

Generalized profiled estimation

Weixin Wang & Jianbin Tan

University of Science and Technology of China

2022.10.19

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems
- 4 A New Efficient Method for ODEs Linear in Parameters
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

Table of Contents

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems
- 4 A New Efficient Method for ODEs Linear in Parameters
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

Generalized profiled estimation

- In this two chapters, we will introduce a profiled estimation for parameter estimation in mechanistic models and discuss its statistical consistency.
- Related algorithms to improve computation efficiency will also be discussed.

Review of current ordinary differential equation parameter estimation strategies

- Data fitting by numerical approximation of an initial value problem.
- Collocation methods or basis function expansions.

Review of parameter estimation strategies

Data fitting by numerical approximation of an initial value problem (NLS method)

- Approximate solutions of ODEs over a range $[t_0, t_1]$ use fixed initial values $x_0 = x(t_0)$.
- One numerical method: Runge–Kutta algorithm.
- Solve the sensitivity differential equations:

$$\frac{d}{dt} \left(\frac{d\mathbf{x}}{d\theta} \right) = \frac{\partial \mathbf{f}}{\partial \theta} + \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \frac{d\mathbf{x}}{d\theta}, \quad \text{with } \left. \frac{d\mathbf{x}}{d\theta} \right|_{t=0} = 0.$$

- In the event that $x(0) = x_0$ must also be estimated, the corresponding sensitivity equations are:

$$\frac{d}{dt} \left(\frac{d\mathbf{x}}{d\mathbf{x}_0} \right) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \frac{d\mathbf{x}}{d\mathbf{x}_0}, \quad \text{with } \left. \frac{d\mathbf{x}}{d\mathbf{x}_0} \right|_{t=0} = \mathbf{I}.$$

NLS procedure has many problems:

- Computationally intensive
- Inaccuracy of the numerical approximation
- Size of the parameter set may be increased by the set of initial conditions that are needed

Collocation methods or basis function expansions

- Express the approximation \hat{x}_i of x_i in terms of a basis function expansion:

$$\hat{x}_i(t) = \sum_k^{K_i} c_{ik} \phi_{ik}(t) = \mathbf{c}_i' \boldsymbol{\phi}_i(t)$$

- Procedure:
 - Each x_i is first estimated by data smoothing methods
 - Minimization of a least squares measure of the fit of $d\hat{\mathbf{x}}/dt$ to $\mathbf{f}(\hat{\mathbf{x}}, \mathbf{u}, t \mid \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$
 - Iteration's roughness penalty: $\|d\hat{\mathbf{x}}/dt - \mathbf{f}(\hat{\mathbf{x}}, \mathbf{u}, t \mid \boldsymbol{\theta})\|$

Table of Contents

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems
- 4 A New Efficient Method for ODEs Linear in Parameters
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

The Parameter Cascading Strategy for Estimating Parameters

The elements of parameter cascading

- Approximate the solution $x(t)$ to the differential equation by the basis function expansion

$$\hat{x}(t) = \sum_k^K c_k \phi_k(t) = \mathbf{c}^T \boldsymbol{\phi}(t).$$

- Constrain the coefficient vector \mathbf{c} to be a smooth function $\mathbf{c}(\boldsymbol{\theta})$ of the parameters in vector $\boldsymbol{\theta}$ that define the differential equation.
- Optimize the fit to observations y_1, \dots, y_n in vector \mathbf{y} by $\hat{\mathbf{x}}(\mathbf{t}) = \mathbf{c}(\boldsymbol{\theta})^T \boldsymbol{\phi}(\mathbf{t})$ with respect to $\boldsymbol{\theta}$ by optimizing a fitting criterion $H(\boldsymbol{\theta} \mid \mathbf{y})$.

Two Classes of Parameters

Main idea

- Replace the coefficients c_k in the basis expansion by values of a smooth function of the parameters in θ
- The status of the K-dimensional coefficient vector c as a set of parameters is eliminated
- The coefficient estimates are actually the values of the function $c(\hat{\theta})$ at the optimal value of the θ

So why we do this?

Two Classes of Parameters

Why parameter cascading?

- Nuisance parameters? Profiled?
- Prevent over-fitting the data
- Reduce computation required to optimize the fit of the model to the data

We will see over and over again that forcing $c(\theta)$ to be a smooth function of the parameters will result in better parameter estimates with smaller confidence intervals without seriously harming the fit to the data that $c(\theta)$ defines.

Defining Coefficients as Functions of Parameters

Approach to defining $c(\theta)$

- Define an inner fitting criterion $J(\mathbf{c} \mid \theta)$ that optimizes \mathbf{c} given any candidate value for the parameters in θ
- Then that is re-optimized each time θ is changed in order to optimize an outer fitting criterion $H(\theta)$
- Higher priority on the estimation θ ; Lower priority on the estimation of \mathbf{c} and the $x(t)$

How can we find a way of blending together the information in the data and the information in the equation?

The Symmetric Relation Between the Data and the Differential Equation

Just capturing the prominent shape features

- Fix a value $\rho \in [0, 1)$, called a smoothing parameter or a bandwidth parameter, that **modulates the relative emphasis on fitting the data as opposed to satisfying the differential equation**
- By choosing ρ judiciously, we can see how much of the variation in the data can be accommodated by what is perhaps a too-simple dynamic model. If this is possible, then the residual variation is apt to give important cues about either how the model is to be elaborated, or about what kind and quantity of data would be more revealing.

Inner Optimization Criterion J

Define the inner optimization function $J(\mathbf{c} \mid \theta, \rho)$ (one dimension):

$$J(\mathbf{c} \mid \theta, \rho) = (1 - \rho) \sum_j^n [y_j - x(t_j)]^2 / n + \rho \int_0^T \{D^m x(t) - f[x(t)]\}^2 dt / T$$

Remark:

- Squared residual fit measures
- The smoothing parameter $\rho \in [0, 1)$
- $\rho = 0 \Rightarrow$ ignore the differential equation entirely and concentrate solely on data-fitting
- $\rho \rightarrow 1 \Rightarrow$ place more and more emphasis on x being close to a solution to the differential equation
- $\rho < 1 \Rightarrow$ not want to ignore the data entirely

A solution to most differential equations is not uniquely defined without some additional information, and the data provides this information.

The Least Squares Cascade Coefficient Function

Some notations

- Differential equation residual function $D^m x(t) - f[x(t) | \theta]$ can be re-expressed as $Lx = 0$ where L is the differential operator form of the differential equation $D^m x = f(x | \theta)$
- The fitting function is $\hat{x}(t) = \mathbf{c}^T \phi(t)$, where the coefficient vector \mathbf{c} and basis function vector ϕ are of length K

Then inner criterion becomes:

$$\begin{aligned} J(\mathbf{c} | \theta) &= (1 - \rho)(\mathbf{y} - \Phi \mathbf{c})^T (\mathbf{y} - \Phi \mathbf{c}) / n + \rho \int_0^T [\mathbf{c}^T L\phi(t)] [\mathbf{c}^T L\phi(t)]^T dt \\ &= (1 - \rho)(\mathbf{y} - \Phi \mathbf{c})^T (\mathbf{y} - \Phi \mathbf{c}) / n + \rho \mathbf{c}^T \left(\int_0^T [L\phi(t)] [L\phi(t)^T] dt \right) \mathbf{c} \\ &= (1 - \rho)(\mathbf{y} - \Phi \mathbf{c})^T (\mathbf{y} - \Phi \mathbf{c}) / n + \rho \mathbf{c}^T \mathbf{R}(\theta) \mathbf{c} / T \end{aligned}$$

where

$$\mathbf{R}(\theta) = \int_0^T [L\phi(t)] [L\phi(t)^T] dt$$

The minimizing value of \mathbf{c}

Inner criterion

$$J(\mathbf{c} \mid \boldsymbol{\theta}) = (1 - \rho)(\mathbf{y} - \boldsymbol{\Phi}\mathbf{c})^T(\mathbf{y} - \boldsymbol{\Phi}\mathbf{c})/n + \rho\mathbf{c}^T\mathbf{R}(\boldsymbol{\theta})\mathbf{c}/T$$

In matrix calculus to work out that the minimizing value of \mathbf{c} that satisfies the equation

$$\mathbf{c}(\boldsymbol{\theta}) = [(1 - \rho)\boldsymbol{\Phi}^T\boldsymbol{\Phi}/n + \rho\mathbf{R}(\boldsymbol{\theta})/T]^{-1} (1 - \rho)\boldsymbol{\Phi}^T\mathbf{y}/n$$

More smooth as we increase ρ

This figure shows how the coefficients $c_3(\rho)$, $c_6(\rho)$, and $c_9(\rho)$ vary as β ranges from 0.5 to 1.5 for the relatively low smoothing parameters $\rho = 0.27$, and for the high value $\rho = 0.98$

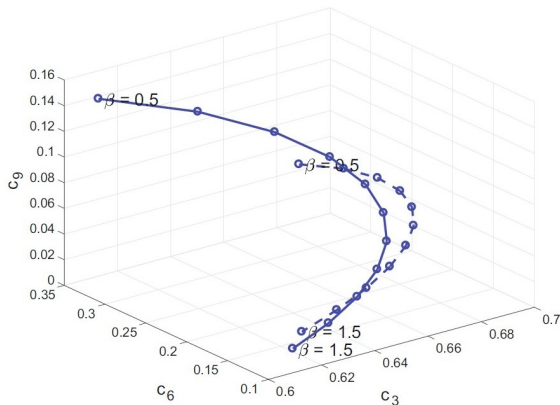


Fig. 9.2 The curves display three of the 13 coefficient functions $c_3(\beta)$, $c_6(\beta)$, and $c_9(\beta)$ defining $x(t)$ varying over values of β between 0.5 and 1.5 for the model and data in Fig. 9.1, and for values of ρ of 0.27 (dashed) and 0.98 (solid)

The Outer Fitting Criterion H

To complete the parameter cascading strategy, we specify an outer fitting criterion

$$H(\boldsymbol{\theta} \mid \rho) = G_0[\mathbf{y}, \mathbf{x}(\mathbf{t}) \mid \boldsymbol{\theta}, \rho].$$

If we stick with least squares as our criterion, then

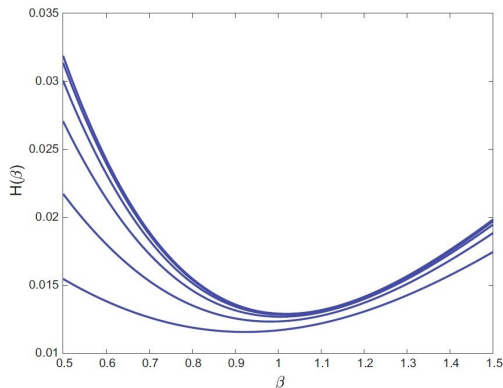
$$H(\boldsymbol{\theta} \mid \rho) = \sum_j [y_j - \hat{x}(t_j)]^2 = \sum_j [y_j - \mathbf{c}^T(\boldsymbol{\theta} \mid \rho)\boldsymbol{\phi}(t_j)]^2$$

Remark:

- Use the Gauss-Newton algorithm
- The shape of the criterion function also varies with the level of ρ

The shape of the criterion function also varies with the level of ρ

- As we increase ρ the height of H increases, its asymmetry increases, as does its curvature, but we also see that the minimizing values moves closer and closer to the true value $\beta = 1$.



Choosing the Smoothing Parameter ρ

A simple and effective answer is to estimate both θ and x for a wide range of ρ values so as to calculate the functional flow $x(\rho)$ or parametric flow $\theta(\rho)$.

Degrees of freedom

$$df(\theta) = \text{trace} [2\mathbf{M}(\theta) - \mathbf{M}(\theta)\mathbf{M}(\theta)^T] .$$

where

$$\mathbf{M}(\theta) = (1 - \rho)\Phi [(1 - \rho)\Phi^T\Phi/n + \rho\mathbf{R}(\theta)/T]^{-1} \Phi^T/n$$

Remark:

- Refers to the effective dimensionality of the fitting function as a function of ρ .
- For $\rho \rightarrow 0$, it turns out that this measure converges to K , the number of basis functions.
- But as $\rho \rightarrow 1$, the measure converges to $K - \text{rank}(\mathbf{M}) = m$, which corresponds to the dimensionality of the solution space for the differential equation.

GCV criterion

$$GCV(\rho) = \frac{n}{[n - df(\rho)]^2} SSE(\rho)$$

where $SSE(\rho)$ is the residual sum of squares

- the value of ρ that minimizes this GCV measure will provide a nearly optimal estimate of x most of the time if the residuals are reasonably close to being uncorrelated.

The general formulation:

$$D^{M_i} x_i(t) = \sum_{k=1}^d \sum_{j=0}^{M_k-1} b_{ijk} \beta_{ijk}(t | \boldsymbol{\theta}) D^j x_k(t) + \sum_{\ell}^* a_{i\ell} \alpha_{i\ell}(t | \boldsymbol{\theta}) u_{i\ell}(t), \quad M_i \geq 0, \quad i = 1, \dots, d.$$

The inner criterion J will now be

$$J(\mathbf{c} | \boldsymbol{\theta}) = (\mathbf{y} - \boldsymbol{\Phi} \mathbf{c})^T \mathbf{W}_1 (\mathbf{y} - \boldsymbol{\Phi} \mathbf{c}) + \mathbf{c}^T \mathbf{W}_2 \mathbf{R}(\boldsymbol{\theta}) \mathbf{c} + \rho \mathbf{c}^T \mathbf{W}_2 \mathbf{S}(\boldsymbol{\theta}),$$

Then we have that

$$\mathbf{c}(\boldsymbol{\theta}) = [\mathbf{W}_1 \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \mathbf{W}_2 \mathbf{R}(\boldsymbol{\theta})]^{-1} [\mathbf{W}_1 \boldsymbol{\Phi}^T \mathbf{y} + \mathbf{W}_2 \mathbf{S}(\boldsymbol{\theta})]$$

Analysis of the Head Impact Data (measure the effects of motorcycle accident on the driver's brain tissue)

- We analyzed the data using the three-parameter forced damped harmonic equation

$$D^2x(t) = -\beta_0x(t) - \beta_1Dx(t) + \alpha u(t)$$

Analysis of the Head Impact Data

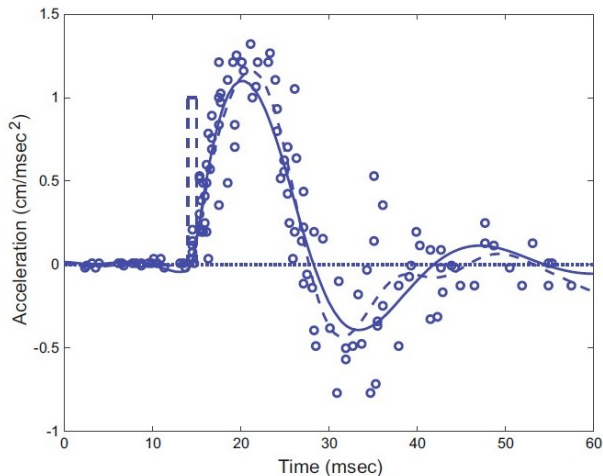


Fig. 9.5 The *circles* indicate observations of acceleration of brain tissue from five replications of an experiment involving striking the cranium of a corpse with a blunt object. The *box* was constructed to represent the impact itself spread over one time unit. The *dashed* and *solid* lines indicate the fits from the parameter cascading analysis of the data in Fig. 9.5 using the Eq. (9.22) and corresponding to $\rho = 0.5$ and $\rho = 0.998$, respectively

Analysis of the Head Impact Data

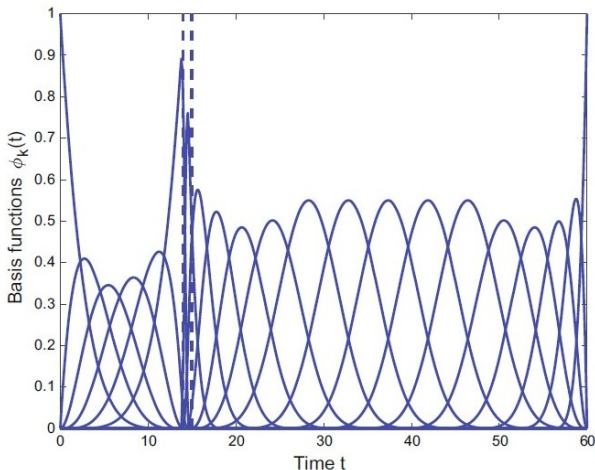


Fig. 9.6 The *solid lines* are the 21 order 6 B-spline basis functions active in a region around the box forcing function for the head impact data. The *vertical dashed lines* indicate the two locations of three coincident knots

Analysis of the Head Impact Data

Figure 9.7 shows how the parameter estimates vary over ρ values 0.50, 0.73, 0.88, 0.95, 0.98, 0.99 and 0.998. The values of parameters β_0 , β_1 , and α , along with two standard errors at $\rho = 0.998$ are 0.056 ± 0.011 , 0.128 ± 0.065 and 0.383 ± 0.101 , respectively. The GCV criterion was minimized at 0.556 for $\rho = 0.98$, corresponding to 10.3 degrees of freedom, and for $\rho = 0.998$, $\text{GCV} = 0.574$ corresponding to 6.4 degrees of freedom.

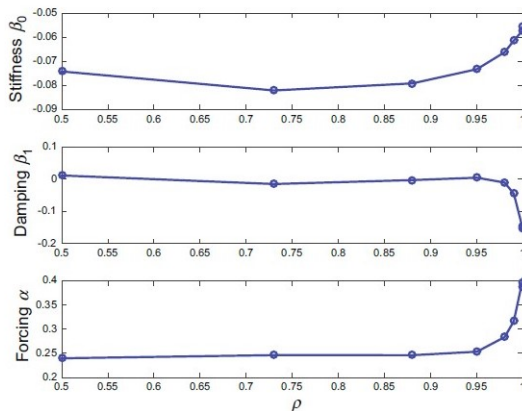


Fig. 9.7 The three parameter values as a function ρ , the largest value of which is 0.998

Table of Contents

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems**
- 4 A New Efficient Method for ODEs Linear in Parameters
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

The Setup for Parameter Cascading

- In the least squares data-fitting case

$$J(\mathbf{c} \mid \boldsymbol{\theta}, \lambda) = \sum_j^n [y_j - x(t_j)]^2 + \lambda \int_0^T \{D^m x(t) - f[x(t), \mathbf{u}(t) \mid \boldsymbol{\theta}]\}^2 Dt$$

where the bandwidth or smoothing parameter $\lambda > 0$ controls the emphasis on fitting the equation but the data-fitting first term remains fixed

The outer criterion

- The outer criterion $H(\theta \mid \lambda)$ defines fit in terms of only the parameter vector θ . The total derivative of H with respect to θ must allow for the fact that H may depend both directly on θ and indirectly via its dependency on $\mathbf{c}(\theta)$:

$$\frac{dH}{d\theta} = \frac{\partial H}{\partial \theta} + \left(\frac{\partial H}{\partial \mathbf{c}} \right) \left(\frac{d\mathbf{c}}{d\theta} \right)$$

- Although for most fitting criteria, including least squares, H has no direct dependency at all on θ , so that $\partial H / \partial \theta = 0$ and can be dropped from the total derivative equation. It is the derivative $d\mathbf{c} / d\theta$ that is the missing link, since we no longer have an explicit expression for \mathbf{c} .

Numerical methods must be used to achieve the optimum of J (implicit differentiation or the Implicit Function Theorem)

- Suppose that we have optimized J using a high quality numerical optimization strategy to obtain an optimum value \mathbf{c}
- $\partial J / \partial \mathbf{c} \approx 0$ to a high level of accuracy
- We now compute the total derivative $d/d\theta$ of the gradient of J with respect to θ , which, since $\partial J / \partial \mathbf{c} = 0$, will also be 0 :

$$\frac{d}{d\theta} \left(\frac{\partial J}{\partial \mathbf{c}} \right) = \frac{\partial^2 J}{\partial \mathbf{c} \partial \theta} + \left(\frac{\partial^2 J}{\partial \mathbf{c}^2} \right) \left(\frac{d\mathbf{c}}{d\theta} \right) = 0$$

Implicit Function Theorem

$$\frac{d}{d\theta} \left(\frac{\partial J}{\partial \mathbf{c}} \right) = \frac{\partial^2 J}{\partial \mathbf{c} \partial \theta} + \left(\frac{\partial^2 J}{\partial \mathbf{c}^2} \right) \left(\frac{d\mathbf{c}}{d\theta} \right) = 0$$

so that

$$\frac{d\mathbf{c}}{d\theta} = - \left(\frac{\partial^2 J}{\partial \mathbf{c}^2} \right)^{-1} \left(\frac{\partial^2 J}{\partial \mathbf{c} \partial \theta} \right)$$

and, substituting this into (10.3) we see that

$$\frac{dH}{d\theta} = \frac{\partial H}{\partial \theta} - \left(\frac{\partial H}{\partial \mathbf{c}} \right) \left(\frac{\partial^2 J}{\partial \mathbf{c}^2} \right)^{-1} \left(\frac{\partial^2 J}{\partial \mathbf{c} \partial \theta} \right).$$

Then we can use gradient descent method or Gauss-Newton method to optimize.

Nonlinear Systems and Other Fitting Criteria

- We use the following expression of the inner fitting criterion involving a general data-fitting function G_J , in the special case where all x_i trajectories are expanded using the same set ϕ of K basis functions:

$$J(\mathbf{c} \mid \boldsymbol{\theta}, \lambda) = G_J[\mathbf{Y}, \mathbf{C}(\boldsymbol{\theta} \mid \lambda)\phi] + \sum_i^d \lambda_i \int_0^T \{Dx_i(t) - f[\mathbf{x}(t), \mathbf{u}_i(t) \mid \boldsymbol{\theta}]\}^2 D$$

- The coefficient matrix \mathbf{C} will, in this simple case, have K rows and d columns. In this formulation, also for simplicity, we retain the integrated squared residuals in the penalty term, but we certainly envisage other penalty structures that might be more appropriate.

- The corresponding expression for the outer criterion H is

$$H(\boldsymbol{\theta} \mid \lambda) = G_H[\mathbf{Y}, \mathbf{C}(\boldsymbol{\theta} \mid \lambda)\boldsymbol{\phi}].$$

- However, the two data fitting functions G_J and G_H may in many applications **be the same**.

Approximating the sampling variation of $\hat{\theta}$ and \hat{c}

Let Σ be the variance-covariance matrix for \mathbf{y} . Making explicit the dependence of H on the data \mathbf{y} by using the notation $H(\theta | \mathbf{y})$, the estimate $\hat{\theta}(\mathbf{y})$ of θ is the solution of the stationary equation $\partial H(\theta | \mathbf{y}) / \partial \theta = 0$.

The usual δ -method that is employed in non-linear least squares produces a variance estimate of the form

$$\text{var}_{\text{GN}}\{\hat{\theta}(\mathbf{y})\} \approx \sigma^2 \left\{ \left(\frac{d\hat{\mathbf{x}}}{d\theta} \right)' \left(\frac{d\hat{\mathbf{x}}}{d\theta} \right) \right\}^{-1}$$

by making use of the approximation

$$\frac{d^2 H}{d\theta^2} \approx \left(\frac{d\hat{\mathbf{x}}}{d\theta} \right)' \left(\frac{d\hat{\mathbf{x}}}{d\theta} \right)$$

Approximating the sampling variation of $\hat{\theta}$ and $\hat{\mathbf{c}}$

By applying the implicit function theorem to $\partial H / \partial \theta$ as a function of \mathbf{y} , we may say that for any \mathbf{y} in \mathcal{N} there is a value $\hat{\theta}(\mathbf{y})$ satisfying $\partial H / \partial \theta = 0$. By taking the \mathbf{y} -derivative of this relation, we obtain

$$\frac{d}{d\mathbf{y}} \left(\frac{dH}{d\theta} \Big|_{\hat{\theta}(\mathbf{y})} \right) = \frac{d^2 H}{d\theta d\mathbf{y}} \Big|_{\hat{\theta}(\mathbf{y})} + \frac{d^2 H}{d\theta^2} \Big|_{\hat{\theta}(\mathbf{y})} \frac{d\hat{\theta}}{d\mathbf{y}} = 0,$$

where

$$\frac{d^2 H}{d\theta^2} = \frac{\partial^2 H}{\partial \theta^2} + \frac{\partial^2 H}{\partial \theta \partial \hat{\mathbf{c}}} \frac{\partial \hat{\mathbf{c}}}{\partial \theta} + \left(\frac{\partial \hat{\mathbf{c}}}{\partial \theta} \right)' \frac{\partial^2 H}{\partial \hat{\mathbf{c}} \partial \theta} + \left(\frac{\partial \hat{\mathbf{c}}}{\partial \theta} \right)' \frac{\partial^2 H}{\partial \hat{\mathbf{c}}^2} \frac{\partial \hat{\mathbf{c}}}{\partial \theta} + \frac{\partial H}{\partial \hat{\mathbf{c}}} \frac{\partial^2 \hat{\mathbf{c}}}{\partial \theta^2}$$

and

$$\frac{d^2 H}{d\theta d\mathbf{y}} = \frac{\partial^2 H}{\partial \theta \partial \mathbf{y}} + \frac{\partial^2 H}{\partial \hat{\mathbf{c}} \partial \mathbf{y}} \frac{\partial \hat{\mathbf{c}}}{\partial \theta} + \frac{\partial^2 H}{\partial \theta \partial \hat{\mathbf{c}}} \frac{\partial \hat{\mathbf{c}}}{\partial \mathbf{y}} + \left(\frac{\partial \hat{\mathbf{c}}}{\partial \theta} \right)' \frac{\partial^2 H}{\partial \hat{\mathbf{c}}^2} \frac{\partial \hat{\mathbf{c}}}{\partial \mathbf{y}} + \frac{\partial H}{\partial \hat{\mathbf{c}}} \frac{\partial^2 \hat{\mathbf{c}}}{\partial \theta \partial \mathbf{y}}$$

Approximating the sampling variation of $\hat{\theta}$ and \hat{c}

Solve this equation

$$\frac{d}{d\mathbf{y}} \left(\frac{dH}{d\boldsymbol{\theta}} \Big|_{\hat{\boldsymbol{\theta}}(\mathbf{y})} \right) = \frac{d^2 H}{d\boldsymbol{\theta} d\mathbf{y}} \Big|_{\hat{\boldsymbol{\theta}}(\mathbf{y})} + \frac{d^2 H}{d\boldsymbol{\theta}^2} \Big|_{\hat{\boldsymbol{\theta}}(\mathbf{y})} \frac{d\hat{\boldsymbol{\theta}}}{d\mathbf{y}} = 0,$$

We obtain the first derivative of $\hat{\boldsymbol{\theta}}$ with respect to \mathbf{y} :

$$\frac{d\hat{\boldsymbol{\theta}}}{d\mathbf{y}} = - \left(\frac{\partial^2 H}{\partial \boldsymbol{\theta}^2} \Big|_{\hat{\boldsymbol{\theta}}(\mathbf{y})} \right)^{-1} \left(\frac{\partial^2 H}{\partial \boldsymbol{\theta} \partial \mathbf{y}} \Big|_{\hat{\boldsymbol{\theta}}(\mathbf{y})} \right)$$

Approximating the sampling variation of $\hat{\theta}$ and \hat{c}

Let $\mu = E(\mathbf{y})$; the first-order Taylor series expansion for $d\hat{\theta}/d\mathbf{y}$ is

$$\frac{d\hat{\theta}}{d\mathbf{y}} \approx \frac{d\hat{\theta}}{d\mu} + \frac{d^2\hat{\theta}}{d^2\mu}(\mathbf{y} - \mu).$$

When $d^2\hat{\theta}/d^2\mu$ is uniformly bounded, we can take the expectation on both sides of approximation and derive $E(d\hat{\theta}/d\mu) \approx E(d\hat{\theta}/d\mathbf{y})$. We can also approximate $\hat{\theta}(\mathbf{y})$ by using the first-order Taylor series expansion:

$$\hat{\theta}(\mathbf{y}) \approx \hat{\theta}(\mu) + \frac{d\hat{\theta}}{d\mu}(\mathbf{y} - \mu)$$

Approximating the sampling variation of $\hat{\theta}$ and \hat{c}

We derive

$$\text{var}\{\hat{\theta}(\mathbf{y})\} \approx \left(\frac{d\hat{\theta}}{d\mu} \right) \Sigma \left(\frac{d\hat{\theta}}{d\mu} \right)' \approx \left(\frac{d\hat{\theta}}{d\mathbf{y}} \right) \Sigma \left(\frac{d\hat{\theta}}{d\mathbf{y}} \right)',$$

since

$$E \left(\frac{d\hat{\theta}}{d\mu} \right) \approx E \left(\frac{d\hat{\theta}}{d\mathbf{y}} \right)$$

Approximating the sampling variation of $\hat{\theta}$ and \hat{c}

Similarly, the sampling variance of $\hat{c}\{\hat{\theta}(\mathbf{y})\}$ is estimated by

$$\text{var}[\hat{c}\{\hat{\theta}(\mathbf{y})\}] = \left(\frac{d\hat{c}}{d\mathbf{y}} \right) \Sigma \left(\frac{d\hat{c}}{d\mathbf{y}} \right)',$$

where

$$\frac{d\hat{c}}{d\mathbf{y}} = \frac{d\hat{c}}{d\hat{\theta}} \frac{d\hat{\theta}}{d\mathbf{y}} + \frac{\partial \hat{c}}{\partial \mathbf{y}}$$

Table of Contents

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems
- 4 A New Efficient Method for ODEs Linear in Parameters**
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

A new efficient method for ODEs linear in parameters

Description of the algorithm

- We iteratively generate a sequence of estimates $(b^{(m)}, \alpha^{(m)})$ with $\hat{u}^{(m)}$ given by

$$\hat{u}^{(m)}(t) = \sum_{j=-1}^{p+1} \mathbf{b}_j^{(m)} B_j(t).$$

- (i) The iteration is initialized with a spline $\hat{u}^{(0)}$ obtained by smoothing the data. More precisely, $\hat{u}^{(0)}$ is given by equation (3.3) with $b^{(0)}$ the minimum of $E_D(b)$.
- (ii) For each $m = 1, 2, \dots$, we obtain $(\mathbf{b}^{(m)}, \alpha^{(m)})$ from $\hat{u}^{(m-1)}$ by minimizing $E(b, \alpha, \hat{u}^{(m-1)})$ with respect to (b, α) . This is a linear least squares problem which is carried out in one step using QR decomposition.
- (iii) We define $\hat{u}^{(m)}$ from $\mathbf{b}^{(m)}$ using equation (3.3). (iv) If for some preset tolerance δ we have $\|\mathbf{b}^{(m)} - \mathbf{b}^{(m-1)}\| < \delta$ then terminate, otherwise return to step (ii).

Properties of the solution

- When the iteration terminates, we have $b^{(m)} = b^{(m-1)}$ to some pre-specified numerical tolerance. Thus

$$\hat{E}_M \left(\mathbf{b}^{(m)}, \alpha, \hat{\mathbf{u}}^{(m-1)} \right) = \hat{E}_M \left(\mathbf{b}^{(m)}, \alpha, \hat{\mathbf{u}}^{(m)} \right) = E_M \left(\mathbf{b}^{(m)}, \alpha \right)$$

and hence

$$\hat{E} \left(\mathbf{b}^{(m)}, \alpha, \hat{\mathbf{u}}^{(m-1)} \right) = \tilde{E} \left(\mathbf{b}^{(m)}, \alpha \right).$$

- Since $\left(\mathbf{b}^{(m)}, \alpha^{(m)} \right)$ minimizes $\hat{E} \left(\mathbf{b}, \alpha, \hat{\mathbf{u}}^{(m-1)} \right)$ with respect to (\mathbf{b}, α) , we see that $\alpha^{(m)}$ is a minimum of $\tilde{E} \left(\mathbf{b}^{(m)}, \alpha \right)$ with respect to α . In general, it does not appear to be possible to ensure that $\tilde{E}(b, \alpha)$ is also minimized with respect to b but in practice the distinction between \hat{E}_M and E_M appears to have negligible effect.

Table of Contents

- 1 Introduction and Review
- 2 Profiled Estimation for Linear Systems Estimated by Least Squares Fitting
- 3 Profiled Estimation for Nonlinear Systems
- 4 A New Efficient Method for ODEs Linear in Parameters
- 5 Fast Inference in Nonlinear Dynamical Systems using Gradient Matching

We model the unknown concentrations of the s th component of the dynamical system at time t with a linear combination of kernels $k(.,.)$ from some function family \mathcal{F} :

$$g_s(t; \mathbf{b}_s) = \sum_{j=1}^n b_{sj} k(t, t_j)$$

We denote by \mathbf{b}_s the vector of kernel regression coefficients b_{sk} and define $\mathbf{B} = (\mathbf{b}_1^\top, \dots, \mathbf{b}_r^\top)$, where r denotes the total number of components in the system.

Objective function

We estimate \mathbf{B} along with θ by minimisation of the following objective function

$$E(\theta, \mathbf{B}) = \sum_{s=1}^r \left(\sum_{i=1}^n [g_s(t_i; \mathbf{b}_s) - y_{si}]^2 \right) \\ + \rho \sum_{s=1}^r \left(\sum_{i=1}^n [\dot{g}_s(t_i; \mathbf{b}_s) - f_s(\mathbf{g}(t_i, \mathbf{B}), \theta)]^2 \right)$$

where $\mathbf{g}(t_i, \mathbf{B}) = (g_1(t_i; \mathbf{b}_1), \dots, g_r(t_i; \mathbf{b}_r))^{\top}$, and $\rho \geq 0$ is a regularization parameter.

The gradient matching term

$$\dot{g}_s(t; b_s) = \sum_{i=1}^n b_{si} \frac{dk(t, t_i)}{dt} = \sum_{i=1}^n b_{si} \dot{k}(t, t_i)$$

We need to solve this inference problem

$$\{\hat{\theta}, \hat{B}\} = \operatorname{argmin}_{\theta, B} E(\theta, B)$$

Step 1

Initialisation of regression parameters and optimisation of kernel hyperparameters

- Regularised loss function

$$\mathcal{L}(\mathbf{b}_s, \varphi_s; \lambda_s) = \sum_{i=1}^n (g_s(t_i; \mathbf{b}_s) - y_{si})^2 + \|\mathbf{g}_s\|^2$$

where the dependence on φ_s is via k_s (which has not been made explicit in the notation), and the regularisation term $\|\mathbf{g}_s\|^2$ is the squared norm in \mathcal{H}_s , $\|\mathbf{g}_s\|^2 = \lambda_s \mathbf{b}_s^T \mathbf{K}_s \mathbf{b}_s$, which contains regularisation parameter $\lambda_s \geq 0$, and \mathbf{K}_s is the Gram matrix with entries $k_s(t_k, t_i)$.

- The minimisation of $\mathcal{L}(\mathbf{b}_s, \varphi_s; \lambda_s)$ with respect to \mathbf{b}_s for given φ_s and λ_s is a convex optimisation problem with solution

$$\mathbf{b}_s = (\mathbf{K}_s + \lambda_s \mathbf{I})^{-1} \mathbf{y}_s$$

Initialisation of ODE parameters using gradient matching

Setting B fixed at the values obtained from Step 1 , the ODE parameters θ are optimised by minimising the objective function E using a standard optimisation routine (e.g. trust region or quasi Newton).

$$E(\theta, \mathbf{B}) = \sum_{s=1}^r \left(\sum_{i=1}^n [g_s(t_i; \mathbf{b}_s) - y_{si}]^2 \right) \\ + \rho \sum_{s=1}^r \left(\sum_{i=1}^n [\dot{g}_s(t_i; \mathbf{b}_s) - f_s(\mathbf{g}(t_i, \mathbf{B}), \theta)]^2 \right)$$

Minimisation of the combined objective function with convergence acceleration

- First, we define the following modified objective function

$$\begin{aligned}\tilde{E}(\boldsymbol{\theta}, \mathbf{B}, \tilde{\mathbf{B}}) &= \sum_{s=1}^r \left(\sum_{i=1}^n [g_s(t_i; \mathbf{b}_s) - y_{si}]^2 \right) \\ &+ \rho \sum_{s=1}^r \left(\sum_{i=1}^n \left[\dot{g}_s(t_i; \mathbf{b}_s) - f_s(\mathbf{g}(t_i, \tilde{\mathbf{B}}), \boldsymbol{\theta}) \right]^2 \right)\end{aligned}$$

- Note that $\tilde{E}(\boldsymbol{\theta}, \mathbf{B}, \mathbf{B}) = E(\boldsymbol{\theta}, \mathbf{B})$.

We now carry out the following iteration until reaching a zero-gradient point:

- Given B and θ , minimize $\tilde{E}(\theta, B^*, B)$ with respect to B^* , i.e. find $B_{\text{new}} = \operatorname{argmin}_{B^*} \tilde{E}(\theta, B^*, B)$
- Set $\mathbf{B} = \mathbf{B}_{\text{new}}$ and minimise $\tilde{E}(\theta, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}})$ wrt θ , i.e. find $\theta_{\text{new}} = \operatorname{argmin}_{\theta} \tilde{E}(\theta, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}})$

Theorem 1

Let \mathcal{G} denote the set of functions defined by equation (12). Assume \mathcal{G} is contained in the solution space of the ODEs in the sense that $\forall g \in \mathcal{G} \exists \theta \in \mathbb{R} : \dot{g} = f(g, \theta)$. Then each parameter adaptation step of the algorithm described above, $(\mathbf{B}, \theta) \rightarrow (\mathbf{B}_{new}, \theta_{new})$, implies that $E(\theta_{new}, \mathbf{B}_{new}) \leq E(\theta, \mathbf{B})$, and the iteration converges to a zero gradient point of $E(\theta, \mathbf{B})$.

Proof. The first step of the algorithm implies that

$$\tilde{E}(\boldsymbol{\theta}, \mathbf{B}_{\text{new}}, \mathbf{B}) \leq \tilde{E}(\boldsymbol{\theta}, \mathbf{B}, \mathbf{B}) = E(\boldsymbol{\theta}, \mathbf{B}) \quad (18)$$

For the second step, note that $\exists \boldsymbol{\theta}^* \in \mathbb{R}$ such that

$\dot{\mathbf{g}}(t, \mathbf{B}_{\text{new}}) = f(\mathbf{g}(t, \mathbf{B}_{\text{new}}), \boldsymbol{\theta}^*)$, by assumption of the theorem. This implies that

$$\|\dot{\mathbf{g}}(t, \mathbf{B}_{\text{new}}) - f(\mathbf{g}(t, \mathbf{B}_{\text{new}}), \boldsymbol{\theta}^*)\|^2 = 0 \quad \forall t$$

$$\begin{aligned}
 \text{Hence } \tilde{E}(\boldsymbol{\theta}, \mathbf{B}_{\text{new}}, \mathbf{B}) &= \\
 &= \sum_{i=1}^N \left\{ \|\mathbf{y}(t_i) - \mathbf{g}(t_i, \mathbf{B}_{\text{new}})\|^2 + \right. \\
 &\quad \left. \rho \|\dot{\mathbf{g}}(t_i, \mathbf{B}_{\text{new}}) - f[\mathbf{g}(t_i, \mathbf{B}), \boldsymbol{\theta}]\|^2 \right\} \\
 &\geq \sum_{i=1}^N \left\{ \|\mathbf{y}(t_i) - \mathbf{g}(t_i, \mathbf{B}_{\text{new}})\|^2 \right\} + 0 \\
 &= \sum_{i=1}^N \left\{ \|\mathbf{y}(t_i) - \mathbf{g}(t_i, \mathbf{B}_{\text{new}})\|^2 + \right. \\
 &\quad \left. \rho \|\dot{\mathbf{g}}(t_i, \mathbf{B}_{\text{new}}) - f[\mathbf{g}(t_i, \mathbf{B}_{\text{new}}), \boldsymbol{\theta}^*]\|^2 \right\} \\
 &= \tilde{E}(\boldsymbol{\theta}^*, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}})
 \end{aligned}$$

This implies that on completion of the second step of the algorithm,
 $\boldsymbol{\theta}_{\text{new}} = \operatorname{argmin}_{\boldsymbol{\theta}} \tilde{E}(\boldsymbol{\theta}, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}})$

We have $\tilde{E}(\theta_{\text{new}}, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}}) = \tilde{E}(\theta^*, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}})$, and hence

$$\begin{aligned} E(\theta_{\text{new}}, \mathbf{B}_{\text{new}}) &= \tilde{E}(\theta_{\text{new}}, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}}) \\ &= \tilde{E}(\theta^*, \mathbf{B}_{\text{new}}, \mathbf{B}_{\text{new}}) \leq \tilde{E}(\theta, \mathbf{B}_{\text{new}}, \mathbf{B}) \quad (19) \end{aligned}$$

Combining equations (18) and (19), we get:

$$E(\theta_{\text{new}}, \mathbf{B}_{\text{new}}) \leq \tilde{E}(\theta, \mathbf{B}_{\text{new}}, \mathbf{B}) \leq E(\theta, \mathbf{B})$$

which completes the first part of the proof. The iteration is continued until $\nabla_{\theta} \tilde{E}(\theta, \mathbf{B}, \mathbf{B}) = \nabla_{\mathbf{B}} \tilde{E}(\theta, \mathbf{B}, \mathbf{B}) = 0$. Since $\tilde{E}(\theta, \mathbf{B}, \mathbf{B}) = E(\theta, \mathbf{B})$, this is a zero-gradient point of the original objective function $E(\theta, \mathbf{B})$.

Table 1. Comparison of computational costs. The table shows the computational costs for a single iteration of two optimisation algorithms compared in our study, and an alternative method discussed in Section 2, using the data generated from eq. (20).

Method	CPU time
Direct minimisation of $E(\boldsymbol{\theta}, \boldsymbol{B})$	599.8s
Proposed acceleration algorithm (RKG3)	6.7s
RKHS method by Gonzalez et al. (GON)	4.2s

- **James Ramsay, Giles Hooker:** "Dynamic Data Analysis - Modeling Data with Differential Equations"
- **Mu Niu, Simon Rogers, Maurizio Filippone, Dirk Husmeier:** "Fast Inference in Nonlinear Dynamical Systems using Gradient Matching"
- **J. O. Ramsay, G. Hooker, D. Campbell and J. Cao :** "Parameter estimation for differential equations: a generalized smoothing approach"
- **Daniel Brewer, Martino Barenco:** "Fitting ordinary differential equations to short time course data"

END

Thank you very much for your valuable time!