



Dynamic Learning with the EM Algorithm for Neural Networks

J.F.G. de FREITAS

Computer Science Division, 387 Soda Hall, University of California, Berkeley, CA 94720-1776, USA

M. NIRANJAN

Department of Computer Science, University of Sheffield, Regent Court, Portabello Street, Sheffield S1 4DP

A.H. GEE

Cambridge University Engineering Department, Trumpington Street, Cambridge CB2 1PZ, England

Received March 25, 1999; Revised September 9, 1999

Abstract. In this paper, we derive an EM algorithm for nonlinear state space models. We use it to estimate jointly the neural network weights, the model uncertainty and the noise in the data. In the E-step we apply a forward-backward Rauch-Tung-Striebel smoother to compute the network weights. For the M-step, we derive expressions to compute the model uncertainty and the measurement noise. We find that the method is intrinsically very powerful, simple and stable.

1. Introduction

In environments where data is available in batches, it is possible to address the general problem of Bayesian learning with Gaussian approximations, in a principled way, using the expectation maximisation (EM) algorithm [1] and dynamical models. Moreover, such an approach allows us to treat non-stationary data sets. This paper will focus on this learning strategy. In particular, it aims to extend the current work on EM learning for dynamical linear systems to the problem of computing the weights of a multi-layer perceptron (MLP), the initial conditions and the noise variances jointly.

The application of the EM algorithm to learning and inference in linear dynamical systems has occupied the attention of several researchers in the past. Chen [2] was one of the pioneers in this field. In particular, he applied the EM algorithm to linear state space models known in the statistics literature as MIMIC models. In these models one observes multiple indicators and multiple causes of a single latent variable. Chen's MIMIC model was implemented in a simulation study relating social status and participation.

Watson and Engle [3] have suggested using the EM algorithm, in conjunction with the method of scoring, for the estimation of linear dynamic factor, MIMIC and varying coefficient regression models. They evaluated their paradigm experimentally by estimating common factors in wage rate data from several industries in Los Angeles, USA.

In 1982, Shumway and Stoffer [4] proposed the use of the EM algorithm and linear state space models for time series smoothing and forecasting with missing observations. To demonstrate their method, they considered a health series representing total expenditures for physician services as measured by two different sources. The time series produced by each source have similar values but exhibit missing observations at different periods. In Shumway and Stoffer's approach, the two series are automatically merged into an overall expenditure series, which is then used for forecasting. Nine years later, Shumway and Stoffer [5] extended their work to switching linear dynamic models. In essence, they derived a state space representation with measurement matrices that switch according to a time varying independent random process. They illustrate

their method on an application involving the tracking of multiple targets.

The method of learning and inference in linear state space models via the EM algorithm has also played a role in the fields of speech analysis and computer vision. Digalakis, Rohlicek and Ostendorf [6] applied it to the speech recognition problem. They made a connection between this method and the Baum-Welch estimation algorithm for hidden Markov models (HMMs). North and Blake [7] have implemented the method to learn linear dynamic state space models used for tracking contours in images. Rao and Ballard [8] have also explored the relevance of the EM algorithm together with state space estimation in the field of vision. They have developed a hierarchical network model of visual recognition that encapsulates these concepts.

Ghahramani [9] has embedded the EM method for learning dynamic linear systems in a graphical models framework. He treats computationally intractable models, such as factorial HMMs and switching state space models, by resorting to Gibbs sampling and variational approximations. In another paper, Roweis and Ghahramani [10] make use of the EM algorithm and linear state space representations to present a unified view of linear Gaussian models including factor analysis, mixtures of Gaussians, standard and probabilistic versions of principal component analysis, vector quantisation, Kalman smoothing and linear hidden Markov models.

This paper is organised as follows. Section 2 introduces the nonlinear state space modelling scheme adopted in the paper. The application of extended Kalman smoothing to estimate the weights of an MLP is discussed in Section 3. Section 4 presents a brief derivation of the EM algorithm, which is used as a step towards the derivation of the EM algorithm for nonlinear state space models in Section 5. Section 6 examines some of the results obtained with experiments on synthetic and real data.

2. Nonlinear State Space Model

To investigate the application of the EM algorithm to state space learning, we shall focus on the following nonlinear state space representation:

$$\begin{aligned}\theta_{t+1} &= A\theta_t + \mathbf{u}_t \\ \mathbf{y}_t &= \hat{\mathbf{f}}(\mathbf{x}_t, \theta_t) + \mathbf{v}_t\end{aligned}\quad (1)$$

where $\mathbf{x} \in \mathbb{R}^d$ denotes the input data, $\mathbf{y} \in \mathbb{R}^c$ denotes the output data and $\theta \in \mathbb{R}^m$ denotes the model states.

The measurements nonlinear mapping $\hat{\mathbf{f}}(\mathbf{x}_t, \theta_t)$ corresponds to a multi-layer perceptron (MLP) whose weights are the model states θ . The framework may be easily extended to encompass recurrent networks, radial basis networks and many other approximation techniques. We assume the measurement (\mathbf{v}_t) and process (\mathbf{u}_t) noise terms to be zero mean Gaussian with covariances R and Q respectively. The matrix A contains information about how the states evolve. It is particularly useful in tracking applications. However, when the above model is employed merely for parameter estimation in neural network models with stationary data, there is no need for the matrix A .

Despite the fact that the data is processed in batches, the model of Eq. (1) allows the weights to be time varying. It is, therefore, possible to deal with non-stationary data sets. In the event of the data being stationary, we should expect the process noise term to vanish. Consequently, if we know that the data is stationary, the estimate of the process noise can be used to determine how well the model explains the data. In Section 6, we demonstrate this method on a few stationary data sets.

Our objective is to estimate the model states (MLP weights) $\hat{\theta}_t$ and the set of parameters $\varphi \triangleq \{R, Q, A, \mu, \Pi\}$ given the measurements $\{\mathbf{x}_{1:N}, \mathbf{y}_{1:N}\}$,¹ where μ and Π denote the mean and covariance of the Gaussian prior $p(\theta_0 | \varphi)$.

3. The Extended Kalman Smoother

One of the earliest implementations of the extended Kalman filter (EKF) to train MLPs is due to Singhal and Wu [11]. The algorithm's computational complexity is of the order $\mathcal{O}(cm^2)$ multiplications per time step. Shah, Palmieri and Datum [12] and Puskorius and Feldkamp [13] have proposed various approximations to the weights covariance so as to simplify this problem. Here, we extend the work in this area by proposing an algorithm to estimate the noise covariances R and Q and the initial conditions μ and Π . The method is also more accurate as it involves smoothing instead of plain filtering.

Smoothing often entails forward and backward filtering over a segment of data so as to obtain improved averaged estimates. Various techniques have been proposed to accomplish this goal [14, 15]. In our work, we make use of the well known Rauch-Tung-Striebel smoother [16]. The forward filtering stage involves computing the estimates $\hat{\theta}_t$ and P_t , over a segment of

N samples, with the following EKF recursions:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{t+1|t} &= A\hat{\boldsymbol{\theta}}_t \\ P_{t+1|t} &= AP_tA' + Q \\ K_{t+1} &= P_{t+1|t}G'_{t+1}(R + G_{t+1}P_{t+1|t}G'_{t+1})^{-1} \\ \hat{\boldsymbol{\theta}}_{t+1} &= \hat{\boldsymbol{\theta}}_{t+1|t} + K_{t+1}(\mathbf{y}_{t+1} - \hat{\mathbf{f}}(\mathbf{x}_{t+1}, \hat{\boldsymbol{\theta}}_{t+1|t})) \\ P_{t+1} &= P_{t+1|t} - K_{t+1}G_{t+1}P_{t+1|t}\end{aligned}$$

where K denotes the Kalman gain, A' the transpose of A and G the Jacobian matrix:

$$G = \left. \frac{\partial \hat{\mathbf{f}}(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}} \right|_{(\boldsymbol{\theta}=\hat{\boldsymbol{\theta}})}$$

$$= \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}_1(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_1} & \frac{\partial \hat{\mathbf{f}}_2(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_1} & \dots & \frac{\partial \hat{\mathbf{f}}_c(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_1} \\ \frac{\partial \hat{\mathbf{f}}_1(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_2} & & & \\ \vdots & & & \vdots \\ \frac{\partial \hat{\mathbf{f}}_1(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_m} & \dots & & \frac{\partial \hat{\mathbf{f}}_c(\boldsymbol{\theta}, \mathbf{x})}{\partial \theta_m} \end{bmatrix}'$$

Subsequently, the Rauch-Tung-Striebel smoother makes use of the following backward recursions:

$$\begin{aligned}J_{t-1} &= P_{t-1}A'P_{t|t-1}^{-1} \\ \hat{\boldsymbol{\theta}}_{t-1|N} &= \hat{\boldsymbol{\theta}}_{t-1}J_{t-1}(\hat{\boldsymbol{\theta}}_{t|N} - A\hat{\boldsymbol{\theta}}_{t-1}) \\ P_{t-1|N} &= P_{t-1} + J_{t-1}(P_{t|N} - P_{t|t-1})J'_{t-1} \\ P_{t,t-1|N} &= P_tJ'_{t-1} + J_t(P_{t+1,t|N} - AP_t)J'_{t-1}\end{aligned}$$

where the parameters, covariance and cross-covariance are defined as follows:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{t|N} &= \mathbb{E}(\boldsymbol{\theta}_t | \mathbf{y}_{1:N}) \\ P_{t|N} &= \mathbb{E}((\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_t)(\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_t)' | \mathbf{y}_{1:N}) \\ P_{t,t-1|N} &= \mathbb{E}((\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_t)(\boldsymbol{\theta}_{t-1} - \hat{\boldsymbol{\theta}}_{t-1})' | \mathbf{y}_{1:N})\end{aligned}$$

They may be initialised with the following values:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{N|N} &= \hat{\boldsymbol{\theta}}_N \\ P_{N|N} &= P_N \\ P_{N,N-1|N} &= (I - K_NG'_N)AP_{N-1}\end{aligned}$$

The extended Kalman smoother provides a minimum variance Gaussian approximation to the posterior probability density function $p(\boldsymbol{\theta}_{0:t} | \mathbf{y}_{1:t})$ [17]. In many

cases, however, the posterior is multi-modal. This problem can be circumvented by implementing a mixture of Kalman smoothers, where each individual smoother approximates a particular mode.

4. The EM Algorithm

So far, we have shown that given a set of parameters $\boldsymbol{\varphi} = \{R, Q, A, \boldsymbol{\mu}, \Pi\}$ and a matrix of N measurements $\mathbf{y}_{1:N}$, it is possible to compute the expected values of the states with an extended Kalman smoother. In this section, we present a treatment of the EM algorithm that will allow us to learn the parameters $\boldsymbol{\varphi}$ of nonlinear state space models.

The EM algorithm is an iterative method for finding a mode of the likelihood function $p(\mathbf{y}_{1:N} | \boldsymbol{\varphi})$. Roughly speaking, it proceeds as follows: (E-step 1) estimate the states $\boldsymbol{\theta}_{0:N}$ given a set of parameters $\boldsymbol{\varphi}$, (M-step 1) estimate the parameters given the new states, (E-step 2) re-estimate the states with the new parameters, and so forth. The most remarkable attribute of the EM algorithm is that it ensures an increase in the likelihood function at each iteration. However, as the EKF can only provide an approximation to the true states $\boldsymbol{\theta}_{0:N}$ in the E step, the EM algorithm to train MLPs is not necessarily guaranteed to converge.

To gain more insight into the EM method, let us express the likelihood function as follows:

$$\begin{aligned}p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}) &= p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}) \frac{p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi})}{p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi})} \\ &= \frac{p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})}{p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi})}\end{aligned}$$

Taking the logarithms of both sides yields the following identity:

$$\begin{aligned}\ln p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}) &= \ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}) \\ &\quad - \ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi})\end{aligned}$$

Let us treat $\boldsymbol{\theta}_{0:N}$ as a random variable with distribution $p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})$, where $\boldsymbol{\varphi}^{\text{old}}$ is the current guess. If we then take expectations on both sides of the previous identity, while remembering that the left hand side does not depend on $\boldsymbol{\theta}_{0:N}$, we get:

$$\begin{aligned}\ln p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}) &= \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\ &\quad - \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}))\end{aligned}\quad (2)$$

where the expectations involve averaging over the matrix $\boldsymbol{\theta}_{0:N}$ under $p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})$. For example:

$$\begin{aligned} & \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\ &= \int (\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}}) d\boldsymbol{\theta}_{0:N} \end{aligned}$$

It is well known that the second term on the right side of Eq. (2) is maximised for $\boldsymbol{\varphi}^{\text{old}}$. That is:

$$\mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})) \geq \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}))$$

for any $\boldsymbol{\varphi}$. To apply the EM algorithm, we need to compute the first term on the right hand side of Eq. (2) repeatedly. The aim is to maximise this term at each iteration. One method of maximising it is discussed in detail in the next section. For the time being, let us assume that we can maximise it, that is:

$$\mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{new}})) \geq \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{old}}))$$

Then, it follows that the likelihood function also increases at every iteration. To demonstrate this important result, consider the change in likelihood for a single iteration:

$$\begin{aligned} & \ln p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{new}}) - \ln p(\mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{old}}) \\ &= (\mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{new}})) \\ &\quad - \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}^{\text{old}}))) \\ &\quad - (\mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{new}})) \\ &\quad - \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}}))) \end{aligned}$$

The right hand side of the above equation is positive because we are averaging under the distribution $p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})$. Consequently, the likelihood function is guaranteed to increase at each iteration. The EM algorithm's name originates from the steps that are required to increase $\mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}))$, namely compute the Expectation and then Maximise it. The EM algorithm thus involves the following steps:

Initialisation: Start with a guess for $\boldsymbol{\varphi}^0$.

E-step: Determine the expected log-likelihood density function of the complete data given the current estimate $\boldsymbol{\varphi}^{\text{old}}$:

$$\begin{aligned} & \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\ &= \int (\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}}) d\boldsymbol{\theta}_{0:N} \end{aligned}$$

M-step: Compute a new value of $\boldsymbol{\varphi}$ that maximises the expected log-likelihood of the complete data. The maximum can be found by simple differentiation of the expected log-likelihood with respect to $\boldsymbol{\varphi}$.

5. The EM Algorithm for Nonlinear State Space Models

To derive the EM algorithm for nonlinear state space models, we need to develop an expression for the likelihood of the completed data. We assume that the likelihood of the data given the states, the initial conditions and the evolution of the states can be represented by Gaussian distributions. In particular, if the initial mean and covariance of the states is given by $\boldsymbol{\mu}$ and $\boldsymbol{\Pi}$, then:

$$\begin{aligned} p(\boldsymbol{\theta}_0 | \boldsymbol{\varphi}) &= \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Pi}|^{1/2}} \exp \left[-\frac{1}{2} (\boldsymbol{\theta}_0 - \boldsymbol{\mu})' \boldsymbol{\Pi}^{-1} (\boldsymbol{\theta}_0 - \boldsymbol{\mu}) \right] \\ p(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}, \boldsymbol{\varphi}) &= \frac{1}{(2\pi)^{m/2} |\boldsymbol{Q}|^{1/2}} \exp \\ &\quad \times \left[-\frac{1}{2} (\boldsymbol{\theta}_t - \boldsymbol{A}\boldsymbol{\theta}_{t-1})' \boldsymbol{Q}^{-1} (\boldsymbol{\theta}_t - \boldsymbol{A}\boldsymbol{\theta}_{t-1}) \right] \\ p(\mathbf{y}_t | \boldsymbol{\theta}_t, \boldsymbol{\varphi}) &= \frac{1}{(2\pi)^{c/2} |\boldsymbol{R}|^{1/2}} \exp \\ &\quad \times \left[-\frac{1}{2} (\mathbf{y}_t - \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t))' \boldsymbol{R}^{-1} (\mathbf{y}_t - \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)) \right] \end{aligned}$$

Under the model assumptions of uncorrelated noise sources and Markov state evolution, the likelihood of the complete data is given by:

$$\begin{aligned} & p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}) \\ &= p(\boldsymbol{\theta}_0 | \boldsymbol{\varphi}) \prod_{t=1}^N p(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}, \boldsymbol{\varphi}) \prod_{t=1}^N p(\mathbf{y}_t | \boldsymbol{\theta}_t, \boldsymbol{\varphi}) \end{aligned}$$

Hence, the log-likelihood of the complete data is given by the following expression:

$$\begin{aligned} & \ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi}) \\ &= - \sum_{t=1}^N \left[\frac{1}{2} (\mathbf{y}_t - \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t))' \boldsymbol{R}^{-1} (\mathbf{y}_t - \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)) \right] \\ &\quad - \frac{N}{2} \ln |\boldsymbol{R}| - \sum_{t=1}^N \left[\frac{1}{2} (\boldsymbol{\theta}_t - \boldsymbol{A}\boldsymbol{\theta}_{t-1})' \boldsymbol{Q}^{-1} (\boldsymbol{\theta}_t - \boldsymbol{A}\boldsymbol{\theta}_{t-1}) \right] \end{aligned}$$

$$\begin{aligned}
 & -\frac{N}{2} \ln |Q| - \frac{1}{2} (\boldsymbol{\theta}_0 - \boldsymbol{\mu})' \Pi^{-1} (\boldsymbol{\theta}_0 - \boldsymbol{\mu}) \\
 & -\frac{1}{2} \ln |\Pi| - \frac{Nc + (N+1)m}{2} \ln(2\pi) \quad (3)
 \end{aligned}$$

As discussed in the previous section, all we need to do now is to compute the expectation of $\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})$ and then differentiate the result with respect to the parameters $\boldsymbol{\varphi}$ so as to maximise it. The EM algorithm for nonlinear state space models will thus involve computing the expected values of the states and covariances with the extended Kalman smoother and then maximising the parameters $\boldsymbol{\varphi}$ with the formulae obtained by differentiating the expected log-likelihood.

5.1. Computing the Expectation of the Log-Likelihood

If we take the expectation of the log-likelihood for the complete data, by averaging over $\boldsymbol{\theta}_{0:N}$ under the distribution $p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})$, we get the following expression:

$$\begin{aligned}
 & \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
 & = -\frac{N}{2} \ln |R| - \frac{N}{2} \ln |Q| \\
 & \quad -\frac{1}{2} \ln |\Pi| - \frac{Nc + (N+1)m}{2} \ln(2\pi) \\
 & \quad - \sum_{t=1}^N \frac{1}{2} \mathbb{E}[\mathbf{y}'_t R^{-1} \mathbf{y}_t - \mathbf{y}'_t R^{-1} \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t) \\
 & \quad \quad - \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)' R^{-1} \mathbf{y}_t + \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)' R^{-1} \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)] \\
 & \quad - \sum_{t=1}^N \frac{1}{2} \mathbb{E}[\boldsymbol{\theta}'_t Q^{-1} \boldsymbol{\theta}_t - \boldsymbol{\theta}'_t Q^{-1} A \boldsymbol{\theta}_{t-1} \\
 & \quad \quad - \boldsymbol{\theta}'_{t-1} A' Q^{-1} \boldsymbol{\theta}_t + \boldsymbol{\theta}'_{t-1} A' Q^{-1} A \boldsymbol{\theta}_{t-1}] \\
 & \quad - \frac{1}{2} \mathbb{E}[\boldsymbol{\theta}'_0 \Pi^{-1} \boldsymbol{\theta}_0 - \boldsymbol{\theta}'_0 \Pi^{-1} \boldsymbol{\mu} \\
 & \quad \quad - \boldsymbol{\mu}' \Pi^{-1} \boldsymbol{\theta}_0 + \boldsymbol{\mu}' \Pi^{-1} \boldsymbol{\mu}]
 \end{aligned}$$

We need to digress briefly to compute the expectation of the measurements mapping $\hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)$. We should recall that the EKF approximation to this mapping is

given by:

$$\begin{aligned}
 \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t) & = \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) + \left. \frac{\partial \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)}{\partial \boldsymbol{\theta}_t} \right|_{(\boldsymbol{\theta}_t = \hat{\boldsymbol{\theta}}_{t|N})} \\
 & \quad \times (\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_{t|N}) + \dots
 \end{aligned}$$

Consequently, if we take expectations on both sides of the equation, we get:

$$\mathbb{E}(\hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)) \approx \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)$$

and

$$\begin{aligned}
 & \mathbb{E}((\hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t) - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))(\hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t) - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))') \\
 & \approx \mathbb{E} \left[\left(\hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) + \left. \frac{\partial \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)}{\partial \boldsymbol{\theta}_t} \right|_{(\boldsymbol{\theta}_t = \hat{\boldsymbol{\theta}}_{t|N})} \right. \right. \\
 & \quad \times (\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_{t|N}) - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) \left. \left. \times \left(\hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) + \left. \frac{\partial \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)}{\partial \boldsymbol{\theta}_t} \right|_{(\boldsymbol{\theta}_t = \hat{\boldsymbol{\theta}}_{t|N})} \right. \right. \right. \\
 & \quad \left. \left. \times (\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_{t|N}) - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) \right) \right]' = G_t P_{t|N} G_t'
 \end{aligned}$$

Hence, under the distribution $p(\boldsymbol{\theta}_{0:N} | \mathbf{y}_{1:N}, \boldsymbol{\varphi}^{\text{old}})$, it follows that:

$$\begin{aligned}
 & \mathbb{E}(\hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t) \hat{\mathbf{f}}(\boldsymbol{\theta}_t, \mathbf{x}_t)') \\
 & \approx G_t P_{t|N} G_t' + \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)'
 \end{aligned}$$

Using this approximation and the fact that the trace and expectation operators are linear, the expectation of the log-likelihood becomes:

$$\begin{aligned}
 & \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
 & \approx -\frac{N}{2} \ln |R| - \frac{N}{2} \ln |Q| \\
 & \quad -\frac{1}{2} \ln |\Pi| - \frac{Nc + (N+1)m}{2} \ln(2\pi) \\
 & \quad - \sum_{t=1}^N \frac{1}{2} \text{tr}(R^{-1} [\mathbf{y}_t \mathbf{y}'_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) \mathbf{y}'_t - \mathbf{y}_t \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)' \\
 & \quad \quad + \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t) \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)' + G_t P_{t|N} G_t']) \\
 & \quad - \sum_{t=1}^N \frac{1}{2} \text{tr}(Q^{-1} [\hat{\boldsymbol{\theta}}_{t|N} \hat{\boldsymbol{\theta}}'_{t|N} + P_{t|N} - 2A(\hat{\boldsymbol{\theta}}_{t|N} \hat{\boldsymbol{\theta}}'_{t-1|N})
 \end{aligned}$$

$$\begin{aligned}
& + P_{t,t-1|N})' + A(\hat{\boldsymbol{\theta}}_{t-1}\hat{\boldsymbol{\theta}}'_{t-1|N} + P_{t-1|N})A') \\
& - \frac{1}{2}\text{tr}(\Pi^{-1}[\hat{\boldsymbol{\theta}}_{0|N}\hat{\boldsymbol{\theta}}'_{0|N} + P_{0|N} - 2\hat{\boldsymbol{\theta}}_{0|N}\boldsymbol{\mu}' + \boldsymbol{\mu}\boldsymbol{\mu}'])
\end{aligned}$$

Completing squares and using the following abbreviations:

$$\begin{aligned}
\Gamma &= \sum_{t=1}^N \hat{\boldsymbol{\theta}}_{t|N}\hat{\boldsymbol{\theta}}'_{t|N} + P_{t|N} \\
\Delta &= \sum_{t=1}^N \hat{\boldsymbol{\theta}}_{t-1|N}\hat{\boldsymbol{\theta}}'_{t-1|N} + P_{t-1|N} \\
\Upsilon &= \sum_{t=1}^N \hat{\boldsymbol{\theta}}_{t|N}\hat{\boldsymbol{\theta}}'_{t-1|N} + P_{t,t-1|N}
\end{aligned}$$

we get our final expression for the approximate expectation of the log-likelihood:

$$\begin{aligned}
& \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
& \approx -\frac{N}{2} \ln |R| - \frac{N}{2} \ln |Q| \\
& - \frac{1}{2} \ln |\Pi| - \frac{Nc + (N+1)m}{2} \ln(2\pi) \\
& - \sum_{t=1}^N \frac{1}{2} \text{tr}(R^{-1}[(\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)) \\
& \times (\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))' + G_t P_{t|N} G_t']) \\
& - \frac{1}{2} \text{tr}(Q^{-1}[\Gamma - 2A\Upsilon' + A\Delta A']) \\
& - \frac{1}{2} \text{tr}(\Pi^{-1}[(\hat{\boldsymbol{\theta}}_{0|N} - \boldsymbol{\mu})(\hat{\boldsymbol{\theta}}_{0|N} - \boldsymbol{\mu})' + P_{0|N}])
\end{aligned} \tag{4}$$

5.2. Differentiating the Expected Log-Likelihood

To maximise the expected value of the log-likelihood with respect to the parameters $\boldsymbol{\varphi}$, we need to compute the derivatives with respect to each parameter individually. This is done in the subsequent sections, where we make use of some results of matrix differentiation [18].

5.2.1. Maximum with Respect to A. Differentiating the expected log-likelihood with respect to A yields:

$$\begin{aligned}
& \frac{\partial}{\partial A} \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
& \approx -\frac{1}{2} \frac{\partial}{\partial A} \text{tr}(Q^{-1}[\Gamma - 2A\Upsilon' + A\Delta A']) \\
& = -\frac{1}{2} (-2Q^{-1}\Upsilon + 2Q^{-1}A\Delta)
\end{aligned}$$

Equating this result to zero yields the value of A that maximises the approximate log-likelihood:

$$A = \Upsilon \Delta^{-1} \tag{5}$$

5.2.2. Maximum with Respect to R. Differentiating the expected log-likelihood with respect to R^{-1} gives:

$$\begin{aligned}
& \frac{\partial}{\partial R^{-1}} \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
& \approx \frac{\partial}{\partial R^{-1}} \left(\frac{N}{2} \ln |R^{-1}| - \sum_{t=1}^N \frac{1}{2} \text{tr}(R^{-1}[G_t P_{t|N} G_t' \right. \\
& \quad \left. + (\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))(\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))'] \right) \\
& = \frac{N}{2} R - \sum_{t=1}^N \frac{1}{2} (G_t P_{t|N} G_t' \\
& \quad + (\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))(\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))')
\end{aligned}$$

Hence, by equating the above result to zero, the approximate maximum of the log-likelihood with respect to R is given by:

$$\begin{aligned}
R &= \frac{1}{N} \sum_{t=1}^N (G_t P_{t|N} G_t' + (\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t)) \\
& \quad \times (\mathbf{y}_t - \hat{\mathbf{f}}(\hat{\boldsymbol{\theta}}_{t|N}, \mathbf{x}_t))')
\end{aligned} \tag{6}$$

5.2.3. Maximum with Respect to Q. Following the same steps, the derivative of the expected log-likelihood with respect to Q^{-1} is given by:

$$\begin{aligned}
& \frac{\partial}{\partial Q^{-1}} \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\
& \approx \frac{N}{2} Q - \frac{1}{2} (\Gamma - 2A\Upsilon' + A\Delta A')
\end{aligned}$$

Hence, equating to zero and using the result that $A = \Upsilon \Delta^{-1}$, the approximate maximum of the log-likelihood with respect to Q is given by:

$$Q = \frac{1}{N} (\Gamma - \Upsilon \Delta^{-1} \Upsilon') \tag{7}$$

5.2.4. Maximum with Respect to $\boldsymbol{\mu}$. It is also possible to treat the initial conditions as parameters and improve their estimates in the M-step of the EM algorithm. Finding the derivative of the expected log-likelihood with respect to the initial states gives:

$$\frac{\partial}{\partial \boldsymbol{\mu}} \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \approx \frac{1}{2} \Pi^{-1} (-2\hat{\boldsymbol{\theta}}_{0|N} + 2\boldsymbol{\mu})$$

Hence, the initial value for the states should be:

$$\boldsymbol{\mu} = \hat{\boldsymbol{\theta}}_{0|N} \quad (8)$$

5.2.5. Maximum with Respect to Π . The derivative of the expected log-likelihood with respect to the inverse of the initial covariance gives:

$$\begin{aligned} & \frac{\partial}{\partial \Pi^{-1}} \mathbb{E}(\ln p(\boldsymbol{\theta}_{0:N}, \mathbf{y}_{1:N} | \boldsymbol{\varphi})) \\ & \approx \frac{\Pi}{2} - \frac{1}{2}((\hat{\boldsymbol{\theta}}_{0|N} - \boldsymbol{\mu})(\hat{\boldsymbol{\theta}}_{0|N} - \boldsymbol{\mu})' + P_{0|N}) \end{aligned}$$

Therefore, the initial covariance should be updated as follows:

$$\Pi = P_{0|N} \quad (9)$$

5.3. The E and M Steps for Nonlinear State Space Models

We can now prescribe the EM algorithm for nonlinear state space models as follows:

Initialisation: Start with a guess for $\boldsymbol{\varphi} = \{R, Q, A, \boldsymbol{\mu}, \Pi\}$.

E-step: Determine the expected values $\hat{\boldsymbol{\theta}}_{t|N}$, $P_{t|N}$ and $P_{t,t-1|N}$, given the current parameter estimate $\boldsymbol{\varphi}^{\text{old}}$, using the extended Kalman smoothing equations described in Section 4.2.

M-step: Compute new values of the parameters $\boldsymbol{\varphi} = \{R, Q, A, \boldsymbol{\mu}, \Pi\}$ using equations (5) to (9).

The complexity of this algorithm is $\mathcal{O}(m^3 N)$ operations per iteration.

6. Experiments

6.1. Simple Regression Example

For the purposes of demonstrating the method, we address the problem of learning the following nonlinear mapping from (x_1, x_2) to y :

$$y = 4 \sin(x_1 - 2) + 2x_2 + 5 + \eta$$

where x_1 and x_2 were chosen to be two normal random sequences of 700 samples each. The noise process η was sampled from a zero mean Gaussian

distribution with variance $R = 0.5$. An MLP with four sigmoidal neurons in the hidden layer and a linear neuron in the output layer was used to approximate the measurements mapping. After 50 iterations, as shown in Fig. 1, the estimate of observation variance R converges to the true value. In addition, the trace of the process noise covariance Q goes to zero. Note that since the data is stationary, the trace of Q should tend to zero. That is, the trace of Q can be used to provide an estimate of how well the model fits the data. The innovations covariance (variance of the evidence function $p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \hat{\boldsymbol{\theta}}_{t|t-1}, Q_{t-1}, R_{t-1})$) tends to R over the entire data set, as shown in Fig. 2. The top plot of this figure shows that the MLP approximates the true function without fitting the noise. That is, it generalises well. Figure 1 also shows how the log-likelihood increases at each step, thereby demonstrating that the algorithm converges well.

6.2. Robot Arm Mapping

This data set is often used as a benchmark to compare neural network algorithms.² It involves implementing a model to map the joint angle of a robot arm (x_1, x_2) to the position of the end of the arm (y_1, y_2) . The data were generated from the following model:

$$y_1 = 2.0 \cos(x_1) + 1.3 \cos(x_1 + x_2) + \epsilon_1$$

$$y_2 = 2.0 \sin(x_1) + 1.3 \sin(x_1 + x_2) + \epsilon_2$$

where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, $\sigma = 0.05$. We use the first 200 observations of the data set to train our models and the last 200 observations to test them.

Figure 3 shows the 3D plots of the training data and the contours of the training and test data. The contour plots also include the typical approximations that were obtained using our algorithm and an MLP with 2 linear output neurons and 20 sigmoidal hidden neurons. Figure 4 illustrates the convergence of the algorithm. In this particular run the training and test mean square errors were 0.0057 and 0.0081 (the minimum bound being $2\sigma^2 = 0.005$). Our mean square errors are of the same magnitude as the ones reported by other researchers [19–23]. Figure 4 also shows the two diagonal entries of the measurements noise covariance and the trace of the process noise covariance. They behave as expected.

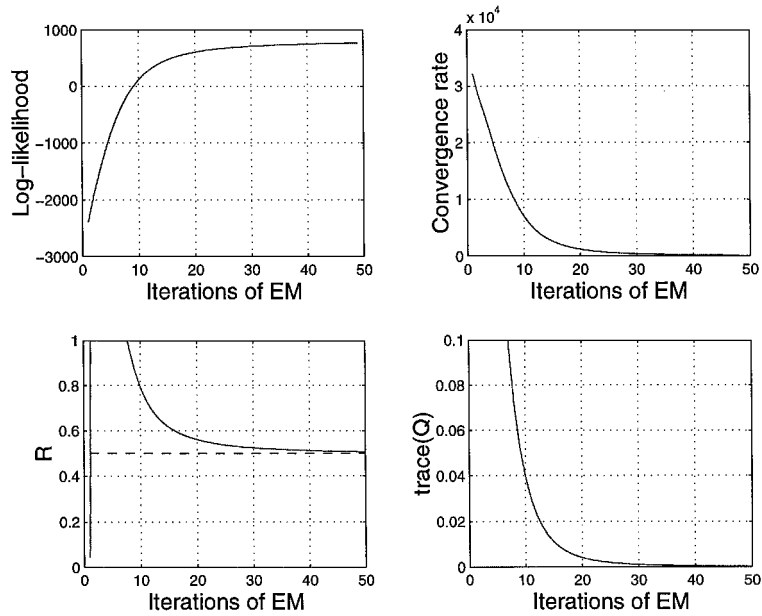


Figure 1. The top plots show the log-likelihood function and the convergence rate (log-likelihood slope) for the simple regression problem. The bottom plots show the convergence of the measurements noise covariance R and the trace of the process noise covariance Q .

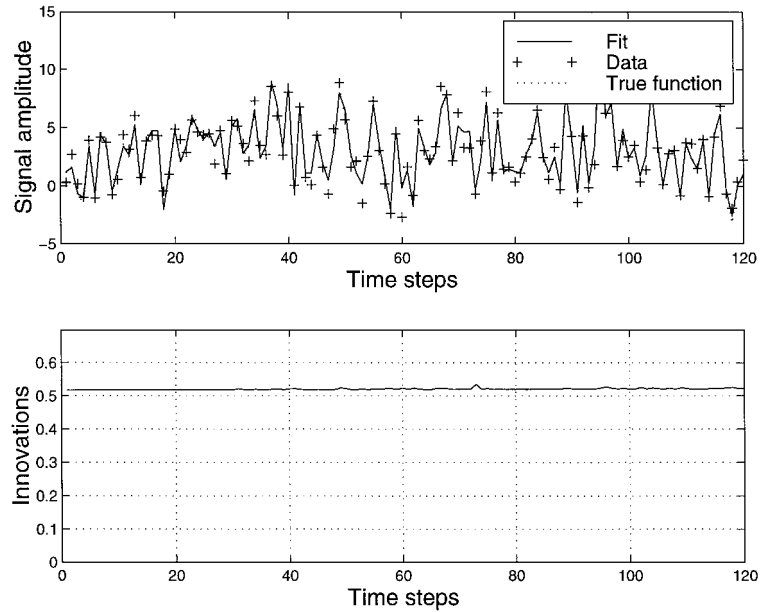


Figure 2. The top plot shows that the MLP fit, for the regression example, approximates the true function (the former is almost exactly on top of the latter); it does not fit the noise. The bottom plot shows that the uncertainty in the predictions (innovations) converges to the uncertainty engendered by the measurement noise.

6.3. Classification with Medical Data

Here, we consider an interesting nonlinear classification data set³ collected as part of a study to iden-

tify patients with muscle tremor [24, 25]. The data was gathered from a group of patients (9 with, primarily, Parkinson’s disease or multiple sclerosis) and from a control group (not exhibiting the disease). Arm

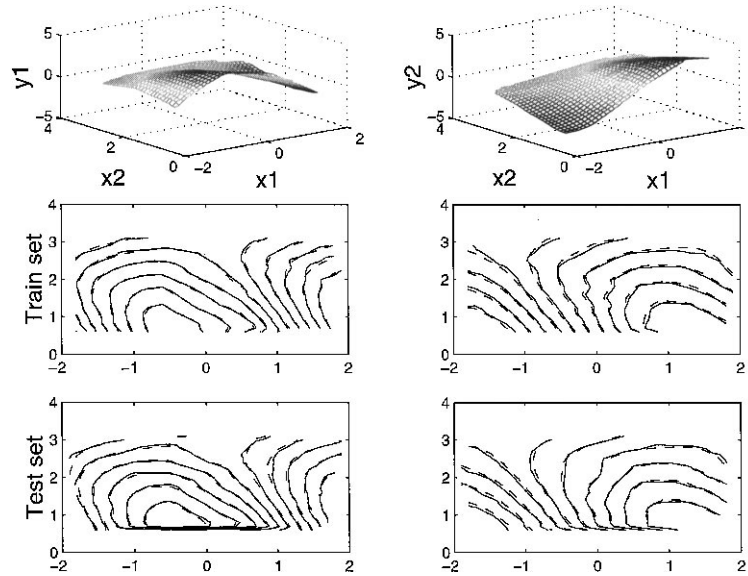


Figure 3. The top plots show the training data surfaces corresponding to each coordinate of the robot arm’s position. The Middle and bottom plots show the training and validation data [- -] and the respective MLP mappings [—].

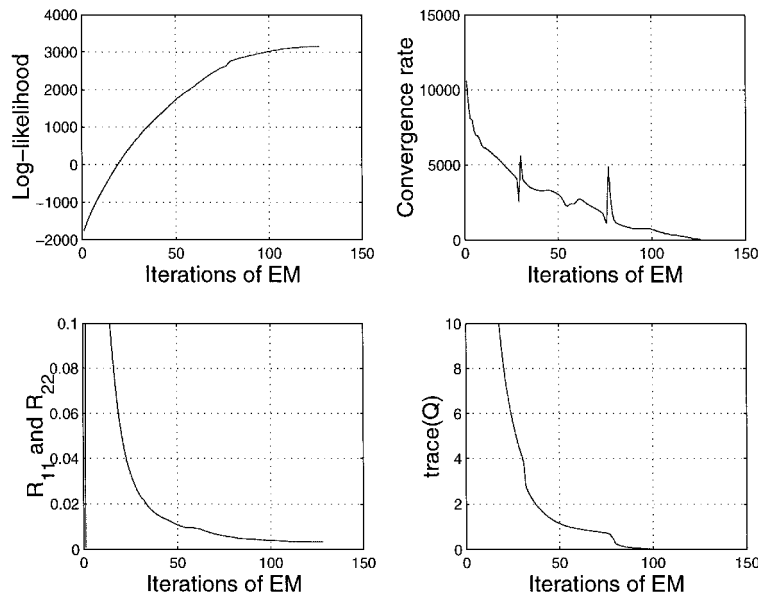


Figure 4. The top plots show the log-likelihood function and the convergence rate (log-likelihood slope) for the robot arm problem. The bottom plots show the convergence of the diagonal entries of the measurements noise covariance R (almost identical) and the trace of the process noise covariance Q .

muscle tremor was measured with a 3-D mouse and a movement tracker in three linear and three angular directions. The time series of the measurements were parameterised using a set of autoregressive models. The number of features was then reduced to two [24]. Figure 5 shows a plot of these features for patient (o)

and control groups (+). The figure also shows the decision boundaries (solid lines) and confidence intervals (dashed lines) obtained with an MLP, consisting of 10 sigmoidal hidden neurons and an output linear neuron. We should point out, however, that having an output linear neuron leads to a classification framework based

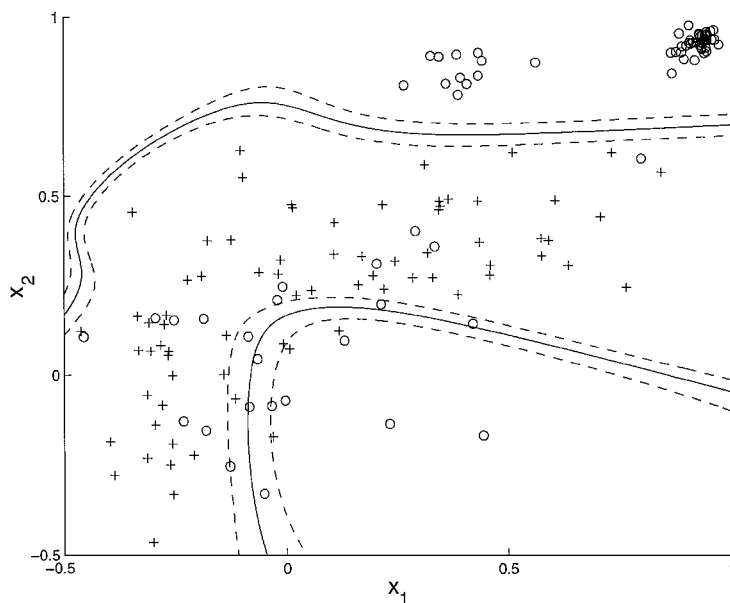


Figure 5. Classification boundaries (—) and confidence intervals (---) for the MLP classifier. The circles indicate patients, while the crosses represent the control group.

on discriminants. An alternative and more principled approach, which we do not pursue here, is to use a logistic output neuron so that the classification scheme is based on probabilities of class membership.

The size of the confidence intervals for the decision boundary is given by the noise variance (σ^2). These intervals are a measure of uncertainty on the

threshold that we apply to the linear output neuron. Our confidence of correctly classifying a sample occurring within these intervals should be very low. The receiver operating characteristic (ROC) curve, shown in Fig. 6, indicates that we can expect to detect patients with a 70% confidence without making any mistakes. The percentage of classification errors in the test set

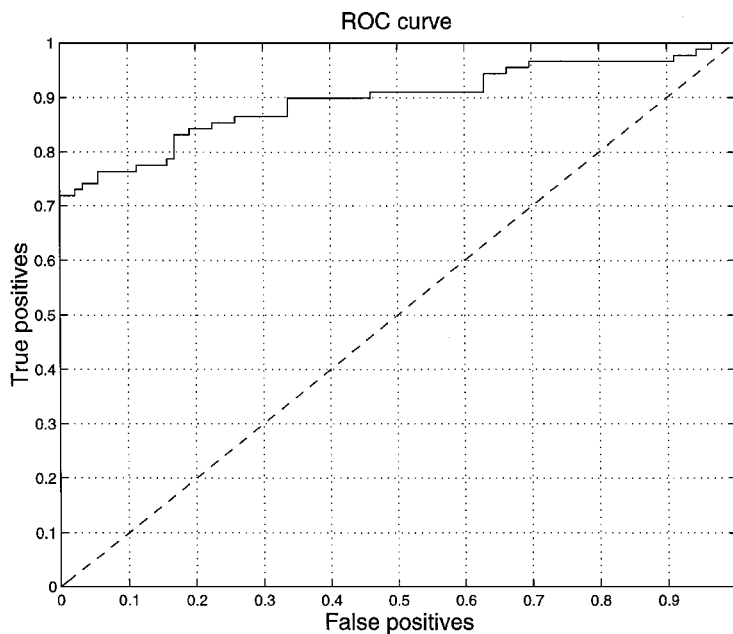


Figure 6. Receiver operating characteristic (ROC) of the classifier for the tremor test data.

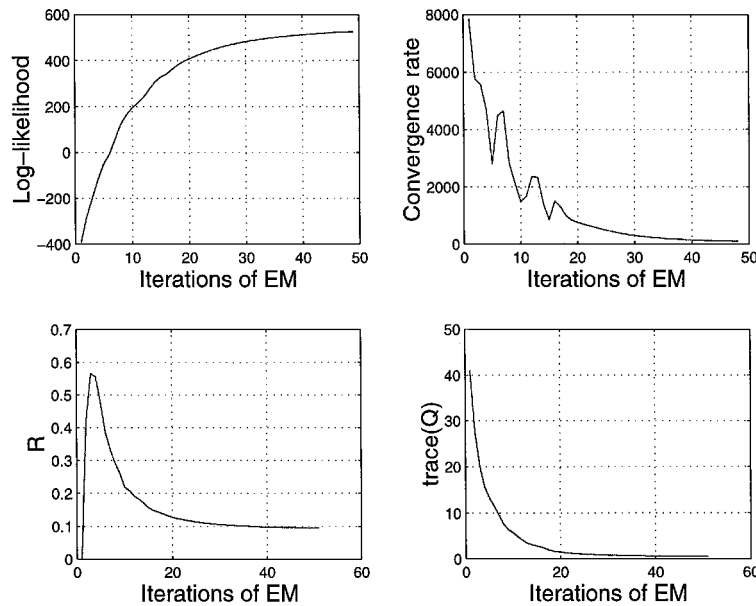


Figure 7. The top plots show the log-likelihood function and the convergence rate (log-likelihood slope) for the tremor data classification problem. The bottom plots show the convergence of the measurements noise covariance R and the trace of the process noise covariance Q .

was found to be 15.17. This error is of the same magnitude as previous results [26]. Finally, the convergence properties of the EM algorithm for this application are illustrated in Fig. 7.

7. Conclusions

In this paper, we derived an EM algorithm to estimate the neural network weights, measurement noise and model uncertainty (in the case of stationary data) jointly. We found that it performs well in terms of model accuracy and generalisation ability. Further research avenues include extending the method to other types of noise processes, establishing theoretical convergence bounds and investigating ways of efficiently initialising the algorithm so as to avoid local minima. The latter problem can be circumvented, to a certain extent, by stochastic global optimisation algorithms. Finally, we have made the software available at <http://www.cs.berkeley.edu/~jfgf>.

Acknowledgments

We would like to thank Mari Ostendorf (Boston University), Vassilis Digalakis (Technical University of Crete), David Melvin (Cambridge Clinical School),

Ben North (Oxford University) and the anonymous reviewers for their valuable comments. We are very grateful to Zoubin Ghahramani (University of Toronto) for his help. We would also like to thank David Stoffer (University of Pittsburgh) for making available some of his technical reports and Stephen Roberts and Will Penny (Imperial College of London) for the tremor data.

João FG de Freitas was financially supported by two University of the Witwatersrand Merit Scholarships, a Foundation for Research Development Scholarship (South Africa), an ORS Award (UK) and a Trinity College External Research Studentship (Cambridge).

Notes

1. We adopt the notation $\mathbf{z}_{1:N} \triangleq \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$. In addition, we suppress the input variables \mathbf{x} in the arguments of the probability distributions.
2. The data set can be found at <http://wol.ra.phy.cam.ac.uk/mackay/>.
3. The data is available at <http://www.ee.ic.ac.uk/hp/staff/sroberts.html>.

References

1. A.P. Dempster, N.M. Laird, and D.B. Rubin, "Maximum Likelihood from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society Series B*, vol. 39, 1977, pp. 1–38.

2. C.F. Chen, "The EM Algorithm to the Multiple Indicators and Multiple Causes Model via the Estimation of the Latent Variable," *Journal of the American Statistical Association*, vol. 76, no. 375, 1981, pp. 704–708.
3. M.W. Watson and R.F. Engle, "Alternative Algorithms for the Estimation of Dynamic Factor, MIMIC and Varying Coefficient regression Models," *Journal of Econometric*, vol. 23, no. 3, 1983, pp. 385–400.
4. R.H. Shumway and D.S. Stoffer, "An Approach to Time Series Smoothing and Forecasting Using the EM Algorithm," *Journal of Time Series Analysis*, vol. 3, no. 4, 1982, pp. 253–264.
5. R.H. Shumway and D.S. Stoffer, "Dynamic Linear Models with Switching," *Journal of the American Statistical Association*, vol. 86, no. 415, 1991, pp. 763–769.
6. V. Digalakis, J.R. Rohlicek, and M. Ostendorf, "ML Estimation of a Stochastic Linear System with the EM Algorithm and its Application to Speech Recognition," *IEEE Transactions on Speech and Audio Processing*, vol. 1, no. 4, 1993, pp. 431–442.
7. B. North and A. Blake, "Learning Dynamical Models using Expectation-Maximisation," in *International Conference on Computer Vision*, Mumbai, India, 1998, pp. 384–389.
8. R.P.N. Rao and D.H. Ballard, "Dynamic Model of Visual Recognition Predicts Neural Response Properties in the Visual Cortex," *Neural Computation*, vol. 9, no. 4, 1997, pp. 721–763.
9. Z. Ghahramani, "Learning Dynamic Bayesian Networks," in *Adaptive Processing of Temporal information*, C.L. Giles and M. Gori (Eds.), Springer-Verlag. Lecture Notes in Artificial Intelligence, vol. 1387.
10. S. Roweis and Z. Ghahramani, "A Unifying Review of Linear Gaussian Models," *Neural Computation*, vol. 11, no. 2, 1999, pp. 305–345.
11. S. Singhal and L. Wu, "Training Multilayer perceptrons with the Extended Kalman Algorithm," in *Advances in Neural Information Processing Systems*, vol. 1, D.S. Touretzky (Ed.), San Mateo, CA, 1988, pp. 133–140.
12. S. Shah, F. Palmieri, and M. Datum, "Optimal Filtering Algorithms for Fast Learning in Feedforward Neural Networks," *Neural Networks*, vol. 5, no. 5, 1992, pp. 779–787.
13. G.V. Puskorius and L.A. Feldkamp, "Decoupled Extended Kalman Filter Training of Feedforward Layered Networks," in *International Joint Conference on Neural Networks*, Seattle, 1991, pp. 307–312.
14. A. Gelb (Ed.), *Applied Optimal Estimation*, MIT Press, 1974.
15. A.H. Jazwinski, *Stochastic processes and Filtering Theory*, Academic Press, 1970.
16. H.E. Rauch, F. Tung, and C.T. Striebel, "Maximum Likelihood Estimates of Linear Dynamic Systems," *AIAA Journal*, vol. 3, no. 8, 1965, pp. 1445–1450.
17. J.F.G. de Freitas, M. Niranjan, and A.H. Gee, "Hierarchical bayesian-Kalman Models for Regularisation and ARD in Sequential Learning," Technical Report CUED/F-INFENG/TR 307, Cambridge University Engineering Department, 1997.
18. A. Graham, *Kronecker Products and Matrix Calculus with Applications*, Ellis Horwood Limited, 1981.
19. C. Andrieu, J.F.G. de Freitas, and A. Doucet, "Robust Full Bayesian Learning for Neural Networks," Technical Report CUED/F-INFENG/TR 343, Cambridge University Engineering Department, 1999.
20. C.C. Holmes and B.K. Mallick, "Bayesian Radial Basis Functions of Variable Dimension," *Neural Computation*, vol. 10, no. 5, 1998, pp. 1217–1233.
21. D.J.C. Mackay, "A Practical bayesian Framework for Backpropagation Networks," *Neural Computation*, vol. 4, no. 3, 1992, pp. 448–472.
22. R.M. Neal, *Bayesian Learning for Neural Networks*, Springer-Verlag, New York, 1996. Lecture Notes in Statistics, vol. 118.
23. D. Rios Insua and P. Müller, "Feedforward Neural Networks for Nonparametric Regression," *Technical Report 98-02*, Institute of Statistics and Decision Sciences, Duke University, 1998.
24. S.J. Roberts, W.D. Penny, and D. Pillot, "Novelty, Confidence and Errors in Connectionist Systems," in *IEE Colloquium on Intelligent Sensors and Fault Detection*, vol. 261, 1996, pp. 10/1–10/6.
25. J.M. Spyers-Ashby, P. Bain, and S.J. Roberts, "A Comparison of Fast Fourier Transform (FFT) and Autoregressive (AR) Spectral Estimation Techniques for the Analysis of Tremor Data," *Journal of Neuroscience Methods*, vol. 83, no. 1, 1998, pp. 35–43.
26. S.J. Roberts and W.D. Penny, "Bayesian Neural Network for Classification: How Useful is the Evidence Framwork?" *Neural Networks*, vol. 12, 1999, pp. 877–892.



Nando de Freitas was awarded an honours B.S.c. degree (with distinction) in electrical engineering at the University of the Witwatersrand (Johannesburg) in 1994. He received an M.S.c. degree (with distinction) at the same place in 1996. In 1999, he obtained a PhD degree in information engineering at the University of Cambridge. He is currently a post-doctoral research scholar at the computer science division of the University of California, Berkeley. His general research interests are in machine learning, Bayesian inference, pattern recognition, non-stationary signal processing, particle filters, Markov chain Monte Carlo simulation and neural networks. jfgf@cs.berkeley.edu



Mahesan Niranjan received a B.Sc. degree from the University of Peradeniya, Sri Lanka (1982) and a Masters from the Netherlands Universities Foundation for International Cooperation, Eindhoven,

The Netherlands (1985), and a Ph.D. degree from the University of Cambridge, England (1990). He was appointed to a Lectureship at Cambridge University and a Fellowship at Robinson College, Cambridge in October 1990. In January 1999 he was appointed to a Professorship in Computer Science at the University of Sheffield, England. His research interests are in the development and critical evaluation of algorithms for machine learning, with application to problems in nonstationary signal processing, medical informatics and computational finance.

M.Niranjan@dcs.shef.ac.uk



Andrew Gee obtained a BA in Electrical and Information Science at the University of Cambridge in 1990. He was awarded the PhD

degree in 1993 from the same university for his thesis on the use of feedback neural networks for combinatorial optimization. He is currently a Lecturer in the Engineering Department of the University of Cambridge, and also a Fellow of Queens' College, Cambridge. He teaches undergraduate courses on all aspects of computing: in 1998 he was awarded a Pilkington Prize by the Trustees of the Cambridge Foundation for excellence in teaching at the University of Cambridge. His current research interests are in medical imaging, computer vision and neural networks, and include three dimensional ultrasound imaging, nonlinear sequential learning and document image processing. He has authored and co-authored around 70 papers. ahg@eng.cam.ac.uk