beta}_\lambda^*(\mathbf{u})$, where $S_\lambda(\beta | \mathbf{u})$ and $\hat{\beta}_\lambda(\mathbf{u})$ are the function and minimizer by replacing \mathbf{T}_N^* with \mathbf{u} in $S_\lambda(\beta | \mathbf{T}_N^*)$ and $\hat{\beta}_\lambda(\mathbf{T}_N^*)$, respectively. The following lemma indicates that the minimizer $\hat{\beta}_\lambda(\mathbf{u})$ exists.

Lemma 3.2. $\hat{\beta}_\lambda(\mathbf{u})$, the minimizer of $S_\lambda(\beta | \mathbf{u})$, exists for any $\mathbf{u} \in \mathbf{R}^{mp}$ and any $\lambda > 0$ if the following condition holds:

(i) $\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i)$, and $n_i = N$, $i = 1, \dots, m$.

The above lemma implies that there exists a function defined on \mathbf{R}^{mp} taking value $\hat{\beta}_\lambda(\mathbf{u})$ at \mathbf{u} . If the minimizers $\hat{\beta}_\lambda(\mathbf{u})$ at \mathbf{u} are not unique, the function takes one of the values. For succinctness, the function is also denoted as $\hat{\beta}_\lambda(\mathbf{u})$. The asymptotic normality of the V_f -valued estimator $\hat{F}_\lambda = \sum_{l=1}^q \hat{\beta}_{\lambda l}(\mathbf{T}_N^*) \psi_l(\boldsymbol{\mu})$ can be established based on the above lemma, the mapping theorem (see Billingsley, 1999, Theorem 2.7) and the result for functions of asymptotically normal vectors (see Serfling, 1980, Chapter 3.3). For a V_f -valued Radon Gaussian variable g (see Ledoux and Talagrand, 1991, Chapter 3), let

$$\Sigma(g) = \sup_{\|L\|_{V_f^*} \leq 1, L \in V_f^*} \{E\{[L(g)]^2\}\}^{1/2},$$

where V_f^* is the topological dual space of V_f .

Theorem 3.2. Let $\hat{F}_\lambda = \sum_{l=1}^q \hat{\beta}_{\lambda l}(\mathbf{T}_N^*) \psi_l(\boldsymbol{\mu})$, where $\psi_l \in V_f$. $\{\mathcal{F}_N = \sqrt{N}(\hat{F}_\lambda - F) : N = 1, 2, \dots\}$ converges in distribution to a V_f -valued centered Radon Gaussian random variable \mathcal{F} with

$$\Sigma(\mathcal{F}) = \sup_{\|L\|_{V_f^*} \leq 1, L \in V_f^*} \{[\mathbf{v}(L)]^t \mathbf{D} \Sigma \mathbf{D}^t \mathbf{v}(L)\}^{1/2}$$

and $\mathbf{v}(L) = [L(\psi_1), \dots, L(\psi_q)]^t$, i.e., $\text{Var}[L(\mathcal{F})] = [\mathbf{v}(L)]^t \mathbf{D} \Sigma \mathbf{D}^t \mathbf{v}(L)$, if the following conditions hold:

- (i) The condition given in Lemma 3.2 holds.
- (ii) Every element of $\hat{\beta}_\lambda(\mathbf{u})$ has a nonzero differential at $\boldsymbol{\mu}_y$ and the associated matrix \mathbf{D} is a $q \times mp$ matrix with the (l, j) th element equal to the first derivative of the function $\hat{\beta}_{\lambda l}(\mathbf{u})$ with respect to the j th element of \mathbf{u} at $\boldsymbol{\mu}_y$.
- (iii) $\{\sqrt{N}(\mathbf{T}_N^* - \boldsymbol{\mu}_y) : N = 1, 2, \dots\}$ converges in distribution to the multivariate normal random variable with zero mean vector and variance-covariance matrix Σ .
- (iv) V_f is separable.

As $\mathbf{T}_i(\mathbf{y}_i) = \sum_{j=1}^N \mathbf{y}_{ij}/N$ and $(\mathbf{y}_{1j}^t, \dots, \mathbf{y}_{mj}^t)^t, j = 1, \dots$, are i.i.d. with a variance-covariance matrix, condition (iii) for the asymptotical normality of the statistics \mathbf{T}_N^* in the above theorem holds by the central limit theorem in \mathbf{R}^{mp} . Note that the only required condition imposed on ψ_l in Theorem 3.2 is $\psi_l \in V_f$, i.e., the continuity assumption imposed on ψ_l being not necessary.

The following two theorems provide the sufficient and necessary optimality conditions for the ordinary equation models. These conditions can be considered as the

random versions of the well-known Karush-Kuhn-Tucker (KKT) conditions (Karush, 1939; Kuhn and Tucker, 1951).

Theorem 3.3. *There exists a subsequence $\{\hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) : j = 1, 2, \dots\}$ of the sequence of random vectors*

$$\left\{ \hat{\beta}_{\lambda}(\mathbf{T}_n) : \lambda = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m \right\}$$

such that

$$\hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) \xrightarrow[p \rightarrow \infty]{} \beta^*,$$

and

$$[\Psi(\mathbf{T}_{n_j})]^t \Psi(\mathbf{T}_{n_j}) \hat{\beta}_{\lambda_j}(\mathbf{T}_{n_j}) \xrightarrow[p \rightarrow \infty]{} \mathbf{0},$$

if the conditions given in Theorem 3.1 hold.

The sufficient conditions for the ordinary equation models can be established, as indicated by the following theorem.

Theorem 3.4. $\tilde{\beta}(\mu_y) = \beta^*$ if the following conditions hold:

(i) *There exists a sequence of estimators $\{\tilde{\beta}(\mathbf{T}_n) : n_i = 1, 2, \dots, i = 1, \dots, m\}$ such that*

$$[\Psi(\mathbf{T}_n)]^t \Psi(\mathbf{T}_n) \tilde{\beta}(\mathbf{T}_n) \xrightarrow[n_1, \dots, n_m \rightarrow \infty]{p} \mathbf{0},$$

and

$$\tilde{\beta}(\mathbf{T}_n) \xrightarrow[n_1, \dots, n_m \rightarrow \infty]{p} \tilde{\beta}(\mu_y),$$

where $\tilde{\beta}(\mu_y)$ is the unit vector.

(ii)

$$\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i) \xrightarrow[n_i \rightarrow \infty]{p} \mu_{y_i}$$

for every i .

As $F \in V_f^0$, the following theorem indicates that the subsequence of the sequence of minimizers of the objective function given in expression (1) generated by the Newton's method converges to the true coefficient vector β^* . Let $\{\hat{\beta}_{\lambda,k}(\mathbf{T}_n) : k = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m\}$ and $\{\hat{\beta}_{\lambda,k}^*(\mu_y) : k = 1, 2, \dots\}$ be the sequences of minimizers of the objective functions given in expressions (1) and (2) generated by the Newton's method at the k th iteration, respectively.

Theorem 3.5. *Assume that $[\mathbf{H}_{\beta}(\beta, \mathbf{u})]^{-1}$ exists for any $\beta \in \mathbf{R}^q$ and any $\mathbf{u} \in \mathbf{R}^{mp}$, where $\mathbf{H}_{\beta}(\beta, \mathbf{u}) = \partial^2 S_{\lambda}^*(\beta | \mathbf{u}) / \partial \beta \partial \beta^t$ is the Hessian matrix. There exists a subsequence $\{\hat{\beta}_{\lambda,k_j}(\mathbf{T}_{n_j}) : j = 1, 2, \dots\}$ of the sequence of random vectors $\{\hat{\beta}_{\lambda,k}(\mathbf{T}_n) : k =$*

$1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m\}$ converging to either β^* or $\hat{\beta}_\lambda^*(\mu_y)$ in probability as $j \rightarrow \infty$ if the following conditions hold:

(i) For each β contained in a compact subset C_1 of the set $\{\beta : \|\beta - \beta^*\|_{R^q} \leq \|\hat{\beta}_{\lambda,1}^*(\mu_y) - \beta^*\|_{R^q}, (\|\beta\|_{R^q}^2 - 1)(\|\beta^* - \beta\|_{R^q} - 1) \geq 0\}$, the corresponding Hessian matrix of the objective function $S_\lambda^*(\beta | \mu_y)$ given in expression (2), i.e., $\partial^2 S_\lambda^*(\beta | \mu_y) / \partial \beta \partial \beta^t$, is positive definite and

$$\sigma_q \left\{ (4\lambda)^{-1} [\Psi^*(\mu_y)]^t \Psi^*(\mu_y) + \beta \beta^t \right\} > \|\beta\|_{R^q}^2,$$

where $\sigma_q(\mathbf{A})$ is the smallest singular value of a $q \times q$ matrix \mathbf{A} .

(ii) $\hat{\beta}_{\lambda,k}^*(\mu_y) \in C_1$ for every k .

(iii)

$$\hat{\beta}_{\lambda,1}(T_n) \xrightarrow[n_1, \dots, n_m \rightarrow \infty]{p} \hat{\beta}_{\lambda,1}^*(\mu_y).$$

(iv)

$$T_{ij}(\mathbf{y}_i) = T_i(\mathbf{y}_i) \xrightarrow[n_i \rightarrow \infty]{p} \mu_{y_i}$$

for every i

If the linear nonhomogeneous differential equations are of interest and the initial and boundary conditions are linear in β , e.g. the illustrative example given in Section 2.2, the following theorem indicates that the subsequence of minimizers of the objective function given in expression (4) generated by the Newton's method converges to $\hat{\beta}^*(\mu_y, \tilde{\mu}_y)$, which is the minimizer of the objective function given in expression (7). Let $\{\hat{\beta}_k(T_n, \tilde{T}_n) : k = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m\}$ and $\{\hat{\beta}_k^*(\mu_y, \tilde{\mu}_y) : k = 1, 2, \dots\}$ be the sequences of minimizers of the objective functions given in expressions (4) and (7) generated by the Newton's method, respectively. Note that $\hat{\beta}_k(T_n, \tilde{T}_n)$ and $\hat{\beta}_k^*(\mu_y, \tilde{\mu}_y)$ independent of λ are the k th iterated vectors generated by the Newton's method. Let $\tilde{T}_{n_j} = [\tilde{T}_{11}^t(\mathbf{y}_1), \dots, \tilde{T}_{1n_{1j}}^t(\mathbf{y}_1), \dots, \tilde{T}_{m1}^t(\mathbf{y}_m), \dots, \tilde{T}_{mn_{mj}}^t(\mathbf{y}_m)]^t$, where $\{\tilde{T}_{in_{ij}} : j = 1, 2, \dots\}$ is the subsequence of the sequence $\{\tilde{T}_{in_i} : n_i = 1, 2, \dots\}$ for every i .

Theorem 3.6. *There exists a subsequence $\{\hat{\beta}_{k_j}(T_{n_j}, \tilde{T}_{n_j}) : j = 1, 2, \dots\}$ of the sequence of random vectors $\{\hat{\beta}_k(T_n, \tilde{T}_n) : k = 1, 2, \dots, n_i = 1, 2, \dots, i = 1, \dots, m\}$ converging to $\hat{\beta}^*(\mu_y, \tilde{\mu}_y)$ in probability as $j \rightarrow \infty$ if the following conditions hold:*

(i) $S^*(\beta | \mu_y, \tilde{\mu}_y) = \beta^t \mathbf{A}_0^*(\mu_y, \tilde{\mu}_y) \beta + [v_0^*(\mu_y, \tilde{\mu}_y)]^t \beta + c_0^*(\mu_y, \tilde{\mu}_y)$, where $\mathbf{A}_0^*(\mathbf{u}, \tilde{\mathbf{u}})$ is nonsingular for any $\mathbf{u}, \tilde{\mathbf{u}} \in \mathbf{R}^{mp}$, is positive definite at $(\mu_y, \tilde{\mu}_y)$ and its elements are continuous at $(\mu_y, \tilde{\mu}_y)$, the elements of v_0^* are continuous at $(\mu_y, \tilde{\mu}_y)$, and c_0^* is a real-valued function.

(ii)

$$\hat{\beta}_1(T_n, \tilde{T}_n) \xrightarrow[n_1, \dots, n_m \rightarrow \infty]{p} \hat{\beta}_1^*(\mu_y, \tilde{\mu}_y).$$

(iii) For every i ,

$$\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i) \xrightarrow[n_i \rightarrow \infty]{p} \boldsymbol{\mu}_{y_i},$$

and

$$\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i) = \tilde{\mathbf{T}}_i(\mathbf{y}_i) \xrightarrow[n_i \rightarrow \infty]{p} \tilde{\boldsymbol{\mu}}_{y_i}.$$

3.2 Error Analysis and Choices of Numbers of Basis Functions

If the orthogonal polynomials such as Chebyshev polynomials or Hermite polynomials are used for the ordinary equation models, it is suggested that a large number of basis functions should be avoided owing to the computational inefficiency and great complexity of the estimated equations. Theoretically, the choice of the basis function for the ordinary equation models mainly depends on the space which F might fall in (see Kreyszig, 1978, Chapter 3.7). For example, if $\mu_j \in (-\infty, \infty)$, Hermite polynomials might be a good choice. On the other hand, the choices of the basis functions and the number of basis functions for the differential equation models mainly depend on the PDEs of interest. The solutions of some PDEs such as the heat and wave equations subject to the specific initial conditions and boundary conditions can be expressed as infinite Fourier series and trigonometric functions can be used as the basis functions. In such cases, an error analysis might be required in order to develop some criteria to choose a sensible value of q , the number of the basis functions. Suppose that V_f is a separable Hilbert space. Then, the true function can be expressed as $F = \sum_{l=1}^{\infty} \beta_l^* \psi_l$, where ψ_l are orthonormal basis functions. Let $\hat{F} = \sum_{l=1}^q \hat{\beta}_l \psi_l$ be the estimator based on the finite basis representation. If $\sigma_{\hat{\beta}_l}^2 = \text{Var}(\hat{\beta}_l)$ exist, the mean integrated squared risk is

$$\begin{aligned} & E \left(\left\| \hat{F} - F \right\|_{V_f}^2 \right) \\ &= E \left(\left\| \hat{F} - E_{\hat{F}} \right\|_{V_f}^2 \right) + \left\| E_{\hat{F}} - F \right\|_{V_f}^2 \\ &= \sum_{l=1}^q \sigma_{\hat{\beta}_l}^2 + \sum_{i=l}^q \left[E(\hat{\beta}_l) - \beta_l^* \right]^2 + \sum_{l=q+1}^{\infty} (\beta_l^*)^2, \end{aligned}$$

where $E_{\hat{F}} = \sum_{l=1}^q E(\hat{\beta}_l) \psi_l$. The risk depends on the variance (the propagated noise error) reflected by the first term and the bias (the approximation error) reflected by the last two terms. The variance term can measure the stability of the estimator, while the bias term can measure how well the finite basis representation approximates the true function in the average sense (see Cavalier, 2011, p. 34, p. 37). As q tends to infinity, the last term $\sum_{l=q+1}^{\infty} (\beta_l^*)^2$ tends to zero. However, as q increases, the number

of summands in the first term increases and the value of this term might possibly increase. Moreover, the computational burden is heavy for obtaining the estimator based on the large number of basis functions. Intuitively, the sensible choice of q is a trade-off between the first term and the sum of the last two terms, $\|E_{\hat{F}} - F\|_{V_f}^2$, i.e., the trade-off between the variance and bias. For illustrations, consider the wave equation $\partial^2 F(\mu_1, \mu_2)/\partial \mu_2^2 = c^2 \partial^2 F(\mu_1, \mu_2)/\partial \mu_1^2$, $0 < \mu_1 < 1$, $0 < \mu_2 < 2T$, with the boundary conditions $F(0, \mu_2) = F(1, \mu_2) = 0$ and the initial conditions $F(\mu_1, 0) = f_1(\mu_1, 0)$, $\partial F(\mu_1, 0)/\partial \mu_2 = f_2(\mu_1, 0)$, where T is a positive integer and f_1, f_2 are real functions. The solution is

$$F(\mu_1, \mu_2) = \sum_{n=1}^{\infty} 2 \sin(n\pi\mu_1) \left\{ \left[\int_0^1 f_1(\mu_1, 0) \sin(n\pi\mu_1) d\mu_1 \right] \cos(cn\pi\mu_2) + \frac{1}{cn\pi} \left[\int_0^1 f_2(\mu_1, 0) \sin(n\pi\mu_1) d\mu_1 \right] \sin(cn\pi\mu_2) \right\}.$$

As $c = 1$, $f_1(\mu_1, 0) = \sum_{n=1}^{\infty} a_n \sin(n\pi\mu_1)$, and $f_2(\mu_1, 0) = \sum_{n=1}^{\infty} n\pi b_n \sin(n\pi\mu_1)$, the solution can be expressed as

$$F(\mu_1, \mu_2) = \sum_{n=1}^{\infty} a_n \sin(n\pi\mu_1) \cos(n\pi\mu_2) + \sum_{n=1}^{\infty} b_n \sin(n\pi\mu_1) \sin(n\pi\mu_2)$$

(see Tveito and Winther, 1998, Chapter 5), where $a_n, b_n \in R$ are some constants such that $\sum_{n=1}^{\infty} a_n^2 < \infty$ and $\sum_{n=1}^{\infty} n^2 b_n^2 < \infty$. Suppose that V_f is the Hilbert space with the inner product given by

$$\langle F_1, F_2 \rangle_{V_f} = \int_0^{2T} \int_0^1 F_1(\mu_1, \mu_2) F_2(\mu_1, \mu_2) d\mu_1 d\mu_2$$

and the orthonormal basis functions

$$\begin{aligned} & \{(T/2)^{-1/2} \sin(n\pi\mu_1) \cos(n\pi\mu_2) : n = 1, \dots\} \\ & \cup \{(T/2)^{-1/2} \sin(n\pi\mu_1) \sin(n\pi\mu_2) : n = 1, \dots\} \\ & \cup \{(2T)^{-1/2}\}. \end{aligned}$$

Then, $\beta_{2n-1}^* = a_n$, $\beta_{2n}^* = b_n$, and the last bias term $\sum_{l=q+1}^{\infty} (\beta_l^*)^2 \xrightarrow{q \rightarrow \infty} 0$ since $\sum_{n=1}^{\infty} (a_n^2 + b_n^2) < \infty$.

It might be reasonable to assume that an “optimal” value of q should result in the smallest value of the mean integrated square risk. Thus, the unbiased or nearly unbiased estimators of the mean integrated squared risk or its corresponding empirical risks can be used as selection criteria. For example, analogous to the one employed in nonparametric regression setting, one estimator of the empirical risk

$$\frac{\|\hat{\mathbf{f}}(q) - \mathbf{f}\|_{R^n}^2}{n} = \frac{[\hat{\mathbf{f}}(q) - \mathbf{f}]^t [\hat{\mathbf{f}}(q) - \mathbf{f}]}{n},$$

is

$$\frac{S_2(\hat{\beta} \mid \mathbf{T}_n)}{\phi(q)} = \frac{\left\| \mathbf{y}_{(p+1)} - \hat{\mathbf{f}}(q) \right\|_{R^n}^2}{\phi(q)} = \frac{[\mathbf{y}_{(p+1)} - \hat{\mathbf{f}}(q)]^t [\mathbf{y}_{(p+1)} - \hat{\mathbf{f}}(q)]}{\phi(q)},$$

where $\phi(q)$ is the penalty function of q ,

$$\mathbf{f} = [F(\boldsymbol{\mu}_{y_1}), \dots, F(\boldsymbol{\mu}_{y_1}), \dots, F(\boldsymbol{\mu}_{y_m}), \dots, F(\boldsymbol{\mu}_{y_m})]^t,$$

$$\hat{\mathbf{f}}(q) = \left(\hat{F}_{11}, \dots, \hat{F}_{1n_1}, \dots, \hat{F}_{m1}, \dots, \hat{F}_{mn_m} \right)^t,$$

$\mathbf{y}_{(p+1)}$ is an $n \times 1$ vector with the l th elements equal to $m^{-1/2} n_i^{-1/2} y_{ij(p+1)}$ and where $\hat{F}_{ij} = m^{-1/2} n_i^{-1/2} \sum_{l=1}^q \hat{\beta}_l \psi_l [\mathbf{T}_{ij}(\mathbf{y}_i)]$ and $l = \sum_{k=1}^{i-1} n_k + j$. The choice of the penalty function ϕ might be different from the one in nonparametric regression setting since $\hat{\mathbf{f}}(q)$ might not be linear in term of $\mathbf{y}_{(p+1)}$. If $\hat{\mathbf{f}}(q)$ can be approximated by $\mathbf{H}(q) \mathbf{y}_{(p+1)}$, several choices of ϕ based on $\text{Tr}[\mathbf{H}(q)]/n$ could be employed (see Eubank, 1988, pp. 38-40), where $\mathbf{H}(q)$ is an $n \times n$ matrix function of q and $\text{Tr}[\mathbf{H}(q)]$ is the trace of the matrix $\mathbf{H}(q)$. This approach is analogous to the one used for statistical model selection problems. Note that the estimators of the empirical risks involve the information provided by the derivatives could be also used for the choices of the numbers of basis functions, for example, $S(\hat{\beta} \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n)/\phi(q)$.

4. Numerical Illustrations

4.1 Simulations

The Nelder-Mead algorithm, the Newton's method, the quasi-Newton methods, and the conjugate gradient methods are commonly used nonlinear programming methods. Therefore, in the first simulation study, several methods, including the Nelder-Mead algorithm, the Newton's method, the BFGS method (Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970), and the conjugate gradient (CG) method proposed by Fletcher and Reeves (1964), were employed for illustrations. Note that the BFGS method is a quasi-Newton method (see Nocedal and Wright, 1999, Chapter 8). A large value of the penalty parameter was pre-specified. Further, $\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i) = \sum_{j=1}^{n_i} \mathbf{y}_{ij}/n_i$.

In the first simulation study, the conic equation in Section 1 was used. The observed data were generated by

$$y_{ij1} = \mu_{y_{i1}} + \epsilon_{ij1}, \quad y_{ij2} = \mu_{y_{i2}} + \epsilon_{ij2}, \quad i = 1, \dots, I, \quad j = 1, \dots, n_i = J,$$

where

$$\mu_{y_{i1}} = \frac{\sqrt{6}}{2} \sin(s_i) - \frac{9\sqrt{2}}{2} \cos(s_i), \quad \mu_{y_{i2}} = \frac{\sqrt{6}}{2} \sin(s_i) + \frac{9\sqrt{2}}{2} \cos(s_i),$$

Table 1. Relative errors for the conic equation model.

$J = 1$	$\sigma = 0.2$	$\sigma = 1$	$\sigma = 2$
Nelder-Mead	0.02674	0.08226	0.25016
CG	0.01259	0.02906	0.08748
BFGS	0.05865	0.33665	0.38301
Newton	0.05866	0.33666	0.38301
$J = 10$	$\sigma = 0.2$	$\sigma = 1$	$\sigma = 2$
Nelder-Mead	0.02577	0.02972	0.04977
CG	0.01214	0.01381	0.01854
BFGS	0.01401	0.11757	0.25144
Newton	0.01394	0.11762	0.25150
$J = 30$	$\sigma = 0.2$	$\sigma = 1$	$\sigma = 2$
Nelder-Mead	0.02391	0.02768	0.03485
CG	0.01225	0.01291	0.01465
BFGS	0.01064	0.05959	0.13804
Newton	0.01056	0.05956	0.13819

and where ϵ_{ij1} and ϵ_{ij2} were independent normal random errors with zero means and standard deviations σ equal to 0.2, 1, and 2, and s_i were the spacing with three settings: $s_i = 0.01i$ corresponding to $I = 629, J = 1$, $s_i = 0.1i$ corresponding to $I = 63, J = 10$, and $s_i = 0.2i$ corresponding to $I = 32, J = 30$. These settings corresponded to different numbers of the repeated observations at each site. The data discussed in Section 1 were generated by the first setting. 1000 replicates of random errors were generated. The Chebyshev polynomials were used for generating the basis functions. The basis functions were

$$\psi_{3(j-1)+i}(\boldsymbol{\mu}) = \phi_i(\mu_1)\phi_j(\mu_2), \quad i = 1, 2, 3, \quad j = 1, 2, 3,$$

where the functions ϕ_i were $\phi_1(x) = 1$, $\phi_2(x) = x$, and $\phi_3(x) = 2x^2 - 1$. Let V_f be the space consisting of the functions defined on M , some subset of \mathbf{R}^2 . The uniqueness of the ordinary equation model relies on both the domain M and the space V_f . If V_f is the space generated by the basis functions $\psi_{3(j-1)+i}(\boldsymbol{\mu})$, M_0 is the set of the vectors (μ_1, μ_2) satisfying the conic equation $28\mu_1^2 + 28\mu_2^2 + 52\mu_1\mu_2 - 162 = 0$, and $M = M_0$, the ordinary equation model is unique. With the same basis functions and the same set M_0 , the ordinary equation model is not unique provided that $M = \mathbf{R}^2$. Nevertheless, the ordinary equation model is unique provided that $M = \mathbf{R}^2$ and V_f is the space spanned by the basis functions $\{\psi_l(\boldsymbol{\mu}) : l = 1, 2, 3, 4, 5, 7\}$.

The underlying equation is equivalent to the normalized equation

$$\sum_{l=1}^9 \beta_l^* \psi_l(\boldsymbol{\mu}) = \frac{1}{\sqrt{21052}} [-134\psi_1(\boldsymbol{\mu}) + 14\psi_3(\boldsymbol{\mu}) + 52\psi_5(\boldsymbol{\mu}) + 14\psi_7(\boldsymbol{\mu})] = 0,$$

where $\sum_{l=1}^9 (\beta_l^*)^2 = 1$. The relative error $\|\hat{\boldsymbol{\beta}}_\lambda(\mathbf{T}_n) - \boldsymbol{\beta}^*\|_{R^9} / \|\boldsymbol{\beta}^*\|_{R^9}$ (see Nocedal and Wright, 1999, p. 606), the ratio of the error to the size of the true coefficient vector, was employed in the simulation study.

The means of the relative errors for different nonlinear programming methods under different settings are provided in Table 1. The medians of the relative errors for these methods and settings are similar to the ones given in Table 1. These methods are quite stable if the variations of the data are small and can result in accurate fits. If the variations of the data are large and $J = 1$, these methods might fail to approximate the true coefficients well and might result in poor fits, as displayed in Figure 4. In such case, sensible approximations of the optimal solutions by the minimizers could be achieved by increasing the number of repeated observations at each site. $\mathbf{T}_i(\mathbf{y}_i) = \sum_{j=1}^{n_i} \mathbf{y}_{ij} / n_i$ might be accurate estimates of $\boldsymbol{\mu}_{y_i}$ for large n_i and hence the estimators obtained by employing large number of repeated observations might approximate the true coefficients well by Theorem 3.1. In fact, the improvements on the fits are significant by using the number of repeated observations equal to 30. These methods can approximate the true coefficients well and can result in accurate fits as $J = 30$. As presented in Table 1, for large data variations, the larger the number of repeated observations is at each site, the smaller the relative errors are. Since the fits based on the BFGS and Newton's methods, respectively, are quite close, only the fits based on the BFGS method are given in these plots.

Since the unconstrained objective function is differentiable in the study, it is not surprising that the conjugate gradient method using the gradient information performs well. Note that the sequence of solutions generated by the methods using conjugate directions can converge to the true solutions in finite steps for quadratic objective functions (see Bazaraa and Shetty, 1979, Theorem 8.6.3). Since the constrained objective function is quadratic and the penalty function is the square of a quadratic function, the quadratic approximation for the unconstrained objective function might be quite well. This provides an explanation for the good performance of the conjugate gradient method. As shown in Figure 4, the conjugate gradient method provides a slightly better fit than the ones by other methods. On the other hand, since the Newton's method and the quasi-Newton method involve the inverses of the Hessian matrices and associated approximated matrices, large data variations might result in the instability of these methods due to possibly significant changes on these inverse matrices. As displayed in Figure 4, the BFGS method might not result in a sensible fit.

In the second simulation study, the one-dimensional wave function in Section 1 was

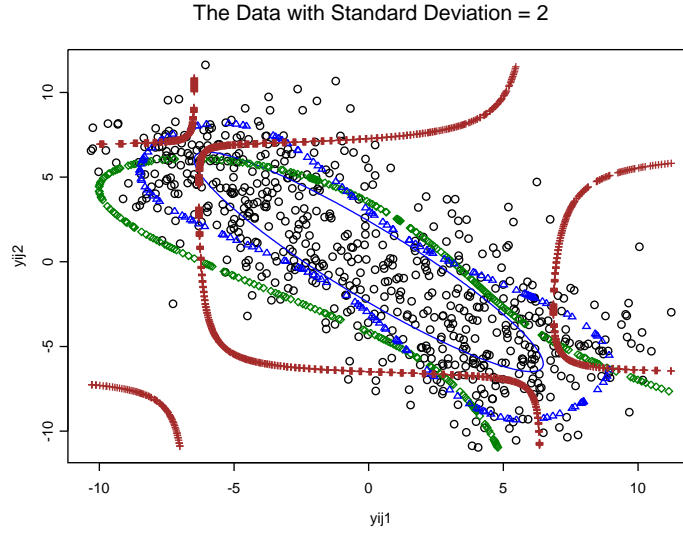


Fig. 4. Different nonlinear programming methods for the data with standard deviation equal to 2: True equation (blue line); Nelder-Mead (green \diamond); CG (blue \triangle); BFGS (brown $+$).

employed and the observed data were generated by

$$y_{ij1} = \mu_{y_{i1}} + \epsilon_{ij1}, \quad y_{ij2} = \mu_{y_{i2}} + \epsilon_{ij2}, \quad y_{ij(3)} = F(\mu_{y_{i1}}, \mu_{y_{i2}}) + \epsilon_{ij3},$$

$i = 1, \dots, 63$, $j = 1, \dots, n_i = J$, where $\mu_{y_{i1}}$ and $\mu_{y_{i2}}$ were sequences starting from 0 to 6.2 with spacing equal to 0.1, $J = 1$ or $J = 30$ was the number of repeated observations at each site, and ϵ_{ij1} , ϵ_{ij2} , and ϵ_{ij3} were independent normal random errors with zero means and standard deviations σ equal to 0.2, 1, and 2. Let the initial conditions be $F(\mu_{y_{i1}}, 0) = 7.5 \cos(\mu_{y_{i1}})$ and $\partial F(\mu_{y_{i1}}, 0) / \partial \mu_{y_{i1}} = -15 \sin(\mu_{y_{i1}})$. 1000 replicates of random errors were generated. The trigonometric functions were used for generating the basis functions. The basis functions were

$$\psi_{5(j-1)+i}(\boldsymbol{\mu}) = \phi_i(\mu_1) \phi_j(\mu_2), \quad i = 1, \dots, 5, \quad j = 1, \dots, 5,$$

where the functions ϕ_i were $\phi_1(x) = 1$, $\phi_2(x) = \cos(x)$, $\phi_3(x) = \cos(2x)$, $\phi_4(x) = \sin(x)$, and $\phi_5(x) = \sin(2x)$. The solution can be expressed as

$$F(\boldsymbol{\mu}) = \sum_{l=1}^{25} \beta_l^* \psi_l(\boldsymbol{\mu}) = 7.5 \psi_{12}(\boldsymbol{\mu}) + 7.5 \psi_{24}(\boldsymbol{\mu}).$$

$\mathbf{T}_{ij}(\mathbf{y}_i) = \mathbf{T}_i(\mathbf{y}_i) = \sum_{j=1}^{n_i} \mathbf{y}_{ij} / n_i$ and $\tilde{\mathbf{T}}_{ij}(\mathbf{y}_i) = \tilde{\mathbf{T}}_i(\mathbf{y}_i) = (\sum_{j=1}^{n_i} y_{ij1} / n_i, 0)^t$ were employed. In the above setting, $M = (0, 6.2) \times (0, 6.2)$ could be the domain (not including

Table 2. Relative errors for the wave equation model.

$J = 1$	$\sigma = 0.2$	$\sigma = 1$	$\sigma = 2$
Data	1.53139	1.31266	1.31326
PDE+Data	0.21631	0.99435	1.03351
PDE+IC	$2.20 \cdot 10^{-14}$	$1.03 \cdot 10^{-15}$	$9.64 \cdot 10^{-16}$
$J = 30$	$\sigma = 0.2$	$\sigma = 1$	$\sigma = 2$
Data	3.24145	1.57171	1.35731
PDE+Data	0.04394	0.19255	0.44219
PDE+IC	$3.93 \cdot 10^{-13}$	$2.60 \cdot 10^{-14}$	$6.61 \cdot 10^{-15}$

the initial points and boundaries) for the solution F and $V_{\partial Y}$ would be the collection of random vectors with the means $\partial M = \{(\mu_1, 0) : 0 < \mu_1 < 6.2\}$.

The minimizer of the objective function given in expression (6) and the minimizers of the objective function given in expression (4) with $S_t(\beta \mid \tilde{\mathbf{T}}_n)$ equal to the ones given in expressions (5) and (6), respectively, were computed. The objective function $S_2(\beta \mid \tilde{\mathbf{T}}_n)$ given in expression (6) only contains the information provided by the data, while the objective function given in expression (4) with $S_t(\beta \mid \tilde{\mathbf{T}}_n)$ equal to the one given in expression (6), i.e., $S(\beta \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n) = S_2(\beta \mid \tilde{\mathbf{T}}_n) + S_0(\beta \mid \mathbf{T}_n)$, contains the information provided by both the data and the partial differential equation. On the other hand, the objective function given in expression (4) with $S_t(\beta \mid \tilde{\mathbf{T}}_n)$ equal to the one given in expression (5), i.e., $S(\beta \mid \mathbf{T}_n, \tilde{\mathbf{T}}_n) = S_1(\beta \mid \tilde{\mathbf{T}}_n) + S_0(\beta \mid \mathbf{T}_n)$, contains the information provided by both the partial differential equation and associated initial and boundary conditions. Note that the partial derivative in $S_1(\beta \mid \tilde{\mathbf{T}}_n)$ needs to be modified in this example. As indicated by Table 2, the means of the relative errors corresponding to the minimizers of the objective function based on the data information only are significantly larger than the ones based on the other objective functions. Even for the data with small variations, the wave function can not be fitted well based on the data information only, as shown in the upper right part of Figure 3. The fits based on the other objective functions are quite consistent with the true wave function, as displayed by the lower left part of Figure 3 and the lower left part of Figure 5. However, for the data with large variations, the fits based on the information provided by the data and the partial differential equation are not consistent with the true wave function, as shown in the upper left and lower right parts of Figure 5. On the other hand, the minimizers of the objective function incorporating with the information from the partial differential equation and associated initial conditions are very close to the true coefficient vector and the fits are very consistent with the true wave function, as displayed in the upper right and lower parts of Figure 5. The setting for the data used

in Figure 3 and Figure 5 was $J = 1$. As $J = 30$ and $\sigma = 2$, the fits based on the data information only or both the data information and the partial differential equation were still not consistent with the true wave function. Nevertheless, for large or medium data variations, the larger the number of repeated observations is at each site, the smaller the relative errors corresponding to the minimizers of the objective function containing both the information from the data and the partial differential equation are, as presented in Table 2.

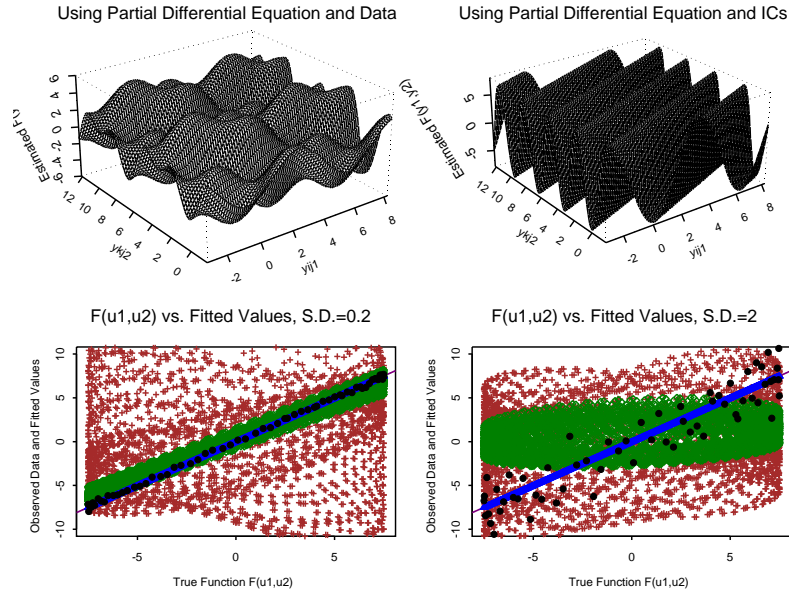


Fig. 5. The true wave Function $F(\mu_1, \mu_2)$ versus the fitted values based on different objective functions: Observed data (black \bullet); Fit using PDE+Data (green \diamond); Fit using PDE+IC (blue \triangle); Fit using the data information only (brown $+$).

4.2 A Modified Example in Thermodynamics

In thermodynamics, the equation of state for an ideal gas is $PV = nRT$ or $PV = Nk_B T$ (see Serway and Jewett, 2004, Chapter 19, Chapter 21), where P is the pressure, V is the volume, T is the temperature, n is the number of moles of gas in the sample, N is the total number of molecules, \mathcal{R} is the universal gas constant, and k_B is Boltzmann's constant. Britt and Luecke (1973) employed nonlinear implicit models for the estimation of parameters in the equation which approximated the modified equation of state for a gas, i.e., the equation (35) given in Britt and Luecke (1973).

In this subsection, based on the data given in Blancett, Hall, and Canfield (1970) and Hoover (1965), ordinary equation models were employed. The pressure measurements were scaled to fall in the interval $(0, 1)$ and a large value of the penalty parameter was pre-specified for data analyses. The objective equation is equation (39) given in Britt and Luecke (1973). However, the goal is to discover the underlying equations based on the pressure data satisfying

$$F(\boldsymbol{\mu}_{y_i}) = \beta_1(1 - N^*)\mu_{y_{i1}}\mu_{y_{i2}} + \beta_2(\mu_{y_{i2}} - N^*\mu_{y_{i1}})\mu_{y_{i1}}\mu_{y_{i2}} + \mu_{y_{i1}} - N^*\mu_{y_{i2}} = 0,$$

$i = 1, \dots, k$, for helium at $273.15^\circ K$ and methane at $131.93^\circ K$ rather than estimating the parameters β_1, β_2 , and N^* , where k is the number of expansions, β_1, β_2 are some parameters, N^* is the volume ratio, and $\mu_{y_{i1}}$ and $\mu_{y_{i2}}$ are the mean pressures of the $(i - 1)$ th and the i th expansions, respectively. The basis functions were

$$\psi_{3(j-1)+i}(\boldsymbol{\mu}) = \phi_i(\mu_1)\phi_j(\mu_2), \quad i = 1, 2, 3, \quad j = 1, 2, 3,$$

where the functions ϕ_i were $\phi_1(x) = 1$, $\phi_2(x) = x$, and $\phi_3(x) = x^2$. The conjugate gradient method was employed. Based on the results given in Table 4 of Britt and Luecke (1973), i.e., for helium at $273.15^\circ K$, $\beta_1 = 11.9517622$, $\beta_2 = 113.9619475$, $N^* = 1.564881$, the estimated equation using the nonlinear implicit models is

$$\begin{aligned} F_{he}(\boldsymbol{\mu}_{y_i}) = & -6.751323\mu_{y_{i1}}\mu_{y_{i2}} + 113.9619475(\mu_{y_{i2}} - 1.564881\mu_{y_{i1}})\mu_{y_{i1}}\mu_{y_{i2}} \\ & + \mu_{y_{i1}} - 1.564881\mu_{y_{i2}} = 0, \end{aligned}$$

which is equivalent to the equation

$$\begin{aligned} \sum_{l=1}^9 \beta_l^* \psi_l(\boldsymbol{\mu}) = & \frac{1}{1.000053} [-0.004722\psi_2(\boldsymbol{\mu}) + 0.00739\psi_4(\boldsymbol{\mu}) + 0.03188\psi_5(\boldsymbol{\mu}) \\ & + 0.8422\psi_6(\boldsymbol{\mu}) - 0.5382\psi_8(\boldsymbol{\mu})] = 0. \end{aligned}$$

The data for 100% helium at $273.15^\circ K$ given in Table I^* of Blancett, Hall, and Canfield (1970) consists of twelve pressure measurements for each run. The relative error $\|\hat{\boldsymbol{\beta}}_\lambda(\mathbf{T}_n) - \boldsymbol{\beta}^*\|_{R^9} / \|\boldsymbol{\beta}^*\|_{R^9}$ for the data in each run was computed. For the two runs, the relative errors were 3.94% and 3.46%, respectively. Similarly, based on the results given in Table 5 of Britt and Luecke (1973), i.e., for methane at $131.93^\circ K$, $\beta_1 = -222.9$, $\beta_2 = -24358$, $N^* = 1.14962$, the estimated equation using the nonlinear implicit models is

$$\begin{aligned} F_{me}(\boldsymbol{\mu}_{y_i}) = & 33.3503\mu_{y_{i1}}\mu_{y_{i2}} - 24358(\mu_{y_{i2}} - 1.14962\mu_{y_{i1}})\mu_{y_{i1}}\mu_{y_{i2}} \\ & + \mu_{y_{i1}} - 1.14962\mu_{y_{i2}} = 0. \end{aligned}$$

The relative error for the data given in Table 2 of Hoover (1965) was 5%. This indicated that the estimated equations using the ordinary equation models were quite consistent with the ones using the nonlinear implicit models.

In addition, to evaluate the performance of the proposed models and methods under different data variations, $F_{he}(\boldsymbol{\mu}_{y_i}) = 0$ and $F_{me}(\boldsymbol{\mu}_{y_i}) = 0$ were assumed to be the true equations and the normal random errors were considered. The pressure data with means satisfying the above equations corresponding to different coefficients of variation, 10%, 30% and 50%, were generated. For each coefficient of variation, 500 simulated data sets with the number of sites equal to 11 and the number of repeated observations equal to 1 were generated and the average relative error was computed. The average relative errors corresponding to helium at $273.15^\circ K$ were smaller than 4%, while the ones corresponding to methane at $131.93^\circ K$ were smaller than 7%. As shown in the right part of Figure 2, a sensible fit can be still obtained for the data with large coefficient of variation corresponding to helium at $273.15^\circ K$.

5. Concluding Discussions

For the data with means satisfying unknown equations, the mathematical equation models can be employed to obtain sensible equations for fitting the data. On the other hand, for the known deterministic equations, the proposed models can be still useful for fitting the experimental data subject to random variations.

The equation estimation turns out to be associated with the nonlinear programming problems subject to the randomness of the coefficients. However, the nonlinear programming problems involving the random coefficients, for instances, the random vectors in the mathematical equation models or some measurements subject to random variations in other situations, have not attracted much attention in the literature. The results concerning the optimality conditions and the convergence of the methods provide the basic theoretical facts for these problems.

If the variations of the data are large, the minimizers of the objective functions given in expressions (1) and (4) may not be accurate estimates of the coefficient vector $\boldsymbol{\beta}^*$. As indicated by Theorem 3.1 and the simulation study, the improvement on the accuracy can be made by increasing the number of repeated observations at each site. It is possible that the repeated observations are not available. In such situation, one possible solution is to “cluster” the data nearly, i.e., the distances among them being small, and then to consider the clustered data as the repeated observations at one site.

If the variations of the data are small, the algorithms employed in Section 4 perform well. If the variations of the data or the number of basis functions are large (see Nocedal and Wright, 1999, Chapter 5) and the objective function is differentiable, the CG method might be a sensible choice. In particular, as a good quadratic approximation for the objective function in the neighborhood of the true solution might exist, the CG method might perform well, e.g. the mathematical equation models with the objective function given in expression (1) or the linear nonhomogeneous partial differential equation models with the objective function given in expression (10). On the other hand,

if the objective function for the mathematical equation models is not differentiable or the gradient of the objective function or its approximation might be difficult to compute, the gradient free methods, such as the Nelder-Mead algorithm or the method of Rosenbrock (1960), could be used.

References

- Adams, R. A. and Fournier, J. J. F. (2003). *Sobolev spaces, 2nd edn.* Academic Press, Boston.
- Aronszajn, N. (1950). Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, **68**, 337-404.
- Aubin, J. P. (2000). *Applied functional analysis, 2nd edn.* Wiley, New York.
- Bazaraa, M. S. and Shetty, C. M. (1979). *Nonlinear programming: Theory and algorithms.* Wiley, New York.
- Berlinet, A. and Thomas-Agnan, C. (2004). *Reproducing kernel Hilbert spaces in probability and statistics.* Kluwer Academic, Boston.
- Billingsley, P. (1999). *Convergence of probability measures, 2nd edn.* Wiley, New York.
- Blancett, A. L., Hall, K. R. and Canfield, F. B. (1970). Isotherms for the He-Ar system at 50°C, 0°C, and -50°C up to 700 atm. *Physica*, **47**, 75-91.
- Britt, H. I. and Luecke, R. H. (1973). The estimation of parameters in nonlinear implicit models. *Technometrics*, **15**, 233-247.
- Broyden, C. G. (1967). Quasi-Newton methods and their application to function minimisation. *Math. Computn.*, **21**, 368-381.
- Broyden, C. G. (1970). The convergence of a class of double-rank minimization algorithms, II: The new algorithm. *J. Inst. Math. Appl.*, **6**, 222-231.
- Cavalier, L. (2011). Inverse problems in statistics. *Inverse problems and high-dimensional estimation: Stats in the Château summer school, August 31 - September 4, 2009, Lecture Notes in Statistics 203.* Springer, Berlin, Heidelberg.
- Chester, C. R. (1970). *Techniques in partial differential equations.* McGraw-Hill, New York.
- Da Prato, G. and Zabczyk, J. (1992). *Stochastic equations in infinite dimensions.* Cambridge University Press, Cambridge.

- Davidon, W. C. (1959). Variable metric method for minimization. Technical Report ANL-5990 (revised), Argonne National Laboratory.
- Ekeland, I. and T  mam, R. (1999). *Convex analysis and variational problems*. SIAM, Philadelphia.
- Eubank, R. L. (1988). *Spline smoothing and nonparametric regression*. Dekker, New York.
- Evans, L. C. (1998). *Partial differential equations*. AMS, Providence.
- Fletcher, R. (1970). A new approach to variable metric algorithms. *Compu. J.*, **13**, 317-322.
- Fletcher, R. and Powell, M. J. D. (1963). A rapidly convergent descent method for minimization. *Compu. J.*, **6**, 163-168.
- Fletcher, R. and Reeves, C. M. (1964). Function minimization by conjugate gradients. *Compu. J.*, **7**, 149-154.
- Goldfarb, D. (1970). A family of variable-metric methods derived by variational means. *Math. Computn.*, **24**, 23-26.
- Hoover, A. E. (1965). Virial coefficients of methane and ethane. Ph.D. Thesis, Department of Chemical Engineering, Rice University.
- Jost, J. (2002). *Partial differential equations*. Springer, New York.
- Karush, W. (1939). Minima of functions of several variables with inequalities as side conditions. Master Thesis, Department of Mathematics, University of Chicago.
- Kiefer, J. (1953). Sequential minimax search for a maximum. *Proc. Amer. Math. Soc.*, **4**, 502-506.
- Kreyszig, E. (1978). *Introductory functional analysis with applications*. Wiley, New York.
- Kuhn, H. W. and Tucker, A. W. (1951). Nonlinear programming. *Proc. 2nd Berkeley Symposium on Math. Statist. and Probab.*, 481-492.
- Ledoux M. and Talagrand M. (1991). *Probability in Banach spaces*. Springer-Verlag, Berlin Heidelberg.
- Meyers, N. and Serrin, J. (1964). H=W. *Proc. Nat. Acad. Sci. U.S.A.*, **51**, 1055-1056.
- Nelder, J. A. and Mead, R. (1965). A simplex method for function minimization. *Compu. J.*, **7**, 308-313.

- Nocedal, J. and Wright, S. J. (1999). *Numerical optimization*. Springer, New York.
- Rheinboldt, W. C. (1998). *Methods for solving systems of nonlinear equations, 2nd edn.* SIAM, Philadelphia.
- Rosenbrock, H. H. (1960). An automatic method for finding the greatest or least value of a function. *Compu. J.*, **3**, 175-184.
- Serfling, R. J. (1980). *Approximation theorems of mathematical statistics*. Wiley, New York.
- Serway, R. A. and Jewett, J. W. (2004). *Physics for scientists and engineers with modern physics, 6th edn.* Thomson, Belmont.
- Shanno, D. F. (1970). Conditioning of quasi-Newton methods for function minimization. *Math. Computn.*, **24**, 647-656.
- Thisted, R. A. (1988). *Elements of statistical computing*. Chapman and Hall, New York.
- Tveito, A. and Winther, R. (1998). *Introduction to partial differential equations: A computational approach*. Springer, New York.

[Received February 2014; accepted August 2014.]

*Journal of the Chinese
Statistical Association*
Vol. 52, (2014) 497-532

數學方程式模式

魏文翔

東海大學統計系

摘 要

本文提出一類利用數學方程式來配適資料之模式。這類模式包含了常被使用之統計模式, 例如線性迴歸模式 (Linear regression models)、無母數迴歸模式 (Nonparametric regression models)、線性混合隨機效應模式 (Linear mixed-effects models) 以及測量誤差模式 (Measurement error models)。這類模式所牽涉之數學方程式亦可是一偏微分方程式 (Partial differential equation)。非線性規劃方法 (Nonlinear programming) 可用來估計這類模式所牽涉之數學方程式。相關估計方法之理論結果亦被建立。模擬研究以及一修正之熱力學 (Thermodynamics) 實例被用來闡明所提出之模式及其相關之估計方法。

關鍵詞: 數學方程式模式 (Mathematical equation models), 非線性規劃, 偏微分方程式, 懲罰函數法 (Penalty function methods), 複製核希爾伯特空間 (Reproducing kernel Hilbert space)。

JEL classification: C61.