

An approach to parameter estimation and model selection in differential equations

M.R.Osborne

Centre for Mathematics and its Applications

School of Mathematical Sciences

<http://www.maths.anu.edu.au/~mike/index.html>

September 23, 2003

The estimation problem In parameterized systems of differential equations this starts with data acquired through observations of system trajectories made in the presence of noise and seeks to estimate the parameter values by solving an optimization problem which matches solutions of the differential equation to the observed data. Thus there is an explicit stochastic component to the problem. Typically, two classes of method are considered:

Class [1]: Explicitly computed solution trajectories are compared directly with the observations in an unconstrained optimization procedure; and

Class [2]: the system of differential equations is imposed as explicit constraints on the optimization problem. The resulting mathematical program typically is solved by a variant of sequential quadratic programming.

Basic data

The differential equation:

$$\frac{d\mathbf{x}}{dt} = \mathbf{w}(t, \mathbf{x}, \boldsymbol{\beta})$$

where $\mathbf{x}, \mathbf{w} \in R^m$, $\boldsymbol{\beta} \in R^p$. Important case is the linear equation

$$\mathbf{w} = M(t, \boldsymbol{\beta})\mathbf{x} + \mathbf{f}(t)$$

which provides the “enabling technology”.

The observed data:

$$y_i = \boldsymbol{\phi}^T \mathbf{x}(t_i) + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

where $\boldsymbol{\phi}$ defines ‘observation functional’, ε_i iid, (normal).

Class [1] algorithm (*embedding method*): Typically this has two main components

1. Given trial β , plus *auxiliary information* $\mathbf{b} \in R^m$, generate trial solution $\mathbf{x}(t, \beta)$.
2. Using trial solution make adjustments to β *and auxiliary information* to improve estimate of β and $\mathbf{x}(t, \beta)$. Measure goodness of fit by

$$F(\beta) = \sum_{i=1}^n r_i(t_i, \beta)^2$$
$$r_i = y_i - \phi^T \mathbf{x}(t_i, \beta).$$

Auxilliary information

Example: Explicit embedding. Idea is to

- (a) select boundary data B_1, B_2, \mathbf{b} ;
- (b) guess variable part of auxiliary information \mathbf{b} ;
- (c) Solve boundary (or initial) value problem

$$B_1 \mathbf{x}(0) + B_2 \mathbf{x}(1) = \mathbf{b},$$
$$\frac{d\mathbf{x}}{dt} = M(t, \beta) \mathbf{x} + \mathbf{f}(t).$$

Try to choose B_1, B_2 so the Green's matrix is 'nicely bounded'. 'Always possible' if dichotomy known. Typically want fast solutions pinned down at $t = 1$, and slow solutions at $t = 0$.

Simple shooting corresponds to $B_1 = I, B_2 = 0$. It requires the initial value problem to be stable.

Leads to nonlinear least squares problem:

$$\min_{\mathbf{b}, \boldsymbol{\beta}} \sum_{i=1}^n (y_i - \phi^T \mathbf{x}(t_i, \boldsymbol{\beta}, \mathbf{b}))^2.$$

Compute correction using Gauss-Newton (Scoring). Note \mathbf{b} occurs as an extra parameter vector to be estimated. Requires the integration of the variational equations:

$$\left\{ \begin{array}{l} \frac{d\Delta_{\boldsymbol{\beta}}}{dt} = M\Delta_{\boldsymbol{\beta}} + \nabla_{\boldsymbol{\beta}} M \mathbf{x}, \\ B_1 \Delta_{\boldsymbol{\beta}}(0) + B_2 \Delta_{\boldsymbol{\beta}}(1) = 0, \end{array} \right. ,$$
$$\left\{ \begin{array}{l} \frac{d\Delta_{\mathbf{b}}}{dt} = M\Delta_{\mathbf{b}}, \\ B_1 \Delta_{\mathbf{b}}(0) + B_2 \Delta_{\mathbf{b}}(1) = I. \end{array} \right. ,$$

where

$$\Delta_u = \frac{\partial \mathbf{x}}{\partial \mathbf{u}}.$$

Can explicit embedding be avoided?:

The class [2] approach (*simultaneous methods*) is formulated as

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \phi^T \mathbf{x}(t_i, \boldsymbol{\beta}))^2$$

subject to equality constraints

$$\mathbf{x}_{i+1} - X(t_{i+1}, t_i) \mathbf{x}_i = \mathbf{v}_i, \quad i = 1, 2, \dots, n-1$$

where X fundamental matrix, \mathbf{v} is the particular integral

$$\begin{aligned} \frac{dX}{dt} &= MX, \quad X(\xi, \xi) = I, \\ \mathbf{v}_i &= \int_{t_i}^{t_{i+1}} X(t_{i+1}, u) \mathbf{f}(u) du. \end{aligned}$$

In practice, the differential equation constraints would be replaced by an appropriate discretization. Here the additional information comes from the Lagrange multipliers, but the constraint system grows with n . In contrast the ODE system has only m degrees of freedom. Redundancy!

Cyclic reduction– More than just elimination!

This is an elimination scheme for the block bidiagonal recurrence

$$A_i^0 \mathbf{x}_{i+1} + B_i^0 \mathbf{x}_i = \mathbf{c}_i^0$$

which combines adjacent rows using techniques such as partial pivoting or orthogonal reduction as follows:

$$\begin{bmatrix} B_{i-1}^0 & A_{i-1}^0 & 0 & \mathbf{c}_{i-1}^0 \\ 0 & B_i^0 & A_i^0 & \mathbf{c}_i^0 \end{bmatrix} \rightarrow \begin{bmatrix} B_{i/2}^1 & 0 & A_{i/2}^1 & \mathbf{c}_{i/2}^1 \\ V_i^1 & -I & W_i^1 & \mathbf{w}_i^1 \end{bmatrix}.$$

The procedure can be applied recursively to give

Interpolation equations

$$\mathbf{x}_t = V_t \mathbf{x}(0) + W_t \mathbf{x}(1) + \mathbf{w}_t,$$

Constraint equation

$$G_1^k \mathbf{x}(0) + G_2^k \mathbf{x}(1) = \mathbf{c}_1^k.$$

Process is simplest if $n = 2^k$. Restriction not necessary in bidiagonal case.

These equations are intrinsic properties of the ODE system in the sense that they do not depend on the boundary conditions.

The reduced system: The aim of choosing B_1, B_2 in embedding the estimation problem is to ensure that

$$\begin{bmatrix} B_1 & B_2 \\ G_1^k & G_2^k \end{bmatrix}$$

has a ‘nicely’ bounded inverse. Thus G_1^k, G_2^k must reflect the dichotomy properties of the ODE system.

The cyclic reduction transformation allows the reformulation of the estimation problem:

$$\min_{\beta} \sum_{t=t_i, i=1}^n \left(y_t - \phi^T (V_t \mathbf{x}(0) + W_t \mathbf{x}(1) + \mathbf{w}_t) \right)^2$$

subject to the constraints

$$G_1^k \mathbf{x}(0) + G_2^k \mathbf{x}(1) = \mathbf{c}_1^k.$$

This reduces the Lagrangian form of the problem to solving an optimization problem involving a fixed number of equality constraints.

Properties: Boundary conditions - $V(0) = I$, $V(1) = 0$, $W(0) = 0$, $W(1) = I$, $\mathbf{w}(0) = \mathbf{w}(1) = 0$.

$V_t, W_t, \mathbf{w}_t, G_1, G_2, c$ are not uniquely defined by the cyclic reduction process. Let C be the transformation that combines adjacent block rows. Then there is an equivalence class of transformations:

$$C \leftarrow \begin{bmatrix} R_1 & 0 \\ R_{21} & R_2 \end{bmatrix} C$$

that preserve the basic structure in the elimination tableau.

Freedom in the interpolation is in $R_2^{-1}R_{21}$.

Freedom in the constraint is in R_1 .

$$\left(G_2^k\right)^{-1} G_1^k = X(1, 0), \left(G_2^k\right)^{-1} \mathbf{c}_1^k = \mathbf{v}_1.$$

Governing equations

The simplest transformation is given by

$$C = \begin{bmatrix} I & X(t_{i+1}, t_i) \\ I & -X(t_{i+1}, t_i) \end{bmatrix},$$

assume $\delta = t_{i+1} - t_i$ is small, expand in powers of δ , and equate leading terms. This gives second order system (so extra boundary condition can be satisfied):

$$\begin{aligned} \frac{d^2}{dt^2} \left(X^{-1} \begin{Bmatrix} V \\ W \end{Bmatrix} \right) &= 0, \\ \Rightarrow V &= X(t, 0)(1 - t), \quad W = X(t, 1)t. \end{aligned}$$

Find other possibilities by fixing

$$R_2^{-1} R_{21} = S_1 = \delta S + O(\delta^2).$$

Substituting this into $C \leftarrow RC$ and repeating calculation gives for V (W is similar)

$$\frac{d^2 V}{dt^2} + 2(S - M) \frac{dV}{dt} + \left(M^2 - 2SM - \frac{dM}{dt} \right) V = 0.$$

Orthogonal reduction

Partial pivoting can be unstable for block bidiagonal systems so use orthogonal transformation. This requires

$$C^T R^T R C = I$$

Substituting and expanding in powers of δ gives

$$S = \frac{M + M^T}{2}$$

Substituting in the general equation gives (order important)

$$\left(\frac{d}{dt} + M^T \right) \left(\frac{d}{dt} - M \right) \begin{cases} V \\ W \end{cases} = 0$$

The general equation corresponds to the first order system (write Y for either V, W)

$$\frac{d}{dt} \begin{bmatrix} Y \\ Z \end{bmatrix} = N \begin{bmatrix} Y \\ Z \end{bmatrix}, \quad N = \begin{bmatrix} M & I \\ & -(2S - M) \end{bmatrix}.$$

The constraint equation: need particular integral equation in form

$$\frac{d}{dt} \begin{bmatrix} \mathbf{w} \\ \mathbf{z} \end{bmatrix} = N \begin{bmatrix} \mathbf{w} \\ \mathbf{z} \end{bmatrix} + \begin{bmatrix} \mathbf{f} \\ 0 \end{bmatrix}.$$

The construction gives \mathbf{x} as a combination of solutions of a higher order equation. Function of constraint is to remove unwanted terms

$$\begin{aligned} 0 &= \frac{d\mathbf{x}}{dt} - M\mathbf{x} - \mathbf{f} \\ &= \left(\frac{dV}{dt} - MV \right) \mathbf{x}(t_1) + \left(\frac{dW}{dt} - MW \right) \mathbf{x}(t_n) \\ &\quad + \frac{d\mathbf{w}}{dt} - M\mathbf{w} - \mathbf{f} \\ &= Z_V(t)\mathbf{x}(t_1) + Z_W(t)\mathbf{x}(t_n) + \mathbf{z}(t). \end{aligned}$$

There is really only one condition here!

Optimization methodology When the model is known then the Gauss-Newton or scoring method appears the method of choice for the first class of methods, and there are good reasons for this which are a consequence of the stochastic setting. Similar approximations appear to work well in SQP methods for the second class of methods, and Zengfeng Li's experience here will be summarized.

Scoring - Two main ideas:

[1] *Maximum likelihood for parameter estimation.* This starts with:

events: $\mathbf{y}_t \in R^m, t \in T$

probability density: $f(\mathbf{y}_t, \boldsymbol{\eta}_t(\mathbf{x}), t)$

exact model: $\boldsymbol{\eta}_t(\mathbf{x}) : R^p \times T \rightarrow R^q$

(parameter and covariate information).

This computes parameter estimate

$$\mathbf{x}_T = \arg \min \mathcal{K}_T(\mathbf{x});$$

$$\mathcal{K}_T(\mathbf{x}) = - \sum_{t \in T} \mathcal{L}_t, \quad \mathcal{L}_t = \log f(\mathbf{y}_t, \boldsymbol{\eta}_t(\mathbf{x}), t)$$

Context: (not just sum of squares)

[1] independent events;

[2] $n = |T| \gg m = \dim \mathbf{y}_t$; and

[3] right kind of analytic properties.

Signal in noise model relevant here:

$$\nabla_{\boldsymbol{\eta}} \mathcal{L}_t = V_t(\boldsymbol{\beta})^{-1} \{\mathbf{y}_t - \boldsymbol{\eta}_t(\mathbf{x}_t, \boldsymbol{\beta})\}$$

[2] *Newton's method for function minimization*

Compute:

$$\mathcal{J} = \nabla^2 \mathcal{K}(\mathbf{x}),$$

$$\mathbf{h} = -\mathcal{J}^{-1}(\mathbf{x}) \nabla \mathcal{K}(\mathbf{x})^T.$$

Update:

$$\mathbf{x} \rightarrow \mathbf{x} + \mathbf{h}.$$

Advantages:

- (1) Fast rate of ultimate convergence to $\hat{\mathbf{x}} \ni \nabla \mathcal{K}(\hat{\mathbf{x}}) = 0$ provided $\mathcal{J}(\hat{\mathbf{x}})$ is nonsingular.
- (2) Good transformation invariance properties.

Disadvantages:

- (1) Convergence is local.
- (2) Method requires $\nabla^2 \mathcal{K}(\mathbf{x})$. It is often regarded as uneconomical or inconvenient to compute this.

Improving global behavior : Introduce merit function $Q(\mathbf{x})$ with properties that

(1) It has same minimizer(s) as \mathcal{K} ;

(2) $\nabla Q(\mathbf{x})\mathbf{h} < 0$ whenever \mathbf{h} , $\nabla Q(\mathbf{x}) \neq 0$.

Aim is to reduce Q at each step $\mathbf{x} \rightarrow \mathbf{x} + \lambda\mathbf{h}$ where typically λ might be chosen to satisfy:

$$\begin{aligned} \varrho &< \psi(\mathbf{x}, \mathbf{h}, \lambda) < 1 - \varrho, \\ \psi &= \frac{Q(\mathbf{x} + \lambda\mathbf{h}) - Q(\mathbf{x})}{\lambda \nabla Q(\mathbf{x})\mathbf{h}}, \\ 0 &< \varrho < .5. \end{aligned}$$

Can always find such a λ , σ . Latter usually small.

Harder to find suitable Q .

- $Q = \mathcal{K}$ will do if $\nabla^2\mathcal{K}$ positive definite.
- $Q = \|\nabla\mathcal{K}\|^2$ works but ...!
- Big advantage of scoring is that suitable Q is available

Define the *sample information matrix* by

$$\mathcal{I}_n = \frac{1}{n} \mathcal{E}_f \{ \nabla_x^2 \mathcal{K}_n \} = \frac{1}{n} \sum_t \nabla_x \eta_t^T V_t^{-1}(\mathbf{x}) \nabla_x \eta_t$$

and estimate the Newton correction by

$$\mathbf{h} = -\mathcal{I}_n(\mathbf{x})^{-1} \frac{1}{n} \nabla_x \mathcal{K}_n(\mathbf{x})^T$$

Properties:

- (1) Scoring has the same good transformation properties as Newton's method.
- (2) It requires only first derivative information.
- (3) \mathcal{I}_n is positive (semi) definite so $\nabla_x \mathcal{K}_n \mathbf{h} < (=) 0$.

This last property ensures that the scoring step is necessarily downhill for minimizing \mathcal{K}_n when \mathcal{I}_n is nonsingular, and that $Q = \mathcal{K}_n$ works as a monitor. This has the consequences:

- (1) $\frac{\nabla_x \mathcal{K} \mathbf{h}}{\|\nabla_x \mathcal{K}\| \|\mathbf{h}\|} < -\frac{1}{\text{cond} \mathcal{I}}$,
- (2) Limit points of the iteration are stationary points of \mathcal{K} .
- (3) $\lambda = 1$ will satisfy the ψ -test eventually if n large enough.

Approximating the expectation

In the expected Hessian

$$V_t(\mathbf{x})^{-1} = \mathcal{E} \left\{ \nabla_{\eta} \mathcal{L}_t^T \nabla_{\eta} \mathcal{L}_t \right\}.$$

If the expectation has to be estimated then the law of large numbers can help

$$\begin{aligned} \frac{1}{n} \mathcal{E} \{ \nabla_x^2 \mathcal{K}_n \} &= \frac{1}{n} \sum_i \mathcal{E} \{ \nabla_x \mathcal{L}_i^T \nabla_x \mathcal{L}_i \} \\ &= -\frac{1}{n} \sum_i \left(\nabla_x \mathcal{L}_i^T \nabla_x \mathcal{L}_i - \mathcal{E} \{ \nabla_x \mathcal{L}_i^T \nabla_x \mathcal{L}_i \} \right) \\ &\quad + \frac{1}{n} \sum_i \nabla_x \mathcal{L}_i^T \nabla_x \mathcal{L}_i \\ &\rightarrow \frac{1}{n} \sum_i \nabla_x \mathcal{L}_i^T \nabla_x \mathcal{L}_i \end{aligned}$$

Rate of convergence Can write as a fixed point iteration when $\lambda = 1$

$$\mathbf{x}_{i+1} = F(\mathbf{x}_i); \quad F(\mathbf{x}) = \mathbf{x} - \mathcal{I}_n(\mathbf{x})^{-1} \frac{1}{n} \nabla_x \mathcal{K}_n(\mathbf{x})^T$$

$\hat{\mathbf{x}}_n$ point of attraction provided

$$\varpi(F'(\hat{\mathbf{x}}_n)) < 1$$

As $\nabla_x \mathcal{K}_n(\hat{\mathbf{x}}_n) = 0$ have

$$\begin{aligned} F'(\hat{\mathbf{x}}_n) &= I - \mathcal{I}_n(\hat{\mathbf{x}}_n)^{-1} \frac{1}{n} \nabla_x^2 \mathcal{K}_n(\hat{\mathbf{x}}_n) \\ &= (\mathcal{I}_n(\hat{\mathbf{x}}_n))^{-1} \left((\mathcal{I}_n(\hat{\mathbf{x}}_n) - \frac{1}{n} \nabla_x^2 \mathcal{K}_n(\hat{\mathbf{x}}_n)) \right) \\ &= F'(\mathbf{x}^*) + \mathcal{O}(\|\hat{\mathbf{x}}_n - \mathbf{x}^*\|), \text{ a.s., } n \rightarrow \infty \end{aligned}$$

But $F'(\mathbf{x}^*) = o(1)$, $n \rightarrow \infty$ using strong law

$$\Rightarrow \varpi(F'(\hat{\mathbf{x}}_n)) \rightarrow 0, \quad n \rightarrow \infty$$

Arbitrary fast rate of (first order) convergence provided effective sample size is large enough.

Scoring with linear constraints The problem

$$\min_{\mathbf{x}} \mathcal{K}_n; C\mathbf{x} = \mathbf{d}, C : R^p \rightarrow R^m, \text{rank}(C) = m.$$

The necessary conditions for a minimum give

$$\nabla_x \mathcal{K}_n = \boldsymbol{\lambda}^T C$$

where $\boldsymbol{\lambda}$ is the vector of Lagrange multipliers.

Limiting form as $n \rightarrow \infty$ follows from

$$\frac{1}{n} \{ \nabla_x \mathcal{K}_n - \mathcal{E}_* \{ \nabla_x \mathcal{K}_n \} \} + \frac{1}{n} \mathcal{E}_* \{ \nabla_x \mathcal{K}_n \} = (\boldsymbol{\lambda}/n)^T C$$

The left hand side has the limiting form

$$- \int_0^1 \mathcal{E}_* \{ \nabla_x \mathcal{L}(\mathbf{y}, \mathbf{x}, t) \} d\omega(t)$$

Thus the limiting system is

$$\begin{aligned} - \int_0^1 \mathcal{E}_* \{ \nabla_x \mathcal{L}(\mathbf{y}, \mathbf{x}, t) \} d\omega(t) &= \boldsymbol{\lambda}^{*T} C \\ C\mathbf{x} &= \mathbf{d} \end{aligned}$$

where $\boldsymbol{\lambda}^* = \lim_{n \rightarrow \infty} \boldsymbol{\lambda}/n$. Solution is

$$\mathbf{x} = \mathbf{x}^*, \boldsymbol{\lambda}^* = 0.$$

SQP framework: Problem

$$\min_{\mathbf{x}} \mathcal{K}(\mathbf{x}); \mathbf{c}(\mathbf{x}) = 0.$$

Introduce Lagrangian

$$l(\mathbf{x}, \boldsymbol{\lambda}) = \mathcal{K}(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{c}(\mathbf{x}).$$

Let B_k be an approximation to $\nabla_{\mathbf{x}}^2 l(\mathbf{x}_k, \boldsymbol{\lambda}_k)$ and solve linear subproblem

$$\begin{aligned} \min_{\mathbf{d} \in S} \nabla \mathcal{K}(\mathbf{x}) \mathbf{d} + \frac{1}{2} \mathbf{d}^T B_k \mathbf{d}, \\ S = \{\mathbf{d}; \mathbf{c}(\mathbf{x}_k) + A(\mathbf{x}_k) \mathbf{d} = 0\} \end{aligned}$$

(use cyclic reduction in constraint reduction step)

Take guarded step $\mathbf{x}_{k+1} := \mathbf{x}_k + \gamma \mathbf{d}_k$

(implemented Byrd and Omojokun trust region strategy)

Update Lagrange multiplier $\boldsymbol{\lambda}$

$$\boldsymbol{\lambda}_{k+1} = -A_k^+ (\nabla \mathcal{K}_k^T + B_k \mathbf{d}_k)$$

Use Gauss-Newton approximation to B_k (possibly guarded):

(i) ignore term $\sum_{i=1}^n r_i \frac{\partial^2 r_i}{\partial x_j \partial x_k}$
(justified by the standard argument)

(ii) ignore term $\sum_{i=1}^{n-1} \lambda_i^T \frac{\partial^2 \mathbf{c}_i}{\partial x_j \partial x_k}$
Numerical experience is that G-N approximation works.

Not too difficult to show suitably scaled $\lambda_i \rightarrow 0$
(not quite enough).

Limiting stochastic differential equation

$$d\lambda = -\nabla_{\mathbf{x}} w(t, \mathbf{x}, \beta)^T \lambda dt + \sigma \phi d\omega.$$

Example

$$M(t, \beta) = \begin{bmatrix} 1 - \beta_1 \cos(\beta_2 t) & 0 & 1 + \beta_1 \sin(\beta_2 t) \\ 0 & \beta_1 & 0 \\ 1 + \beta_1 \sin(\beta_2 t) & 0 & 1 + \beta_1 \cos(\beta_2 t) \end{bmatrix}$$

$$f(t) = e^t \begin{bmatrix} -1 + 19(\cos(2t) - \sin(2t)) \\ -18 \\ 1 - 19(\cos(2t) + \sin(2t)) \end{bmatrix}$$

$$x(t) = e^t \mathbf{e}$$

Data has form $x(t) + \sigma \text{rnd}$, $\sigma = 5., 1., .01$
 initial parameter vector 20% larger than true -
 [19, 2].

n	Ne	GN	Ne	GN	Ne	GN
$2^5 + 1$	15	55	6	11	4	4
$2^7 + 1$	16	20	6	10	3	4
$2^{10} + 1$	7	13	4	5	3	3

Model selection It all becomes harder if the only information available is that the model is known to lie within a parameterized class of systems. Presumably one should start the searching with the simpler members of this class (the potentially under-specified systems). However, the scoring method requires a consistency result and thus loses its justification in this case. A stochastic embedding procedure which produces spline-like objects which offer the possibility of overcoming this difficulty is being studied for systems linear in the state variables. Key components include the use of a multiple shooting version of the Kalman filter for problems with nontrivial dichotomy, and the use of information criteria for discriminating between contending models.

Splines as parametric models: Smoothing spline $\eta(t)$ defined by:

$$\min_{\eta} \sum_{i=1}^n (y_i - \eta(t_i))^2 + \tau \int_0^1 \left(\frac{d^k \eta}{dt^k} \right)^2 dt$$

τ provides a compromise between data fit and smoothness.

Stochastic formulation (Wahba)

$$\begin{aligned} \eta(t) &= \mathcal{E} \{y(t) | y_1, y_2, \dots, y_n, \lambda\} \\ \frac{d^k \eta}{dt^k} &= \sigma \sqrt{\lambda} \frac{d\omega}{dt} \end{aligned}$$

Here $\lambda = 1/\tau$. Plus consistency result $\eta(t) \rightarrow \mathcal{E}\{y(t)\}$, $n \rightarrow \infty$ provided λ chosen appropriately.

Generalisation to more general differential operators (g-splines)

$$\min_{\eta} \sum_{i=1}^n (y_i - \eta(t_i))^2 + \tau \int_0^1 (\mathcal{M}_k \eta)^2 dt$$

As τ gets large η forced to null space of \mathcal{M}_k . Possibility of identifying linear model for signal $\eta(t)$ in this case.

First order systems (Wecker, Ansley, Kohn)

Write \mathcal{M}_k in first order system form

$$\frac{d\mathbf{x}}{dt} = M_k \mathbf{x}$$

Stochastic form (here $\mathbf{b} = \mathbf{e}_k$)

$$\begin{aligned} \eta(t) &= \mathcal{E} \{x_1(t) | y_1, y_2, \dots, y_n, \lambda\}, \\ dx &= M_k \mathbf{x} dt + \sigma \sqrt{\lambda} \mathbf{b} d\omega. \end{aligned}$$

Generalise - $y_i = \phi^T \mathbf{x}(t_i)$, $M(t, \beta)$ matrix of general linear system, \mathbf{b} involved in smoothness control. Let $X(t, \xi)$ be fundamental matrices of deterministic equation. Then obtain relations

$$\begin{aligned} \mathbf{x}_{i+1} &= X(t_{i+1}, t_i) \mathbf{x}_i + \sigma \sqrt{\lambda} \mathbf{u}_i, \\ \mathbf{u}_i &= \int_{t_i}^{t_{i+1}} X(t_{i+1}, s) \mathbf{b} d\omega(s), \\ &\sim N(0, \sigma^2 R(t_{i+1}, t_i)), \\ R(t_{i+1}, t_i) &= \lambda \int_{t_i}^{t_{i+1}} X(t_{i+1}, s) \mathbf{b} \mathbf{b}^T X(t_{i+1}, s)^T ds. \end{aligned}$$

Kalman filter: System is in form required for computing $\mathbf{x}(t|n)$ using Kalman filter and interpolation smoother. The filter is a forward recursion giving $\mathbf{x}_{i|i} = \mathcal{E} \{ \mathbf{x}(t_i) | y_1, y_2, \dots, y_i, \lambda \}$, and $\sigma^2 S_{i|i}$, the corresponding covariance. The interpolation smoother gives the dependence on all the data. If $t_i \leq t \leq t_{i+1}$:

$$\begin{aligned} \mathbf{x}(t|n) &= X(t, t_i) \mathbf{x}_{i|i} + A(t, t_i) \left(\mathbf{x}_{i+1|n} - \mathbf{x}_{i+1|i} \right), \\ A(t, t_i) &= \left\{ X(t, t_i) S_{i|i} X_i + \Gamma(t, t_i) \right\} S_{i+1|i}^{-1}, \\ \Gamma(t, t_i) &= R(t, t_i) X(t_{i+1}, t)^T. \end{aligned}$$

Two choices of the initial condition $\mathbf{x}_{1|0}$ - either to choose it as constant or assume diffuse prior ($\mathbf{x}_{1|0} = 0, S_{1|0} \uparrow \infty$). Both have been considered for smoothing splines. The filter is an initial value process, but does involve a correction step. Stability is a legitimate question.

Smoothness - choice of ϕ , \mathbf{b} .

Differentiating the interpolation smoother gives

$$\frac{d\mathbf{x}(t|n)}{dt} = M\mathbf{x}(t|n) + \mathbf{b}\mathbf{b}^T X(t_{i+1}, t)^T \mathbf{v},$$

$$\mathbf{v} = S_{i+1|i}^{-1}(\mathbf{x}_{i+1|n} - \mathbf{x}_{i+1|i})$$

Non-smoothness only at the points t_i . The interesting term is that involving $\mathbf{b}\mathbf{b}^T X(t_{i+1}, t)^T$. Calculations need derivatives of $X(t_i, t)$:

$$\frac{d^j X(t_i, t)}{dt^j} = X(t_i, t) P_j(M)$$

$$P_0 = I, P_j = \frac{dP_{j-1}}{dt} - MP_{j-1}, j = 1, 2, \dots$$

Smoothness of $\frac{d^j \mathbf{x}(t|n)}{dt^j}$ at t_i requires

$$\phi^T P_{j-1}(M)\mathbf{b} = 0, j = 1, 2, \dots$$

Limit to smoothness if $P_{j-1}(M)^T \phi$ linearly independent $\Rightarrow j < k$. Spline results follow by setting

$$\phi = \mathbf{e}_1, \mathbf{b} = \mathbf{e}_k,$$

$$M = \begin{bmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \cdots & \cdots & \\ & & & 0 & 1 \\ -m_k & -m_{k-1} & \cdots & \cdots & -m_1 \end{bmatrix}$$

Smoothness - implications for $R(t_{i+1}, t_i)$.

If δ small then Taylor expansion gives

$$R(t + \delta, t) = \int_t^{t+\delta} \sum_{i,j} \frac{(s - (t + \delta))^{i+j}}{i!j!} W_{ij} ds$$

$$W_{ij} = P_i(M) \mathbf{b} \mathbf{b}^T P_j(M)^T.$$

As $\delta \rightarrow 0$ can use Rayleigh quotient for eigenvalue information.

(1) Largest eigenvalue:

$$\pi_k = \lambda \delta \mathbf{b}^T \mathbf{b} + O(\delta^2).$$

Corresponding eigenvector $\rightarrow \mathbf{b}$.

(2) If orthogonality conditions $\mathbf{b}^T P_{j-1}(M)^T \phi = 0$ satisfied then eigenvector with smallest eigenvalue $\rightarrow \phi$. Corresponding RQ is

$$\pi_1 = \frac{\lambda}{((k-1)!)^2} \frac{(\mathbf{b}^T P_{k-1}(M)^T \phi)^2}{\phi^T \phi} \frac{\delta^{2k-1}}{2k-1} + O(\delta^{2k}).$$

This is an upper bound!

Parameter estimation: Two main approaches - cross validation and GML. The latter involves a “likelihood” approach. Starting point is that innovations $\zeta_i = y_i - \phi^T \mathbf{x}_{i|i-1}$ are independent, normally distributed with variance $\sigma^2 \mathcal{V}_i$ where $\mathcal{V}_i = (1 + \phi^T S_{i|i-1} \phi)$. Idea is to minimize

$$\sum'_i \left\{ \log \sigma^2 + \log \mathcal{V}_i + \frac{\zeta_i^2}{\sigma^2 \mathcal{V}_i} \right\}$$

Minimizing with respect to σ^2 gives

$$\hat{\sigma}^2 = \frac{1}{N} \sum'_i \frac{\zeta_i^2}{\mathcal{V}_i}.$$

Substituting back gives concentrated likelihood

$$GML = \sum'_i \log \mathcal{V}_i + N \log \left(\sum'_i \frac{\zeta_i^2}{\mathcal{V}_i} \right)$$

Generalised cross validation This gives an objective function

$$GCV = \frac{\sum_{i=1}^n (y_i - \phi^T \mathbf{x}_{i|n})^2 / n}{\{\text{trace}\{I - T\}/n\}^2}.$$

Here T is the influence matrix mapping observations y_i into the estimated signal $\phi^T \mathbf{x}_{i|n}$.

Advantages are claimed for its use in estimating λ . Problem is finding an implementation in less than $O(n^2)$ cost that can be used for parameter estimation.

GML readily easy to calculate in $O(n)$ cost. It appears relatively insensitive to λ .