

A Unifying Review of Linear Gaussian Models

Sam Roweis*

Computation and Neural Systems, California Institute of Technology, Pasadena, CA 91125, U.S.A.

Zoubin Ghahramani*

Department of Computer Science, University of Toronto, Toronto, Canada

Factor analysis, principal component analysis, mixtures of gaussian clusters, vector quantization, Kalman filter models, and hidden Markov models can all be unified as variations of unsupervised learning under a single basic generative model. This is achieved by collecting together disparate observations and derivations made by many previous authors and introducing a new way of linking discrete and continuous state models using a simple nonlinearity. Through the use of other nonlinearities, we show how independent component analysis is also a variation of the same basic generative model. We show that factor analysis and mixtures of gaussians can be implemented in autoencoder neural networks and learned using squared error plus the same regularization term. We introduce a new model for static data, known as sensible principal component analysis, as well as a novel concept of spatially adaptive observation noise. We also review some of the literature involving global and local mixtures of the basic models and provide pseudocode for inference and learning for all the basic models.

1 A Unifying Review ---

Many common statistical techniques for modeling multidimensional static data sets and multidimensional time series can be seen as variants of one underlying model. As we will show, these include factor analysis, principal component analysis (PCA), mixtures of gaussian clusters, vector quantization, independent component analysis models (ICA), Kalman filter models (also known as linear dynamical systems), and hidden Markov models (HMMs). The relationships between some of these models has been noted in passing in the recent literature. For example, Hinton, Revow, and Dayan (1995) note that FA and PCA are closely related, and Digalakis, Rohlicek, and Ostendorf (1993) relate the forward-backward algorithm for HMMs to

* Present address: {roweis, zoubin}@gatsby.ucl.ac.uk. Gatsby Computational Neuroscience Unit, University College London, 17 Queen Square, London WC1N 3AR U.K.

Kalman filtering. In this article we unify many of the disparate observations made by previous authors (Rubin & Thayer, 1982; Delyon, 1993; Digalakis et al., 1993; Hinton et al., 1995; Elliott, Aggoun, & Moore, 1995; Ghahramani & Hinton, 1996a,b, 1997; Hinton & Ghahramani, 1997) and present a review of all these algorithms as instances of a single basic generative model. This unified view allows us to show some interesting relations between previously disparate algorithms. For example, factor analysis and mixtures of gaussians can be implemented using autoencoder neural networks with different nonlinearities but learned using a squared error cost penalized by the same regularization term. ICA can be seen as a nonlinear version of factor analysis. The framework also makes it possible to derive a new model for static data that is based on PCA but has a sensible probabilistic interpretation, as well as a novel concept of spatially adaptive observation noise. We also review some of the literature involving global and local mixtures of the basic models and provide pseudocode (in the appendix) for inference and learning for all the basic models.

2 The Basic Model

The basic models we work with are discrete time linear dynamical systems with gaussian noise. In such models we assume that the state of the process in question can be summarized at any time by a k -vector of state variables or causes \mathbf{x} that we cannot observe directly. However, the system also produces at each time step an output or observable p -vector \mathbf{y} to which we do have access.

The state \mathbf{x} is assumed to evolve according to simple first-order Markov dynamics; each output vector \mathbf{y} is generated from the current state by a simple linear observation process. Both the state evolution and the observation processes are corrupted by additive gaussian noise, which is also hidden. If we work with a continuous valued state variable \mathbf{x} , the basic generative model can be written¹ as:

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t = \mathbf{A}\mathbf{x}_t + \mathbf{w}_\bullet \quad \mathbf{w}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (2.1a)$$

$$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_\bullet \quad \mathbf{v}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \quad (2.1b)$$

where \mathbf{A} is the $k \times k$ state transition matrix and \mathbf{C} is the $p \times k$ observation, measurement, or generative matrix.

The k -vector \mathbf{w} and p -vector \mathbf{v} are random variables representing the state evolution and observation noises, respectively, which are independent

¹ All vectors are column vectors. To denote the transpose of a vector or matrix, we use the notation \mathbf{x}^T . The determinant of a matrix is denoted by $|\mathbf{A}|$ and matrix inversion by \mathbf{A}^{-1} . The symbol \sim means “distributed according to.” 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 \mathbf{z} is denoted $\mathcal{N}(\mu, \Sigma) |_{\mathbf{z}}$.

of each other and of the values of \mathbf{x} and \mathbf{y} . Both of these noise sources are temporally white (uncorrelated from time step to time step) and spatially gaussian distributed² with zero mean and covariance matrices, which we denote \mathbf{Q} and \mathbf{R} , respectively. We have written \mathbf{w}_\bullet and \mathbf{v}_\bullet in place of \mathbf{w}_t and \mathbf{v}_t to emphasize that the noise processes do not have any knowledge of the time index. The restriction to zero-mean noise sources is not a loss of generality.³ Since the state evolution noise is gaussian and its dynamics are linear, \mathbf{x}_t is a first-order Gauss-Markov random process. The noise processes are essential elements of the model. Without the process noise \mathbf{w}_\bullet , the state \mathbf{x}_t would always either shrink exponentially to zero or blow up exponentially in the direction of the leading eigenvector of \mathbf{A} ; similarly in the absence of the observation noise \mathbf{v}_\bullet the state would no longer be hidden. Figure 1 illustrates this basic model using both the engineering system block form and the network form more common in machine learning.

Notice that there is degeneracy in the model. All of the structure in the matrix \mathbf{Q} can be moved into the matrices \mathbf{A} and \mathbf{C} . This means that we can, without loss of generality, work with models in which \mathbf{Q} is the identity matrix.⁴ Of course, \mathbf{R} cannot be restricted in the same way since the values \mathbf{y}_t are observed, and hence we are not free to whiten or otherwise rescale them. Finally, the components of the state vector can be arbitrarily reordered; this corresponds to swapping the columns of \mathbf{C} and \mathbf{A} . Typically we choose an ordering based on the norms of the columns of \mathbf{C} , which resolves this degeneracy.

The network diagram of Figure 1 can be unfolded in time to give separate units for each time step. Such diagrams are the standard method of illustrating graphical models, also known as probabilistic independence networks, a category of models that includes Markov networks, Bayesian (or belief) networks, and other formalisms (Pearl, 1988; Lauritzen & Spiegelhalter, 1988; Whittaker, 1990; Smyth et al., 1997). A graphical model is a representation of the dependency structure between variables in a multivariate probability distribution. Each node corresponds to a random variable, and the absence of an arc between two variables corresponds to a particular conditional independence relation. Although graphical models are beyond

² An assumption that is weakly motivated by the central limit theorem but more strongly by analytic tractability.

³ Specifically we could always add a $k + 1$ st dimension to the state vector, which is fixed at unity. Then augmenting \mathbf{A} with an extra column holding the noise mean and an extra row of zeros (except unity in the bottom right corner) takes care of a nonzero mean for \mathbf{w}_\bullet . Similarly, adding an extra column to \mathbf{C} takes care of a nonzero mean for \mathbf{v}_\bullet .

⁴ In particular, since it is a covariance matrix, \mathbf{Q} is symmetric positive semidefinite and thus can be diagonalized to the form $\mathbf{E}\mathbf{\Lambda}\mathbf{E}^T$ (where \mathbf{E} is a rotation matrix of eigenvectors and $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues). Thus, for any model in which \mathbf{Q} is not the identity matrix, we can generate an exactly equivalent model using a new state vector $\mathbf{x}' = \mathbf{\Lambda}^{-1/2}\mathbf{E}^T\mathbf{x}$ with $\mathbf{A}' = (\mathbf{\Lambda}^{-1/2}\mathbf{E}^T)\mathbf{A}(\mathbf{E}\mathbf{\Lambda}^{1/2})$ and $\mathbf{C}' = \mathbf{C}(\mathbf{E}\mathbf{\Lambda}^{1/2})$ such that the new covariance of \mathbf{x}' is the identity matrix: $\mathbf{Q}' = \mathbf{I}$.

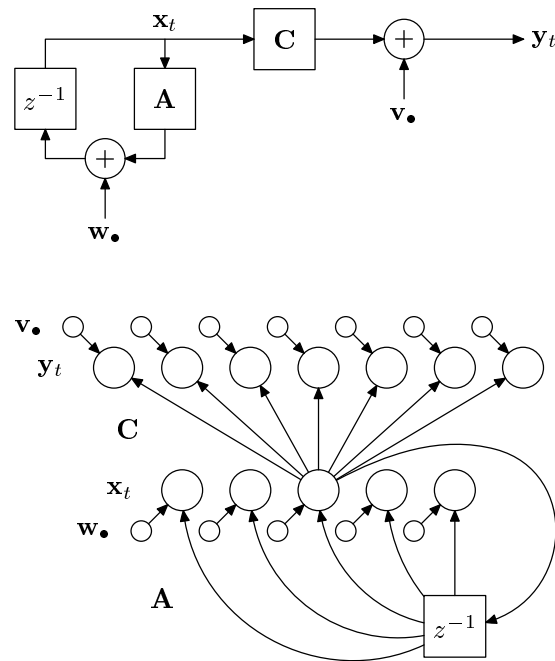


Figure 1: Linear dynamical system generative model. The z^{-1} block is a unit delay. The covariance matrix of the input noise w is Q and the covariance matrix of the output noise v is R . In the network model below, the smaller circles represent noise sources and all units are linear. Outgoing weights have only been drawn from one hidden unit. This model is equivalent to a Kalman filter model (linear dynamical system).

the scope of this review, it is important to point out that they provide a very general framework for working with the models we consider here. In this review, we unify and extend some well-known statistical models and signal processing algorithms by focusing on variations of linear graphical models with gaussian noise.

The main idea of the models in equations 2.1 is that the hidden state sequence x_t should be an informative lower dimensional projection or explanation of the complicated observation sequence y_t . With the aid of the dynamical and noise models, the states should summarize the underlying causes of the data much more succinctly than the observations themselves. For this reason, we often work with state dimensions much smaller than the number of observables—in other words, $k \ll p$.⁵ We assume that both

⁵ More precisely, in a model where all the matrices are full rank, the problem of inferring

\mathbf{A} and \mathbf{C} are of rank k and that \mathbf{Q} , \mathbf{R} , and \mathbf{Q}_1 (introduced below) are always full rank.

3 Probability Computations

The popularity of linear gaussian models comes from two fortunate analytical properties of gaussian processes: the sum of two independent gaussian distributed quantities is also gaussian distributed,⁶ and the output of a linear system whose input is gaussian distributed is again gaussian distributed. This means that if we assume the initial state \mathbf{x}_1 of the system to be gaussian distributed,

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{Q}_1), \quad (3.1)$$

then all future states \mathbf{x}_t and observations \mathbf{y}_t will also be gaussian distributed. In fact, we can write explicit formulas for the conditional expectations of the states and observables:

$$P(\mathbf{x}_{t+1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{A}\mathbf{x}_t, \mathbf{Q})|_{\mathbf{x}_{t+1}}, \quad (3.2a)$$

$$P(\mathbf{y}_t|\mathbf{x}_t) = \mathcal{N}(\mathbf{C}\mathbf{x}_t, \mathbf{R})|_{\mathbf{y}_t}. \quad (3.2b)$$

Furthermore, because of the Markov properties of the model and the gaussian assumptions about the noise and initial distributions, it is easy to write an expression for the joint probability of a sequence of τ states and outputs:

$$P(\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}, \{\mathbf{y}_1 \dots \mathbf{y}_\tau\}) = P(\mathbf{x}_1) \prod_{t=1}^{\tau-1} P(\mathbf{x}_{t+1}|\mathbf{x}_t) \prod_{t=1}^{\tau} P(\mathbf{y}_t|\mathbf{x}_t). \quad (3.3)$$

The negative log probability (cost) is just the sum of matrix quadratic forms:

$$\begin{aligned} & -2 \log P(\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}, \{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}) \\ &= \sum_{t=1}^{\tau} [(\mathbf{y}_t - \mathbf{C}\mathbf{x}_t)^T \mathbf{R}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t) + \log |\mathbf{R}|] \end{aligned}$$

the state from a sequence of τ consecutive observations is well defined as long $k \leq \tau p$ (a notion related to observability in systems theory; Goodwin & Sin, 1984). For this reason, in dynamic models it is sometimes useful to use state-spaces of larger dimension than the observations, $k > p$, in which case a single state vector provides a compact representation of a sequence of observations.

⁶ In other words the convolution of two gaussians is again a gaussian. In particular, the convolution of $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ is $\mathcal{N}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)$. This is not the same as the (false) statement that the sum of two gaussians is a gaussian but is the same as the (Fourier domain equivalent) statement that the multiplication of two gaussians is a gaussian (although no longer normalized).

$$\begin{aligned}
& + \sum_{t=1}^{\tau-1} [(\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \mathbf{Q}^{-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) + \log |\mathbf{Q}|] \\
& + (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \mathbf{Q}_1^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + \log |\mathbf{Q}_1| + \tau(p+k) \log 2\pi. \quad (3.4)
\end{aligned}$$

4 Learning and Estimation Problems

Latent variable models have a wide spectrum of application in data analysis. In some scenarios, we know exactly what the hidden states are supposed to be and just want to estimate them. For example, in a vision problem, the hidden states may be the location or pose of an object; in a tracking problem, the states may be positions and momenta. In these cases, we can often write down a priori the observation or state evolution matrices based on our knowledge of the problem structure or physics. The emphasis is on accurate inference of the unobserved information from the data we do have—for example, from an image of an object or radar observations. In other scenarios, we are trying to discover explanations or causes for our data and have no explicit model for what these causes should be. The observation and state evolution processes are mostly or entirely unknown. The emphasis is instead on robustly learning a few parameters that model the observed data well (assign it high likelihood). Speech modeling is a good example of such a situation; our goal is to find economical models that perform well for recognition tasks, but the particular values of hidden states in our models may not be meaningful or important to us. These two goals—estimating the hidden states given observations and a model and learning the model parameters—typically manifest themselves in the solution of two distinct problems: inference and system identification.

4.1 Inference: Filtering and Smoothing. The first problem is that of inference or filtering and smoothing, which asks: Given fixed model parameters $\{\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}, \boldsymbol{\mu}_1, \mathbf{Q}_1\}$, what can be said about the unknown hidden state sequence given some observations? This question is typically made precise in several ways. A very basic quantity we would like to be able to compute is the total likelihood of an observation sequence:

$$\begin{aligned}
& P(\{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}) \\
& = \int_{\text{all possible } \{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}} P(\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}, \{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}) d\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}. \quad (4.1)
\end{aligned}$$

This marginalization requires an efficient way of integrating (or summing) the joint probability (easily computed by equation 3.4 or similar formulas) over all possible paths through state-space.

Once this integral is available, it is simple to compute the conditional distribution for any one proposed hidden state sequence given the observations

by dividing the joint probability by the total likelihood of the observations:

$$P(\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}|\{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}) = \frac{P(\{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}, \{\mathbf{y}_1, \dots, \mathbf{y}_\tau\})}{P(\{\mathbf{y}_1, \dots, \mathbf{y}_\tau\})}. \quad (4.2)$$

Often we are interested in the distribution of the hidden state at a particular time t . In filtering, we attempt to compute this conditional posterior probability,

$$P(\mathbf{x}_t|\{\mathbf{y}_1, \dots, \mathbf{y}_t\}), \quad (4.3)$$

given all the observations up to and including time t . In smoothing, we compute the distribution over \mathbf{x}_t ,

$$P(\mathbf{x}_t|\{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}), \quad (4.4)$$

given the entire sequence of observations. (It is also possible to ask for the conditional state expectation given observations that extend only a few time steps into the future—partial smoothing—or that end a few time steps before the current time—partial prediction.) These conditional calculations are closely related to the computation of equation 4.1 and often the intermediate values of a recursive method used to compute that equation give the desired distributions of equations 4.3 or 4.4. Filtering and smoothing have been extensively studied for continuous state models in the signal processing community, starting with the seminal works of Kalman (1960; Kalman & Bucy, 1961) and Rauch (1963; Rauch, Tung, & Striebel, 1965), although this literature is often not well known in the machine learning community. For the discrete state models, much of the literature stems from the work of Baum and colleagues (Baum & Petrie, 1966; Baum & Eagon, 1967; Baum, Petrie, Soules, & Weiss, 1970; Baum, 1972) on HMMs and of Viterbi (1967) and others on optimal decoding. The recent book by Elliott and colleagues (1995) contains a thorough mathematical treatment of filtering and smoothing for many general models.

4.2 Learning (System Identification). The second problem of interest with linear gaussian models is the learning or system identification problem: given only an observed sequence (or perhaps several sequences) of outputs $\{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}$ find the parameters $\{\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}, \boldsymbol{\mu}_1, \mathbf{Q}_1\}$ that maximize the likelihood of the observed data as computed by equation 4.1.

The learning problem has been investigated extensively by neural network researchers for static models and also for some discrete state dynamic models such as HMMs or the more general Bayesian belief networks. There is a corresponding area of study in control theory known as system identification, which investigates learning in continuous state models. For linear gaussian models, there are several approaches to system identification

(Ljung & Söderström, 1983), but to clarify the relationship between these models and the others we review in this article, we focus on system identification methods based on the expectation-maximization (EM) algorithm. The EM algorithm for linear gaussian dynamical systems was originally derived by Shumway and Stoffer (1982) and recently reintroduced (and extended) in the neural computation field by Ghahramani and Hinton (1996a,b). Digalakis et al. (1993) made a similar reintroduction and extension in the speech processing community. Once again we mention the book by Elliott et al. (1995), which also covers learning in this context.

The basis of all the learning algorithms presented by these authors is the powerful EM algorithm (Baum & Petrie, 1966; Dempster, Laird, & Rubin, 1977). The objective of the algorithm is to maximize the likelihood of the observed data (equation 4.1) in the presence of hidden variables. Let us denote the observed data by $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_\tau\}$, the hidden variables by $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_\tau\}$, and the parameters of the model by θ . Maximizing the likelihood as a function of θ is equivalent to maximizing the log-likelihood:

$$\mathcal{L}(\theta) = \log P(\mathbf{Y}|\theta) = \log \int_{\mathbf{X}} P(\mathbf{X}, \mathbf{Y}|\theta) d\mathbf{X}. \quad (4.5)$$

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

$$\log \int_{\mathbf{X}} P(\mathbf{Y}, \mathbf{X}|\theta) d\mathbf{X} = \log \int_{\mathbf{X}} Q(\mathbf{X}) \frac{P(\mathbf{X}, \mathbf{Y}|\theta)}{Q(\mathbf{X})} d\mathbf{X} \quad (4.6a)$$

$$\geq \int_{\mathbf{X}} Q(\mathbf{X}) \log \frac{P(\mathbf{X}, \mathbf{Y}|\theta)}{Q(\mathbf{X})} d\mathbf{X} \quad (4.6b)$$

$$= \int_{\mathbf{X}} Q(\mathbf{X}) \log P(\mathbf{X}, \mathbf{Y}|\theta) d\mathbf{X} - \int_{\mathbf{X}} Q(\mathbf{X}) \log Q(\mathbf{X}) d\mathbf{X} \quad (4.6c)$$

$$= \mathcal{F}(Q, \theta), \quad (4.6d)$$

where the middle inequality is known as Jensen's inequality and can be proved using the concavity of the log function. If we define the energy of a global configuration (\mathbf{X}, \mathbf{Y}) to be $-\log P(\mathbf{X}, \mathbf{Y}|\theta)$, then some readers may notice that 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 (Neal & Hinton, 1998). 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 :

$$\text{E-step:} \quad Q_{k+1} \leftarrow \arg \max_Q \mathcal{F}(Q, \theta_k) \quad (4.7a)$$

$$\mathbf{M}\text{-step:} \quad \theta_{k+1} \leftarrow \arg \max_{\theta} \mathcal{F}(Q_{k+1}, \theta). \quad (4.7b)$$

It is easy to show that the maximum in the E-step results when Q is exactly the conditional distribution of \mathbf{X} : $Q_{k+1}(\mathbf{X}) = P(\mathbf{X}|\mathbf{Y}, \theta_k)$, 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 first term in equation 4.6c, since the entropy of Q does not depend on θ :

$$\mathbf{M}\text{-step:} \quad \theta_{k+1} \leftarrow \arg \max_{\theta} \int_{\mathbf{X}} P(\mathbf{X}|\mathbf{Y}, \theta_k) \log P(\mathbf{X}, \mathbf{Y}|\theta) d\mathbf{X}. \quad (4.8)$$

This is the expression most often associated with the EM algorithm, but it obscures the elegant interpretation of EM as coordinate ascent in \mathcal{F} (Neal & Hinton, 1998). 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.

Therefore, at the heart of the EM learning procedure is the following idea: use the solutions to the filtering and 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 linear gaussian models, this typically involves minimizing quadratic forms such as equation 3.4, 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. We shall review the details of particular algorithms as we present the various cases; however, we now touch on one general point that often causes confusion. Our goal is to maximize the *total* likelihood (see equation 4.1) (or equivalently maximize the total log likelihood) of the observed data with respect to the model parameters. This means integrating (or summing) over all ways in which the generative model could have produced the data. 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, where the expectation is taken over the distribution of hidden values predicted by the current model parameters and the observations. Thus, it appears that we are maximizing the incorrect quantity, but doing so is in fact guaranteed to increase (or keep the same) the quantity of interest at each iteration of the algorithm.

5 Continuous-State Linear Gaussian Systems ---

Having described the basic model and learning procedure, we now focus on specific linear instances of the model in which the hidden state variable \mathbf{x} is continuous and the noise processes are gaussian. This will allow us to

elucidate the relationship among factor analysis, PCA, and Kalman filter models. We divide our discussion into models that generate static data and those that generate dynamic data. Static data have no temporal dependence; no information would be lost by permuting the ordering of the data points \mathbf{y}_i ; whereas for dynamic data, the time ordering of the data points is crucial.

5.1 Static Data Modeling: Factor Analysis, SPCA, and PCA. In many situations we have reason to believe (or at least to assume) that each point in our data set was generated independently and identically. In other words, there is no natural (temporal) ordering to the data points; they merely form a collection. In such cases, we assume that the underlying state vector \mathbf{x} has no dynamics; the matrix \mathbf{A} is the zero matrix, and therefore \mathbf{x} is simply a constant (which we take without loss of generality to be the zero vector) corrupted by noise. The new generative model then becomes:

$$\mathbf{A} = \mathbf{0} \quad \Rightarrow \quad \mathbf{x}_\bullet = \mathbf{w}_\bullet \quad \mathbf{w}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (5.1a)$$

$$\mathbf{y}_\bullet = \mathbf{C}\mathbf{x}_\bullet + \mathbf{v}_\bullet \quad \mathbf{v}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{R}). \quad (5.1b)$$

Notice that since \mathbf{x}_t is driven only by the noise \mathbf{w}_\bullet and since \mathbf{y}_t depends only on \mathbf{x}_t , all temporal dependence has disappeared. This is the motivation for the term *static* and for the notations \mathbf{x}_\bullet and \mathbf{y}_\bullet above. We also no longer use a separate distribution for the initial state: $\mathbf{x}_1 \sim \mathbf{x}_\bullet \sim \mathbf{w}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$.

This model is illustrated in Figure 2. We can analytically integrate equation 4.1 to obtain the marginal distribution of \mathbf{y}_\bullet , which is the gaussian,

$$\mathbf{y}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{C}\mathbf{Q}\mathbf{C}^T + \mathbf{R}). \quad (5.2)$$

Two things are important to notice. First, the degeneracy mentioned above persists between the structure in \mathbf{Q} and \mathbf{C} .⁷ This means there is no loss of generality in restricting \mathbf{Q} to be diagonal. Furthermore, there is arbitrary sharing of scale between a diagonal \mathbf{Q} and \mathbf{C} . Typically we either restrict the columns of \mathbf{C} to be unit vectors or make \mathbf{Q} the identity matrix to resolve this degeneracy. In what follows we will assume $\mathbf{Q} = \mathbf{I}$ without loss of generality.

Second, the covariance matrix \mathbf{R} of the observation noise must be restricted in some way for the model to capture any interesting or informative projections in the state \mathbf{x}_\bullet . If \mathbf{R} were not restricted, learning could simply choose $\mathbf{C} = \mathbf{0}$ and then set \mathbf{R} to be the sample covariance of the data, thus trivially achieving the maximum likelihood model by explaining all of the

⁷ If we diagonalize \mathbf{Q} and rewrite the covariance of \mathbf{y}_\bullet , the degeneracy becomes clear: $\mathbf{y}_\bullet \sim \mathcal{N}(\mathbf{0}, (\mathbf{C}\mathbf{E}\mathbf{A}^{1/2})(\mathbf{C}\mathbf{E}\mathbf{A}^{1/2})^T + \mathbf{R})$. To make \mathbf{Q} diagonal, we simply replace \mathbf{C} with $\mathbf{C}\mathbf{E}$.

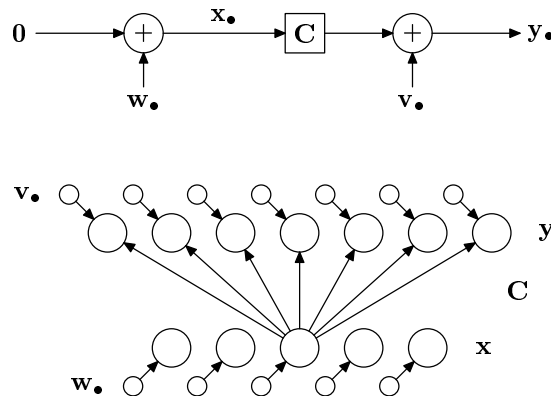


Figure 2: Static generative model (continuous state). The covariance matrix of the input noise w is Q and the covariance matrix of the output noise v is R . In the network model below, the smaller circles represent noise sources and all units are linear. Outgoing weights have only been drawn from one hidden unit. This model is equivalent to factor analysis, SPCA and PCA models depending on the output noise covariance. For factor analysis, $Q = I$ and R is diagonal. For SPCA, $Q = I$ and $R = \alpha I$. For PCA, $Q = I$ and $R = \lim_{\epsilon \rightarrow 0} \epsilon I$.

structure in the data as noise. (Remember that since the model has reduced to a single gaussian distribution for y_{\bullet} , we can do no better than having the covariance of our model equal the sample covariance of our data.) Note that restricting R , unlike making Q diagonal, does constitute some loss of generality from the original model of equations 5.1.

There is an intuitive spatial way to think about this static generative model. We use white noise to generate a spherical ball (since $Q = I$) of density in k -dimensional state-space. This ball is then stretched and rotated into p -dimensional observation space by the matrix C , where it looks like a k -dimensional pancake. The pancake is then convolved with the covariance density of v_{\bullet} (described by R) to get the final covariance model for y_{\bullet} . We want the resulting ellipsoidal density to be as close as possible to the ellipsoid given by the sample covariance of our data. If we restrict the shape of the v_{\bullet} covariance by constraining R , we can force interesting information to appear in both R and C as a result.

Finally, observe that all varieties of filtering and smoothing reduce to the same problem in this static model because there is no time dependence. We are seeking only the posterior probability $P(x_{\bullet} | y_{\bullet})$ over a single hidden state given the corresponding single observation. This inference is easily done by

linear matrix projection, and the resulting density is itself gaussian:

$$P(\mathbf{x}_\bullet | \mathbf{y}_\bullet) = \frac{P(\mathbf{y}_\bullet | \mathbf{x}_\bullet) P(\mathbf{x}_\bullet)}{P(\mathbf{y}_\bullet)} = \frac{\mathcal{N}(\mathbf{C}\mathbf{x}_\bullet, \mathbf{R}) |_{\mathbf{y}_\bullet} \mathcal{N}(\mathbf{0}, \mathbf{I}) |_{\mathbf{x}_\bullet}}{\mathcal{N}(\mathbf{0}, \mathbf{C}\mathbf{C}^T + \mathbf{R}) |_{\mathbf{y}_\bullet}} \quad (5.3a)$$

$$P(\mathbf{x}_\bullet | \mathbf{y}_\bullet) = \mathcal{N}(\beta \mathbf{y}_\bullet, \mathbf{I} - \beta \mathbf{C}) |_{\mathbf{x}_\bullet}, \quad \beta = \mathbf{C}^T (\mathbf{C}\mathbf{C}^T + \mathbf{R})^{-1}, \quad (5.3b)$$

from which we obtain not only the expected value $\beta \mathbf{y}_\bullet$ of the unknown state but also an estimate of the uncertainty in this value in the form of the covariance $\mathbf{I} - \beta \mathbf{C}$. Computing the likelihood of a data point \mathbf{y}_\bullet is merely an evaluation under the gaussian in equation 5.2. The learning problem now consists of identifying the matrices \mathbf{C} and \mathbf{R} . There is a family of EM algorithms to do this for the various cases discussed below, which are given in detail at the end of this review.

5.2 Factor Analysis. If we restrict the covariance matrix \mathbf{R} that controls the observation noise to be diagonal (in other words, the covariance ellipsoid of \mathbf{v}_\bullet is axis aligned) and set the state noise \mathbf{Q} to be the identity matrix, then we recover exactly a standard statistical model known as maximum likelihood factor analysis. The unknown states \mathbf{x} are called the factors in this context; the matrix \mathbf{C} is called the factor loading matrix, and the diagonal elements of \mathbf{R} are often known as the uniquenesses. (See Everitt, 1984, for a brief and clear introduction.) The inference calculation is done exactly as in equation 5.3b. The learning algorithm for the loading matrix and the uniquenesses is exactly an EM algorithm except that we must take care to constrain \mathbf{R} properly (which is as easy as taking the diagonal of the unconstrained maximum likelihood estimate; see Rubin & Thayer, 1982; Ghahramani & Hinton, 1997). If \mathbf{C} is completely free, this procedure is called exploratory factor analysis; if we build a priori zeros into \mathbf{C} , it is confirmatory factor analysis. In exploratory factor analysis, we are trying to model the covariance structure of our data with $p + pk - k(k-1)/2$ free parameters⁸ instead of the $p(p+1)/2$ free parameters in a full covariance matrix.

The diagonality of \mathbf{R} is the key assumption here. Factor analysis attempts to explain the covariance structure in the observed data by putting all the variance unique to each coordinate in the matrix \mathbf{R} and putting all the correlation structure into \mathbf{C} (this observation was first made by Lyttkens, 1966, in response to work by Wold). In essence, factor analysis considers the axis rotation in which the original data arrived to be special because observation noise (often called *sensor noise*) is independent along the coordinates in these axes. However, the original scaling of the coordinates is unimportant. If we were to change the units in which we measured some of the components of \mathbf{y} , factor analysis could merely rescale the corresponding entry in \mathbf{R} and

⁸ The correction $k(k-1)/2$ comes in because of degeneracy in unitary transformations of the factors. See, for example, Everitt (1984).

row in \mathbf{C} and achieve a new model that assigns the rescaled data identical likelihood. On the other hand, if we rotate the axes in which we measure the data, we could not easily fix things since the noise \mathbf{v} is constrained to have axis aligned covariance (\mathbf{R} is diagonal).

EM for factor analysis has been criticized as being quite slow (Rubin & Thayer, 1982). Indeed, the standard method for fitting a factor analysis model (Jöreskog, 1967) is based on a quasi-Newton optimization algorithm (Fletcher & Powell, 1963), which has been found empirically to converge faster than EM. We present the EM algorithm here not because it is the most efficient way of fitting a factor analysis model, but because we wish to emphasize that for factor analysis and all the other latent variable models reviewed here, EM provides a unified approach to learning. Finally, recent work in online learning has shown that it is possible to derive a family of EM-like algorithms with faster convergence rates than the standard EM algorithm (Kivinen & Warmuth, 1997; Bauer, Koller, & Singer, 1997).

5.3 SPCA and PCA. If instead of restricting \mathbf{R} to be merely diagonal, we require it to be a multiple of the identity matrix (in other words, the covariance ellipsoid of \mathbf{v}_\bullet is spherical), then we have a model that we will call *sensible principal component analysis* (SPCA) (Roweis, 1997). The columns of \mathbf{C} span the principal subspace (the same subspace found by PCA), and we will call the scalar value on the diagonal of \mathbf{R} the global noise level. Note that SPCA uses $1 + pk - k(k - 1)/2$ free parameters to model the covariance. Once again, inference is done with equation 5.3b and learning by the EM algorithm (except that we now take the trace of the maximum likelihood estimate for \mathbf{R} to learn the noise level; see (Roweis, 1997)). Unlike factor analysis, SPCA considers the original axis rotation in which the data arrived to be unimportant: if the measurement coordinate system were rotated, SPCA could (left) multiply \mathbf{C} by the same rotation, and the likelihood of the new data would not change. On the other hand, the original scaling of the coordinates is privileged because SPCA assumes that the observation noise has the same variance in all directions in the measurement units used for the observed data. If we were to rescale one of the components of \mathbf{y} , the model could not be easily corrected since \mathbf{v} has spherical covariance ($\mathbf{R} = \epsilon \mathbf{I}$). The SPCA model is very similar to the independently proposed probabilistic principal component analysis (Tipping & Bishop, 1997).

If we go even further and take the limit $\mathbf{R} = \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{I}$ (while keeping the diagonal elements of \mathbf{Q} finite)⁹ then we obtain the standard principal component analysis (PCA) model. The directions of the columns of \mathbf{C} are

⁹ Since isotropic scaling of the data space is arbitrary, we could just as easily take the limit as the diagonal elements of \mathbf{Q} became infinite while holding \mathbf{R} finite or take both limits at once. The idea is that the noise variance becomes infinitesimal compared to the scale of the data.

known as the principal components. Inference now reduces to simple least squares projection:¹⁰

$$P(\mathbf{x}_\bullet | \mathbf{y}_\bullet) = \mathcal{N}(\beta \mathbf{y}_\bullet, I - \beta \mathbf{C}) |_{\mathbf{x}_\bullet}, \quad \beta = \lim_{\epsilon \rightarrow 0} \mathbf{C}^T (\mathbf{C} \mathbf{C}^T + \epsilon \mathbf{I})^{-1} \quad (5.4a)$$

$$\begin{aligned} P(\mathbf{x}_\bullet | \mathbf{y}_\bullet) &= \mathcal{N}\left((\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{y}_\bullet, \mathbf{0}\right) |_{\mathbf{x}_\bullet} \\ &= \delta(\mathbf{x}_\bullet - (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{y}_\bullet). \end{aligned} \quad (5.4b)$$

Since the noise has become infinitesimal, the posterior over states collapses to a single point, and the covariance becomes zero. There is still an EM algorithm for learning (Roweis, 1997), although it can learn only \mathbf{C} . For PCA, we could just diagonalize the sample covariance of the data and take the leading k eigenvectors multiplied by their eigenvalues to be the columns of \mathbf{C} . This approach would give us \mathbf{C} in one step but has many problems.¹¹ The EM learning algorithm amounts to an iterative procedure for finding these leading eigenvectors without explicit diagonalization.

An important final comment is that (regular) PCA does not define a proper density model in the observation space, so we cannot ask directly about the likelihood assigned by the model to some data. We can, however, examine a quantity that is proportional to the negative log-likelihood in the limit of zero noise. This is the sum squared deviation of each data point from its projection. It is this “cost” that the learning algorithm ends up minimizing and is the only available evaluation of how well a PCA model fits new data. This is one of the most critical failings of PCA: translating points by arbitrary amounts inside the principal subspace has no effect on the model error.

¹⁰ Recall that if \mathbf{C} is $p \times k$ with $p > k$ and is rank k , then left multiplication by $\mathbf{C}^T (\mathbf{C} \mathbf{C}^T)^{-1}$ (which appears not to be well defined because $\mathbf{C} \mathbf{C}^T$ is not invertible) is exactly equivalent to left multiplication by $(\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T$. This is the same as the singular value decomposition idea of defining the “inverse” of the diagonal singular value matrix as the inverse of an element unless it is zero, in which case it remains zero. The intuition is that although $\mathbf{C} \mathbf{C}^T$ truly is not invertible, the directions along which it is not invertible are exactly those that \mathbf{C}^T is about to project out.

¹¹ It is computationally very hard to diagonalize or invert large matrices. It also requires an enormous amount of data to make a large sample covariance matrix full rank. If we are working with patterns in a large (thousands) number of dimensions and want to extract only a few (tens) principal components, we cannot naively try to diagonalize the sample covariance of our data. Techniques like the snapshot method (Sirovich, 1987) attempt to address this but still require the diagonalization of an $N \times N$ matrix where N is the number of data points. The EM algorithm approach solves all of these problems, requiring no explicit diagonalization whatsoever and the inversion of only a $k \times k$ matrix. It is guaranteed to converge to the true principal subspace (the same subspace spanned by the principal components). Empirical experiments (Roweis, 1997) indicate that it converges in a few iterations, unless the ratio of the leading eigenvalues is near unity.

5.4 Time-Series Modeling: Kalman Filter Models. We use the term *dynamic data* to refer to observation sequences in which the temporal ordering is important. For such data, we do not want to ignore the state evolution dynamics, which provides the only aspect of the model capable of capturing temporal structure. Systems described by the original dynamic generative model, shown in equations 2.1a and 2.2b, are known as *linear dynamical systems* or *Kalman filter models* and have been extensively investigated by the engineering and control communities for decades. The emphasis has traditionally been on inference problems: the famous discrete Kalman filter (Kalman, 1960; Kalman & Bucy, 1961) gives an efficient recursive solution to the optimal filtering and likelihood computation problems, while the RTS recursions (Rauch, 1963; Rauch et al., 1965) solve the optimal smoothing problem. Learning of unknown model parameters was studied by Shumway and Stoffer (1982) (\mathbf{C} known) and by Ghahramani and Hinton (1996a) and Digalakis et al. (1993) (all parameters unknown). Figure 1 illustrates this model, and the appendix gives pseudocode for its implementation.

We can extend our spatial intuition of the static case to this dynamic model. As before, any point in state-space is surrounded by a ball (or ovoid) of density (described by \mathbf{Q}), which is stretched (by \mathbf{C}) into a pancake in observation space and then convolved with the observation noise covariance (described by \mathbf{R}). However, unlike the static case, in which we always centered our ball of density on the origin in state-space, the center of the state-space ball now “flows” from time step to time step. The flow is according to the field described by the eigenvalues and eigenvectors of the matrix \mathbf{A} . We move to a new point according to this flow field; then we center our ball on that point and pick a new state. From this new state, we again flow to a new point and then apply noise. If \mathbf{A} is the identity matrix (not the zero matrix), then the “flow” does not move us anywhere, and the state just evolves according to a random walk of the noise set by \mathbf{Q} .

6 Discrete-State Linear Gaussian Models

We now consider a simple modification of the basic continuous state model in which the state at any time takes on one of a finite number of discrete values. Many real-world processes, especially those that have distinct modes of operation, are better modeled by internal states that are not continuous. (It is also possible to construct models that have a mixed continuous and discrete state.) The state evolution is still first-order Markovian dynamics, and the observation process is still linear with additive gaussian noise. The modification involves the use of the winner-take-all nonlinearity $\mathbf{WTA}[\cdot]$, defined such that $\mathbf{WTA}[\mathbf{x}]$ for any vector \mathbf{x} is a new vector with unity in the position of the largest coordinate of the input and zeros in all other positions. The discrete-state generative model is now simply:

$$\mathbf{x}_{t+1} = \mathbf{WTA}[\mathbf{A}\mathbf{x}_t + \mathbf{w}_t] = \mathbf{WTA}[\mathbf{A}\mathbf{x}_t + \mathbf{w}_\bullet] \quad (6.1a)$$

$$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_\bullet \quad (6.1b)$$

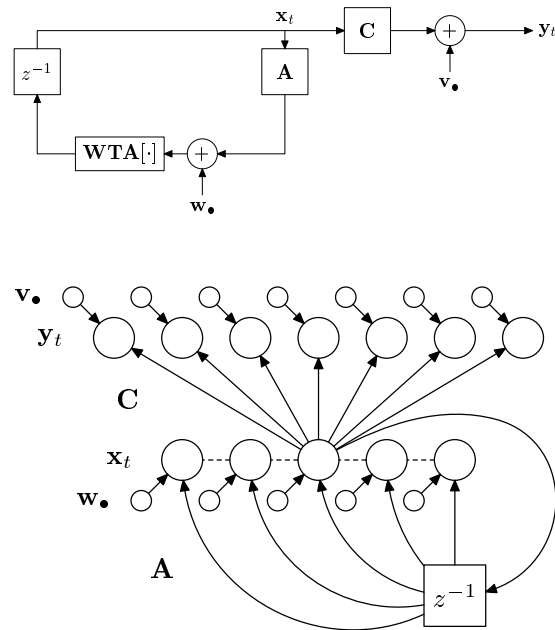


Figure 3: Discrete state generative model for dynamic data. The $\text{WTA}[\cdot]$ block implements the winner-take-all nonlinearity. The z^{-1} block is a unit delay. The covariance matrix of the input noise w is Q and the covariance matrix of the output noise v is R . In the network model below, the smaller circles represent noise sources and the hidden units x have a winner take all behaviour (indicated by dashed lines). Outgoing weights have only been drawn from one hidden unit. This model is equivalent to a hidden Markov model with tied output covariances.

where A is no longer known as the state transition matrix (although we will see that matrix shortly). As before, the k -vector w and p -vector v are temporally white and spatially gaussian distributed noises independent of each other and of x and y . The initial state x_1 is generated in the obvious way:

$$x_1 = \text{WTA}[\mathcal{N}(\mu_1, Q_1)] \quad (6.2)$$

(though we will soon see that without loss of generality Q_1 can be restricted to be the identity matrix). This discrete state generative model is illustrated in Figure 3.

6.1 Static Data Modeling: Mixtures of Gaussians and Vector Quantization. Just as in the continuous-state model, we can consider situations in

which there is no natural ordering to our data, and so set the matrix \mathbf{A} to be the zero matrix. In this discrete-state case, the generative model becomes:

$$\mathbf{A} = \mathbf{0} \quad \Rightarrow \quad \mathbf{x}_\bullet = \mathbf{WTA}[\mathbf{w}_\bullet] \quad \mathbf{w}_\bullet \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}) \quad (6.3)$$

$$\mathbf{y}_\bullet = \mathbf{C}\mathbf{x}_\bullet + \mathbf{v}_\bullet \quad \mathbf{v}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{R}). \quad (6.4)$$

Each state \mathbf{x}_\bullet is generated independently¹² according to a fixed discrete probability histogram controlled by the mean and covariance of \mathbf{w}_\bullet . Specifically, $\pi_j = P(\mathbf{x}_\bullet = \mathbf{e}_j)$ is the probability assigned by the gaussian $\mathcal{N}(\boldsymbol{\mu}, \mathbf{Q})$ to the region of k -space in which the j th coordinate is larger than all others. (Here \mathbf{e}_j is the unit vector along the j th coordinate direction.) Notice that to obtain nonuniform priors π_j with the $\mathbf{WTA}[\cdot]$ nonlinearity, we require a nonzero mean $\boldsymbol{\mu}$ for the noise \mathbf{w}_\bullet . Once the state has been chosen, the corresponding output \mathbf{y}_\bullet is generated from a gaussian whose mean is the j th column of \mathbf{C} and whose covariance is \mathbf{R} . This is exactly the standard mixture of gaussian clusters model except that the covariances of all the clusters are constrained to be the same. The probabilities $\pi_j = P(\mathbf{x}_\bullet = \mathbf{e}_j)$ correspond to the mixing coefficients of the clusters, and the columns of \mathbf{C} are the cluster means. Constraining \mathbf{R} in various ways corresponds to constraining the shape of the covariance of the clusters. This model is illustrated in Figure 4.

To compute the likelihood of a data point, we can explicitly perform the sum equivalent to the integral in equation 4.1 since it contains only k terms,

$$\begin{aligned} P(\mathbf{y}_\bullet) &= \sum_{i=1}^k P(\mathbf{x}_\bullet = \mathbf{e}_i, \mathbf{y}_\bullet) = \sum_{i=1}^k \mathcal{N}(\mathbf{C}_i, \mathbf{R}) |_{\mathbf{y}_\bullet} P(\mathbf{x}_\bullet = \mathbf{e}_i) \\ &= \sum_{i=1}^k \mathcal{N}(\mathbf{C}_i, \mathbf{R}) |_{\mathbf{y}_\bullet} \pi_i, \end{aligned} \quad (6.5)$$

where \mathbf{C}_i denotes the i th column of \mathbf{C} . Again, all varieties of inference and filtering are the same, and we are simply seeking the set of discrete probabilities $P(\mathbf{x}_\bullet = \mathbf{e}_j | \mathbf{y}_\bullet)$ $j = 1, \dots, k$. In other words, we need to do probabilistic classification. The problem is easily solved by computing the responsibilities $\hat{\mathbf{x}}_\bullet$ that each cluster has for the data point \mathbf{y}_\bullet :

$$\begin{aligned} (\hat{\mathbf{x}}_\bullet)_j &= P(\mathbf{x}_\bullet = \mathbf{e}_j | \mathbf{y}_\bullet) = \frac{P(\mathbf{x}_\bullet = \mathbf{e}_j, \mathbf{y}_\bullet)}{P(\mathbf{y}_\bullet)} = \frac{P(\mathbf{x}_\bullet = \mathbf{e}_j, \mathbf{y}_\bullet)}{\sum_{i=1}^k P(\mathbf{x}_\bullet = \mathbf{e}_i, \mathbf{y}_\bullet)} \quad (6.6a) \\ (\hat{\mathbf{x}}_\bullet)_j &= \frac{\mathcal{N}(\mathbf{C}_j, \mathbf{R}) |_{\mathbf{y}_\bullet} P(\mathbf{x}_\bullet = \mathbf{e}_j)}{\sum_{i=1}^k \mathcal{N}(\mathbf{C}_i, \mathbf{R}) |_{\mathbf{y}_\bullet} P(\mathbf{x}_\bullet = \mathbf{e}_i)} \end{aligned}$$

¹² As in the continuous static case, we again dispense with any special treatment of the initial state.

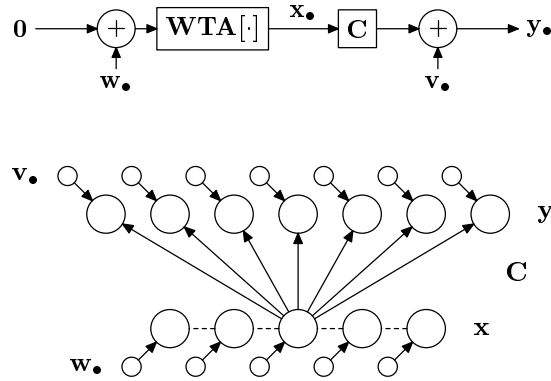


Figure 4: Static generative model (discrete state). The $\text{WTA}[\cdot]$ block implements the winner-take-all nonlinearity. The covariance matrix of the input noise \mathbf{w} is \mathbf{Q} and the covariance matrix of the output noise \mathbf{v} is \mathbf{R} . In the network model below, the smaller circles represent noise sources and the hidden units \mathbf{x} have a winner take all behaviour (indicated by dashed lines). Outgoing weights have only been drawn from one hidden unit. This model is equivalent to a mixture of Gaussian clusters with tied covariances \mathbf{R} or to vector quantization (VQ) when $\mathbf{R} = \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{I}$.

$$= \frac{\mathcal{N}(\mathbf{C}_j, \mathbf{R}) |y_{\bullet}, \pi_j}{\sum_{i=1}^k \mathcal{N}(\mathbf{C}_i, \mathbf{R}) |y_{\bullet}, \pi_i}. \quad (6.6b)$$

The mean $\hat{\mathbf{x}}_{\bullet}$ of the state vector given a data point is exactly the vector of responsibilities for that data point. This quantity defines the entire posterior distribution of the discrete hidden state given the data point. As a measure of the randomness or uncertainty in the hidden state, one could evaluate the entropy or normalized entropy¹³ of the discrete distribution corresponding to $\hat{\mathbf{x}}_{\bullet}$. Although this may seem related to the variance of the posterior in factor analysis, this analogy is deceptive. Since $\hat{\mathbf{x}}_{\bullet}$ defines the entire distribution, no other “variance” measure is needed. Learning consists of finding the cluster means (columns of \mathbf{C}), the covariance \mathbf{R} , and the mixing coefficients π_j . This is easily done with EM and corresponds exactly to maximum likelihood competitive learning (Duda & Hart, 1973; Nowlan, 1991), except

¹³ The entropy of the distribution divided by the logarithm of k so that it always lies between zero and one.

that all the clusters share the same covariance. Later we introduce extensions to the model that remove this restriction.

As in the continuous-state case, we can consider the limit as the observation noise becomes infinitesimal compared to the scale of the data. What results is the standard vector quantization model. The inference (classification) problem is now solved by the one-nearest-neighbor rule, using Euclidean distance if \mathbf{R} is a multiple of the identity matrix, or Mahalanobis distance in the unscaled matrix \mathbf{R} otherwise. Similarly to PCA, since the observation noise has disappeared, the posterior collapses to have all of its mass on one cluster (the closest), and the corresponding uncertainty (entropy) becomes zero. Learning with EM is equivalent to using a batch version of the k-means algorithm such as that proposed by Lloyd (1982). As with PCA, vector quantization does not define a proper density in the observation space. Once again, we examine the sum squared deviation of each point from its closest cluster center as a quantity proportional to the likelihood in the limit of zero noise. Batch k-means algorithms minimize this cost in lieu of maximizing a proper likelihood.

6.2 Time-Series Modeling: Hidden Markov Models. We return now to the fully dynamic discrete-state model introduced in equations 6.1a and 6.2b. Our key observation is that the dynamics described by equation 6.1a are exactly equivalent to the more traditional discrete Markov chain dynamics using a state transition matrix \mathbf{T} , where $\mathbf{T}_{ij} = P(\mathbf{x}_{t+1} = \mathbf{e}_j | \mathbf{x}_t = \mathbf{e}_i)$. It is easy to see how to compute the equivalent state transition matrix \mathbf{T} given \mathbf{A} and \mathbf{Q} above: \mathbf{T}_{ij} is the probability assigned by the gaussian whose mean is the i th column of \mathbf{A} (and whose covariance is \mathbf{Q}) to the region of k -space in which the j th coordinate is larger than all others. It is also true that for any transition matrix \mathbf{T} (whose rows each sum to unity), there exist matrices \mathbf{A} and \mathbf{Q} such that the dynamics are equivalent.¹⁴ Similarly, the initial probability mass function for \mathbf{x}_1 is easily computed from μ_1 and \mathbf{Q}_1 and for any desired histogram over the states for \mathbf{x}_1 there exist a μ_1 and \mathbf{Q}_1 that achieve it.

Similar degeneracy exists in this discrete-state model as in the continuous-state model except that it is now between the structure of \mathbf{A} and \mathbf{Q} . Since for any noise covariance \mathbf{Q} , the means in the columns of \mathbf{A} can be chosen to set any equivalent transition probabilities \mathbf{T}_{ij} , we can without loss of generality restrict \mathbf{Q} to be the identity matrix and use only the means in the columns

¹⁴ Although harder to see. Sketch of proof: Without loss of generality, always set the covariance to the identity matrix. Next, set the dot product of the mean vector with the k -vector having unity in all positions to be zero since moving along this direction does not change the probabilities. Now there are $(k - 1)$ degrees of freedom in the mean and also in the probability model. Set the mean randomly at first (except that it has no projection along the all-unity direction). Move the mean along a line defined by the constraint that all probabilities but two should remain constant until one of those two probabilities has the desired value. Repeat this until all have been set correctly.

of \mathbf{A} to set probabilities. Equivalently, we can restrict $\mathbf{Q}_1 = \mathbf{I}$ and use only the mean μ_1 to set the probabilities for the initial state \mathbf{x}_1 .

Thus, this generative model is equivalent to a standard HMM except that the emission probability densities are all constrained to have the same covariance. Likelihood and filtering computations are performed with the so-called forward (alpha) recursions, while complete smoothing is done with the forward-backward (alpha-beta) recursions. The EM algorithm for learning is exactly the well-known Baum-Welch reestimation procedure (Baum & Petrie, 1966; Baum et al., 1970).

There is an important and peculiar consequence of discretizing the state that affects the smoothing problem. The state sequence formed by taking the most probable state of the posterior distribution at each time (as computed by the forward-backward recursions given the observed data and model parameters) is not the single state sequence most likely to have produced the observed data. In fact, the sequence of states obtained by concatenating the states that individually have maximum posterior probability at each time step may have zero probability under the posterior. This creates the need for separate inference algorithms to find the single most likely state sequence given the observations. Such algorithms for filtering and smoothing are called Viterbi decoding methods (Viterbi, 1967). Why was there no need for similar decoding in the continuous-state case? It turns out that due to the smooth and unimodal nature of the posterior probabilities for individual states in the continuous case (all posteriors are gaussian), the sequence of maximum a posteriori states is exactly the single most likely state trajectory, so the regular Kalman filter and RTS smoothing recursions suffice. It is possible (see, for example, Rabiner & Juang, 1986) to learn the discrete-state model parameters based on the results of the Viterbi decoding instead of the forward-backward smoothing—in other words, to maximize the joint likelihood of the observations and the single most likely state sequence rather than the total likelihood summed over all possible paths through state-space.

7 Independent Component Analysis

There has been a great deal of recent interest in the blind source separation problem that attempts to recover a number of “source” signals from observations resulting from those signals, using only the knowledge that the original sources are independent. In the “square-linear” version of the problem, the observation process is characterized entirely by a square and invertible matrix \mathbf{C} . In other words, there are as many observation streams as sources, and there is no delay, echo, or convolutional distortion. Recent experience has shown the surprising result that for nongaussian distributed sources, this problem can often be solved even with no prior knowledge about the sources or about \mathbf{C} . It is widely believed (and beginning to be proved theo-

retically; see MacKay, 1996) that high kurtosis source distributions are most easily separated.

We will focus on a modified, but by now classic, version due to Bell and Sejnowski (1995) and Baram and Roth (1994) of the original independent component analysis algorithm (Comon, 1994). Although Bell and Sejnowski derived it from an information-maximization perspective, this modified algorithm can also be obtained by defining a particular prior distribution over the components of the vector \mathbf{x}_t of sources and then deriving a gradient learning rule that maximizes the likelihood of the data \mathbf{y}_t in the limit of zero output noise (Amari, Cichocki, & Yang, 1996; Pearlmutter & Parra, 1997; MacKay, 1996). The algorithm, originally derived for unordered data, has also been extended to modeling time series (Pearlmutter & Parra, 1997).

We now show that the generative model underlying ICA can be obtained by modifying slightly the basic model we have considered thus far. The modification is to replace the $\mathbf{WTA}[\cdot]$ nonlinearity introduced above with a general nonlinearity $g(\cdot)$ that operates componentwise on its input. Our generative model (for static data) then becomes:

$$\mathbf{x}_\bullet = g(\mathbf{w}_\bullet) \quad \mathbf{w}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (7.1a)$$

$$\mathbf{y}_\bullet = \mathbf{C}\mathbf{x}_\bullet + \mathbf{v}_\bullet \quad \mathbf{v}_\bullet \sim \mathcal{N}(\mathbf{0}, \mathbf{R}). \quad (7.1b)$$

The role of the nonlinearity is to convert the gaussian distributed prior for \mathbf{w}_\bullet into a nongaussian prior for \mathbf{x}_\bullet . Without loss of generality, we can set $\mathbf{Q} = \mathbf{I}$, since any covariance structure in \mathbf{Q} can be obtained by a linear transformation of a $\mathcal{N}(\mathbf{0}, \mathbf{I})$ random variable, and this linear transformation can be subsumed into the nonlinearity $g(\cdot)$. Assuming that the generative nonlinearity $g(\cdot)$ is invertible and differentiable, any choice of the generative nonlinearity results in a corresponding prior distribution on each source given by the probability density function:

$$p_{\mathbf{x}}(x) = \frac{\mathcal{N}(\mathbf{0}, \mathbf{1})|_{g^{-1}(x)}}{|g'(g^{-1}(x))|}. \quad (7.2)$$

It is important to distinguish this generative nonlinearity from the nonlinearity found in the ICA learning rule. We call this the learning rule nonlinearity, $f(\cdot)$, and clarify the distinction between the two nonlinearities below.

Classic ICA is defined for square and invertible \mathbf{C} in the limit of vanishing noise, $\mathbf{R} = \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{I}$. Under these conditions, the posterior density of \mathbf{x}_\bullet given \mathbf{y}_\bullet is a delta function at $\mathbf{x}_\bullet = \mathbf{C}^{-1}\mathbf{y}_\bullet$, and the ICA algorithm can be defined in terms of learning the recognition (or unmixing) weights $\mathbf{W} = \mathbf{C}^{-1}$, rather than the generative (mixing) weights \mathbf{C} . The gradient learning rule to increase the likelihood is

$$\Delta \mathbf{W} \propto \mathbf{W}^{-T} + f(\mathbf{W}\mathbf{y}_\bullet)\mathbf{y}_\bullet^T, \quad (7.3)$$

where the learning rule nonlinearity $f(\cdot)$ is the derivative of the implicit log prior: $f(x) = \frac{d \log p_x(x)}{dx}$ (MacKay, 1996). Therefore, any generative nonlinearity $g(\cdot)$ results in a nongaussian prior $p_x(\cdot)$, which in turn results in a nonlinearity $f(\cdot)$ in the maximum likelihood learning rule. Somewhat frustratingly from the generative models perspective, ICA is often discussed in terms of the learning rule nonlinearity without any reference to the implicit prior over the sources.

A popular choice for the ICA learning rule nonlinearity $f(\cdot)$ is the $\tanh(\cdot)$ function, which corresponds to a heavy tailed prior over the sources (MacKay, 1996):

$$p_x(x) = \frac{1}{\pi \cosh(x)}. \quad (7.4)$$

From equation 7.2, we obtain a general relationship between the cumulative distribution function of the prior on the sources, $\text{cdf}_x(x)$, and of the zero-mean, unit variance noise w ,

$$\text{cdf}_x(g(w)) = \text{cdf}_w(w) = \frac{1}{2} + \frac{1}{2} \text{erf}(w/\sqrt{2}), \quad (7.5)$$

for monotonic g , where $\text{erf}(z)$ is the error function $2/\sqrt{\pi} \int_0^z e^{-u^2} du$. This relationship can often be solved to obtain an expression for g . For example, if $p_x(x) = \frac{1}{\pi \cosh(x)}$, we find that setting

$$g(w) = \ln \left(\tan \left(\frac{\pi}{4} \left(1 + \text{erf}(w/\sqrt{2}) \right) \right) \right) \quad (7.6)$$

causes the generative model of equations 7.1 to generate vectors \mathbf{x} in which each component is distributed exactly according to $1/(\pi \cosh(x))$. This nonlinearity is shown in Figure 5.

ICA can be seen either as a linear generative model with nongaussian priors for the hidden variables or as a nonlinear generative model with gaussian priors for the hidden variables. It is therefore possible to derive an EM algorithm for ICA, even when the observation noise \mathbf{R} is nonzero and there are fewer sources than observations. The only complication is that the posterior distribution of \mathbf{x}_\bullet given \mathbf{y}_\bullet will be the product of a nongaussian prior and a gaussian likelihood term, which can be difficult to evaluate. Given this posterior, the M step then consists of maximizing the expected log of the joint probability as a function of \mathbf{C} and \mathbf{R} . The M-step for \mathbf{C} is

$$\mathbf{C} \leftarrow \arg \max_{\mathbf{C}} \sum_i \langle \log P(\mathbf{x}) + \log P(\mathbf{y}_i | \mathbf{x}, \mathbf{C}, \mathbf{R}) \rangle_i, \quad (7.7)$$

where i indexes the data points and $\langle \cdot \rangle_i$ denotes expectation with respect to the posterior distribution of \mathbf{x} given \mathbf{y}_i , $P(\mathbf{x} | \mathbf{y}_i, \mathbf{C}, \mathbf{R})$. The first term does not

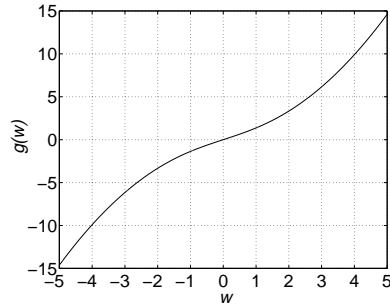


Figure 5: The nonlinearity $g(\cdot)$ from equation 7.6 which converts a gaussian distributed source $w \sim \mathcal{N}(0, 1)$ into one distributed as $x = g(w) \sim 1/(\pi \cosh(x))$.

depend on \mathbf{C} , and the second term is a quadratic in \mathbf{C} , so taking derivatives with respect to \mathbf{C} , we obtain a linear system of equations that can be solved in the usual manner:

$$\mathbf{C} \leftarrow \left(\sum_i \mathbf{y}_i \mathbf{x}_i^T \right) \left(\sum_i \langle \mathbf{x} \mathbf{x}^T \rangle_i \right)^{-1}. \quad (7.8)$$

A similar M-step can be derived for \mathbf{R} . Since, given \mathbf{x}_* , the generative model is linear, the M-step requires only evaluating the first and second moments of the posterior distribution of \mathbf{x} : $\langle \mathbf{x} \rangle_i$ and $\langle \mathbf{x} \mathbf{x}^T \rangle_i$. It is not necessary to know anything else about the posterior if its first two moments can be computed. These may be computed using Gibbs sampling or, for certain source priors, using table lookup or closed-form computation.¹⁵ In particular, Moulines et al. (1997) and Attias and Schreiner (1998) have independently proposed using a gaussian mixture to model the prior for each component of the source, \mathbf{x} . The posterior distribution over \mathbf{x} is then also a gaussian mixture, which can be evaluated analytically and used to derive an EM algorithm for both the mixing matrix and the source densities. The only caveat is that the number of gaussian components in the posterior grows exponentially in the number of sources,¹⁶ which limits the applicability of this method to models with only a few sources.

¹⁵ In the limit of zero noise, $\mathbf{R} = \mathbf{0}$, the EM updates derived in this manner degenerate to $\mathbf{C} \leftarrow \mathbf{C}$ and $\mathbf{R} \leftarrow \mathbf{R}$. Since this does not decrease the likelihood, it does not contradict the convergence proof for the EM algorithm. However, it also does not increase the likelihood, which might explain why no one uses EM to fit the standard zero-noise ICA model.

¹⁶ If each source is modeled as a mixture of k gaussians and there are m sources, then there are k^m components in the mixture.

Alternatively, we can compute the posterior distribution of \mathbf{w}_\bullet given \mathbf{y}_\bullet , which is the product of a gaussian prior and a nongaussian likelihood. Again, this may not be easy, and we may wish to resort to Gibbs sampling (Geman & Geman, 1984) or other Markov chain Monte Carlo methods (Neal, 1993). Another option is to employ a deterministic trick recently used by Bishop, Svenson, and Williams (1998) in the context of the generative topographic map (GTM), which is a probabilistic version of Kohonen's (1982) self-organized topographic map. We approximate the gaussian prior via a finite number (N) of *fixed* points (this is the trick). In other words,

$$P(\mathbf{w}_\bullet) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \approx \tilde{P}(\mathbf{w}_\bullet) = \sum_{j=1}^N \delta(\mathbf{w}_\bullet - \mathbf{w}_j), \quad (7.9)$$

where the \mathbf{w}_j 's are a finite sample from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The generative model then takes these N points, maps them through a fixed nonlinearity g , an adaptable linear mapping \mathbf{C} , and adds gaussian noise with covariance \mathbf{R} to produce \mathbf{y}_\bullet . The generative model is therefore a constrained mixture of N gaussians, where the constraint comes from the fact that the only way the centers can move is by varying \mathbf{C} . Then, computing the posterior over \mathbf{w}_\bullet amounts to computing the responsibility under each of the N gaussians for the data point. Given these responsibilities, the problem is again linear in \mathbf{C} , which means that it can be solved using equation 7.8. For the traditional zero-noise limit of ICA, the responsibilities will select the center closest to the data point in exactly the same manner as standard vector quantization. Therefore, ICA could potentially be implemented using EM for GTMs in the limit of zero output noise.

8 Network Interpretations and Regularization

Early in the modern history of neural networks, it was realized that PCA could be implemented using a linear autoencoder network (Baldi & Hornik, 1989). The data are fed as both the input and target of the network, and the network parameters are learned using the squared error cost function. In this section, we show how factor analysis and mixture of gaussian clusters can also be implemented in this manner, albeit with a different cost function.

To understand how a probabilistic model can be learned using an autoencoder it is very useful to make a recognition-generation decomposition of the autoencoder (Hinton & Zemel, 1994; Hinton, Dayan, & Revow, 1997). An autoencoder takes an input \mathbf{y}_\bullet , produces some internal representation in the hidden layer $\hat{\mathbf{x}}_\bullet$, and generates at its output a reconstruction of the input $\hat{\mathbf{y}}_\bullet$ in Figure 6. We call the mapping from hidden to output layers the generative part of the network since it generates the data from a (usually more compact) representation under some noise model. Conversely, we call the mapping from input to hidden units the recognition part of the network

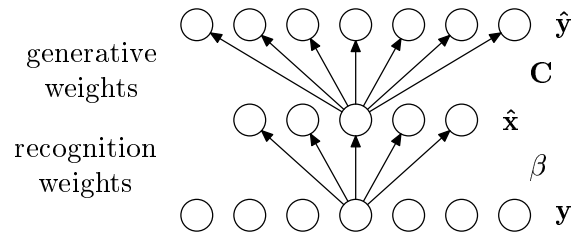


Figure 6: A network for state inference and for learning parameters of a static data model. The input \mathbf{y}_\bullet is clamped to the input units (bottom), and the mean $\hat{\mathbf{x}}_\bullet$ of the posterior of the estimated state appears on the hidden units above. The covariance of the state posterior is constant at $\mathbf{I} - \beta\mathbf{C}$ which is easily computed if the weights β are known. The inference computation is a trivial linear projection, but learning the weights of the inference network is difficult. The input to hidden weights are always constrained to be a function of the hidden to output weights, and the network is trained as an autoencoder using self-supervised learning. Outgoing weights have only been drawn from one input unit and one hidden unit.

since it produces some representation in the hidden variables given the input. Because autoencoders are usually assumed to be deterministic, we will think of the recognition network as computing the posterior mean of the hidden variables given the input.

The generative model for factor analysis assumes that both the hidden states and the observables are normally distributed, from which we get the posterior probabilities for the hidden states in equation 5.3b. If we assume that the generative weight matrix from the hidden units to the outputs is \mathbf{C} and the noise model at the output is gaussian with covariance \mathbf{R} , then the posterior mean of the hidden variables is $\hat{\mathbf{x}}_\bullet = \beta\mathbf{y}_\bullet$, where $\beta = \mathbf{C}^T(\mathbf{C}\mathbf{C}^T + \mathbf{R})^{-1}$. Therefore, the hidden units can compute the posterior mean exactly if they are linear and the weight matrix from input to hidden units is β . Notice that β is tied to \mathbf{C} and \mathbf{R} , so we only need to estimate \mathbf{C} and \mathbf{R} during learning. We denote expectations under the posterior state distribution by $\langle \cdot \rangle$, for example,

$$\langle \mathbf{x}_\bullet \rangle = \int \mathbf{x}_\bullet P(\mathbf{x}_\bullet | \mathbf{y}_\bullet) d\mathbf{x}_\bullet = \hat{\mathbf{x}}_\bullet.$$

From the theory of the EM algorithm (see section 4.2), we know that one way to maximize the likelihood is to maximize the expected value of the log of the joint probability under the posterior distribution of the hidden variables:

$$\langle \log P(\mathbf{x}_\bullet, \mathbf{y}_\bullet | \mathbf{C}, \mathbf{R}) \rangle.$$

Changing signs and ignoring constants, we can equivalently minimize the following cost function:

$$\mathcal{C} = \langle (\mathbf{y}_\bullet - \mathbf{C}\mathbf{x}_\bullet)^T \mathbf{R}^{-1} (\mathbf{y}_\bullet - \mathbf{C}\mathbf{x}_\bullet) \rangle + \log |\mathbf{R}| \quad (8.1a)$$

$$= \mathbf{y}_\bullet^T \mathbf{R}^{-1} \mathbf{y}_\bullet - 2\mathbf{y}_\bullet^T \mathbf{R}^{-1} \mathbf{C} \langle \mathbf{x}_\bullet \rangle + \langle \mathbf{x}_\bullet^T \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \mathbf{x}_\bullet \rangle + \log |\mathbf{R}| \quad (8.1b)$$

$$= (\mathbf{y}_\bullet - \mathbf{C}\hat{\mathbf{x}}_\bullet)^T \mathbf{R}^{-1} (\mathbf{y}_\bullet - \mathbf{C}\hat{\mathbf{x}}_\bullet) + \log |\mathbf{R}| + \text{trace}[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \Sigma]. \quad (8.1c)$$

Here we have defined Σ to be the posterior covariance of \mathbf{x}_\bullet ,

$$\Sigma \equiv \langle \mathbf{x}_\bullet \mathbf{x}_\bullet^T \rangle - \langle \mathbf{x}_\bullet \rangle \langle \mathbf{x}_\bullet \rangle^T = \mathbf{I} - \beta \mathbf{C},$$

and in the last step we have reorganized terms and made use of the fact that $\langle \mathbf{x}_\bullet^T \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \mathbf{x}_\bullet \rangle = \text{trace}[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \langle \mathbf{x}_\bullet \mathbf{x}_\bullet^T \rangle]$.

The first two terms of cost function in equation 8.1c are just a squared error cost function evaluated with respect to a gaussian noise model with covariance \mathbf{R} . They are exactly the terms minimized when fitting a standard neural network with this gaussian noise model. The last term is a regularization term that accounts for the posterior variance in the hidden states given the inputs.¹⁷ When we take derivatives of this cost function, we do not differentiate $\hat{\mathbf{x}}$ and Σ with respect to \mathbf{C} and \mathbf{R} . As is usual for the EM algorithm, we differentiate the cost given the posterior distribution of the hidden variables. Taking derivatives with respect to \mathbf{C} and premultiplying by \mathbf{R} , we obtain a weight change rule,

$$\Delta \mathbf{C} \propto (\mathbf{y}_\bullet - \mathbf{C}\mathbf{x}_\bullet) \mathbf{x}_\bullet^T - \mathbf{C} \Sigma. \quad (8.2)$$

The first term is the usual delta rule. The second term is simply a weight-decay term decaying the columns of \mathbf{C} with respect to the posterior covariance of the hidden variables. Intuitively, the higher the uncertainty in a hidden unit, the more its outgoing weight vector is shrunk toward zero. To summarize, factor analysis can be implemented in an autoassociator by tying the recognition weights to the generative weights and using a particular regularizer in addition to squared reconstruction error during learning.

We now analyze the mixture of gaussians model in the same manner. The recognition network is meant to produce the mean of the hidden variable given the inputs. Since we assume that the discrete hidden variable is represented as a unit vector, its mean is just the vector of probabilities of being in each of its k settings given the inputs, that is, the responsibilities. Assuming equal mixing coefficients, $P(\mathbf{x}_\bullet = \mathbf{e}_j) = P(\mathbf{x}_\bullet = \mathbf{e}_i) \forall ij$, the responsibilities

¹⁷ PCA assumes infinitesimal noise, and therefore the posterior “distribution” over the hidden states has zero variance ($\Sigma \rightarrow \mathbf{0}$) and the regularizer vanishes ($\mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \Sigma \rightarrow \mathbf{I}$).

defined in equation 6.6b are

$$(\hat{\mathbf{x}}_{\bullet})_j = P(\mathbf{x}_{\bullet} = \mathbf{e}_j | \mathbf{y}_{\bullet}) = \frac{\exp\{-\frac{1}{2}(\mathbf{y}_{\bullet} - \mathbf{C}_j)^T \mathbf{R}^{-1}(\mathbf{y}_{\bullet} - \mathbf{C}_j)\}}{\sum_{i=1}^k \exp\{-\frac{1}{2}(\mathbf{y}_{\bullet} - \mathbf{C}_i)^T \mathbf{R}^{-1}(\mathbf{y}_{\bullet} - \mathbf{C}_i)\}} \quad (8.3a)$$

$$= \frac{\exp\{\beta_j \mathbf{y}_{\bullet} - \alpha_j\}}{\sum_{i=1}^k \exp\{\beta_i \mathbf{y}_{\bullet} - \alpha_i\}}, \quad (8.3b)$$

where we have defined $\beta_j = \mathbf{C}_j \mathbf{R}^{-1}$ and $\alpha_j = \frac{1}{2} \mathbf{C}_j^T \mathbf{R}^{-1} \mathbf{C}_j$. Equation 8.3b describes a recognition model that is linear followed by the softmax nonlinearity, Φ , written in full matrix notation: $\hat{\mathbf{x}}_{\bullet} = \Phi(\beta \mathbf{y}_{\bullet} - \alpha)$. In other words, a simple network could do exact inference with linear input to hidden weights β and softmax hidden units.

Appealing again to the EM algorithm, we obtain a cost function that when minimized by an autoassociator will implement the mixture of gaussians.¹⁸ The log probability of the data given the hidden variables can be written as

$$-2 \log P(\mathbf{y}_{\bullet} | \mathbf{x}_{\bullet}, \mathbf{C}, \mathbf{R}) = (\mathbf{y}_{\bullet} - \mathbf{C} \mathbf{x}_{\bullet})^T \mathbf{R}^{-1} (\mathbf{y}_{\bullet} - \mathbf{C} \mathbf{x}_{\bullet}) + \log |\mathbf{R}| + \text{const.}$$

Using this and the previous derivation, we obtain the cost function,

$$\mathcal{C} = (\mathbf{y}_{\bullet} - \mathbf{C} \hat{\mathbf{x}}_{\bullet})^T \mathbf{R}^{-1} (\mathbf{y}_{\bullet} - \mathbf{C} \hat{\mathbf{x}}_{\bullet}) + \log |\mathbf{R}| + \text{trace}[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \Sigma], \quad (8.4)$$

where $\Sigma = \langle \mathbf{x}_{\bullet} \mathbf{x}_{\bullet}^T \rangle - \langle \mathbf{x}_{\bullet} \rangle \langle \mathbf{x}_{\bullet} \rangle^T$. The second-order term, $\langle \mathbf{x}_{\bullet} \mathbf{x}_{\bullet}^T \rangle$, evaluates to a matrix with $\hat{\mathbf{x}}_{\bullet}$ along its diagonal and zero elsewhere. Unlike in factor analysis, Σ now depends on the input.

To summarize, the mixture of gaussians model can also be implemented using an autoassociator. The recognition part of the network is linear, followed by a softmax nonlinearity. The cost function is the usual squared error penalized by a regularizer of exactly the same form as in factor analysis. Similar network interpretations can be obtained for the other probabilistic models.

9 Comments and Extensions

There are several advantages, both theoretical and practical, to a unified treatment of the unsupervised methods reviewed here. From a theoretical viewpoint, the treatment emphasizes that all of the techniques for inference in the various models are essentially the same and just correspond to probability propagation in the particular model variation. Similarly, all the learning procedures are nothing more than an application of the EM

¹⁸ Our derivation assumes tied covariance and equal mixing coefficients. Slightly more complex equations result for the general case.

algorithm to increase the total likelihood of the observed data iteratively. Furthermore, the origin of zero-noise-limit algorithms such as vector quantization and PCA is easy to see. A unified treatment also highlights the relationship between similar questions across the different models. For example, picking the number of clusters in a mixture model or state dimension in a dynamical system or the number of factors or principal components in a covariance model or the number of states in an HMM are all really the same question.

From a practical standpoint, a unified view of these models allows us to apply well-known solutions to hard problems in one area to similar problems in another. For example, in this framework it is obvious how to deal properly with missing data in solving both the learning and inference problems. This topic has been well understood for many static models (Little & Rubin, 1987; Tresp, Ahmad, & Neuneier, 1994; Ghahramani & Jordan, 1994) but is typically not well addressed in the linear dynamical systems literature. As another example, it is easy to design and work with models having a mixed continuous- and discrete-state vector, (for example, hidden filter HMMs (Fraser & Dimitriadis, 1993), which is something not directly addressed by the individual literatures on discrete or continuous models.

Another practical advantage is the ease with which natural extensions to the basic models can be developed. For example, using the hierarchical mixture-of-experts formalism developed by Jordan and Jacobs (1994) we can consider global mixtures of any of the model variants discussed. In fact, most of these mixtures have already been considered: mixtures of linear dynamical systems are known as switching state-space models (see Shumway & Stoffer, 1991; Ghahramani & Hinton, 1996b); mixtures of factor analyzers (Ghahramani and Hinton, 1997) and of pancakes (PCA) (Hinton et al., 1995); and mixtures of HMMs (Smyth, 1997). A mixture of m of our constrained mixtures of gaussians each with k clusters gives a mixture model with mk components in which there are only m possible covariance matrices. This “tied covariance” approach is popular in speech modeling to reduce the number of free parameters. (For $k = 1$, this corresponds to a full “unconstrained” mixture of gaussians model with m clusters.)

It is also possible to consider “local mixtures” in which the conditional probability $P(\mathbf{y}_t|\mathbf{x}_t)$ is no longer a single gaussian but a more complicated density such as a mixture of gaussians. For our (constrained) mixture of gaussians model, this is another way to get a “full” mixture. For HMMs, this is a well-known extension and is usually the standard approach for emission density modeling (Rabiner & Juang, 1986). It is even possible to use constrained mixture models as the output density model for an HMM (see, for example, Saul & Rahim, 1998, which uses factor analysis as the HMM output density). However, we are not aware of any work that considers this variation in the continuous-state cases, for either static or dynamic data.

Another important natural extension is spatially adaptive observation noise. The idea here is that the observation noise \mathbf{v} can have different statis-

tics in different parts of (state or observation) space rather than being described by a single matrix \mathbf{R} . For discrete-mixture models, this idea is well known, and it is achieved by giving each mixture component a private noise model. However, for continuous-state models, this idea is relatively unexplored and is an interesting area for further investigation. The crux of the problem is how to parameterize a positive definite matrix over some space. We propose some simple ways to achieve this. One possibility is replacing the single covariance shape \mathbf{Q} for the observation noise with a conic¹⁹ linear blending of k “basis” covariance shapes. In the case of linear dynamical systems or factor analysis, this amounts to a novel type of model in which the local covariance matrix \mathbf{R} is computed as a conic linear combination of several “canonical” covariance matrices through a tensor product between the current state vector \mathbf{x} (or equivalently the “noiseless” observation $\mathbf{C}\mathbf{x}$) and a master noise tensor \mathcal{R} .²⁰ Another approach would be to drop the conic restriction (allow general linear combinations) and then add a multiple of the identity matrix to the resulting noise matrix in order to make it positive definite. A third approach is to represent the covariance shape as the compression of an elastic sphere due to a spatially varying force field. This representation is easier to work with because the parameterization of the field is unconstrained, but it is hard to learn the local field from measurements of the effective covariance shape. Bishop (1995, sec. 6.3) and others have considered simple nonparametric methods for estimating input-dependent noise levels in regression problems. Goldberg, Williams, and Bishop (1998) have explored this idea in the context of gaussian processes.

It is also interesting to consider what happens to the dynamic models when the output noise tends to zero. In other words, what are the dynamic analogs of PCA and vector quantization? For both linear dynamical systems and HMMs, this causes the state to no longer be hidden. In linear dynamical systems, the optimal observation matrix is then found by performing PCA on the data and using the principal components as the columns of \mathbf{C} ; for HMMs, \mathbf{C} is found by vector quantization of the data (using the codebook vectors as the columns of \mathbf{C}). Given these observation matrices, the state is no longer hidden. All that remains is to identify a first-order Markov dynamics in state-space: this is a simple AR(1) model in the continuous case or a first-order Markov chain in the discrete case. Such zero-noise limits are not only interesting models in their own right, but are also valuable as good choices for initialization of learning in linear dynamical systems and HMMs.

¹⁹ A conic linear combination is one in which all the coefficients are positive.

²⁰ For mixtures of gaussians or hidden Markov models, this kind of linear “blending” merely selects the j th submatrix of the tensor if the discrete state is e_j . This is yet another way to recover the conventional “full” or unconstrained mixture of gaussians or hidden Markov model emission density in which each cluster or state has its own private covariance shape for observation noise.

Appendix

In this appendix we review in detail the inference and learning algorithms for each of the models. The goal is to enable readers to implement these algorithms from the pseudocode provided. For each class of model, we first present the solution to the inference problem, and then the EM algorithm for learning the model parameters. For this appendix only, we adopt the notation that the transpose of a vector or matrix is written as \mathbf{x}' , not \mathbf{x}^T . We use T instead of τ to denote the length of a time series. We also define the binary operator \odot to be the element-wise product of two equal-size vectors or matrices. Comments begin with the symbol %.

A.1 Factor Analysis, SPCA, and PCA.

A.1.1 Inference. For factor analysis and related models, the posterior probability of the hidden state given the observations, $P(\mathbf{x}_\bullet | \mathbf{y}_\bullet)$, is gaussian. The inference problem therefore consists of computing the mean and covariance of this gaussian, $\hat{\mathbf{x}}_\bullet$ and $\mathbf{V} = \text{Cov}[\mathbf{x}_\bullet]$:

```
FactorAnalysisInference( $\mathbf{y}_\bullet, \mathbf{C}, \mathbf{R}$ )
   $\beta \leftarrow \mathbf{C}'(\mathbf{C}\mathbf{C}' + \mathbf{R})^{-1}$ 
   $\hat{\mathbf{x}}_\bullet \leftarrow \beta\mathbf{y}_\bullet$ 
   $\mathbf{V} \leftarrow \mathbf{I} - \beta\mathbf{C}$ 
  return  $\hat{\mathbf{x}}_\bullet, \mathbf{V}$ 
```

Since the observation noise matrix \mathbf{R} is assumed to be diagonal and \mathbf{x}_\bullet is of smaller dimension than \mathbf{y}_\bullet , β can be more efficiently computed using the matrix inversion lemma:

$$\beta = \mathbf{C}'\mathbf{R}^{-1} \left(\mathbf{I} - \mathbf{C}(\mathbf{I} + \mathbf{C}'\mathbf{R}^{-1}\mathbf{C})^{-1}\mathbf{C}'\mathbf{R}^{-1} \right).$$

Computing the (log) likelihood of an observation is nothing more than an evaluation under the gaussian $\mathcal{N}(\mathbf{0}, \mathbf{C}\mathbf{C}' + \mathbf{R})$.

The sensible PCA (SPCA) algorithm is a special case of factor analysis in which the observation noise is assumed to be spherically symmetric: $\mathbf{R} = \alpha\mathbf{I}$. Inference in SPCA is therefore identical to inference for factor analysis.

The traditional PCA algorithm can be obtained as a limiting case of factor analysis: $\mathbf{R} = \lim_{\epsilon \rightarrow 0} \epsilon\mathbf{I}$. The inverse used for computing β in factor analysis is no longer well defined. However, the limit of β is well defined: $\beta = (\mathbf{C}'\mathbf{C})^{-1}\mathbf{C}'$. Also, the posterior collapses to a single point, so $\mathbf{V} = \text{Cov}[\mathbf{x}_\bullet] = \mathbf{I} - \beta\mathbf{C} = \mathbf{0}$.

```
PCAINference( $\mathbf{y}_\bullet, \mathbf{C}$ ) % Projection onto principal components
   $\beta \leftarrow (\mathbf{C}'\mathbf{C})^{-1}\mathbf{C}'$ 
   $\hat{\mathbf{x}}_\bullet \leftarrow \beta\mathbf{y}_\bullet$ 
  return  $\hat{\mathbf{x}}_\bullet$ 
```

A.1.2 Learning. The EM algorithm for learning the parameters of a factor analyzer with k factors from a zero-mean data set $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$ (each column of the $p \times n$ matrix \mathbf{Y} is a data point) is

```

FactorAnalysisLearn( $\mathbf{Y}, k, \epsilon$ )
  initialize  $\mathbf{C}, \mathbf{R}$ 
  compute sample covariance  $\mathbf{S}$  of  $\mathbf{Y}$ 
  while change in log likelihood  $> \epsilon$ 
    % E step
     $\hat{\mathbf{X}}, \mathbf{V} \leftarrow \text{FactorAnalysisInference}(\mathbf{Y}, \mathbf{C}, \mathbf{R})$ 
     $\delta \leftarrow \mathbf{Y}\hat{\mathbf{X}}'$ 
     $\gamma \leftarrow \hat{\mathbf{X}}\hat{\mathbf{X}}' + n\mathbf{V}$ 
    % M step
     $\mathbf{C} \leftarrow \delta\gamma^{-1}$ 
    set diagonal elements of  $\mathbf{R}$  to  $\mathbf{R}_{ii} \leftarrow (\mathbf{S} - \mathbf{C}\delta'/n)_{ii}$ 
  end
  return  $\mathbf{C}, \mathbf{R}$ 

```

Here `FactorAnalysisInference($\mathbf{Y}, \mathbf{C}, \mathbf{R}$)` has the obvious interpretation of the inference function applied to the entire matrix of observations. Since β and \mathbf{V} do not depend on \mathbf{Y} , this can be computed efficiently in matrix form. Since the data appear only in outer products, we can run factor analysis learning with just the sample covariance. Note also that the log-likelihood is computed as $-\frac{1}{2}\mathbf{y}'(\mathbf{C}\mathbf{C}' + \mathbf{R})^{-1}\mathbf{y} + \frac{n}{2}\log|\mathbf{C}\mathbf{C}' + \mathbf{R}| + \text{const}$.

The EM algorithm for SPCA is identical to the EM algorithm for factor analysis, except that since the observation noise covariance is spherically symmetrical, the M-step for \mathbf{R} is changed to $\mathbf{R} \leftarrow \alpha\mathbf{I}$, where $\alpha \leftarrow \sum_{j=1}^p (\mathbf{S} - \mathbf{C}\delta')_{jj}/p$.

The EM algorithm for PCA can be obtained in a similar manner:

```

PCALearn( $\mathbf{Y}, k, \epsilon$ )
  initialize  $\mathbf{C}$ 
  while change in squared reconstruction error  $> \epsilon$ 
    % E step
     $\hat{\mathbf{X}} \leftarrow \text{PCAIInference}(\mathbf{Y}, \mathbf{C})$ 
     $\delta \leftarrow \mathbf{Y}\hat{\mathbf{X}}'$ 
     $\gamma \leftarrow \hat{\mathbf{X}}\hat{\mathbf{X}}'$ 
    % M step
     $\mathbf{C} \leftarrow \delta\gamma^{-1}$ 
  end
  return  $\mathbf{C}$ 

```

Since PCA is not a probability model (it assumes zero noise), the likelihood is undefined, so convergence is assessed by monitoring the squared reconstruction error.

A.2 Mixtures of Gaussians and Vector Quantization.

A.2.1 Inference. We begin by discussing the inference problem for mixtures of gaussians and then discuss the inference problem in vector quantization as a limiting case. The hidden variable in a mixture of gaussians is a discrete variable that can take on one of k values. We represent this variable using a vector \mathbf{x}_\bullet of length k , where each setting of the hidden variable corresponds to \mathbf{x}_\bullet taking a value of unity in one dimension and zero elsewhere. The probability distribution of the discrete hidden variable, which has $k - 1$ degrees of freedom (since it must sum to one), is fully described by the mean of \mathbf{x}_\bullet . Therefore, the inference problem is limited to computing the posterior mean of \mathbf{x}_\bullet given a data point \mathbf{y}_\bullet and the model parameters, which are π (the prior mean of \mathbf{x}_\bullet), \mathbf{C} (the matrix whose k columns are the means of \mathbf{y}_\bullet given each of the k settings of \mathbf{x}_\bullet) and \mathbf{R} (the observation noise covariance matrix).

```
MixtureofGaussiansInference( $\mathbf{y}_\bullet, \mathbf{C}, \mathbf{R}, \pi$ ) % compute
                                         % responsibilities
 $\alpha \leftarrow 0$ 
for  $i = 1$  to  $k$ 
     $\Delta_i \leftarrow (\mathbf{y}_\bullet - \mathbf{C}_i)' \mathbf{R}^{-1} (\mathbf{y}_\bullet - \mathbf{C}_i)$ 
     $\gamma_i \leftarrow \pi_i \exp \{-\frac{1}{2} \Delta_i\}$ 
     $\alpha \leftarrow \alpha + \gamma_i$ 
end
 $\hat{\mathbf{x}}_\bullet \leftarrow \gamma / \alpha$ 
return  $\hat{\mathbf{x}}_\bullet$ 
```

A measure of the randomness of the hidden state can be obtained by evaluating the entropy of the discrete distribution corresponding to $\hat{\mathbf{x}}_\bullet$.

Standard VQ corresponds to the limiting case $\mathbf{R} = \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{I}$ and equal priors $\pi_i = 1/k$. Inference in this case is performed by the well known (1-)nearest-neighbor rule.

```
VQInference( $\mathbf{y}_\bullet, \mathbf{C}$ ) % 1-nearest-neighbor
for  $i = 1$  to  $k$ 
     $\Delta_i \leftarrow (\mathbf{y}_\bullet - \mathbf{C}_i)' (\mathbf{y}_\bullet - \mathbf{C}_i)$ 
end
 $\hat{\mathbf{x}}_\bullet \leftarrow \mathbf{e}_j$  for  $j = \arg \min \Delta_i$ 
return  $\hat{\mathbf{x}}_\bullet$ 
```


As before, \mathbf{e}_j is the unit vector along the j th coordinate direction. The squared distances Δ_i can be generalized to a Mahalanobis metric with respect to some matrix \mathbf{R} , and nonuniform priors π_i can easily be incorporated. As was the case with PCA, the posterior distribution has zero entropy.

A.2.2 Learning. The EM algorithm for learning the parameters of a mixture of gaussian is:

```
MixtureofGaussiansLearn( $\mathbf{Y}, k, \epsilon$ ) % ML soft competitive learning
  initialize  $\mathbf{C}, \mathbf{R}, \boldsymbol{\pi}$ 
  while change in log likelihood >  $\epsilon$ 
    initialize  $\boldsymbol{\delta} \leftarrow \mathbf{0}, \boldsymbol{\gamma} \leftarrow \mathbf{0}, \boldsymbol{\alpha} \leftarrow \mathbf{0}$ 
    % E step
    for  $i = 1$  to  $n$ 
       $\mathbf{x}_i \leftarrow$  MixtureofGaussiansInference( $\mathbf{y}_i, \mathbf{C}, \mathbf{R}, \boldsymbol{\pi}$ )
       $\boldsymbol{\delta} \leftarrow \boldsymbol{\delta} + \mathbf{y}_i \mathbf{x}_i'$ 
       $\boldsymbol{\gamma} \leftarrow \boldsymbol{\gamma} + \mathbf{x}_i$ 
    end
    % M step
    for  $j = 1$  to  $k$ 
       $\mathbf{C}_j \leftarrow \boldsymbol{\delta}_j / \boldsymbol{\gamma}_j$ 
      for  $i = 1$  to  $n$ 
         $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \mathbf{x}_{ij}(\mathbf{y}_i - \mathbf{C}_j)(\mathbf{y}_i - \mathbf{C}_j)'$ 
      end
    end
    end
     $\mathbf{R} \leftarrow \boldsymbol{\alpha} / n$ 
     $\boldsymbol{\pi} \leftarrow \boldsymbol{\gamma} / n$ 
  end
  return  $\mathbf{C}, \mathbf{R}, \boldsymbol{\pi}$ 
```

We have assumed a common covariance matrix \mathbf{R} for all the gaussians; the extension to different covariances for each gaussian is straightforward.

The k -means vector quantization learning algorithm results when we take the appropriate limit of the above algorithm:

```
VQLearn( $\mathbf{Y}, k, \epsilon$ ) % k-means
  initialize  $\mathbf{C}$ 
  while change in squared reconstruction error >  $\epsilon$ 
    % E step
    initialize  $\boldsymbol{\delta} \leftarrow \mathbf{0}, \boldsymbol{\gamma} \leftarrow \mathbf{0}$ 
    for  $i = 1$  to  $n$ 
       $\mathbf{x}_i \leftarrow$  VQInference( $\mathbf{y}_i, \mathbf{C}$ )
       $\boldsymbol{\delta} \leftarrow \boldsymbol{\delta} + \mathbf{y}_i \mathbf{x}_i'$ 
       $\boldsymbol{\gamma} \leftarrow \boldsymbol{\gamma} + \mathbf{x}_i$ 
    end
```

```

end
% M step
for j = 1 to k
    Cj ← δj/γj
end
end
return C

```

A.3 Linear Dynamical Systems.

A.3.1 Inference. Inference in a linear dynamical system involves computing the posterior distributions of the hidden state variables given the sequence of observations. As in factor analysis, all the hidden state variables are assumed gaussian and are therefore fully described by their means and covariance matrices. The algorithm for computing the posterior means and covariances consists of two parts: a forward recursion, which uses the observations from \mathbf{y}_1 to \mathbf{y}_t , known as the Kalman filter (Kalman, 1960), and a backward recursion, which uses the observations from \mathbf{y}_T to \mathbf{y}_{t+1} (Rauch, 1963). The combined forward and backward recursions are known as the Kalman or Rauch-Tung-Streibel (RTS) smoother.

To describe the smoothing algorithm it will be useful to define the following quantities: \mathbf{x}_t^s and \mathbf{V}_t^s are, respectively, the mean and covariance matrix of \mathbf{x}_t given observations $\{\mathbf{y}_1, \dots, \mathbf{y}_s\}$; $\hat{\mathbf{x}}_t \equiv \mathbf{x}_t^T$ and $\hat{\mathbf{V}}_t \equiv \mathbf{V}_t^T$ are the “full smoother” estimates. To learn the \mathbf{A} matrix using EM, it is also necessary to compute the covariance across time between \mathbf{x}_t and \mathbf{x}_{t-1} :

```

LDSInference(Y, A, C, Q, R, x10, V10) % Kalman smoother
for t = 1 to T % Kalman filter (forward pass)
    xtt-1 ← Axt-1t-1 if t > 1
    Vtt-1 ← AVt-1t-1A' + Q if t > 1
    Kt ← Vtt-1C'(CVtt-1C' + R)-1
    xtt ← xtt-1 + Kt(yt - Cxtt-1)
    Vtt ← Vtt-1 - KtCVtt-1
end
initialize VT,T-1hat = (I - KTC)AVT-1T-1
for t = T to 2 % Rauch recursions (backward pass)
    Jt-1 ← Vt-1t-1A'(Vtt-1)-1
    xt-1hat ← xt-1t-1 + Jt-1(xthat - Axt-1t-1)
    Vt-1hat ← Vt-1t-1 + Jt-1(Vthat - Vtt-1)Jt-1'
    Vt,t-1hat ← Vt-1tJt-1' + Jt(Vt+1,that - AVtt)Jt-1' if t < T

```

```

end
return  $\hat{\mathbf{x}}_t, \hat{\mathbf{V}}_t, \hat{\mathbf{V}}_{t,t-1}$  for all  $t$ 

```

A.3.2 Learning. The EM algorithm for learning a linear dynamical system (Shumway & Stoffer, 1982; Ghahramani & Hinton, 1996a) is given below, assuming for simplicity that we have only a single sequence of observations:

```

LDSLearn( $\mathbf{Y}, k, \epsilon$ )
  initialize  $\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}, \mathbf{x}_1^0, \mathbf{V}_1^0$ 
  set  $\alpha \leftarrow \sum_t \mathbf{y}_t \mathbf{y}_t'$ 
  while change in log likelihood  $> \epsilon$ 
    LDSInference( $\mathbf{Y}, \mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}, \mathbf{x}_1^0, \mathbf{V}_1^0$ ) % E step
    initialize  $\delta \leftarrow \mathbf{0}, \gamma \leftarrow \mathbf{0}, \beta \leftarrow \mathbf{0}$ 
    for  $t=1$  to  $T$ 
       $\delta \leftarrow \delta + \mathbf{y}_t \hat{\mathbf{x}}_t'$ 
       $\gamma \leftarrow \gamma + \hat{\mathbf{x}}_t \hat{\mathbf{x}}_t' + \hat{\mathbf{V}}_t$ 
       $\beta \leftarrow \beta + \hat{\mathbf{x}}_t \hat{\mathbf{x}}_{t-1}' + \hat{\mathbf{V}}_{t,t-1}$  if  $t > 1$ 
    end
     $\gamma_1 \leftarrow \gamma - \hat{\mathbf{x}}_T \hat{\mathbf{x}}_T' - \hat{\mathbf{V}}_T$ 
     $\gamma_2 \leftarrow \gamma - \hat{\mathbf{x}}_1 \hat{\mathbf{x}}_1' - \hat{\mathbf{V}}_1$ 
    % M step
     $\mathbf{C} \leftarrow \delta \gamma^{-1}$ 
     $\mathbf{R} \leftarrow (\alpha - \mathbf{C} \delta') / T$ 
     $\mathbf{A} \leftarrow \beta \gamma_1^{-1}$ 
     $\mathbf{Q} \leftarrow (\gamma_2 - \mathbf{A} \beta') / (T - 1)$ 
     $\mathbf{x}_1^0 \leftarrow \hat{\mathbf{x}}_1$ 
     $\mathbf{V}_1^0 \leftarrow \hat{\mathbf{V}}_1$ 
  end
  return  $\mathbf{A}, \mathbf{C}, \mathbf{Q}, \mathbf{R}, \mathbf{x}_1^0, \mathbf{V}_1^0$ 

```

A.4 Hidden Markov Models.

A.4.1 Inference. The forward-backward algorithm computes the posterior probabilities of the hidden states in an HMM and therefore forms the basis of the inference required for EM. We use the following standard definitions,

$$\alpha_t = P(\mathbf{x}_t, \mathbf{y}_1, \dots, \mathbf{y}_t) \quad (\text{A.1a})$$

$$\beta_t = P(\mathbf{y}_{t+1}, \dots, \mathbf{y}_T | \mathbf{x}_t), \quad (\text{A.1b})$$

where both α and β are vectors of the same length as \mathbf{x} . We present the case where the observations \mathbf{y}_t are real-valued p -dimensional vectors and the probability density of an observation given the corresponding state (the “output model”) is assumed to be a single gaussian with mean $\mathbf{C}\mathbf{x}_t$ and covariance \mathbf{R} . The parameters of the model are therefore a $k \times k$ transition matrix \mathbf{T} , initial state prior probability vector π , an observation mean matrix \mathbf{C} , and an observation noise matrix \mathbf{R} that is tied across states:

```
HMMInference(Y, T,  $\pi$ , C, R) % forward-backward algorithm
for t=1 to T % forward pass
  for i=1 to k
     $\mathbf{b}_{t,i} \leftarrow \exp \left\{ -\frac{1}{2}(\mathbf{y}_t - \mathbf{C}_i)' \mathbf{R}^{-1}(\mathbf{y}_t - \mathbf{C}_i) \right\} |\mathbf{R}|^{-1/2} (2\pi)^{-p/2}$ 
  end
  if (t=1)  $\alpha_t \leftarrow \pi \odot \mathbf{b}_1$  else  $\alpha_t \leftarrow [\mathbf{T}' \alpha_{t-1}] \odot \mathbf{b}_t$  end
   $\rho_t \leftarrow \sum_i \alpha_{t,i}$ 
   $\alpha_t \leftarrow \alpha_t / \rho_t$ 
end
 $\beta_T \leftarrow \mathbf{1} / \rho_T$ 
for t=T-1 to 1 % backward pass
   $\beta_t \leftarrow \mathbf{T} [\beta_{t+1} \odot \mathbf{b}_{t+1}] / \rho_t$ 
end
 $\gamma_t \leftarrow \alpha_t \odot \beta_t / (\alpha_t' \beta_t)$ 
 $\xi_t \leftarrow \alpha_t [\beta_{t+1} \odot \mathbf{b}_{t+1}]'$ 
 $\xi_t \leftarrow \xi_t / (\sum_{ij} \xi_{t,ij})$ 
return  $\gamma_t, \xi_t, \rho_t$  for all t
```

The definitions of α and β in equations A.1a and A.1b correspond to running the above algorithm without the scaling factors ρ_t . These factors, however, are essential to the numerical stability of the algorithm; otherwise, for long sequences, both α and β become vanishingly small. Furthermore, from the ρ 's we can compute the log-likelihood of the sequence

$$\log P(\mathbf{y}_1, \dots, \mathbf{y}_T) = \sum_{t=1}^T \log \rho_t,$$

which is why it is useful for the above function to return them.

A.4.2 Learning. Again, we assume for simplicity that we have a single sequence of observations from which we wish to learn the parameters of an HMM. The EM algorithm for learning these parameters, known as the

Baum-Welch algorithm, is:

```

HMMLearn( $\mathbf{Y}, k, \epsilon$ ) % Baum-Welch
  initialize  $\mathbf{T}, \boldsymbol{\pi}, \mathbf{C}, \mathbf{R}$ 
  while change in log likelihood >  $\epsilon$ 
    HMMInference( $\mathbf{Y}, \mathbf{T}, \boldsymbol{\pi}, \mathbf{C}, \mathbf{R}$ ) % E step
    % M step
     $\boldsymbol{\pi} \leftarrow \boldsymbol{\gamma}_1$ 
     $\mathbf{T} \leftarrow \sum_t \boldsymbol{\xi}_t$ 
     $\mathbf{T}_{ij} \leftarrow \mathbf{T}_{ij} / \sum_\ell \mathbf{T}_{i\ell}$  for all  $i, j$ 
     $\mathbf{C}_j \leftarrow \sum_t \gamma_{t,j} \mathbf{y}_t / \sum_t \gamma_{t,j}$  for all  $j$ 
     $\mathbf{R} \leftarrow \sum_{t,j} \gamma_{t,j} (\mathbf{y}_t - \mathbf{C}_j)(\mathbf{y}_t - \mathbf{C}_j)' / T$ 
  return  $\mathbf{T}, \boldsymbol{\pi}, \mathbf{C}, \mathbf{R}$ 

```

Acknowledgments

We thank Carlos Brody, Sanjoy Mahajan, and Erik Winfree for many fruitful discussions in the early stages, the anonymous referees for helpful comments, and Geoffrey Hinton and John Hopfield for providing outstanding intellectual environments and guidance. S.R. was supported in part by the Center for Neuromorphic Systems Engineering as a part of the National Science Foundation Engineering Research Center Program under grant EEC-9402726 and by the Natural Sciences and Engineering Research Council of Canada under an NSERC 1967 Award. Z.G. was supported by the Ontario Information Technology Research Centre.

References

- Amari, S., Cichocki, A., & Yang, H. H. (1996). A new learning algorithm for blind signal separation. In D. S. Touretzky, M. C. Mozer, & M. E. Hasselmo (Eds.), *Advances in neural information processing systems*, 8 (pp. 757–763). Cambridge, MA: MIT Press.
- Attias, H., & Schreiner, C. (1998). Blind source separation and deconvolution: The dynamic component analysis algorithm. *Neural Computation*, 10, 1373–1424.
- Baldi, P., & Hornik, K. (1989). Neural networks and principal components analysis: Learning from examples without local minima. *Neural Networks*, 2, 53–58.
- Baram, Y., & Roth, Z. (1994). *Density shaping by neural networks with application to classification, estimation and forecasting* (Tech. Rep. TR-CIS-94-20). Haifa, Israel: Center for Intelligent Systems, Technion, Israel Institute for Technology.
- Bauer, E., Koller, D., & Singer, Y. (1997). Update rules for parameter estimation in Bayesian networks. In *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence (UAI-97)*. San Mateo, CA: Morgan Kaufmann.
- Baum, L. E. (1972). An inequality and associated maximization technique in sta-

- tistical estimation of probabilistic functions of a Markov process. *Inequalities*, 3, 1–8.
- Baum, L. E., & Eagon, J. A. (1967). An inequality with applications to statistical estimation for probabilistic functions of Markov processes and to a model for ecology. *Bulletin of the American Mathematical Society*, 73, 360–363.
- Baum, L. E., & Petrie, T. (1966). Statistical inference for probabilistic functions of finite state Markov chains. *Annals of Mathematical Statistics*, 37, 1554–1563.
- Baum, L. E., Petrie, T., Soulds, G., & Weiss, N. (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Annals of Mathematical Statistics*, 41(1), 164–171.
- Bell, A. J., & Sejnowski, T. J. (1995). An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6), 1129–1159.
- Bishop, C. (1995). *Neural networks for pattern recognition*. Oxford: Clarendon Press.
- Bishop, C. M., Svensen, M., & Williams, C. K. I. (1998). GTM: A principled alternative to the self-organizing map. *Neural Computation*, 10, 215–234.
- Comon, P. (1994). Independent component analysis: A new concept. *Signal Processing*, 36, 287–314.
- Delyon, B. (1993). *Remarks on filtering of semi-Markov data* (Tech. Rep. 733). Beaulieu, France: Institute de Recherche en Informatique et Systems Aleatoires.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society B*, 39, 1–38.
- Digalakis, V., Rohlicek, J. R., & Ostendorf, M. (1993). ML estimation of a stochastic linear system with the EM algorithm and its application to speech recognition. *IEEE Transactions on Speech and Audio Processing*, 1(4), 431–442.
- Duda, R. O., & Hart, P. E. (1973). *Pattern classification and scene analysis*. New York: Wiley.
- Elliott, R. J., Aggoun, L., & Moore, J. B. (1995). *Hidden Markov models: Estimation and control*. New York: Springer-Verlag.
- Everitt, B. S. (1984). *An introduction to latent variable models*. London: Chapman and Hill.
- Fletcher, R., & Powell, M. J. D. (1963). A rapidly convergent descent method for minimization. *Computing Journal*, 2, 163–168.
- Fraser, A. M., & Dimitriadis, A. (1993). Forecasting probability densities by using hidden Markov models with mixed states. In A. S. Weigend & N. A. Gershenfeld (Eds.), *Time series prediction: Forecasting the future and understanding the past* (pp. 265–282). Reading, MA: Addison-Wesley.
- Geman, S., & Geman, D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6, 721–741.
- Ghahramani, Z., & Hinton, G. (1996a). *Parameter estimation for linear dynamical systems* (Tech. Rep. CRG-TR-96-2). Toronto: Department of Computer Science, University of Toronto. Available from <ftp://ftp.cs.toronto.edu/pub/zoubin/>.
- Ghahramani, Z., & Hinton, G. (1996b). *Switching state-space models* (Tech. Rep.

- CRG-TR-96-3). Toronto: Department of Computer Science, University of Toronto. Submitted for publication.
- Ghahramani, Z., & Hinton, G. (1997). *The EM algorithm for mixtures of factor analyzers* (Tech. Rep. CRG-TR-96-1). Toronto: Department of Computer Science, University of Toronto. Available from <ftp://ftp.cs.toronto.edu/pub/zoubin/>.
- Ghahramani, Z., & Jordan, M. I. (1994). Supervised learning from incomplete data via an EM approach. In J. D. Cowan, G. Tesauro, & J. Alspector (Eds.), *Advances in neural information processing systems, 6* (pp. 120–127). San Mateo, CA: Morgan Kaufmann.
- Goldberg, P., Williams, C., & Bishop, C. (1998). Regression with input-dependent noise: A gaussian process treatment. In M. Kearns, M. Jordan, & T. Solla (Eds.), *Advances in neural information processing systems, 10* (pp. 493–499). Cambridge, MA: MIT Press.
- Goodwin, G. C., & Sin, K. S. (1984). *Adaptive filtering prediction and control*. Englewood Cliffs, NJ: Prentice Hall.
- Hinton, G. E., Dayan, P., & Revow, M. (1997). Modelig the manifolds of images of handwritten digits. *IEEE Transactions on Neural Networks, 8*, 65–74.
- Hinton, G., & Ghahramani, Z. (1997). Generative models for discovering sparse distributed representations. *Philosophical Transactions of the Royal Society B, 352*, 1177–1190.
- Hinton, G. E., Revow, M., & Dayan, P. (1995). Recognizing handwritten digits using mixtures of linear models. In G. Tesauro, D. Touretzky, & T. Leen (Eds.), *Advances in neural information processing systems, 7* (pp. 1015–1022). Cambridge, MA: MIT Press.
- Hinton, G. E., & Zemel, R. S. (1994). Autoencoders, minimum description length, and Helmholtz free energy. In J. D. Cowan, G. Tesauro, & J. Alspector (Eds.), *Advances in neural information processing systems, 6* (pp. 3–10). San Mateo, CA: Morgan Kaufmann.
- Jordan, M. I., & Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation, 6*(2), 181–214.
- Jöreskog, K. G. (1967). Some contributions to maximum likelihood factor analysis. *Psychometrika, 32*, 443–482.
- Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Trans. Am. Soc. Mech. Eng., Series D, Journal of Basic Engineering, 82*, 35–45.
- Kalman, R. E., & Bucy, R. S. (1961). New results in linear filtering and prediction theory. *Trans. Am. Soc. Mech. Eng., Series D, Journal of Basic Engineering, 83*, 95–108.
- Kivinen, J., & Warmuth, M. K. (1997). Exponentiated gradient versus gradient descent for linear predictors. *Journal of Information and Computation, 132*(1), 1–64.
- Kohonen, T. (1982). Self-organized formation of topologically correct feature maps. *Biological Cybernetics, 43*, 59–69.
- Lauritzen, S. L., & Spiegelhalter, D. J. (1988). Local computations with probabilities on graphical structures and their application to expert systems. *J. Royal Statistical Society B, 50*, 157–224.
- Little, R. J. A., & Rubin, D. B. (1987). *Statistical analysis with missing data*. New

- York: Wiley.
- Ljung, L., & Söderström, T. (1983). *Theory and practice of recursive identification*. Cambridge, MA: MIT Press.
- Lloyd, S. P. (1982). Least square quantization in PCM. *IEEE Transactions on Information Theory*, 28, 128–137.
- Lyttkens, E. (1966). On the fixpoint property of Wold's iterative estimation method for principal components. In P. Krishnaiah (Ed.), *Paper in multivariate analysis*. New York: Academic Press.
- MacKay, D. J. C. (1996). *Maximum likelihood and covariant algorithms for independent component analysis* (Tech. Rep. Draft 3.7). Cambridge: Cavendish Laboratory, University of Cambridge.
- Moulines, E., & Cardoso, J.-F., & Gassiat, E. (1997). Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing* (Vol. 5), pp. 3617–3620.
- Neal, R. M. (1993). *Probabilistic inference using Markov chain Monte Carlo methods* (Tech. Rep. CRG-TR-93-1). Toronto: Department of Computer Science, University of Toronto.
- Neal, R. M., & Hinton, G. E. (1998). A view of the EM algorithm that justifies incremental, sparse and other variants. In M. I. Jordan (Ed.), *Learning in graphical models* (pp. 355–368). Dordrecht, MA: Kluwer.
- Nowlan, S. J. (1991). Maximum likelihood competitive learning. In R. P. Lippmann, J. E. Moody, & D. S. Touretzky (Eds.), *Advances in neural information processing systems*, 3 (pp. 574–582). San Mateo, CA: Morgan Kaufmann.
- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. San Mateo, CA: Morgan Kaufmann.
- Pearlmutter, B. A., & Parra, L. C. (1997). Maximum likelihood blind source separation: A context-sensitive generalization of ICA. In M. Mozer, M. Jordan, & T. Petsche (Eds.), *Advances in neural information processing systems*, 9 (pp. 613–619). Cambridge, MA: MIT Press.
- Rabiner, L. R., & Juang, B. H. (1986). An introduction to hidden Markov models. *IEEE ASSP Magazine*, 3(1), 4–16.
- Rauch, H. E. (1963). Solutions to the linear smoothing problem. *IEEE Transactions on Automatic Control*, 8, 371–372.
- Rauch, H. E., Tung, F., & Striebel, C. T. (1965). Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, 3(8), 1445–1450.
- Roweis, S. T. (1998). EM algorithms for PCA and SPCA. In M. Kearns, M. Jordan, and S. Solla (Eds.), *Advances in neural information processing systems*, 10 (626–632). Cambridge, MA: MIT Press. Also Tech Report CNS-TR-97-02, Computation and Neural Systems, Calif. Institute of Technology.
- Rubin, D. B., & Thayer, D. T. (1982). EM algorithms for ML factor analysis. *Psychometrika*, 47(1), 69–76.
- Saul, L. and Rahim, M. (1998). Modeling acoustic correlations by factor analysis. In M. Kearns, M. Jordan, & S. Solla (Eds.), *Advances in neural information processing systems*, 10 (pp. 749–755). Cambridge, MA: MIT Press.
- Shumway, R. H., & Stoffer, D. S. (1982). An approach to time series smoothing and forecasting using the EM algorithm. *Journal of Time Series Analysis*, 3(4),

- 253–264.
- Shumway, R. H., & Stoffer, D. S. (1991). Dynamic linear models with switching. *Journal of the American Statistical Association*, *86*(415), 763–769.
- Sirovich, L. (1987). Turbulence and the dynamics of coherent structures. *Quarterly Applied Mathematics*, *45*(3), 561–590.
- Smyth, P. (1997). Clustering sequences with hidden Markov models. In G. Tesauro, D. Touretzky, & T. Leen (Eds.), *Advances in neural information processing systems*, *9* (pp. 648–654). Cambridge, MA: MIT Press.
- Smyth, P., Heckerman, D., & Jordan, M. (1997). Probabilistic independence networks for hidden Markov probability models. *Neural Computation*, *9*, 227–269.
- Tipping, M. E., & Bishop, C. M. (1999). Mixtures of probabilistic principal component analyzers. *Neural Computation*, *11*, 443–482. Also Tech. Report TR-NCRG/97/003, Neural Computing Research Group, Aston University.
- Tresp, V., Ahmad, S., & Neuneier, R. (1994). Training neural networks with deficient data. In J. D. Cowan, G. Tesauro, & J. Alspector (Eds.), *Advances in neural information processing systems*, *6* (pp. 128–135). San Mateo, CA: Morgan Kaufmann.
- Viterbi, A. J. (1967). Error bounds for convolutional codes and an asymptotically optimal decoding algorithm. *IEEE Transactions on Information Theory*, *IT-13*, 260–269.
- Whittaker, J. (1990). *Graphical models in applied multivariate statistics*. Chichester, England: Wiley.

Received September 5, 1997; accepted April 23, 1998.