

An EM Algorithm for Identification of Nonlinear Dynamical Systems

Sam Roweis, Zoubin Ghahramani

Abstract— We provide a novel solution to the problem of simultaneously estimating the unknown parameters and hidden states of a nonlinear dynamical system. Our solution is based on the expectation–maximization (EM) algorithm, an iterative procedure for maximum likelihood parameter estimation from data sets with missing or hidden variables [1]. EM has been applied to system identification in linear state-space models, where the state variables are hidden from the observer and both the state and the parameters of the model have to be estimated simultaneously [2], [3], [4]. Here we generalize the EM algorithm to estimate parameters of nonlinear dynamical state-space models. The “expectation” step makes use of Extended Kalman Smoothing to estimate the state, while the “maximization” step re-estimates the parameters using these uncertain state estimates. In general, the nonlinear maximization step is difficult because it requires integrating out the uncertainty in the states. However, if Gaussian radial basis function (RBF) approximators are used to model the nonlinearities, the integrals become tractable and the maximization step can be solved via systems of linear equations. We derive an online version of this EM-EKS algorithm, as well as a version for non-stationary time series. We also consider the identifiability and expressive power of nonlinear dynamical systems and relate our learning algorithm with more traditional system identification procedures based on dual and joint Extended Kalman Filtering. Finally, we demonstrate our algorithm on several synthetic problems and one real time series.

I. LEARNING STOCHASTIC NONLINEAR DYNAMICS

SINCE the advent of cybernetics, dynamical systems have been an important modeling tool in fields ranging from engineering to the physical and social sciences. Most realistic dynamical systems models have two essential features. First, they are stochastic—the observed outputs are a noisy function of the inputs, and the dynamics itself may be driven by some unobserved noise process. Second, they can be characterized by some finite-dimensional internal state which, while not directly observable, summarizes at any time all information about the past behaviour of the process relevant to predicting its future evolution.

From a modeling standpoint, stochasticity is essential to allow a model with a few fixed parameters to generate a rich variety of time-series outputs.¹ Explicitly modeling the internal state makes it possible to decouple the internal dynamics from the observation process. For example, to model a sequence of video images of a balloon floating

in the wind, it would be computationally very costly to directly predict the array of camera pixel intensities from a sequence of arrays of previous pixel intensities. It seems much more sensible to attempt to infer the true state of the balloon (its position, velocity, and orientation) and decouple the process which governs the balloon dynamics from the observation process which maps the actual balloon state to an array of measured pixel intensities.

Often we are able to write down equations governing these dynamical systems directly, based on prior knowledge of the problem structure and the sources of noise—for example, from the physics of the situation. In such cases, we may want to infer the hidden state of the system from a sequence of observations of the system’s inputs and outputs. Solving this *inference* or *state estimation* problem is essential for tasks such as tracking or the design of state-feedback controllers, and there exist well-known algorithms for this.

However in many cases the exact parameter values, or even the gross structure of the dynamical system itself may be unknown. In such cases the dynamics of the system have to be *learned* or *identified* from sequences of observations only. Learning may be a necessary precursor if the ultimate goal is effective state inference. But learning nonlinear state-based models is also useful in its own right, even when we are not explicitly interested in the internal states of the model, for tasks such as prediction (extrapolation), time-series classification, outlier detection, and filling-in of missing observations (interpolation). This chapter addresses the problem of learning time series models when the internal state is hidden. Below, we briefly review the two fundamental algorithms that form the basis of our learning procedure. In section II we introduce our algorithm and derive its learning rules. Section III presents results of using the algorithm to identify nonlinear dynamical systems. Finally, we present some conclusions and potential extensions to the algorithm in sections IV and V.

A. State inference and model learning

Two remarkable algorithms from the 1960’s, one developed in engineering and the other in statistics, form the basis of modern techniques in state estimation and model learning. The Kalman filter, introduced by Kalman and Bucy in 1961 [5], was developed in a setting where the physical model of the dynamical system of interest was readily available; its goal is optimal state estimation in systems with known parameters. The expectation-maximization (EM) algorithm, pioneered by Baum and colleagues [6] and later generalized and named by Dempster *et.al* [1], was de-

Gatsby Computational Neuroscience Unit, University College London, London WC1N 3AR, U.K. <http://www.gatsby.ucl.ac.uk>

¹There are, of course, completely deterministic but *chaotic* systems with this property. If we separate the noise processes in our models from the deterministic portions of the dynamics and observations, we can think of the noises as another deterministic (but highly chaotic) system which depends on initial conditions and exogenous inputs that we do not know. Indeed when we run simulations using a pseudo-random number generator started with a particular seed this is precisely what we are doing.

veloped to learn parameters of statistical models in the presence of incomplete data or hidden variables.

In this chapter we bring together these two algorithms in order to learn the dynamics of stochastic nonlinear systems with hidden states. Our goal is two-fold: both to develop a method for identifying the dynamics of nonlinear systems whose hidden states we wish to infer, and to develop a general nonlinear time-series modeling tool. We examine inference and learning in discrete-time² stochastic nonlinear dynamical systems with hidden states x_t , external inputs u_t , and noisy outputs y_t . (All lowercase characters (except indices) denote vectors. Matrices are represented by uppercase characters.) The systems are parametrized by a set of tunable matrices, vectors and scalars which we shall collectively denote as θ . The inputs, outputs and states are related to each other by

$$x_{t+1} = f(x_t, u_t) + w_t \quad (1a)$$

$$y_t = g(x_t, u_t) + v_t \quad (1b)$$

where w_t and v_t are zero mean Gaussian noise processes. The state vector x evolves according to a nonlinear but stationary Markov dynamics³ driven by the inputs u and by the noise source w . The outputs y are nonlinear, noisy but stationary and instantaneous functions of the current state and current input. The vector-valued nonlinearities f and g are assumed to be differentiable, but otherwise arbitrary.

Models of this kind have been examined for decades in systems and control engineering. They can also be viewed within the framework of *probabilistic graphical models* which use graph theory to represent the conditional dependencies between a set of variables [7], [8]. A graphical model is a diagram corresponding to a probabilistic generative model; it has a node for each possibly vector valued random variable, with directed arcs representing stochastic dependencies. Absent connections indicate conditional independence. In particular, nodes are conditionally independent from their non-descendants, given their parents—where parents, children, descendants, etc, are defined with respect to the directionality of the arcs (i.e. from parent to child). We can capture the dependencies in equations (1) compactly by drawing its graphical model as shown in figure 1.

Graphical models have helped clarify the relationship between dynamical systems and other probabilistic models such as hidden Markov models and factor analysis [9]. Graphical models have also made it possible to develop

²Continuous time dynamical systems (in which *derivatives* are specified as functions of the current state and inputs) can be converted into discrete time systems by sampling their outputs and using “zero-order holds” on their inputs. In particular, for a continuous-time linear system $\dot{x}(t) = A_c x(t) + B_c u(t)$ sampled at interval τ the corresponding dynamics and input driving matrices so that $x_{t+1} = A x_t + B u_t$ are $A = \sum_{k=0}^{\infty} A_c^k \tau^k / k! = \exp(A_c \tau)$ and $B = A_c^{-1} (A - I) B_c$.

³Stationarity means here that neither f , nor the covariance of the noise process w_t depend on time; that is, the dynamics are *time invariant*. Markov refers to the fact that given the current state the next state does not depend on the past history of the states.

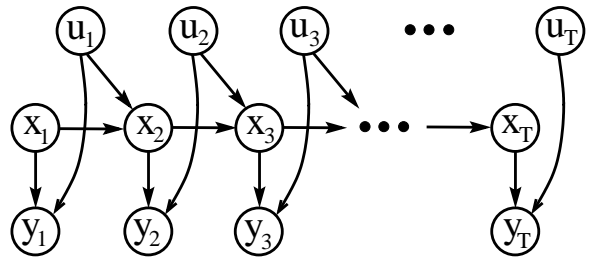


Fig. 1. A probabilistic graphical model for stochastic dynamical systems with hidden states x_t , inputs u_t and observables y_t .

probabilistic inference algorithms that are vastly more general than the Kalman filter.

Since our learning algorithm for nonlinear dynamical systems combines the key ideas behind Kalman filtering and EM, we review these in the following two sections. The goal is to develop an algorithm which can be used to model the probability density of output sequences (or the conditional density of outputs given inputs) using only a finite number of example time-series. The crux of the problem is that both the hidden state trajectory and the parameters are unknown.

If we knew the parameters, the operation of interest would be to infer the hidden state sequence. The uncertainty in this sequence would be encoded by computing the posterior distributions of the hidden state variables given the sequence of observations. The Kalman filter provides a solution to this problem in the case where f and g are linear. If, on the other hand, we had access to the hidden state trajectories as well as to the observables, then the problem would be one of model fitting, i.e. estimating the parameters of f and g and the noise covariances. Given observations of the (no longer hidden) states and outputs, f and g can be obtained as the solution to a possibly nonlinear regression problem, and the noise covariances can be obtained from the residuals of the regression. How should we proceed when *both* the system model and the hidden states are unknown?

The classical approach to solving this problem is to treat the parameters θ as “extra” hidden variables, and to apply an Extended Kalman Filtering algorithm (described below) to the nonlinear system with the state vector augmented by the parameters [10], [11]. For stationary models, the dynamics of the parameter portion of this extended state vector are set to the identity function. This approach can be made inherently on-line, which may be important in certain applications. Furthermore, it provides an estimate of the covariance of the parameters at each time step. Finally, its objective, probabilistically speaking, is to find an optimum in the joint space of parameters and hidden state sequences.

In contrast, the algorithm we present is a batch algorithm (although as we discuss in section IV-B, online extensions are possible) and does not attempt to estimate the covariance of the parameters. Like other instances of the EM algorithm, which we describe below, its goal is to integrate over the uncertain estimates of the unknown hid-

den states and optimize the resulting marginal likelihood of the parameters given the observed data. An Extended Kalman Smoother (EKS) is used to estimate the approximate state distribution in the E-step, and a radial basis function (RBF) network [12] is used for nonlinear regression in the M-step. It is important not to confuse this use of the Extended Kalman algorithm, to estimate just the hidden state as part of the E-step of EM, with the use we described in the previous paragraph, to simultaneously estimate parameters and hidden states.

B. The Kalman Filter

Linear dynamical systems with additive white Gaussian noises are the most basic models to examine when considering the state estimation problem because they admit exact and efficient inference. (Here, and in what follows we call a system linear if both the state evolution function and the state-to-output observation function are linear, and nonlinear otherwise.) The linear dynamics and observation processes correspond to matrix operations which we denote with A, B and C, D respectively, giving the classic state-space formulation of input-driven linear dynamical systems:

$$x_{t+1} = Ax_t + Bu_t + w_t \quad (2a)$$

$$y_t = Cx_t + Du_t + v_t \quad (2b)$$

The Gaussian noise vectors w and v have zero mean and covariances Q and R respectively. If the prior probability distribution $p(x_1)$ over initial states is taken to be Gaussian, then the joint probabilities of all states and outputs at future times are also Gaussian, since the Gaussian distribution is closed under the linear operations applied by state evolution and output mapping and under the convolution applied by additive Gaussian noise. Thus, all distributions over hidden state variables are fully described by their means and covariance matrices. The algorithm for exactly computing the posterior mean and covariance for x_t given some sequence of observations consists of two parts: a forward recursion which uses the observations from y_1 to y_t , known as the *Kalman filter* [13], and a backward recursion which uses the observations from y_T to y_{t+1} . The combined forward and backward recursions are known as the Kalman or Rauch-Tung-Streifel (RTS) *smoother* [14].

There are three key insights to understanding the Kalman filter.

The first insight is that the Kalman filter is simply a method for implementing Bayes rule. Consider the very general setting where we have a prior $p(x)$ on some state variable and an observation model $p(y|x)$ for the noisy outputs given the state. Bayes rule gives us the state inference procedure:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x)}{Z} \quad (3a)$$

$$Z = p(y) = \int_x p(y|x)p(x)dx \quad (3b)$$

where the normalizer Z is the unconditional density of the observation. All we need to do to convert our prior on the

state into a posterior is multiply by the likelihood from the observation equation and renormalize.

The second insight is that there is no need to invert the output or dynamics functions, as long as we work with easily normalizable distributions over hidden states. We see this by applying Bayes rule to the linear Gaussian case for a single time step.⁴ We start with a Gaussian belief $\mathcal{N}(x_{t-1}, V_{t-1})$ on the current hidden state, use the dynamics to convert this to a prior $\mathcal{N}(x^+, V^+)$ on the next state, and then condition on the observation to convert this prior into a posterior $\mathcal{N}(x_t, V_t)$. This gives the classic Kalman filtering equations:

$$p(x_{t-1}) = \mathcal{N}(x^+, V^+) \quad (4a)$$

$$x^+ = Ax_{t-1}, \quad V^+ = AV_{t-1}A^\top + Q \quad (4b)$$

$$p(y_t|x_t) = \mathcal{N}(Cx_t, R) \quad (4c)$$

$$p(x_t|y_t) = \mathcal{N}(x_t, V_t) \quad (4d)$$

$$x_t = x^+ + K(y_t - Cx^+), \quad V_t = (I - KC)V^+ \quad (4e)$$

$$K = V^+C^\top (CV^+C^\top + R)^{-1} \quad (4f)$$

The posterior is again Gaussian and analytically tractable. Notice that neither the dynamics matrix A nor the observation matrix C needed to be inverted.

The third insight is that the state estimation procedures can be implemented recursively. The posterior from the previous time step is run through the dynamics model and becomes our prior for the current time step. We then convert this prior into a new posterior by using the current observation.

For the general case of a nonlinear system with non-Gaussian noise, state estimation is much more complex. In particular, mapping through arbitrary nonlinearities f and g can result in arbitrary state distributions, and the integrals required for Bayes rule can become intractable. Several methods have been proposed to overcome this intractability, each providing a distinct approximate solution to the inference problem. Assuming f and g are differentiable and the noise is Gaussian, one approach is to locally linearize the nonlinear system about the current state estimate so that applying the Kalman filter to the linearized system the approximate state distribution remains Gaussian. Such algorithms are known as *extended Kalman filters* (EKF) [15], [16]. The EKF has been used both in the classical setting of state estimation for nonlinear dynamical systems, and also as a basis for online learning algorithms for feedforward neural networks [17] and radial basis function networks [18], [19].

Another possibility is to propagate a set of discrete samples in state space through f and g , and to re-weight them using the likelihood $p(y|x)$. Algorithms which use this general strategy are known as *particle filters* [20]; a particular form of which is known as the CONDENSATION algorithm.

⁴Some notation: A multivariate normal (Gaussian) distribution with mean μ and covariance matrix Σ is written as $\mathcal{N}(\mu, \Sigma)$. The same Gaussian evaluated at the point z is denoted $\mathcal{N}(\mu, \Sigma)|_z$. The determinant of a matrix is denoted by $|A|$ and matrix inversion by A^{-1} . The symbol \sim means “distributed according to”.

m [21]. (These filters have also been called Monte Carlo filters [22], bootstrap filters [23], and dynamic mixture models [24], [25]. See the book [26] for a recent survey.) A third approximation called the *unscented* filter [27], [28], [29] deterministically chooses a set of balanced points and propagates them through the nonlinearities in order to recursively approximate a Gaussian state distribution. Finally, there are algorithms for approximate inference and learning based on mean field theory and variational methods [30], [31]. Although we have chosen to make local linearization (EKS) the basis of our algorithms below, it is possible to formulate the same learning algorithms using any approximate inference method, for example, the unscented filter.

C. The EM Algorithm

The EM or *expectation-maximization* algorithm [32], [1] is a widely applicable iterative parameter re-estimation procedure. The objective of the EM algorithm is to maximize the likelihood of the observed data $P(Y|\theta)$ in the presence of hidden⁵ variables x . (We will denote the entire sequence of observed data by $Y = \{y_1 \dots y_\tau\}$, of hidden variables by $X = \{x_1 \dots x_\tau\}$, and the parameters of the model by θ .) Maximizing the likelihood as a function of θ is equivalent to maximizing the log likelihood:

$$\mathcal{L}(\theta) = \log P(Y|\theta) = \log \int_X P(X, Y|\theta) dX \quad (5)$$

Using *any* distribution $Q(X)$ over the hidden variables, we can obtain a lower bound on \mathcal{L} :

$$\log \int_X P(Y, X|\theta) dX = \log \int_X Q(X) \frac{P(X, Y|\theta)}{Q(X)} dX \quad (6a)$$

$$\geq \int_X Q(X) \log \frac{P(X, Y|\theta)}{Q(X)} dX \quad (6b)$$

$$= \int_X Q(X) \log P(X, Y|\theta) dX - \int_X Q(X) \log Q(X) dX \quad (6c)$$

$$= \mathcal{F}(Q, \theta) \quad (6d)$$

where the middle inequality is known as Jensen's inequality and can be proven using the concavity of the log function. If we define the *energy* of a global configuration (X, Y) to be $-\log P(X, Y|\theta)$, then the lower bound $\mathcal{F}(Q, \theta) \leq \mathcal{L}(\theta)$ is the negative of a quantity known in statistical physics as the *free energy*: the expected energy under Q minus the entropy of Q [33]. The EM algorithm alternates between maximizing \mathcal{F} with respect to the *distribution* Q and the *parameters* θ , respectively, holding the other fixed. Starting from some initial parameters θ_0 we alternately apply:

$$\text{E-step:} \quad Q_{k+1} \leftarrow \arg \max_Q \mathcal{F}(Q, \theta_k) \quad (7a)$$

$$\text{M-step:} \quad \theta_{k+1} \leftarrow \arg \max_\theta \mathcal{F}(Q_{k+1}, \theta) \quad (7b)$$

⁵Often called *missing data* or *auxiliary parameters*.

It is easy to show that the maximum in the E-step results when Q is exactly the conditional distribution of X , $Q_{k+1}^*(X) = P(X|Y, \theta_k)$, at which point the bound becomes an equality: $\mathcal{F}(Q_{k+1}^*, \theta_k) = \mathcal{L}(\theta_k)$. The maximum in the M-step is obtained by maximizing the the first term in (6c), since the entropy of Q does not depend on θ :

$$\text{M-step:} \quad \theta_{k+1}^* \leftarrow \arg \max_\theta \int_X P(X|Y, \theta_k) \log P(X, Y|\theta) dX. \quad (8)$$

This is the expression most often associated with the EM algorithm, but it obscures the elegant interpretation [33] of EM as coordinate ascent in \mathcal{F} (see figure 2). Since $\mathcal{F} = \mathcal{L}$ at the beginning of each M-step, and since the E-step does not change θ , we are guaranteed not to decrease the likelihood after each combined EM step. (While this is obviously true of "complete" EM algorithms as described above, it may also be true for "incomplete" or "sparse" variants in which approximations are used during the E- and/or M-steps so long as \mathcal{F} always goes up; see also the earlier work of [34].) For example, this can take the form of a gradient M-step algorithm (where we increase $P(Y|\theta)$ with respect to θ but do not strictly maximize it), or any E-step which improves the bound \mathcal{F} without saturating it [33].)

In dynamical systems with hidden states, the E-step corresponds exactly to solving the smoothing problem: estimating the hidden state trajectory given both the observations/inputs and the parameter values. The M-step involves system identification using the state estimates from the smoother. Therefore, at the heart of the EM learning procedure is the following idea: use the solutions to the filtering/smoothing problem to estimate the unknown hidden states given the observations and the current model parameters. Then use this fictitious complete data to solve for new model parameters. Given the estimated states obtained from the inference algorithm, it is usually easy to solve for new parameters. For example, when working with linear Gaussian models this typically involves minimizing

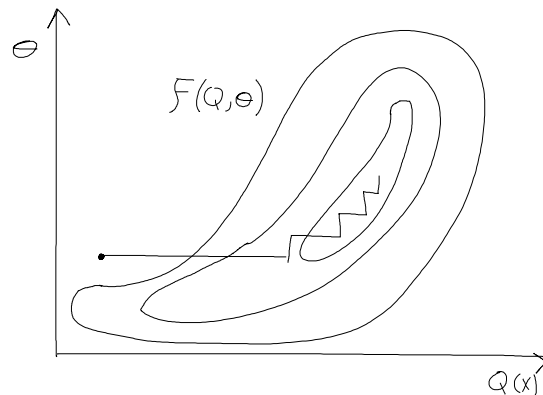


Fig. 2. The EM algorithm can be thought of as coordinate ascent in the functional $\mathcal{F}(Q(X), \theta)$ (see text). The E-step maximizes \mathcal{F} with respect to $Q(X)$ given fixed θ (horizontal moves) while the M-step maximizes \mathcal{F} with respect to θ given fixed $Q(X)$ (vertical moves).

quadratic forms which can be done with linear regression. This process is repeated, using these new model parameters to infer the hidden states again, and so on. Keep in mind that our goal is to maximize the log-likelihood (5) (or equivalently maximize the total likelihood) of the observed data with respect to the model parameters. This means integrating (or summing) over all ways in which the model could have produced the data (i.e. hidden state sequences). As a consequence of using the EM algorithm to do this maximization we find ourselves needing to compute (and maximize) the *expected* log likelihood of the *joint* data (8), where the expectation is taken over the distribution of hidden values predicted by the current model parameters and the observations.

In the past, the EM algorithm has been applied to learning linear dynamical systems in specific cases (such as multiple-indicator multiple-cause (MIMC) models with a single latent variable [2] or state space models with the observation matrix known [3]) as well more generally [4]. This chapter is an extension of our earlier work [35]; since then, there has been other similar work applying EM to nonlinear dynamical systems [36], [37]. Whereas other work uses sampling for the E-step and gradient M-steps, our algorithm uses the RBF networks to obtain a computationally efficient and exact M-step.

There are four important advantages the EM algorithm has over classical approaches. First, the EM algorithm provides a straightforward and principled method for handling missing inputs or outputs. (Indeed this was the original motivation for Shumway and Stoffer's application of the EM algorithm to learning of partially unknown linear dynamical systems [3].) Second, EM generalizes readily to more complex models with combinations of discrete and real-valued hidden variables. For example, one can formulate EM for a *mixture* of nonlinear dynamical systems [38], [39]. Third, whereas it is often very difficult to prove or analyze stability within the classical on-line approach, the EM algorithm is always attempting to maximize the likelihood, which acts as a Lyapunov function for stable learning. Fourth, the EM framework facilitates Bayesian extensions to learning, for example through the use of variational approximations [31].

II. COMBINING EKS AND EM

In the next sections we will describe the basic components of our EM learning algorithm. For the expectation step of the algorithm, we infer an approximate conditional distribution of the hidden states using Extended Kalman Smoothing (section II-A). For the maximization step we first discuss the general case (section II-B) and then describe the particular case where the nonlinearities are represented using Gaussian radial basis function (RBF) networks (section II-C). Since, as with all EM or likelihood ascent algorithms, our algorithm is not guaranteed to find the globally optimum solutions, good initialization is a key factor in practical success. We typically use a variant of factor analysis followed by estimation of a purely linear dynamical system as the starting point for training our nonlinear

models (section II-D).

A. Extended Kalman smoothing (E-step)

Given a system described by equations (1), the E-step of an EM learning algorithm needs to infer the hidden states from a history of observed inputs and outputs. The quantities at the heart of this inference problem are two conditional densities

$$P(x_t | u_1, \dots, u_T, y_1, \dots, y_T), \quad 1 \leq t \leq T \quad (9)$$

$$P(x_t, x_{t+1} | u_1, \dots, u_T, y_1, \dots, y_T), \quad 1 \leq t \leq T - 1 \quad (10)$$

For nonlinear systems these conditional densities are in general non-Gaussian and can in fact be quite complex. For all but a very few nonlinear systems exact inference equations cannot be written down in closed form. Furthermore, for many nonlinear systems of interest exact inference is intractable (even numerically) meaning that in principle the amount of computation required grows exponentially in the length of the time series observed. The intuition behind all extended Kalman algorithms is that they approximate a stationary nonlinear dynamical system with a *non-stationary* (time-varying) but linear system. In particular, Extended Kalman Smoothing (EKS) simply applies regular Kalman smoothing to a local linearization of the nonlinear system. At every point \hat{x} in x -space, the derivatives of the vector-valued functions f and g define the matrices, $A_{\hat{x}} \equiv \left. \frac{\partial f}{\partial x} \right|_{x=\hat{x}}$ and $C_{\hat{x}} \equiv \left. \frac{\partial g}{\partial x} \right|_{x=\hat{x}}$, respectively. The dynamics are linearized about \hat{x}_t , the mean of the current filtered (not smoothed) state estimate at time t . The output equation can be similarly linearized. These linearizations yield:

$$x_{t+1} = f(\hat{x}_t, u_t) + A_{\hat{x}_t} (x_t - \hat{x}_t) + w \quad (11)$$

$$y_t = g(\hat{x}_t, u_t) + C_{\hat{x}_t} (x_t - \hat{x}_t) + v. \quad (12)$$

If the noise distributions and the prior distribution of the hidden state at $t = 1$ are Gaussian, then, in this progressively linearized system, the conditional distribution of the hidden state at any time t given the history of inputs and outputs will also be Gaussian. Thus, Kalman smoothing can be used on the linearized system to infer this conditional distribution; this is illustrated in figure 3.

Notice that although the algorithm performs *smoothing* (in other words, it takes into account all observations including future ones when inferring the state at any time) the linearization is only done in the forward direction. Why not re-linearize about the backwards estimates during the Rauch recursions? While in principle this approach might give better results it is difficult to implement in practice because it requires the dynamics function to be uniquely invertible, which is often not true.

Unlike the normal (linear) Kalman smoother, in the EKS the error covariances for the state estimates and the Kalman gain matrices *do* depend on the observed data, not just on the time index t . Furthermore, it is no longer true that if the system is stationary the Kalman gain will converge to a value which makes the smoother act as the

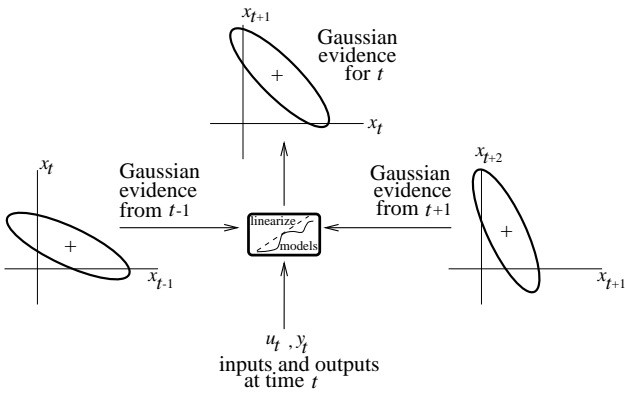


Fig. 3. Illustration of the information used in Extended Kalman Smoothing (EKS), which infers the hidden state distribution during the E-step of our algorithm. The nonlinear model is linearized about the current state estimate at each time, and then Kalman Smoothing is used on the linearized system to infer Gaussian state estimates.

optimal Wiener filter in the steady state. (In fact, the Kalman gain need not converge to a fixed value at all.)

B. Learning model parameters (M-step)

The M-step of our EM algorithm re-estimates the parameters of the model given the observed inputs, outputs, and the conditional distributions over the hidden states. For the model we have described, the parameters define the nonlinearities f and g , and the noise covariances Q and R (as well as the mean and covariance of the initial state, x_1).

Two complications can arise in the M-step. First, fully re-estimating f and g in each M-step may be computationally expensive. For example, if they are represented by neural network regressors, a single full M-step would be a lengthy training procedure using backpropagation, conjugate gradients, or some other optimization method. To avoid this, one could use partial M-steps which increase but do not maximize the expected log likelihood (8), for example each consisting of one or a few gradient steps. However, this will in general make the fitting procedure much slower.

The second complication is that f and g have to be trained using the *uncertain* state estimates output by the EKS algorithm. This makes it difficult to apply standard curve-fitting or regression techniques. Consider fitting f , which takes as inputs x_t and u_t and outputs x_{t+1} . For each t , the conditional density estimated by EKS is a full-covariance Gaussian in (x_t, x_{t+1}) -space. So f has to be fit not to a set of data points but instead to a mixture of full-covariance Gaussians in input-output space (Gaussian “clouds” of data). Ideally, to follow the EM framework, this conditional density should be *integrated over* during the fitting process. Integrating over this type of data is non-trivial for almost any form of f . One simple but inefficient approach to bypass this problem is to draw a large sample from these Gaussian clouds of data and then fit f to these samples in the usual way. A similar situation occurs with the fitting of the output function g .

We present an alternative approach, which is to choose the form of the function approximator to make the integra-

tion easier. As we will show, using Gaussian radial basis function (RBF) networks [12] to model f and g allows us to do the integrals exactly and efficiently. With this choice of representation both of the above complications vanish.

C. Fitting Radial Basis Functions to Gaussian Clouds

We will present a general formulation of an RBF network from which it should be clear how to fit special forms for f and g . Consider the following nonlinear mapping from input vectors x and u to an output vector z :

$$z = \sum_{i=1}^I h_i \rho_i(x) + Ax + Bu + b + w, \quad (13)$$

where w is a zero-mean Gaussian noise variable with covariance Q , and ρ_i are RBFs defined below. This general mapping can be used in several ways to represent dynamical systems, depending on which of the input to hidden to output mappings are assumed to be nonlinear. Three examples are: (1) representing f using (13) with the substitutions $x \leftarrow x_t$, $u \leftarrow u_t$, and $z \leftarrow x_{t+1}$; (2) representing f using $x \leftarrow (x_t, u_t)$, $u \leftarrow \emptyset$, and $z \leftarrow x_{t+1}$; and (3) representing g using the substitutions $x \leftarrow x_t$, $u \leftarrow u_t$, and $z \leftarrow y_t$. (Indeed for different simulations we will use different forms.) The parameters are: the I coefficients h_i of the RBFs; the matrices A and B multiplying inputs x and u , respectively; and an output bias vector b , and the noise covariance Q . Each RBF is assumed to be a Gaussian in x -space, with center c_i and width given by the covariance matrix S_i :

$$\rho_i(x) = |2\pi S_i|^{-1/2} \exp \left\{ -\frac{1}{2} (x - c_i)^\top S_i^{-1} (x - c_i) \right\}. \quad (14)$$

For now we assume the centers and widths of the RBFs are fixed, although we discuss learning their locations in section IV

The goal is to fit this model to data (u, x, z) . The complication is that the data set comes in the form of a mixture of Gaussian distributions. Here we show how to analytically integrate over this mixture distribution to fit the RBF model.

Assume the data set is:

$$P(x, z, u) = \frac{1}{J} \sum_j \mathcal{N}_j(x, z) \delta(u - u_j). \quad (15)$$

That is, we observe samples from the u variables, each paired with a Gaussian “cloud” of data, \mathcal{N}_j , over (x, z) . The Gaussian \mathcal{N}_j has mean μ_j and covariance matrix C_j .

Let $\hat{z}_\theta(x, u) = \sum_{i=1}^I h_i \rho_i(x) + Ax + Bu + b$, where θ is the set of parameters. The log likelihood of a single fully observed data point under the model would be:

$$-\frac{1}{2} [z - \hat{z}_\theta(x, u)]^\top Q^{-1} [z - \hat{z}_\theta(x, u)] - \frac{1}{2} \ln |Q| + \text{const.}$$

Since the (x, z) values in the data set are uncertain, the maximum expected log likelihood RBF fit to the mixture

of Gaussian data is obtained by minimizing the following integrated quadratic form:

$$\min_{\theta, Q} \left\{ \sum_j \int_{x, z} \mathcal{N}_j(x, z) [z - \hat{z}_\theta(x, u_j)]^\top Q^{-1} [z - \hat{z}_\theta(x, u_j)] dx dz + J \ln |Q| \right\}. \quad (16)$$

We rewrite this in a slightly different notation, using angled brackets $\langle \cdot \rangle_j$ to denote expectation over \mathcal{N}_j , and defining

$$\begin{aligned} \theta &\equiv [h_1 \ h_2 \ \dots \ h_I \ A \ B \ b] \\ \Phi &\equiv [\rho_1(x) \ \rho_2(x) \ \dots \ \rho_I(x) \ x^\top \ u^\top \ 1]^\top. \end{aligned}$$

Then, the objective can be written

$$\min_{\theta, Q} \left\{ \sum_j \langle (z - \theta - \Phi)^\top Q^{-1} (z - \theta - \Phi) \rangle_j + J \ln |Q| \right\}. \quad (17)$$

Taking derivatives with respect to θ , premultiplying by $-Q^{-1}$, and setting to zero gives the linear equations $\sum_j \langle (z - \theta - \Phi) \Phi^\top \rangle_j = 0$, which we can solve for θ and Q :

$$\hat{\theta} = \left(\sum_j \langle z \Phi^\top \rangle_j \right) \left(\sum_j \langle \Phi \Phi^\top \rangle_j \right)^{-1}, \quad (18a)$$

$$\hat{Q} = \frac{1}{J} \left(\sum_j \langle z z^\top \rangle_j - \hat{\theta} \sum_j \langle \Phi z^\top \rangle_j \right). \quad (18b)$$

In other words, given the expectations in the angled brackets, the optimal parameters can be solved for via a set of linear equations. In appendix A we show that these expectations can be computed analytically and efficiently, which means that we can take full and exact M-steps. The derivation is somewhat laborious, but the intuition is very simple: the Gaussian RBFs multiply the Gaussian densities \mathcal{N}_j to form new unnormalized Gaussians in (x, y) -space. Expectations under these new Gaussians are easy to compute. This fitting algorithm is illustrated in figure 4.

Note that of the four advantages we mentioned previously for the EM algorithm—ability to handle missing observations, generalizability to extensions of the basic model, Bayesian approximations, and guaranteed stability through a Lyapunov function—we have had to forgo one. There is no guarantee that extended Kalman smoothing increases the lower bound on the true likelihood and therefore stability cannot be assured. Although in practice the algorithm rarely became unstable and the approximation worked well—in our experiments the likelihood increased monotonically and good density models were learned—it may be desirable to derive guaranteed-stable algorithms for certain special cases using lower-bound preserving variational approximations [31].

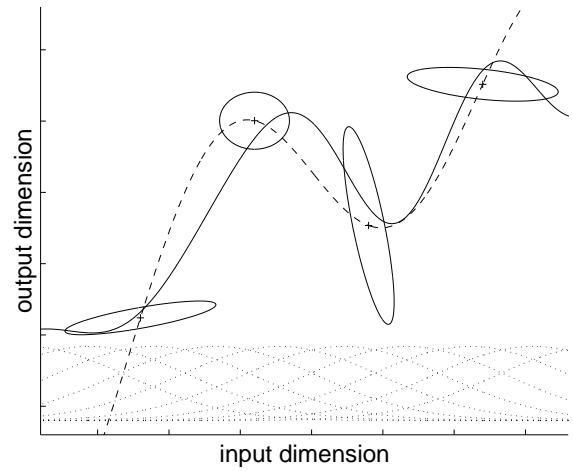


Fig. 4. Illustration of the regression technique employed during the M-step. A fit to a mixture of Gaussian densities is required; if Gaussian RBF networks are used then this fit can be solved analytically. The dashed line shows a regular RBF fit to the centres of the four Gaussian densities while the solid line shows the analytic RBF fit using the covariance information. The dotted lines below show the support of the RBF kernels.

The ability to fully integrate over uncertain state estimates provides practical benefits as well as being theoretically pleasing. We have compared fitting our RBF networks using only the means of the state estimates to performing the full integration as derived above. When using only the means, we found it necessary to introduce a ridge regression (weight decay) parameter in the M-step to penalize the very large coefficients that would otherwise occur based on precise cancellations between inputs. This ridge regression regularizer is like adding white noise to the radial basis outputs (i.e. *after* the RBF kernels have been applied). This is equivalent to Gaussian noise at the inputs with a covariance determined by the derivatives of the RBFs at the input locations. The uncertain state estimates provide exactly this sort of noise, and thus automatically regularize the RBF fit in the M-step. This naturally avoids the need to introduce a penalty on large coefficients and improves generalization.

D. Initialization of models and choosing locations for RBF kernels

The practical success of our algorithm depends on two design choices which need to be made at the beginning of the training procedure. The first is to judiciously select the placement of the RBF kernels in the representation of the the state dynamics and/or output function. The second is to sensibly initialize the parameters of the model so that iterative improvement with the EM algorithm (which finds only local maxima of the likelihood function) finds a good solution.

In models with low-dimensional hidden states, placement of RBF kernel centres can be done by gridding the state space and placing one kernel on each grid point. Since the scaling of the state variables is given by the covariance matrix of the state dynamics noise w_t (1), which without loss

of generality we have set to I , it is possible to determine both a suitable size for the gridding region over the state-space, and a suitable scaling of the RBF kernels themselves. However, the number of kernels in a grid increased exponentially with the grid dimension, so for more than 3 or 4 state variables gridding the state-space is impractical. In these cases we first use a simple initialization, such as a linear dynamical system, to infer the hidden states, and then place RBF kernels on a randomly chosen subset of the inferred state means.⁶ We set the widths (variances) of the RBF kernels once we have the spacing of their centres by attempting to make neighbouring kernels cross when their outputs are half of their peak value. This ensures that with all the coefficients set approximately equal the RBF network will have an almost “flat” output across the space.⁷

These heuristics can be used both for fixed assignments of centres and widths, and as initialization to an adaptive RBF placement procedure. In section IV-A we discuss techniques for adapting both the positions of the RBF centres and their widths during training of the model.

For systems with nonlinear dynamics but approximately linear output functions we initialize using maximum likelihood factor analysis (FA) [9] trained on the collection of output observations (or conditional factor analysis for models with inputs). The loading matrix of the factor analysis solution is used to initialize the observation matrix C in the dynamical system. By doing time-independent inference through the factor analysis model we can obtain approximate estimates for the state at each time (independently). These estimates can be used to initialize the nonlinear RBF regressor by fitting the estimates at one time step as a function of those at the previous time step. (We also sometimes do a few iterations of training using a purely linear dynamical system before initializing the nonlinear RBF network.) Since such systems are nonlinear flows embedded in linear manifolds, this initialization estimates the embedding manifold using a linear statistical technique (FA) and the flow using a nonlinear regression based on projections into the estimated manifold.

If the output function is nonlinear but the dynamics are approximately linear then a mixture of factor analyzers (MFA) can be trained on the output observations [40], [41]. Such systems are linear flows in a nonlinear embedding manifold. The MFA initialization captures the nonlinear shape of the output manifold. Estimating the dynamics is difficult (since the hidden states of the individual analyzers in the mixture cannot be combined easily into a single

⁶In order to properly cover the portions of the state space which are most frequently used, we require a minimum distance between RBF kernel centres. Thus in practice we reject centres that fall too close together.

⁷One way to see this is to consider Gaussian RBFs in a n -dimensional grid, i.e. a square lattice, all with heights 1. The RBF centres define a hypercube, the distance between neighboring RBFs being $2d$ where d is chosen such that $e^{-d^2/(2\sigma^2)} = 1/2$. At the centres of the hypercubes, there are 2^n contributions from neighboring Gaussians, each of which is a distance $\sqrt{n}d$, and so contributes $(1/2)^n$ to the height. Therefore the height at the interiors is approximately equal to the height at the corners.

internal state representation) but still possible.⁸

Ideally Bayesian methods would be used to control the complexity of the model by estimating the internal state dimension and optimal number of RBF centres. However, in general only approximate techniques such as cross validation or variational approximations can be implemented in practice (see section IV-D). Currently we have set these complexity parameters either by hand or with cross validation.

III. RESULTS

We tested how well our algorithm could learn the dynamics of a nonlinear system by observing only the system inputs and outputs. We investigated the behaviour on simple one and two-dimensional state space problems whose nonlinear dynamics were known as well as on a weather time-series problem involving real temperature data.

A. One and two-dimensional nonlinear state space models

In order to be able to compare our algorithm’s learned internal state representation with a ground truth state representation, we first tested it on synthetic data generated by nonlinear dynamics whose form was known. The systems we considered consisted of three inputs and four observables at each time, with either one or two hidden state variables. The relation of the state from one time step to the next was given by a variety of nonlinear functions followed by Gaussian noise. The outputs were a linear function of the state and inputs plus Gaussian noise. The inputs affected the state only through a linear driving function. The true and learned state transition functions for these systems as well as sample outputs in response to Gaussian noise inputs and internal driving noise are shown in figures 5, 6 and 7 (right columns).

We initialized each nonlinear model with a linear dynamical model trained with EM, which in turn we initialized with a variant of factor analysis (see section II-D). The one-dimensional state space models were given 11 RBFs in x -space, which were uniformly spaced within a range which was automatically determined from the density of inferred points. Two-dimensional state space models were given 25 RBFs spaced in a 5×5 grid uniformly over the range of inferred states. After the initialization was over, the algorithm discovered the nonlinearities in the dynamics within less than 5 iterations of EM (see figures 5, 6 and 7, left and middle panels).

After training the models on input-output observations from the dynamics we examined the learned internal state representation and compared it to the known structure of the generating system. As the figures show, the algorithm

⁸As an approximate solution to the problem of getting a single hidden state from a MFA, we can use the following procedure. (1) Estimate the “similarity” between analyzer centres using average separation in time between data points for which they are active. (2) Use standard embedding techniques such as multidimensional scaling (MDS) [42] to place the MFA centres in a Euclidean space of dimension k . (3) Time-independent state inference for each observation now consists of the responsibility weighted low-dimensional MFA centres, where the responsibilities are the posterior probabilities of each analyzer given the observation under the MFA.

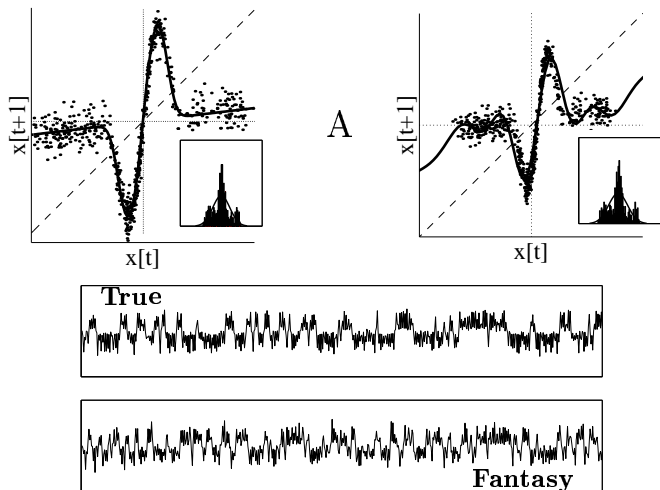


Fig. 5. Example of fitting a system with nonlinear dynamics and linear observation function. The panels show the fitting of a nonlinear system with a one-dimensional hidden state and 4 noisy outputs driven by Gaussian noise inputs and internal state noise. **(left)** The true dynamics function (line) and states (dots) used to generate the training data. Inset: histogram of internal states. **(right)** The learned dynamics function and states inferred on the training data. Inset: histogram of inferred internal states. **(middle)** The first component of the observable time series from the training data. **(bottom)** First component of fantasy data generated from the learned model (same scale as middle plot).

recovers the form of the nonlinear dynamics quite well. We are also able to generate “fantasy” data from the models once they have been learned by exciting them with Gaussian noise of similar variance to that applied during training. The resulting observation streams look qualitatively very similar to the time series from the true systems.

We can quantify this quality of fit by comparing the log-likelihood of the training sequences and novel test sequences under our nonlinear model to the likelihood under a basic linear dynamical system model or a static model such as factor analysis. Figure 8 presents this comparison. The nonlinear dynamical system had significantly superior likelihood on both training and test data for all the example systems. (Notice that for system **E**, the linear dynamical system is much better than factor analysis because of the strong hysteresis (mode-locking) in the system. Thus, the output at the previous time step is an excellent predictor of the current output.)

B. Weather data

As an example of a real system with a nonlinear output function as well as important dynamics, we trained our model on records of the daily maximum and minimum temperatures in Melbourne, Australia, over the period 1981–1990.⁹ We used a model with two internal state variables, three outputs and no inputs. During the training phase, the three outputs were the minimum and maximum daily temperature as well as a real valued output indicating the

⁹This data is available on the world wide web from the Australian Bureau of Meteorology.

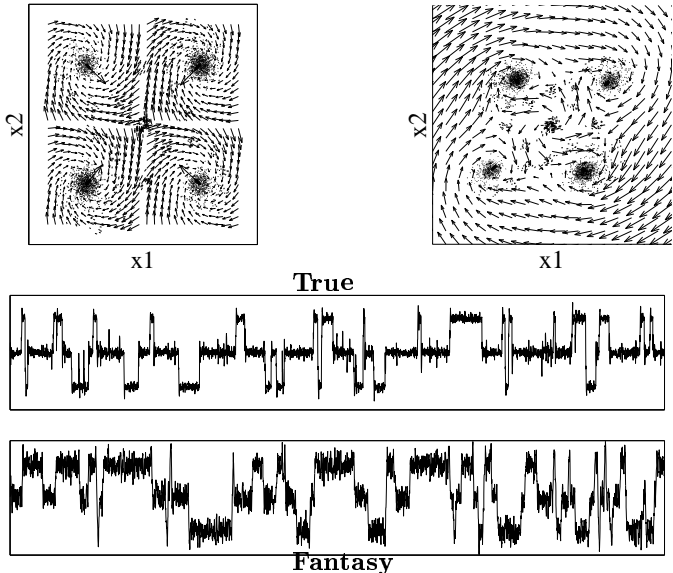


Fig. 7. Multidimensional example of fitting a system with nonlinear dynamics and linear observation functions. The true system is piecewise linear across the state space. The plots show the fitting of a nonlinear system with a two-dimensional hidden state and 4 noisy outputs driven by Gaussian noise inputs and internal state noise. **(top left)** The true dynamics vector field (arrows) and states (dots) used to generate the training data. **(top right)** The learned dynamics vector field and states inferred on the training data. **(middle/bottom)** The first component of the observable time series: training data above and fantasy data generated from the learned model below.

time of the year (month) in the range $[0,12]$. The model was trained on 1500 days of temperature records, or just over four seasons. We tested on the remaining 2150 days by showing the model only the minimum and maximum daily temperatures and attempting to predict the time of year (month). The prediction was performed by using the EKS algorithm to do state inference given only the two available observation streams. Once state inference was performed, the learned output function of the model could be used to predict the time of year. This prediction problem inherently requires the use of information from previous and/or future times since the static relationship between temperature and season is ambiguous during spring/fall. Figure 9 shows the results of this prediction after training; the algorithm has discovered a relationship between the hidden state and the observations which allows it to perform reasonable prediction for this task. Also show are the model predictions of minimum and maximum temperatures given the inferred state.

Although not explicitly part of the generative model, the learned system implicitly parameterizes a relationship between time of year and temperature. We can discover this relationship by evaluating the nonlinear output function at many points in the state space. Each evaluation yields a triple of month, minimum and maximum temperature. These triples can then be plotted against each other as in figure 10 to show that the model has discovered the southern hemisphere’s seasonal temperature variations.

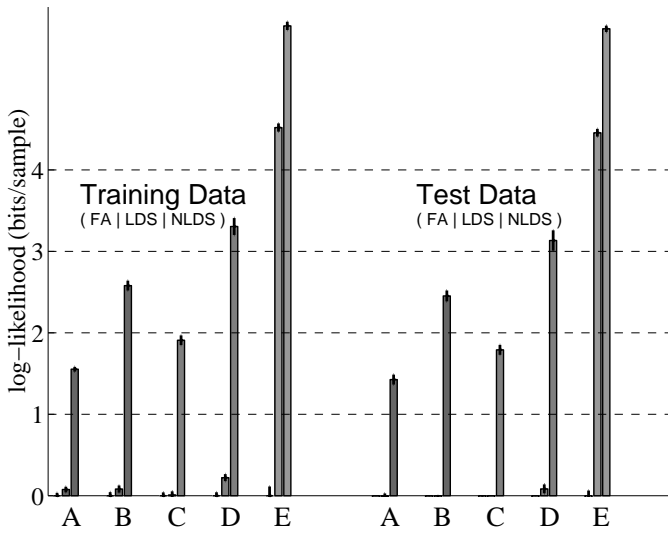


Fig. 8. Differences in log-likelihood assigned by various models to training and test data from the systems in figure 6. Each adjacent group of three bars shows the log-likelihood under factor analysis (FA), linear dynamical systems (LDS), and the nonlinear dynamical systems (NLDS). The letters on the ordinate match the labels in previous figures. Results on training data appear on the left and on test data on the right; taller bars represent better models. Log-likelihoods are offset so that FA on the training data is zero. Error bars represent the 68% quantile about the median across 100 repetitions of training or testing.

IV. EXTENSIONS

A. Learning the means and widths of the RBFs

It is possible to relax the assumption that the Gaussian radial basis functions have fixed centers and widths, although this results in a somewhat more complicated and slower fitting algorithm. To derive learning rules for the RBF centers c_i and width matrices S_i we need to consider how they play into the cost function (17) through the RBF kernel, (14). We take derivatives with respect to the expectation of the cost function \mathcal{C} and exchange the order of the expectation and the derivative:

$$\begin{aligned} \left\langle \frac{\partial \mathcal{C}}{\partial c_i} \right\rangle &= \left\langle \frac{\partial \mathcal{C}}{\partial \rho_i} \frac{\partial \rho_i}{\partial c_i} \right\rangle \\ &= 2 \langle (\theta \Phi - z)^\top Q^{-1} h_i \rho_i(x) S_i^{-1} (x - c_i) \rangle \end{aligned} \quad (19)$$

Recalling that $\Phi = [\rho_1(x) \rho_2(x) \dots \rho_I(x) x^\top u^\top 1]^\top$, it is clear that c_i figures nonlinearly in several places in this equation, and therefore it is not possible to solve for c_i in closed form. We can however, use the above gradient to move the center c_i to decrease the cost, which corresponds to taking a *partial* M-step w.r.t. c_i . Equation (19) requires computing three third order expectations in addition to the first and second order expectations needed to optimize θ and Q : $\langle \rho_i(x) \rho_k(x) x_\ell \rangle_j$, $\langle \rho_i(x) x_k x_\ell \rangle_j$, and $\langle \rho_i(x) z_k x_\ell \rangle_j$. Similarly, differentiating the cost with respect to S_i^{-1} gives:

$$\begin{aligned} \left\langle \frac{\partial \mathcal{C}}{\partial S_i^{-1}} \right\rangle &= \left\langle \frac{\partial \mathcal{C}}{\partial \rho_i} \frac{\partial \rho_i}{\partial S_i^{-1}} \right\rangle = \\ &\langle [(\theta \Phi - z)^\top Q^{-1} h_i] \rho_i(x) [S_i - (x - c_i)(x - c_i)^\top] \rangle. \end{aligned} \quad (20)$$

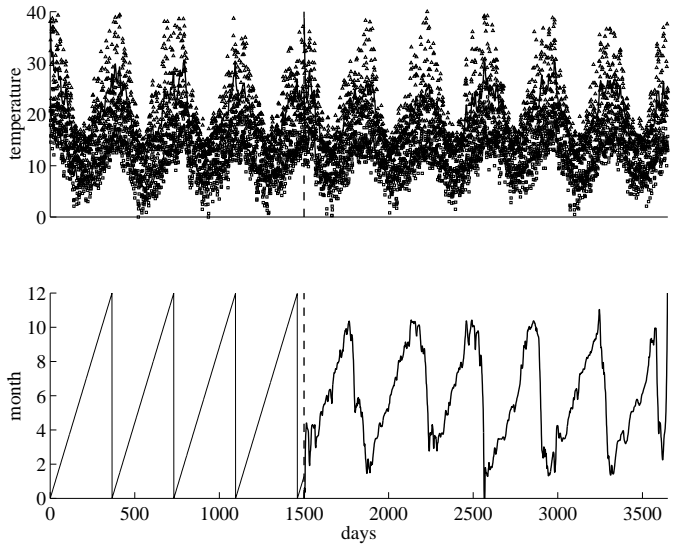


Fig. 9. Model of maximum and minimum daily temperatures in Melbourne, Australia from 1981–1990. **(left of vertical line)** A system with two hidden states governed by linear dynamics and a nonlinear output function was trained on observation vectors of a three-dimensional time-series consisting of maximum and minimum temperatures for each day as well as the (real-valued) month of the year. Training points are shown as triangles (max temp), squares (min temp) and a solid line (sawtooth wave below). **(right of vertical line)** After training, the system can infer its internal state from only the temperature observations. Having inferred its internal state it can predict the month of the year as a missing output (line below). The solid lines in the upper plots show the model's prediction of minimum and maximum temperature given the inferred state at the time.

We now need three *fourth* order expectations as well: $\langle \rho_i(x) \rho_k(x) x_\ell x_m \rangle_j$, $\langle \rho_i(x) x_k x_\ell x_m \rangle_j$, and $\langle \rho_i(x) z_k x_\ell x_m \rangle_j$.

These additional expectations increase both the storage and computation time of the algorithm, a cost which may not be compensated by the added advantage of moving of centers and widths by small gradient steps. One heuristic is to place centers and widths using unsupervised techniques like the EM algorithm for Gaussian mixtures, which considers solely the input density and not the output nonlinearity. Alternatively, some of these higher-order expectations can be approximated using, for example $\langle \rho_i(x) \rangle \approx \rho_i(\langle x \rangle)$.

B. Online learning

One of the major limitations of the algorithm we have presented in this paper is that it is a *batch* algorithm, i.e. it assumes that we use the entire sequence of observations to estimate the model parameters. Fortunately, it is relatively straightforward to derive an online version of the algorithm, which updates parameters as it receives observations. This is achieved using the recursive least squares (RLS) algorithm [11].

The key observation is that the cost minimized in the M-step of the algorithm (17) is a quadratic function of the parameters θ . RLS is simply a way of solving quadratic problems online. Using t to index time step, the resulting

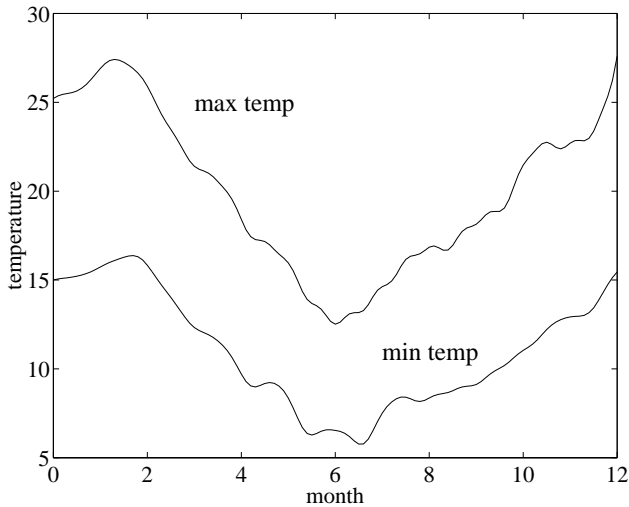


Fig. 10. Prediction of maximum and minimum daily temperatures based on time of year. The model from figure 9 implicitly learns a relationship between time of year and minimum/maximum temperature. This relationship is not directly invertible, but the temporal information used by Extended Kalman Smoothing correctly infers month given temperature as shown in figure 9.

algorithm for scalar z is as follows:

$$\theta_t = \theta_{t-1} + (\langle z \Phi \rangle_t - \theta_{t-1} \langle \Phi \Phi^\top \rangle_t) P_t \quad (21a)$$

$$P_t = P_{t-1} - \frac{P_{t-1} \langle \Phi \Phi^\top \rangle_t P_{t-1}}{1 + \langle \Phi^\top P_{t-1} \Phi \rangle_t} \quad (21b)$$

$$Q_t = Q_{t-1} + \frac{1}{t} [\langle z^2 \rangle_t - \theta_t \langle \Phi z \rangle_t - Q_{t-1}] \quad (21c)$$

Let us ignore the expectations for now. Initializing $\theta_0 = 0$, $Q_0 = 1$ and P_0 very large, it is easy to show that after a few iterations the estimates of θ_t will rapidly converge to the exact values obtained by the least square solution. The estimate of Q will converge to the correct values plus a bias incurred by the fact that the early estimates of Q were based on residuals from θ_t rather than $\lim_{t \rightarrow \infty} \theta_t$. P_t is a recursive estimate of $(\sum_{j=1}^t \langle \Phi \Phi \rangle_j)^{-1}$, obtained by using the matrix inversion lemma.

There is an important way in which this online algorithm is an approximation to the batch EM algorithm we have described for nonlinear state-space models. The expectations $\langle \cdot \rangle_t$ in the online algorithm are computed by running a single step of the extended Kalman filter using the previous parameters θ_{t-1} . In the batch EM algorithm, the expectations are computed by running an extended Kalman *smoother* over the entire sequence using the current parameter estimate. Moreover, these expectations are used to re-estimate the parameters, the smoother is then re-run, the parameters are re-re-estimated, and so on, to perform the usual iterations of EM. In general, we can expect that, unless the time series is non-stationary, the parameter estimates obtained by the batch algorithm after convergence will model the data better than those obtained by the online algorithm.

Interestingly, the updates for the RLS online algorithm described here are very similar to the parameter updates

used in a dual extended Kalman filter approach to system identification [43] (discussed in section V-E). This similarity is not coincidental, since the Kalman filter can be derived as a generalization of the RLS algorithm. In fact this similarity can be exploited in an elegant manner to derive an online algorithm for parameter estimation for non-stationary nonlinear dynamical systems.

C. Nonstationarity

To handle nonstationary time series we assume that the parameters can drift according to a Gaussian random walk with covariance Σ_θ :

$$\theta_t = \theta_{t-1} + \epsilon_t \quad \text{where } \epsilon_t \sim \mathcal{N}(0, \Sigma_\theta)$$

As before, we have the following function relating the z variables to the parameters θ and nonlinear kernels Φ :

$$z_t = \theta_t \Phi_t + w_t \quad \text{where } w_t \sim \mathcal{N}(0, Q)$$

which we can view as the observation model for a “state variable” θ_t with time varying “output matrix” Φ_t . Since both the dynamics and observation models are linear in θ and the noise is Gaussian we can apply the following Kalman filter to recursively compute the distribution of drifting parameters θ :

$$\hat{\theta}_t = \hat{\theta}_{t-1} + \frac{(\langle z \Phi \rangle_t - \hat{\theta}_{t-1} \langle \Phi \Phi^\top \rangle_t) P_{t|t-1}}{Q_{t-1} + \langle \Phi^\top P_{t|t-1} \Phi \rangle_t} \quad (22a)$$

$$P_{t|t-1} = P_{t-1} + \Sigma_\theta \quad (22b)$$

$$P_t = P_{t|t-1} - \frac{P_{t|t-1} \langle \Phi \Phi^\top \rangle_t P_{t|t-1}}{Q_{t-1} + \langle \Phi^\top P_{t|t-1} \Phi \rangle_t} \quad (22c)$$

$$Q_t = Q_{t-1} + \lambda [\langle z^2 \rangle_t - \hat{\theta}_t \langle \Phi z \rangle_t - Q_{t-1}] \quad (22d)$$

There are two important things to note. First, these equations describe an ordinary Kalman filter except that both the “output” z and “output matrix” Φ_t are jointly uncertain with a Gaussian distribution. Second, we have also assumed that the output noise covariance can drift by introducing a forgetting factor λ in its re-estimation equation. As before, the expectations are computed by running one step of the EKF over the hidden variables using θ_{t-1} .

While we derived this online algorithm starting from the batch EM algorithm, what we have ended up with appears almost identical to the dual extended Kalman filter. Indeed, we have two Kalman filters, one extended and one ordinary, running in parallel, estimating the hidden states and parameters, respectively.

We can also view this online algorithm as an approximation to the Bayesian posterior over parameters and hidden variables. The true posterior would be some complicated distribution over the x , z and θ parameters. Here we have recursively approximated it with two independent Gaussians, one over (x, z) and one over θ . The approximated posterior for θ_t has mean $\hat{\theta}_t$ and covariance P_t .

D. Using Bayesian methods for model selection and complexity control

Like any other maximum likelihood procedure, the EM algorithm described in this chapter has the potential to overfit the data set, that is, to find spurious patterns in noise in the data thereby generalizing poorly. In our implementation we used some ridge regression, i.e. a weight decay regularizer on the h_i parameters, which seemed to work well in practice but required some heuristics for setting regularization parameters. (Although, as mentioned previously, integrating over the hidden variables acts as a sort of modulated input noise and so in effect performs ridge regression which can eliminate the need for explicit regularization.)

A second closely related problem faced by maximum likelihood methods is that there is no built-in procedure for doing model selection. That is, the value of the maximum of the likelihood is not a suitable way to choose between different model structures. For example, consider the problems of choosing the dimensionality of the state space x and choosing the number of basis functions I . Higher dimensions of x and more basis functions should always, in principle, result in higher maxima of the likelihood, which means that more complex models will always be preferred to simpler ones. But this, of course, leads to overfitting.

Bayesian methods provide a very general framework for simultaneously handling the overfitting and model selection problems in a consistent manner. The key idea of the Bayesian approach is to avoid maximization wherever possible. Instead, possible models, structures, parameters—in short, all settings of unknown quantities—should be weighted by their posterior probabilities and predictions should be made according to this weighted posterior.

For our nonlinear dynamical system, we can, for example, treat the parameters θ as an unknown. Then the model's prediction of the output at time $t + 1$ is:

$$\begin{aligned} & p(y_{t+1}|u_{1:t+1}, y_{1:t}) \\ &= \int d\theta p(y_{t+1}|u_{t+1}, y_{1:t}, u_{1:t}, \theta) p(\theta|y_{1:t}, u_{1:t}) \\ &= \int d\theta p(\theta|y_{1:t}, u_{1:t}) \int dx_{t+1} p(y_{t+1}|u_{t+1}, x_{t+1}, \theta) p(x_{t+1}|u_{1:t+1}, y_{1:t}, \theta) \end{aligned} \quad (23)$$

where the first integral on the last line is over the posterior distribution of the parameters and the second integral is over the posterior distribution of the hidden variables.

The posterior distribution over parameters can be obtained recursively from Bayes rule:

$$p(\theta|y_{1:t}, u_{1:t}) = \frac{p(y_t|u_{1:t}, y_{1:t-1}, \theta) p(\theta|u_{1:t-1}, y_{1:t-1})}{p(y_t|u_{1:t}, y_{1:t-1})}$$

The dual extended Kalman filter, the joint extended Kalman filter, and the non-stationary online algorithm from section IV-C are all coarse approximations of these Bayesian recursions.

The above equations are all implicitly conditioned on some choice of model structure \mathcal{S}_m , i.e. dimension of x and number of basis functions. Although the Bayesian modeling philosophy advocates averaging predictions of different model structures, if necessary it is also possible to use Bayes rule to *choose* between model structures according to their probabilities:

$$P(\mathcal{S}_m|y_{1:t}, u_{1:t}) = \frac{p(y_{1:t}|u_{1:t}, \mathcal{S}_m)P(\mathcal{S}_m)}{\sum_n p(y_{1:t}|u_{1:t}, \mathcal{S}_n)P(\mathcal{S}_n)}$$

Tractable approximations to the required integrals can be obtained in several ways. We highlight three ideas, without going into much detail; an adequate solution to this problem for nonlinear dynamical systems requires further research. The first idea is the use of Markov chain Monte Carlo techniques to sample over both parameters and hidden variables. MCMC methods such as Gibbs sampling have been used for linear dynamical systems [44], [45] while a promising method for nonlinear systems is particle filtering [20], [26]. The second idea is the use of so-called “automatic relevance determination” (ARD; [46], [47]). This consists of using a zero-mean Gaussian prior on each parameter with tunable variances. Optimizing these variance hyperparameters results in “irrelevant” being eliminated from the model. ARD for RBF networks with a center on each data point has been used by Tipping [48] successfully for nonlinear regression, and given the name “relevance vector machine” in analogy to support vector machines. The third idea is the use of variational methods to lower bound the model structure posterior probabilities. Variational Bayesian methods have been used by [49] to infer the structure of linear dynamical systems, although the generalization to nonlinear systems of the kind described in this chapter is not straightforward.

Of course, in principle the Bayesian approach would advocate averaging over all possible choices of c_i , S_i , I , Q , etc. It is easy to see how this can rapidly get very unwieldy.

V. DISCUSSION

A. Identifiability and Expressive Power

As we saw in the experiments above, the algorithm we have presented is capable of learning good density models for a variety of nonlinear time series. Specifying the class of nonlinear systems which our algorithm can model well defines its *expressive power*. A related question is, what is the ability of this model, in principle, to recover the actual parameters of specific nonlinear systems? This is the question of *model identifiability*. These two questions are intimately tied as they both describe the mapping between actual nonlinear systems and model parameter settings.

There are three trivial degeneracies which make our model technically unidentifiable, but should not concern us. First, it is always possible to permute the dimensions in the state space and, by permuting the domain of the output mapping and dynamics in the corresponding fashion, obtain an exactly equivalent model. Second, the state variables can be rescaled or in fact transformed by any invertible linear mapping. This transformation can be absorbed

by the output and dynamics functions yielding a model with identical input-output behaviour. Without loss of generality, we always set the covariance of the state evolution noise to be the identity matrix which both sets the scale of the state space, and disallows certain state transformations without reducing the expressive power of the model. Third, we take the observation noise to be uncorrelated with the state noise and both noises to be zero mean since again, without loss of generality these can be absorbed into the f and g functions¹⁰

There exist other forms of unidentifiability which are more difficult to overcome. For example, if both f and g are nonlinear then at least in the noise free case, for any arbitrary invertible transformation of the state, there exist transformations of f and g which result in identical input-output behavior. In this case, it would be very hard to detect that the recovered model is indeed a good model of the actual system since the estimated and actual states would appear to be unrelated.

Clearly, not all systems can be modeled by assuming that f is linear and g is nonlinear. Similarly, not all systems can be modeled by assuming that f is nonlinear and g is linear. For example, consider the case where the observations y_t and y_{t+n} are statistically independent, but each observation lies on a curved low-dimensional manifold in a high dimensional space. Modeling this would require a nonlinear g as in nonlinear factor analysis, but an $f = 0$. Therefore, choosing either f or g to be linear restricts the expressive power of the model.

Unlike for the state noise covariance Q , assuming that the observation noise covariance R is diagonal *does* restrict the expressive power of the model. This is easy to see for the case where the dimension of the state space is small and the dimension of the observation vector is large. A full covariance R can capture all correlations between observations at a single time step, while a diagonal R model cannot.

For nonlinear dynamical systems, the Gaussian noise assumption is not as restrictive as it may initially appear. This is because the nonlinearity can be used to turn Gaussian noise into non-Gaussian noise [9].

Of course, we have restricted our expressive power by using an RBF network, especially once the means and centers of the RBFs are fixed. One could try to appeal to universal approximation theorems to make the claim that we could in principle model any nonlinear dynamical system. But this would be misleading in the light of the noise assumptions and the fact that only a finite and usually small number of RBFs are going to be used in practice.

B. Embedded Flows

There are two ways to think about the dynamical models we have investigated, shown in figure 11. One is as a non-

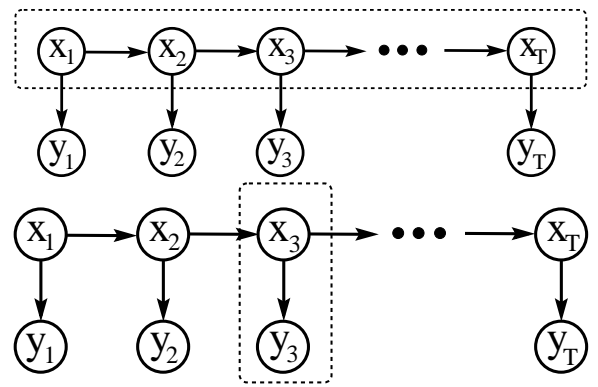


Fig. 11. Two interpretations of the graphical model for stochastic (non)linear dynamical systems (see text). (top) A Markov process embedded in a manifold. (bottom) Nonlinear factor analysis through time.

linear Markov process (flow) x_t which has been embedded (or potentially projected) into a manifold y_t . From this perspective, the function f controls the evolution of the stochastic process and the function g specifies the nonlinear embedding (or projection) operation.¹¹

Another way to think of the same model is as a nonlinear version of a latent variable model such as factor analysis (but possibly with external inputs as well) in which the latent variables or factors evolves through time rather than being drawn independently for each observation. The nonlinear factor analysis model is represented by g and the time evolution of the latent variables by f .

If the state space is of lower dimension than the observation space and the observation noise is additive, then a useful geometrical intuition applies. In such cases, we have observed a flow inside an embedded manifold. The observation function g specifies the structure (shape) of the manifold, while the dynamics f specifies the flow within the manifold. Armed with this intuition, the learning problem looks like it might be decoupled into two separate stages: first find the manifold by doing some sort of density modeling on the collection of observed outputs (ignoring their time order); second, find the flow (dynamics) by projecting the observations into the manifold and doing nonlinear regression from one time step to the next. This intuition is partly true, and indeed provides the basis for many of the practical and effective initialization schemes we have tried. However, the crucial point as far as the design of learning algorithms is concerned, is that the two learning problems interact in a way that makes the problem *easier*. Once we know something about the dynamics, this information gives some prior knowledge when trying to learn the manifold shape. For example, if the dynamics suggest that the next state will be near a certain point we can use this information to do better than naive projection when we locate a noisy observation on the manifold. Conversely, knowing something about the manifold allows us to estimate the

¹⁰Imagine that the joint noise covariance was nonzero: $\langle w_t v_t^T \rangle = S$. Replacing A with $A' = A - SR^{-1}C$ gives a new noise process w' with covariance $Q' = Q - SR^{-1}S^T$ that is uncorrelated with v , leaving the input output behavior invariant. Similarly any non-zero noise means can be absorbed into the b terms in the functions f and g .

¹¹To simplify presentation we'll neglect driving inputs u_t in this section, although the arguments extend as well to systems with inputs.

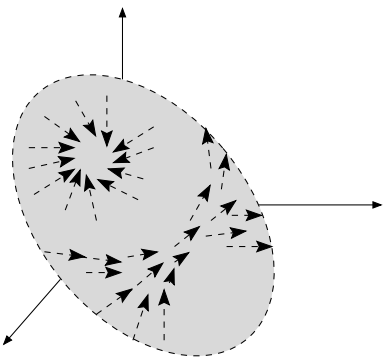


Fig. 12. Linear and nonlinear dynamical systems represents flow fields embedded in manifolds. For systems with linear output functions, such as the one illustrated, the manifold is a hyper-plane while the dynamics may be complex. For systems with nonlinear output functions the shape of the embedding manifold is also curved.

dynamics more effectively.

We discuss separately two special cases of flows in manifolds: systems with linear output functions but nonlinear dynamics and systems with linear dynamics but nonlinear output function.

When the output function g is linear and the dynamics f is nonlinear (figure 12) the observed sequence forms a nonlinear flow in a linear subspace of the observation space. The manifold estimation is made easier, even with high levels of observation noise, by the fact that its shape is known to be a hyperplane. All that is required is to find its orientation and the character of the output noise. Time-invariant analysis of the observations by algorithms such as factor analysis is an excellent way to initialize estimates of the hyperplane and noises. However during learning we made have cause to tilt the hyperplane to make the dynamics fit better, or conversely to modify the dynamics to make the hyperplane model better.

This setting is actually more expressive than it might seem initially. Consider a nonlinear output function $g(x)$ that is “invertible” in the sense that it be written in the form $g(x) = C\tilde{g}(x)$ for invertible \tilde{g} and non-square matrix C . Any such nonlinear output function can be made strictly linear if we transform to a new state variable \tilde{x} :

$$\begin{aligned} \tilde{x} &= \tilde{g}(x) \quad \Rightarrow \\ \tilde{x}_{t+1} &= \tilde{f}(\tilde{x}_t, w_t) = \tilde{g}(f(\tilde{g}^{-1}(\tilde{x}))) + w_t \\ y_t &= C\tilde{x}_t + v_t = g(x_t) + v_t \end{aligned}$$

which gives an equivalent model but with a purely linear output process, and potentially non-additive dynamics noise.

For nonlinear output functions g paired with linear dynamics f , the observation sequence forms a matrix (linear) flow in a nonlinear manifold:

$$x_{t+1} = Ax_t + w_t \quad (24a)$$

$$y_t = g(x_t) + v_t \quad (24b)$$

The manifold learning is harder now, because we must estimate a thin, curved subspace of the observation space in

the presence of noise. However, once we have learned this manifold approximately, we project the observations into it and learn only linear dynamics. The win comes from the following fact: in the locations where the projected dynamics do not look linear, we know that we should bend the manifold to make the dynamics more linear. Thus not only the shape of the outputs (ignoring time) but also the linearity of the dynamics give us clues to learning the manifold.

C. Stability

Stability is a key issue in the study of any dynamical system. Here we have to consider stability at two levels: the stability of the learning procedure, and the stability of the learned nonlinear dynamical system.

Since every step of the EM algorithm is guaranteed to increase the log likelihood until convergence, it has a built-in Lyapunov function for stable learning. However, as we pointed out, our use of extended Kalman smoothing in the E-step of the algorithm represents an approximation to the exact E-step, and therefore we have to forego any guarantees of stability of learning. While we rarely had problems with stability of learning, this is sure to be problem-specific, depending both of the quality of the EKS approximation and on how close the true system dynamics are to the boundary of stability. In contrast to the EKS approximations, certain variational approximations [31] transform the intractable Lyapunov function into a tractable one, and therefore preserve stability of learning. It is not clear how to apply these variational approximations to nonlinear dynamics, although this would clearly be an interesting area of research.

Stability of the learned nonlinear dynamical system can be analyzed by making use of some linear systems theory. We know that for discrete time linear dynamical systems if all eigenvalues of the A matrix lie inside the unit circle, the system is globally stable. The nonlinear dynamics of our RBF network f can be decomposed into two parts (cf eq (13)): a linear component given by A , and a nonlinear component given by $\sum_i h_i \rho_i(x)$. Clearly for the system to be globally stable, A has to satisfy the eigenvalue criterion for linear systems. Moreover, if the RBF coefficients for both f and g are bounded, i.e. $\max_i |h_i| < \bar{h}$, and the RBF support is bounded in the following way: $\min_i \det(S_i) > 0$ and $\max_{ij} |c_i - c_j| < \bar{c}$, then the nonlinear system is stable in a bounded-input bounded-output sense. That is, for any sequence of bounded inputs the output sequence of the noise free system will be bounded with probability 1. Intuitively, although unstable behavior might occur in the region of RBF support, one x leaves this region it is drawn back in by A .

For the online EM learning algorithm, the hidden state dynamics and the parameter re-estimation dynamics will interact, and therefore a stability analysis would be quite challenging. However, since there is no stability guarantee for the batch EKS-EM algorithm, it seems very unlikely that a simple form of the online algorithm could be provably stable.

It has long been known that for linear systems, there is an equivalence between so called state-space formulations which involve hidden variables and direct vector autoregressive models of the time-series. In 1980, Takens proved a remarkable theorem [50] which tells us that for almost any deterministic nonlinear dynamical system with a d -dimensional state space, the state can be effectively reconstructed by observing $2d + 1$ time lags of any one of its outputs. In particular, Takens showed that such a lag vector will be a smooth embedding (diffeomorphism) of the true state if one exists. This notion of finding an ‘‘embedding’’ for the state has been used to justify a nonlinear regression approach to learning nonlinear dynamical systems. That is, if you suspect that the system is nonlinear and that it has d state dimensions, instead of building a state-space model, you can do away with representing states and just build an autoregressive (AR) model directly on the observations that nonlinearly relates previous outputs and the current output. This view begs the question, do we need our models to have hidden states at all?

While no constructive realization for Takens' theorem exists in general, for linear systems there are very strong results. For purely linear systems, we can appeal to the Cayley-Hamilton theorem¹² to show that the hidden state can always be eliminated to obtain a equivalent vector autoregressive model by taking only d time lags of the output. Furthermore there is a construction which allows this conversion to be performed explicitly.¹³ Takens' theorem offers us a similar guarantee for elimination of hidden states in nonlinear dynamical systems as long as we take $2d + 1$ output lags. However, no similar recipe exists for explicitly converting to an autoregressive form. These results make hidden states seem unnecessary.

The problem with this view is that it does not generalize well to many realistic high-dimensional and noisy scenarios. Consider the example mentioned in the introduction. While it is mathematically true that the pixels in the video frame of a balloon floating in the wind are a (highly nonlinear) function of the pixels in the previous video frames, it would be ludicrous from the modeling perspective to build an AR model of the video images. This would require a number of parameters on the order of the number of pixels squared. Furthermore, unlike the noise-free case of Takens' theorem, when the dynamics are noisy the optimal prediction of the observation would have to depend on the entire history of past observations. Any truncation of this his-

¹²Any square matrix A of size n satisfies its own characteristic equation. Equivalently, any matrix power A^m for $m \geq n$ can be written as a linear combination of lower matrix powers $I, A, A^2, \dots, A^{n-1}$.

¹³Start with the system $x_{t+1} = Ax_t + w_t$, $y_t = Cx_t + v_t$. Create a d -dimensional lag vector $z_t = [y_t; y_{t+1}; \dots; y_{t+d-1}]$ which holds the current and $d - 1$ future outputs. Write $z_t = Gx_t + n_t$ for $G = [CI; CA; CA^2; \dots; CA^{d-1}]$ and Gaussian noise n (although with non-diagonal covariance). The Cayley-Hamilton theorem assures us that G is full rank and thus we need not take any more lags. Given the lag vector z_t we can solve the system $z_t = Gx_t$ for x_t ; write this solution as G^+z_t . Using the original observation equation d times, to solve for y_t, \dots, y_{t+d-1} in terms of z_t we can write an autoregression for z_t as $z_{t+1} = G^+AGz_t + m_t$ for Gaussian noise m .

tory throws away potentially valuable information about the unobserved state. The state-space formulation of nonlinear dynamical systems allows us to overcome both of these limitations of nonlinear autoregressive models. That is, it allows us to have compact representations of dynamics, and to integrate uncertain information over time. The price paid for this is that it requires having to do inference over the hidden state.

E. Should parameters and hidden states be treated differently?

The maximum likelihood framework on which the EM algorithm is based makes a distinction between parameters and hidden variables: it attempts to *integrate* over hidden variables to *maximize* the likelihood as a function of parameters. This leads to the two-step approach which computes sufficient statistics over the hidden variables in the E-step and optimizes parameters in the M-step. In contrast, a fully Bayesian approach to learning nonlinear dynamical state-space models would treat both hidden variables and parameters as unknown and attempt to compute or approximate the joint posterior distribution over them, in effect integrating over both.

It is important to compare these approaches to system identification with more traditional ones. We highlight two such approaches: *joint EKF* approaches, and *dual EKF* approaches.

In joint EKF approaches [10], [11], an augmented hidden state space is constructed which comprises the original hidden state space and the parameters. Since parameters and hidden states interact, even for linear dynamical systems this approach results in nonlinear dynamics over the augmented hidden states. Initializing a Gaussian prior distribution both over parameters and states, an extended Kalman filter is then used to recursively update the joint distribution over states and parameters based on the observations, $p(X, \theta|Y)$. This approach has the advantage that it can model uncertainties in the parameters and correlations between parameters and hidden variables. In fact this approach treats parameters and state variables completely symmetrically and can be thought of as iteratively implementing a Gaussian approximation to the recursive Bayes rule computations. Nonstationarity can be easily built in by giving the parameters e.g. random-walk dynamics. Although it has some very appealing properties, this approach is known to suffer from instability problems, which is the reason why dual EKF approaches have been proposed.

In dual EKF approaches [43], two interacting but distinct extended Kalman filters operate simultaneously. One computes a Gaussian approximation of the state posterior given a parameter estimate and the observations: $p(X|\hat{\theta}_{old}, Y)$, while the other computes a Gaussian approximation of the parameter posterior given the estimated states $p(\theta|\hat{X}_{old}, Y)$. The two EKFs interact by each feeding its estimate (i.e. the posterior means \hat{X} and $\hat{\theta}$) into the other. One can think of the dual EKF as performing approximate coordinate ascent in $p(X, \theta|Y)$ by itera-

tively maximizing $p(X|\hat{\theta}_{old}, Y)$ and $p(\theta|\hat{X}_{old}, Y)$ under the assumption that each conditional is Gaussian. Since the only interaction between parameters and hidden variables occurs through their respective means, the procedure has the flavor of mean-field methods in physics and neural networks [51]. Like these methods it is also likely to suffer from the overconfidence problem—namely, since the parameter estimate does not take into account the uncertainty in the states, the parameter covariance will be overly narrow, likewise for the states.

For large systems, both joint and dual EKF methods suffer from the fact that the parameter covariance matrix is quadratic in the number of parameters. This problem is more pronounced for the joint EKF since it considers the concatenated state space. Furthermore, both joint and dual EKF methods rely on Gaussian approximations to parameter distributions. This can sometimes be problematic, e.g. consider retaining positive definiteness of a noise covariance matrix under the assumption that its parameters are Gaussian distributed.

VI. CONCLUSION

This chapter brings together two classic algorithms, one from statistics and another from systems engineering, to address the learning of stochastic nonlinear dynamical systems. We have shown that by pairing the Extended Kalman Smoothing algorithm for approximate state estimation in the E-step, with a radial basis function learning model that permits exact analytic solution of the M-step, the EM algorithm is capable of learning a nonlinear dynamical model from data. As a side effect we have derived an algorithm for training a radial basis function network to fit data in the form of a mixture of Gaussians. We have also derived an online version of the algorithm and a version for dealing with non-stationary time series.

We have demonstrated the algorithm on a series of synthetic and realistic nonlinear dynamical systems and shown that it is able to learn accurate models from only observations of inputs and outputs. Initialization of model parameters and placement of the radial basis kernels are important to the practical success of the algorithm. We have discussed techniques for making these choices as well as provided gradient rules for adapting the centres and widths of the basis functions.

The belief network literature has recently been dominated by two methods for approximate inference, Markov chain Monte Carlo [52] and variational approximations [31]. To our knowledge [35] and [43] were the first instances where extended Kalman smoothing had been used to perform approximate inference in the E-step of EM. While EKS does not have the theoretical guarantees of variational methods (which are also approximate but monotonically optimize a computable objective function during learning), its simplicity has gained it wide acceptance in the estimation and control literatures as a method for doing inference in nonlinear dynamical systems. Our practical success in modeling a variety of nonlinear time series suggests that the combination of Extended Kalman algorithms and the

EM algorithm can provide powerful tools for learning nonlinear dynamical systems.

ACKNOWLEDGEMENTS

The authors acknowledge support from the Gatsby Charitable Fund, and an NSF PDF Fellowship to SR.

APPENDIX

I. EXPECTATIONS REQUIRED TO FIT THE RBFs

The expectations we need to compute for equation (18) are $\langle x \rangle_j$, $\langle z \rangle_j$, $\langle xx^\top \rangle_j$, $\langle zz^\top \rangle_j$, $\langle xz^\top \rangle_j$, and $\langle \rho_i(x) \rangle_j$, $\langle x \rho_i(x) \rangle_j$, $\langle z \rho_i(x) \rangle_j$, $\langle \rho_i(x) \rho_\ell(x) \rangle_j$.

Starting with some of the easier ones that do not depend on the RBF kernel ρ :

$$\begin{aligned} \langle x \rangle_j &= \mu_j^x & \langle z \rangle_j &= \mu_j^z \\ \langle xx^\top \rangle_j &= \mu_j^x \mu_j^{x,T} + C_j^{xx} & \langle zz^\top \rangle_j &= \mu_j^z \mu_j^{z,T} + C_j^{zz} \\ \langle xz^\top \rangle_j &= \mu_j^x \mu_j^{z,T} + C_j^{xz} \end{aligned}$$

Observe that when we multiply the Gaussian RBF kernel $\rho_i(x)$ (equation 14) and \mathcal{N}_j we get a Gaussian density over (x, z) with mean and covariance

$$\mu_{ij} = C_{ij} \left(C_j^{-1} \mu_j + \begin{bmatrix} S_i^{-1} c_i \\ 0 \end{bmatrix} \right) \quad (25)$$

and

$$C_{ij} = \left(C_j^{-1} + \begin{bmatrix} S_i^{-1} & 0 \\ 0 & 0 \end{bmatrix} \right)^{-1}, \quad (26)$$

and an extra constant (due to lack of normalization),

$$\beta_{ij} = (2\pi)^{-d_x/2} |S_i|^{-1/2} |C_j|^{-1/2} |C_{ij}|^{1/2} \exp\{-\delta_{ij}/2\}$$

where $\delta_{ij} = c_i^\top S_i^{-1} c_i + \mu_j^\top C_j^{-1} \mu_j - \mu_{ij}^\top C_{ij}^{-1} \mu_{ij}$. Using β_{ij} and μ_{ij} , we can evaluate the other expectations:

$$\langle \rho_i(x) \rangle_j = \beta_{ij}, \quad \langle x \rho_i(x) \rangle_j = \beta_{ij} \mu_{ij}^x, \quad \text{and} \quad \langle z \rho_i(x) \rangle_j = \beta_{ij} \mu_{ij}^z.$$

Finally,

$$\langle \rho_i(x) \rho_\ell(x) \rangle_j = (2\pi)^{-d_x} |C_j|^{-1/2} |S_i S_\ell|^{-1/2} |C_{i\ell j}|^{1/2} \exp\{-\gamma_{i\ell j}/2\}, \quad (27)$$

where

$$C_{i\ell j} = \left(C_j^{-1} + \begin{bmatrix} S_i^{-1} + S_\ell^{-1} & 0 \\ 0 & 0 \end{bmatrix} \right)^{-1} \quad (28)$$

and

$$\mu_{i\ell j} = C_{i\ell j} \left(C_j^{-1} \mu_j + \begin{bmatrix} S_i^{-1} c_i + S_\ell^{-1} c_\ell \\ 0 \end{bmatrix} \right), \quad (29)$$

and $\gamma_{i\ell j} = c_i^\top S_i^{-1} c_i + c_\ell^\top S_\ell^{-1} c_\ell + \mu_j^\top C_j^{-1} \mu_j - \mu_{i\ell j}^\top C_{i\ell j}^{-1} \mu_{i\ell j}$.

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, pp. 1–38, 1977.
- [2] C. F. Chen, "The EM approach to the multiple indicators and multiple causes model via the estimation of the latent variables," *Journal of the American Statistical Association*, vol. 76, no. 375, pp. 704–708, 1981.

- [3] 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, pp. 253–264, 1982.
- [4] Zoubin Ghahramani and Geoffrey Hinton, "Parameter estimation for linear dynamical systems," Tech. Rep. CRG-TR-96-2, Dept. of Computer Science, University of Toronto, February 1996.
- [5] R. E. Kalman and R. S. Bucy, "New results in linear filtering and prediction," *Trans. American Society of Mechanical Engineers, Series D, Journal of Basic Engineering*, vol. 83D, pp. 95–108, 1961.
- [6] L. E. Baum and J. A. Eagon, "An inequality with applications to statistical estimation for probabilistic functions of markov processes and to a model for ecology," *Bulletin of American Mathematical Society*, vol. 73, pp. 360–363, 1967.
- [7] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann, San Mateo, CA, 1988.
- [8] S. L. Lauritzen and D. J. Spiegelhalter, "Local computations with probabilities on graphical structures and their application to expert systems," *Journal of the Royal Statistical Society B*, vol. 50, no. 2, pp. 157–224, 1988.
- [9] Sam Roweis and Zoubin Ghahramani, "A unifying review of linear Gaussian models," *Neural Computation*, vol. 11, no. 2, 1999.
- [10] L. Ljung and T. Söderström, *Theory and Practice of Recursive Identification*, MIT Press, Cambridge, MA, 1983.
- [11] G.C. Goodwin and K.S. Sin, *Adaptive filtering prediction and control*, Prentice-Hall, 1984.
- [12] J. Moody and C. Darken, "Fast learning in networks of locally-tuned processing units," *Neural Computation*, vol. 1, no. 2, pp. 281–294, 1989.
- [13] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Trans. American Society of Mechanical Engineers, Series D, Journal of Basic Engineering*, vol. 82D, pp. 35–45, March 1960.
- [14] H. E. Rauch, "Solutions to the linear smoothing problem," *IEEE Transactions on Automatic Control*, vol. 8, pp. 371–372, 1963.
- [15] R. E. Kopp and R. J. Orford, "Linear regression applied to system identification and adaptive control systems," *Journal of the American Institute of Aeronautics and Astronautics*, vol. 1, no. 10, pp. 2300–2306, 1963.
- [16] H. Cox, "On the estimation of state variables and parameters for noisy dynamic systems," *IEEE Transactions on Automatic Control*, vol. 9, pp. 5–12, 1964.
- [17] S. Singhal and L. Wu, "Training multilayer perceptrons with the extended Kalman algorithm," in *Advances in Neural Information Processing Systems*, vol. 1, pp. 133–140. Morgan Kaufmann, 1989.
- [18] V. Kadirkamanathan and M. Niranjan, "A function estimation approach to sequential learning with neural networks," *Neural Computation*, vol. 5, pp. 954–975, 1993.
- [19] I. T. Nabney, A. McLachlan, and D. Lowe, "Practical methods of tracking of nonstationary time series applied to real-world data," in *AeroSense '96: Applications and Science of Artificial Neural Networks II*, S. K. Rogers and D. W. Ruck, Eds. 1996, pp. 152–163, SPIE Proceedings.
- [20] J.E. Handschin and D. Q. Mayne, "Monte Carlo techniques to estimate the conditional expectation in multi-stage non-linear filtering," *International Journal of Control*, vol. 9, no. 5, pp. 547–559, 1969.
- [21] M. Isard and A. Blake, "CONDENSATION – conditional density propagation for visual tracking," *International Journal of Computer Vision*, vol. 29, no. 1, pp. 5–28, 1998.
- [22] Genshiro Kitagawa, "Monte Carlo filter and smoother for non-Gaussian nonlinear state space models," *Journal of Computer Graphics and Graphical Statistics*, vol. 5, no. 1, pp. 1–25, Mar. 1996.
- [23] N. J. Gordon, D. J. Salmond, and A. F. M. Smith, "A novel approach to nonlinear/non-Gaussian Bayesian state space estimation," *IEE Proceedings F: Radar and Signal Processing*, vol. 140, pp. 107–113, 1993.
- [24] Mike West, "Approximating posterior distributions by mixtures," *Journal of the Royal Statistical Society (Ser. B)*, vol. 54, pp. 553–568, 1993.
- [25] Mike West, "Mixture models, monte carlo, bayesian updating and dynamic models," *Computing Science and Statistics*, vol. 24, pp. 325–333, 1993.
- [26] A. Doucet, J. F. G. de Freitas, and N.J. Gordon, *Sequential Monte Carlo Methods in Practice*, Springer-Verlag, New York, 2000.
- [27] S. J. Julier, J. K. Uhlmann, and H. F. Durrant-Whyte, "A new approach for filtering nonlinear systems," in *Proceedings of the 1995 American Control Conference*, Seattle, 1995, pp. 1628–1632.
- [28] S. J. Julier and J. K. Uhlmann, "A new extension of the Kalman filter to nonlinear systems," in *Proceedings of AeroSense: The 11th International Symposium on Aerospace/Defense Sensing, Simulation and Controls*, Orlando, 1997.
- [29] Eric A. Wan, R. van der Merwe, and Alex T. Nelson, "Dual estimation and the unscented transformation," in *Advances in Neural Information Processing Systems*. 1999, vol. 12, MIT Press.
- [30] Zoubin Ghahramani and Geoffrey E. Hinton, "Variational learning for switching state-space models," *Neural Computation*, vol. 12, no. 4, pp. 831–864, 2000.
- [31] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, "An Introduction to variational methods in graphical models," *Machine Learning*, 1999.
- [32] L. E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains," *Annals of Mathematical Statistics*, vol. 37, pp. 1554–1563, 1966.
- [33] R. M. Neal and G. E. Hinton, "A view of the EM algorithm that justifies incremental, sparse and other variants," in *Learning in Graphical Models*, M. I. Jordan, Ed. 1998, pp. 355–368, Kluwer Academic Press.
- [34] I. Csizsár and G. Tusnády, "Information geometry and alternating minimization procedures," *Statistics & Decisions, Supplement Issue 1*, pp. 205–237, 1984.
- [35] Zoubin Ghahramani and Sam Roweis, "Learning nonlinear dynamical systems using an em algorithm," in *Advances in Neural Information Processing Systems*. 1999, vol. 11, pp. 431–437, MIT Press.
- [36] J. F. G. de Freitas, M. Niranjan, and A. H. Gee, "Nonlinear state space estimation with neural networks and the EM algorithm," Tech. Rep., Cambridge University Engineering Dept., 1999.
- [37] T. Briegel and V. Tresp, "Fisher scoring and a mixture of modes approach for approximate inference and learning in nonlinear state space models," in *Advances in Neural Information Processing Systems*, vol. 11. MIT Press, 1999.
- [38] Zoubin Ghahramani and Geoffrey Hinton, "Switching state-space models," Tech. Rep. CRG-TR-96-3, Dept. of Computer Science, University of Toronto, July 1996.
- [39] Kevin Murphy, "Switching Kalman filters," Tech. Rep., Dept. of Computer Science, University of California Berkeley, August 1998.
- [40] G. E. Hinton, P. Dayan, and M. Revow, "Modeling the manifolds of Images of handwritten digits," *IEEE Transactions on Neural Networks*, vol. 8, pp. 65–74, 1997.
- [41] Zoubin Ghahramani and Geoffrey Hinton, "The EM algorithm for mixtures of factor analyzers," Tech. Rep. CRG-TR-96-1, Dept. of Computer Science, University of Toronto, May 1996 (revised Feb. 1997).
- [42] W. S. Torgerson, "Multidimensional scaling I. Theory and method," *Psychometrika*, vol. 17, pp. 401–419, 1952.
- [43] Eric A. Wan and Alex T. Nelson, "Dual Kalman filtering methods for nonlinear prediction," in *Advances in Neural Information Processing Systems*. 1997, vol. 9, MIT Press.
- [44] C. K. Carter and R. Kohn, "On Gibbs sampling for state space models," *Biometrika*, vol. 81, pp. 541–553, 1994.
- [45] S. Frühwirth-Schnatter, "Bayesian model discrimination and Bayes factors for linear Gaussian state space models," *Journal of the Royal Statistical Society, Series B*, vol. 57, pp. 237–246, 1995.
- [46] D. J. C. MacKay, "Bayesian non-linear modelling for the prediction competition," in *ASHRAE Transactions*, Atlanta Georgia, 1994, vol. 100 pt.2, pp. 1053–1062, ASHRAE.
- [47] R.M. Neal, "Assessing relevance determination methods using DELVE," in *Neural Networks and Machine Learning*, C.M. Bishop, Ed. 1998, pp. 97–129, Springer-Verlag.
- [48] M. E. Tipping, "The relevance vector machine," in *Advances in Neural Information Processing Systems*. 2000, vol. 12, pp. 652–658, MIT Press.
- [49] Z. Ghahramani and M. J. Beal, "Variational inference for Bayesian mixtures of factor analysers," in *Advances in Neural Information Processing Systems*. 2000, vol. 12, MIT Press.

- [50] F. Takens, "Detecting strange attractors in turbulence," in *Dynamical Systems and Turbulence*, D. A. Rand and L.-S. Young, Eds., Proceedings of 1980 Warwick conference, 1981, vol. 898 of *Lecture Notes in Mathematics*, pp. 366–381, Springer-Verlag, Berlin.
- [51] J. Hertz, A. Krogh, and R.G. Palmer, *Introduction to the Theory of Neural Computation*, Addison-Wesley, Redwood City, CA, 1991.
- [52] R. M. Neal, "Probabilistic inference using Markov chain monte carlo methods," Tech. Rep. CRG-TR-93-1, Department of Computer Science, University of Toronto, 1993.

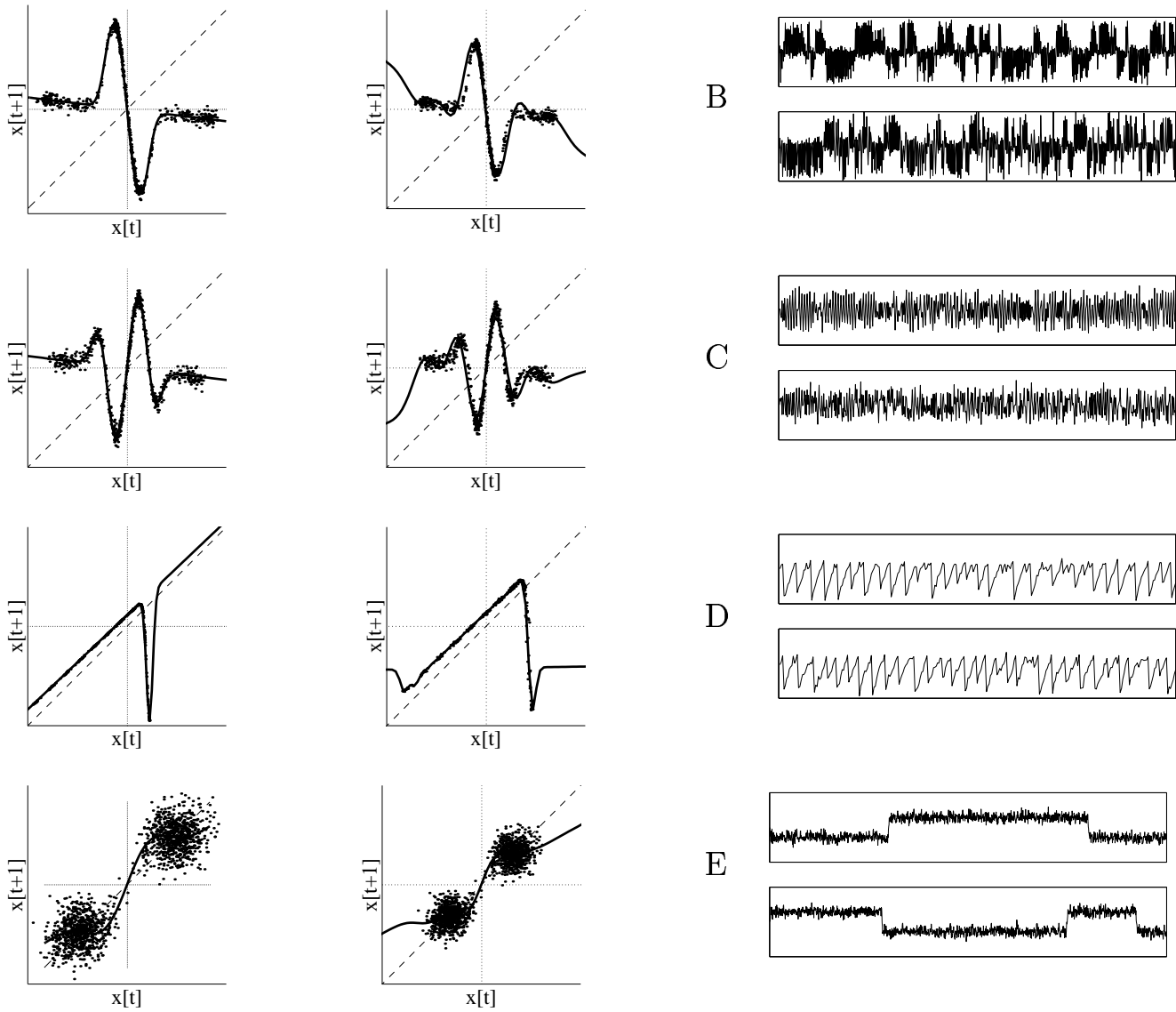


Fig. 6. More examples of fitting systems with nonlinear dynamics and linear observation functions. Each of the five rows shows the fitting of a nonlinear system with a one-dimensional hidden state and 4 noisy outputs driven by Gaussian noise inputs and internal state noise. **(left)** The true dynamics function (line) and states (dots) used to generate the training data. **(middle)** The learned dynamics function and states inferred on the training data. **(right)** The first component of the observable time series: training data on the top and fantasy data generated from the learned model on the bottom. The nonlinear dynamics can produce quasi-periodic outputs in response to white driving noise.