

Inference and Characterization of Planar Trajectories

Zhanglong Cao

a thesis submitted for the degree of
Doctor of Philosophy
at the University of Otago, Dunedin,
New Zealand.

30 August 2017

Contents

0.1	Introduction	1
0.2	Gaussian Process Regression	2
0.3	A Reproducing Kernel in Space \mathbb{H}	4
0.4	Covariance Matrix and Posterior Mean	6
0.5	A 1-D Gaussian Process Spline Construction	7
0.5.1	Tractor Spline	7
0.5.2	Tractor Spline Estimated by GP	11
0.6	Cross Validation	13
0.6.1	K-Fold Cross Validation	15
1.7	Writing Something	17
1.8	Introduction	17
1.9	Tractor Spline	20
1.9.1	Objective Function	20
1.9.2	Basis Functions	20
1.9.3	Solution to The Objective Function	22
1.9.4	Adjusted Penalty Term and Parameter Function	24
1.10	Parameter Selection and Cross Validation	25
1.11	Simulation	27
1.11.1	Numerical Examples	27
1.11.2	Evaluation	33
1.12	Conclusion and Discussion	33
1.13	Introduction	36
1.13.1	Spline	36
1.13.2	Gaussian Process Regression	37
1.13.3	The Smoothing Spline as Bayes Estimates	38
1.14	A reproducing kernel on $\mathcal{C}_{p.w.}^2[0, 1]$	39
1.15	Computation of Polynomial Smoothing Splines	40
1.16	Polynomial Smoothing Splines as Bayes Estimates	41
1.17	Numeric Simulation of Smoothing Spline and GPR	43
3.18	State Space Models	45
3.19	Sequential Monte Carlo Method	46
3.19.1	Filtering Problem and Estimation	47
3.19.2	Sampling Methods	48
3.20	Bayesian Parameter Estimation	52
3.20.1	Off-line Methods	53
3.20.2	On-line Methods	53

3.21	A Sequential Monte Carlo Algorithm for Parameter Estimation	55
3.22	States Estimation	56
4.23	State Space Models	57
4.24	Sequential Monte Carlo Method	58
4.24.1	Filtering Problem and Estimation	59
4.24.2	Sampling Methods	60
4.25	Bayesian Parameter Estimation	65
4.25.1	Off-line Methods	65
4.25.2	On-line Methods	66
4.26	Combined State and Parameters Estimation of Sequential Monte Carlo Algorithm	67
4.26.1	General Linear Space	68
4.26.2	High Dimension Parameters Space of OU-Process	73
4.27	Prior Distribution for Variance Parameters	81
4.27.1	Priors Discussion	82
4.27.2	Discussion two	83
References		85

List of Tables

- 1.1 MSE. Mean square errors of different methods. The star sign (*) marks the smallest error among these methods under the same level. The difference is not significant. 35
- 1.2 TMSE. True mean square errors of different methods. The star sign (*) marks the smallest error among these methods under the same level. The proposed tractor spline returns the smallest TMSE among all the methods under the same level except for *Doppler* with SNR=7. The differences are significant. 36

List of Figures

1	The two basis functions N_{2k+1} and N_{2k+2} on interval $[t_k, t_{k+2}]$. It is apparently that these basis functions are continuous on this interval and have continuous first and second derivatives.	9
1.2	The two basis functions N_{2k+1} and N_{2k+2} on interval $[t_k, t_{k+2}]$. It is apparently that these basis functions are continuous on this interval and have continuous first derivatives.	22
1.3	Numerical example: <i>Blocks</i> . (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.	29
1.4	Numerical example: <i>Bumps</i> . (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.	30
1.5	Numerical example: <i>HeaviSine</i> . (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.	31
1.6	Numerical example: <i>Doppler</i> . (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.	32

1.7	Estimated penalty functions. Left side shows how the value of $\lambda(t)$ changes on the interval. Right side projects $\lambda(t)$ into reconstructions. The bigger the blacks dots present, the larger the penalty values are.	34
1.8	Estimated velocity functions by taking the first derivative of tractor spline. (a) Fitted <i>Blocks</i> . (b) Fitted <i>Bumps</i> . (c) Fitted <i>HeaviSine</i> . (d) Fitted <i>Doppler</i>	35
1.9	(a) Comparing two methods under the same parameters $\lambda = 0.01$ and $\gamma = 0.1$. In this graph, the blue line is reconstruction from tractor spline, the red line is the mean of Gaussian Process, which is the posterior $\mathbb{E}(\eta(x) \mathbf{Y}, \mathbf{V})$. (b) The differences between two methods under the same parameters.	44

All knowledge is, in the final analysis, history.
All science are, in the abstract, mathematics.
All judgments are, in their rationale, statistics.
– C. Radhakrishna Rao.

0.1 Introduction

In regression problem, linear regression, linear discriminant analysis, logistic regression and separating hyperplanes all rely on a linear model. With the good property of linear model, easy to be interpreted and first order Taylor approximation to $f(t)$, it is more convenient to represent $f(t)$ by linear model. However, the true function $f(t)$ is unlikely to be an actual linear function in space \mathbb{R} . Researchers found some methods for moving beyond linearity. One of them is replacing the vector of inputs \mathbf{T} with its transformations as new variables, and then use linear models in this new space of derived input features.

Denote by $h_m(t) : \mathbb{R} \mapsto \mathbb{R}$ the m th transformation of t , $m = 1, \dots, M$. We then model

$$f(t) = \sum_{m=1}^M \beta_m h_m(t). \quad (1)$$

a linear basis expansion of \mathbf{t} in \mathbb{R} , where $h_m(t)$ are named basis functions, β_m are coefficients. Once the basis functions h_m have been determined, the models are linear in these new variables, and the fitting proceeds as before.

Suppose we are given observed data t_1, t_2, \dots, t_n on interval $[0, 1]$, satisfying $0 \leq t_1 < t_2 < \dots < t_n \leq 1$. A piecewise polynomial function $f(t)$ can be obtained by dividing the interval into contiguous intervals $(t_1, t_2), \dots, (t_{n-1}, t_n)$, and representing f by a separate polynomial in each interval. The points t_i are called knots. For example,

$$f(t) = d_i(t - t_i)^3 + c_i(t - t_i)^2 + b_i(t - t_i) + a_i, \quad (2)$$

for given coefficients d_i, c_i, b_i and a_i , where $t_i \leq t \leq t_{i+1}$, $i = 1, 2, \dots, n$. f is a cubic spline on $[0, 1]$ if (1) on each intervals f is a polynomial; (2) the polynomial pieces fit together at knots t_i in such a way that f itself and its first and second derivatives are continuous at each t_i . If the second and third derivatives of f are zero at 0 and 1,

f is said to be a natural cubic spline. These conditions are called natural boundary conditions.

Over all spline functions $f(t)$ with two continuous derivatives fitting these observed data, the curve estimate $\hat{f}(t)$ will be defined to be the minimizer the following penalized residual sum of squares, [\[edited expressions\]](#)

$$\text{MSE}(f, \lambda) = \frac{1}{n} \sum_{i=1}^n (f(t_i) - y_i)^2 + \lambda \int_0^1 (f''(t))^2 dt \quad (3)$$

where λ is a fixed smoothing parameter, (t_i, y_i) , $i = 1, \dots, n$ are observed data and $0 \leq t_1 < t_2 < \dots < t_n \leq 1$. In equation (1.123), the smoothing parameter λ controls the trade-off between over-fitting and bias,

$$\begin{cases} \lambda = 0 : & f \text{ can be any function that interpolates the data,} \\ \lambda = \infty : & \text{the simple least squares line fit since no second derivative can be tolerated.} \end{cases} \quad (4)$$

In our case, the velocity data set v_i with some independent Gaussian distributed errors $\varepsilon_i \sim N(0, \frac{\sigma_n^2}{\gamma})$ are used to estimate $f(t)$ simultaneously. f is a linear combination of basis functions, as shown in equation (1.122), in the meantime, f' is a linear combination of the first derivative of these basis functions

$$f'(t) = \sum_{m=1}^M \alpha_m h'_m(t). \quad (5)$$

The velocity information is incorporated into MSE equation (1.123) by the addition of velocity term $(f'(t_i) - v_i)^2$. Then it becomes

$$\text{MSE}(f, \lambda, \gamma) = \frac{1}{n} \sum_{i=1}^n (f(t_i) - y_i)^2 + \frac{\gamma}{n} \sum_{i=1}^n (f'(t_i) - v_i)^2 + \lambda \int_0^1 (f''(t))^2 dt, \quad (6)$$

and \hat{f} is the minimizer of the MSE equation (6).

In the model $y = f(t) + \varepsilon$, it is reasonable to assume that the observed data y_i is Gaussian distribution with mean $f(t_i)$ and variance σ_n^2 . In a similar way, the velocity is estimated as $v = f'(t) + \frac{\varepsilon}{\gamma}$, where v_i is Gaussian distribution with mean $f'(t_i)$ and variance $\frac{\sigma_n^2}{\gamma}$. Then the joint distribution of $\mathbf{y}, \mathbf{v}, f(t)$ and $f'(t)$ is normal with zero mean and a covariance matrix, which can be estimated through Gaussian Process Regression.

0.2 Gaussian Process Regression

A Gaussian Process is a collection of random variables, any finite number of which have a joint Gaussian distribution, Rasmussen and Williams (2006).

A GP is fully defined by its mean $m(t)$ and covariance $K(s, t)$ functions as

$$m(t) = \mathbb{E}[f(t)] \quad (7)$$

$$K(s, t) = \mathbb{E}[(f(s) - m(s))(f(t) - m(t))], \quad (8)$$

where s and t are two variables, and a function f distributed as such is denoted in form of

$$f \sim GP(m(t), K(s, t)). \quad (9)$$

Usually the mean function is assumed to be zero everywhere.

Given a set of input variables \mathbf{T} for function $f(t)$ and the output $\mathbf{y} = f(\mathbf{T}) + \varepsilon$ with independent identically distributed Gaussian noise ε with variance σ_n^2 , we can use the above definition to predict the value of the function $f_* = f(t_*)$ at a particular input t_* . As the noisy observations becoming

$$\text{cov}(y_p, y_q) = K(t_p, t_q) + \sigma_n^2 \delta_{pq} \quad (10)$$

where δ_{pq} is a Kronecker delta which is one iff $p = q$ and zero otherwise, the joint distribution of the observed outputs \mathbf{y} and the estimated output f_* according to prior is

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim N \left(0, \begin{bmatrix} K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I & K(\mathbf{T}, t_*) \\ K(t_*, \mathbf{T}) & K(t_*, t_*) \end{bmatrix} \right). \quad (11)$$

The posterior distribution over the predicted value is obtained by conditioning on the observed data

$$f_* | \mathbf{y}, \mathbf{T}, t_* \sim N(\bar{f}_*, \text{cov}(f_*)) \quad (12)$$

where

$$\bar{f}_* = \mathbb{E}[f_* | \mathbf{y}, \mathbf{T}, t_*] = K(t_*, \mathbf{T})[K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I]^{-1} \mathbf{y}, \quad (13)$$

$$\text{cov}(f_*) = K(t_*, t_*) - K(t_*, \mathbf{T})[K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I]^{-1} K(\mathbf{T}, t_*). \quad (14)$$

We now add velocity information $\mathbf{v} = f'(\mathbf{T}) + \varepsilon'$, where ε' is independent distributed Gaussian noise with variance $\frac{\sigma_n^2}{\gamma}$.

It is expected that a position point y_i and velocity point v_i are all effected by other points \mathbf{y} and \mathbf{v} . So the covariance matrix for \mathbf{y} and \mathbf{v} is

$$\Sigma(\mathbf{y}, \mathbf{v}) = \begin{bmatrix} \text{cov}(\mathbf{y}, \mathbf{y}) & \text{cov}(\mathbf{y}, \mathbf{v}) \\ \text{cov}(\mathbf{v}, \mathbf{y}) & \text{cov}(\mathbf{v}, \mathbf{v}) \end{bmatrix}, \quad (15)$$

where obviously $\text{cov}(\mathbf{y}, \mathbf{v}) = \text{cov}(\mathbf{v}, \mathbf{y})$. Then the joint distribution is

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{v} \end{bmatrix} \sim N(\mu_{y,v}, \Sigma_{y,v}). \quad (16)$$

Define f_* and f'_* the estimated position and velocity values at point t_* . From equation (15) and using similar idea, it is easily to get the covariance matrices

$$\begin{aligned} \Sigma(f_*, \mathbf{v}) &= \begin{bmatrix} \text{cov}(f_*, f_*) & \text{cov}(f_*, \mathbf{v}) \\ \text{cov}(\mathbf{v}, f_*) & \text{cov}(\mathbf{v}, \mathbf{v}) \end{bmatrix}, \\ \Sigma(\mathbf{y}, f'_*) &= \begin{bmatrix} \text{cov}(\mathbf{y}, \mathbf{y}) & \text{cov}(\mathbf{y}, f'_*) \\ \text{cov}(f'_*, \mathbf{y}) & \text{cov}(f'_*, f'_*) \end{bmatrix}, \\ \Sigma(f_*, f'_*) &= \begin{bmatrix} \text{cov}(f_*, f_*) & \text{cov}(f_*, f'_*) \\ \text{cov}(f'_*, f_*) & \text{cov}(f'_*, f'_*) \end{bmatrix}, \end{aligned} \quad (17)$$

[will need to give the form of these covariances at some point. in an appendix? i think you need discussion of how f' is related to f for a GP]

0.3 A Reproducing Kernel in Space \mathbb{H}

$N_1(t), \dots, N_n(t)$ denote n basis function having first derivative in space \mathbb{H} . For any continuous function $f \in \mathbb{H}$, it is a combination of these basis functions

$$f(t) = \sum_{i=1}^n \alpha_i N_i(t), \quad (18)$$

where $\alpha_i (i = 1, \dots, n)$ are coefficients. With an inner product

$$\langle f, g \rangle = \left\langle \sum_{i=1}^n \alpha_i N_i(t), \sum_{i=1}^n \beta_i N_i(t) \right\rangle = \sum_{i=1}^n \alpha_i \beta_i, \quad (19)$$

it can be shown that the representer of evaluation $[s](\cdot)$ is

$$R_s(t) = \sum_{i=1}^n N_i(s) N_i(t), \quad (20)$$

Then we can prove that the space \mathbb{H} is a Reproducing Kernel Hilbert Space. In fact,

$$\langle f(t), R(s, t) \rangle = \left\langle \sum_{i=1}^n \alpha_i N_i(t), \sum_{i=1}^n N_i(s) N_i(t) \right\rangle = \sum_{i=1}^n \alpha_i N_i(s) = f(s). \quad (21)$$

The term $R(s, t) = R_s(t)$ is called the reproducing kernel function.

We now introduce a new notation $\dot{R}(s, t)$ in the following and use it to find the covariance matrix Σ of the joint distribution of $\mathbf{y}, \mathbf{v}, f$ and f' .

Define $\dot{R}(s, t)$ and $R'(s, t)$ are the first partial derivative of $R(s, t)$ with respect to the first and second argument respectively

$$\dot{R}(s, t) = \frac{\partial R(s, t)}{\partial s} = \sum_{i=1}^n \frac{dN_i(s)}{ds} N_i(t) = \sum_{i=1}^n N'_i(s) N_i(t), \quad (22)$$

$$R'(s, t) = \frac{\partial R(s, t)}{\partial t} = \sum_{i=1}^n N_i(s) \frac{dN_i(t)}{dt} = \sum_{i=1}^n N_i(s) N'_i(t), \quad (23)$$

Then $\dot{R}'(s, t)$ is the second partial derivative of $R(s, t)$ with respect to both arguments

$$\dot{R}'(s, t) = \frac{\partial^2 R(s, t)}{\partial s \partial t} = \sum_{i=1}^n \frac{dN_i(s)}{ds} \frac{dN_i(t)}{dt} = \sum_{i=1}^n N'_i(s) N'_i(t). \quad (24)$$

It is easy to prove that $\dot{R}(s, t) = R'(t, s)$ and

$$\begin{aligned} \langle f(t), R'(s, t) \rangle &= \langle \sum_i \alpha_i N_i(t), \sum_i N_i(s) N'_i(t) \rangle = \sum_i \alpha_i N_i(t) = f(t), \\ \langle f(t), \dot{R}(s, t) \rangle &= \langle \sum_i \alpha_i N_i(t), \sum_i N'_i(s) N_i(t) \rangle = \sum_i \alpha_i N'_i(s) = f'(s), \end{aligned} \quad (25)$$

Given the sample points $t_i, i = 1, \dots, n$ and noting that the space

$$\mathbb{A} = \{f : f = \sum_{i=1}^n \alpha_i R(t_i, \cdot)\} \quad (26)$$

is one linear subspace of \mathbb{H} . Then $f \in \mathbb{H}$ can be written as

$$f(t) = \sum_{i=1}^n c_i R(t_i, t) + \rho(t) \quad (27)$$

where c_i are coefficients, $\rho(t) \in \mathbb{H} \ominus \mathbb{A}$, and

$$f'(t) = \sum_{i=1}^n c_i R'(t_i, t) + \rho'(t). \quad (28)$$

The equation (6) can be written inform of

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{j=1}^n c_j R(t_j, t_i) - \rho(t_i))^2 &+ \frac{\gamma}{n} \sum_{i=1}^n (v_i - \sum_{j=1}^n c_j R'(t_j, t_i) - \rho'(t_i))^2 \\ &+ \lambda \int_0^1 (\sum_{j=1}^n c_j R''(t_j, t) + \rho''(t))^2 dt \end{aligned} \quad (29)$$

As $\dot{R}(t_i, \cdot) = \sum_{j=1}^n N_j'(t_i) N_j(t) \in \mathbb{A}$, then by orthogonality and property of reproducing kernel functions, $\rho(t_i) = \langle R(t_i, \cdot), \rho \rangle = 0$, and $\rho'(t_i) = \langle \rho, \dot{R}(t_i, \cdot) \rangle = 0$, where $i = 1, \dots, n$.

Denoting by Q the $n \times n$ matrix with the (i, j) th entry $R(t_i, t_j)$, by P the $n \times n$ matrix with the (i, j) th entry $\dot{R}(t_i, t_j)$ the equation (6) can be written as

$$(\mathbf{y} - Q\mathbf{c})^\top (\mathbf{y} - Q\mathbf{c}) + \gamma(\mathbf{v} - P\mathbf{c})^\top (\mathbf{v} - P\mathbf{c}) + n\lambda\Omega + \lambda(\rho, \rho). \quad (30)$$

Note that ρ only appears in the third term in (1.145), which is minimized at $\rho = 0$. Hence, a polynomial smoothing spline resides in the space \mathbb{A} of finite dimension. Then, following the method given from Gu (2013), the solution can be computed via minimization of term in (30) with respect to \mathbf{c}

0.4 Covariance Matrix and Posterior Mean

Consider f and f' in \mathbb{H} , having Gaussian priors with zero mean. By equation (1.127), their covariance functions are

$$\begin{aligned} \text{cov}(f(s), f(t)) &= \tau^2 R(s, t) + \sigma_n^2 I \\ \text{cov}(f(s), f'(t)) &= \tau^2 R'(s, t) + \frac{\sigma_n^2}{\sqrt{\gamma}} I \\ \text{cov}(f'(s), f(t)) &= \tau^2 \dot{R}(s, t) + \frac{\sigma_n^2}{\sqrt{\gamma}} I \\ \text{cov}(f'(s), f'(t)) &= \tau^2 \dot{R}'(s, t) + \frac{\sigma_n^2}{\gamma} I \end{aligned} \quad (31)$$

Observing $y_i \sim N(f(t_i), \sigma_n^2)$ and $v_i \sim N(f'(t_i), \frac{\sigma_n^2}{\gamma})$, the joint distribution of $\mathbf{y}, \mathbf{v}, f(t)$ and $f'(t)$ is normal with zero mean and covariance matrix

$$\begin{aligned} \text{cov}(\mathbf{y}, \mathbf{v}, f, f') &= \begin{bmatrix} \tau^2 R(t_i, t_j) + \sigma_n^2 I & \tau^2 R'(t_i, t_j) + \frac{\sigma_n^2}{\sqrt{\gamma}} I & \tau^2 R(t_i, t) & \tau^2 R'(t_i, t) \\ \tau^2 \dot{R}(t_i, t_j) + \frac{\sigma_n^2}{\sqrt{\gamma}} I & \tau^2 \dot{R}'(t_i, t_j) + \frac{\sigma_n^2}{\gamma} I & \tau^2 \dot{R}(t_i, t) & \tau^2 \dot{R}'(t_i, t) \\ \tau^2 R^\top(t_i, t) & \tau^2 \dot{R}^\top(t_i, t) & \tau^2 R(t, t) & \tau^2 R'(t, t) \\ \tau^2 R'^\top(t_i, t) & \tau^2 \dot{R}'^\top(t_i, t) & \tau^2 \dot{R}(t, t) & \tau^2 \dot{R}'(t, t) \end{bmatrix} \\ &= \begin{bmatrix} \tau^2 Q + \sigma_n^2 I & \tau^2 O + \frac{\sigma_n^2}{\sqrt{\gamma}} I & \tau^2 \xi & \tau^2 \xi' \\ \tau^2 O + \frac{\sigma_n^2}{\sqrt{\gamma}} I & \tau^2 P + \frac{\sigma_n^2}{\gamma} I & \tau^2 \dot{\xi} & \tau^2 \dot{\xi}' \\ \tau^2 \xi^\top & \tau^2 \dot{\xi}^\top & \tau^2 R(t, t) & \tau^2 R'(t, t) \\ \tau^2 \xi'^\top & \tau^2 \dot{\xi}'^\top & \tau^2 \dot{R}(t, t) & \tau^2 \dot{R}'(t, t) \end{bmatrix} \end{aligned} \quad (32)$$

where $\{Q\}_{ij}$ is the matrix with elements $R(t_i, t_j)$, $\{O\}_{ij}$ is the matrix with elements $\dot{R}(t_i, t_j) = R'(t_j, t_i)$, $\{P\}_{ij}$ is the matrix with elements $\dot{R}'(t_i, t_j)$, ξ is a $n \times 1$ matrix with i th elements $R(x_i, x)$, and $\dot{\xi}$ is a $n \times 1$ matrix with i th elements $\dot{R}(x_i, x)$. Then

$$\begin{aligned} E \begin{bmatrix} f|\mathbf{y} \\ f'|\mathbf{v} \end{bmatrix} &= \begin{bmatrix} \xi^\top & \dot{\xi}^\top \\ \xi'^\top & \dot{\xi}'^\top \end{bmatrix} \begin{bmatrix} Q + n\lambda I & O + \frac{n\lambda}{\sqrt{\gamma}} I \\ O + \frac{n\lambda}{\sqrt{\gamma}} I & P + \frac{n\lambda}{\gamma} I \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ \gamma\mathbf{v} \end{bmatrix} \\ &\triangleq \begin{bmatrix} \xi^\top & \dot{\xi}^\top \\ \xi'^\top & \dot{\xi}'^\top \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \gamma\mathbf{v} \end{bmatrix} \\ &= \begin{bmatrix} \xi^\top (A\mathbf{y} + B\gamma\mathbf{v}) + \dot{\xi}^\top (C\mathbf{y} + D\gamma\mathbf{v}) \\ \xi'^\top (A\mathbf{y} + B\gamma\mathbf{v}) + \dot{\xi}'^\top (C\mathbf{y} + D\gamma\mathbf{v}) \end{bmatrix} \end{aligned} \quad (33)$$

where $n\lambda = \sigma_n^2/\tau^2$. The posterior mean $E(f|\mathbf{y}, \mathbf{v})$ is a linear combination of basis functions $N_i(t)$, and both ξ and $\dot{\xi}$ contain $N_i(t)$, thus the posterior mean is of the form $\xi^\top \mathbf{c} + \dot{\xi}^\top \mathbf{d}$. Similarly, $E(f'|\mathbf{y}, \mathbf{v})$ is of the form $\xi'^\top \mathbf{c} + \dot{\xi}'^\top \mathbf{d}$, with the same coefficients given by

$$\mathbf{c} = A\mathbf{y} + B\gamma\mathbf{v} \quad (34)$$

$$\mathbf{d} = C\mathbf{y} + D\gamma\mathbf{v} \quad (35)$$

0.5 A 1-D Gaussian Process Spline Construction

Trajectories are represented by a series of 2D position points (x_t, y_t) and velocity points (u_t, v_t) corresponding to measurements taken at discrete time steps t , where x_t and u_t represented longitude, y_t and v_t represented latitude position and velocity respectively Ellis *et al.* (2009). For now, we just focus on the problem of fitting trajectories in 1 Dimension situation.

For any $t \in [t_1, t_n]$, we wish to estimate the latitude position $y(t)$ and velocity $v(t)$ with model

$$y(t) = f(t) + \varepsilon, \quad (36)$$

$$v(t) = f'(t) + \frac{\varepsilon}{\gamma}, \quad (37)$$

where ε is zero-mean Gaussian noise. A Gaussian process prior over $f \sim GP(m(t), K(s, t))$ leading to the approximate estimation model

$$p(y_t, v_t|\mathbf{y}, \mathbf{v}) \sim N(GP_\mu(\mathbf{y}, \mathbf{v}), GP_\Sigma(\mathbf{y}, \mathbf{v})). \quad (38)$$

0.5.1 Tractor Spline

Suppose we have observed dataset $t_1 < t_2 < \dots < t_n$. The function $f(t)$ defined on this interval $[t_1, t_n]$ is called tractor spline, if on each interval (t_i, t_{i+1}) , $i = 2, \dots, n-2$, $f(t)$ is a cubic polynomial, but on interval (t_1, t_2) and (t_{n-1}, t_n) can be a linear function; $f(t)$ fits together at each point t_i in such a way that $f(t)$ itself and its first and second derivatives are continuous at each t_i , $i = 2, \dots, n-2$.

On an arbitrary interval $[t_i, t_{i+1}]$, we have Hermite Spline basis functions as following

$$h_{00}^{(i)}(t) = \begin{cases} 2(\frac{t-t_i}{t_{i+1}-t_i})^3 - 3(\frac{t-t_i}{t_{i+1}-t_i})^2 + 1, & t_i \leq t \leq t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (39)$$

$$h_{10}^{(i)}(t) = \begin{cases} \frac{(t-t_i)^3}{(t_{i+1}-t_i)^2} - 2\frac{(t-t_i)^2}{t_{i+1}-t_i} + (t-t_i), & t_i \leq t \leq t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (40)$$

$$h_{01}^{(i)}(t) = \begin{cases} -2(\frac{t-t_i}{t_{i+1}-t_i})^3 + 3(\frac{t-t_i}{t_{i+1}-t_i})^2, & t_i \leq t \leq t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (41)$$

$$h_{11}^{(i)}(t) = \begin{cases} \frac{(t-t_i)^3}{(t_{i+1}-t_i)^2} - \frac{(t-t_i)^2}{t_{i+1}-t_i}, & t_i \leq t \leq t_{i+1} \\ 0 & \text{otherwise} \end{cases}. \quad (42)$$

Construct new basis functions on entire interval $[t_1, t_n]$ in such way, that $N_1 = h_{00}^{(1)}$, $N_2 = h_{10}^{(1)}$, $N_{2n-1} = h_{01}^{(n)}$, $N_{2n} = h_{11}^{(n)}$. For all $k = 1, 2, \dots, n-2$ define N_{2k+1} by

$$N_{2k+1}(t) = \begin{cases} h_{01}^{(k)} + h_{00}^{(k+1)} & t \neq t_{k+1} \\ 1 & t = t_{k+1}. \end{cases}$$

and $N_{2k+2} = h_{11}^{(k)} + h_{10}^{(k+1)}$. Then $N_1(t), \dots, N_{2n}(t)$ are the new basis functions on $[t_1, t_n]$.

We now prove that N_1, N_2, \dots, N_{2n} are linear independent.

Lemma 1. *Peng (1983) Functions $x_1(t), x_2(t), \dots, x_n(t)$ on interval $[a, b]$, if they are linear dependent, the necessary and sufficient condition is for any $c_1, c_2, \dots, c_n \in [a, b]$, the determinant $D(c_1, c_2, \dots, c_n) = 0$; if they are linear independent, the necessary and sufficient condition is that there exist $c_1, c_2, \dots, c_n \in [a, b]$, so that the determinant $D(c_1, c_2, \dots, c_n) \neq 0$, where*

$$D(c_1, c_2, \dots, c_n) = \begin{vmatrix} x_1(c_1) & x_1(c_2) & \dots & x_1(c_n) \\ x_2(c_1) & x_2(c_2) & \dots & x_2(c_n) \\ \vdots & \vdots & \ddots & \vdots \\ x_n(c_1) & x_n(c_2) & \dots & x_n(c_n) \end{vmatrix} \quad (43)$$

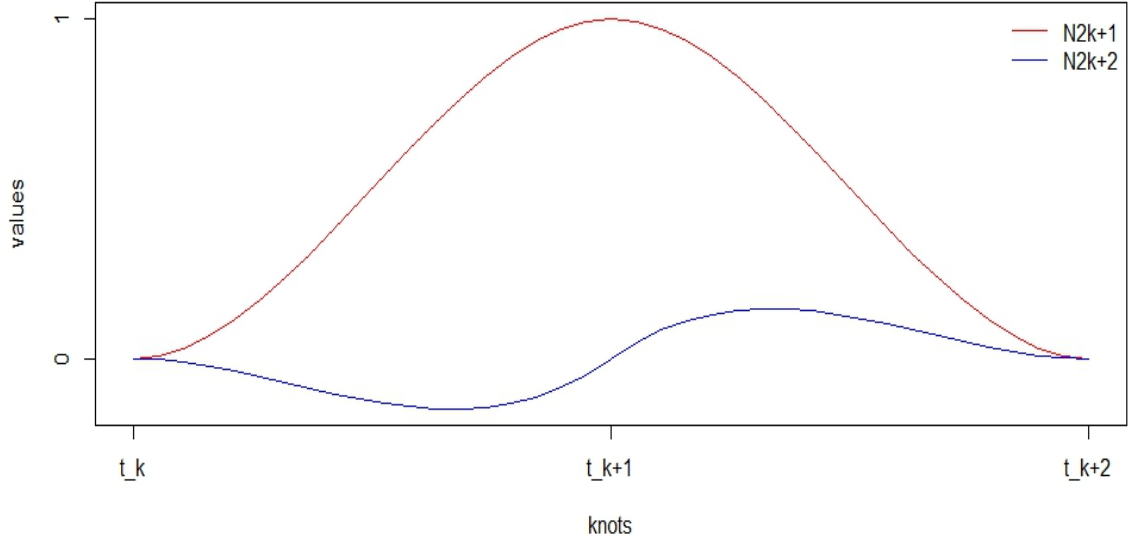


Figure 1: The two basis functions N_{2k+1} and N_{2k+2} on interval $[t_k, t_{k+2}]$. It is apparently that these basis functions are continuous on this interval and have continuous first and second derivatives.

Theorem 1. *The functions N_1, \dots, N_{2n} provide a basis for the set of functions on $[t_1, t_n]$ which are continuous, have continuous first derivatives and which are cubic on each open interval (t_i, t_{i+1}) .*

The proof of theorem 4 is in appendices.

As independent basis functions, $N_1(t), \dots, N_{2n}(t)$ span a $2n$ dimensional space \mathbb{H} . For any $f \in \mathbb{H}$, it is represented in the form of

$$f = \sum_{i=1}^{2n} \theta_i N_i(t). \quad (44)$$

Suppose that we have observations y_1, \dots, y_n and v_1, \dots, v_n . $f(t)$ can be found by minimizing equation (6), which reduces to

$$\text{MSE}(\theta, \lambda, \gamma) = (\mathbf{y} - \mathbf{B}\theta)^\top (\mathbf{y} - \mathbf{B}\theta) + \gamma(\mathbf{v} - \mathbf{C}\theta)^\top (\mathbf{v} - \mathbf{C}\theta) + n\lambda\theta^\top \Omega \theta \quad (45)$$

where $\{\mathbf{B}\}_{ij} = N_j(t_i)$, $\{\mathbf{C}\}_{ij} = N'_j(t_i)$ and $\{\Omega_{2n}\}_{jk} = \int N''_j(t)N''_k(t)dt$. After substituting the series observation t_1, \dots, t_n into basis functions, we get $N_1(t_1) = 1, N_1(t_2) = 0, \dots, N_{2k-1}(t_k) = 1, N_{2k}(t_k) = 0, \dots, N_{2n-1}(t_n) = 1, N_{2n}(t_n) = 0$; and into first derivative of basis functions, we get $N'_1(t_1) = 0, N'_1(t_2) = 1, \dots, N'_{2k-1}(t_k) = 0, N'_{2k}(t_k) =$

$1, \dots, N'_{2n-1}(t_n) = 0, N'_{2n}(t_n) = 1$. That means the matrices \mathbf{B} and \mathbf{C} in MSE equation (1.92) are $n \times 2n$ dimensional and the elements are

$$\mathbf{B} = \{B\}_{ij} = \begin{cases} 1, & j = 2i - 1 \\ 0, & \text{otherwise} \end{cases} \quad (46)$$

$$\mathbf{C} = \{C\}_{ij} = \begin{cases} 1, & j = 2i \\ 0, & \text{otherwise} \end{cases} \quad (47)$$

where $i = 1, \dots, n$. Elements of penalty matrix $\{\Omega_{2n}\}_{jk}$ is given in appendices.

The solution to (1.92) is easily seen to be

$$\hat{\theta} = (\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\lambda\Omega)^{-1}(\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \quad (48)$$

a generalized ridge regression. Then the fitted smoothing spline is given by

$$\hat{f}(t) = \sum_{i=1}^{2n} N_i(t) \hat{\theta}_i \quad (49)$$

A smoothing spline with parameters λ and γ is an example of a linear smoother Trevor Hastie (2009). This is because the estimated parameters in (1.96) are a linear combination of y_i and v_i . Denote by $\hat{\mathbf{f}}$ the $2n$ vector of fitted values $\hat{f}(t_i)$ and $\hat{\mathbf{f}}'$ the $2n$ vector of fitted values $\hat{f}'(t_i)$ at the training points t_i . Then

$$\begin{aligned} \hat{\mathbf{f}} &= \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\lambda\Omega)^{-1}(\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \\ &\triangleq \mathbf{S}_{\lambda,\gamma} \mathbf{y} + \gamma \mathbf{T}_{\lambda,\gamma} \mathbf{v} \end{aligned} \quad (50)$$

$$\begin{aligned} \hat{\mathbf{f}}' &= \mathbf{C}(\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\lambda\Omega)^{-1}(\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \\ &\triangleq \mathbf{U}_{\lambda,\gamma} \mathbf{y} + \gamma \mathbf{V}_{\lambda,\gamma} \mathbf{v} \end{aligned} \quad (51)$$

The fitted $\hat{\mathbf{f}}$ and $\hat{\mathbf{f}}'$ are linear in \mathbf{y} and \mathbf{v} , and the finite linear operators $\mathbf{S}_{\lambda,\gamma}$, $\mathbf{T}_{\lambda,\gamma}$, $\mathbf{U}_{\lambda,\gamma}$ and $\mathbf{V}_{\lambda,\gamma}$ are known as the smoother matrices. One consequence of this linearity is that the recipe for producing $\hat{\mathbf{f}}$ and $\hat{\mathbf{f}}'$ from \mathbf{y} and \mathbf{v} , do not depend on \mathbf{y} and \mathbf{v} themselves; $\mathbf{S}_{\lambda,\gamma}$, $\mathbf{T}_{\lambda,\gamma}$, $\mathbf{U}_{\lambda,\gamma}$ and $\mathbf{V}_{\lambda,\gamma}$ depend only on the t_i , λ and γ .

Suppose in a traditional least squares fitting, \mathbf{B}_ξ is $N \times M$ matrix of M cubic-spline basis functions evaluated at the N training points x_i , with knot sequence ξ and $M \ll N$. Then the vector of fitted spline values is given by

$$\hat{\mathbf{f}} = \mathbf{B}_\xi (\mathbf{B}_\xi^\top \mathbf{B}_\xi)^{-1} \mathbf{B}_\xi \mathbf{y} = \mathbf{H}_\xi \mathbf{y} \quad (52)$$

Here the linear operator \mathbf{H}_ξ is a symmetric, positive semidefinite matrices, and $\mathbf{H}_\xi \mathbf{H}_\xi = \mathbf{H}_\xi$ (idempotent). In our case, it is easily seen that $\mathbf{S}_{\lambda,\gamma}$, $\mathbf{T}_{\lambda,\gamma}$, $\mathbf{U}_{\lambda,\gamma}$ and $\mathbf{V}_{\lambda,\gamma}$ are symmetric, positive semidefinite matrices as well. However, only when $\lambda = \gamma = 0$, the matrix $\mathbf{S}_{\lambda=0,\gamma=0}$ is idempotent.

0.5.2 Tractor Spline Estimated by GP

A tractor spline on interval $[t_1, t_n]$ has $2n$ basis functions $N_1(t), \dots, N_{2n}(t)$, which are linear independent. So the space \mathbb{H} , spanned by these basis functions, is a $2n$ dimensional space. Following the definition in section 0.3, it can be proved that the space \mathbb{H} is a Reproducing Kernel Hilbert Space with inner product given in (19), and kernel function $R(s, t)$ defined in (20).

Noticing the definition of Hermite Spline from equation (1.82) to (42), prior status of y_i and v_i will only affect the status y_{i+1} and v_{i+1} in the following time period t_i . Define a covariance matrix Λ_i as

$$\Lambda_i = \text{cov}\left(\begin{bmatrix} y_i \\ v_i \end{bmatrix}, \begin{bmatrix} y_{i+1} \\ v_{i+1} \end{bmatrix}\right) \quad (53)$$

Then f and f' in \mathbb{H} with zero mean Gaussian priors, have covariance functions

$$\begin{aligned} \text{cov}(f(s), f(t)) &= \tau^2 R(s, t) + \sigma_n^2 \Lambda \\ \text{cov}(f(s), f'(t)) &= \tau^2 R'(s, t) + \frac{\sigma_n^2}{\sqrt{\gamma}} \Lambda \\ \text{cov}(f'(s), f(t)) &= \tau^2 \dot{R}(s, t) + \frac{\sigma_n^2}{\sqrt{\gamma}} \Lambda \\ \text{cov}(f'(s), f'(t)) &= \tau^2 \dot{R}'(s, t) + \frac{\sigma_n^2}{\gamma} \Lambda \end{aligned} \quad (54)$$

Observing $y_i \sim N(f(t_i), \sigma_n^2)$ and $v_i \sim N(f'(t_i), \frac{\sigma_n^2}{\gamma})$, the joint distribution of $\mathbf{y}, \mathbf{v}, f(t)$ and $f'(t)$ is normal with zero mean and covariance matrix

$$\text{cov}(\mathbf{y}, \mathbf{v}, f, f') = \begin{bmatrix} \tau^2 Q + \sigma_n^2 \Lambda & \tau^2 O + \frac{\sigma_n^2}{\sqrt{\gamma}} \Lambda & \tau^2 \xi & \tau^2 \xi' \\ \tau^2 O + \frac{\sigma_n^2}{\sqrt{\gamma}} \Lambda & \tau^2 P + \frac{\sigma_n^2}{\gamma} \Lambda & \tau^2 \dot{\xi} & \tau^2 \dot{\xi}' \\ \tau^2 \xi^\top & \tau^2 \dot{\xi}^\top & \tau^2 R(t, t) & \tau^2 R'(t, t) \\ \tau^2 \xi'^\top & \tau^2 \dot{\xi}'^\top & \tau^2 \dot{R}(t, t) & \tau^2 \dot{R}'(t, t) \end{bmatrix} \quad (55)$$

where $\{Q\}_{ij}$, $\{O\}_{ij}$, $\{P\}_{ij}$, ξ and $\dot{\xi}$ are the same as that in (32). Then

$$\begin{aligned} E \begin{bmatrix} f \\ f' \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{v} \end{bmatrix} &= \begin{bmatrix} \xi^\top & \dot{\xi}^\top \\ \xi'^\top & \dot{\xi}'^\top \end{bmatrix} \begin{bmatrix} Q + n\lambda\Lambda & O + \frac{n\lambda}{\sqrt{\gamma}}\Lambda \\ O + \frac{n\lambda}{\sqrt{\gamma}}\Lambda & P + \frac{n\lambda}{\gamma}\Lambda \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ \gamma\mathbf{v} \end{bmatrix} \\ &\triangleq \begin{bmatrix} \xi^\top & \dot{\xi}^\top \\ \xi'^\top & \dot{\xi}'^\top \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \gamma\mathbf{v} \end{bmatrix} \\ &= \begin{bmatrix} \xi^\top (A\mathbf{y} + B\gamma\mathbf{v}) + \dot{\xi}^\top (C\mathbf{y} + D\gamma\mathbf{v}) \\ \xi'^\top (A\mathbf{y} + B\gamma\mathbf{v}) + \dot{\xi}'^\top (C\mathbf{y} + D\gamma\mathbf{v}) \end{bmatrix} \end{aligned} \quad (56)$$

where $n\lambda = \sigma_n^2/\tau^2$. The posterior mean is of the form $\xi^\top \mathbf{c} + \dot{\xi}^\top \mathbf{d}$, and $E(f'|\mathbf{y}, \mathbf{v})$ is of the form $\xi'^\top \mathbf{c} + \dot{\xi}'^\top \mathbf{d}$, with the same coefficients given by

$$\mathbf{c} = A\mathbf{y} + B\gamma\mathbf{v} \quad (57)$$

$$\mathbf{d} = C\mathbf{y} + D\gamma\mathbf{v} \quad (58)$$

Following the procedure in section 0.3 and 0.4, we use observations $t_i, i = 1, \dots, n$ to construct a subspace $\mathbb{A} \subset \mathbb{H}$, which is a linear combination of kernel functions $R(s, t)$, as

$$\mathbb{A} = \{f : f = \sum_{i=1}^n \alpha_i R(t_i, \cdot)\}. \quad (59)$$

The covariance matrix and posterior mean all given above, and the solution can be computed via finding \mathbf{c} and \mathbf{d} .

In fact, for a tractor spline, the space \mathbb{A} only contains terms of odd basis functions N_{2k-1} , which can be seen from matrix \mathbf{B} in (1.92). So we construct another subspace

$$\mathbb{B} = \{f : f = \sum_{i=1}^n \beta_i \dot{R}(t_i, \cdot) = \sum_{i=1}^n \alpha_i \sum_{j=1}^{2n} N'_j(t_i) N_j(\cdot)\} \quad (60)$$

where \dot{R} is defined in equation (22). This subspace contains even terms basis functions N_{2k} . So $\mathbb{A} \cap \mathbb{B} = \emptyset$.

Thus $f \in \mathbb{H}$ can be written as

$$f(t) = \sum_{i=1}^n c_i R(t_i, t) + \sum_{i=1}^n d_i \dot{R}(t_i, t) + \rho(t) \quad (61)$$

where c_i, d_i are coefficients, $\rho(t) \in \mathbb{H} \ominus (\mathbb{A} \oplus \mathbb{B})$, and

$$f'(t) = \sum_{i=1}^n c_i R'(t_i, t) + \sum_{i=1}^n d_i \dot{R}'(t_i, t) + \rho'(t). \quad (62)$$

The equation (6) can be written as

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{j=1}^n c_j R(t_j, t_i) - \sum_{j=1}^n d_j \dot{R}(t_j, t_i) - \rho(t_i))^2 \\ & + \frac{\gamma}{n} \sum_{i=1}^n (v_i - \sum_{j=1}^n c_j R'(t_j, t_i) - \sum_{j=1}^n d_j \dot{R}'(t_j, t_i) - \rho'(t_i))^2 \\ & + \lambda \int_{t_1}^{t_n} (\sum_{j=1}^n c_j R''(t_j, t) + \sum_{j=1}^n d_j \dot{R}''(t_j, t) + \rho''(t))^2 dt \end{aligned} \quad (63)$$

By orthogonality, $\rho(t_i) = \langle R(t_i, \cdot), \rho \rangle = 0$, and $\rho'(t_i) = \langle \dot{R}(t_i, \cdot), \rho \rangle = 0$, where $i = 1, \dots, n$.

Denoting by Q the $n \times n$ matrix with the (i, j) th entry $R(t_i, t_j)$, by P the $n \times n$ matrix with the (i, j) th entry $\dot{R}(t_i, t_j)$ the equation (1.92) can be written as

$$(\mathbf{y} - Q\mathbf{c} - P\mathbf{d})^\top (\mathbf{y} - Q\mathbf{c} - P\mathbf{d}) + \gamma \left(\mathbf{v} - \frac{\partial Q}{\partial t} \mathbf{c} - \frac{\partial P}{\partial t} \mathbf{d} \right)^\top \left(\mathbf{v} - \frac{\partial Q}{\partial t} \mathbf{c} - \frac{\partial P}{\partial t} \mathbf{d} \right) + n\lambda\Omega + \lambda(\rho, \rho). \quad (64)$$

The elements of penalty matrix Ω is in appendices. Note that ρ only appears in the third term in (1.145), which is minimized at $\rho = 0$. Hence, a polynomial smoothing spline resides in the space $\mathbb{A} \oplus \mathbb{B}$ of finite dimension. Then the solution could be computed via minimization of term in (1.145) with respect to \mathbf{c} and \mathbf{d} .

0.6 Cross Validation

The coefficients can be calculated by minimizing MSE function. While another problem is how to choose smoothing parameter. There are two different philosophical approaches to the question of choosing the smoothing parameter. The first approach is to regard the free choice of smoothing parameter as an advantageous feature of the procedure. The other is a need for an automatic method whereby the smoothing parameter values is chose by the data, Green and Silverman (1993).

Assuming that the random error has zero mean, the true regression curve f has the property that, if an observation y is taken at a point t , the value $f(t)$ is the best predictor of y in terms of returning a small value of $(y - f(t))^2$.

Now we focus on an observation y_i at point t_i as being a new observation by omitting it from the set of data, which are used to estimate \hat{f} . Denote by $\hat{f}^{(-i)}(t, \lambda)$ the estimated function from the remaining data, where λ is the smoothing parameter. Then $\hat{f}^{(-i)}(t, \lambda)$ minimizes

$$\frac{1}{n} \sum_{j \neq i} (y_j - f(t_j))^2 + \lambda \int f''^2 dt \quad (65)$$

and λ can be quantified by cross-validation score function

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}^{(-i)}(t_i, \lambda)\}^2. \quad (66)$$

The basis idea of cross-validation is to choose the value of λ that minimizes $CV(\lambda)$.

An efficient way to calculate cross validation score is given by Green and Silverman (1993). Through the equation (1.100), we know that the value of the smoothing spline

\hat{f} depend linearly on the data y_i . Define the matrix $A(\lambda)$, which is a map vector of observed values y_i to predicted values $\hat{f}(t_i)$. Then we have

$$\mathbf{f} = A(\lambda)\mathbf{y} \quad (67)$$

and the following lemma.

Lemma 2. *The cross validation score satisfies*

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}(t_i)}{1 - A_{ii}(\lambda)} \right)^2 \quad (68)$$

where \hat{f} is the spline smoother calculated from the full data set $\{(t_i, y_i)\}$ with smoothing paramter λ .

For a tractor spline and its MSE function, there are two parameters need to be estimated λ and γ . Thus the objective function becomes

$$\frac{1}{n} \sum_{j \neq i} (y_j - f(t_j))^2 + \frac{\gamma}{n} \sum_{j \neq i} (v_j - f'(t_j))^2 + \lambda \int f''^2 dt, \quad (69)$$

and the cross-validation score function is

$$CV(\lambda, \gamma) = \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}^{(-i)}(t_i, \lambda, \gamma)\}^2. \quad (70)$$

For a tractor spline, the parameter $\hat{\theta} = (B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}(B^\top \mathbf{y} + \gamma C^\top \mathbf{v})$, then

$$\begin{aligned} \hat{f} &= B\hat{\theta} = B(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}B^\top \mathbf{y} + B(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}C^\top \mathbf{v} \\ &= S\mathbf{y} + \gamma T\mathbf{v}, \end{aligned} \quad (71)$$

$$\begin{aligned} \hat{f}' &= C\hat{\theta} = C(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}B^\top \mathbf{y} + C(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}C^\top \mathbf{v} \\ &= U\mathbf{y} + \gamma V\mathbf{v}. \end{aligned} \quad (72)$$

Then

Theorem 2. *The cross validation score of a tractor spline satisfies*

$$CV(\lambda, \gamma) = \frac{1}{n} \sum_{i=1}^n \frac{\hat{f}(t_i) - y_i + \gamma \frac{T_{ii}}{1 - \gamma V_{ii}} (\hat{f}'(t_i) - v_i)}{1 - S_{ii} - \gamma \frac{T_{ii}}{1 - \gamma V_{ii}} U_{ii}} \quad (73)$$

where \hat{f} is the tractor spline smoother calculated from the full data set $\{(t_i, y_i, v_i)\}$ with smoothing paramter λ and γ .

The proof of Theorem 5 follows immediately from a lemma, and gives an expression for the deleted residuals $y_i - \hat{f}^{(-i)}(t_i)$ and $v_i - \hat{f}'^{(-i)}(t_i)$ in terms of $y_i - \hat{f}(t_i)$ and $v_i - \hat{f}'(t_i)$ respectively.

Lemma 3. For fixed λ, γ and i , denote $\mathbf{f}^{(-i)}$ by the vector with components $f_j^{(-i)} = \hat{f}^{(-i)}(t_j, \lambda, \gamma)$, $\mathbf{f}'^{(-i)}$ by the vector with components $f_j'^{(-i)} = \hat{f}'^{(-i)}(t_j, \lambda, \gamma)$, and define vectors \mathbf{y}^* and \mathbf{v}^* by

$$\begin{cases} y_j^* = y_j & j \neq i \\ y_i^* = \hat{f}^{(-i)}(t_i) & \text{otherwise} \end{cases}, \quad (74)$$

$$\begin{cases} v_j^* = v_j & j \neq i \\ v_i^* = \hat{f}'^{(-i)}(t_i) & \text{otherwise} \end{cases}. \quad (75)$$

Then

$$\hat{\mathbf{f}}^{(-i)} = S\mathbf{y}^* + \gamma T\mathbf{v}^* \quad (76)$$

$$\hat{\mathbf{f}}'^{(-i)} = U\mathbf{y}^* + \gamma V\mathbf{v}^* \quad (77)$$

0.6.1 K-Fold Cross Validation

Based on the procedure given by Wahba and Wold (1975), we follow the improved steps to calculate a K-fold cross validation.

Step 1. Remove the first data t_1 and last date t_n from the dataset.

Step 2. Divide dataset into k groups:

Group 1 : t_2, t_{2+k}, \dots

Group 2 : t_3, t_{3+k}, \dots

\vdots

Group k : t_{k+1}, t_{2k+1}, \dots

Step 3. Guess values of $\lambda_{down}, \lambda_{up}$ and γ .

Step 4. Delete the first group of data. Fit a smoothing spline to the first data, the rest groups of dataset and the last data, with $\lambda_{down}, \lambda_{up}$ and γ in step 3. Compute the sum of squared deviations of this smoothing spline from the deleted data points.

Step 5. Delete instead the second group of data. Fit a smoothing spline to the remaining data with $\lambda_{down}, \lambda_{up}$ and γ . Compute the sum of squared deviations of the spline from deleted data points.

Step 6. Repeat Step 5 for the 3rd, 4th, \dots , k th group of data.

Step 7. Add the sums of squared deviations from steps 4 to 6 and divide by k . This is the cross validation score of three parameters λ_{down} , λ_{up} and γ .

Step 8. Vary λ_{down} , λ_{up} and γ systematically and repeat steps 4-7 until CV shows a minimum.

1.7 Writing Something

I'm thinking, I should start writing from Spline method, without adaptive terms on λ . An advanced Cross Validation method is given. Then when I tried to use this method to construct trajectories (application on real data), there appears some issues. So I brought adaptive terms.

It's a batch case method. Trajectory is reconstructed from a batch of data and some issues appeared. The numeric simulations proved that this new method is better. However, some issues still exist in real data application (long-gap-curve).

Then I will introduce Gaussian Process Regression and how it could estimate Smoothing Spline. Based on some references and tractor spline, I will prove some relationships between Tractor Spline and GPR. this is a spin-off results.

These are all batch method. Even if researchers could do some online smoothing, the computation time cost a lot.

Why don't we do online estimation by using some efficiency method? Introduce filters and other methods (Probably Kalman filter, particle filter and other filter) and dynamic linear regression models. Then the proposed method, parameters estimation and something else. It's an online case.

Parameter estimation methods: Metropolis-Hastings methods, better than MLE because the latter doesn't give the distribution of the estimates.

Moreover, I need to figure out how these different methods working and the advantages and disadvantages they have.

1.8 Introduction

In a vehicular system, the simplest way of getting the trajectory of a moving-object is connecting position points by a sequence of lines (line-based trajectory representation) in a 3-dimensional or 4-dimensional space-time Agarwal *et al.* (2003). Due to the noises generated from observation units, one can use regression method to find the best fitting returning the smallest sum square errors among all the sequences. Consider a regression model $y_i = f(t_i) + \epsilon_i$, where $a \leq t_1 < \dots < t_n \leq b$ and $f \in C^2[a, b]$ is an unknown smooth function, $(\epsilon_i)_{i=1}^n \sim N(0, \sigma^2)$ are random errors. In a classical parametric regression, f is assumed having the form $f(x, \beta)$, which is known up to the data estimated parameters β Kim and Gu (2004). When $f(x, \beta)$ is linear in β , we will have a standard linear model. However, most of the natural moving objects return

smooth trajectories without any angles. Therefore, a spline method is used to construct such trajectories. A curved-base method uses a parametric cubic function $P(t) = a_0 + a_1t + a_2t^2 + a_3t^3$ to obtain a spline that passes through any given sequence of joint position-velocity paired points $(x_1, v_1), (x_2, v_2), \dots, (x_n, v_n)$ Yu *et al.* (2004). In Yang and Sukkarieh (2010), an efficient and analytical continuous curvature path-smoothing algorithm based on parametric cubic Bézier curves is proposed. It can fit ordered sequential points smoothly. In computer (or computerized) numerical control (CNC), Altintas and Erkorkmaz Erkorkmaz and Altintas (2001) presented a quintic spline trajectory generation algorithm connecting a series of reference knots that produces continuous position, velocity and acceleration profiles.

However, a parametric approach only captures features contained in the preconceived class of functions Yao *et al.* (2005) and increases model bias. To avoid this, an alternative approach called nonparametric model was invented. Rather than giving a specified parameter, it is desired to reconstruct f from the data $y(t_i) \equiv y_i, i = 1, \dots, n$ Craven and Wahba (1978). Smoothing spline estimate of the f function appears as a solution to the following minimization problem: Find $\hat{f} \in C^2[a, b]$ that minimizes the penalized residual sum of squares:

$$\text{RSS} = \sum_{j=1}^n (y_j - f(t_j))^2 + \lambda \int_a^b f''(t)^2 dt \quad (1.78)$$

for pre-specified value $\lambda > 0$ Aydin and Tuzemen (2012). In equation (1.78), the first part is residual sum square and it penalizes the lack of fit. The second part is roughness penalty term weighted by a smoothing parameter λ , which varies from 0 to $+\infty$ and establishes a trade-off between interpolation and a linear model. The motivation of roughness penalty term is from a formalization of a mechanical device: if a thin piece of flexible wood, called a spline, is bent to the shape of the graph g , then the leading term in the strain energy is proportional to $\int f''^2$ Green and Silverman (1993). The cost of equation (1.78) is determined not only by its goodness-of-fit to the data quantified by the residual sum of squares, but also by its roughness Schwarz (2012). For a given λ , minimizing equation (1.78) will give the best compromise between smoothness and goodness-of-fit. Notice that the first term in equation (1.78) depends only on the values of f at knots $t_i, i = 1, \dots, n$. In the book, the authors show that the function that minimizes the roughness penalty for fixed values of $f(t_i)$ is a cubic spline: an interpolation of points via a continuous piecewise cubic function, with continuous first and second derivatives Green and Silverman (1993). The continuity requirements uniquely determine the interpolating spline, except at the boundaries Sealton *et al.* (2005).

A conventional smoothing spline is controlled by one single parameter, which controls the smoothness of a spline on the whole domain. A natural extension is to allow the smoothing parameter to vary as a penalty function of the independent variable, adapting to the change of roughness in different domains Silverman (1985), Donoho *et al.* (1995). In this way, a new objective function is formulated in the form of

$$\sum_{j=1}^n (y_j - f(x_j))^2 + \int_T \lambda(t) f''(t)^2 dt, \quad (1.79)$$

by minimizing which, the best estimation \hat{f} can be found. This approach makes adaptive smoothing as a minimization problem with a new penalty term.

Similar to conventional smoothing spline problem, researchers are wondering how to choose the penalty function $\lambda(t)$. The fundamental idea of nonparametric smoothing is to let the data choose the amount of smoothness, which consequently decides the model complexity Gu (1998). Most of them focus on data driven criteria, such as cross validation (CV), generalized cross validation (GCV) Craven and Wahba (1978) and generalized maximum likelihood (GML) Wahba (1985). A new challenge is posed that the smoothing parameter becomes a function and is varying in domains. The structure of this penalty function controls the complexity on each domain and the whole final model. Liu and Guo proposed to approximate the penalty function with an indicator and extended the generalized likelihood to the adaptive smoothing spline Liu and Guo (2010).

In this paper, we propose an adaptive smoothing spline method based on Hermite Spline basis functions to get reconstruction of f and f' from noisy data \mathbf{y} and \mathbf{v} . Rather than only using residuals of $f(t_i) - y_i$ term in the objective function in equation (1.79), we added the residuals of $f'(t_i) - v_i$ as a new term containing a new parameter γ . In this way, the spline keeps a balance on both observed \mathbf{y} and \mathbf{v} . Using new generated basis functions, we reconstruct a smoothing spline on the whole interval $[a, b]$. Derived from the new objective function, an advanced cross validation formula of $f(t)$ and $f'(t)$ is given. This method can be used in either getting true signal from noisy data or moving-object database.

1.9 Tractor Spline

1.9.1 Objective Function

In a 2D curve nonparametric regression, consider n time points $t_{1:n}$, such that $a \leq t_1, \dots, t_n \leq b$. Let $z_i = (x_i, y_i)$ and $w_i = (u_i, v_i)$ for $i = 1, \dots, n$. We define a positive piecewise constant function $\lambda(t)$:

$$\lambda(t) = \lambda_i > 0, \quad (1.80)$$

where $t_i \leq t < t_{i+1}$, $t_0 = a$, $t_{n+1} = b$, that will control the curvature penalty of each interval. For a function $f : [a, b] \rightarrow \mathbb{R}^2$ and $\gamma > 0$, define the objective function

$$J[f] = \frac{1}{n} \sum_{i=1}^n (f(t_i) - z_i)^2 + \frac{\gamma}{n} \sum_{i=1}^n (f'(t_i) - w_i)^2 + \sum_{i=0}^n \lambda_i \int_{t_i}^{t_{i+1}} f''(t)^2 dt, \quad (1.81)$$

where γ is a coefficient of the velocity information \mathbf{v} and it weights the residuals between \mathbf{f}' and \mathbf{v} , and $\lambda(t)$ is the smoothing parameter function.

Theorem 3. *For $n \geq 2$, the objective function $J[f]$ is minimized by a cubic spline that is linear outside the knots.*

The solution to the objective function (1.81) is called tractor spline.

In the following, we divide the 2D function $f(x, y)$ into two sub functions $f_x(t)$ on x -axis and $f_y(t)$ on y -axis. Compared with other parameters, choosing time t to be the parameter has some advantages: 1. The expressions of all the constraints are simpler Zhang *et al.* (2013); 2. It can be simply applied from 2-dimension to 3-dimension by adding an extra z -axis.

1.9.2 Basis Functions

Suppose we have a time series sequence of observed dataset $a = t_1 < t_2 < \dots < t_n = b$. The function $f(t)$ (stands for $f_y(t)$) defined on this interval $[t_1, t_n]$ is called tractor spline, if it is the solution to the objective function (1.81). Then it has the following property: on each interior interval (t_i, t_{i+1}) , $i = 2, \dots, n-2$, $f(t)$ is a cubic polynomial, but on interval (t_1, t_2) and (t_{n-1}, t_n) can be linear; $f(t)$ fits together at each point t_i in such a way that $f(t)$ itself and its first derivatives are continuous at each t_i , $i = 2, \dots, n-2$.

Using Hermite interpolation on an arbitrary interval $[t_i, t_{i+1}]$, the cubic spline basis functions can be constructed as follows

$$h_{00}^{(i)}(t) = \begin{cases} 2\left(\frac{t-t_i}{t_{i+1}-t_i}\right)^3 - 3\left(\frac{t-t_i}{t_{i+1}-t_i}\right)^2 + 1 & t_i \leq t < t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (1.82)$$

$$h_{10}^{(i)}(t) = \begin{cases} \frac{(t-t_i)^3}{(t_{i+1}-t_i)^2} - 2\frac{(t-t_i)^2}{t_{i+1}-t_i} + (t-t_i) & t_i \leq t < t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (1.83)$$

$$h_{01}^{(i)}(t) = \begin{cases} -2\left(\frac{t-t_i}{t_{i+1}-t_i}\right)^3 + 3\left(\frac{t-t_i}{t_{i+1}-t_i}\right)^2 & t_i \leq t < t_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad (1.84)$$

$$h_{11}^{(i)}(t) = \begin{cases} \frac{(t-t_i)^3}{(t_{i+1}-t_i)^2} - \frac{(t-t_i)^2}{t_{i+1}-t_i} & t_i \leq t < t_{i+1} \\ 0 & \text{otherwise} \end{cases}. \quad (1.85)$$

Then a Hermite spline $f^{(i)}(t)$ on interval $[t_i, t_{i+1})$ with points $p_i = \{y_i, v_i\}$ and $p_{i+1} = \{y_{i+1}, v_{i+1}\}$ can be expressed as

$$f^{(i)}(t) = h_{00}^{(i)}(t)y_i + h_{10}^{(i)}(t)v_i + h_{01}^{(i)}(t)y_{i+1} + h_{11}^{(i)}(t)v_{i+1}. \quad (1.86)$$

To construct a tractor spline on the entire interval $[t_1, t_n]$, the new basis functions are defined in such way, that $N_1 = h_{00}^{(1)}$, $N_2 = h_{10}^{(1)}$, and for all $k = 1, 2, \dots, n-2$,

$$N_{2k+1} = \begin{cases} h_{01}^{(k)} + h_{00}^{(k+1)} & \text{if } t < t_n \\ 2\left(\frac{t-t_{n-1}}{t_n-t_{n-1}}\right)^3 - 3\left(\frac{t-t_{n-1}}{t_n-t_{n-1}}\right)^2 + 1 & \text{if } t = t_n \end{cases}, \quad (1.87)$$

$$N_{2k+2} = \begin{cases} h_{11}^{(k)} + h_{10}^{(k+1)} & \text{if } t < t_n \\ \frac{(t-t_{n-1})^3}{(t_n-t_{n-1})^2} - 2\frac{(t-t_{n-1})^2}{t_n-t_{n-1}} + (t-t_{n-1}) & \text{if } t = t_n \end{cases}, \quad (1.88)$$

and

$$N_{2n-1} = \begin{cases} h_{01}^{(n-1)} & \text{if } t < t_n \\ -2\left(\frac{t-t_{n-1}}{t_n-t_{n-1}}\right)^3 + 3\left(\frac{t-t_{n-1}}{t_n-t_{n-1}}\right)^2 & \text{if } t = t_n \end{cases}, \quad (1.89)$$

$$N_{2n} = \begin{cases} h_{11}^{(n-1)} & \text{if } t < t_n \\ \frac{(t-t_{n-1})^3}{(t_n-t_{n-1})^2} - \frac{(t-t_{n-1})^2}{t_n-t_{n-1}} & \text{if } t = t_n \end{cases}. \quad (1.90)$$

Theorem 4. On $[t_1, t_n]$, the functions N_1, \dots, N_{2n} provide a basis for the set of functions which are continuous, have continuous first derivatives and are cubic on each open interval (t_i, t_{i+1}) , where $i = 1, \dots, n-1$.

As independent basis functions, $N_1(t), \dots, N_{2n}(t)$ span a $2n$ dimensional function space \mathbb{H} . For any $f \in \mathbb{H}$, it can be represented in the form of

$$f = \sum_{k=1}^{2n} \theta_k N_k(t), \quad (1.91)$$

where $\{\theta_k\}_{k=1}^{2n}$ are parameters.

Figure (1.2) presents two basis functions on an arbitrary interval $[t_k, t_{k+2})$ where they are continuous and differential. At the interior joint knot t_k , basis functions in the previous and following interval share the same position y_k and velocity v_k .

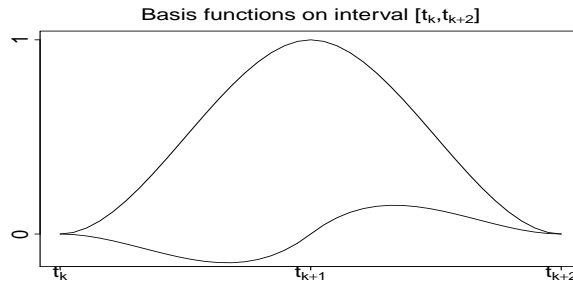


Figure 1.2: The two basis functions N_{2k+1} and N_{2k+2} on interval $[t_k, t_{k+2})$. It is apparently that these basis functions are continuous on this interval and have continuous first derivatives.

1.9.3 Solution to The Objective Function

Basis functions have been defined in the previous subsection, therefore the tractor spline $f(t)$ on $[a, b]$, where $a \leq t_1 < t_2 < \dots < t_{n-1} < t_n \leq b$, can be found by minimizing objective function (1.81), which reduces to

$$\text{MSE}(\theta, \lambda, \gamma) = (\mathbf{y} - \mathbf{B}\theta)^\top (\mathbf{y} - \mathbf{B}\theta) + \gamma(\mathbf{v} - \mathbf{C}\theta)^\top (\mathbf{v} - \mathbf{C}\theta) + n\theta^\top \Omega_\lambda \theta, \quad (1.92)$$

where $\{\mathbf{B}\}_{ij} = N_j(t_i)$, $\{\mathbf{C}\}_{ij} = N'_j(t_i)$ and $\{\Omega_{2n}^{(k)}\}_{jk} = \int_{t_k}^{t_{k+1}} \lambda_k N_j''(t) N_k''(t) dt$. After substituting the series observation t_1, \dots, t_n into basis functions, we get $N_1(t_1) = 1, N_1(t_2) = 0, \dots, N_{2k-1}(t_k) = 1, N_{2k}(t_k) = 0, \dots, N_{2n-1}(t_n) = 1, N_{2n}(t_n) = 0$; and into first derivative of basis functions, we get $N'_1(t_1) = 0, N'_1(t_2) = 1, \dots, N'_{2k-1}(t_k) = 0, N'_{2k}(t_k) = 1, \dots, N'_{2n-1}(t_n) = 0, N'_{2n}(t_n) = 1$. That means the matrices \mathbf{B} and \mathbf{C} in

MSE equation (1.92) are $n \times 2n$ dimensional and the elements are

$$\mathbf{B} = \{B\}_{ij} = \begin{cases} 1, & j = 2i - 1 \\ 0, & \text{otherwise} \end{cases} \quad (1.93)$$

$$\mathbf{C} = \{C\}_{ij} = \begin{cases} 1, & j = 2i \\ 0, & \text{otherwise} \end{cases} \quad (1.94)$$

where $i = 1, \dots, n$. The k -th $\Omega_{\lambda_k}^{(k)}$ is a $2n \times 2n$ matrix and its details is in appendix. Then the penalty term is

$$\Omega_\lambda = \sum_{k=1}^{n-1} \Omega_{\lambda_k}^{(k)}, \quad (1.95)$$

which is a bandwidth four matrix.

The solution to (1.92) is easily seen to be

$$\hat{\theta} = (\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\Omega_\lambda)^{-1} (\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \quad (1.96)$$

a generalized ridge regression. Then the fitted smoothing spline is given by

$$\hat{f}(t) = \sum_{i=1}^{2n} N_i(t) \hat{\theta}_i \quad (1.97)$$

A smoothing spline with parameters $\lambda(t)$ and γ is an example of a linear smoother Trevor Hastie (2009). This is because the estimated parameters in equation (1.96) are a linear combination of y_i and v_i . Denote by $\hat{\mathbf{f}}$ and $\hat{\mathbf{f}}'$ the $2n$ vector of fitted values $\hat{f}(t_i)$ and $\hat{f}'(t_i)$ at the training points t_i . Then

$$\begin{aligned} \hat{\mathbf{f}} &= \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\Omega_\lambda)^{-1} (\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \\ &\triangleq \mathbf{S}_{\lambda, \gamma} \mathbf{y} + \gamma \mathbf{T}_{\lambda, \gamma} \mathbf{v} \end{aligned} \quad (1.98)$$

$$\begin{aligned} \hat{\mathbf{f}}' &= \mathbf{C}(\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\Omega_\lambda)^{-1} (\mathbf{B}^\top \mathbf{y} + \gamma \mathbf{C}^\top \mathbf{v}) \\ &\triangleq \mathbf{U}_{\lambda, \gamma} \mathbf{y} + \gamma \mathbf{V}_{\lambda, \gamma} \mathbf{v} \end{aligned} \quad (1.99)$$

The fitted $\hat{\mathbf{f}}$ and $\hat{\mathbf{f}}'$ are linear in \mathbf{y} and \mathbf{v} , and the finite linear operators $\mathbf{S}_{\lambda, \gamma}$, $\mathbf{T}_{\lambda, \gamma}$, $\mathbf{U}_{\lambda, \gamma}$ and $\mathbf{V}_{\lambda, \gamma}$ are known as the smoother matrices. One consequence of this linearity is that the recipe for producing $\hat{\mathbf{f}}$ and $\hat{\mathbf{f}}'$ from \mathbf{y} and \mathbf{v} , do not depend on \mathbf{y} and \mathbf{v} themselves; $\mathbf{S}_{\lambda, \gamma}$, $\mathbf{T}_{\lambda, \gamma}$, $\mathbf{U}_{\lambda, \gamma}$ and $\mathbf{V}_{\lambda, \gamma}$ depend only on t_i , $\lambda(t)$ and γ .

Suppose in a traditional least squares fitting, \mathbf{B}_ξ is $N \times M$ matrix of M cubic-spline basis functions evaluated at the N training points x_i , with knot sequence ξ and $M \ll N$. Then the vector of fitted spline values is given by

$$\hat{\mathbf{f}} = \mathbf{B}_\xi (\mathbf{B}_\xi^\top \mathbf{B}_\xi)^{-1} \mathbf{B}_\xi \mathbf{y} = \mathbf{H}_\xi \mathbf{y} \quad (1.100)$$

Here the linear operator \mathbf{H}_ξ is a symmetric, positive semidefinite matrices, and $\mathbf{H}_\xi \mathbf{H}_\xi = \mathbf{H}_\xi$ (idempotent) Trevor Hastie (2009). In our case, it is easily seen that $\mathbf{S}_{\lambda,\gamma}$, $\mathbf{T}_{\lambda,\gamma}$, $\mathbf{U}_{\lambda,\gamma}$ and $\mathbf{V}_{\lambda,\gamma}$ are symmetric, positive semidefinite matrices as well. Additionally, by Cholesky decomposition

$$(\mathbf{B}^\top \mathbf{B} + \gamma \mathbf{C}^\top \mathbf{C} + n\Omega_\lambda)^{-1} = \mathbf{R}\mathbf{R}^\top, \quad (1.101)$$

it is easily to prove that $\mathbf{T}_{\lambda,\gamma} = \mathbf{B}\mathbf{R}\mathbf{R}^\top \mathbf{C}^\top$ and $\mathbf{U}_{\lambda,\gamma} = \mathbf{C}\mathbf{R}\mathbf{R}^\top \mathbf{B}^\top$, then we will have $\mathbf{T}_{\lambda,\gamma} = \mathbf{U}_{\lambda,\gamma}^\top$. When $\lambda = \gamma = 0$, the matrix $\mathbf{S}_{\lambda_0,\gamma_0} = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$ is idempotent.

Additionally, by playing with $\lambda(t)$ and γ , tractor spline has the following property:

- if $\lambda(t)$ is a piecewise constant and $\gamma \neq 0$, then f and f' are continuous, f'' is piecewise linear but not continuous at knots;
- if $\lambda(t)$ is a piecewise constant and $\gamma = 0$, the same as above;
- if $\lambda(t) = \lambda$ is a constant and $\gamma \neq 0$, the same as above;
- if $\lambda(t) = \lambda$ is a constant and $\gamma = 0$, then f , f' are continuous, f'' is piecewise linear and continuous at knots.

1.9.4 Adjusted Penalty Term and Parameter Function

To get the reconstructed trajectory in a multi-dimensional space, one can use tractor spline to find the trajectory in each dimensions simultaneously can combine them together at the end. Additionally, sometimes the data is not recorded in equal space. Due to the property of Hermite spline, the combination of multi-dimensional reconstructions and non-equal space data will bring some issues. Imagining this situation that a vehicle is moving along the x -axis, its x position changes consequently, but its y position might stay the same. By fitting \mathbf{x} and \mathbf{u} on x -axis, the tractor spline $f_x(t)$ will give us a best fit which returns smallest errors to the objective function. While with the same parameter $\lambda(t)$ and γ , $f_y(t)$ will return a cubic curve. However, it should give us a straight line as we expected. Moreover, in some circumstances, the time mark increases but \mathbf{f} and \mathbf{f}' keep the same, or changes slightly. Facing this situation, the Hermite spline will return a wiggle in one dimension and a curve in two dimensions. To get a reliable reconstruction, we introduce an adjusted term $\frac{(\Delta t_i)^\alpha}{(\Delta d_i)^\beta}$, where $\alpha \geq 0$ and $\beta \geq 0$, to the penalty function $\lambda(t)$, which means that the tractor spline should be penalized by its real difference of Δd_i and Δt_i between two points p_i and p_{i+1} . With this term, when \mathbf{u} goes down or equals to 0, it will make sure that the penalty function

will be large enough and return a straight line rather than a curve in each dimension of x and y . Because of the unit of the penalty term is m^2/t^3 , to keep the same scale in the space, α and β in the adjusted penalty term are chosen as 3 and 2. Then the final form of the penalty function is

$$\lambda(t) = \frac{(\Delta t_i)^3}{(\Delta d_i)^2} \lambda, \quad (1.102)$$

where $t_i \leq t < t_{i+1}$. Eventually in objective function there is one parameter λ controlling the curvature of tractor spline in different status, and another one parameter γ controlling the residuals of velocity.

1.10 Parameter Selection and Cross Validation

The problem of choosing the smoothing parameter is ubiquitous in curve estimation. And there are two different philosophical approaches to this question. The first one is to regard the free choice of smoothing parameter as an advantageous feature of the procedure. The other one is to find the parameter automatically by the data Green and Silverman (1993). We more prefer the latter one, use data to train our model and find the best parameters. The most well known method is cross-validation.

Assuming that the random errors has zero mean, the true regression curve $f(t)$ has the property that, if an observation y is taken at a point t , the value $f(t)$ is the best predictor of y in terms of returning a small value of $(y - f(t))^2$.

Now we focus on an observation y_i at point t_i as being a new observation by omitting it from the set of data, which are used to estimate \hat{f} . Denote by $\hat{f}^{(-i)}(t, \lambda)$ the estimated function from the remaining data, where λ is the smoothing parameter. Then $\hat{f}^{(-i)}(t, \lambda)$ minimizes

$$\frac{1}{n} \sum_{j \neq i} (y_j - f(t_j))^2 + \lambda \int f''^2 dt \quad (1.103)$$

and λ can be quantified by cross-validation score function

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}^{(-i)}(t_i, \lambda)\}^2. \quad (1.104)$$

The basis idea of cross-validation is to choose the value of λ that minimizes $CV(\lambda)$.

An efficient way to calculate cross validation score is given by Green and Silverman (1993). Through the equation (1.100), we know that the value of the smoothing spline

\hat{f} depend linearly on the data y_i . Define the matrix $A(\lambda)$, which is a map vector of observed values y_i to predicted values $\hat{f}(t_i)$. Then we have

$$\hat{\mathbf{f}} = A(\lambda)\mathbf{y} \quad (1.105)$$

and the following lemma.

Lemma 4. *The cross validation score satisfies*

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}(t_i)}{1 - A_{ii}(\lambda)} \right)^2 \quad (1.106)$$

where \hat{f} is the spline smoother calculated from the full data set $\{(t_i, y_i)\}$ with smoothing paramter λ .

For a tractor spline and its MSE function, there are two parameters need to be estimated λ and γ . Then the objective function (1.103) becomes

$$\frac{1}{n} \sum_{j \neq i} (y_j - f(t_j))^2 + \frac{\gamma}{n} \sum_{j \neq i} (v_j - f'(t_j))^2 + \int \lambda(t) f''^2 dt, \quad (1.107)$$

and the cross-validation score function is

$$CV(\lambda, \gamma) = \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}^{(-i)}(t_i, \lambda, \gamma)\}^2. \quad (1.108)$$

For a tractor spline, the parameter $\hat{\theta} = (B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}(B^\top \mathbf{y} + \gamma C^\top \mathbf{v})$ gives us

$$\begin{aligned} \hat{\mathbf{f}} &= B\hat{\theta} = B(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}B^\top \mathbf{y} + B(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}C^\top \mathbf{v} \\ &= S\mathbf{y} + \gamma T\mathbf{v}, \end{aligned} \quad (1.109)$$

$$\begin{aligned} \hat{\mathbf{f}}' &= C\hat{\theta} = C(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}B^\top \mathbf{y} + C(B^\top B + \gamma C^\top C + n\Omega_\lambda)^{-1}C^\top \mathbf{v} \\ &= U\mathbf{y} + \gamma V\mathbf{v}. \end{aligned} \quad (1.110)$$

From lemma 4, we can prove the following theorem.

Theorem 5. *The cross validation score of a tractor spline satisfies*

$$CV(\lambda, \gamma) = \frac{1}{n} \sum_{i=1}^n \left(\frac{\hat{f}(t_i) - y_i + \gamma \frac{T_{ii}}{1 - \gamma V_{ii}} (\hat{f}'(t_i) - v_i)}{1 - S_{ii} - \gamma \frac{T_{ii}}{1 - \gamma V_{ii}} U_{ii}} \right)^2 \quad (1.111)$$

where \hat{f} is the tractor spline smoother calculated from the full data set $\{(t_i, y_i, v_i)\}$ with smoothing parameter λ and γ .

The proof of Theorem 5 follows immediately from a lemma, and gives an expression for the deleted residuals $y_i - \hat{f}^{(-i)}(t_i)$ and $v_i - \hat{f}'^{(-i)}(t_i)$ in terms of $y_i - \hat{f}(t_i)$ and $v_i - \hat{f}'(t_i)$ respectively.

Lemma 5. For fixed λ, γ and i , denote $\mathbf{f}^{(-i)}$ by the vector with components $f_j^{(-i)} = \hat{f}^{(-i)}(t_j, \lambda, \gamma)$, $\mathbf{f}'^{(-i)}$ by the vector with components $f_j'^{(-i)} = \hat{f}'^{(-i)}(t_j, \lambda, \gamma)$, and define vectors \mathbf{y}^* and \mathbf{v}^* by

$$\begin{cases} y_j^* = y_j & j \neq i \\ y_i^* = \hat{f}^{(-i)}(t_i) & \text{otherwise} \end{cases}, \quad (1.112)$$

$$\begin{cases} v_j^* = v_j & j \neq i \\ v_i^* = \hat{f}'^{(-i)}(t_i) & \text{otherwise} \end{cases}. \quad (1.113)$$

Then

$$\hat{\mathbf{f}}^{(-i)} = S\mathbf{y}^* + \gamma T\mathbf{v}^* \quad (1.114)$$

$$\hat{\mathbf{f}}'^{(-i)} = U\mathbf{y}^* + \gamma V\mathbf{v}^* \quad (1.115)$$

1.11 Simulation

1.11.1 Numerical Examples

In this section, we examine the visual quality of the proposed method with four functions: Blocks, Bumps, HeaviSine and Doppler, which have been used in Donoho and Johnstone (1994), Donoho and Johnstone (1995) and Abramovich *et al.* (1998) because of their caricature features in imaging, spectroscopy and other scientific signal processing. However it is unfair for tractor spline fitting "jump" position in Blocks and Bumps function, because it fits position and velocity simultaneously and these points imply infinite first derivative in original functions, which are impossible for vehicles or individuals. In terms of this issue, we treat these functions as velocity, and use noise free points to generate accurate position data, then add noises back to them.

For calculating consideration, we use $n = 1024$ Nason (2010). Because all noises are randomly generated, for convenience of reinitialization and repetition of comparing, we set random seed as 2016. The noises are independent Gaussian distribution $\epsilon \sim N(0, 1)$ and signal-to-noise ratio (SNR) is 7. These data are treated as velocity (first derivative). By setting initial position $y_0 = 0$, acceleration $a_0 = 0$ and using the following formula

to calculate position

$$y_{i+1} = y_i + (v_i + v_{i+1}) \frac{t_{i+1} - t_i}{2}, \quad (1.116)$$

we can easily generate position data. Then we add some noises, which are independent Gaussian distribution $\epsilon \sim N(0, 1)$ and SNR is 7. For wavelet reconstruction, we use the threshold policy of "sure" and "BayesThresh" with levels $j = 4, \dots, 9$. Penalized B-spline is added in comparison. For tractor spline we have two parameters λ and γ . To evaluate the performance of the velocity term in objective function (1.81) and the adjusted penalty term in (1.102), the parameter γ is set as 0 in one reconstruction of tractor spline, whose objective function and solution become

$$J[f]_{\gamma=0} = \frac{1}{n} \sum_{i=1}^n (f(t_i) - y_i)^2 + \sum_{i=1}^{n-1} \lambda_i \int_{t_i}^{t_{i+1}} f''(t)^2 dt, \quad (1.117)$$

and

$$\hat{\theta}_{\gamma=0} = (\mathbf{B}^\top \mathbf{B} + n\Omega_\lambda)^{-1} \mathbf{B}^\top \mathbf{y} \quad (1.118)$$

and the adjusted penalty term in (1.102) was removed from another reconstruction, noted as "tractor spline without APT". Figure (1.3) to (1.6) display the original (velocity), generated position, wavelet with two different threshold methods, P-spline and three kinds of tractor spline fitted functions. The parameters λ and γ of a tractor spline are automatically selected from formula (1.111) by *optim()* function in *R*.

By comparing, we can see that all these methods can rebuild up the skeleton of generated trajectory. Wavelets(sure) method have more wiggles in interior interval than Wavelet(BayesThresh), and the latter one becomes fluctuation near boundary knots. P-spline gives much smoother fitting than wavelets, but the drawback is losing more specific details. Tractor spline without velocity might give us less fitting, as can be seen from Blocks and Bumps where there should be a straight line. Tractor spline without adjusted penalty term could also get over fitting when the direction changes more frequently than normal, although it catches specific feature in HeaviSine. The proposed tractor spline performs much better than other methods and returns the near-true trajectory reconstruction.

Figure 1.7 shows the estimated penalty function

$$\lambda(t) = \frac{(\Delta t)^3}{(\Delta d)^2} \lambda. \quad (1.119)$$

The left column illustrates the value of penalty function on different intervals and the right column is the projection on position. Bigger black dots present larger penalty values. It can be seen that $\lambda(t)$ adapts to the smoothness pattern of position and will

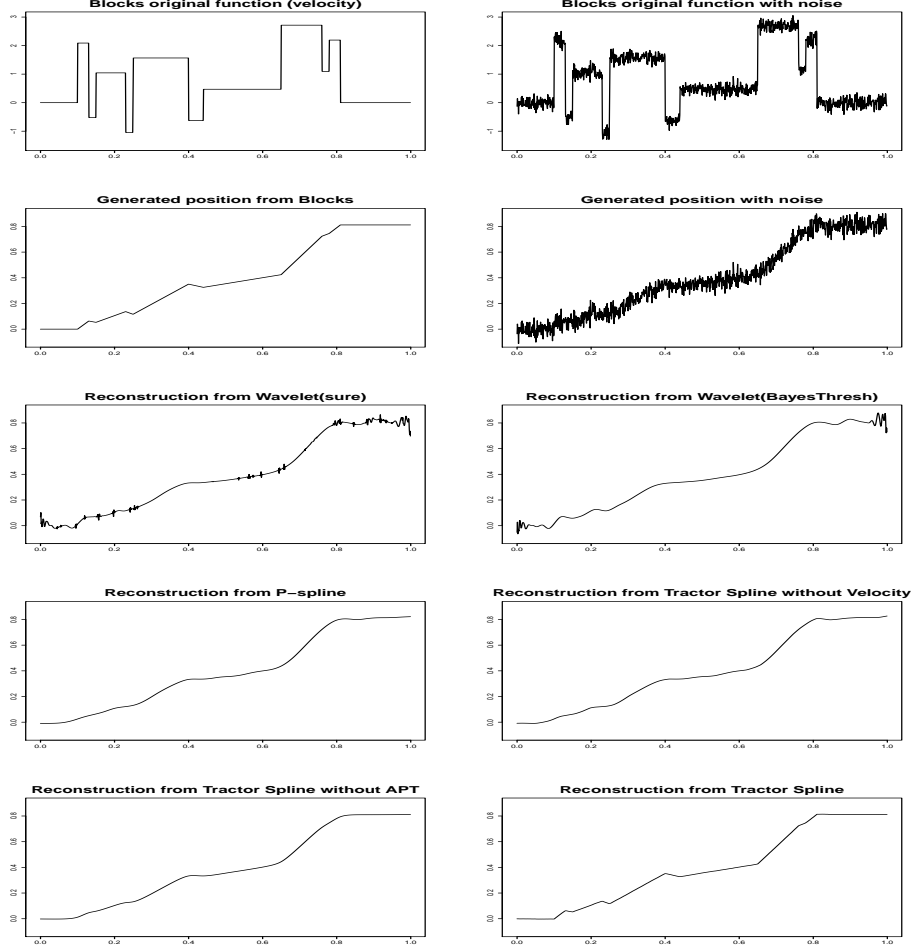


Figure 1.3: Numerical example: *Blocks*. (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.

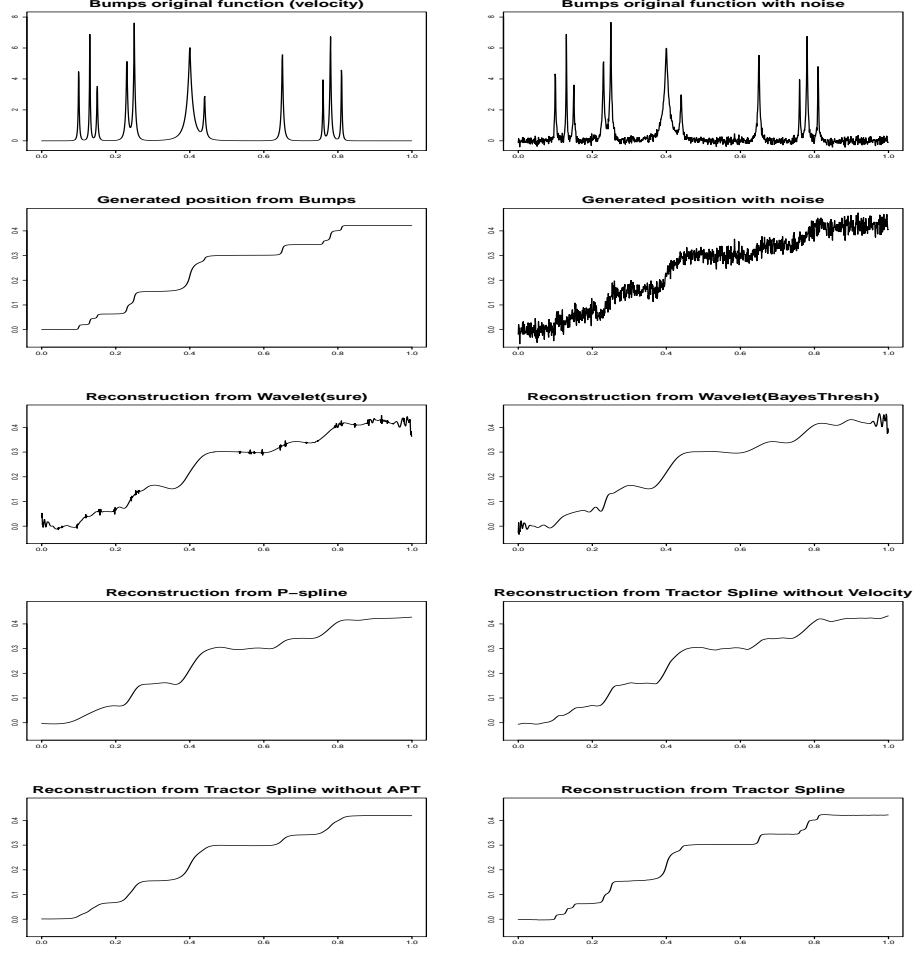


Figure 1.4: Numerical example: *Bumps*. (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.

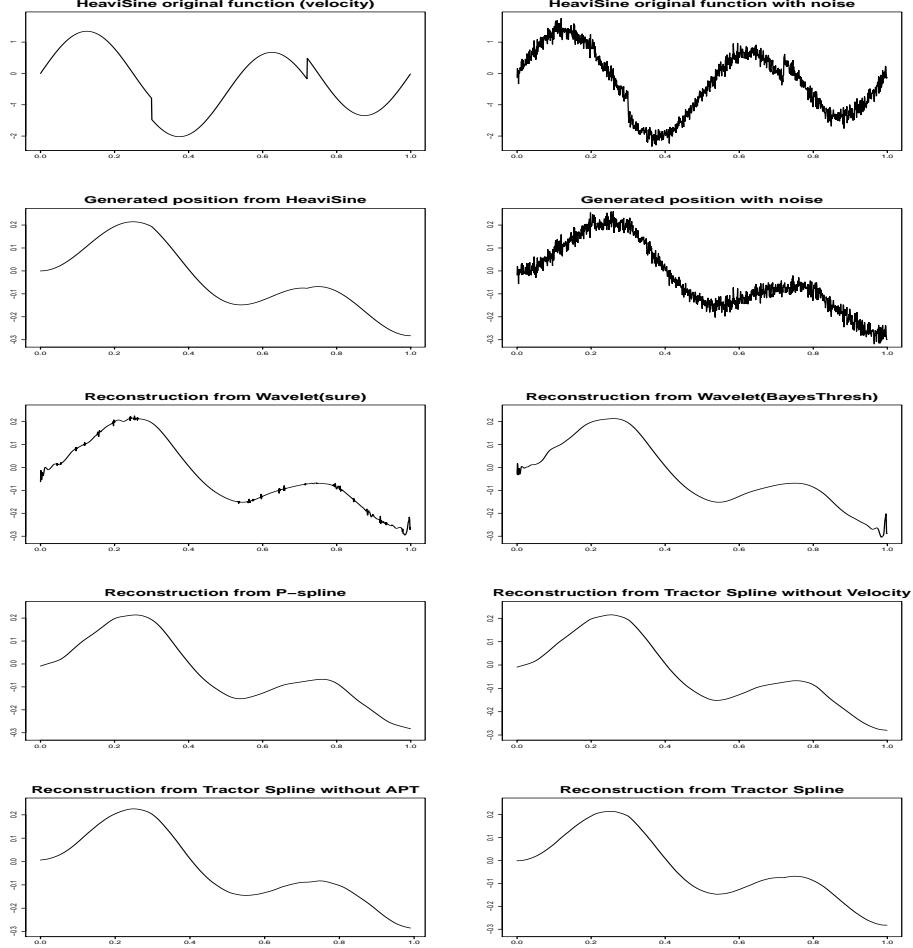


Figure 1.5: Numerical example: *HeaviSine*. (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.

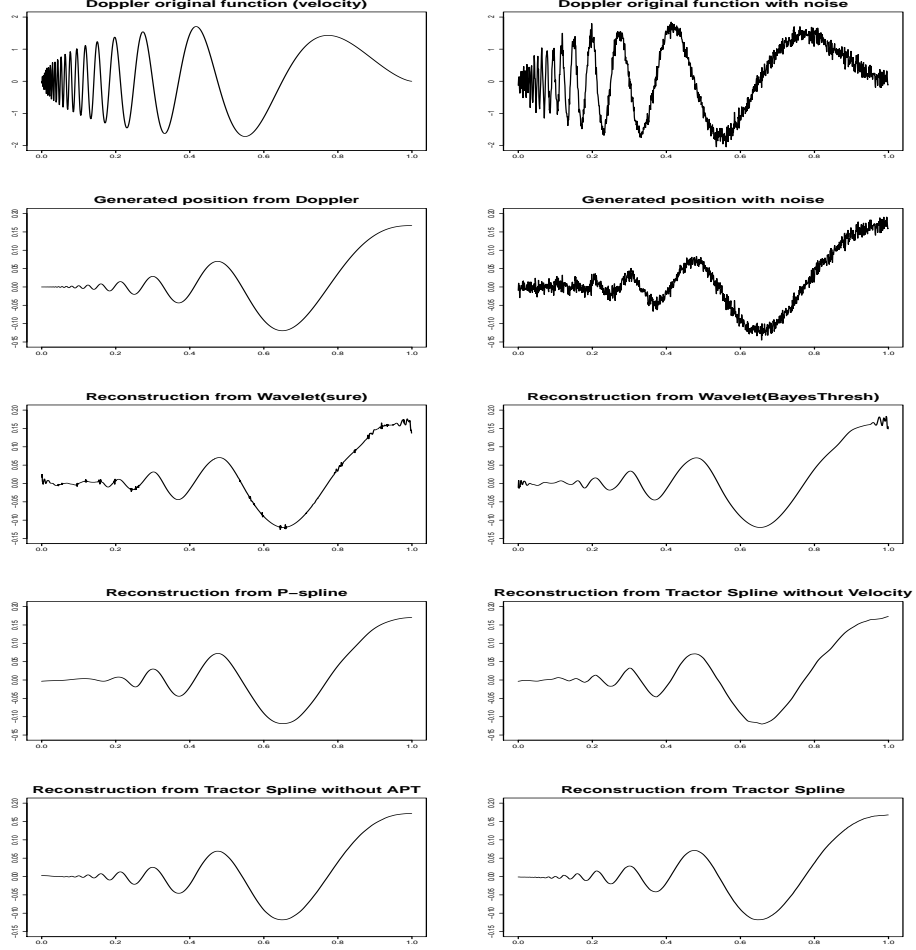


Figure 1.6: Numerical example: *Doppler*. (a) The true velocity function. (b) Velocity with Gaussian noise at SNR=7. (c) Generated position function. (d) Position with Gaussian noise at SNR=7. (e) Reconstruction from Wavelet with sure threshold. (f) Reconstruction from Wavelet with BayesThresh approach. (g) Reconstruction by P-spline. (h) Reconstruction by tractor spline setting $\gamma = 0$. (i) Reconstruction by tractor spline with normal penalty term. (j) Reconstruction by proposed tractor spline.

be large where a long time gap may occur. The details of how this penalty function works will be explained in next subsection.

Figure 1.8 demonstrates the estimated velocity functions. By taking the first derivative of fitted tractor spline, it is easy to get the original four velocity functions. The fitting of velocity is not as smooth as that in position, because we only care about the smoothness of position rather than velocity in our cross-validation formula (5). However, velocity information does help us reconstructing trajectory.

1.11.2 Evaluation

To examine the performance of tractor spline, we conducted an evaluation by comparing the mean square errors and true mean square errors, which are respectively calculated in

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{\lambda, \gamma}(t_i))^2, \quad (1.120)$$

$$\text{TMSE} = \frac{1}{n} \sum_{i=1}^n (f(t_i) - \hat{f}_{\lambda, \gamma}(t_i))^2. \quad (1.121)$$

The results are shown in table 1.1 and 1.2. All of these methods have good performances in fitting noisy data. The differences of mean square error between these methods are not significant, as can be seen from table 1.1. The proposed method is not the best among these simulations according to MSE. However, from table 1.2, tractor spline returns the smallest true mean square errors. The difference is significant, that means the reconstruction from tractor spline is closer to the true trajectory.

1.12 Conclusion and Discussion

In this paper, we proposed a tractor spline model with first derivative and adaptive penalty terms to reconstruct trajectory. This method performs better when we know \mathbf{z} and \mathbf{w} information than other methods. Additionally, the reconstruction of a tractor spline contains $4 \times (n - 1)$ parameters if we have n knots. By adding $2 \times (n - 2)$ constraints, the original function and its first derivative are continuous on each interior knots, the degrees of freedom will be $4 \times (n - 1) - 2 \times (n - 2) = 2n$. Because there are n location and n velocity data, thus we don't need to specify more parameters or add more constraints on the model. Although the mean square error of tractor spline is not the smallest comparing with other methods, the true mean square error is the smallest

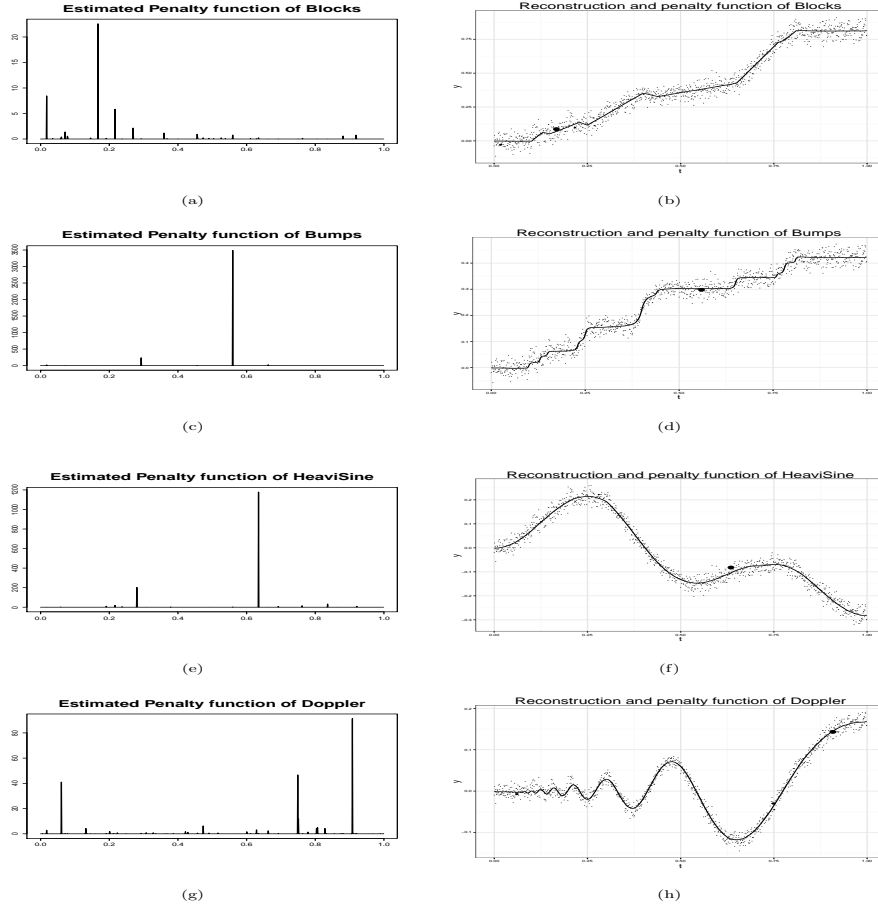


Figure 1.7: Estimated penalty functions. Left side shows how the value of $\lambda(t)$ changes on the interval. Right side projects $\lambda(t)$ into reconstructions. The bigger the blacks dots present, the larger the penalty values are.

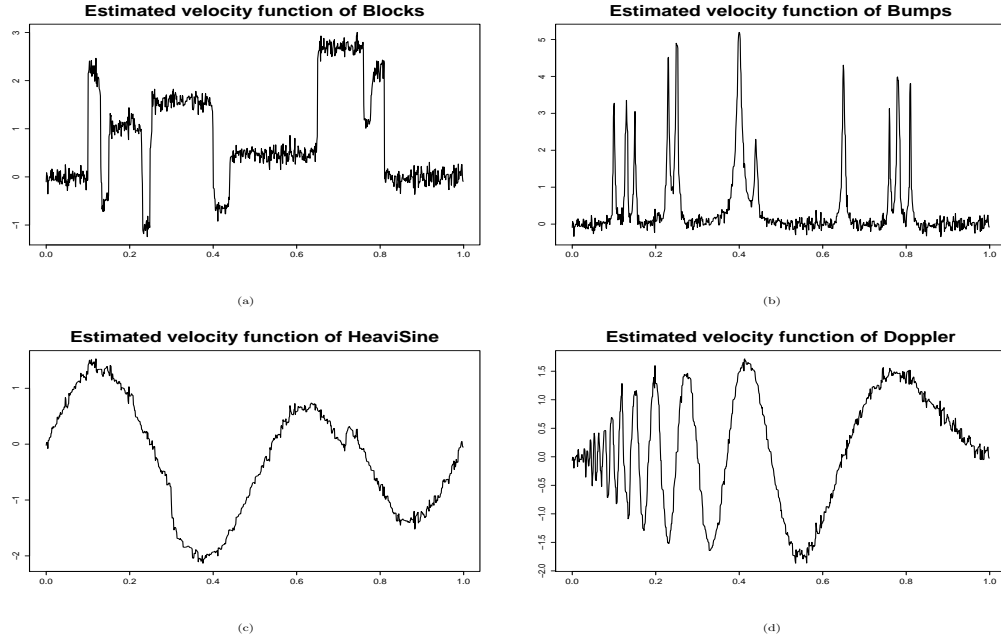


Figure 1.8: Estimated velocity functions by taking the first derivative of tractor spline. (a) Fitted *Blocks*. (b) Fitted *Bumps*. (c) Fitted *HeaviSine*. (d) Fitted *Doppler*.

Table 1.1: MSE. Mean square errors of different methods. The star sign (*) marks the smallest error among these methods under the same level. The difference is not significant.

MSE (10^{-4})	SNR	TS	$TS_{\gamma=0}$	$TS_{APT=0}$	P-spline	Wavelet(sure)	Wavelet(Bayes)
<i>Blocks</i>	7	16.53	15.99	16.69	16.14	*15.39	16.68
<i>Blocks</i>	3	89.79	*87.64	89.94	88.27	98.35	90.24
<i>Bumps</i>	7	4.40	4.19	4.55	4.33	*4.18	4.59
<i>Bumps</i>	3	23.93	*23.19	24.10	23.55	26.23	23.74
<i>HeaviSine</i>	7	4.16	4.01	4.16	4.02	*3.79	4.19
<i>HeaviSine</i>	3	22.63	*22.19	22.65	22.02	23.53	22.07
<i>Doppler</i>	7	1.15	*1.07	1.10	1.15	*1.07	1.13
<i>Doppler</i>	3	6.27	*5.94	6.28	6.05	6.85	6.29

Table 1.2: TMSE. True mean square errors of different methods. The star sign (*) marks the smallest error among these methods under the same level. The proposed tractor spline returns the smallest TMSE among all the methods under the same level except for *Doppler* with SNR=7. The differences are significant.

TMSE (10^{-6})	SNR	TS	$TS_{\gamma=0}$	$TS_{APT=0}$	P-spline	Wavelet(sure)	Wavelet(Bayes)
<i>Blocks</i>	7	*1.75	54.25	28.68	54.76	201.02	182.12
<i>Blocks</i>	3	*16.44	152.5	30.76	171.59	1138.08	712.36
<i>Bumps</i>	7	*1.64	23.44	21.10	24.21	71.71	69.26
<i>Bumps</i>	3	*8.51	77.78	37.12	77.52	330.77	238.79
<i>HeaviSine</i>	7	*1.53	7.80	1.56	9.54	55.37	44.88
<i>HeaviSine</i>	3	*8.21	33.56	8.49	34.26	240.72	110.49
<i>Doppler</i>	7	1.51	6.67	*1.08	8.26	14.87	12.01
<i>Doppler</i>	3	*8.10	22.14	8.25	19.95	81.48	50.33

most of the time. It means that the reconstruction is closer to the truth. In cross validation, we only focus on the errors of f ignoring that in f' . So the reconstruction of f' is not as smooth as that in f , which does not affect trajectory reconstruction. A drawback of tractor spline is the cost in finding local minimal parameters, where the cross validation algorithm returns a smaller score. So we optimized our code to make it runs as faster as it can.

1.13 Introduction

1.13.1 Spline

In interpolation and curve fitting, piecewise linear approximation may not have the practical significance of cubic spline, or even higher order, approximation. These "broken lines" are neither very smooth nor very efficient approximation. Researchers can go to piecewise polynomial approximation with higher order pieces De Boor *et al.* (1978), which is called spline method. A spline is a numeric function that is piecewise-defined by polynomial functions, and which possesses a high degree of smoothness at the places where the polynomial pieces connect (known as knots) Judd (1998)Chen (2009). Suppose we are given observed data t_1, t_2, \dots, t_n on interval $[0, 1]$, satisfying $0 \leq t_1 < t_2 < \dots < t_n \leq 1$. A piecewise polynomial function $f(t)$ can be obtained by

dividing the interval into contiguous intervals $(t_1, t_2), \dots, (t_{n-1}, t_n)$ and represented by a separate polynomial in each interval. For any continuous $f \in \mathbb{C}^{(m)}[0, 1]$, it can be represented in a linear combination of basis functions $h_m(t)$, just as every vector in a vector space can be represented as a linear combination of basis vectors. So we have

$$f(t) = \sum_{m=1}^M \beta_m h_m(t), \quad (1.122)$$

where β_m are coefficients Ellis *et al.* (2009).

Suppose we were attempt to fit a model of $f(t)$ by least squares without any restrictions, the best fitting $\hat{f}(t)$ would go through every given data to reduce sum of squares to zero. Most of the time, the results are unsatisfactory as explanations of the given data. The roughness penalty approach is introduced to quantify the notion of a rapidly fluctuating curve and then to pose the estimation problem in a way that makes explicit the necessary compromise between varying rapidly fluctuation and slowly trend in curve estimation Green and Silverman (1993). By adding a penalty term $\int_0^1 (f^{(m)}(t))^2 dt$, the curve estimate $\hat{f}(t)$ exists and is unique over all spline functions $f(t)$ with $m - 1$ continuous derivatives fitting observed data in the space $\mathbb{C}^{(m)}[0, 1]$, and it can be found by minimizing the following penalized mean residual sum squares

$$\text{MSE}(f, \lambda) = \frac{1}{n} \sum_{i=1}^n (f(t_i) - y_i)^2 + \lambda \int_0^1 (f^{(m)}(t))^2 dt, \quad (1.123)$$

where λ is a fixed smoothing parameter, (t_i, y_i) , $i = 1, \dots, n$ are observed data and $0 \leq t_1 < t_2 < \dots < t_n \leq 1$. In equation (1.123), the smoothing parameter λ controls the trade-off between over-fitting and bias. Smoothing spline provides a powerful tool to estimate nonparametric functions Hastie and Tibshirani (1990).

1.13.2 Gaussian Process Regression

A Gaussian Process is a collection of random variables, any finite number of which have a joint Gaussian distribution Rasmussen and Williams (2006). It is fully defined by its mean $m(t)$ and covariance $K(s, t)$ functions as

$$m(t) = \mathbb{E}[f(t)] \quad (1.124)$$

$$K(s, t) = \mathbb{E}[(f(s) - m(s))(f(t) - m(t))], \quad (1.125)$$

where s and t are two variables, and a function f distributed as such is denoted in form of

$$f \sim GP(m(t), K(s, t)). \quad (1.126)$$

Usually the mean function is assumed to be zero everywhere.

Given a set of input variables \mathbf{T} for function $f(t)$ and the output $\mathbf{y} = f(\mathbf{T}) + \varepsilon$ with independent identically distributed Gaussian noise ε with variance σ_n^2 , we can use the above definition to predict the value of the function $f_* = f(t_*)$ at a particular input t_* . As the noisy observations becoming

$$\text{cov}(y_p, y_q) = K(t_p, t_q) + \sigma_n^2 \delta_{pq} \quad (1.127)$$

where δ_{pq} is a Kronecker delta which is one iff $p = q$ and zero otherwise, the joint distribution of the observed outputs \mathbf{y} and the estimated output f_* according to prior is

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim N \left(0, \begin{bmatrix} K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I & K(\mathbf{T}, t_*) \\ K(t_*, \mathbf{T}) & K(t_*, t_*) \end{bmatrix} \right). \quad (1.128)$$

The posterior distribution over the predicted value is obtained by conditioning on the observed data

$$f_* | \mathbf{y}, \mathbf{T}, t_* \sim N(\bar{f}_*, \text{cov}(f_*)) \quad (1.129)$$

where

$$\bar{f}_* = \mathbb{E}[f_* | \mathbf{y}, \mathbf{T}, t_*] = K(t_*, \mathbf{T})[K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I]^{-1} \mathbf{y}, \quad (1.130)$$

$$\text{cov}(f_*) = K(t_*, t_*) - K(t_*, \mathbf{T})[K(\mathbf{T}, \mathbf{T}) + \sigma_n^2 I]^{-1} K(\mathbf{T}, t_*). \quad (1.131)$$

1.13.3 The Smoothing Spline as Bayes Estimates

It is possible to interpret the smoothing spline regression estimator as a Bayesian estimate when the mean function $r(\cdot)$ is given an improper prior distribution. Berline and Thomas-Agnan (2011) Wahba (1990)

Consider the model

$$y_i = f(t_i) + \varepsilon_i, \quad (1.132)$$

where $i = 1, \dots, n$, ε_i are i.i.d. Gaussian distributed noise with variance σ^2 . Assume $f \in \mathbb{H}^{(m)}[0, 1]$, where

$$\mathbb{H}^{(m)}[0, 1] = \{f : f^{(\nu)} \text{ absolutely continuous}, \nu = 0, \dots, m-1, \int_0^1 (f^{(m)}(t))^2 dt < \infty\}. \quad (1.133)$$

A smoothing spline \hat{f}_λ is the minimizer of objective function (1.123) Wang (1998). Equipped with an appropriate inner product

$$\langle f, g \rangle = \sum_{\nu=0}^{m-1} f^{(\nu)}(0)g^{(\nu)}(0) + \int_0^1 f^{(m)}g^{(m)}dt, \quad (1.134)$$

the space $\mathbb{H}^{(m)}[0, 1]$ becomes a reproducing kernel Hilbert space.

Let $\phi_\nu(t) = \frac{t^{\nu-1}}{(\nu-1)!}$ where $\nu = 1, \dots, m$ and $R_1(s, t) = \int_0^1 \frac{(s-u)_+^{m-1}}{(m-1)!} \frac{(t-u)_+^{m-1}}{(m-1)!} du$. Denote $S = \{\phi_\nu(t_i)\}_{n \times m}$ where $i = 1, \dots, n, \nu = 1, \dots, m$ and $Q = \{R_1(t_i, t_j)\}_{n \times n}$ where $i = 1, \dots, n, j = 1, \dots, n$. Kimeldorf and Wahba (1971) and Kimeldorf and Wahba (1970) proved that \hat{f}_λ has the form

$$\hat{f}(t) = \sum_{\nu=1}^m d_\nu \phi_\nu(t) + \sum_{i=1}^n c_i R_1(t, t_i). \quad (1.135)$$

By denoting $M = Q + n\lambda I$, Gu (2013) found that the coefficients will be given by

$$\mathbf{c} = (M^{-1} - M^{-1}S(S^\top M^{-1}S)^{-1}S^\top M^{-1})\mathbf{Y}, \quad (1.136)$$

$$\mathbf{d} = (S^\top M^{-1}S)^{-1}S^\top M^{-1}\mathbf{Y}. \quad (1.137)$$

1.14 A reproducing kernel on $\mathcal{C}_{p.w.}^2[0, 1]$

The minimizer $\eta(x)$ of

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \eta(x_i))^2 + \frac{\gamma}{n} \sum_{i=1}^n (v_i - \eta'(x_i))^2 + \lambda \int_0^1 \eta''^2 dx \quad (1.138)$$

in the space $\mathcal{C}_{p.w.}^2[0, 1] = \{f : f, f' \text{ are continuous and } f'' \text{ is piecewise continuous on } [0, 1]\}$ is a tractor spline. Equipped with an appropriate inner product

$$(f, g) = f(0)g(0) + f'(0)g'(0) + \int_0^1 f''g'' dx, \quad (1.139)$$

the space $\mathcal{C}_{p.w.}^2[0, 1]$ is made a reproducing kernel Hilbert space. In fact, the representer $R_x(\cdot)$ is

$$R_x(y) = 1 + xy + \int_0^1 (x-u)_+(y-u)_+ du. \quad (1.140)$$

It can be seen that $R_x(0) = 1$, $R'_x(0) = x$, and $R''_x(y) = (x-y)_+$.

The two terms of the reproducing kernel $R(x, y) = R_x(y) = R_0(x, y) + R_1(x, y)$, where

$$R_0(x, y) = 1 + xy \quad (1.141)$$

$$R_1(x, y) = \int_0^1 (x-u)_+(y-u)_+ du \quad (1.142)$$

are both non-negative definite themselves.

Theorem 6. *If the reproducing kernel R of a space \mathcal{H} on domain X can be decomposed into $R = R_0 + R_1$, where R_0 and R_1 are both non-negative definite, $R_0(x, \cdot), R_1(x, \cdot) \in \mathcal{H}$, $\forall x \in X$, and $(R_0(x, \cdot), R_1(y, \cdot)) = 0$, $\forall x, y \in X$, then the spaces \mathcal{H}_0 and \mathcal{H}_1 corresponding respectively to R_0 and R_1 form a tensor sum decomposition of \mathcal{H} . Conversely, if R_0 and R_1 are both non-negative definite and $\mathcal{H}_0 \cap \mathcal{H}_1 = \{0\}$, then $\mathcal{H} = \mathcal{H}_0 \oplus \mathcal{H}_1$ has a reproducing kernel $R = R_0 + R_1$.*

According to Theorem 1, R_0 can correspond the space of polynomials $\mathcal{H}_0 = \{f : f'' = 0\}$ with an inner product $(f, g)_0 = f(0)g(0) + f'(0)g'(0)$, and R_1 corresponds the orthogonal complement of \mathcal{H}_0

$$\mathcal{H}_1 = \{f : f(0) = 0, f'(0) = 0, \int_0^1 f''^2 dx < \infty\}$$

with inner product $(f, g)_1 = \int_0^1 f''g'' dx$. Thus, \mathcal{H}_0 and \mathcal{H}_1 are two subspaces of the $\mathcal{C}_{p.w.}^2[0, 1]$, and the reproducing kernel is $R_x(\cdot) = R_0(x, \cdot) + R_1(x, \cdot)$.

Define a new notation $\dot{R}(x, y) = \frac{\partial R}{\partial x}(x, y) = \frac{\partial R_0}{\partial x}(x, y) + \frac{\partial R_1}{\partial x}(x, y) = y + \int_0^x (y-u)_+ du$. Obviously $\dot{R}_x(y) \in \mathcal{C}_{p.w.}^2[0, 1]$. Additionally, we have $\dot{R}_x(0) = 0$, $\dot{R}'_x(0) = \frac{\partial \dot{R}_x}{\partial y}(0) = 1$, and $\dot{R}_x''(y) = \begin{cases} 0 & x \leq y \\ 1 & x > y \end{cases}$. Then, for any $f \in \mathcal{C}_{p.w.}^2[0, 1]$, it gives us

$$(\dot{R}_x, f) = \dot{R}_x(0)f(0) + \dot{R}'_x(0)f'(0) + \int_0^1 \dot{R}_x'' f'' du = f'(0) + \int_0^y f'' du = f'(y).$$

It can be seen that the first term $\dot{R}_0 = y \in \mathcal{H}_0$, and the space spanned by the second term $\dot{R}_1 = \int_0^x (y-u)_+ du$, denoted as $\dot{\mathcal{H}}$, is not in \mathcal{H}_1 , but $\dot{\mathcal{H}} \cap \mathcal{H}_1 \neq \emptyset$. Then we have a new space $\mathcal{H}_* = \dot{\mathcal{H}} \cup \mathcal{H}_1$. Thus the two new sub spaces in $\mathcal{C}_{p.w.}^2[0, 1]$ are \mathcal{H}_0 and \mathcal{H}_* .

1.15 Computation of Polynomial Smoothing Splines

Given the sample points $x_j, j = 1, \dots, n$ in equation (1.138) and noting that the space

$$\mathcal{A} = \{f : f = \sum_{j=1}^n \alpha_j R_1(x_j, \cdot) + \sum_{j=1}^n \beta_j \dot{R}_1(x_j, \cdot)\} \quad (1.143)$$

is a closed linear subspace of \mathcal{H}_* . Then $\eta \in \mathcal{C}_{p.w.}^2[0, 1]$ can be written as

$$\eta(x) = d_1 + d_2 x + \sum_{j=1}^n c_j R_1(x_j, x) + \sum_{i=j}^n b_j \dot{R}_1(x_j, \cdot) + \rho(x) \quad (1.144)$$

where \mathbf{d}, \mathbf{c} and \mathbf{b} are coefficients, and $\rho(x) \in \mathcal{H}_* \ominus \mathcal{A}$.

The equation (1.138) can be written as

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \left(Y_i - d_1 - d_2 x - \sum_{j=1}^n c_j R_1(x_j, x_i) - \sum_{j=1}^n b_j \dot{R}_1(x_j, x_i) - \rho(x_i) \right)^2 \\ & \frac{\gamma}{n} \sum_{i=1}^n \left(V_i - d_2 - \sum_{j=1}^n c_j R'_1(x_j, x_i) - \sum_{j=1}^n b_j \dot{R}'_1(x_j, x_i) - \rho'(x_i) \right)^2 \\ & + \lambda \int_0^1 \left(\sum_{j=1}^n c_j R''_1(x_j, x) + \sum_{j=1}^n c_j \dot{R}''_1(x_j, x) + \rho''(x) \right)^2 dx \end{aligned}$$

By orthogonality, $\rho(x_i) = (R_1(x_i, \cdot), \rho) = 0$, $\rho'(x_i) = (\dot{R}_1(x_i, \cdot), \rho') = 0$, $i = 1, \dots, n$. Denoting by

$$\begin{aligned} S &= \{S_{ij}\}_{n \times 2} = \begin{bmatrix} 1 & x_i \end{bmatrix}, \quad Q = \{Q_{ij}\}_{n \times n} = R_1(x_j, x_i), \quad P = \{P_{ij}\}_{n \times n} = \dot{R}_1(x_j, x_i), \\ S' &= \{S'_{ij}\}_{n \times 2} = \begin{bmatrix} 0 & 1 \end{bmatrix}, \quad Q' = \{Q'_{ij}\}_{n \times n} = R'_1(x_j, x_i), \quad P' = \{P'_{ij}\}_{n \times n} = \dot{R}'_1(x_j, x_i). \end{aligned}$$

and noting that $\int_0^1 R''_1(x_i, x) R''_1(x_j, x) dx = R_1(x_i, x_j)$, $\int_0^1 R''_1(x_i, x) \dot{R}''_1(x_j, x) dx = \int_0^v (x_i - x) dx = \dot{R}_1(x_j, x_i)$, and $\int_0^1 \dot{R}''_1(x_i, x) \dot{R}''_1(x_j, x) dx = \int_0^v 1 dx = \dot{R}'_1(x_i, x_j)$, where $v = \min(x_i, x_j)$, the above equation can be written as

$$\begin{aligned} & (\mathbf{Y} - S\mathbf{d} - Q\mathbf{c} - P\mathbf{b})^\top (\mathbf{Y} - S\mathbf{d} - Q\mathbf{c} - P\mathbf{b}) + \gamma (\mathbf{V} - S'\mathbf{d} - Q'\mathbf{c} - P'\mathbf{b})^\top (\mathbf{V} - S'\mathbf{d} - Q'\mathbf{c} - P'\mathbf{b}) \\ & + n\lambda (\mathbf{c}^\top Q\mathbf{c} + 2\mathbf{c}^\top P\mathbf{b} + \mathbf{b}^\top P'\mathbf{b}) + n\lambda(\rho, \rho). \end{aligned} \quad (1.145)$$

Note that ρ only appears in the third term in (1.145), which is minimised at $\rho = 0$. Hence, a polynomial smoothing spline resides in the space $\mathcal{H}_0 \oplus \mathcal{A}$ of finite dimension. Then the solution to (1.138) could be computed via minimization of the first two terms in (1.145) with respect to \mathbf{d} , \mathbf{c} and \mathbf{b} .

1.16 Polynomial Smoothing Splines as Bayes Estimates

Now in the model $Y = \eta(x) + \epsilon$ and $V = \eta'(x) + \frac{\epsilon}{\gamma}$ where $\epsilon \sim N(0, \sigma^2)$, according to equation (1.144), for $\eta(x) \in \mathcal{C}_{p.w.}^2[0, 1]$ and $x \in X$, we have

$$\eta(x) = (d_1 + d_2 x) + \sum_{i=1}^n c_i R_1(x_i, x) + \sum_{i=1}^n b_i \dot{R}_1(x_i, x). \quad (1.146)$$

The covariance functions for Y, V and η, η' are

$$\begin{aligned}
\mathbb{E}(\eta(x)\eta(y)) &= \tau^2 R_0(x, y) + \beta R_1(x, y) & \mathbb{E}(\eta(x)\eta'(y)) &= \tau^2 R'_0(x, y) + \beta R'_1(x, y) \\
\mathbb{E}(\eta'(x)\eta(y)) &= \tau^2 \dot{R}_0(x, y) + \beta \dot{R}_1(x, y) & \mathbb{E}(\eta'(x)\eta'(y)) &= \tau^2 \dot{R}'_0(x, y) + \beta \dot{R}'_1(x, y) \\
\mathbb{E}(y_i, y_j) &= \tau^2 R_0(x_i, x_j) + \beta R_1(x_i, x_j) & \mathbb{E}(v_i, v_j) &= \tau^2 \dot{R}'_0(x_i, x_j) + \beta \dot{R}'_1(x_i, x_j) \\
&+ \sigma^2 \delta_{ij} & &+ \frac{\sigma^2}{\gamma} \delta_{ij} \\
\mathbb{E}(v_i, y_j) &= \tau^2 \dot{R}_0(x_i, x_j) + \beta \dot{R}_1(x_i, x_j) & \mathbb{E}(y_i, v_j) &= \tau^2 R'_0(x_i, x_j) + \beta R'_1(x_i, x_j) \\
\mathbb{E}(y_i, \eta(x)) &= \tau^2 R_0(x_i, x) + \beta R_1(x_i, x) & \mathbb{E}(y_i, \eta'(x)) &= \tau^2 R'_0(x_i, x) + \beta R'_1(x_i, x) \\
\mathbb{E}(v_i, \eta(x)) &= \tau^2 \dot{R}_0(x_i, x) + \beta \dot{R}_1(x_i, x) & \mathbb{E}(v_i, \eta'(x)) &= \tau^2 \dot{R}'_0(x_i, x) + \beta \dot{R}'_1(x_i, x)
\end{aligned}$$

where $R(x, y)$ is taken from (1.140).

Observing $Y_i \sim N(\eta(x_i), \sigma^2)$ and $V_i \sim N(\eta(x_i), \frac{\sigma^2}{\gamma})$, the joint distribution of Y, V and $\eta(x)$ is normal with mean zero and a covariance matrix can be found. The posterior mean of $\eta(x)$ is

$$\begin{aligned}
\mathbb{E}(\eta|Y, V) &= \begin{bmatrix} \text{cov}(Y, \eta) & \text{cov}(\eta, V) \end{bmatrix} \begin{bmatrix} \text{var}(Y) & \text{cov}(Y, V) \\ \text{cov}(V, Y) & \text{var}(V) \end{bmatrix}^{-1} \begin{bmatrix} Y \\ V \end{bmatrix} \\
&= \begin{bmatrix} \tau^2 \phi^\top S^\top + \beta \xi^\top & \tau^2 \phi^\top S'^\top + \beta \psi^\top \end{bmatrix} \begin{bmatrix} \tau^2 S S^\top + \beta Q + \sigma^2 I & \tau^2 S S'^\top + \beta P \\ \tau^2 S' S^\top + \beta Q' & \tau^2 S' S'^\top + \beta P' + \frac{\sigma^2}{\gamma} I \end{bmatrix}^{-1} \begin{bmatrix} Y \\ V \end{bmatrix} \\
&= \begin{bmatrix} \rho \phi^\top S^\top + \xi^\top & \rho \phi^\top S'^\top + \psi^\top \end{bmatrix} \begin{bmatrix} \rho S S^\top + Q + n\lambda I & \rho S S'^\top + P \\ \rho S' S^\top + Q' & \rho S' S'^\top + P' + \frac{n\lambda}{\gamma} I \end{bmatrix}^{-1} \begin{bmatrix} Y \\ V \end{bmatrix} \\
&= (\phi^\top \rho \begin{bmatrix} S \\ S' \end{bmatrix}^\top + \begin{bmatrix} \xi^\top & \psi^\top \end{bmatrix}) \left(\rho \begin{bmatrix} S \\ S' \end{bmatrix}^\top \begin{bmatrix} S \\ S' \end{bmatrix} + \begin{bmatrix} Q + n\lambda I & P \\ Q' & P' + \frac{n\lambda}{\gamma} I \end{bmatrix} \right)^{-1} \begin{bmatrix} Y \\ V \end{bmatrix} \\
&\triangleq \phi^\top \rho T^\top (\rho T^\top T + M)^{-1} \begin{bmatrix} Y \\ V \end{bmatrix} + \begin{bmatrix} \xi^\top & \psi^\top \end{bmatrix} (\rho T^\top T + M)^{-1} \begin{bmatrix} Y \\ V \end{bmatrix}
\end{aligned} \tag{1.147}$$

where ϕ is 2×1 matrix with entry 1 and x , ξ is $n \times 1$ matrix with i th entry $R(x_i, x)$ and ψ is $n \times 1$ matrix with i th entry $\dot{R}(x_i, x)$, $\rho = \tau^2/\beta$ and $n\lambda = \sigma^2/\beta$.

Lemma 6. Suppose M is symmetric and nonsingular and S is of full column rank.

$$\begin{aligned}
\lim_{\rho \rightarrow \infty} (\rho T T^\top + M)^{-1} &= M^{-1} - M^{-1} T (T^\top M^{-1} T)^{-1} T^\top M^{-1}, \\
\lim_{\rho \rightarrow \infty} \rho T^\top (\rho T T^\top + M)^{-1} &= (T^\top M^{-1} T)^{-1} T^\top M^{-1}.
\end{aligned}$$

Setting $\rho \rightarrow \infty$ in equation (1.147) and applying Lemma 6, the posterior mean $\mathbb{E}(\eta(x)|\mathbf{Y}, \mathbf{V})$ is of the form $\eta = \phi^\top \mathbf{d} + \xi^\top \mathbf{c} + \psi^\top \mathbf{b}$, with the coefficients given by

$$\mathbf{d} = (T^\top M^{-1} T)^{-1} T^\top M^{-1} \begin{bmatrix} Y \\ V \end{bmatrix},$$

$$\begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix} = (M^{-1} - M^{-1} T (T^\top M^{-1} T)^{-1} T^\top M^{-1}) \begin{bmatrix} Y \\ V \end{bmatrix},$$

where $T = \begin{bmatrix} S \\ S' \end{bmatrix}$ and $M = \begin{bmatrix} Q + n\lambda I & P \\ Q' & P' + \frac{n\lambda}{\gamma} I \end{bmatrix}$.

It is easy to verify that d, c, b are the solutions to

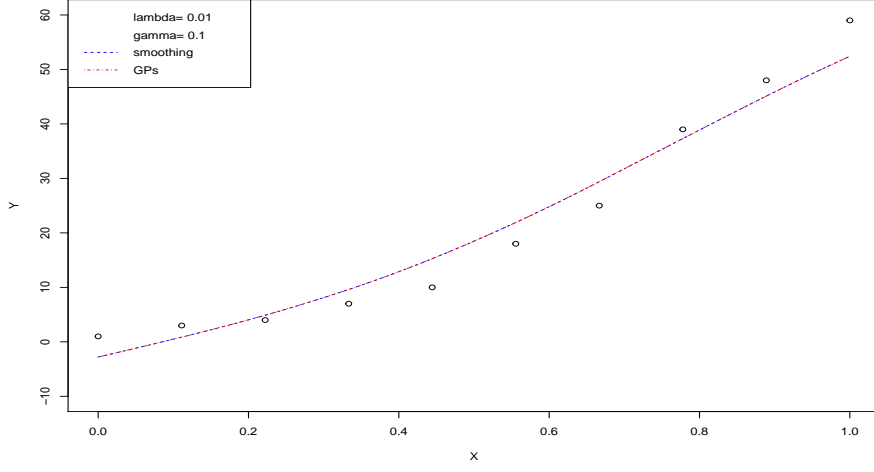
$$\begin{cases} S^\top (Sd + Qc + Pb - Y) + \gamma S'^\top (S'd + P^\top c + S'^\top P'b - V) = 0, \\ Q(Sd + (Q + n\lambda I)c + Pb - Y) + P(\gamma S'd + \gamma P^\top c + (\gamma P' + n\lambda I)b - \gamma V) = 0, \\ P^\top (Sd + (Q + n\lambda I)c + Pb - Y) + P'(\gamma S'b + P^\top c + (\gamma P' + n\lambda I)b - \gamma V) = 0. \end{cases}$$

1.17 Numeric Simulation of Smoothing Spline and GPR

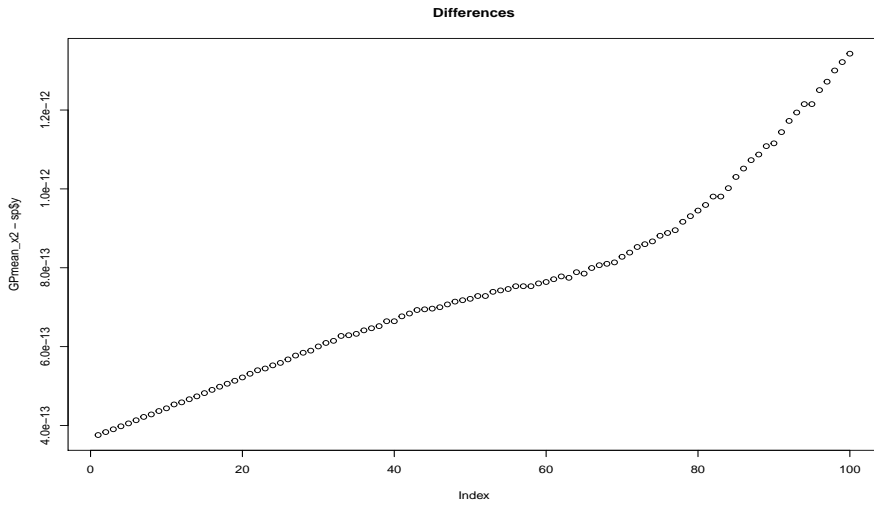
Given \mathbf{X} a length-10 sequence from 0 to 1,

$$\mathbf{Y} = [1, 3, 4, 7, 10, 18, 25, 39, 48, 59],$$

$$v_i = \begin{cases} \frac{y_{i+1} - y_i}{x_{i+1} - x_i} & \text{if } 1 \leq i \leq 9 \\ 0 & \text{if } i = 10 \end{cases}. \text{ A simulated result is in figure 1.}$$



(a)



(b)

Figure 1.9: (a) Comparing two methods under the same parameters $\lambda = 0.01$ and $\gamma = 0.1$. In this graph, the blue line is reconstruction from tractor spline, the red line is the mean of Gaussian Process, which is the posterior $\mathbb{E}(\eta(x)|\mathbf{Y}, \mathbf{V})$. (b) The differences between two methods under the same parameters.

3.18 State Space Models

State space models are the natural form of system models relying on the general concept of state. If we describe a system as an operator mapping from the space of inputs to the space of outputs, then we may need the entire input-output history of the system together with the planned input in order to compute the future output values Hangos *et al.* (2006). In an alternative way, by using new information at time t containing all the past information up to the current state and initial conditions to get the current output is possible, that is known as a sequential method. A genetic state space model consists of two sets of equations: state equation and output equation. The state equation describes the evolution of the true input and state variables sequentially as a function and passes the variable one after one, generally, with some noises. The output equation catches the input values and interprets it out by an algebraic equation. A general state space model looks like the following form

$$\text{State equation } x_t = G_t(x_{t-1}) + w_t, \quad (3.148)$$

$$\text{Output equation } y_t = F_t(x_t) + \epsilon_t \quad (3.149)$$

with an initial state x_0 , where ϵ_t and w_t are noises passing through the process G_t and F_t . x_t are true status variables and y_t are output values. Many researchers have been interested in this model and its application because of its good property. It can be used to model univariate or multivariate time series, also in the presence of non-stationarity, structural changes, and irregular patterns Petris *et al.* (2009).

The most simple and important system is given by Gaussian linear state space models, also known by dynamic linear models (DLM), which defines a very general class of non-stationary time series models. Firstly, the model is linear, which means G_t and F_t are linear processes and satisfying linearity property. Secondly, the it is specified by a normal prior distribution for the p -dimensional state vector at initial state $t = 0$,

$$x_0 \sim N_p(m, 0, C_0)$$

and two independent zero mean normal distributed noises $\epsilon_t \sim N_p(0, V_t)$ and $w_t \sim N_p(0, W_t)$ Petris *et al.* (2009). The well known Kalman Filter is a particular algorithm that is used to solve state space models in the linear case. This was first derived by Kalman Kalman *et al.* (1960).

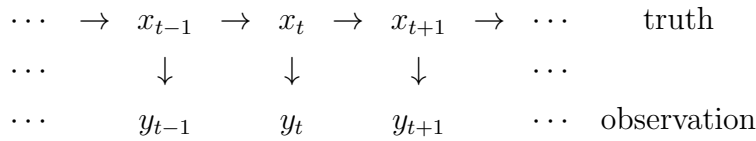
In a nonlinear state space model, the process G_t and F_t are no longer linear functions and the situation becomes more complicated. Here gives a simple nonlinear example

of such a model, which has been used extensively in the literature for benchmarking numerical filtering techniques Kitagawa (1996) West (1993) Gordon *et al.* (1993) assuming the sequence is Markovian.

$$x_t = \frac{x_{t-1}}{2} + 25 \frac{x_{t-1}}{1 + x_{t-1}^2} + 8 \cos(1.2t) + u_t$$

$$y_t = \frac{x_t^2}{20} + v_t,$$

where $u_t \sim N(0, \sigma_u^2)$, $v_t \sim N(0, \sigma_v^2)$, $\sigma_u^2 = 10$ and $\sigma_v^2 = 1$ are considered fixed and known. The initial state $x_0 \sim N(0, 10)$. The assumption Markovian keeps the current state x_t only depending on the previous one step x_{t-1} and the observed y_t depending on x_t . A state-space is shown in the diagram below:



In applications, the process function G_t and F_t contain unknown parameters to be estimated De Jong (1988) and the target is to estimate the true states on sequential observations y_1, \dots, y_t . Then it becomes to estimate a joint density of $p(x_{1:t}, \theta \mid y_{1:t})$, where $x_{1:t} = \{x_1, x_2, \dots, x_t\}$ are the hidden states and $y_{1:t} = \{y_1, y_2, \dots, y_t\}$ are the observed outcomes and θ is a set of unknown parameters.

3.19 Sequential Monte Carlo Method

The use of Monte Carlo methods for non-linear filtering can be traced back to the pioneering contributions of Handschin and Mayne (1969) Handschin and Mayne (1969) and Handschin (1970) Handschin (1970). These researchers tried to use an importance sampling paradigm to approximate the target distributions. Later on, an importance sampling algorithms were implemented sequentially in the non-linear filtering context. This algorithm is called sequential importance sampling, often abbreviated SIS, and has been known since the early 1970s. Limited by the power of computers and suffering from sample impoverishment or weight degeneracy, the SIS didn't develop very well until 1993. A particle filter algorithm was proposed to allow rejuvenation of the set of samples by duplicating the samples with high importance weights and, on the contrary, removing samples with low weights Cappé *et al.* (2009). Since then, sequential Monte Carlo (SMC) methods have been applied in many different fields including but not limited to computer vision, signal processing, control, econometrics, finance, robotics, and statistics Arnaud Doucet (2011) Ristic *et al.* (2004).

3.19.1 Filtering Problem and Estimation

Sequential Monte Carlo method, also known as particle filter, is a technique based on sampling and importance sampling methods to find the best state estimation given by Gordon in 1993 Gordon *et al.* (1993) and was the first successful application of sequential Monte Carlo techniques to the field of non-linear filtering Cappé *et al.* (2009). In the state space model, a generic particle filter estimates the posterior distribution of the hidden states using the observation measurement process. The filtering problem is to estimate sequentially the values of the hidden states x_k given the values of the observation process $y_{1:k}$ at any time step k . In another word, it is to find the value of $p(x_k | y_{1:k})$. The process is divided into two steps: prediction and updating. In the prediction step, the assumption of Markov Chain is the current status x_k only depends on the previous one x_{k-1} . Then we can calculate the probability of x_k by

$$\begin{aligned} p(x_k | y_{1:k-1}) &= \int p(x_k, x_{k-1} | y_{1:k-1}) dx_{k-1} \\ &= \int p(x_k | x_{k-1}, y_{1:k-1}) p(x_{k-1} | y_{1:k-1}) dx_{k-1} \\ &= \int p(x_k | x_{k-1}) p(x_{k-1} | y_{1:k-1}) dx_{k-1}. \end{aligned}$$

In the updating step, once $p(x_k | y_{1:k-1})$ is known, $p(x_k | y_{1:k})$ can be found by

$$\begin{aligned} p(x_k | y_{1:k}) &= \frac{p(y_k | x_k, y_{1:k-1}) p(x_k | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_k) p(x_k | y_{1:k-1})}{p(y_k | y_{1:k-1})}, \end{aligned}$$

where the normalization $p(y_k | y_{1:k-1}) = \int p(y_k | x_k) p(x_k | y_{1:k-1}) dx_k$ Arulampalam *et al.* (2002).

Imagine that the stat space is partitioned as many parts, in which the particles are filled according to some probability measure. The higher probability, the denser the particles are concentrated. Suppose the particles x_1, \dots, x_N are drawn from the target probability density function $p(x)$, then these particles are used to estimate the expectation and variance of $f(x)$ by

$$\begin{aligned} E(f(x)) &= \int_a^b f(x) p(x) dx \\ Var(f(x)) &= E(f(x) - E(f(x)))^2 p(x) dx. \end{aligned}$$

Back to our target, the posterior distribution or density is empirically represented by

a weighted sum of samples x_1, \dots, x_N

$$\hat{p}(x_n | y_{1:k}) = \frac{1}{N} \sum_{i=1}^N \delta(x_n - x_n^{(i)}) \approx p(x_n | y_{1:k}),$$

where $f(x) = \delta(x_n - x_n^{(i)})$ is Dirac delta function. When N is sufficiently large, $\hat{p}(x_n | y_{1:k})$ approximates the true posterior $p(x_n | y_{1:k})$. By this approximation, the filtering problem becomes to get the expectation of current status

$$\begin{aligned} E(f(x_n)) &\approx \int f(x_n) \hat{p}(x_n | y_{1:k}) dx_n \\ &= \frac{1}{N} \sum_{i=1}^N \int f(x_n) \delta(x_n - x_n^{(i)}) dx_n \\ &= \frac{1}{N} \sum_{i=1}^N f(x_n^{(i)}). \end{aligned}$$

The expectation is the mean of the status of all particles x_1, \dots, x_N .

However, the posterior distribution is unknown and impossible to sample from the true posterior. So some sampling methods are introduced.

3.19.2 Sampling Methods

Importance sampling

It is common to sample from an easy-to-implement distribution, the so-called proposal distribution $q(x | y)$, hence

$$\begin{aligned} E(f(x)) &= \int f(x_k) \frac{p(x_k | y_{1:k})}{q(x_k | y_{1:k})} q(x_k | y_{1:k}) dx_x \\ &= \int f(x_k) \frac{p(x_k) p(y_{1:k} | x_k)}{p(y_{1:k}) q(x_k | y_{1:k})} q(x_k | y_{1:k}) dx_x \\ &= \int f(x_k) \frac{W_k(x_k)}{p(y_{1:k})} q(x_k | y_{1:k}) dx_x, \end{aligned}$$

where $W_k(x_k) = \frac{p(x_k)p(y_{1:k}|x_k)}{q(x_k|y_{1:k})} \propto \frac{p(x_k|y_{1:k})}{q(x_k|y_{1:k})}$. Because $p(y_{1:k}) = \int p(y_{1:k} | x_k)p(x_k)dx_k$, so the above equation can be rewritten as

$$\begin{aligned} E(f(x)) &= \frac{1}{p(y_{1:k})} \int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k \\ &= \frac{\int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k}{\int p(y_{1:k} | x_k)p(x_k)dx_k} \\ &= \frac{\int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k}{\int W_k(x_k)q(x_k | y_{1:k})dx_k} \\ &= \frac{E_{q(x_k|y_{1:k})}[W_k(x_k)f(x_k)]}{E_{q(x_k|y_{1:k})}[W_k(x_k)]}. \end{aligned}$$

To solve the above equation, we can use Monte Carlo method by drawing samples $\{x_k^{(i)}\}$ from $q(x_k | y_{1:k})$ and get their expectation, which approximate by

$$\begin{aligned} E(f(x_k)) &\approx \frac{\frac{1}{N} \sum_{i=1}^N W_k(x_k^{(i)})f(x_k^{(i)})}{\frac{1}{N} \sum_{i=1}^N W_k(x_k^{(i)})} \\ &= \sum_{i=1}^N \tilde{W}_k(x_k^{(i)})f(x_k^{(i)}), \end{aligned}$$

where $\tilde{W}_k(x_k^{(i)}) = \frac{W_k(x_k^{(i)})}{\sum_{i=1}^N W_k(x_k^{(i)})}$ is factorized weight. Each particles has its own weighted value, so the expectation is a weighted mean. However, the drawback of this method is that the computation is quite expensive. A smarter way is to update $W_k^{(i)}$ recursively. Suppose the proposal distribution

$$q(x_{0:k} | y_{1:k}) = q(x_{0:k-1} | y_{1:k-1})q(x_k | x_{0:k-1}, y_{1:k}),$$

then the recursive form of the posterior distribution is

$$\begin{aligned} p(x_{0:k} | y_{1:k}) &= \frac{p(y_k | x_{0:k}, y_{1:k-1})p(x_{0:k} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_{0:k}, y_{1:k-1})p(x_k | x_{0:k-1}, y_{1:k-1})p(x_{0:k-1} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_k)p(x_k | x_{k-1})p(x_{0:k-1} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &\propto p(y_k | x_k)p(x_k | x_{k-1})p(x_{0:k-1} | y_{1:k-1}), \end{aligned}$$

the recursive form of the weights are

$$\begin{aligned}
W_k^{(i)} &\propto \frac{p(x_{0:k}^{(i)} | y_{1:k})}{q(x_{0:k}^{(i)} | y_{1:k})} \\
&= \frac{p(y_{1:k} | x_{0:k}^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})p(x_{0:k-1}^{(i)} | y_{1:k-1})}{q(x_k^{(i)} | x_{0:k-1}^{(i)}, y_k)q(x_{0:k-1}^{(i)} | y_{1:k-1})} \\
&= W_{k-1}^{(i)} \frac{p(y_{1:k} | x_{0:k}^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})}{q(x_k^{(i)} | x_{0:k-1}^{(i)}, y_k)}.
\end{aligned}$$

Sequential Importance Sampling and Resampling

In practice, we are interested in the current filtered estimate $p(x_k | y_{1:k})$ instead of $p(x_{0:k} | y_{1:k})$. Provided

$$q(x_k | x_{0:k-1}, y_{1:k}) = q(x_k | x_{k-1}, y_k),$$

the importance weights $W_k^{(i)}$ can be updated recursively

$$W_k^{(i)} \propto W_{k-1}^{(i)} \frac{p(y_k | x_k^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})}{q(x_k^{(i)} | x_{k-1}^{(i)}, y_k)}.$$

The problem of SIS filter is that the distribution of importance weights becomes more and more skewed as time increases. Hence, after some iterations, only very few particles have non-zero importance weights. This phenomenon is called *weight degeneracy* or *sample impoverishment* Arnaud Doucet (2011).

The effective sample size N_{eff} is suggested to monitor how bad the degeneration is, which is

$$N_{eff} = \frac{N}{1 + \text{var}(w_k^{*(i)})},$$

where $w_k^{*(i)} = \frac{p(x_k^{(i)} | y_{1:k})}{q(x_k^{(i)} | x_{k-1}^{(i)}, y_{1:k})}$. The more different between the biggest weight and smallest weight, the worse the degeneration is. In practice, the effective sample size is approximated by

$$\hat{N}_{eff} \approx \frac{1}{\sum_{i=1}^N (w_k^{(i)})^2}.$$

If the value of N_{eff} is less than some threshold, some procedure should be used to avoid a worse degeneration. There are two ways one can do: choose an appropriate PDF for importance sampling, or use re-sampling after SIS.

The idea of resampling is keeping the same size of particles, replacing the low weights particles with new ones. As discussed before,

$$p(x_k \mid y_{1:k}) = \sum_{i=1}^N w_k^{(i)} \delta(x_k - x_k^{(i)}).$$

After resampling, it becomes

$$\tilde{p}(x_k \mid y_{1:k}) = \sum_{j=1}^N \frac{1}{N} \delta(x_k - x_k^{(j)}) = \sum_{i=1}^N \frac{n_i}{N} \delta(x_k - x_k^{(i)}),$$

where n_i represents how many times the new particles $x_k^{(j)}$ were duplicated from $x_k^{(i)}$.

Then the process of SIS particle filter with re-sampling is:

- Initial particles when $k = 0$. For $i = 1, \dots, N$, draw samples $\{x_0^{(i)}\}$ from $p(x_0)$.
- For $k = 1, 2, \dots$, run the process recursively
 - Importance sampling: draw sample $\{\tilde{x}_k^{(i)}\}_{i=1}^N$ from $q(x_k \mid y_{1:k})$, calculate their weights $\tilde{w}_k^{(i)}$ and normalize them.
 - Re-sampling: Re-sample $\{\tilde{x}_k^{(i)}, \tilde{w}_k^{(i)}\}$ and get a new set $\{x_k^{(i)}, \frac{1}{N}\}$.
 - Output the status at time k : $\hat{x}_k = \sum_{i=1}^N \tilde{x}_k^{(i)} \tilde{w}_k^{(i)}$.

In SIR, if we choose

$$q(x_k^{(i)} \mid x_{k-1}^{(i)}, y_k) = p(x_k^{(i)} \mid x_{k-1}^{(i)}),$$

the weights become

$$\begin{aligned} w_k^{(i)} &\propto w_{k-1}^{(i)} \frac{p(y_k \mid x_k^{(i)}) p(x_k^{(i)} \mid x_{k-1}^{(i)})}{q(x_k^{(i)} \mid x_{k-1}^{(i)}, y_k)} \\ &\propto w_{k-1}^{(i)} p(y_k \mid x_k^{(i)}). \end{aligned}$$

Because $w_{k-1}^{(i)} = \frac{1}{N}$, thus we have $w_k^{(i)} \propto p(y_k \mid x_k^{(i)})$ and

$$w = \frac{1}{\sqrt{2\pi\Sigma}} \exp\left(-\frac{1}{2}(y_{true} - y)\Sigma^{-1}(y_{true} - y)\right).$$

Metropolis-Hastings Algorithm with Delayed Acceptance

Importance sampling works well only if the proposal density $Q(x)$ is similar to $P(x)$. In large and complex problems it is difficult to create a single density $Q(x)$ that has this property MacKay (2003). Here, we introduce the Metropolis-Hastings algorithm, which makes use of a proposal density $Q(x)$ depending on the current state x_t instead. We assume that we can evaluate $P^*(\theta)$ for any θ . A tentative new state θ' is generated from the proposal density $Q(\theta'; \theta^{(t)})$. To decide whether to accept the new state, we compute the quantity

$$\alpha = \frac{P^*(\theta')}{P^*(\theta^{(t)})} \frac{Q(\theta^{(t)}; \theta')}{Q(\theta'; \theta^{(t)})}.$$

If $\alpha \geq 1$, then the new state is accepted. Otherwise, the new state is accepted with probability α .

In our case, the proposal $\theta' \sim N(\theta^{(t)}, \sigma)$, and the density Q is symmetric, so

$$\alpha = \frac{P^*(\theta')}{P^*(\theta^{(t)})} = \frac{P(y_{1:t} | \theta') P(\theta')}{P(y_{1:t} | \theta^{(t)}) P(\theta^{(t)})}.$$

Delayed acceptance random walk Metropolis Sherlock *et al.* (2016)

1. Standard MH acceptance formula:

$$\alpha_1 = \min \left\{ 1, \frac{\hat{\pi}(\theta^*) q(\theta | \theta^*)}{\hat{\pi}(\theta) q(\theta^* | \theta)} \right\},$$

where $\hat{\pi}(\cdot) = N(\cdot | \hat{\theta}, \sigma)$, $q(\theta^* | \theta) = N(\theta^* | \theta, \epsilon^2)$ and $q(\theta | \theta^*) = N(\theta | \theta^*, \epsilon^2)$.

2. The acceptance probability is

$$\alpha_2 = \min \left\{ 1, \frac{\pi(\theta^*) \hat{\pi}(\theta)}{\pi(\theta) \hat{\pi}(\theta^*)} \right\}.$$

The overall acceptance probability $\alpha_1 \alpha_2$ ensures that detailed balance is satisfied with respect to π ; however if a rejection occurs at Stage 1 then the expensive evaluation of $\pi(\theta)$ at Stage 2 is unnecessary.

3.20 Bayesian Parameter Estimation

The state transition density and the conditional likelihood function depend not only upon the dynamic state x_t , but also on a static parameter vector θ , which will be stressed by use of the notations $f(x_t | x_{t-1}, \theta)$ and $g(y_t | x_t, \theta)$. To estimate θ , we would consider a Bayesian method in the following two situations: off-line, estimating the parameters by a batch of data, and on-line, by an instant updated sequential

data stream. Specifically, the advantage of Bayesian than maximum likelihood method is that the unknown parameter is considered random and assigned a suitable prior distribution, which is addressed from the experiences of researchers or a learning process and easily to be implemented in the algorithm of machine learning.

Generally, in the Bayesian setting, we choose a suitable prior density $p(\theta)$ for θ and compute the joint posterior density $p(x_{0:t}, \theta \mid y_{0:t})$ in the off-line case, or the sequence of posterior densities $\{p(x_{0:n}, \theta \mid y_{0:n})\}$ in the on-line setting Kantas *et al.* (2009).

3.20.1 Off-line Methods

In the off-line setting, the parameters can be estimated with non-sequential Monte Carlo methods, such as Markov Chain Monte Carlo Robert (2004). However, it is recognized that the sequential MC methods have some significant advantages in some certain cases, like Cappé *et al.* (2009) and Del Moral *et al.* (2006). Additionally, it is difficult to design an efficient MCMC sampling algorithm for a nonlinear non-Gaussian state space model. A Particle MCMC method is proposed by Andrieu *et al.* (2010), which is a new class of MCMC techniques relying on Standard MC methods to build efficient high dimensional proposal distributions.

PMMH jointly updates θ and $x_{0:t}$ for state space models. It proposes a new θ^* from a proposal density function $q(\theta^* \mid \theta)$, and then generates $x_{0:t}^*$ by running bootstrap particle filter with θ^* . The acceptance ratio of this sampler is

$$\begin{aligned} \alpha &= \min \left\{ 1, \frac{p(x_{0:t}^*, \theta^* \mid y_{0:t})q((x_{0:t}, \theta) \mid (x_{0:t}^*, \theta^*))}{p(x_{0:t}, \theta \mid y_{0:t})q((x_{0:t}^*, \theta^*) \mid (x_{0:t}, \theta))} \right\} \\ &= \min \left\{ 1, \frac{p_{\theta^*}(y_{0:t})p(\theta^*)q(\theta \mid \theta^*)}{p_{\theta}(y_{0:t})p(\theta)q(\theta^* \mid \theta)} \right\}. \end{aligned}$$

The PMMH sampler is an approximation of the ideal MMH sampler for sampling from $p(x^t, \theta \mid y^t)$. Apparently, the higher number of particles N the better the mixing properties of the algorithm, in contrast, the lower efficiency of computation.

3.20.2 On-line Methods

Putting the algorithms on-line means to update the parameters and states instantly as new observations coming into the data stream. For Bayesian dynamic models, however, the most natural option consists in treating the unknown parameter θ , using the state space representation, as a component of the state which has no dynamic evolution, also referred to as a static parameter Cappé *et al.* (2007).

The standard SMC is deficiency for on-line estimation. As a result of the successive resampling steps, after a certain time n , the approximation $\hat{p}(\theta | y^{1:t})$ will only contain a single unique value for θ . In other words, SMC approximation of the marginalized parameter posterior distribution is represented by a single Dirac delta function. It also causes error accumulation in successive Monte Carlo (MC) steps grows exponentially or polynomially in time.

The target is to estimate $p(\theta | y_{1:t})$ by

$$p(\theta | y_{1:t}) \propto p(y_{1:t} | \theta)p(\theta) \quad (3.150)$$

without introducing any bias or controlling the bias in states propagation. A pragmatic approach to reduce parameter sample degeneracy and error accumulation in successive MC approximations is to adding an artificial dynamic equation on θ Higuchi (2001) Kitagawa (1998), which gives

$$\theta_{n+1} = \theta_n + \varepsilon_{n+1}.$$

With a small artificial noise, SMC can now be applied to approximate $p(x^t, \theta | y^t)$. A related kernel density estimation method proposes a kernel density estimate of the target Liu and West (2001)

$$\hat{p}(\theta | y^t) = \frac{1}{N} \sum M(\theta - \theta_n^{(i)}).$$

Both of these methods require a significant amount of tuning.

A fixed-lag practical filtering is used to approximate

$$p(x_{0:n-L}, \theta | y_{0:n-1}) \approx p(x_{0:n-L}, \theta | y^n)$$

for L large enough in reference Polson *et al.* (2008). $x_{0:n-L}$ has very little influence on observations coming after n . The choice of the lag L is difficult and ther is a non-vanishing bias which is difficult to quantify.

A MCMC kernel with invariant density $p(x^t, \theta | y^t)$ is used in SMC algorithm. This method was firstly used in an on-line Bayesian parameter estimation, where the author in Andrieu *et al.* (1999) were using

$$K_n(x'_{1:t}, \theta' | x_{1:t}, \theta) = \delta_{x_{1:t}}(x'_{1:t})p(\theta' | x_{1:t}, y_{1:t}),$$

where $p(y^t | \theta, x^t) = p(\theta | s_t(x^t, y^t))$ and $s_t(x^t, y^t)$ is a fixed-dimensional vector of sufficient statistics. MCMC can be used to maintain the diversity of the samples of θ . Here the stationary distribution for the MCMC will be the full joint posterior distribution of states and parameters and apply MH or Gibbs sampling separately to $p(\theta | x^t, y^t)$ and $p(x^t | \theta, y^t)$. However, this method is not feasible for large dataset.

3.21 A Sequential Monte Carlo Algorithm for Parameter Estimation

Several methods can be used for parameter estimation, cross validation, EM method, Gibbs sampling and Metropolis-Hastings algorithm.

The process is

$$\begin{aligned} y_t \mid x_t &\sim N(x_t, \sigma^2) \\ x_t \mid x_{t-1} &\sim N(\phi x_{t-1}, \tau^2). \end{aligned}$$

The joint distribution of $x_{0:t}$ and $y_{1:t}$ is

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N(0, \Sigma),$$

where Σ^{-1} is

$$\begin{bmatrix} \frac{1}{L^2} + \frac{\phi^2}{\tau^2} & \frac{-\phi}{\tau^2} & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \frac{-\phi}{\tau^2} & \frac{1+\phi^2}{\tau^2} + \frac{1}{\sigma^2} & \cdots & 0 & -\frac{1}{\sigma^2} & 0 & \cdots & 0 \\ 0 & \frac{-\phi}{\tau^2} & \cdots & 0 & 0 & -\frac{1}{\sigma^2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\tau^2} + \frac{1}{\sigma^2} & 0 & 0 & \cdots & -\frac{1}{\sigma^2} \\ 0 & -\frac{1}{\sigma^2} & \cdots & 0 & \frac{1}{\sigma^2} & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \frac{1}{\sigma^2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -\frac{1}{\sigma^2} & 0 & 0 & \cdots & \frac{1}{\sigma^2} \end{bmatrix}.$$

The block matrix $\Sigma^{-1} = \begin{bmatrix} A(\theta) & -\frac{1}{\sigma^2}I \\ -\frac{1}{\sigma^2}I & \frac{1}{\sigma^2}I \end{bmatrix}$ and

$$\Sigma = \begin{bmatrix} (A - B)^{-1} & (A - B)^{-1} \\ (A - B)^{-1} & (I - A^{-1}B)^{-1}B^{-1} \end{bmatrix},$$

where $\theta = \{\phi, \tau, \sigma\}$ and $B = \frac{1}{\sigma^2}I$. The posterior of ϕ is

$$\begin{aligned} P(\theta \mid y_{1:t}) &= \frac{P(y_{1:t} \mid \theta)P(\theta)}{P(y_{1:t})} \propto P(y_{1:t} \mid \theta)P(\theta) \\ &= e^{-\frac{1}{2}y^\top B(I - A^{-1}B)y} \sqrt{\det B(I - A^{-1}B)} P(\theta). \end{aligned}$$

Because $\det B(I - A^{-1}B) = \det(B) \det(A^{-1}) \det(A - B) = \frac{1}{\sigma^{2n} \tau^{2n} \det A}$. Taking the Cholesky decomposition of $A = LL^\top$ and natural logarithm of the posterior will give us

$$\ln P(\theta \mid y_{1:t}) = \frac{1}{2\sigma^4} u^\top u - \frac{1}{2\sigma^2} y^\top y - n \ln \sigma - n \ln \tau - \sum \ln \text{tr}(L) - \frac{n}{2} \ln(2\pi).$$

3.22 States Estimation

4.23 State Space Models

State space models are the natural form of system models relying on the general concept of state. If we describe a system as an operator mapping from the space of inputs to the space of outputs, then we may need the entire input-output history of the system together with the planned input in order to compute the future output values Hangos *et al.* (2006). In an alternative way, by using new information at time t containing all the past information up to the current state and initial conditions to get the current output is possible, that is known as a sequential method. A genetic state space model consists of two sets of equations: state equation and output equation. The state equation describes the evolution of the true input and state variables sequentially as a function and passes the variable one after one, generally, with some noises. The output equation catches the input values and interprets it out by an algebraic equation. A general state space model looks like the following form

$$\text{State equation } x_t = G_t(x_{t-1}) + w_t, \quad (4.151)$$

$$\text{Output equation } y_t = F_t(x_t) + \epsilon_t \quad (4.152)$$

with an initial state x_0 , where ϵ_t and w_t are noises passing through the process G_t and F_t . x_t are true status variables and y_t are output values. Many researchers have been interested in this model and its application because of its good property. It can be used to model univariate or multivariate time series, also in the presence of non-stationarity, structural changes, and irregular patterns Petris *et al.* (2009).

The most simple and important system is given by Gaussian linear state space models, also known by dynamic linear models (DLM), which defines a very general class of non-stationary time series models. Firstly, the model is linear, which means G_t and F_t are linear processes and satisfying linearity property. Secondly, the it is specified by a normal prior distribution for the p -dimensional state vector at initial state $t = 0$,

$$x_0 \sim N_p(m, 0, C_0)$$

and two independent zero mean normal distributed noises $\epsilon_t \sim N_p(0, V_t)$ and $w_t \sim N_p(0, W_t)$ Petris *et al.* (2009). The well known Kalman Filter is a particular algorithm that is used to solve state space models in the linear case. This was first derived by Kalman Kalman *et al.* (1960).

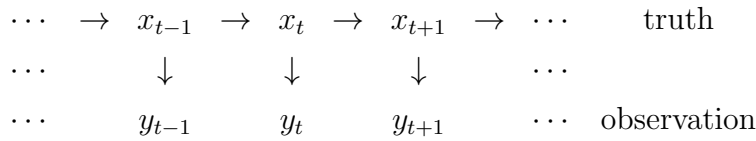
In a nonlinear state space model, the process G_t and F_t are no longer linear functions and the situation becomes more complicated. Here gives a simple nonlinear example

of such a model, which has been used extensively in the literature for benchmarking numerical filtering techniques Kitagawa (1996) West (1993) Gordon *et al.* (1993) assuming the sequence is Markovian.

$$x_t = \frac{x_{t-1}}{2} + 25 \frac{x_{t-1}}{1 + x_{t-1}^2} + 8 \cos(1.2t) + u_t$$

$$y_t = \frac{x_t^2}{20} + v_t,$$

where $u_t \sim N(0, \sigma_u^2)$, $v_t \sim N(0, \sigma_v^2)$, $\sigma_u^2 = 10$ and $\sigma_v^2 = 1$ are considered fixed and known. The initial state $x_0 \sim N(0, 10)$. The assumption Markovian keeps the current state x_t only depending on the previous one step x_{t-1} and the observed y_t depending on x_t . A state-space is shown in the diagram below:



In applications, the process function G_t and F_t contain unknown parameters to be estimated De Jong (1988) and the target is to estimate the true states on sequential observations y_1, \dots, y_t . Then it becomes to estimate a joint density of $p(x_{1:t}, \theta \mid y_{1:t})$, where $x_{1:t} = \{x_1, x_2, \dots, x_t\}$ are the hidden states and $y_{1:t} = \{y_1, y_2, \dots, y_t\}$ are the observed outcomes and θ is a set of unknown parameters.

4.24 Sequential Monte Carlo Method

The use of Monte Carlo methods for non-linear filtering can be traced back to the pioneering contributions of Handschin and Mayne (1969) Handschin and Mayne (1969) and Handschin (1970) Handschin (1970). These researchers tried to use an importance sampling paradigm to approximate the target distributions. Later on, an importance sampling algorithms were implemented sequentially in the non-linear filtering context. This algorithm is called sequential importance sampling, often abbreviated SIS, and has been known since the early 1970s. Limited by the power of computers and suffering from sample impoverishment or weight degeneracy, the SIS didn't develop very well until 1993. A particle filter algorithm was proposed to allow rejuvenation of the set of samples by duplicating the samples with high importance weights and, on the contrary, removing samples with low weights Cappé *et al.* (2009). Since then, sequential Monte Carlo (SMC) methods have been applied in many different fields including but not limited to computer vision, signal processing, control, econometrics, finance, robotics, and statistics Arnaud Doucet (2011) Ristic *et al.* (2004).

4.24.1 Filtering Problem and Estimation

Sequential Monte Carlo method, also known as particle filter, is a technique based on sampling and importance sampling methods to find the best state estimation given by Gordon in 1993 Gordon *et al.* (1993) and was the first successful application of sequential Monte Carlo techniques to the field of non-linear filtering Cappé *et al.* (2009). In the state space model, a generic particle filter estimates the posterior distribution of the hidden states using the observation measurement process. The filtering problem is to estimate sequentially the values of the hidden states x_k given the values of the observation process $y_{1:k}$ at any time step k . In another word, it is to find the value of $p(x_k | y_{1:k})$. The process is divided into two steps: prediction and updating. In the prediction step, the assumption of Markov Chain is the current status x_k only depends on the previous one x_{k-1} . Then we can calculate the probability of x_k by

$$\begin{aligned} p(x_k | y_{1:k-1}) &= \int p(x_k, x_{k-1} | y_{1:k-1}) dx_{k-1} \\ &= \int p(x_k | x_{k-1}, y_{1:k-1}) p(x_{k-1} | y_{1:k-1}) dx_{k-1} \\ &= \int p(x_k | x_{k-1}) p(x_{k-1} | y_{1:k-1}) dx_{k-1}. \end{aligned}$$

In the updating step, once $p(x_k | y_{1:k-1})$ is known, $p(x_k | y_{1:k})$ can be found by

$$\begin{aligned} p(x_k | y_{1:k}) &= \frac{p(y_k | x_k, y_{1:k-1}) p(x_k | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_k) p(x_k | y_{1:k-1})}{p(y_k | y_{1:k-1})}, \end{aligned}$$

where the normalization $p(y_k | y_{1:k-1}) = \int p(y_k | x_k) p(x_k | y_{1:k-1}) dx_k$ Arulampalam *et al.* (2002).

Imagine that the stat space is partitioned as many parts, in which the particles are filled according to some probability measure. The higher probability, the denser the particles are concentrated. Suppose the particles x_1, \dots, x_N are drawn from the target probability density function $p(x)$, then these particles are used to estimate the expectation and variance of $f(x)$ by

$$\begin{aligned} E(f(x)) &= \int_a^b f(x) p(x) dx \\ Var(f(x)) &= E(f(x) - E(f(x)))^2 p(x) dx. \end{aligned}$$

Back to our target, the posterior distribution or density is empirically represented by

a weighted sum of samples x_1, \dots, x_N

$$\hat{p}(x_n | y_{1:k}) = \frac{1}{N} \sum_{i=1}^N \delta(x_n - x_n^{(i)}) \approx p(x_n | y_{1:k}),$$

where $f(x) = \delta(x_n - x_n^{(i)})$ is Dirac delta function. When N is sufficiently large, $\hat{p}(x_n | y_{1:k})$ approximates the true posterior $p(x_n | y_{1:k})$. By this approximation, the filtering problem becomes to get the expectation of current status

$$\begin{aligned} E(f(x_n)) &\approx \int f(x_n) \hat{p}(x_n | y_{1:k}) dx_n \\ &= \frac{1}{N} \sum_{i=1}^N \int f(x_n) \delta(x_n - x_n^{(i)}) dx_n \\ &= \frac{1}{N} \sum_{i=1}^N f(x_n^{(i)}). \end{aligned}$$

The expectation is the mean of the status of all particles x_1, \dots, x_N .

However, the posterior distribution is unknown and impossible to sample from the true posterior. So some sampling methods are introduced.

4.24.2 Sampling Methods

Importance sampling

It is common to sample from an easy-to-implement distribution, the so-called proposal distribution $q(x | y)$, hence

$$\begin{aligned} E(f(x)) &= \int f(x_k) \frac{p(x_k | y_{1:k})}{q(x_k | y_{1:k})} q(x_k | y_{1:k}) dx_x \\ &= \int f(x_k) \frac{p(x_k) p(y_{1:k} | x_k)}{p(y_{1:k}) q(x_k | y_{1:k})} q(x_k | y_{1:k}) dx_x \\ &= \int f(x_k) \frac{W_k(x_k)}{p(y_{1:k})} q(x_k | y_{1:k}) dx_x, \end{aligned}$$

where $W_k(x_k) = \frac{p(x_k)p(y_{1:k}|x_k)}{q(x_k|y_{1:k})} \propto \frac{p(x_k|y_{1:k})}{q(x_k|y_{1:k})}$. Because $p(y_{1:k}) = \int p(y_{1:k} | x_k)p(x_k)dx_k$, so the above equation can be rewritten as

$$\begin{aligned} E(f(x)) &= \frac{1}{p(y_{1:k})} \int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k \\ &= \frac{\int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k}{\int p(y_{1:k} | x_k)p(x_k)dx_k} \\ &= \frac{\int f(x_k)W_k(x_k)q(x_k | y_{1:k})dx_k}{\int W_k(x_k)q(x_k | y_{1:k})dx_k} \\ &= \frac{E_{q(x_k|y_{1:k})}[W_k(x_k)f(x_k)]}{E_{q(x_k|y_{1:k})}[W_k(x_k)]}. \end{aligned}$$

To solve the above equation, we can use Monte Carlo method by drawing samples $\{x_k^{(i)}\}$ from $q(x_k | y_{1:k})$ and get their expectation, which approximate by

$$\begin{aligned} E(f(x_k)) &\approx \frac{\frac{1}{N} \sum_{i=1}^N W_k(x_k^{(i)})f(x_k^{(i)})}{\frac{1}{N} \sum_{i=1}^N W_k(x_k^{(i)})} \\ &= \sum_{i=1}^N \tilde{W}_k(x_k^{(i)})f(x_k^{(i)}), \end{aligned}$$

where $\tilde{W}_k(x_k^{(i)}) = \frac{W_k(x_k^{(i)})}{\sum_{i=1}^N W_k(x_k^{(i)})}$ is factorized weight. Each particles has its own weighted value, so the expectation is a weighted mean. However, the drawback of this method is that the computation is quite expensive. A smarter way is to update $W_k^{(i)}$ recursively. Suppose the proposal distribution

$$q(x_{0:k} | y_{1:k}) = q(x_{0:k-1} | y_{1:k-1})q(x_k | x_{0:k-1}, y_{1:k}),$$

then the recursive form of the posterior distribution is

$$\begin{aligned} p(x_{0:k} | y_{1:k}) &= \frac{p(y_k | x_{0:k}, y_{1:k-1})p(x_{0:k} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_{0:k}, y_{1:k-1})p(x_k | x_{0:k-1}, y_{1:k-1})p(x_{0:k-1} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &= \frac{p(y_k | x_k)p(x_k | x_{k-1})p(x_{0:k-1} | y_{1:k-1})}{p(y_k | y_{1:k-1})} \\ &\propto p(y_k | x_k)p(x_k | x_{k-1})p(x_{0:k-1} | y_{1:k-1}), \end{aligned}$$

the recursive form of the weights are

$$\begin{aligned}
W_k^{(i)} &\propto \frac{p(x_{0:k}^{(i)} | y_{1:k})}{q(x_{0:k}^{(i)} | y_{1:k})} \\
&= \frac{p(y_{1:k} | x_{0:k}^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})p(x_{0:k-1}^{(i)} | y_{1:k-1})}{q(x_k^{(i)} | x_{0:k-1}^{(i)}, y_k)q(x_{0:k-1}^{(i)} | y_{1:k-1})} \\
&= W_{k-1}^{(i)} \frac{p(y_{1:k} | x_{0:k}^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})}{q(x_k^{(i)} | x_{0:k-1}^{(i)}, y_k)}.
\end{aligned}$$

Sequential Importance Sampling and Resampling

In practice, we are interested in the current filtered estimate $p(x_k | y_{1:k})$ instead of $p(x_{0:k} | y_{1:k})$. Provided

$$q(x_k | x_{0:k-1}, y_{1:k}) = q(x_k | x_{k-1}, y_k),$$

the importance weights $W_k^{(i)}$ can be updated recursively

$$W_k^{(i)} \propto W_{k-1}^{(i)} \frac{p(y_k | x_k^{(i)})p(x_k^{(i)} | x_{k-1}^{(i)})}{q(x_k^{(i)} | x_{k-1}^{(i)}, y_k)}.$$

The problem of SIS filter is that the distribution of importance weights becomes more and more skewed as time increases. Hence, after some iterations, only very few particles have non-zero importance weights. This phenomenon is called *weight degeneracy* or *sample impoverishment* Arnaud Doucet (2011).

The effective sample size N_{eff} is suggested to monitor how bad the degeneration is, which is

$$N_{eff} = \frac{N}{1 + var(w_k^{*(i)})},$$

where $w_k^{*(i)} = \frac{p(x_k^{(i)} | y_{1:k})}{q(x_k^{(i)} | x_{k-1}^{(i)}, y_{1:k})}$. The more different between the biggest weight and smallest weight, the worse the degeneration is. In practice, the effective sample size is approximated by

$$\hat{N}_{eff} \approx \frac{1}{\sum_{i=1}^N (w_k^{(i)})^2}.$$

If the value of N_{eff} is less than some threshold, some procedure should be used to avoid a worse degeneration. There are two ways one can do: choose an appropriate PDF for importance sampling, or use re-sampling after SIS.

The idea of resampling is keeping the same size of particles, replacing the low weights particles with new ones. As discussed before,

$$p(x_k \mid y_{1:k}) = \sum_{i=1}^N w_k^{(i)} \delta(x_k - x_k^{(i)}).$$

After resampling, it becomes

$$\tilde{p}(x_k \mid y_{1:k}) = \sum_{j=1}^N \frac{1}{N} \delta(x_k - x_k^{(j)}) = \sum_{i=1}^N \frac{n_i}{N} \delta(x_k - x_k^{(i)}),$$

where n_i represents how many times the new particles $x_k^{(j)}$ were duplicated from $x_k^{(i)}$.

Then the process of SIS particle filter with re-sampling is:

- Initial particles when $k = 0$. For $i = 1, \dots, N$, draw samples $\{x_0^{(i)}\}$ from $p(x_0)$.
- For $k = 1, 2, \dots$, run the process recursively
 - Importance sampling: draw sample $\{\tilde{x}_k^{(i)}\}_{i=1}^N$ from $q(x_k \mid y_{1:k})$, calculate their weights $\tilde{w}_k^{(i)}$ and normalize them.
 - Re-sampling: Re-sample $\{\tilde{x}_k^{(i)}, \tilde{w}_k^{(i)}\}$ and get a new set $\{x_k^{(i)}, \frac{1}{N}\}$.
 - Output the status at time k : $\hat{x}_k = \sum_{i=1}^N \tilde{x}_k^{(i)} \tilde{w}_k^{(i)}$.

In SIR, if we choose

$$q(x_k^{(i)} \mid x_{k-1}^{(i)}, y_k) = p(x_k^{(i)} \mid x_{k-1}^{(i)}),$$

the weights become

$$\begin{aligned} w_k^{(i)} &\propto w_{k-1}^{(i)} \frac{p(y_k \mid x_k^{(i)}) p(x_k^{(i)} \mid x_{k-1}^{(i)})}{q(x_k^{(i)} \mid x_{k-1}^{(i)}, y_k)} \\ &\propto w_{k-1}^{(i)} p(y_k \mid x_k^{(i)}). \end{aligned}$$

Because $w_{k-1}^{(i)} = \frac{1}{N}$, thus we have $w_k^{(i)} \propto p(y_k \mid x_k^{(i)})$ and

$$w = \frac{1}{\sqrt{2\pi\Sigma}} \exp\left(-\frac{1}{2}(y_{true} - y)\Sigma^{-1}(y_{true} - y)\right).$$

Delayed Acceptance Metropolis-Hastings Algorithm

Importance sampling works well only if the proposal density $q(x)$ is similar to $p(x)$. In large and complex problems it is difficult to create a single density $q(x)$ that has this property MacKay (2003). Here, we introduce the Metropolis-Hastings algorithm, which makes use of a proposal density $q(x)$ depending on the current state x_t instead. We assume that we can evaluate $p(\theta)$ for any θ . The transition probabilities should satisfy the detailed balance condition

$$\pi(\theta)p(\theta' | \theta) = \pi(\theta')p(\theta | \theta'),$$

that means transition from $\pi(\theta)$ to $\pi(\theta')$ has the same probability as that from $\pi(\theta')$ to $\pi(\theta)$. In sampling method, drawing θ_i first and then drawing θ_j should have the same probability as drawing θ_j and then drawing θ_i . However, in most situations, the detailed balance condition is not satisfied. Therefore, we introduce a function $\alpha(x, y)$ satisfying

$$p(\theta_i)q(\theta_i, \theta_j)\alpha(\theta_i, \theta_j) = p(\theta_j)q(\theta_j, \theta_i)\alpha(\theta_j, \theta_i).$$

A tentative new state θ' is generated from the proposal density $q(\theta'; \theta^{(t)})$. To decide whether to accept the new state, we compute the quantity

$$\alpha = \frac{p(\theta')}{p(\theta^{(t)})} \frac{q(\theta^{(t)}; \theta')}{q(\theta'; \theta^{(t)})}.$$

If $\alpha \geq 1$, then the new state is accepted. Otherwise, the new state is accepted with probability α . A drawback of MH algorithm is a large time consuming in calculating $p(\theta)$ if it's in an irregular structure. A delayed acceptance MH algorithm introduces a cheap approximation $\hat{p}(\theta)$ for $p(\theta)$ in two stages. In stage one, the quantity α_1 is found by a standard MH acceptance formula

$$\alpha_1 = \min \left\{ 1, \frac{\hat{p}(\theta^*)q(\theta | \theta^*)}{\hat{p}(\theta)q(\theta^* | \theta)} \right\}, \quad (4.153)$$

where $\hat{p}(\cdot)$ is a cheap estimation for θ and a simple form is $\hat{p}(\cdot) = N(\cdot | \hat{\theta}, \sigma)$. Once α_1 is accepted, the process goes into stage two and the acceptance probability α_2 is

$$\alpha_2 = \min \left\{ 1, \frac{p(\theta^*)\hat{p}(\theta)}{p(\theta)\hat{p}(\theta^*)} \right\}, \quad (4.154)$$

where the overall acceptance probability $\alpha_1\alpha_2$ ensures that detailed balance is satisfied with respect to $p(\cdot)$; however if a rejection occurs at stage one then the expensive evaluation of $p(\theta)$ at stage two is unnecessary.

In a random walk, the proposal density function $q(\cdot)$ can be chosen for some suitable normal distribution, and hence $q(\theta^* | \theta) = N(\theta^* | \theta, \epsilon^2)$ and $q(\theta | \theta^*) = N(\theta | \theta^*, \epsilon^2)$ cancel in the Delayed Acceptance MH process stage one Sherlock *et al.* (2016). Then in our case, the proposal $\theta' \sim N(\theta^{(t)}, \sigma)$ and the density q is symmetric, so it becomes

$$\alpha = \frac{p^*(\theta')}{p^*(\theta^{(t)})} = \frac{p(y_{1:t} | \theta')p(\theta')}{p(y_{1:t} | \theta^{(t)})p(\theta^{(t)})}. \quad (4.155)$$

4.25 Bayesian Parameter Estimation

The state transition density and the conditional likelihood function depend not only upon the dynamic state x_t , but also on a static parameter vector θ , which will be stressed by use of the notations $f(x_t | x_{t-1}, \theta)$ and $g(y_t | x_t, \theta)$. To estimate θ , we would consider a Bayesian method in the following two situations: off-line, estimating the parameters by a batch of data, and on-line, by an instant updated sequential data stream. Specifically, the advantage of Bayesian than maximum likelihood method is that the unknown parameter is considered random and assigned a suitable prior distribution, which is addressed from the experiences of researchers or a learning process and easily to be implemented in the algorithm of machine learning.

Generally, in the Bayesian setting, we choose a suitable prior density $p(\theta)$ for θ and compute the joint posterior density $p(x_{0:t}, \theta | y_{0:t})$ in the off-line case, or the sequence of posterior densities $\{p(x_{0:n}, \theta | y_{0:n})\}$ in the on-line setting Kantas *et al.* (2009).

4.25.1 Off-line Methods

In the off-line setting, the parameters can be estimated with non-sequential Monte Carlo methods, such as Markov Chain Monte Carlo Robert (2004). However, it is recognized that the sequential MC methods have some significant advantages in some certain cases, like Cappé *et al.* (2009) and Del Moral *et al.* (2006). Additionally, it is difficult to design an efficient MCMC sampling algorithm for a nonlinear non-Gaussian state space model. A Particle MCMC method is proposed by Andrieu *et al.* (2010), which is a new class of MCMC techniques relying on Standard MC methods to build efficient high dimensional proposal distributions.

PMMH jointly updates θ and $x_{0:t}$ for state space models. It proposes a new θ^* from a proposal density function $q(\theta^* | \theta)$, and then generates $x_{0:t}^*$ by running bootstrap

particle filter with θ^* . The acceptance ratio of this sampler is

$$\begin{aligned}\alpha &= \min \left\{ 1, \frac{p(x_{0:t}^*, \theta^* | y_{0:t})q((x_{0:t}, \theta) | (x_{0:t}^*, \theta^*))}{p(x_{0:t}, \theta | y_{0:t})q((x_{0:t}^*, \theta^*) | (x_{0:t}, \theta))} \right\} \\ &= \min \left\{ 1, \frac{p_{\theta^*}(y_{0:t})p(\theta^*)q(\theta | \theta^*)}{p_{\theta}(y_{0:t})p(\theta)q(\theta^* | \theta)} \right\}.\end{aligned}$$

The PMMH sampler is an approximation of the ideal MMH sampler for sampling from $p(x^t, \theta | y^t)$. Apparently, the higher number of particles N the better the mixing properties of the algorithm, in contrast, the lower efficiency of computation.

4.25.2 On-line Methods

Putting the algorithms on-line means to update the parameters and states instantly as new observations coming into the data stream. For Bayesian dynamic models, however, the most natural option consists in treating the unknown parameter θ , using the state space representation, as a component of the state which has no dynamic evolution, also referred to as a static parameter Cappé *et al.* (2007).

The standard SMC is deficiency for on-line estimation. As a result of the successive resampling steps, after a certain time n , the approximation $\hat{p}(\theta | y^{1:t})$ will only contain a single unique value for θ . In other words, SMC approximation of the marginalized parameter posterior distribution is represented by a single Dirac delta function. It also causes error accumulation in successive Monte Carlo (MC) steps grows exponentially or polynomially in time.

The target is to estimate $p(\theta | y_{1:t})$ by

$$p(\theta | y_{1:t}) \propto p(y_{1:t} | \theta)p(\theta) \quad (4.156)$$

without introducing any bias or controlling the bias in states propagation. A pragmatic approach to reduce parameter sample degeneracy and error accumulation in successive MC approximations is to adding an artificial dynamic equation on θ Higuchi (2001) Kitagawa (1998), which gives

$$\theta_{n+1} = \theta_n + \varepsilon_{n+1}.$$

With a small artificial noise, SMC can now be applied to approximate $p(x^t, \theta | y^t)$. A related kernel density estimation method proposes a kernel density estimate of the target Liu and West (2001)

$$\hat{p}(\theta | y^t) = \frac{1}{N} \sum M(\theta - \theta_n^{(i)}).$$

Both of these methods require a significant amount of tuning.

A fixed-lag practical filtering is used to approximate

$$p(x_{0:n-L}, \theta \mid y_{0:n-1}) \approx p(x_{0:n-L}, \theta \mid y^n)$$

for L large enough in reference Polson *et al.* (2008). $x_{0:n-L}$ has very little influence on observations coming after n . The choice of the lag L is difficult and there is a non-vanishing bias which is difficult to quantify.

A MCMC kernel with invariant density $p(x^t, \theta \mid y^t)$ is used in SMC algorithm. This method was firstly used in an on-line Bayesian parameter estimation, where the author in Andrieu *et al.* (1999) were using

$$K_n(x'_{1:t}, \theta' \mid x_{1:t}, \theta) = \delta_{x_{1:t}}(x'_{1:t})p(\theta' \mid x_{1:t}, y_{1:t}),$$

where $p(y^t \mid \theta, x^t) = p(\theta \mid s_t(x^t, y^t))$ and $s_t(x^t, y^t)$ is a fixed-dimensional vector of sufficient statistics. MCMC can be used to maintain the diversity of the samples of θ . Here the stationary distribution for the MCMC will be the full joint posterior distribution of states and parameters and apply MH or Gibbs sampling separately to $p(\theta \mid x^t, y^t)$ and $p(x^t \mid \theta, y^t)$. However, this method is not feasible for large dataset.

4.26 Combined State and Parameters Estimation of Sequential Monte Carlo Algorithm

To work out the best estimations for x and u , one has to solve the target function

$$p(X \mid Y) = \int p(X \mid Y, \theta)p(\theta \mid Y)d\theta. \quad (4.157)$$

The main work need to be done is finding an efficient way to sort out the integration in the above equation. Several methods can be used, such as cross validation, Expectation Maximization algorithm, Gibbs sampling and Metropolis-Hastings algorithm and so on. A Monte Carlo method is popular in research area solving this problem. Monte Carlo method is an algorithm that relies on repeated random sampling to obtain numerical results. To compute an integration of $\int f(x)dx$, one has to sampling as many independent x_i ($i = 1, \dots, N$) as possible and numerically to find $\frac{1}{N} \sum_i f(x_i)$ to approximate the target function.

In our target function, we have to sampling θ and use a numerical way to calculate

its integration. Here are two ways of solving the sampling problem sequentially:

$$\begin{cases} \text{M1 : } p(\theta \mid Y_{1:t}, Y_{t+1}) \propto p(Y_{1:t}, Y_{t+1} \mid \theta)p(\theta) \\ \text{M2 : } p(\theta \mid Y_{1:t}, Y_{t+1}) \propto p(Y_{t+1} \mid \theta, Y_{1:t})p(\theta \mid Y_{1:t}) \end{cases}.$$

NOTES: add more.....

4.26.1 General Linear Space

In one dimensional state space model, we consider the hidden state process $\{x_t, t \geq 1\}$ is a stationary and ergodic Markov process and transited by $f(x' \mid x)$. In this paper, we assume that the current state x_t only depends on the previous one step x_{t-1} , which is known as *AR(1)* model. As indicated by its name, the states are not observed directly but by another process $\{y_t, t \geq 1\}$, which is assumed depending on $\{x_t\}$ by the process $g(y \mid x)$ only and independent with each other. If the transition processes f and g are linear and normal distributed, we call this model *Linear Gaussian Model*, that can be written as

$$\begin{aligned} y_t \mid x_t &\sim N(\gamma x_t, \sigma^2) \\ x_t \mid x_{t-1} &\sim N(\phi x_{t-1}, \tau^2), \end{aligned}$$

where σ and τ are errors occurring in processes, γ and ϕ are static process parameters.

In a simple scenario, by assuming $\gamma = 1$ gives us the joint distribution for $x_{0:t}$ and $y_{1:t}$ as following

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N(0, \Sigma),$$

where Σ^{-1} looks like

$$\begin{bmatrix} \frac{1}{L^2} + \frac{\phi^2}{\tau^2} & \frac{-\phi}{\tau^2} & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \frac{-\phi}{\tau^2} & \frac{1+\phi^2}{\tau^2} + \frac{1}{\sigma^2} & \cdots & 0 & -\frac{1}{\sigma^2} & 0 & \cdots & 0 \\ 0 & \frac{-\phi}{\tau^2} & \cdots & 0 & 0 & -\frac{1}{\sigma^2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\tau^2} + \frac{1}{\sigma^2} & 0 & 0 & \cdots & -\frac{1}{\sigma^2} \\ 0 & -\frac{1}{\sigma^2} & \cdots & 0 & \frac{1}{\sigma^2} & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \frac{1}{\sigma^2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -\frac{1}{\sigma^2} & 0 & 0 & \cdots & \frac{1}{\sigma^2} \end{bmatrix},$$

is the general procedure matrix denoted as $\Sigma^{-1} = \begin{bmatrix} A & -B \\ -B & B \end{bmatrix}$. Its inverse is

$$\Sigma = \begin{bmatrix} (A - B)^{-1} & (A - B)^{-1} \\ (A - B)^{-1} & (I - A^{-1}B)^{-1}B^{-1} \end{bmatrix} \triangleq \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix} \quad (4.158)$$

the covariance matrix, where B is a $t \times t$ diagonal matrix with elements $\frac{1}{\sigma^2}$. The covariance matrices $\Sigma_{XX} = (A - B)^{-1}$ and $\Sigma_{YY} = (I - A^{-1}B)^{-1}B^{-1}$ are easily found. Here the parameter θ represents for the unknown parameters ϕ, σ, τ .

To find a recursive way of calculating the log likelihood posteriors, we introduce the Sherman-Morrison-Woodbury formula here first. In the late 1940s and the 1950s, Sherman and Morrison (1950), Woodbury (1950), Bartlett (1951) and Bodewig (1959) discovered the following result. The original Sherman-Morrison-Woodbury (for short SMW) formula has been used to consider the inverse of matrices Deng (2011). In this paper, we will consider the more generalized case.

Theorem 1.1 (Sherman-Morrison-Woodbury). Let $A \in B(H)$ and $G \in B(K)$ both be invertible, and $Y, Z \in B(K, H)$. Then $A + YGZ^*$ is invertible if and only if $G^{-1} + ZA^{-1}Y$ is invertible. In which case,

$$(A + YGZ^*)^{-1} = A^{-1} - A^{-1}Y(G^{-1} + ZA^{-1}Y)^{-1}ZA^{-1}. \quad (4.159)$$

A simple form of SMW formula is Sherman-Morrison formula represented in the following statement Bartlett (1951): Suppose $A \in R^{n \times n}$ is an invertible square matrix and $u, v \in R^n$ are column vectors. Then $A + uv^\top$ is invertible $\iff 1 + u^\top A^{-1}v \neq 0$. If $A + uv^\top$ is invertible, then its inverse is given by

$$(A + uv^\top)^{-1} = A^{-1} - \frac{A^{-1}uv^\top A^{-1}}{1 + v^\top A^{-1}u}. \quad (4.160)$$

The Forecast Distribution $p(y_{t+1} \mid y_{1:t}, \theta)$

The joint distribution for y_{t+1} and $y_{1:t}$ is $p(y_{1:t+1} \mid \theta) \sim N(0, \Sigma_{YY})$, where $\Sigma_{YY} = (I - A^{-1}B)^{-1}B^{-1}$ is the covariance matrix given above. One may find the inverse of the covariance matrix

$$\Sigma_{YY}^{-1} = B(I - A^{-1}B) = \frac{1}{\sigma^4}(\sigma^2 I - A^{-1}) \triangleq \frac{1}{\sigma^4} \begin{bmatrix} Z_{t+1} & b_{t+1} \\ b_{t+1}^\top & K_{t+1} \end{bmatrix}.$$

Therefore, the original form of this covariance is

$$\Sigma_{YY} = \sigma^4 \begin{bmatrix} (Z - bK^{-1}b^\top)^{-1} & -Z^{-1}b(K - b^\top Z^{-1}b)^{-1} \\ -K^{-1}b^\top(Z - bK^{-1}b^\top)^{-1} & (K - b^\top Z^{-1}b)^{-1} \end{bmatrix}.$$

For sake of simplicity, here we are using Z to represent the $t \times t$ matrix Z_{t+1} , b to represent the $t \times 1$ vector b_{t+1} and K to represent the 1×1 constant K_{t+1} .

By denoting $C_{t+1} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$ and multiplying Σ_{YY}^{-1} , it gives us

$$\Sigma_{YY}^{-1}C_{t+1} = \frac{1}{\sigma^4}(\sigma^2 I - A^{-1})C_{t+1} = \frac{1}{\sigma^4} \begin{bmatrix} b_{t+1} \\ K_{t+1} \end{bmatrix}.$$

In order to find b and K easily, one has to use Sherman-Morrison formula in the following way, that

$$A_{t+1}^{-1}C_{t+1} = \left(I - \frac{M_{t+1}^{-1}u_{t+1}u_{t+1}^\top}{1 + u_{t+1}^\top M_{t+1}^{-1}u_{t+1}} \right) M_{t+1}^{-1}C_{t+1}, \quad (4.161)$$

in which

$$M_{t+1}^{-1}C_{t+1} = \begin{bmatrix} A_t^{-1} & 0 \\ 0 & \sigma^2 \end{bmatrix} C_{t+1} = \sigma^2 C_{t+1},$$

$$u_{t+1}^\top C_{t+1} = \begin{bmatrix} 0 & \cdots & 0 & \frac{-\phi}{\tau} & \frac{1}{\tau} \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \frac{1}{\tau}.$$

Then the above equation becomes

$$A_{t+1}^{-1}C_{t+1} = \sigma^2 C_{t+1} - \frac{M_{t+1}^{-1}u_{t+1}\frac{\sigma^2}{\tau}}{1 + u_{t+1}^\top M_{t+1}^{-1}u_{t+1}}. \quad (4.162)$$

Moreover,

$$M_{t+1}^{-1}u_{t+1} = \begin{bmatrix} A_t^{-1} & 0 \\ 0 & \sigma^2 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ -\frac{\phi}{\tau} \\ \frac{1}{\tau} \end{bmatrix} = \begin{bmatrix} A_t^{-1} & 0 \\ 0 & \sigma^2 \end{bmatrix} \begin{bmatrix} -\frac{\phi}{\tau}C_t \\ \frac{1}{\tau} \end{bmatrix} = \begin{bmatrix} -\frac{\phi}{\tau}A_t^{-1}C_t \\ \frac{\sigma^2}{\tau} \end{bmatrix},$$

$$u_{t+1}^\top M_{t+1}^{-1}u_{t+1} = \begin{bmatrix} 0 & \cdots & 0 & -\frac{\phi}{\tau} & \frac{1}{\tau} \end{bmatrix} \begin{bmatrix} -\frac{\phi}{\tau}A_t^{-1}C_t \\ \frac{\sigma^2}{\tau} \end{bmatrix} = \begin{bmatrix} -\frac{\phi}{\tau}C_t^\top & \frac{1}{\tau} \end{bmatrix} \begin{bmatrix} -\frac{\phi}{\tau}A_t^{-1}C_t \\ \frac{\sigma^2}{\tau} \end{bmatrix} = \frac{\phi^2}{\tau^2}C_t^\top A_t^{-1}C_t + \frac{\sigma^2}{\tau^2}.$$

Thus

$$\begin{aligned}
A_{t+1}^{-1}C_{t+1} &= \begin{bmatrix} -b_{t+1} \\ \sigma^2 - K_{t+1} \end{bmatrix} = \sigma^2 C_{t+1} - \frac{1}{1 + \frac{\phi^2}{\tau^2} C_t^\top A_t^{-1} C_t + \frac{\sigma^2}{\tau^2}} \begin{bmatrix} -\frac{\phi\sigma^2}{\tau^2} A_t^{-1} C_t \\ \frac{\sigma^4}{\tau^2} \end{bmatrix} \\
&= \sigma^2 C_{t+1} - \frac{1}{\tau^2 + \phi^2 C_t^\top A_t^{-1} C_t + \sigma^2} \begin{bmatrix} -\phi\sigma^2 A_t^{-1} C_t \\ \sigma^4 \end{bmatrix}
\end{aligned} \tag{4.163}$$

and

$$\sigma^2 - K_{t+1} = \sigma^2 - \frac{\sigma^4}{\tau^2 + \phi^2 C_t^\top A_t^{-1} C_t + \sigma^2} = \sigma^2 - \frac{\sigma^4}{\tau^2 + \sigma^2 + \phi^2(\sigma^2 - K_t)},$$

therefore

$$K_{t+1} = \frac{\sigma^4}{\tau^2 + \sigma^2 + \phi^2(\sigma^2 - K_t)}, \tag{4.164}$$

and

$$b_{t+1} = \begin{bmatrix} \frac{b_t \phi K_{t+1}}{\sigma^2} \\ \frac{K_{t+1}(\sigma^2 + \tau^2) - \sigma^4}{\phi \sigma^2} \end{bmatrix}, \tag{4.165}$$

$$\begin{aligned}
\bar{\mu}_{t+1} &= 0 - \sigma^4 K^{-1} b^\top (Z - b K^{-1} b^\top)^{-1} \sigma^{-4} (Z - b K^{-1} b^\top) y_{1:t} \\
&= -K^{-1} b^\top y_{1:t} \\
&= \frac{\phi}{\sigma^2} K_t \bar{\mu}_t + \phi \left(1 - \frac{K_t}{\sigma^2}\right) y_t, \\
\bar{\Sigma}_{t+1} &= \sigma^4 (K - b^\top Z^{-1} b)^{-1} - \sigma^4 K^{-1} b^\top (Z - b K^{-1} b^\top)^{-1} (Z - b K^{-1} b^\top) Z^{-1} b (K - b^\top Z^{-1} b)^{-1} \\
&= \sigma^4 (I - K^{-1} b^\top Z^{-1} b) (K - b^\top Z^{-1} b)^{-1} \\
&= \sigma^4 K_{t+1}^{-1},
\end{aligned}$$

where $K_1 = \frac{\sigma^4}{\frac{\phi^2}{\tau^2} + \frac{1}{L^2}}$.

The Estimation Distribution $p(x_{t+1} \mid y_{1:t+1}, \theta)$

The joint distribution for x_{t+1} and $y_{1:t+1}$ is $p(x_{t+1}, y_{1:t+1} \mid \theta) \sim N(0, \Gamma)$, where

$$\Gamma = \begin{bmatrix} C_{t+1}^\top (A - B)^{-1} C_{t+1} & C_{t+1}^\top (A - B)^{-1} \\ (A - B)^{-1} C_{t+1} & (I - A^{-1} B)^{-1} B^{-1} \end{bmatrix},$$

where C_{t+1}^\top is a $1 \times t + 1$ vector $\begin{bmatrix} 0 & \dots & 0 & 1 \end{bmatrix}_{1 \times t+1}$ retrieving the last column of a matrix. Because of

$$C_{t+1}^\top A_{t+1}^{-1} = \begin{bmatrix} -b_{t+1}^\top & \sigma^2 - K_{t+1} \end{bmatrix},$$

thus $x_{t+1} \mid y_{1:t+1} \sim N(\bar{\mu}_{t+1}^{(x)}, \bar{\sigma}_{t+1}^{(x)2})$, where

$$\begin{aligned}
\bar{\mu}_{t+1}^{(x)} &= \phi \hat{x}_t + C_{t+1}^\top (A - B)^{-1} B (I - A^{-1} B) y_{1:t+1} \\
&= \phi \hat{x}_t + C_{t+1}^\top A^{-1} B y_{1:t+1} \\
&= \phi \hat{x}_t + \frac{1}{\sigma^2} C_{t+1}^\top A^{-1} y_{1:t+1} \\
&= 0 + \frac{1}{\sigma^2} \begin{bmatrix} -b_{t+1}^\top & \sigma^2 - K_{t+1} \end{bmatrix} \begin{bmatrix} y_{1:t} \\ y_{t+1} \end{bmatrix} \\
&= -\frac{1}{\sigma^2} b_{t+1}^\top y_{1:t} + (1 - \frac{K_{t+1}}{\sigma^2}) y_{t+1} \\
&= \frac{K_{t+1} \bar{\mu}_t}{\sigma^2} + (1 - \frac{K_{t+1}}{\sigma^2}) y_{t+1} \\
\bar{\sigma}_{t+1}^{(x)2} &= C_{t+1}^\top (A - B)^{-1} C_{t+1} - C_{t+1}^\top (A - B)^{-1} B (I - A^{-1} B) (A - B)^{-1} C_{t+1} \\
&= C_{t+1}^\top (A - B)^{-1} C_{t+1} - C_{t+1}^\top A^{-1} B (A - B)^{-1} C_{t+1} \\
&= C_{t+1}^\top A^{-1} C_{t+1} \\
&= \sigma^2 - K_{t+1}.
\end{aligned}$$

Approximations of The Parameters Posterior

Because of the covariance $\Sigma_{YY} = (I - A^{-1} B)^{-1} B^{-1}$, therefore the inverse is

$$\Sigma_{YY}^{-1} = B(I - A^{-1} B) = B A^{-1} \Sigma_{XX}^{-1}.$$

Given the Choleski decomposition $LL^\top = A$, we have

$$\begin{aligned}
\Sigma_{YY}^{-1} &= B L^{-\top} L^{-1} \Sigma_{XX}^{-1} \\
&= (L^{-1} B)^\top (L^{-1} \Sigma_{XX}^{-1}) \\
&= \text{solve}(L, B)^\top \text{solve}(L, \Sigma_{XX}^{-1}).
\end{aligned}$$

More usefully, by given another Choleski decomposition $RR^\top = A - B = \Sigma_{XX}^{-1}$,

$$\begin{aligned}
Y^\top \Sigma_{YY}^{-1} Y &= \text{solve}(L, BY)^\top \text{solve}(L, \Sigma_{XX}^{-1} Y) \\
&\triangleq W^\top \text{solve}(L, \Sigma_{XX}^{-1} Y)
\end{aligned} \tag{4.166}$$

$$\begin{aligned}
\det \Sigma_{YY}^{-1} &= \det B \det L^{-\top} \det L^{-1} \det R \det R^\top \\
&= \det B (\det L^{-1})^2 (\det R)^2.
\end{aligned} \tag{4.167}$$

From the objective function, the posterior distribution of θ is

$$p(\theta \mid Y) \propto p(Y \mid \theta) p(\theta) \propto e^{-\frac{1}{2} Y \Sigma_{YY}^{-1} Y} \sqrt{\det \Sigma_{YY}^{-1}} p(\theta).$$

Then by taking natural logarithm on the posterior of θ and using the useful solutions in equations (4.182) and (4.183), we will have

$$\ln L(\theta) = -\frac{1}{2}Y^\top \Sigma_{YY}^{-1}Y + \frac{1}{2} \sum \ln \text{tr}(B) - \sum \ln \text{tr}(L) + \sum \ln \text{tr}(R) + \ln p(\theta). \quad (4.168)$$

4.26.2 High Dimension Parameters Space of OU-Process

The Brownian motion is used to construct the Ornstein Uhlenbeck (OU) process, which has become a popular tool for modeling interest rates and vehicle moving. The derivative of the Brownian motion x_t does not exist at any point in time. Thus, if x_t represents the position of a particle, we might be interested in obtaining its velocity, which is the derivative of the motion. The OU process is an alternative model to the Brownian motion that overcomes the preceding problem. It does this by considering the velocity u_t of a Brownian motion at time t . Over a small time interval, two factors affect the change in velocity: the frictional resistance of the surrounding medium whose effect is proportional to u_t and the random impact of neighboring particles whose effect can be represented by a standard Wiener process. Thus, because mass times velocity equals force, we have that

$$mdu_t = -\omega u_t dt + dW_t,$$

where $\omega > 0$ is called the friction coefficient and $m > 0$ is the mass. If we define $\gamma = \omega/m$ and $\lambda = 1/m$, we obtain the OU process with the following differential equation:

$$du_t = -\gamma u_t dt + \lambda dW_t. \quad (4.169)$$

The OU process is used to describe the velocity of a particle in a fluid and is encountered in statistical mechanics. It is the model of choice for random movement toward a concentration point. It is sometimes called a continuous-time Gauss Markov process, where a Gauss Markov process is a stochastic process that satisfies the requirements for both a Gaussian process and a Markov process. Because a Wiener process is both a Gaussian process and a Markov process, in addition to being a stationary independent increment process, it can be considered a Gauss-Markov process with independent increments Kijima (1997).

An OU-process model combining states and velocity is in the form of

$$\begin{cases} du_t = -\gamma u_t dt + \lambda dW_t, \\ dx_t = u_t dt + \xi dW'_t. \end{cases} \quad (4.170)$$

The solution can be found by integrating dt out, that gives us

$$\begin{cases} u_t &= u_{t-1}e^{-\gamma t} + \int_0^t \lambda e^{-\gamma(t-s)} dW_s, \\ x_t &= x_{t-1} + \frac{u_{t-1}}{\gamma}(1 - e^{-\gamma t}) + \int_0^t \frac{\lambda}{\gamma} e^{\gamma s} (1 - e^{-\gamma t}) dW_s + \int_0^t \xi dW'_s. \end{cases} \quad (4.171)$$

Therefore, the joint distribution is

$$\begin{bmatrix} x_t \\ u_t \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_t^{(x)} \\ \mu_t^{(u)} \end{bmatrix}, \begin{bmatrix} \sigma_t^{(x)2} & \rho_t \sigma_t^{(x)} \sigma_t^{(u)} \\ \rho_t \sigma_t^{(x)} \sigma_t^{(u)} & \sigma_t^{(u)2} \end{bmatrix} \right), \quad (4.172)$$

where $\mu_t^{(x)}$ and $\mu_t^{(u)}$ are from the forward map process

$$\begin{bmatrix} \mu_t^{(x)} \\ \mu_t^{(u)} \end{bmatrix} = \begin{bmatrix} 1 & \frac{1-e^{-\gamma\Delta_t}}{\gamma} \\ 0 & e^{-\gamma\Delta_t} \end{bmatrix} \begin{bmatrix} x_{t-1} \\ u_{t-1} \end{bmatrix} \triangleq \Phi \begin{bmatrix} x_{t-1} \\ u_{t-1} \end{bmatrix}, \quad (4.173)$$

and

$$\begin{cases} \sigma_t^{(x)2} &= \frac{\lambda^2(e^{2\gamma\Delta_t}-1)(1-e^{-\gamma\Delta_t})^2}{2\gamma^3} + \xi^2\Delta_t \\ \sigma_t^{(u)2} &= \frac{\lambda^2(1-e^{-2\gamma\Delta_t})}{2\gamma} \\ \rho_t \sigma_t^{(x)} \sigma_t^{(u)} &= \frac{\lambda^2(e^{\gamma\Delta_t}-1)(1-e^{-2\gamma\Delta_t})}{2\gamma^2} \end{cases}$$

In the above equations, $\Delta_t = T_t - T_{t-1}$, $\Delta_1 = 0$, $x_0 \sim N(0, L_x^2)$, $u_0 \sim N(0, L_u^2)$, $\rho_t^2 = 1 - \frac{\xi^2\Delta_t}{\sigma_t^{(x)2}}$. To be useful, $1 - \rho_t^2 = \frac{\sigma_t^{(x)2}}{\xi^2\Delta_t}$.

Moreover, the independent observation processes are

$$\begin{cases} y_t = x_t + \epsilon_t, \\ v_t = u_t + \epsilon'_t, \end{cases}$$

where $\epsilon_t \sim N(0, \sigma)$, $\epsilon'_t \sim N(0, \tau)$ are normally distributed independent errors. Thus, the joint distribution of observations is

$$\begin{bmatrix} y_t \\ v_t \end{bmatrix} \sim N \left(\begin{bmatrix} x_t \\ u_t \end{bmatrix}, \begin{bmatrix} \sigma^2 & 0 \\ 0 & \tau^2 \end{bmatrix} \right). \quad (4.174)$$

Consequently, the parameter θ of an entire Ornstein-Uhlenbeck process is a set of five parameters from both hidden status and observation process, which is represented as $\theta = \{\gamma, \xi^2, \lambda, \sigma^2, \tau^2\}$.

Starting from the joint distribution of $x_{0:t}, u_{0:t}$ and $y_{1:t}, v_{1:t}$ by given θ , it can be found that

$$\begin{bmatrix} \tilde{X} \\ \tilde{Y} \end{bmatrix} \Big| \theta \sim N(0, \tilde{\Sigma}), \quad (4.175)$$

where \tilde{X} represents for the hidden statues $\{x, u\}$, \tilde{Y} represents for observed $\{y, v\}$, θ is the set of five parameters. The inverse of the covariance matrix $\tilde{\Sigma}^{-1}$ is the procedure matrix in the form of

$$\tilde{\Sigma}^{-1} = \begin{bmatrix} Q_{xx} & Q_{xu} & -\frac{1}{\sigma^2}I & 0 \\ Q_{ux} & Q_{uu} & 0 & -\frac{1}{\tau^2}I \\ -\frac{1}{\sigma^2}I & 0 & \frac{1}{\sigma^2}I & 0 \\ 0 & -\frac{1}{\tau^2}I & 0 & \frac{1}{\tau^2}I \end{bmatrix}.$$

To make the covariance matrix a more beautiful form and convenient computing, \tilde{X} , \tilde{Y} and $\tilde{\Sigma}$ can be rearranged in a time series order, that makes $X = \{x_1, u_1, x_2, u_2, \dots, x_t, u_t\}$, $Y = \{y_1, v_1, y_2, v_2, \dots, y_t, v_t\}$ and the new procedure matrix Σ^{-1} looks like

$$\Sigma^{-1} = \begin{bmatrix} \sigma_{11}^{(x)2} + \frac{1}{\sigma^2} & \sigma_{11}^{(xu)2} & \dots & \sigma_{1t}^{(x)2} & \sigma_{1t}^{(xu)2} & -\frac{1}{\sigma^2} & 0 & \dots & 0 & 0 \\ \sigma_{11}^{(ux)2} & \sigma_{11}^{(u)2} + \frac{1}{\tau^2} & \dots & \sigma_{1t}^{(ux)2} & \sigma_{1t}^{(u)2} & 0 & -\frac{1}{\tau^2} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \sigma_{t1}^{(x)2} & \sigma_{t1}^{(xu)2} & \dots & \sigma_{tt}^{(x)2} + \frac{1}{\sigma^2} & \sigma_{tt}^{(xu)2} & 0 & 0 & \dots & -\frac{1}{\sigma^2} & 0 \\ \sigma_{t1}^{(ux)2} & \sigma_{t1}^{(u)2} & \dots & \sigma_{tt}^{(ux)2} & \sigma_{tt}^{(u)2} + \frac{1}{\tau^2} & 0 & 0 & \dots & 0 & -\frac{1}{\tau^2} \\ -\frac{1}{\sigma^2} & 0 & \dots & 0 & 0 & \frac{1}{\sigma^2} & 0 & \dots & 0 & 0 \\ 0 & -\frac{1}{\tau^2} & \dots & 0 & 0 & 0 & \frac{1}{\tau^2} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & -\frac{1}{\sigma^2} & 0 & 0 & 0 & \dots & \frac{1}{\sigma^2} & 0 \\ 0 & 0 & \dots & 0 & -\frac{1}{\tau^2} & 0 & 0 & \dots & 0 & \frac{1}{\tau^2} \end{bmatrix} \triangleq \begin{bmatrix} A_t \\ -B_t^\top \end{bmatrix}$$

where B_t is a $2t \times 2t$ diagonal matrix of observation errors at time t in the form of

$$\begin{bmatrix} \frac{1}{\sigma^2} & \cdot & \cdot & \cdot & \cdot \\ \cdot & \frac{1}{\tau^2} & \cdot & \cdot & \cdot \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \cdot & \cdot & \cdot & \frac{1}{\sigma^2} & \cdot \\ \cdot & \cdot & \cdot & \cdot & \frac{1}{\tau^2} \end{bmatrix}. \text{ In fact, the matrix } A_t \text{ is a } 2t \times 2t \text{ bandwidth six sparse matrix}$$

at time t in the process. Temporally, we are using A and B to represent the matrices A_t and B_t here. Then we may find the covariance matrix by calculating the inverse of

the procedure matrix as

$$\begin{aligned}
\Sigma &= \begin{bmatrix} (A - B^\top B^{-1} B)^{-1} & -(A - B^\top B^{-1} B)^{-1} B^\top B^{-1} \\ -B^{-1} B (A - B^\top B^{-1} B)^{-1} & (B - B^\top A^{-1} B)^{-1} \end{bmatrix} \\
&= \begin{bmatrix} (A - B)^{-1} & (A - B)^{-1} \\ (A - B)^{-1} & (I - A^{-1} B)^{-1} B^{-1} \end{bmatrix} \\
&\triangleq \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}.
\end{aligned}$$

The Forecast Distribution $p(Y_{t+1}|Y_{1:t}, \theta)$

We are now using the capital letter Y to represent the joint $\{y, v\}$ and $Y_{1:t} = \{y_1, v_1, y_2, v_2, \dots, y_t, v_t\}$, $Y_{t+1} = \{y_{t+1}, v_{t+1}\}$. It is known that

$$\begin{aligned}
p(Y_{1:t}, \theta) &\sim N(0, \Sigma_{YY}^{(t)}) \\
p(Y_{t+1}, Y_{1:t}, \theta) &\sim N(0, \Sigma_{YY}^{(t+1)}) \\
p(Y_{t+1} | Y_{1:t}, \theta) &\sim N(\bar{\mu}_{t+1}, \bar{\Sigma}_{t+1})
\end{aligned}$$

where the covariance matrix of the joint distribution is $\Sigma_{YY}^{(t+1)} = (I_{t+1} - A_{t+1}^{-1} B_{t+1})^{-1} B_{t+1}^{-1}$.

Then, by taking its inverse, we will get

$$\Sigma_{YY}^{(t+1)(-1)} = B_{t+1}(I_{t+1} - A_{t+1}^{-1} B_{t+1}).$$

To be clear, the matrix B_t is short for the matrix $B_t(\sigma^2, \tau^2)$, which is $2t \times 2t$ diagonal matrix with elements $\frac{1}{\sigma^2}, \frac{1}{\tau^2}$ repeating for t times on its diagonal. For instance, the very simple $B_1(\sigma^2, \tau^2) = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\tau^2} \end{bmatrix}_{2 \times 2}$ is a 2×2 matrix.

Because of A is symmetric and invertible, B is the diagonal matrix defined as above, then they have the following property

$$\begin{aligned}
AB &= A^\top B^\top = (BA)^\top, \\
A^{-1}B &= A^{-\top} B^\top = (BA^{-1})^\top.
\end{aligned}$$

Followed up the form of $\Sigma_{YY}^{(t+1)(-1)}$, we can find out that

$$\begin{aligned}
\Sigma_{YY}^{(t+1)(-1)} &= B_{t+1}(I_{t+1} - A_{t+1}^{-1} B_{t+1}) \\
&= B_{t+1}(B_{t+1}^{-1} - A_{t+1}^{-1}) B_{t+1} \\
&\triangleq \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} \begin{bmatrix} Z_{t+1} & b_{t+1} \\ b_{t+1}^\top & K_{t+1} \end{bmatrix} \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix}
\end{aligned}$$

where Z_{t+1} is a $2t \times 2t$ matrix, b_{t+1} is a $2t \times 2$ matrix and K_{t+1} is a 2×2 matrix. Thus by taking its inverse again, we will get

$$\Sigma_{YY}^{(t+1)} = \begin{bmatrix} B_t^{-1}(Z_{t+1} - b_{t+1}K_{t+1}^{-1}b_{t+1}^\top)^{-1}B_t^{-1} & -B_t^{-1}Z_{t+1}^{-1}b_{t+1}(K_{t+1} - b_{t+1}^\top Z_{t+1}^{-1}b_{t+1})^{-1}B_1^{-1} \\ -B_1^{-1}K_{t+1}^{-1}b_{t+1}^\top(Z_{t+1} - b_{t+1}K_{t+1}^{-1}b_{t+1}^\top)^{-1}B_t^{-1} & B_1^{-1}(K_{t+1} - b_{t+1}^\top Z_{t+1}^{-1}b_{t+1})^{-1}B_1^{-1} \end{bmatrix}.$$

It is easy to find the relationship between A_{t+1} and A_t in the Sherman-Morrison-Woodbury form is

$$A_{t+1} = \begin{bmatrix} A_t & \cdot & \cdot \\ \cdot & \frac{1}{\sigma^2} & \cdot \\ \cdot & \cdot & \frac{1}{\tau^2} \end{bmatrix} + U_{t+1}U_{t+1}^\top \triangleq M_{t+1} + U_{t+1}U_{t+1}^\top,$$

where $M_{t+1} = \begin{bmatrix} A_t & \cdot & \cdot \\ \cdot & \frac{1}{\sigma^2} & \cdot \\ \cdot & \cdot & \frac{1}{\tau^2} \end{bmatrix} = \begin{bmatrix} A_t & 0 \\ 0 & B_1 \end{bmatrix}$ and its inverse is $M_{t+1}^{-1} = \begin{bmatrix} A_t^{-1} & 0 \\ 0 & B_1^{-1} \end{bmatrix}$.

Additionally, U is a $2t + 2 \times 2$ matrix in the following form

$$U_{t+1} = \frac{1}{\sqrt{1-\rho_{t+1}^2}} \begin{bmatrix} \mathbf{0}_{2t-2} & \mathbf{0}_{2t-2} \\ \frac{1}{\sigma_{t+1}^{(x)}} & 0 \\ \frac{1-e^{-\gamma\Delta_{t+1}}}{\gamma\sigma_{t+1}^{(x)}} - \frac{\rho_{t+1}e^{-\gamma\Delta_{t+1}}}{\sigma_{t+1}^{(u)}} & \frac{\sqrt{1-\rho_{t+1}^2}e^{-\gamma\Delta_{t+1}}}{\sigma_{t+1}^{(u)}} \\ -\frac{1}{\sigma_{t+1}^{(x)}} & 0 \\ \frac{\rho_{t+1}}{\sigma_{t+1}^{(u)}} & -\frac{\sqrt{1-\rho_{t+1}^2}}{\sigma_{t+1}^{(u)}} \end{bmatrix} \triangleq \begin{bmatrix} C_t S_{t+1} \\ D_{t+1} \end{bmatrix},$$

denoted by $S_{t+1} = \frac{1}{\sqrt{1-\rho_{t+1}^2}} \begin{bmatrix} \frac{1}{\sigma_{t+1}^{(x)}} & 0 \\ \frac{1-e^{-\gamma\Delta_{t+1}}}{\gamma\sigma_{t+1}^{(x)}} - \frac{\rho_{t+1}e^{-\gamma\Delta_{t+1}}}{\sigma_{t+1}^{(u)}} & \frac{\sqrt{1-\rho_{t+1}^2}e^{-\gamma\Delta_{t+1}}}{\sigma_{t+1}^{(u)}} \end{bmatrix}$, $D_{t+1} = \frac{1}{\sqrt{1-\rho_{t+1}^2}} \begin{bmatrix} -\frac{1}{\sigma_{t+1}^{(x)}} & 0 \\ \frac{\rho_{t+1}}{\sigma_{t+1}^{(u)}} & -\frac{\sqrt{1-\rho_{t+1}^2}}{\sigma_{t+1}^{(u)}} \end{bmatrix}$

and $C_{t+1} = \begin{bmatrix} 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{0}_t \\ I_2 \end{bmatrix}.$

By post-multiplying $\Sigma_{YY}^{(t+1)(-1)}$ with C_{t+1} , it gives us

$$\begin{aligned}
\Sigma_{YY}^{(t+1)(-1)} C_{t+1} &= B_{t+1}(I_{t+1} - A_{t+1}^{-1} B_{t+1}) C_{t+1} \\
&= \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} \left(\begin{bmatrix} B_t^{-1} & 0 \\ 0 & B_1^{-1} \end{bmatrix} - A_{t+1}^{-1} \right) \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} C_{t+1} \\
&= \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} \begin{bmatrix} Z_{t+1} & b_{t+1} \\ b_{t+1}^\top & K_{t+1} \end{bmatrix} \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} C_{t+1} \\
&= \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} \begin{bmatrix} b_{t+1} B_1 \\ K_{t+1} B_1 \end{bmatrix}.
\end{aligned}$$

and the property of A_{t+1}^{-1} is

$$A_{t+1}^{-1} C_{t+1} = \begin{bmatrix} -b_{t+1} \\ B_1^{-1} - K_{t+1} \end{bmatrix}.$$

Moreover, by pre-multiplying C_{t+1}^\top on the left side of the above equation, we will have

$$C_{t+1}^\top A_{t+1}^{-1} C_{t+1} = B_1^{-1} - K_{t+1}, \quad (4.176)$$

$$K_{t+1} = B_1^{-1} - C_{t+1}^\top A_{t+1}^{-1} C_{t+1}. \quad (4.177)$$

We may use Sherman-Morrison-Woodbury formula to find the inverse of A_{t+1} in a recursive way, which is

$$A_{t+1}^{-1} = (M_{t+1} + U_{t+1} U_{t+1}^\top)^{-1} = M_{t+1}^{-1} - M_{t+1}^{-1} U_{t+1} (I + U_{t+1}^\top M_{t+1}^{-1} U_{t+1})^{-1} U_{t+1}^\top M_{t+1}^{-1}. \quad (4.178)$$

Consequently, it is easy to find that $M_{t+1}^{-1} C_{t+1} = \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix}$ and

$$\begin{aligned}
A_{t+1}^{-1} C_{t+1} &= \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix} - \begin{bmatrix} A_t^{-1} & 0 \\ 0 & B_1^{-1} \end{bmatrix} \begin{bmatrix} C_t S_{t+1} \\ D \end{bmatrix} (I + U_{t+1}^\top M_{t+1}^{-1} U_{t+1})^{-1} \begin{bmatrix} S_{t+1}^\top C_t^\top & D_{t+1}^\top \end{bmatrix} \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix} \\
&= \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix} - \begin{bmatrix} A_t^{-1} C_t S_{t+1} \\ B_1^{-1} D_{t+1} \end{bmatrix} (I + U_{t+1}^\top M_{t+1}^{-1} U_{t+1})^{-1} D_{t+1}^\top B_1^{-1} \\
&= \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix} - \begin{bmatrix} A_t^{-1} C_t S_{t+1} \\ B_1^{-1} D_{t+1} \end{bmatrix} (I + S_{t+1}^\top C_t^\top A_t^{-1} C_t S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1} \\
&= \begin{bmatrix} 0 \\ B_1^{-1} \end{bmatrix} - \begin{bmatrix} A_t^{-1} C_t S_{t+1} \\ B_1^{-1} D_{t+1} \end{bmatrix} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1}.
\end{aligned}$$

Thus, by using the equation (4.176), we will get

$$K_{t+1} = B_1^{-1} D_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1}, \quad (4.179)$$

and

$$\begin{aligned} b_{t+1} &= A_t^{-1} C_t S_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1} \\ &= \begin{bmatrix} -b_t \\ B_1^{-1} - K_t \end{bmatrix} S_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1}. \end{aligned}$$

To achieve the recursive updating formula, firstly we need to find the form of $b_{t+1}^\top B_t^2 Y_{1:t}$.

In fact, it is

$$\begin{aligned} b_{t+1}^\top B_t Y_{1:t} &= B_1^{-\top} D_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-\top} S_{t+1}^\top \begin{bmatrix} -b_t^\top & B_1^{-1} - K_t \end{bmatrix} B_t \begin{bmatrix} Y_{1:t-1} \\ Y_t \end{bmatrix} \\ &= B_1^{-\top} D_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-\top} S_{t+1}^\top (-b_t^\top B_{t-1} Y_{1:t-1} + (B_1^{-1} - K_t) Y_t) \\ &= B_1^{-\top} D_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-\top} S_{t+1}^\top (K_t B_1 \bar{\mu}_t + (I - K_t B_1) Y_t), \end{aligned}$$

By using equation (4.181) and simplifying the above equation, one can achieve a recursive updating form of the mean, which is

$$\begin{aligned} \bar{\mu}_{t+1} &= -B_1 K_{t+1}^{-1} b_{t+1}^\top B_t Y_{1:t} \\ &= -D_{t+1}^{-\top} S_{t+1}^\top (K_t B_1 \bar{\mu}_t + (I - K_t B_1) Y_t) \\ &= -D_{t+1}^{-\top} S_{t+1}^\top (Y_t + K_t B_1 (\bar{\mu}_t - Y_t)), \end{aligned}$$

where by simplifying $D^{-\top} S^\top$, one may find

$$D_{t+1}^{-\top} S_{t+1}^\top = \begin{bmatrix} -1 & -\frac{1-e^{-\gamma\Delta_{t+1}}}{\gamma} \\ 0 & -e^{-\gamma\Delta_{t+1}} \end{bmatrix} = -\Phi_{t+1},$$

which is the negative of forward process. Then the final form of recursive updating formula are

$$\begin{cases} \bar{\mu}_{t+1} &= \Phi_{t+1} K_t B_1 \bar{\mu}_t + \Phi_{t+1} (I - K_t B_1) Y_t \\ \bar{\Sigma}_{t+1} &= (B_1 K_{t+1} B_1)^{-1} \end{cases}. \quad (4.180)$$

The matrix K_{t+1} is updated via

$$K_{t+1} = B_1^{-1} D_{t+1} (I + S_{t+1}^\top (B_1^{-1} - K_t) S_{t+1} + D_{t+1}^\top B_1^{-1} D_{t+1})^{-1} D_{t+1}^\top B_1^{-1}, \quad (4.181)$$

or updating its inverse in the following form makes the computation faster, that is

$$\begin{cases} K_{t+1}^{-1} &= B_1 D_{t+1}^{-\top} D_{t+1}^{-1} B_1 + B_1 \Phi_{t+1} (B_1^{-1} - K_t) \Phi_{t+1}^\top B_1 + B_1, \\ \bar{\Sigma}_{t+1} &= D_{t+1}^{-\top} D_{t+1}^{-1} + \Phi_{t+1} (B_1^{-1} - K_t) \Phi_{t+1}^\top + B_1^{-1} \end{cases}$$

$$\text{and } K_1 = B_1^{-1} - A_1^{-1} = \begin{bmatrix} \frac{\sigma^4}{\sigma^2 + L_x^2} & 0 \\ 0 & \frac{\tau^4}{\tau^2 + L_u^2} \end{bmatrix}.$$

The Estimation Distribution $p(X_{t+1}|Y_{1:t+1}, \theta)$

The filtering distribution of the state given parameters is $p(X_t | Y_{1:t}, \theta)$. To find its form, one can use the joint distribution of X_{t+1} and $Y_{1:t+1}$, which is $p(X_{t+1}, Y_{1:t+1} | \theta) \sim N(0, \Gamma)$, where

$$\Gamma = \begin{bmatrix} C_{t+1}^\top (A - B)^{-1} C_{t+1} & C_{t+1}^\top (A - B)^{-1} \\ (A - B)^{-1} C_{t+1} & (I - A^{-1}B)^{-1} B^{-1} \end{bmatrix}.$$

Because of

$$C_{t+1}^\top A_{t+1}^{-1} = \begin{bmatrix} -b_{t+1}^\top & B_1^{-1} - K_{t+1} \end{bmatrix},$$

then $X_{t+1} | Y_{1:t+1}, \theta \sim N(\bar{\mu}_{t+1}^{(X)}, \bar{\sigma}_{t+1}^{(X)2})$, where

$$\begin{aligned} \bar{\mu}_{t+1}^{(X)} &= \Phi \hat{x}_t + C_{t+1}^\top (A - B)^{-1} B (I - A^{-1}B) Y_{1:t+1} \\ &= \Phi \hat{x}_t + C_{t+1}^\top A^{-1} B Y_{1:t+1} \\ &= 0 + \begin{bmatrix} -b_{t+1}^\top & B_1^{-1} - K_{t+1} \end{bmatrix} \begin{bmatrix} B_t & 0 \\ 0 & B_1 \end{bmatrix} \begin{bmatrix} Y_{1:t} \\ Y_{t+1} \end{bmatrix} \\ &= -b^\top B_t Y_{1:t} + (I - B_1 K_{t+1}) Y_{t+1} \\ &= K_{t+1} B_1 \bar{\mu}_{t+1} + (I - B_1 K_{t+1}) Y_{t+1} \\ \bar{\sigma}_{t+1}^{(X)2} &= C_{t+1}^\top (A - B)^{-1} C_{t+1} - C_{t+1}^\top (A - B)^{-1} B (I - A^{-1}B) (A - B)^{-1} C_{t+1} \\ &= C_{t+1}^\top (A - B)^{-1} C_{t+1} - C_{t+1}^\top A^{-1} B (A - B)^{-1} C_{t+1} \\ &= C_{t+1}^\top A^{-1} C_{t+1} \\ &= B_1^{-1} - K_{t+1}. \end{aligned}$$

Approximations of The Parameters Posterior

From the joint distribution of X and Y , we can find that

$$\Sigma_{YY}^{-1} = B(I - A^{-1}B) = BA^{-1}\Sigma_{XX}^{-1}.$$

Given the Choleski decomposition $LL^\top = A$, we have

$$\begin{aligned} \Sigma_{YY}^{-1} &= BL^{-\top} L^{-1} \Sigma_{XX}^{-1} \\ &= (L^{-1}B)^\top (L^{-1} \Sigma_{XX}^{-1}) \\ &= \text{solve}(L, B)^\top \text{solve}(L, \Sigma_{XX}^{-1}). \end{aligned}$$

More usefully, by given another Choleski decomposition $RR^\top = A - B = \Sigma_{XX}^{-1}$,

$$\begin{aligned} Y^\top \Sigma_{YY}^{-1} Y &= \text{solve}(L, BY)^\top \text{solve}(L, \Sigma_{XX}^{-1} Y) \\ &\triangleq W^\top \text{solve}(L, \Sigma_{XX}^{-1} Y) \end{aligned} \quad (4.182)$$

$$\begin{aligned} \det \Sigma_{YY}^{-1} &= \det B \det L^{-\top} \det L^{-1} \det R \det R^\top \\ &= \det B (\det L^{-1})^2 (\det R)^2. \end{aligned} \quad (4.183)$$

From the objective function, the second term in the integral is

$$p(\theta | Y) \propto p(Y | \theta) p(\theta) \propto e^{-\frac{1}{2} Y^\top \Sigma_{YY}^{-1} Y} \sqrt{\det \Sigma_{YY}^{-1}} P(\theta).$$

Then by taking natural logarithm on the posterior of θ and using the useful solutions in equations (4.182) and (4.183), we will have

$$\ln L(\theta) = -\frac{1}{2} Y^\top \Sigma_{YY}^{-1} Y + \frac{1}{2} \sum \ln \text{tr}(B) - \sum \ln \text{tr}(L) + \sum \ln \text{tr}(R). \quad (4.184)$$

4.27 Prior Distribution for Variance Parameters

The well known Hierarchical Linear Model, where the parameters vary at more than one level, was firstly introduced by Lindley and Smith in 1972 and 1973 Lindley and Smith (1972) Smith (1973). An extension of these models is non-linear Hierarchical Model. Hierarchical Model can be used on data with many levels, although 2-level models are the most common ones. The state space model in equations (4.151) and (4.152) is one of Hierarchical Linear Model if G_t and F_t are linear and non-linear model if G_t and F_t are non-linear processes. Researchers have made a few discussions and works on these both linear and non-linear models. In this section, we only discuss on the prior for variance parameters in these models.

Jonathan and Thomas in Stroud and Bengtsson (2007) have discussed a model, which is slightly different with a Gaussian state-space model in equations (4.151) and (4.152) from section one. The two errors ω_t and ϵ_t are assumed normally distributed as

$$\begin{aligned} \omega_t &\sim N(0, \alpha Q), \\ \epsilon_t &\sim N(0, \alpha R), \end{aligned}$$

where the two matrices R and Q are known and α is an unknown scale factor to be estimated. (Note that a perfect model is obtained by setting $Q = 0$.) Therefore, the

density of Gaussian state-space model is

$$p(y_t | x_t, \alpha) = N(F(x_t), \alpha R),$$
$$p(x_t | x_{t-1}, \alpha) = N(G(x_{t-1}), \alpha Q).$$

The parameter α is assumed *Inverse Gamma* distribution.

Various non-informative and weakly-informative prior distributions have been suggested for scale parameters in hierarchical models. Andrew Gelman gave a discussion on prior distributions for variance parameters in hierarchical models in 2006 Gelman *et al.* (2006). General considerations include using invariance Jeffries (1961), maximum entropy Jaynes (1983) and agreement with classical estimators Box and Tiao (2011).

4.27.1 Priors Discussion

http://andrewgelman.com/2007/07/18/informative_and/

Informative and noninformative priors Posted by Andrew on 18 July 2007, 8:04 am
Neal writes,

As I start your Bayesian stuff, can I ask you the same question I asked Boris a few years ago, namely, as you note, noninf priors simply represent the situation where we know very little and want the data to speak (so in the end not too far from the classical view). Can you point me to any social science (closer to ps is better) where people actually update, so that the prior in a second study is the posterior of the first (whether or not the two studies done by same person or not).

Equivalently point me to a study which uses non-inf priors. (as more than a toy i know the piece by gill and his student).

Btw do you know the old piece by Harry Roberts, saying that as a scientist all we can report is the likelihood, and that everyone should put their own prior in and then produce their own posterior. so all articles would just be a computer program which takes as input my prior and produces my posterior given the likelihood surface estimated by the author?

My reply: now I like weakly informative priors. But thats new since our books. Regarding informative priors in applied research, we can distinguish three categories:

(1) Prior distributions giving numerical information that is crucial to estimation of the model. This would be a traditional informative prior, which might come from a literature review or explicitly from an earlier data analysis.

(2) Prior distributions that are not supplying any controversial information but are strong enough to pull the data away from inappropriate inferences that are consistent

with the likelihood. This might be called a weakly informative prior.

(3) Prior distributions that are uniform, or nearly so, and basically allow the information from the likelihood to be interpreted probabilistically. These are noninformative priors, or maybe, in some cases, weakly informative.

I have examples of (1), (2), and (3) in my own applied research. Category (3) is the most common for me, but an example of (2) is my 1990 paper with King on seats-votes curves, where we fit a mixture model and used an informative prior to constrain the locations, scales, and masses of the three components. An example of (3) is my 1996 paper with Bois and Jiang where we used an informative prior distribution for several parameters in a toxicology model. We were careful to parameterize the model so that these priors made sense, and the model also had an interesting two-level structure which we discuss in that paper and also in Section 9.1 of *Bayesian Data Analysis*.

Regarding your question about models where people actually update: we did this in our radon analysis (see here) where the posterior distribution from a national data analysis (based on data from over 80,000 houses) gives inference for each county in the U.S., which is in turn used as the prior distribution for the radon level in your house, which in turn can be updated if you have information from a measurement in your house.

One of the convenient things about doing applied statistics is that eventually I can come up with an example for everything from my own experience. (This also makes it fun to write books.)

Regarding your last comment: yes, there is an idea that a Bayesian wants everyone else to be non-Bayesian so that he or she can do cleaner analyses. I discuss that idea in this talk from 2003 which I've been too lazy to write up as a paper.

4.27.2 Discussion two

<http://andrewgelman.com/2007/05/11/weaklyinformative/>

Weakly informative priors Posted by Andrew on 11 May 2007, 1:01 pm

Bayesians traditionally consider prior distributions that (a) represent the actual state of subject-matter knowledge, or (b) are completely or essentially noninformative. We consider an alternative strategy: choosing priors that convey some generally useful information but clearly less than we actually have for the particular problem under study. We give some examples, including the Cauchy (0, 2.5) prior distribution for logistic regression coefficients, and then briefly discuss the major unsolved problem in Bayesian inference: the construction of models that are structured enough to learn

from data but weak enough to learn from data.

Im speaking Monday on this at Jun Lius workshop on Monte Carlo methods at Harvard (my talk is 9:45-10:30am Monday at 104 Harvard Hall).

Heres the presentation. I think this is potentially a huge advance in how we think about Bayesian models.

References

- Abramovich, F., Sapatinas, T., and Silverman, B. W. (1998). Wavelet thresholding via a Bayesian approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(4), 725–749.
- Agarwal, P. K., Arge, L., and Erickson, J. (2003). Indexing moving points. *Journal of Computer and System Sciences*, 66(1), 207–243.
- Andrieu, C., De Freitas, N., and Doucet, A. (1999). Sequential MCMC for Bayesian model selection. In *Higher-Order Statistics, 1999. Proceedings of the IEEE Signal Processing Workshop on*, 130–134. IEEE.
- Andrieu, C., Doucet, A., and Holenstein, R. (2010). Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3), 269–342.
- Arnaud Doucet, Nando de Freitas, N. G. (Ed.) (2011). *Sequential Monte Carlo Methods in Practice*. Springer-Verlag New York.
- Arulampalam, M. S., Maskell, S., Gordon, N., and Clapp, T. (2002). A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE Transactions on signal processing*, 50(2), 174–188.
- Aydin, D. and Tuzemen, M. S. (2012). Smoothing parameter selection problem in nonparametric regression based on smoothing spline: A simulation study. *Journal of Applied Sciences*, 12(7), 636.
- Bartlett, M. S. (1951). An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, 22(1), 107–111.
- Berlinet, A. and Thomas-Agnan, C. (2011). *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media.

- Bodewig, E. (1959). *Matrix Calculus*, North.
- Box, G. E. and Tiao, G. C. (2011). *Bayesian inference in statistical analysis*, Volume 40. John Wiley & Sons.
- Cappé, O., Godsill, S. J., and Moulines, E. (2007). An overview of existing methods and recent advances in sequential Monte Carlo. *Proceedings of the IEEE*, 95(5), 899–924.
- Cappé, O., Moulines, E., and Rydén, T. (2009). Inference in hidden markov models. In *Proceedings of EUSFLAT Conference*, 14–16.
- Chen, W.-K. (2009). *Feedback, nonlinear, and distributed circuits*. CRC Press.
- Craven, P. and Wahba, G. (1978). Smoothing noisy data with spline functions. *Numerische Mathematik*, 31(4), 377–403.
- De Boor, C., De Boor, C., Mathématicien, E.-U., De Boor, C., and De Boor, C. (1978). *A practical guide to splines*, Volume 27. Springer-Verlag New York.
- De Jong, P. (1988). The likelihood for a state space model. *Biometrika*, 165–169.
- Del Moral, P., Doucet, A., and Jasra, A. (2006). Sequential monte carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3), 411–436.
- Deng, C. Y. (2011). A generalization of the Sherman–Morrison–Woodbury formula. *Applied Mathematics Letters*, 24(9), 1561–1564.
- Donoho, D. L. and Johnstone, I. M. (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the american statistical association*, 90(432), 1200–1224.
- Donoho, D. L., Johnstone, I. M., Kerkycharian, G., and Picard, D. (1995). Wavelet shrinkage: asymptopia? *Journal of the Royal Statistical Society. Series B (Methodological)*, 301–369.
- Donoho, D. L. and Johnstone, J. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3), 425–455.

- Ellis, D., Sommerlade, E., and Reid, I. (2009). Modelling pedestrian trajectory patterns with gaussian processes. In *Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on*, 1229–1234. IEEE.
- Erkorkmaz, K. and Altintas, Y. (2001). High speed CNC system design. Part I: jerk limited trajectory generation and quintic spline interpolation. *International Journal of machine tools and manufacture*, 41(9), 1323–1345.
- Gelman, A. *et al.* (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian analysis*, 1(3), 515–534.
- Gordon, N. J., Salmond, D. J., and Smith, A. F. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In *IEE Proceedings F (Radar and Signal Processing)*, Volume 140, 107–113. IET.
- Green, P. J. and Silverman, B. W. (1993). *Nonparametric regression and generalized linear models: a roughness penalty approach*. CRC Press.
- Gu, C. (1998). Model indexing and smoothing parameter selection in nonparametric function estimation. *Statistica Sinica*, 607–623.
- Gu, C. (2013). *Smoothing Spline ANOVA Models Second Edition*. Springer New York Heidelberg Dordrecht London.
- Handschin, J. (1970). Monte Carlo techniques for prediction and filtering of non-linear stochastic processes. *Automatica*, 6(4), 555–563.
- Handschin, J. E. and Mayne, D. Q. (1969). Monte Carlo techniques to estimate the conditional expectation in multi-stage non-linear filtering. *International journal of control*, 9(5), 547–559.
- Hangos, K. M., Bokor, J., and Szederkényi, G. (2006). *Analysis and control of nonlinear process systems*. Springer Science & Business Media.
- Hastie, T. J. and Tibshirani, R. J. (1990). *Generalized additive models*, Volume 43. CRC Press.
- Higuchi, T. (2001). Self-organizing time series model. In *Sequential Monte Carlo Methods in Practice*, 429–444. Springer.

- Jaynes, E. T. (1983). Papers On Probability. *Statistics and Statistical Physics*.
- Jeffries, H. (1961). Theory of probability.
- Judd, K. L. (1998). *Numerical methods in economics*. MIT press.
- Kalman, R. E. *et al.* (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1), 35–45.
- Kantas, N., Doucet, A., Singh, S., and Maciejowski, J. (2009). An overview of sequential Monte Carlo methods for parameter estimation. In *in General State-Space Models, in IFAC System Identification, no. Ml*. Citeseer.
- Kijima, M. (1997). *Markov processes for stochastic modeling*, Volume 6. CRC Press.
- Kim, Y.-J. and Gu, C. (2004). Smoothing spline Gaussian regression: more scalable computation via efficient approximation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(2), 337–356.
- Kimeldorf, G. and Wahba, G. (1971). Some results on Tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1), 82–95.
- Kimeldorf, G. S. and Wahba, G. (1970). A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2), 495–502.
- Kitagawa, G. (1996). Monte Carlo filter and smoother for non-Gaussian nonlinear state space models. *Journal of computational and graphical statistics*, 5(1), 1–25.
- Kitagawa, G. (1998). A self-organizing state-space model. *Journal of the American Statistical Association*, 1203–1215.
- Lindley, D. V. and Smith, A. F. (1972). Bayes estimates for the linear model. *Journal of the Royal Statistical Society. Series B (Methodological)*, 1–41.
- Liu, J. and West, M. (2001). Combined parameter and state estimation in simulation-based filtering. In *Sequential Monte Carlo methods in practice*, 197–223. Springer.
- Liu, Z. and Guo, W. (2010). Data driven adaptive spline smoothing. *Statistica Sinica*, 1143–1163.

- MacKay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge university press.
- Nason, G. (2010). *Wavelet methods in statistics with R*. Springer Science & Business Media.
- Peng, F. (1983). A Necessary and Sufficient Condition to Judge Function Linearly Independent. *JOURNAL OF JISHOU UNIVERSITY(NATURAL SCIENCE EDITION)*, (0), 25–28.
- Petris, G., Petrone, S., and Campagnoli, P. (2009). Dynamic linear models. *Dynamic Linear Models with R*, 31–84.
- Polson, N. G., Stroud, J. R., and Müller, P. (2008). Practical filtering with sequential parameter learning. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(2), 413–428.
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. The MIT Press.
- Ristic, B., Arulampalam, S., and Gordon, N. (2004). Beyond the Kalman filter. *IEEE Aerospace and Electronic Systems Magazine*, 19(7), 37–38.
- Robert, C. P. (2004). *Monte carlo methods*. Wiley Online Library.
- Schwarz, K.-P. (2012). *Geodesy Beyond 2000: The Challenges of the First Decade, IAG General Assembly Birmingham, July 19–30, 1999*, Volume 121. Springer Science & Business Media.
- Sealfon, C., Verde, L., and Jimenez, R. (2005). Smoothing spline primordial power spectrum reconstruction. *Physical Review D*, 72(10), 103520.
- Sherlock, C., Golightly, A., and Henderson, D. A. (2016). Adaptive, delayed-acceptance MCMC for targets with expensive likelihoods. *Journal of Computational and Graphical Statistics*, (just-accepted).
- Sherman, J. and Morrison, W. J. (1950). Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *The Annals of Mathematical Statistics*, 21(1), 124–127.

- Silverman, B. W. (1985). Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society. Series B (Methodological)*, 1–52.
- Smith, A. F. (1973). A general Bayesian linear model. *Journal of the Royal Statistical Society. Series B (Methodological)*, 67–75.
- Stroud, J. R. and Bengtsson, T. (2007). Sequential state and variance estimation within the ensemble Kalman filter. *Monthly Weather Review*, 135(9), 3194–3208.
- Trevor Hastie, Robert Tibshirani, J. F. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second Edition*. Springer-Verlag.
- Wahba, G. (1985). A comparison of GCV and GML for choosing the smoothing parameter in the generalized spline smoothing problem. *The Annals of Statistics*, 1378–1402.
- Wahba, G. (1990). *Spline models for observational data*, Volume 59. Siam.
- Wahba, G. and Wold, S. (1975). A completely automatic French curve: Fitting spline functions by cross validation. *Communications in Statistics-Theory and Methods*, 4(1), 1–17.
- Wang, Y. (1998). Smoothing spline models with correlated random errors. *Journal of the American Statistical Association*, 93(441), 341–348.
- West, M. (1993). Mixture models, Monte Carlo, Bayesian updating, and dynamic models. *Computing Science and Statistics*, 325–325.
- Woodbury, M. A. (1950). Inverting modified matrices. *Memorandum report*, 42(106), 336.
- Yang, K. and Sukkarieh, S. (2010). An analytical continuous-curvature path-smoothing algorithm. *Robotics, IEEE Transactions on*, 26(3), 561–568.
- Yao, F., Müller, H.-G., Wang, J.-L., *et al.* (2005). Functional linear regression analysis for longitudinal data. *The Annals of Statistics*, 33(6), 2873–2903.
- Yu, B., Kim, S. H., Bailey, T., and Gamboa, R. (2004). Curve-based representation of moving object trajectories. In *Database Engineering and Applications Symposium, 2004. IDEAS'04. Proceedings. International*, 419–425. IEEE.

Zhang, K., Guo, J.-X., and Gao, X.-S. (2013). Cubic spline trajectory generation with axis jerk and tracking error constraints. *International Journal of Precision Engineering and Manufacturing*, 14(7), 1141–1146.

