

CAMBRIDGE UNIVERSITY
ENGINEERING DEPARTMENT

**GENERALISED
LINEAR GAUSSIAN MODELS**

A-V.I. Rosti & M.J.F. Gales
CUED/F-INFENG/TR.420

November 23, 2001

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

E-mail: {avir2, mjfg}@eng.cam.ac.uk

Abstract

This paper addresses the time-series modelling of high dimensional data. Currently, the hidden Markov model (HMM) is the most popular and successful model especially in speech recognition. However, there are well known shortcomings in HMMs particularly in the modelling of the correlation between successive observation vectors; that is, inter-frame correlation. Standard diagonal covariance matrix HMMs also lack the modelling of the spatial correlation in the feature vectors; that is, intra-frame correlation. Several other time-series models have been proposed recently especially in the segment model framework to address the inter-frame correlation problem such as Gauss-Markov and dynamical system segment models. The lack of intra-frame correlation has been compensated for with transform schemes such as semi-tied full covariance matrices (STC). All these models can be regarded as belonging to the broad class of generalised linear Gaussian models. Linear Gaussian models (LGM) are popular as many forms may be trained efficiently using the expectation maximisation algorithm. In this paper, several LGMs and generalised LGMs are reviewed. The models can be roughly categorised into four combinations according to two different state evolution and two different observation processes. The state evolution process can be based on a discrete finite state machine such as in the HMMs or a linear first-order Gauss-Markov process such as in the traditional linear dynamical systems. The observation process can be represented as a factor analysis model or a linear discriminant analysis model. General HMMs and schemes proposed to improve their performance such as STC can be regarded as special cases in this framework.

1 Introduction

Currently, the most popular and successful time-series model is the hidden Markov model (HMM). HMMs can be applied in a broad range of areas such as speech recognition, bioinformatics and stock market analysis. For example HMM based large vocabulary speech recognition systems have dominated the standard evaluation tasks such as Broadcast News Transcription and Switchboard. The benefits of using HMMs include efficient training and recognition algorithms in which the automatic segmentation using Viterbi algorithm [35] is in very significant role. On the other hand, the state conditional independence assumption in HMMs is a major drawback in case of strongly correlated feature vectors such as in speech recognition [6]. This correlation between successive feature vectors is often called inter-frame correlation. To overcome this deficiency of HMMs, several schemes have been proposed from segmental features to segment models [34]. In general, the use of segment models is restricted by the complicated algorithms due to additional duration models.

Another drawback in the standard diagonal covariance matrix HMMs is its weak spatial correlation modelling. Despite the attempts to reduce this intra-frame correlation by front-end transforms such as discrete cosine transform (DCT) in case of standard Mel-frequency cepstral coefficients [42] there is always some correlation present due to the fixed basis functions in DCT. An optimal front-end should employ data dependent transform bases such as in Karhunen-Loève transform [10] but the complexity of the system becomes equivalent to a full covariance matrix HMMs. Several multiple class based transform schemes have been proposed; e.g., semi-tied full covariance matrices (STC) [11]. Linear discriminant analysis and factor analysis have been also proposed to overcome this problem [29, 40]. In addition to correlation modelling they address the problem of high dimensionality of the feature vectors by allowing lower dimensional subspaces to be used.

The machine learning community has been interested in linear Gaussian models (LGM) for some time now [38] due to the efficiency and applicability of the expectation maximisation (EM) algorithm [5] which provides a consistent framework in supervised learning. Despite the attempts to unify the field of LGMs, several interesting models have been omitted; e.g., independent factor analysis [1], linear discriminant analysis and its extensions [12, 29, 26]. The forms of models reviewed in this paper generalise some of the currently used techniques in the HMM framework such as semi-tied full covariance matrices.

Linear Gaussian models are a subset of more general state-space models that consist of state evolution process and observation process. Strictly speaking linear Gaussian models are state-space models in which the state evolution and observation equations are linear and the distributions are Gaussians [38]. Some models presented in this paper deviate from the strict definition in that mixtures of Gaussians are allowed as the distributions. Therefore they are called generalised LGMs. This generalisation is important in case of feature vectors which tend to have multi-modal distributions due to source and environment differences such as in speech recognition.

Linear Gaussian models also allow the investigation of new subspaces of the feature vectors which often tend to be very high in dimensionality. It is well known that the inclusion of dynamic features such as first and second-order differences to the vectors of cepstral coefficients increase the performance of many recognition systems. That is often thought to handle partially the inter-frame correlation. It is not guaranteed, though, that the standard feature space is optimal [12] and one might argue that there exists some optimal subspace where the modelling assumptions hold better. Determining the optimal state-space dimensionality for different models would be an interesting topic by itself but it is not included in this paper.

This paper is organised as follows. In Section 2, the general modelling framework is described by introducing the stochastic processes involved. Bayesian networks are briefly introduced as a way of illustrating the conditional independencies in probabilistic models. Also, the expectation maximisation algorithm used in estimating the model parameters is outlined. Static models based on factor analysis and mixtures of factor analysers are described in Section 3. In Section 4, linear dynamical systems are presented. Various modelling assumptions used to either decrease the number of model parameters or modify them to better suit the nature of complex signals are introduced as well. Hidden Markov models and an extension of HMMs combining it with a generic

factor analyser are introduced in Section 5. Section 6 is devoted to linear discriminant analysis based observation processes in the both cases of discrete and continuous state evolution processes. In section 7, implementation issues are discussed. Finally, Section 8 concludes this paper and some future directions are proposed.

2 Generalised Linear Gaussian Models

State-space models are generally based on a k dimensional state vector¹ \mathbf{x}_t and a p dimensional observation vector \mathbf{o}_t which satisfy the following generative model

$$\mathbf{x}_{t+1} = f(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{w}_t) \quad (1)$$

$$\mathbf{o}_t = g(\mathbf{x}_t, \mathbf{v}_t) \quad (2)$$

where the function $f(\cdot)$ describes the state evolution process and the function $g(\cdot)$ the observation process. Random vectors \mathbf{w}_t and \mathbf{v}_t are called state evolution noise and observation noise, respectively. Although the models described by Eqs. 1 and 2 can exhibit any linear or non-linear functions, the models presented in this paper are restricted to linear ones. In strict definition of linear Gaussian models the noise sources \mathbf{w}_t and \mathbf{v}_t are distributed according to Gaussian distributions $\mathcal{N}(\boldsymbol{\mu}_t^{(x)}, \boldsymbol{\Sigma}_t^{(x)})$ and $\mathcal{N}(\boldsymbol{\mu}_t^{(o)}, \boldsymbol{\Sigma}_t^{(o)})$. Since many feature vectors generally do not have unimodal distributions due to source and environment variability, the unimodal assumption is relaxed by including mixture models in this paper. Often the time dependence of the noise source statistics is also omitted so that the distributions have static parameters, $\boldsymbol{\mu}^{(o)}$ and $\boldsymbol{\mu}^{(x)}$. In the same way the noise vectors can be written without the time indices as \mathbf{v} and \mathbf{w} .

The state evolution process may be viewed as some underlying phenomenon which may be inherent for the process to be modelled. Alternatively, it may be viewed only as a compact representation of the high dimensional observation vectors. For example in speech recognition, the state vector may be viewed as representing positions of the articulatory organs and the state evolution process, $f(\cdot)$, describes their movement in time. Although the above generic model fits into this interpretation, the linear Gaussian models are a crude approximation since the movement of articulators is rather non-linear [4]. Due to this non-linear nature, the articulatory interpretation is not often stressed.

The observation equation, $g(\cdot)$, describes the function mapping the current state vector \mathbf{x}_t mixed with the observation noise \mathbf{v}_t onto the current observation \mathbf{o}_t which is usually higher in dimensionality; i.e. $p > k$. For example in speech recognition, the observation process may be viewed as mapping the positions of the articulatory organs onto the measurable acoustic representation where the observation noise \mathbf{v}_t represents all the noise induced by the environment and the recording equipment. As in case of state evolution process, this interpretation is only valid for non-linear models [4]. In this paper, the observation process is only viewed as a scheme to carry out dimensionality reduction.

2.1 State Evolution Process

Firstly, all the state evolution processes used in the following sections are presented. Although the dynamic models are more interesting, the static counterparts provide a general framework to

¹In this paper, bold capital letters are used to denote matrices e.g. \mathbf{A} , bold letters refer to vectors e.g. \mathbf{a} and plain letters represent scalars e.g. c . All vectors are column vectors. Prime is used to denote the transpose of a matrix or a vector e.g. \mathbf{A}' , \mathbf{a}' . The determinant of a matrix is denoted by $|\mathbf{A}|$. Gaussian distributed vectors e.g. \mathbf{x} with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ are denoted by $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. The likelihood of a vector \mathbf{z} being generated by the above Gaussian; i.e., the Gaussian evaluated at the point \mathbf{z} , is represented as $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$. Vectors distributed according to a mixture of Gaussians are denoted shortly by $\mathbf{x} \sim \sum_m c_m \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$. As opposed to the lower case letter p used to represent a continuous distribution, a capital letter P is used to denote a probability mass function of a discrete variables. The probability that a discrete random variable ω equals m is denoted shortly as $P(\omega = m) = P_\omega(m)$ and the subscript is also omitted if not ambiguous.

begin with. The relationship between the static models and the traditional covariance modelling techniques are discussed in the following sections.

The simplest non-trivial stochastic state evolution process is a static multivariate Gaussian with a mean vector $\boldsymbol{\mu}^{(x)}$ and a covariance matrix $\boldsymbol{\Sigma}^{(x)}$. In the case of static models, the time index t in the subscript is omitted. The generative model can be represented as follows

$$\boldsymbol{x} = \boldsymbol{w}, \quad \boldsymbol{w} \sim \mathcal{N}(\boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \quad (3)$$

It is often possible to make further assumptions about the distribution of the noise depending on the situation. Especially diagonal covariance matrices can be used if the state evolution noise is assumed to be spatially uncorrelated.

The mixture of Gaussians state process adds a new discrete hidden layer in the model where the individual Gaussians are selected by a hidden discrete mixture indicator variable ω . Generally, mixture of Gaussians model can be represented as follows

$$\boldsymbol{x} = \boldsymbol{w}, \quad \boldsymbol{w} \sim \sum_n c_n^{(x)} \mathcal{N}(\boldsymbol{\mu}_n^{(x)}, \boldsymbol{\Sigma}_n^{(x)}) \quad (4)$$

where the mixture priors $c_n^{(x)} = P_\omega(n)$ must sum to unity to guarantee $p(\boldsymbol{x})$ being a true distribution. It is also possible to use factorial mixture of Gaussians where individual elements of the noise vector are drawn from a univariate mixture of Gaussians; e.g., see [1]. The factorial mixture includes a stronger assumption on the independence of the individual state vector components but provides more flexible modelling of the component variances. The difference between mixture and factorial mixture models is the same as between vector quantisation and cooperative vector quantisation [43] as well as hidden Markov models and factorial hidden Markov models [18].

The simplest dynamic state evolution process is a linear first-order Gauss-Markov random process. A Markov assumption is made to the general state evolution process in Eq. 1 so that the current state vector \boldsymbol{x}_t depends only on the previous state vector \boldsymbol{x}_{t-1} . The new state vector is generated by a linear state evolution process from the current state vector as follows

$$\boldsymbol{x}_{t+1} = \boldsymbol{A}_t \boldsymbol{x}_t + \boldsymbol{w}_t, \quad \boldsymbol{w}_t \sim \mathcal{N}(\boldsymbol{\mu}_t^{(x)}, \boldsymbol{\Sigma}_t^{(x)}) \quad (5)$$

where \boldsymbol{A}_t is a k by k state transition matrix and \boldsymbol{w} the state evolution noise. Since the initial state vector \boldsymbol{x}_1 is assumed to be a Gaussian with a mean vector $\boldsymbol{\mu}^{(i)}$ and a covariance matrix $\boldsymbol{\Sigma}^{(i)}$, all the subsequent state vectors are Gaussian distributed as follows

$$p(\boldsymbol{x}_t | \boldsymbol{x}_{t-1}) = \mathcal{N}(\boldsymbol{x}_t; \boldsymbol{A}_t \boldsymbol{x}_{t-1} + \boldsymbol{\mu}_t^{(x)}, \boldsymbol{\Sigma}_t^{(x)}) \quad (6)$$

and their joint likelihood is a Gaussian. The state evolution noise is often assumed to be zero mean and spatially uncorrelated; i.e., $\boldsymbol{\Sigma}^{(x)}$ is diagonal. Non-zero mean vector corresponds to a constant drift of the state vector in the direction of $\boldsymbol{\mu}^{(x)}$ and diagonality of $\boldsymbol{\Sigma}^{(x)}$ is often used to reduce the number of model parameters. If the state evolution noise is distributed as a mixture of Gaussians, the state vector is also distributed according to a mixture of Gaussians with the number of components growing exponentially. The implications of this growth are discussed later in this paper.

A Gauss-Markov random process with parameters

$$\boldsymbol{\mu}^{(i)} = \begin{bmatrix} -1 \\ -2 \end{bmatrix}, \quad \boldsymbol{\Sigma}^{(i)} = \begin{bmatrix} 1.3 & 0 \\ 0 & 2.5 \end{bmatrix}, \quad \boldsymbol{A} = \begin{bmatrix} 0.9 & -0.2 \\ 0.6 & 0.6 \end{bmatrix}, \quad \boldsymbol{\mu}^{(x)} = \begin{bmatrix} 1.8 \\ 1.3 \end{bmatrix}, \quad \boldsymbol{\Sigma}^{(x)} = \begin{bmatrix} 0.3 & 0 \\ 0 & 0.5 \end{bmatrix}$$

is illustrated on the left-hand side of Figure 1. At $t = 1$, the main axes of the k dimensional ellipsoid ($k = 2$ in the figure) representing the covariance matrix of the state vector are parallel to the coordinate axes since a diagonal initial state covariance matrix is used. At every subsequent time step, the covariance ellipsoid of the preceding state vector is stretched and rotated according to the state evolution matrix and the resulting ellipsoid is convolved with the ellipsoid representing the state evolution noise covariance. The state evolution noise mean vector causes the linear drift

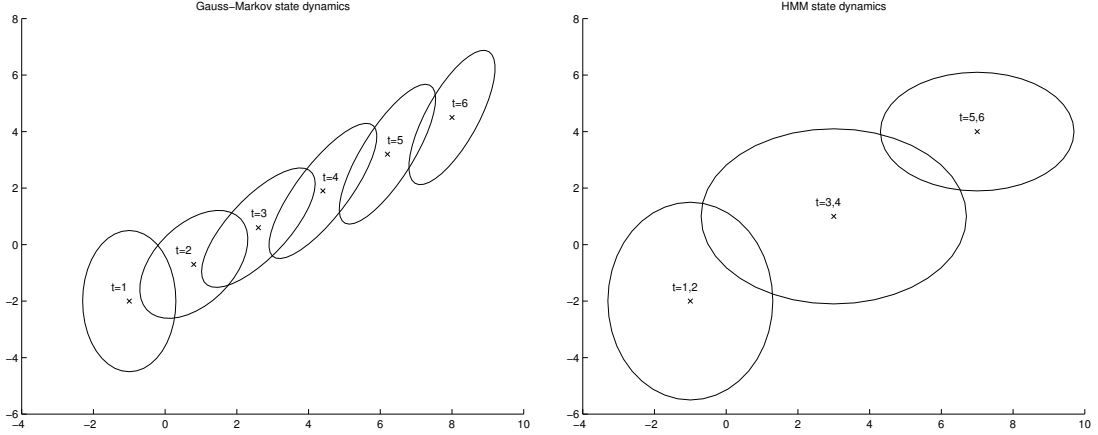


Figure 1: State evolution processes of 6 time steps in length. 1) Linear first-order Gauss-Markov process. 2) Hidden Markov model.

of the ellipsoid centres (marked by x). In modelling of zero mean signals, state evolution noise mean can be omitted.

The last type of the state evolution can be assumed to be generated by a discrete state sequence rather than continuous as above. The most commonly used dynamic discrete state model is the hidden Markov model (HMM). The hidden Markov models are based on N_s hidden states with state conditional densities $b_j(\mathbf{x}_t) = p(\mathbf{x}_t|j)$, a set of transition probabilities $a_{ij} = P_{q_t|q_{t-1}}(j|i)$ and initial state probabilities $\pi_j = P_{q_1}(j)$. The joint likelihood of the state vector sequence \mathbf{X} and the discrete state sequence Q is

$$p(\mathbf{X}, Q) = \pi_{q_1} b_{q_1}(\mathbf{x}_1) a_{q_1 q_2} b_{q_2}(\mathbf{x}_2) a_{q_2 q_3} b_{q_3}(\mathbf{x}_3) \dots \quad (7)$$

Since, in practise, the underlying discrete state sequence is unknown (hidden) the likelihood is estimated by summing over all possible discrete state sequences

$$p(\mathbf{X}) = \sum_{\{Q_T\}} \pi_{q_1} \prod_{t=2}^T b_{q_t}(\mathbf{x}_t) a_{q_t-1 q_t} \quad (8)$$

where $\{Q_T\}$ denotes the set of all possible state sequences of length T . In HMMs the output vectors are generated by one of the N_s discrete states and the transition from one discrete state to the next is assumed to be abrupt. The hidden Markov model as a generative model for \mathbf{x}_{t_i} is represented later in this paper as follows

$$\mathbf{x}_t \sim \mathcal{M}^{hmm} \quad (9)$$

A three state hidden Markov model state evolution process with the following state distribution parameters

$$\begin{aligned} \boldsymbol{\mu}_1^{(x)} &= \begin{bmatrix} -1 \\ -2 \end{bmatrix} & \boldsymbol{\mu}_2^{(x)} &= \begin{bmatrix} 3 \\ 1 \end{bmatrix} & \boldsymbol{\mu}_3^{(x)} &= \begin{bmatrix} 7 \\ 4 \end{bmatrix} \\ \boldsymbol{\Sigma}_1^{(x)} &= \begin{bmatrix} 2.3 & 0 \\ 0 & 3.5 \end{bmatrix} & \boldsymbol{\Sigma}_2^{(x)} &= \begin{bmatrix} 3.7 & 0 \\ 0 & 3.1 \end{bmatrix} & \boldsymbol{\Sigma}_3^{(x)} &= \begin{bmatrix} 2.7 & 0 \\ 0 & 2.1 \end{bmatrix} \end{aligned} \quad (10)$$

is depicted on the right-hand side of Figure 1. Since diagonal state conditional densities are used, the main axes of the k dimensional ellipsoid ($k = 2$ in the figure) at a given time are parallel to the coordinate axes. The figure represents a three state HMM and one possible assignment of the state vectors into the states in case of a sequence of six output vectors. The HMM performs clustering

of the state vector space in case of full transition probability matrix but different constraints can be set for the transition probabilities according to the nature of the signals to be modelled. In case of causal time-series modelling, so called left to right HMMs are often used so that the discrete states represent parts of the modelling units that can occur only sequentially.

2.2 Observation Process

Two different observation processes are presented here. Both processes can perform dimension reduction in that the state-space is lower in dimensionality than the observation space; $k < p$. The first observation equation is called factor analysis (FA) since in a static case, the model reduces to a standard FA model with some additional assumptions about the state and noise distributions. The factor analysis observation process can be represented as

$$\mathbf{o}_t = \mathbf{C}_t \mathbf{x}_t + \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(\boldsymbol{\mu}_t^{(o)}, \boldsymbol{\Sigma}_t^{(o)}) \quad (11)$$

where \mathbf{C} is a p by k observation matrix and \mathbf{v} the observation noise. The observation noise is independent of the state vector and its covariance matrix is often assumed to be diagonal to capture the correlation into the observation matrix. In case of static data modelling, the time indices in the observation equation can be omitted. If the current state vector is Gaussian distributed, the current observation is a Gaussian with the following distribution

$$p(\mathbf{o}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}_t \mathbf{x}_t + \boldsymbol{\mu}_t^{(o)}, \boldsymbol{\Sigma}_t^{(o)}) \quad (12)$$

and all the possible joint distributions are Gaussians. Factor analysis and different mixture assumptions are described in more detail in Section 3.

The second type of the observation equation is called linear discriminant analysis (LDA) because in a restricted static case, the model reduces to a standard LDA. The meaningful dimensions are modelled by the state vector and the nuisance dimensions are modelled by a single Gaussian distributed observation noise as follows

$$\mathbf{o}_t = \mathbf{C}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{v}_t \end{bmatrix}, \quad \mathbf{v}_t \sim \mathcal{N}(\boldsymbol{\mu}_t^{(o)}, \boldsymbol{\Sigma}_t^{(o)}) \quad (13)$$

where \mathbf{C}_t is a p by p matrix and \mathbf{v} is a $p - k$ dimensional vector. There are several ways to express the posterior of the observations. The most commonly used is

$$p(\mathbf{o}_t | \mathbf{x}_t) = \text{abs}|\mathbf{C}_t^{-1}| \mathcal{N}(\mathbf{C}_t^{-1} \mathbf{o}_t; \begin{bmatrix} \mathbf{x}_t \\ \boldsymbol{\mu}_t^{(o)} \end{bmatrix}, \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_t^{(o)} \end{bmatrix}) \quad (14)$$

where the first k columns and k rows of the covariance matrix are filled with zeroes because the first k elements of the observation vector are deterministic given the current state vector \mathbf{x}_t . The first term in the equation is needed to scale the likelihood to be valid and the absolute value is needed to guarantee non-negative likelihoods. Since the deterministic elements do not have any influence on the likelihoods, another way of representing the posterior is

$$p(\mathbf{o}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}_t \begin{bmatrix} \mathbf{x}_t \\ \boldsymbol{\mu}_t^{(o)} \end{bmatrix}, \mathbf{C}_{t[p-k]} \boldsymbol{\Sigma}_t^{(o)} \mathbf{C}'_{t[p-k]}) \quad (15)$$

where $\mathbf{C}_{t[p-k]}$ denotes a matrix consisting of the $p - k$ rightmost columns of matrix \mathbf{C}_t . In this representation, the covariance matrix is clearly non-singular provided the transformation and the observation noise covariance are non-singular as well. Linear discriminant analysis is reviewed in more detail in Section 6.

2.3 Basic Models

Different combinations of the state and observation processes presented in this paper are depicted in Figures² 2 and 3. Figure 2 depicts models based on static state evolution process. In the top of

²These diagrams do not illustrate all the possible forms of linear Gaussian models. E.g., several combinations using LDA observation process are not present.

the diagram is a single static multivariate Gaussian as the basis for all the subsequent models. A mixture of Gaussians can be regarded as vector quantisation with the Gaussian mean vectors as the cluster centres. Vector quantisation is often used in the initialisation of, e.g., HMM parameters [42]. Instead of using only one code-book to assign a vector into a single cluster, cooperative vector quantisation [43] or factorial mixture of Gaussians uses several independent code-books in a distributed manner. All the three models above correspond to state-space model with identity matrix as the observation matrix and zero observation noise.

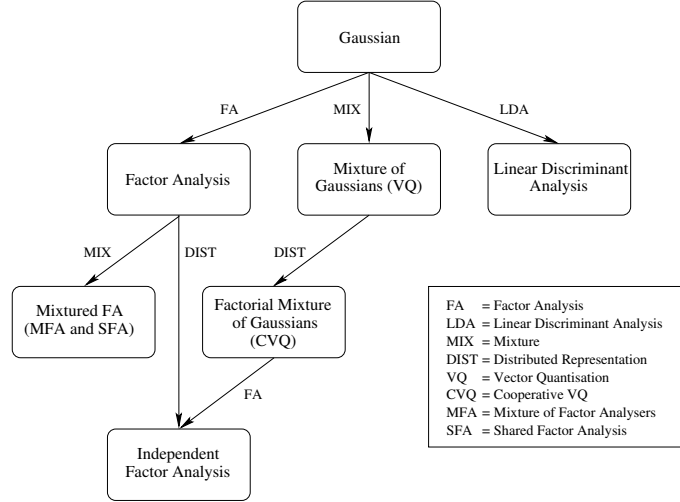


Figure 2: Diagram of static linear Gaussian models. The arrows represent additional properties to the model they are attached to.

Factor analysis [20, 26, 39] is based on a static multivariate Gaussian state process and a factor analysis observation process. In standard factor analysis the state vectors are assumed to be distributed according to a standard Gaussian density, $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Independent factor analysis [1], mixture of factor analysers [15] and shared factor analysis [22] are based on the factor analysis model with different mixture assumptions. These mixture assumptions are described in a unified way in Section 3. Linear discriminant analysis [26] is another static data modelling scheme with linear discriminant analysis observation process.

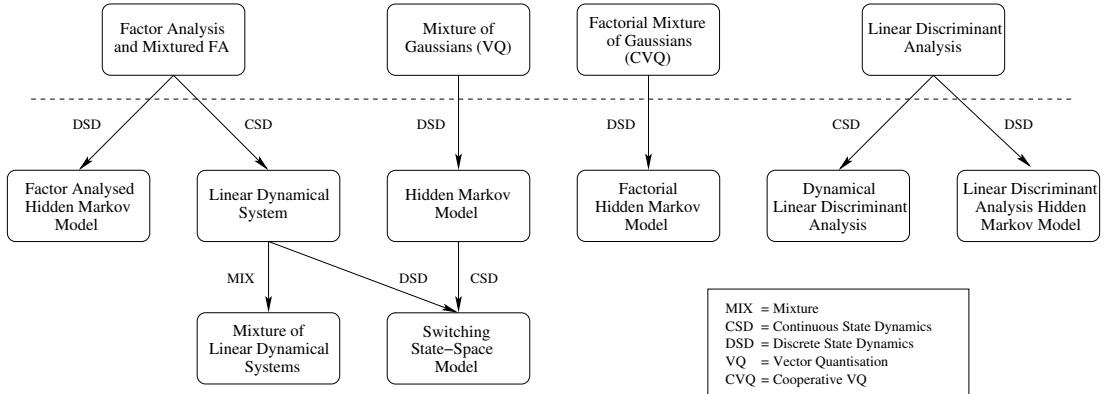


Figure 3: Dynamic linear Gaussian models and how they relate to some of the static models.

Dynamic linear Gaussian models and the corresponding static models are illustrated in Figure 3. Dynamic models with factor analysis observation process include linear dynamical systems [6, 7, 16, 32, 34, 38], mixture of linear dynamical systems and switching state-space model [17, 33]

as well as different variations of factor analysed HMMs presented later in this paper and its restricted version in [40]. The linear discriminant observation process is illustrated in case of HMM based [12, 21, 29] and linear first-order Gauss-Markov based state evolution processes. The dynamical linear discriminant analysis is presented later in this paper.

Standard hidden Markov model [35, 42] can be considered as a special case of both the observation processes by just omitting the observation noise and setting the observation matrix to an identity matrix; i.e., $\mathbf{C} = \mathbf{I}$. Also semi-tied covariance matrix HMMs (STC) [11] can be described by both observation processes when $k = p$ and $\mathbf{v} = \mathbf{0}$. Factorial hidden Markov models [18] use distributed representation of the discrete state-space so that several independent HMMs can be viewed to have produced the observation vectors.

2.4 Bayesian Networks

In this paper, Bayesian networks [14] are used to illustrate the statistical independencies between different random variables in the probabilistic models. Bayesian networks are directed acyclic graphs, also known as graphical models. The notation is adopted from [33] where round nodes were used to denote continuous and squared nodes discrete random variables. The observable variables are shaded and the lack of an arrow between nodes represents conditional independence.

As an example, consider three continuous random variables \mathbf{z}, \mathbf{x} and \mathbf{o} of which \mathbf{o} is observed and the others are hidden. The joint likelihood can be factored as a product of conditional likelihoods as follows

$$p(\mathbf{z}, \mathbf{x}, \mathbf{o}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})p(\mathbf{o}|\mathbf{x}, \mathbf{z}) \quad (16)$$

This factorisation is perfectly valid in every case. If no assumptions of conditional independence can be made, the corresponding Bayesian network must be illustrated as the first network in Figure 4. There are two arrows pointing the node representing \mathbf{o} since in the above factorisation \mathbf{o} depends on both \mathbf{z} and \mathbf{x} .

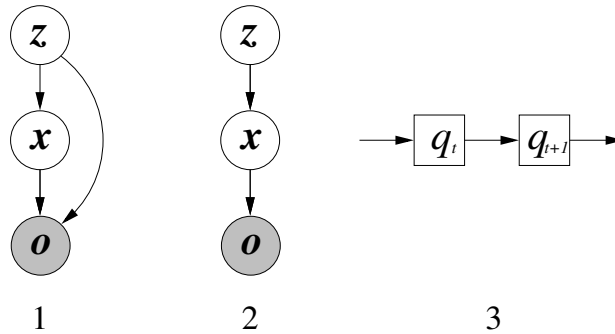


Figure 4: Examples of Bayesian networks representing different assumptions on conditional independence. 1) Continuous random variables \mathbf{z} , \mathbf{x} and \mathbf{o} (shading denotes observable) are fully dependent, 2) \mathbf{o} is conditionally independent of \mathbf{z} given \mathbf{x} , 3) discrete random variable q_{t+1} is conditionally independent of all its predecessors given q_t (discrete Markov chain).

If the random variable \mathbf{o} is assumed to be conditionally independent of \mathbf{z} given \mathbf{x} the joint likelihood can be rewritten as follows

$$p(\mathbf{z}, \mathbf{x}, \mathbf{o}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})p(\mathbf{o}|\mathbf{x}) \quad (17)$$

The corresponding Bayesian network is depicted as the second graph in Figure 4. The arrow between the nodes representing \mathbf{z} and \mathbf{o} can be deleted.

As another example, consider an ordered set of discrete random variables, $Q = q_1, \dots, q_T$, which are encountered several times in the following sections. The joint likelihood of the variables

up to time instant $t + 1$ can be written using conditional likelihoods as follows

$$P(q_1, \dots, q_{t+1}) = P(q_1)P(q_2|q_1)P(q_3|q_1, q_2) \dots P(q_{t+1}|q_1, \dots, q_t) \quad (18)$$

Often a simplification is achieved by using a Markov assumption which says that the likelihood of the variable q_{t+1} is conditionally independent of all the other previous variables given the immediate predecessor q_t ; i.e.,

$$P(q_{t+1}|q_1, \dots, q_t) = P(q_{t+1}|q_t) \quad (19)$$

This is often called a discrete Markov chain and is illustrated as the third graph of Figure 4.

2.5 The EM Algorithm

Estimation of the parameters in probabilistic models involving latent (hidden) variables is often carried out using the expectation maximisation (EM) algorithm [5]. To find the maximum likelihood (ML) estimates for the model parameters, the joint log-likelihood of the data $\log p(\mathbf{O})$ has to be maximised. This can be done in a single run if there are no hidden variables present. The hidden variables, often a collection of hidden state vectors \mathbf{X} , can be regarded as missing data. In the EM framework the objective is to maximise the log-likelihood of the complete data, $\log p(\mathbf{O}, \mathbf{X}|\hat{\mathcal{M}})$, given the new model parameters, $\hat{\mathcal{M}}$, via an iterative two step process of estimating the ML solution for the hidden state vectors given the observed data, and then maximising the joint log-likelihood.

By using Jensen's inequality, the joint log-likelihood of the data given the new model parameters is bound from below as follows

$$\log p(\mathbf{O}|\hat{\mathcal{M}}) = \log \int p(\mathbf{O}, \mathbf{X}|\hat{\mathcal{M}})d\mathbf{X} \geq \int p^\circ(\mathbf{X}|\mathbf{O}) \log \frac{p(\mathbf{O}, \mathbf{X}|\hat{\mathcal{M}})}{p^\circ(\mathbf{X}|\mathbf{O})}d\mathbf{X} \quad (20)$$

where $p^\circ(\mathbf{X}|\mathbf{O})$ is an arbitrary conditional distribution and $\hat{\mathcal{M}}$ is the set of new parameters. The equality in Eq. 20 is attained exactly when the conditional distribution is chosen to be the posterior of the state sequence given the observation sequence and the set of old model parameters, $p^\circ(\mathbf{X}|\mathbf{O}) = p(\mathbf{X}|\mathbf{O}, \mathcal{M})$. Since the denominator on the right hand side of Eq. 20 is independent of the new model parameters, the function to be maximised reduces to the expectation of the log-likelihood of the complete data given the new model parameters with respect to the sequence of state vectors given the observation sequence and the old model parameters; i.e.,

$$\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) = E\{\log p(\mathbf{O}, \mathbf{X}|\hat{\mathcal{M}})|\mathbf{O}, \mathcal{M}\} = \int p(\mathbf{X}|\mathbf{O}, \mathcal{M}) \log p(\mathbf{O}, \mathbf{X}|\hat{\mathcal{M}})d\mathbf{X} \quad (21)$$

which is often called as an auxiliary function and is denoted by $\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}})$. The selection of the arbitrary conditional distribution guarantees that in any iteration, the likelihood never decreases.

The derivation of the EM algorithm for a given model always starts off with finding the expected state statistics given the observation sequence and finding the log-likelihood function for the complete data. With linear Gaussian models the statistics of the posterior needed in the maximisation of the auxiliary function are just the mean vector and covariance matrix since the posteriors are also Gaussian. In case of dynamic models based on Gauss-Markov dynamics, the cross covariance matrix between two successive state vectors has to be obtained but it is easy to obtain since the state vectors are also jointly Gaussian.

In the appendices, the auxiliary function is written as a function of the new model parameters marked with caps; e.g. $\hat{\mathcal{C}}$. It has to be noted that in the M step, $\hat{\mathcal{C}}$ is treated as a variable for which the value that maximises the auxiliary function is determined using standard optimisation methods.

3 Factor Analysis

Factor analysis is a statistical method for modelling the covariance structure of high dimensional static data using a small number of latent (hidden) variables [26]. Traditionally, the latent variables (state vector), \mathbf{x} , are assumed to be distributed according to a standard Gaussian density, $\mathcal{N}(\mathbf{0}, \mathbf{I})$, so that the covariance structure would be captured into the factor loading matrix (observation matrix) \mathbf{C} . If the state vectors had a non-zero mean and its covariance matrix was not an identity matrix, the model would be degenerate since the state vector mean can always be subsumed to the observation noise mean vector and the covariance structure to the observation matrix [38]. The observation noise is assumed to be a Gaussian with mean vector $\boldsymbol{\mu}^{(o)}$ and diagonal covariance matrix $\boldsymbol{\Sigma}^{(o)}$. Diagonality of the observation noise covariance matrix is important since otherwise the parameter estimation could choose the loading matrix to be zero and estimate the observation noise by the sample averages. This would be a maximum likelihood estimate but not very interesting since all the information would be regarded as a p dimensional noise and no gain would be attained by the covariance model.

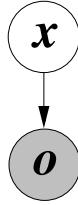


Figure 5: Bayesian network representing factor analysis model.

The generative model of factor analysis can be written as follows

$$\mathbf{x} = \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (22)$$

$$\mathbf{o} = \mathbf{C}\mathbf{x} + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (23)$$

where the state vector \mathbf{x} is k dimensional, the observation noise \mathbf{v} as well as the observation \mathbf{o} are p dimensional and \mathbf{C} is a p by k observation matrix. The generative model of factor analysis can be represented as the Bayesian network in Figure 5. The likelihood of an observation given the state is simply

$$p(\mathbf{o}_j | \mathbf{x}) = \mathcal{N}(\mathbf{o}_j; \mathbf{C}\mathbf{x} + \boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (24)$$

since when the state vector \mathbf{x} is given, the product $\mathbf{C}\mathbf{x}$ is a constant vector added to the observation noise vector \mathbf{v} . Furthermore, the joint likelihood of the state and observation is Gaussian with the following distribution

$$p(\mathbf{o}_j, \mathbf{x}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{o}_j \\ \mathbf{x} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}^{(o)} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{C}\mathbf{C}' + \boldsymbol{\Sigma}^{(o)} & \mathbf{C} \\ \mathbf{C}' & \mathbf{I} \end{bmatrix}\right) \quad (25)$$

According to the generative model in Eqs. 22 and 23, the k dimensional factors are distributed as a standard Gaussian distribution as illustrated on the upper left-hand side of Figure 6 where $k = 1$. The covariance matrix of the factors is therefore a k dimensional ball with unit radius. The k dimensional ball is then stretched and rotated according to the observation matrix \mathbf{C} . The observation noise covariance can be represented as a p dimensional ellipsoid with the main axes parallel to the coordinate axes due to the diagonal covariance matrix assumption. Since the sum of two independent random variables is distributed according to the convolution of the individual distributions, the covariance matrix of the transformed factors is convolved with the ellipsoid

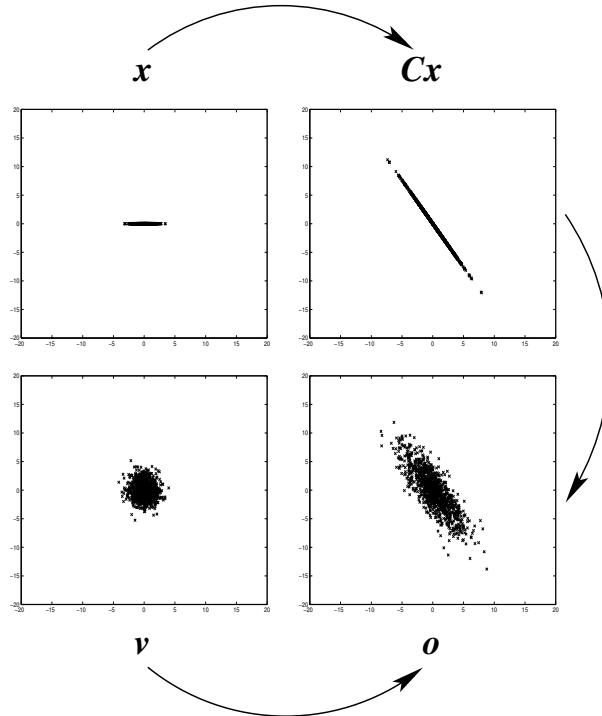


Figure 6: Factor analysis with 1-dimensional state-space and 2-dimensional observation space. $\mathcal{N}(0, 1)$ distributed state vectors, x , are stretched and rotated according to the transform C , and $\mathcal{N}(\mathbf{0}, \Sigma^{(o)})$ distributed noise vectors are added to form the observation vectors. It should be noted that a sum of two random variables is distributed according to the convolution of the individual distributions; in this example $\mathcal{N}(\mathbf{0}, CC' + \Sigma^{(o)})$.

representing the observation noise covariance as depicted on the bottom right-hand side of Figure 6.

The number of parameters needed to model the observations as a single p dimensional multivariate Gaussian is $p(p + 1)/2 + p$ where the first term in the sum corresponds to the symmetric covariance matrix and the second to the mean vector. A factor analysis model requires $pk + 2p$ parameters where the first term of the sum corresponds to the factor loading matrix and the second to the diagonal covariance matrix and the mean vector of the observation noise. A reduction in the number of model parameters by using factor analysis can be attained by choosing the state-space dimensionality according to $k < (p - 1)/2$.

There are two different approaches in solving the factor analysis problem [26]. The direct method is computationally very expensive since it requires spectral factorisation of the factor loading matrix. The maximum likelihood (ML) solution suits better the linear Gaussian model framework presented here since the efficient EM algorithm can be used.

3.1 EM Algorithm for Factor Analysis

Let $\mathbf{O} = \mathbf{o}_1, \dots, \mathbf{o}_N$ be a set of N independent observation vectors. In the case of factor analysis the auxiliary function can be simplified since the log-likelihood of the complete data depends only on the observation vectors due to the distribution assumption on the factors. The auxiliary function can be expressed as follows

$$\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) = E\{\log p(\mathbf{O}|\mathbf{x}, \hat{\mathcal{M}})|\mathbf{O}, \mathcal{M}\} = \sum_{j=1}^N \int p(\mathbf{x}|\mathbf{o}_j, \mathcal{M}) \log p(\mathbf{o}_j|\mathbf{x}, \hat{\mathcal{M}}) d\mathbf{x} \quad (26)$$

where $\hat{\mathcal{M}}$ is the set of new model parameters. Since the observation vectors are independent, the posterior of the data given the factors becomes a sum after taking the logarithm.

The E step requires estimating the statistics of the posterior of the state vectors. Since the posterior is also a Gaussian, only the first and second-order statistics are non-zero as follows

$$\hat{\mathbf{x}}(j) = \mathbf{K}(\mathbf{o}_j - \boldsymbol{\mu}^{(o)}) \quad (27)$$

$$\hat{\mathbf{R}}(j) = \mathbf{I} - \mathbf{K}\mathbf{C} + \hat{\mathbf{x}}(j)\hat{\mathbf{x}}'(j) \quad (28)$$

where $\mathbf{K} = \mathbf{C}'(\mathbf{C}\mathbf{C}' + \boldsymbol{\Sigma}^{(o)})^{-1}$ which can be computed in advance for the current model set.

The M step is also very straightforward. However it requires some manipulation to obtain the new parameters using convenient matrix operations. The new model parameters can be obtained by the following two re-estimation formulae

$$\left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] = \left(\sum_{j=1}^N [\mathbf{o}_j \hat{\mathbf{x}}'(j) \mathbf{o}_j] \right) \left(\sum_{j=1}^N \begin{bmatrix} \hat{\mathbf{R}}(j) & \hat{\mathbf{x}}(j) \\ \hat{\mathbf{x}}'(j) & 1 \end{bmatrix} \right)^{-1} \quad (29)$$

$$\hat{\boldsymbol{\Sigma}}^{(o)} = \frac{1}{N} \sum_{j=1}^N \text{diag} \left(\mathbf{o}_j \mathbf{o}_j' - \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] [\mathbf{o}_j \hat{\mathbf{x}}'(j) \mathbf{o}_j]' \right) \quad (30)$$

where $\text{diag}(\cdot)$ denotes setting the elements outside the main diagonal to zeroes. The details of the derivation are presented in Appendix B.

3.2 Mixture of Factor Analysers

The standard factor analysis model works well for correlated data with Gaussian distribution provided the number of factors (state-space dimensionality) is chosen well. In reality the data are not always Gaussian distributed; e.g., speech feature vectors may have bimodal distributions due to the gender variations. A generic mixture of factor analysers (MFA) is based on three random indicator variables ω^x , ω^o and ω^c with priors $c_n^{(x)} = P_{\omega^x}(m)$, $c_m^{(o)} = P_{\omega^o}(m)$ and $c_l^{(c)} = P_{\omega^c}(l)$. The Bayesian network corresponding to this model is illustrated as the first graph in Figure 7. The state-space indicator, ω^x , chooses the state noise parameters $(\boldsymbol{\mu}_n^{(x)}, \boldsymbol{\Sigma}_n^{(x)})$ from $M^{(x)}$ different sets, the observation noise indicator, ω^o , chooses the state noise parameters $(\boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)})$ from $M^{(o)}$ different sets and the observation matrix indicator chooses a constant matrix \mathbf{C}_l from $M^{(c)}$ alternatives.

It should be noted that in this generic notation, the distribution of the observation matrix is crucial for the model to be a generalised linear Gaussian model. If the observation matrix is Gaussian distributed, the observations are not anymore distributed as Gaussians or even mixtures of Gaussians. To guarantee the observations being Gaussians, the observation matrix is chosen by the mixture indicator from a set of possible constant matrices. Therefore, the observation matrix is distributed according to a weighted delta distribution as presented in the generative model below.

Several parameters of MFA can be assumed to depend on the same mixture indicator. For example, the second graph in Figure 7 represents a model where all the observation parameters $(\mathbf{C}_m, \boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)})$ depend on the same indicator, ω^o . The third graph depicts a case where all the parameters, $(\boldsymbol{\mu}_m^{(x)}, \boldsymbol{\Sigma}_m^{(x)}, \mathbf{C}_m, \boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)})$, depend on a single indicator, ω . In this paper only the latter two cases are considered.

The generative model for the generic case without parameter tying can be written as

$$\mathbf{x} = \mathbf{w}, \quad \mathbf{w} \sim \sum_n c_n^{(x)} \mathcal{N}(\boldsymbol{\mu}_n^{(x)}, \boldsymbol{\Sigma}_n^{(x)}) \quad (31)$$

$$\begin{aligned} \mathbf{v} &\sim \sum_m c_m^{(o)} \mathcal{N}(\boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)}) \\ \mathbf{o} = \mathbf{C}\mathbf{x} + \mathbf{v}, \quad \mathbf{C} &\sim \sum_l c_l^{(c)} \delta(\mathbf{C} - \mathbf{C}_l) \end{aligned} \quad (32)$$

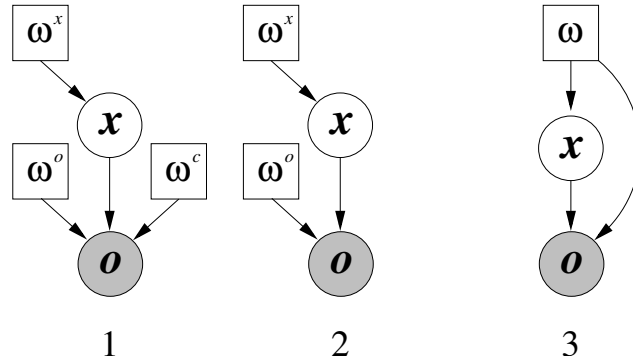


Figure 7: Bayesian networks representing mixture of factor analysers. 1) Generic MFA with three mixture indicators. 2) Observation parameters ($\mathbf{C}_m, \boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)}$) depend on the same indicator. 3) All parameters ($\boldsymbol{\mu}_m^{(x)}, \boldsymbol{\Sigma}_m^{(x)}, \mathbf{C}_m, \boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)}$) depend on the same indicator.

which can be easily modified to represent any of the other mixture models just by modifying the corresponding superscripts. Only difference between the training of these models comes with the posterior estimation. Using the conditional independence assumptions, the joint likelihood of the observation, state and the mixture components is simply

$$p(\mathbf{o}_j, \mathbf{x}, m, n) = c_n^{(x)} p(\mathbf{x}|n) c_m^{(o)} p(\mathbf{o}_j|\mathbf{x}, m) \quad (33)$$

In general, learning the model parameters of mixture of factor analysers can be carried out by the EM algorithm. The normal factor analysis E step has to be modified to take component posteriors into account. Regardless of the model assumptions, the joint component posteriors can be represented as follows

$$\gamma_{mn}(j) = P(m, n|\mathbf{o}_j, \mathcal{M}) = \frac{c_m^{(o)} c_n^{(x)} p(\mathbf{o}_j|m, n, \mathcal{M})}{\sum_{l=1}^{M^{(o)}} c_l^{(o)} \sum_{i=1}^{M^{(x)}} c_i^{(x)} p(\mathbf{o}_j|l, i, \mathcal{M})} \quad (34)$$

where $p(\mathbf{o}_j|m, n, \mathcal{M})$ varies according to the choice of the scheme. In the second case of the figure it can be written as

$$p(\mathbf{o}_j|m, n, \mathcal{M}) = \mathcal{N}(\mathbf{o}_j; \mathbf{C}_m \boldsymbol{\mu}_n^{(x)} + \boldsymbol{\mu}_m^{(o)}, \mathbf{C}_m \boldsymbol{\Sigma}_n^{(x)} \mathbf{C}_m + \boldsymbol{\Sigma}_m^{(o)}) \quad (35)$$

The estimation of the first and second-order component dependent expected state statistics, $\hat{\mathbf{x}}_{mn}(j) = E\{\mathbf{x}|\mathbf{o}_j, m, n, \mathcal{M}\}$ and $\hat{\mathbf{R}}_{mn}(j) = E\{\mathbf{x}\mathbf{x}'|\mathbf{o}_j, m, n, \mathcal{M}\}$, has to be modified accordingly.

The re-estimation of the model parameters in the third case of mixture of factor analysers is the same as in case of a standard factor analyser in Eqs. 29 and 30 apart from the inclusion of the mixture posteriors given above. The M steps for the other mixture assumptions are described in the following two sections. The estimation of the new component priors, $\hat{c}_m^{(o)}$ and $\hat{c}_n^{(x)}$, can be carried out as follows

$$\hat{c}_m^{(o)} = \frac{1}{N} \sum_{j=1}^N \gamma_m^{(o)}(j) \quad (36)$$

$$\hat{c}_n^{(x)} = \frac{1}{N} \sum_{j=1}^N \gamma_n^{(x)}(j) \quad (37)$$

The derivation for this is presented in Appendix C.

3.3 Mixtures of Factors

In this section a mixture of factor analysers with a mixture of Gaussians as the state process is presented. Similar approach is taken in the independent factor analysis (IFA) [1] apart from the independence assumption. The main difference is that IFA is based on factorial mixture of Gaussians where the elements of the state vector (factors) are assumed independent and generated individually by univariate mixtures of Gaussians. Bayesian networks corresponding to a factor analyser with mixtures of factors and IFA are depicted in Figure 8.

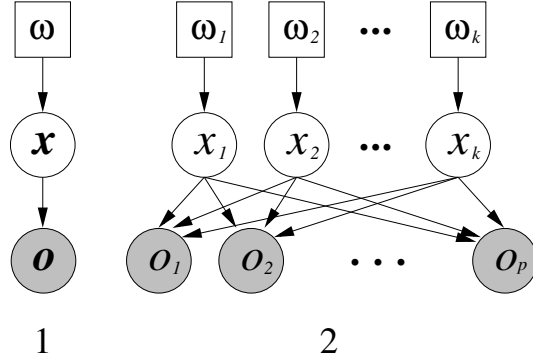


Figure 8: Bayesian networks representing factor analysers with mixtures of factors. 1) One mixture indicator chooses the factors (state vector). 2) In IFA, every factor has its own mixture indicator.

It should be noted that only $(M - 1)(2k + 1)$ additional parameters are effectively used in this model due to the degeneracy mentioned with the standard factor analysis model. The mean vector and the covariance matrix of one state noise component can always be subsumed to the observation noise mean vector and the observation matrix, respectively. Also, only $M - 1$ parameters are needed to define the component priors.

Using a mixture of Gaussians as the state evolution noise makes the E step a bit more complicated. The conditioning of multivariate Gaussians presented in Appendix A.2 becomes useful. Detailed derivation of the E step is presented in Appendix C. The sufficient statistics of the state posterior given the observation and the mixture component are

$$\hat{\mathbf{x}}_n^{(x)}(j) = \boldsymbol{\mu}_n^{(x)} + \mathbf{K}_n(\mathbf{o}_j - \mathbf{C}\boldsymbol{\mu}_n^{(x)} - \boldsymbol{\mu}^{(o)}) \quad (38)$$

$$\hat{\mathbf{R}}_n^{(x)}(j) = \boldsymbol{\Sigma}_n^{(x)} - \mathbf{K}_n \mathbf{C} \boldsymbol{\Sigma}_n^{(x)} + \hat{\mathbf{x}}_n^{(x)}(j) \hat{\mathbf{x}}_n^{(x)}(j)' \quad (39)$$

where $\mathbf{K}_n = \boldsymbol{\Sigma}_n^{(x)} \mathbf{C}' (\mathbf{C} \boldsymbol{\Sigma}_n^{(x)} \mathbf{C}' + \boldsymbol{\Sigma}^{(o)})^{-1}$. Since all the observation parameters are tied over the mixture indicator, the state posteriors given the observation are needed. They can be obtained easily as follows

$$\hat{\mathbf{x}}(j) = \sum_{n=1}^M \gamma_n^{(x)}(j) \hat{\mathbf{x}}_n^{(x)}(j) \quad (40)$$

$$\hat{\mathbf{R}}(j) = \sum_{n=1}^M \gamma_n^{(x)}(j) \hat{\mathbf{R}}_n^{(x)}(j) \quad (41)$$

where $\gamma_n^{(x)}(j)$ corresponds to the posterior of the state-space mixture component. These statistics can be used in the re-estimation of the observation parameters as in normal factor analysis case; see Eqs. 29 and 30.

The re-estimation formulae for the state parameters are

$$\hat{\boldsymbol{\mu}}_n^{(x)} = \frac{\sum_{j=1}^N \gamma_n^{(x)}(j) \hat{\boldsymbol{x}}_n^{(x)}(j)}{\sum_{j=1}^N \gamma_n^{(x)}(j)} \quad (42)$$

$$\hat{\boldsymbol{\Sigma}}_n^{(x)} = \text{diag} \left(\frac{\sum_{j=1}^N \gamma_n^{(x)}(j) \hat{\boldsymbol{R}}_n^{(x)}(j)}{\sum_{j=1}^N \gamma_n^{(x)}(j)} - \hat{\boldsymbol{\mu}}_n^{(x)} \hat{\boldsymbol{\mu}}_n^{(x)'} \right) \quad (43)$$

assuming all the model parameters are updated in the same run. If only covariance matrices are to be re-estimated, Eq. 43 will be more complicated since the substitution of the new mean vector simplifies the formula as seen in Appendix C.

The update formulae for independent factor analysis can be derived in the same way. The advantage of IFA is to have effectively more different state distributions with fewer underlying Gaussian mixture components. The usefulness of IFA depends on the validity of the independence assumption which is not generally true in case of speech feature vectors.

3.4 Mixture of Observation Processes

Mixturing the observation process increases the number of model parameters significantly and therefore tying becomes important. In [15] the noise covariance matrices were chosen to be tied so that the state to observation space transforms are distinct to each mixture component. This model was called mixture of factor analysers (MFA) but it should not be mixed with the generic mixture model in this paper. In [22] the observation matrices were tied and it is called shared factor analysis (SFA) in this paper. Both the models can be illustrated by the same Bayesian network in Figure 9.

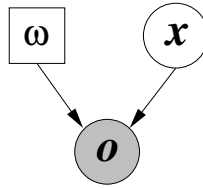


Figure 9: Bayesian network representing factor analyser with mixture of observation processes.

The MFA model provides more freedom in choosing the number of factors, k , in each analyser but the number of model parameters is much higher since M full \mathbf{C}_m matrices have to be estimated. In SFA framework the observation noise is modelled as a mixture of Gaussians and the observation matrix is shared so that less model parameters ($M2p$ additional) have to be estimated due to the diagonal covariance matrix assumption.

The parameters for the mixture of factor analysers can be obtained after finding the component and state posteriors given the observation and the mixture component. The derivation for the M step can be found in [15]. The re-estimation formulae for individual parameters of the MFA model can be represented as follows

$$\left[\hat{\mathbf{C}}_m \hat{\boldsymbol{\mu}}_m^{(o)} \right] = \left(\sum_{j=1}^N \gamma_m^{(o)}(j) \left[\mathbf{o}_j \hat{\boldsymbol{x}}_m^{(o)'}(j) \ \mathbf{o}_j \right] \right) \left(\sum_{j=1}^N \gamma_m^{(o)}(j) \begin{bmatrix} \hat{\mathbf{R}}_m^{(o)}(j) & \hat{\boldsymbol{x}}_m^{(o)}(j) \\ \hat{\boldsymbol{x}}_m^{(o)'}(j) & 1 \end{bmatrix} \right)^{-1} \quad (44)$$

The derivation follows very closely to that of the normal factor analysis since the observation matrix is shared and cancels nicely during the re-estimation of the other parameters. The re-estimation of the shared observation noise covariance matrix is simple and can be represented as follows

$$\hat{\Sigma}^{(o)} = \frac{1}{N} \sum_{j=1}^N \sum_{m=1}^M \gamma_m^{(o)}(j) \text{diag} \left(\mathbf{o}_j \mathbf{o}_j' - \left[\hat{\mathbf{C}}_m \hat{\boldsymbol{\mu}}_m^{(o)} \right] \left[\mathbf{o}_j \hat{\mathbf{x}}_m^{(o)'}(j) \mathbf{o}_j \right]' \right) \quad (45)$$

This is easy to verify since the re-estimation formulae for the other parameters are consistent with the standard factor analysis framework.

Optimisation of the model parameters for the shared factor analysis is more difficult since the re-estimation formulae for any of the model parameters cannot be represented using only the sufficient statistics. Therefore the derivation follows generalised EM algorithm where in addition to the sufficient statistics, the old observation noise parameters are used in updating the observation matrix. This optimisation strategy is the same as used for the optimisation of maximum likelihood linear regression [30] transform matrix in HMM based speaker adaptation. To make the optimisation one row at a time easy, diagonal observation noise covariance matrices have to be used.

To be able to stick to the notation, let $\hat{\mathbf{B}} = \hat{\mathbf{C}}'$. The rows of the observation matrix can now be optimised through the columns of the $\hat{\mathbf{B}}$ matrix. In order to do that, the following statistics have to be defined

$$\mathbf{G}_l = \sum_{m=1}^M \frac{1}{\sigma_{ml}^{(o)2}} \sum_{j=1}^N \gamma_m^{(o)}(j) \hat{\mathbf{R}}_m^{(o)}(j) \quad (46)$$

$$\mathbf{k}_l = \sum_{m=1}^M \frac{1}{\sigma_{ml}^{(o)2}} \sum_{j=1}^N \gamma_m^{(o)}(j) (o_{jl} - \mu_{ml}^{(o)}) \hat{\mathbf{x}}_m^{(o)}(j) \quad (47)$$

where $\sigma_{ml}^{(o)2}$ is the l th diagonal element of the observation noise covariance matrix $\Sigma_m^{(o)}$, o_{jl} is the l th element of the observation vector \mathbf{o}_j and $\mu_{ml}^{(o)}$ is the l th element of the observation noise mean vector $\boldsymbol{\mu}_m^{(o)}$. The columns of the $\hat{\mathbf{B}}$ matrix can be obtained as follows

$$\hat{\mathbf{b}}_l = \mathbf{G}_l^{-1} \mathbf{k}_l \quad (48)$$

over all the p rows and the new observation matrix is $\hat{\mathbf{C}} = [\hat{\mathbf{b}}_1 \dots \hat{\mathbf{b}}_p]'$.

The observation noise parameters can be re-estimated in the usual way apart from the simplified formula for the observation noise covariance matrix. The individual component parameters can be re-estimated as follows

$$\hat{\boldsymbol{\mu}}_m^{(o)} = \frac{\sum_{j=1}^N \gamma_m^{(o)}(j) (\mathbf{o}_j - \hat{\mathbf{C}} \hat{\mathbf{x}}_m^{(o)}(j))}{\sum_{j=1}^N \gamma_m^{(o)}(j)} \quad (49)$$

$$\begin{aligned} \hat{\Sigma}_m^{(o)} = & \frac{1}{N} \sum_{j=1}^N \gamma_m^{(o)}(j) \text{diag} \left(\mathbf{o}_j \mathbf{o}_j' - \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right] \left[\mathbf{o}_j \hat{\mathbf{x}}_m^{(o)'}(j) \mathbf{o}_j \right]' \right. \\ & \left. - \left[\mathbf{o}_j \hat{\mathbf{x}}_m^{(o)'}(j) \mathbf{o}_j \right] \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right]' + \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right] \begin{bmatrix} \hat{\mathbf{R}}_m^{(o)}(j) & \hat{\mathbf{x}}_m^{(o)}(j) \\ \hat{\mathbf{x}}_m^{(o)'}(j) & 1 \end{bmatrix} \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right]' \right) \quad (50) \end{aligned}$$

3.5 Decoding Cost

When using any of the factor analysis models presented above, the likelihood of the observation being generated by a given model can be obtained by using Eq. 35. In case of mixture of factor analysers it has to be weighted by the corresponding mixture priors and the resulting joint likelihoods of the observation and the mixture components have to be summed over the components. The majority of the cost involved in obtaining the likelihood is due to inverting the corresponding full covariance matrices in the observation space, $\mathbf{C}_m \boldsymbol{\Sigma}_n^{(x)} \mathbf{C}_m + \boldsymbol{\Sigma}_m^{(o)}$. Instead of inverting a p by p matrix there is a useful result from matrix algebra that converts the problem to inverting a k by k matrix as discussed later in Section 7. Furthermore, it should be noted that all the full inverse covariance matrices in the observation space can be pre-computed before starting off with the decoding.

4 Linear Dynamical Systems

Linear dynamical systems (LDS) are the simplest dynamical models with continuous state vectors. The state evolution process is a linear first-order Gauss-Markov random process and the observation process is a factor analyser. The benefits of using linear dynamical systems in speech recognition are the smooth trajectory modelling, improved covariance modelling and it allows subspace modelling [34].

Linear dynamical systems are based on continuous hidden k dimensional state vectors, \mathbf{x}_t , which evolve according to first-order Markov dynamics. A p dimensional observation vector, \mathbf{o}_t , is generated from the current state by a linear observation process. LDS can be described by the following two equations

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \quad (51)$$

$$\mathbf{o}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (52)$$

where \mathbf{A} is a k by k state transition matrix and \mathbf{C} is a p by k observation matrix. Vectors \mathbf{w} and \mathbf{v} are called the state evolution noise and observation noise respectively, and they are independent of each other and of the values \mathbf{x}_t and \mathbf{o}_t . Both of these noise sources are Gaussian distributed and temporally uncorrelated (white) processes; therefore, \mathbf{x}_t is a first-order Gauss-Markov random vector process. It should be noted that the observation process is a dynamic version of the factor analysis. The model parameters are $\mathcal{M} = \{\boldsymbol{\mu}^{(i)}, \boldsymbol{\Sigma}^{(i)}, \mathbf{A}, \boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}, \mathbf{C}, \boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}\}$. Traditionally, linear dynamical systems do not contain noise mean vectors but they have been included in the generic notation. On top of being more general, mean vectors prove to be essential when dealing with mixture models as already seen in the case of the factor analysis.

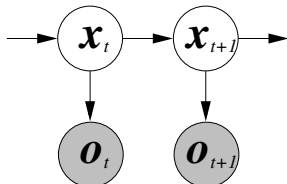


Figure 10: Dynamic Bayesian network representing a linear dynamical system.

The generative model above can be represented as a dynamic Bayesian network depicted in Figure 10. Since the noise sources are Gaussian distributed, \mathbf{x}_t and \mathbf{o}_t are Gaussian distributed and their conditional densities are

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \mathbf{A}\mathbf{x}_{t-1} + \boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \quad (53)$$

$$p(\mathbf{o}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}\mathbf{x}_t + \boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (54)$$

which is consistent with the conditional independencies illustrated in the figure. By assuming a Gaussian initial state density, $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}^{(i)}, \boldsymbol{\Sigma}^{(i)})$, and using the Markov property, the joint density of an observation sequence, $\mathbf{O} = \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T$, and state sequence, $\mathbf{X} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$, can be represented as

$$p(\mathbf{X}, \mathbf{O}) = p(\mathbf{x}_1) \prod_{t=2}^T p(\mathbf{x}_t | \mathbf{x}_{t-1}) \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{x}_t) \quad (55)$$

The number of model parameters is large in the standard LDS and reliable estimation of the model parameters becomes hard. Noting the same degeneracy regarding the state noise parameters as with the standard factor analysis, the effective number of parameters is $p + p(p+1)/2 + pk + k^2 + k + k(k+1)/2$ where the first two terms correspond to the observation noise statistics, the next two to the observation matrix and the state evolution matrix, respectively, and the last two to the initial state statistics. The generic model can be modified by introducing several restrictions to the model such as using diagonal noise covariance matrices and tying parameters especially when using mixture models. These restrictions and mixture models are reviewed after the parameter estimation has been presented.

4.1 EM Algorithm for LDS

Learning the parameters (system identification) of a linear dynamical system given only a sequence of observations can be carried out using again the powerful expectation maximisation algorithm [7]. The auxiliary function to be maximised can be represented as follows

$$\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) = E\{\log p(\mathbf{X}, \mathbf{O} | \hat{\mathcal{M}}) | \mathbf{O}, \mathcal{M}\} = \int p(\mathbf{X} | \mathbf{O}, \mathcal{M}) \log p(\mathbf{X}, \mathbf{O} | \hat{\mathcal{M}}) d\mathbf{X} \quad (56)$$

Due to the Markov assumption in the state dynamics, the observations are dependent on the respective states and the state posteriors become dependent on the entire observation sequence. Therefore, it is essential in estimating the sufficient statistics for the M step to take the entire observation sequence into account. This must be done in two passes by first estimating the state posteriors given the observations up to the current time index; i.e., $E\{\mathbf{x}_t | \mathbf{o}_1, \dots, \mathbf{o}_t\}$ and $E\{\mathbf{x}_t \mathbf{x}_t' | \mathbf{o}_1, \dots, \mathbf{o}_t\}$. Using these statistics, the state posteriors given the entire observation sequence can be obtained; i.e., $E\{\mathbf{x}_t | \mathbf{O}\}$ and $E\{\mathbf{x}_t \mathbf{x}_t' | \mathbf{O}\}$. Since both of these posteriors are Gaussians, it is rather easy to derive the recursions for the estimates using matrix algebra and the generalised forward-backward algorithm as described in Appendix E.1.

By defining $\mathbf{x}^{(\tau)}(t) = E\{\mathbf{x}_t | \mathbf{o}_1, \dots, \mathbf{o}_\tau\}$, $\mathbf{R}^{(\tau)}(t) = E\{\mathbf{x}_t \mathbf{x}_t' | \mathbf{o}_1, \dots, \mathbf{o}_\tau\}$ and $\boldsymbol{\Sigma}^{(\tau)}(t) = \mathbf{R}^{(\tau)}(t) - \mathbf{x}^{(\tau)}(t) \mathbf{x}^{(\tau)'}(t)$ the forward recursion can be represented as follows

$$\mathbf{x}^{(t-1)}(t) = \mathbf{A} \mathbf{x}^{(t-1)}(t-1) + \boldsymbol{\mu}^{(x)} \quad (57)$$

$$\boldsymbol{\Sigma}^{(t-1)}(t) = \mathbf{A} \boldsymbol{\Sigma}^{(t-1)}(t-1) \mathbf{A}' + \boldsymbol{\Sigma}^{(x)} \quad (58)$$

$$\boldsymbol{\Sigma}^{(e)}(t) = \mathbf{C} \boldsymbol{\Sigma}^{(t-1)}(t) \mathbf{C}' + \boldsymbol{\Sigma}^{(o)} \quad (59)$$

$$\mathbf{K}(t) = \boldsymbol{\Sigma}^{(t-1)}(t) \mathbf{C}' \boldsymbol{\Sigma}^{(e)-1}(t) \quad (60)$$

$$\mathbf{e}(t) = \mathbf{o}_t - \mathbf{C} \mathbf{x}^{(t-1)}(t) - \boldsymbol{\mu}^{(o)} \quad (61)$$

$$\mathbf{x}^{(t)}(t) = \mathbf{x}^{(t-1)}(t) + \mathbf{K}(t) \mathbf{e}(t) \quad (62)$$

$$\boldsymbol{\Sigma}^{(t)}(t) = \boldsymbol{\Sigma}^{(t-1)}(t) - \mathbf{K}(t) \mathbf{C} \boldsymbol{\Sigma}^{(t-1)}(t) \quad (63)$$

where $\mathbf{x}^{(0)}(1) = \boldsymbol{\mu}^{(i)}$ and $\boldsymbol{\Sigma}^{(0)}(1) = \boldsymbol{\Sigma}^{(i)}$. The matrix $\mathbf{K}(t)$ is traditionally called the Kalman gain matrix and the vector $\mathbf{e}(t)$ is called the prediction error or the innovation vector [27, 28]. The innovation vector contains the “new” information that cannot be predicted, hence the name.

The generic scaled forward algorithm allows the likelihood of the observation sequence to be obtained via the scaling factors. The scaling factors which represent the posterior of the current observation given all the previous observations, $\kappa_t = p(\mathbf{o}_t | \mathbf{o}_1, \dots, \mathbf{o}_{t-1})$, can be obtained by

$$\kappa_t = \mathcal{N}(\mathbf{e}(t); \mathbf{0}, \boldsymbol{\Sigma}^{(e)}(t)) \quad (64)$$

and the joint likelihood of an observation sequence is simply $p(\mathbf{O}) = \prod_{t=1}^T \kappa_t$.

The parallel backward recursion, also known as Kalman or Rauch-Tung-Streibel smoother [36, 37], is derived in Appendix E.2. By defining $\hat{\mathbf{x}}(t) = E\{\mathbf{x}_t|\mathbf{O}\}$ and $\hat{\mathbf{R}}(t) = E\{\mathbf{x}_t\mathbf{x}_t'|\mathbf{O}\}$ the estimates of the required statistics can now be initialised as $\hat{\mathbf{x}}(T) = \mathbf{x}^{(T)}(T)$, $\hat{\mathbf{\Sigma}}(T) = \mathbf{\Sigma}^{(T)}(T)$ and $\hat{\mathbf{R}}(T) = \hat{\mathbf{\Sigma}}(T) + \hat{\mathbf{x}}(T)\hat{\mathbf{x}}'(T)$. The rest of the estimates can be obtained by using the following Kalman smoother recursions,

$$\mathbf{J}(t-1) = \mathbf{\Sigma}^{(t-1)}(t-1)\mathbf{A}'(\mathbf{\Sigma}^{(t-1)}(t))^{-1} \quad (65)$$

$$\hat{\mathbf{x}}(t-1) = \mathbf{x}^{(t-1)}(t-1) + \mathbf{J}(t-1)(\hat{\mathbf{x}}(t) - \mathbf{x}^{(t-1)}(t)) \quad (66)$$

$$\hat{\mathbf{\Sigma}}(t-1) = \mathbf{\Sigma}^{(t-1)}(t-1) + \mathbf{J}(t-1)(\hat{\mathbf{\Sigma}}(t) - \mathbf{\Sigma}^{(t-1)}(t))\mathbf{J}'(t-1) \quad (67)$$

$$\hat{\mathbf{\Sigma}}^{(t-1)}(t) = \hat{\mathbf{\Sigma}}(t)\mathbf{J}'(t-1) \quad (68)$$

for time indices $t = T, \dots, 2$. Then $\hat{\mathbf{R}}(t) = \hat{\mathbf{\Sigma}}(t) + \hat{\mathbf{x}}(t)\hat{\mathbf{x}}'(t)$ and $\hat{\mathbf{R}}^{(t-1)}(t) = \hat{\mathbf{\Sigma}}^{(t-1)}(t) + \hat{\mathbf{x}}(t)\hat{\mathbf{x}}'(t-1)$. This recursion is considerably more efficient than the traditional recursion presented in [38] since the estimate of the cross covariance of two consecutive state vectors, $\hat{\mathbf{\Sigma}}^{(t-1)}(t)$, can be obtained via one matrix multiplication in contrast to four in the traditional method. This simplification follows directly from the joint Gaussianity of the state posteriors [32].

The derivation of the re-estimation formulae for linear dynamical systems is outlined in Appendix E.3. Same re-estimation formulae can be found in [34] and can be nicely expressed by the following matrix operations

$$\left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] = \left(\sum_{t=1}^T \begin{bmatrix} \mathbf{o}_t \hat{\mathbf{x}}'(t) & \mathbf{o}_t \end{bmatrix} \right) \left(\sum_{t=1}^T \begin{bmatrix} \hat{\mathbf{R}}(t) & \hat{\mathbf{x}}(t) \\ \hat{\mathbf{x}}'(t) & 1 \end{bmatrix} \right)^{-1} \quad (69)$$

$$\hat{\mathbf{\Sigma}}^{(o)} = \frac{1}{T} \sum_{t=1}^T \left(\mathbf{o}_t \mathbf{o}_t' - \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] \begin{bmatrix} \mathbf{o}_t \hat{\mathbf{x}}'(t) & \mathbf{o}_t \end{bmatrix}' \right) \quad (70)$$

$$\left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] = \left(\sum_{t=2}^T \begin{bmatrix} \hat{\mathbf{R}}^{(t-1)}(t) & \hat{\mathbf{x}}(t) \end{bmatrix} \right) \left(\sum_{t=2}^T \begin{bmatrix} \hat{\mathbf{R}}(t-1) & \hat{\mathbf{x}}(t-1) \\ \hat{\mathbf{x}}'(t-1) & 1 \end{bmatrix} \right)^{-1} \quad (71)$$

$$\hat{\mathbf{\Sigma}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \left(\hat{\mathbf{R}}(t) - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \begin{bmatrix} \hat{\mathbf{R}}^{(t-1)}(t) & \hat{\mathbf{x}}(t) \end{bmatrix}' \right) \quad (72)$$

$$\hat{\boldsymbol{\mu}}^{(i)} = \hat{\mathbf{x}}(1) \quad (73)$$

$$\hat{\mathbf{\Sigma}}^{(i)} = \hat{\mathbf{R}}(1) - \hat{\boldsymbol{\mu}}^{(i)} \hat{\boldsymbol{\mu}}^{(i)'} \quad (74)$$

where the noise covariance matrices are full.

4.2 Diagonal Covariance and State Transition Matrices

To reduce the number of model parameters and to guarantee that the state transition matrix and observation matrix model the covariance structure of the data, it is possible to restrict the noise source covariance matrices to be diagonal. The only modifications to the re-estimation formulae are just taking the main diagonal elements of the re-estimated matrices as follows

$$\hat{\mathbf{\Sigma}}^{(o)} = \frac{1}{T} \sum_{t=1}^T \text{diag} \left(\mathbf{o}_t \mathbf{o}_t' - \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] \begin{bmatrix} \mathbf{o}_t \hat{\mathbf{x}}'(t) & \mathbf{o}_t \end{bmatrix}' \right) \quad (75)$$

$$\hat{\mathbf{\Sigma}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \text{diag} \left(\hat{\mathbf{R}}(t) - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \begin{bmatrix} \hat{\mathbf{R}}^{(t-1)}(t) & \hat{\mathbf{x}}(t) \end{bmatrix}' \right) \quad (76)$$

$$\hat{\mathbf{\Sigma}}^{(i)} = \text{diag} \left(\hat{\mathbf{R}}(1) - \hat{\boldsymbol{\mu}}^{(i)} \hat{\boldsymbol{\mu}}^{(i)'} \right) \quad (77)$$

In practise, it is, of course, not necessary to calculate the full covariance matrices but just the diagonal elements. The number of parameters for a diagonal covariance LDS is $2p + pk + k^2 + 2k$.

It is possible to further reduce the number of model parameters by assuming diagonal state transition matrix which corresponds to an independent state trajectory assumption provided that the covariance matrices are diagonal as well. The Kalman filter and smoother recursions do not need any modifications since they give the maximum likelihood estimates for the required statistics regardless the structure of the matrices. The covariance matrices obtained by Kalman smoother will be full matrices even if the state transition matrix is diagonal due to the full observation matrix. Therefore the re-estimation formula for the state transition matrix without state evolution noise mean vector has to be modified as follows

$$\hat{\mathbf{A}} = \left(\text{diag} \left(\sum_{t=2}^T \hat{\mathbf{R}}^{(t-1)}(t) \right) \right) \left(\text{diag} \left(\sum_{t=2}^T \hat{\mathbf{R}}(t-1) \right) \right)^{-1} \quad (78)$$

Unfortunately, the state evolution noise mean vector cannot be re-estimated simultaneously with the diagonal state transition matrix. The simultaneous re-estimation would require major changes in the notation and an inefficient form of Eq. 71 should be used. The re-estimation has to be done separately as follows

$$\hat{\mathbf{A}} = \left(\text{diag} \left(\sum_{t=2}^T \hat{\mathbf{R}}^{(t-1)}(t) - \frac{1}{T-1} \sum_{t=2}^T \hat{\mathbf{x}}(t) \sum_{t=2}^T \hat{\mathbf{x}}'(t-1) \right) \right) \left(\text{diag} \left(\sum_{t=2}^T \hat{\mathbf{R}}(t-1) - \frac{1}{T-1} \sum_{t=2}^T \hat{\mathbf{x}}(t-1) \sum_{t=2}^T \hat{\mathbf{x}}'(t-1) \right) \right)^{-1} \quad (79)$$

$$\hat{\boldsymbol{\mu}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \left(\hat{\mathbf{x}}(t) - \hat{\mathbf{A}} \hat{\mathbf{x}}(t-1) \right) \quad (80)$$

It is easy to verify that Eq. 76 can be used to re-estimate the state evolution noise covariance if diagonal covariance matrices are used as well. In case of full covariance matrices, all the terms in Eq. 255 need to be used.

4.3 Mixture of Linear Dynamical Systems

In linear dynamical systems, the observation vectors are also Gaussian distributed. A Gaussian density is unimodal whereas often the feature vectors are multi-modal as noted already. Therefore, mixture models are more interesting. There are three distinct noise models present in the linear dynamical systems; the initial state, the state evolution noise and the observation noise. If mixtures were used in the state evolution noise, there would be exponential growth in the number of observation components. Propagating M mixture components through the state evolution equations leads to M^2 components after first time instant and in the end of an observation sequence there are M^T components. However, there are several approximate approaches for dealing with this exponential growth as discussed in [33]. Most of the approximate methods work well only in the filtering and the smoothing is more prone to errors. Similar exponential growth occurs in the state posteriors if mixture models are used as the observation noise or the initial state distribution.

The most generic form of a mixture of linear dynamical systems is depicted on the left-hand side of Figure 11. The observation matrix can be distributed according to a mixture model like in the case of a mixture of factor analysers in Section 3.2. The same applies for the state transition matrix for which a mixture indicator ω_t^a is used. The generative model for this model can be written as

$$\begin{aligned}
\mathbf{x}_{t+1} &= \mathbf{A}\mathbf{x}_t + \mathbf{w}, & \mathbf{w} &\sim \sum_n c_n^{(x)} \mathcal{N}(\boldsymbol{\mu}_n^{(x)}, \boldsymbol{\Sigma}_n^{(x)}) \\
& & \mathbf{A} &\sim \sum_i c_i^{(a)} \delta(\mathbf{A} - \mathbf{A}_i) & (81) \\
\mathbf{o}_t &= \mathbf{C}\mathbf{x}_t + \mathbf{v}, & \mathbf{v} &\sim \sum_m c_m^{(o)} \mathcal{N}(\boldsymbol{\mu}_m^{(o)}, \boldsymbol{\Sigma}_m^{(o)}) \\
& & \mathbf{C} &\sim \sum_l c_l^{(c)} \delta(\mathbf{C} - \mathbf{C}_l) & (82)
\end{aligned}$$

As discussed above, the filtering and smoothing are not tractable with this model. Nevertheless, this model is very interesting and some approximate algorithms are mentioned below. Another interesting intractable model is illustrated on the right-hand side of Figure 11. It is often called switching Kalman filter model [33] and it combines a hidden Markov model and a LDS. Generally, the combinations are called switching state-space models [17] since the HMM states can be viewed as switching variables which choose the active parameters of the LDS. All the parameters of switching models are time dependent via the discrete state variable q_t .

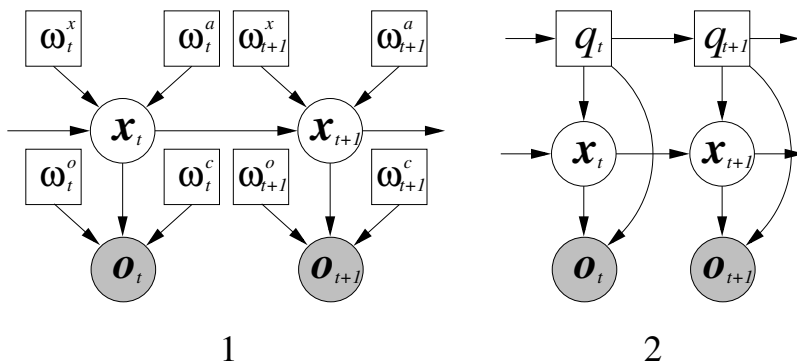


Figure 11: Dynamic Bayesian networks representing mixtures of linear dynamical systems. 1) Generic mixture of LDSs. 2) Switching Kalman filter.

There is only one way to restrict the model to be tractable but still employ a kind of mixture distribution. If the HMM states are restricted to choose the parameters from M sets and the consecutive states are forced to be equal through the entire observation sequence, the learning becomes tractable. This can be viewed as M independent LDSs working in parallel with prior likelihoods $c_m = P_\omega(m)$. Since the number of parameters grows with the factor of M , tying becomes essential to obtain reliable estimates in the re-estimation. If all parameters other than the observation noise are tied, the observation noise is crudely a mixture model. The difference to the standard mixture model is due to the mixture indicator having the same value throughout the observation sequence and therefore not permitting as flexible a model.

To optimise the parameters of mixture of M LDSs with component likelihoods, $c_m = P_\omega(m)$, the following auxiliary function has to be maximised

$$\begin{aligned}
\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) &= & (83) \\
E\{\log p(\mathbf{O}, \mathbf{X} | \hat{\mathcal{M}} | \mathbf{O}, \mathcal{M})\} &= \sum_{m=1}^M P_\omega(m | \mathbf{O}, \mathcal{M}) E\{\log p(\mathbf{O}, \mathbf{X} | m, \hat{\mathcal{M}} | \mathbf{O}, m, \mathcal{M})\}
\end{aligned}$$

The E step requires M passes of Kalman filter and smoother, and the likelihoods of the observation

sequence given the model, $p(\mathbf{O}|\mathcal{M}^m)$, are obtained. The component posteriors are then

$$\gamma_m = P_\omega(m|\mathbf{O}, \mathcal{M}^1, \dots, \mathcal{M}^M) = \frac{c_m p(\mathbf{O}|\mathcal{M}^m)}{\sum_{l=1}^M c_l p(\mathbf{O}|\mathcal{M}^l)} \quad (84)$$

The observation process is very similar to the shared factor analysis discussed in Section 3.2 apart from the component posteriors which are conditioned by the entire observation sequence. The re-estimation formula for the observation matrix can be derived in the usual fashion by first defining $\hat{\mathbf{B}} = \hat{\mathbf{C}}'$ and the following statistics

$$\mathbf{G}_l = \sum_{m=1}^M \frac{1}{\sigma_{ml}^{(o)2}} \gamma_m \sum_{t=1}^T \hat{\mathbf{R}}_m(t) \quad (85)$$

$$\mathbf{k}_l = \sum_{m=1}^M \frac{1}{\sigma_{ml}^{(o)2}} \gamma_m \sum_{t=1}^T (o_{tl} - \mu_{ml}^{(o)}) \hat{\mathbf{x}}_m(t) \quad (86)$$

where $\sigma_{ml}^{(o)2}$ is the l th diagonal element of the observation noise covariance matrix $\Sigma_m^{(o)}$, o_{tl} is the l th element of the observation vector \mathbf{o}_t and $\mu_{ml}^{(o)}$ is the l th element of the observation noise mean vector $\boldsymbol{\mu}_m^{(o)}$. Then the columns of the $\hat{\mathbf{B}}$ matrix can be obtained as follows

$$\hat{\mathbf{b}}_l = \mathbf{G}_l^{-1} \mathbf{k}_l \quad (87)$$

over all the p rows and the new observation matrix is $\hat{\mathbf{C}} = [\hat{\mathbf{b}}_1 \dots \hat{\mathbf{b}}_p]'$. The component posteriors are also the maximum likelihood estimates for the new component prior likelihoods, $\hat{c}_m = \gamma_m$. Since the priors are subject to constraint, $\sum_m c_m = 1$, this can be derived using Lagrange multipliers as in Appendix C.2.

The re-estimation formulae for the observation noise parameters can be obtained easily noting that the component posteriors cancel since they do not depend on the time, t . The formulae can be represented as follows

$$\hat{\boldsymbol{\mu}}_m^{(o)} = \frac{1}{T} \sum_{t=1}^T (\mathbf{o}_t - \hat{\mathbf{C}} \hat{\mathbf{x}}_m(t)) \quad (88)$$

$$\begin{aligned} \hat{\Sigma}_m^{(o)} &= \frac{1}{T} \sum_{t=1}^T \text{diag}(\mathbf{o}_t \mathbf{o}_t' - [\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)}] [\mathbf{o}_t \hat{\mathbf{x}}_m'(t) \mathbf{o}_t]') \\ &\quad - [\mathbf{o}_t \hat{\mathbf{x}}_m'(t) \mathbf{o}_t] [\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)}]' + [\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)}] \begin{bmatrix} \hat{\mathbf{R}}_m(t) & \hat{\mathbf{x}}_m(t) \\ \hat{\mathbf{x}}_m'(t) & 1 \end{bmatrix} [\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)}]' \end{aligned} \quad (89)$$

The re-estimation formulae for the state evolution process is simply obtained using the component posteriors as weights

$$\left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] = \left(\sum_{t=2}^T \sum_{m=1}^M \gamma_m \begin{bmatrix} \hat{\mathbf{R}}_m^{(t-1)}(t) & \hat{\mathbf{x}}_m(t) \end{bmatrix} \right) \left(\sum_{t=2}^T \sum_{m=1}^M \gamma_m \begin{bmatrix} \hat{\mathbf{R}}_m(t-1) & \hat{\mathbf{x}}_m(t-1) \\ \hat{\mathbf{x}}_m'(t-1) & 1 \end{bmatrix} \right)^{-1} \quad (90)$$

$$\hat{\Sigma}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \sum_{m=1}^M \gamma_m \left(\hat{\mathbf{R}}_m(t) - [\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)}] \begin{bmatrix} \hat{\mathbf{R}}_m^{(t-1)}(t) & \hat{\mathbf{x}}_m(t) \end{bmatrix}' \right) \quad (91)$$

$$\hat{\boldsymbol{\mu}}^{(i)} = \sum_{m=1}^M \gamma_m \hat{\mathbf{x}}_m(1) \quad (92)$$

$$\hat{\Sigma}^{(i)} = \sum_{m=1}^M \gamma_m \hat{\mathbf{R}}_m(1) - \hat{\boldsymbol{\mu}}^{(i)} \hat{\boldsymbol{\mu}}^{(i)'} \quad (93)$$

where the noise covariance matrices are full but can be diagonalised as usual.

4.4 Sampling Approaches for Mixture Models

As noted above, using mixture models as the noise processes of a linear dynamical system leads to exponential growth in the state posteriors. Recently, in statistical signal processing [8, 9, 19] and computer vision [25, 24] Monte Carlo sampling methods have been employed in filtering problems with non-linear non-Gaussian models. These methods are based on a “particle cloud” representation in the posterior propagation. A set of N points sampled from the prior distribution are propagated through the system and the resulting N^2 posterior samples are then re-sampled and weighted according to importance functions obtained from the model. Particle filtering can be used instead of Kalman filtering but the smoothing approaches presented in [8, 9] perform poorly in estimating $p(\mathbf{x}_t|\mathbf{O})$ when $T - t$ is significant. More stable smoothing estimates have to be studied and developed.

An alternative way to deal with the mixture models is to propagate the posterior as usual resulting into M^2 mixture components and then apply Gaussian selection [13] to project these M^2 components back to M components. This is also an approximate method but it takes into account the fact that all the posteriors are always mixtures of Gaussians. Unfortunately, it is not possible to guarantee that the state posteriors obtained this way would increase the likelihood of the data after re-estimation of the model parameters but the validity of the approximation should be evaluated experimentally.

5 Hidden Markov Models

Hidden Markov models are based on a finite N_s state machine with discrete initial state and transition probabilities, π_j and a_{ij} respectively, as well as state dependent observation distributions $b_j(\mathbf{o}_t) = p(\mathbf{o}_t|j)$. Mixtures of M Gaussians with mixture weights $c_{jm} = P_\omega(m|j)$ are used as the generative model for each state in this paper so that a HMM can be described by the set of parameters $\mathcal{M} = \{\forall j, i \in (1, N_s), m \in (1, M) : \pi_j, a_{ij}, c_{jm}, \boldsymbol{\mu}_{jm}, \boldsymbol{\Sigma}_{jm}\}$. In speech recognition the initial state and the transition probabilities are restricted to form so called left to right models but in this paper a general formulation is used.

One of the main advantages of HMMs is the efficient segmentation which is carried out using Viterbi algorithm [41]. Training of HMMs is also very efficient using the EM framework and it is often referred to as the Baum-Welch algorithm [2, 3]. The major drawbacks are discontinuous state transitions which prevent trajectory interpretation of the feature elements, the state conditional independence assumption which states that given the state, the subsequent observation vectors are independent, and rather weak covariance modelling if diagonal covariance matrices are used.

Using mixture models as the state dependent observation distributions has yielded better performance in HMM based systems especially in speech recognition. There are three major reasons for that. Firstly, mixtures can model multi-modal distributions that might result from source classes with more than one distinct acoustic characteristics. Secondly, mixtures can model spatial correlation in the feature vectors and therefore several distributions with diagonal covariance matrices can approximate one distribution with full covariance matrix with fewer parameters. Thirdly, mixtures can be used to model non-symmetric distributions. In the limit Gaussian mixture models can represent arbitrary distributions that fulfil certain continuity requirements.

The p dimensional observation vector likelihood $b_j(\mathbf{o}_t)$ is defined as the following mixture of Gaussians

$$b_j(\mathbf{o}_t) = p(\mathbf{o}_t|j) = \sum_{m=1}^M c_{jm} \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{jm}, \boldsymbol{\Sigma}_{jm}) \quad (94)$$

$$(95)$$

where $\boldsymbol{\mu}_{jm}$ and $\boldsymbol{\Sigma}_{jm}$ are the mean vector and covariance matrix associated with the state $q_t = j$ and the mixture component $\omega = m$. The simplest continuous single mixture HMM is illustrated as a dynamic Bayesian network on the left-hand side of Figure 12. On the right-hand side is a mixture of Gaussians HMM.

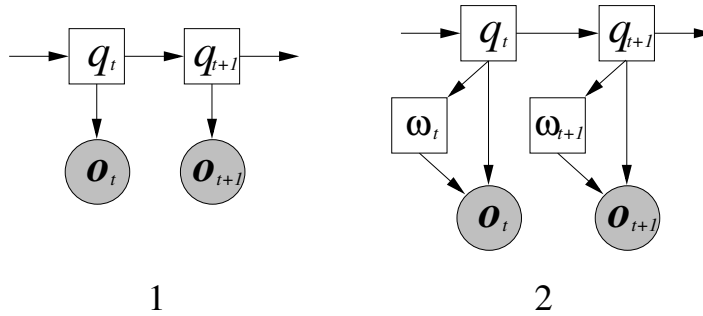


Figure 12: Bayesian networks representing hidden Markov models. 1) Single Gaussian continuous density HMM. 2) Mixture of Gaussians continuous density HMM.

The joint likelihood of an observation sequence, $\mathbf{O} = \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T$, and a HMM state sequence, $Q = q_1, q_2, \dots, q_T$, can be represented as

$$p(\mathbf{O}, Q) = P(q_1) \prod_{t=2}^T P(q_t | q_{t-1}) \prod_{t=1}^T p(\mathbf{o}_t | q_t) \quad (96)$$

where $P(q_1) = \pi_{q_1}$ are the initial state probabilities and $P(q_t | q_{t-1}) = a_{q_{t-1}q_t}$ are the state transition probabilities.

5.1 EM Algorithm for HMM

Learning the model parameters of hidden Markov models given only an observation sequence or a set of sequences can be carried out using the powerful expectation maximisation algorithm also known as Baum-Welch algorithm. The auxiliary function to be maximised can be represented as follows

$$\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) = E\{\log p(\mathbf{O}, Q | \hat{\mathcal{M}}) | \mathbf{O}, \mathcal{M}\} = \sum_{\{Q_T\}} P(Q | \mathbf{O}, \mathcal{M}) \log p(\mathbf{O}, Q | \hat{\mathcal{M}}) \quad (97)$$

Since the state posteriors depend on discrete state transition probabilities, they are not Gaussians any more. Due to the finite number of possible state sequences, efficient likelihood calculation requires the use of dynamic programming. The sufficient statistics $\gamma_j(t) = P_{q_t}(j | \mathbf{O}, \mathcal{M})$ and $\xi_{ij}(t) = P_{q_{t-1}, q_t}(i, j | \mathbf{O}, \mathcal{M})$ for the M step can be obtained using the familiar forward-backward algorithm for HMMs [35]. The forward algorithm starts off as follows

$$\alpha_j(1) = \pi_j b_j(\mathbf{o}_1) \quad (98)$$

where $\alpha_j(t) = p(j, \mathbf{o}_1, \dots, \mathbf{o}_t)$. The rest of the forward variables are obtained by the following recursion

$$\alpha_j(t) = b_j(\mathbf{o}_t) \sum_{i=1}^{N_s} a_{ij} \alpha_i(t-1) \quad (99)$$

$$(100)$$

The joint likelihood of the observation sequence can be obtained using the forward variables at time T ; i.e., $p(\mathbf{O}) = \sum_{j=1}^{N_s} \alpha_j(T)$.

The backward algorithm starts off with $\beta_i(T) = 1$ and the rest of the backward variables are obtained recursively as follows

$$\beta_i(t-1) = \sum_{j=1}^{N_s} a_{ij} b_j(\mathbf{o}_t) \beta_j(t) \quad (101)$$

where $\beta_i(t) = p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | i)$.

Now, the sufficient statistics can be obtained easily as

$$\gamma_j(t) = \frac{1}{p(\mathbf{O})} \alpha_j(t) \beta_j(t) \quad (102)$$

$$\xi_{ij}(t) = \frac{1}{p(\mathbf{O})} \alpha_i(t-1) a_{ij} b_j(\mathbf{o}_t) \beta_j(t) \quad (103)$$

To update the mixture components the mixture posteriors $\gamma_{jm}(t) = P_{q_t, \omega}(j, m | \mathbf{O}, \mathcal{M})$ are needed. Since the mixture components can be regarded as additional states [31], the joint likelihood of being in state j and mixture component m is

$$\gamma_{jm}(t) = \frac{1}{p(\mathbf{O})} c_{jm} b_{jm}(\mathbf{o}_t) \sum_{i=1}^{N_s} a_{ij} \alpha_i(t-1) \beta_j(t) \quad (104)$$

where $b_{jm}(\mathbf{o}_t)$ is the posterior of the observation \mathbf{o}_t given the state j and the mixture component m .

The traditional derivation of the Baum-Welch algorithm can be found in [35]. A derivation more consistent with the notation used in this paper can be found in Appendix F by replacing the state dependent observation posteriors properly, and using \mathbf{o}_t and $\mathbf{o}_t \mathbf{o}_t'$ instead of $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{R}}_t$, respectively. The re-estimation formulae can be represented as follows

$$\hat{\pi}_j = \frac{\gamma_j(1)}{\sum_{i=1}^{N_s} \gamma_i(1)} \quad (105)$$

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T \xi_{ij}(t)}{\sum_{t=2}^T \gamma_i(t-1)} \quad (106)$$

$$\hat{c}_{jm} = \frac{\sum_{t=1}^T \gamma_{jm}(t)}{\sum_{t=1}^T \gamma_j(t)} \quad (107)$$

$$\hat{\boldsymbol{\mu}}_{jm} = \frac{\sum_{t=1}^T \gamma_{jm}(t) \mathbf{o}_t}{\sum_{t=1}^T \gamma_{jm}(t)} \quad (108)$$

$$\hat{\boldsymbol{\Sigma}}_{jm} = \frac{\sum_{t=1}^T \gamma_{jm}(t) \mathbf{o}_t \mathbf{o}_t'}{\sum_{t=1}^T \gamma_{jm}(t)} - \hat{\boldsymbol{\mu}}_j \hat{\boldsymbol{\mu}}_j' \quad (109)$$

provided all the parameters are update in the same iteration. Otherwise, the update formula 109 for the covariance matrices has to be modified slightly.

To reduce the number of model parameters and the computational complexity, the covariance matrices can be assumed to be diagonal. The validity of the assumption has been argued since there is often some spatial correlation between the feature vectors obtained via currently popular

parameterisation methods. Instead of full covariance matrices the correlation of the feature vectors can be modelled, e.g., using semi-tied full covariance matrices [12, 11]. The number of parameters per state in a standard diagonal covariance matrix HMM in speech recognition is $(M-1)+M2p+1$. The first term corresponds to the component priors, the second term to the state distribution parameters and the last to the state transition probability of a left to right HMM.

5.2 Factor Analysed Hidden Markov Models

Factor analysis was combined with HMMs in [40] where all the covariance matrices are factored according to a standard factor analysis model presented in Section 3. This model and the corresponding mixture model can be illustrated as the dynamic Bayesian networks in Figure 13. It was argued that fewer model parameters than with full covariance systems can be used without significantly lowering the system performance. Nevertheless, the evaluation was done using simple alpha-digit and town name tasks. The number of model parameters in this case is huge since a separate observation matrix is attached to every component in the HMM system.

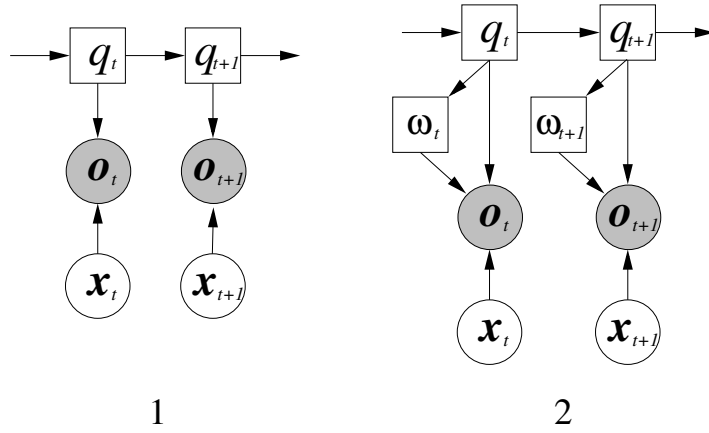


Figure 13: Bayesian networks representing hidden Markov models with factor analysed covariance matrices. 1) Single Gaussian continuous density case. 2) Mixture of Gaussians continuous density case.

Factor analysed HMMs (FAHMM) described in this paper are dynamic versions of the factor analysers with mixtures of factors and shared factor analysis presented in Sections 3.3 and 3.4. FAHMM is based on a HMM as a state evolution process and a factor analyser as an observation process. All the benefits of conventional HMMs are present and a better covariance modelling is achieved since the assumption of spatially uncorrelated feature vectors is removed. FAHMM allows also subspace modelling. The drawbacks of FAHMM include virtually the same as HMMs apart from the weak spatial covariance modelling. The generative model of FAHMM can be represented as follows

$$\begin{aligned} \mathbf{x}_t &\sim \mathcal{M}^{hmm}, & \mathcal{M}^{hmm} &= \{a_{ij}, c_{jn}^{(x)}, \boldsymbol{\mu}_{jn}^{(x)}, \boldsymbol{\Sigma}_{jn}^{(x)}\} & (110) \\ \mathbf{v}_t &\sim \sum c_{jm}^{(o)} \mathcal{N}(\boldsymbol{\mu}_{jm}^{(o)}, \boldsymbol{\Sigma}_{jm}^{(o)}) \\ \mathbf{o}_t &= \mathbf{C}_t \mathbf{x}_t + \mathbf{v}_t, & \mathbf{C}_t &\sim \sum_l c_{jl}^{(c)} \delta(\mathbf{C}_t - \mathbf{C}_{jl}) & (111) \end{aligned}$$

where the observation noise distributions and the observation matrices are state dependent. The observation matrices are delta distributed to guarantee the model to be a generalised linear Gaussian model as discussed in the case of a mixture of factor analysers in Section 3.2. Only a single

component observation matrix distribution is used later in this paper. Therefore the observation matrix prior cancels out in the formulation. The dynamic Bayesian network corresponding to this model is illustrated in Figure 14.

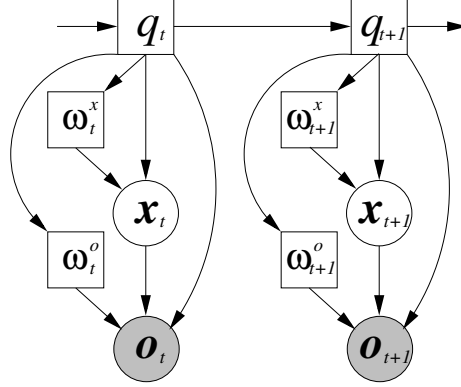


Figure 14: Bayesian network representing a generic factor analysed hidden Markov model.

Since there are now two mixtures of Gaussians, the mixture weights for the state vector distributions are denoted as $c_{jn}^{(x)} = P_{\omega^x|q_t}(n|j)$ and for the observation noise distributions as $c_{jm}^{(o)} = P_{\omega^o|q_t}(m|j)$. The model parameters of FAHMM are $\mathcal{M} = \{\forall j, i \in (1, N_s), n \in (1, M^{(x)}), m \in (1, M^{(o)}) : \pi_j, a_{ij}, c_{jn}^{(x)}, \boldsymbol{\mu}_{jn}^{(x)}, \boldsymbol{\Sigma}_{jn}^{(x)}, \mathbf{C}_j, c_{jm}^{(o)}, \boldsymbol{\mu}_{jm}^{(o)}, \boldsymbol{\Sigma}_{jm}^{(o)}\}$. Since the reliable estimation of all the parameters requires a lot of data, sharing becomes important. The general model allows arbitrary sharing of the parameters to be used in many different levels. For example the observation matrix can be tied between several states, several models or even globally.

The joint likelihood of an observation sequence, $\mathbf{O} = \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T$, state vector sequence, $\mathbf{X} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$, and HMM state sequence, $Q = q_1, q_2, \dots, q_T$, can be represented as

$$p(\mathbf{O}, \mathbf{X}, Q) = P(q_1) \prod_{t=2}^T P(q_t|q_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t|q_t)p(\mathbf{o}_t|\mathbf{x}_t) \quad (112)$$

where $P(q_1) = \pi_{q_1}$ is the initial state probability and $P(q_t|q_{t-1}) = a_{q_{t-1}q_t}$ are the state transition probabilities.

5.3 EM Algorithm for FAHMM

In the expectation maximisation algorithm for factor analysed HMM the following auxiliary function has to be maximised

$$\begin{aligned} \mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) &= E\{\log p(\mathbf{O}, \mathbf{X}, Q|\hat{\mathcal{M}})|\mathbf{O}, \mathcal{M}\} \\ &= \sum_{\{Q_T\}} \int p(\mathbf{X}|\mathbf{O}, Q, \mathcal{M})P(Q|\mathbf{O}, \mathcal{M}) \log p(\mathbf{O}, \mathbf{X}, Q|\hat{\mathcal{M}})d\mathbf{X} \end{aligned} \quad (113)$$

where the posteriors for both the state vector and the discrete HMM state are needed.

The posterior probability of being in state j at time t $\gamma_j(t) = P_{q_t}(j|\mathbf{O}, \mathcal{M})$, and being in state j at time t and in state i at time $t-1$ $\xi_{ij}(t) = P_{q_{t-1}, q_t}(i, j|\mathbf{O}, \mathcal{M})$ for the transition parameter re-estimation can be obtained using the same forward-backward algorithm as for HMMs presented in the previous section and using the following posterior of an observation

$$b_j(\mathbf{o}_t) = \sum_{m=1}^{M^{(o)}} c_{jm}^{(o)} \sum_{n=1}^{M^{(x)}} c_{jn}^{(x)} \mathcal{N}(\mathbf{o}_t; \mathbf{C}_j \boldsymbol{\mu}_{jn}^{(x)} + \boldsymbol{\mu}_{jm}^{(o)}, \mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' + \boldsymbol{\Sigma}_{jm}^{(o)}) \quad (114)$$

The posterior of being in state j , state-space component n and observation space component m at time t $\gamma_{jmn}(t)$ is needed for the re-estimation of the other parameters. This can be obtained as in the standard HMMs using the forward and backward variables as follows

$$\gamma_{jmn}(t) = \frac{1}{p(\mathbf{O})} c_{jm}^{(o)} c_{jn}^{(x)} b_{jmn}(\mathbf{o}_t) \sum_{i=1}^{N_s} a_{ij} \alpha_i(t-1) \beta_j(t) \quad (115)$$

where $b_{jmn}(\mathbf{o}_t)$ is the posterior of an observation given both the mixture components; i.e., the Gaussian in Eq. 114.

The first and second order state vector posteriors are needed in the update formulae. The posteriors given both mixture components can be obtained as follows

$$\hat{\mathbf{x}}_{jmn}(t) = \boldsymbol{\mu}_{jn}^{(x)} + \mathbf{K}_{jmn}(\mathbf{o}_t - \mathbf{C}_j \boldsymbol{\mu}_{jn}^{(x)} - \boldsymbol{\mu}_{jm}^{(o)}) \quad (116)$$

$$\hat{\mathbf{R}}_{jmn}(t) = \boldsymbol{\Sigma}_{jn}^{(x)} - \mathbf{K}_{jmn} \mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} + \hat{\mathbf{x}}_{jmn}(t) \hat{\mathbf{x}}_{jmn}'(t) \quad (117)$$

where $\mathbf{K}_{jmn} = \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' (\mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' + \boldsymbol{\Sigma}_{jm}^{(o)})^{-1}$. Now, $\hat{\mathbf{x}}_{jmn}(t)$ corresponds to the expected state vector given the state, the mixture components, the observation sequence and the old model parameters $E\{\mathbf{x}_t | j, m, n, \mathbf{O}, \mathcal{M}\}$. The similarity to the mixtures of factors in Section 3.3 should be noted.

Re-estimation formulae are derived in Appendix F. The formulae for the state parameters are very similar to the normal HMM parameter estimation and can be written as follows

$$\hat{\pi}_j = \frac{\gamma_j(1)}{\sum_{i=1}^{N_s} \gamma_i(1)} \quad (118)$$

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T \xi_{ij}(t)}{\sum_{t=2}^T \gamma_i(t-1)} \quad (119)$$

$$\hat{c}_{jn}^x = \frac{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t)}{\sum_{t=1}^T \gamma_j(t)} \quad (120)$$

$$\hat{\boldsymbol{\mu}}_{jn}^{(x)} = \frac{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) \hat{\mathbf{x}}_{jmn}(t)}{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t)} \quad (121)$$

$$\hat{\boldsymbol{\Sigma}}_{jn}^{(x)} = \text{diag} \left(\frac{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) \hat{\mathbf{R}}_{jmn}(t)}{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t)} - \hat{\boldsymbol{\mu}}_{jn}^{(x)} \hat{\boldsymbol{\mu}}_{jn}^{(x)'} \right) \quad (122)$$

The observation parameter update formulae are the dynamic version of shared factor analysis and

a new matrix $\hat{\mathbf{B}}_j = \hat{\mathbf{C}}_j'$ can be obtained defining the following statistics

$$\mathbf{G}_{jl} = \sum_{m=1}^{M^{(o)}} \frac{1}{\sigma_{jml}^{(o)2}} \sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t) \hat{\mathbf{R}}_{jmn}(t) \quad (123)$$

$$\mathbf{k}_{jl} = \sum_{m=1}^{M^{(o)}} \frac{1}{\sigma_{jml}^{(o)2}} \sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t) (o_{tl} - \mu_{jml}^{(o)}) \hat{\mathbf{x}}_{jmn}(t) \quad (124)$$

where $\sigma_{jml}^{(o)2}$ is the l th diagonal element of the observation noise covariance matrix $\Sigma_{jm}^{(o)}$, o_{tl} is the l th element of the observation vector \mathbf{o}_t and $\mu_{jml}^{(o)}$ is the l th elements of the observation noise mean vector $\boldsymbol{\mu}_{jm}^{(o)}$. The columns of the $\hat{\mathbf{B}}$ matrix can be obtained as follows

$$\hat{\mathbf{b}}_{jl} = \mathbf{G}_{jl}^{-1} \mathbf{k}_{jl} \quad (125)$$

over all the p rows and the new observation matrix is $\hat{\mathbf{C}}_j = [\hat{\mathbf{b}}_{j1} \dots \hat{\mathbf{b}}_{jp}]'$.

The re-estimation formulae for the observation noise parameters can be obtained in the usual way as follows

$$\hat{c}_{jm}^{(o)} = \frac{\sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t)}{\sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_j(t)} \quad (126)$$

$$\hat{\boldsymbol{\mu}}_{jm}^{(o)} = \frac{\sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t) (\mathbf{o}_t - \hat{\mathbf{C}}_j \hat{\mathbf{x}}_{jmn}(t))}{\sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t)} \quad (127)$$

$$\hat{\boldsymbol{\Sigma}}_{jm}^{(o)} = \frac{1}{\sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t)} \sum_{t=1}^T \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t) \text{diag} \left(\mathbf{o}_t \mathbf{o}_t' - \left[\hat{\mathbf{C}}_j \hat{\boldsymbol{\mu}}_{jm}^{(o)} \right] \left[\mathbf{o}_t \hat{\mathbf{x}}'_{jmn}(t) \ \mathbf{o}_t \right]' \right. \\ \left. - \left[\mathbf{o}_t \hat{\mathbf{x}}'_{jmn}(t) \ \mathbf{o}_t \right] \left[\hat{\mathbf{C}}_j \hat{\boldsymbol{\mu}}_{jm}^{(o)} \right]' + \left[\hat{\mathbf{C}}_j \hat{\boldsymbol{\mu}}_{jm}^{(o)} \right] \left[\begin{array}{c} \hat{\mathbf{R}}_{jmn}(t) \ \hat{\mathbf{x}}_{jmn}(t) \\ \hat{\mathbf{x}}'_{jmn}(t) \ 1 \end{array} \right] \left[\hat{\mathbf{C}}_j \hat{\boldsymbol{\mu}}_{jm}^{(o)} \right]' \right) \quad (128)$$

Now, diagonal covariance matrices have to be used to avoid a vast amount of parameters to be optimised and to make the easy observation parameter update possible. The diagonality of the HMM state dependent covariance matrices is not such a bad assumption anymore since the observation matrix is used to model the spatial covariance structure of the observations. The number of parameters per state in a FAHMM in speech recognition is $(M^{(x)} - 1)(1 + 2k) + pk + (M^{(o)} - 1) + M^{(o)}2p + 1$. The first term correspond to the state-space mixture priors and HMM state distribution statistics. One state-space distribution component can be again omitted due to the degeneracy of the factor analysis model as discussed in Section 3.3. The second term corresponds to the state dependent observation matrix, the third to the observation noise priors and the fourth for the distribution. The last term corresponds to the transition probability of a left to right HMM.

5.4 Relationship Between FAHMM and STC

As noted in the section about factor analysis the observation noise spans the $p - k$ nuisance dimensions making it a significant part of the model when $k < p$. If the dimensionality of the

observation and the state-space are the same, the observation noise can be cancelled without losing the identifiability of the system. However, in that case the observation equation becomes deterministic and it does not have any influence on the joint likelihood of the observation and state sequences. The re-estimation formulae change their form dramatically as well.

Since the observation equation is deterministic, the transform \mathbf{C}_j can be subsumed to the state dependent HMM output covariance matrix and the usual semi-tied full covariance matrix [11] system results with a transform class for every unique state in the system. This model can also be regarded as a noiseless independent factor analysis [1] with the factorial mixture model replaced by a standard mixture of Gaussians. Since arbitrary sharing is possible in the FAHMM framework, the transform classes can be defined more reasonably. The optimisation and further characteristics of semi-tied covariance matrices is discussed in the following section dealing with linear discriminant analysis.

6 Linear Discriminant Analysis

An alternative dimension reduction scheme to factor analysis is linear discriminant analysis [26]. LDA is based on a linear transformation which projects the samples so that the within class variance is minimised and the between class variance maximised. Provided the classes are separable in some lower dimensional space, LDA finds the optimal projection and thus can be regarded as a dimension reduction scheme. However, not much is known about the generalisation of LDA from two dimensional case and the performance of such a scheme must be evaluated experimentally.

In traditional linear discriminant analysis, the within class covariance matrices are assumed to be the same. In [29] an extension to LDA was presented in which the restriction of the equal covariance matrices was removed and it is called heteroscedastic LDA (HLDA). The generative model for LDA can be represented in the framework of this paper by allowing more general densities at the state process than a single Gaussian. For example, mixture of Gaussians can represent multiple classes with mean vectors as the cluster centres.

The traditional linear discriminant analysis is illustrated in Figure 15 where two classes in p dimensional space ($p = 2$ in the figure) are not easy to separate along the axis in the original space. Fortunately, there exists a transformation to lower dimensional space ($k = 1$ in the figure) where the classes are easily separable. Dimension \mathbf{x}_1 in the figure represents the new space and dimension \mathbf{x}_2 represents the nuisance dimension where both classes are represented by a single Gaussian.

The optimisation of a static HLDA model parameters culminates into finding the new transformation matrix $\hat{\mathbf{C}}$. Defining $\hat{\mathbf{B}} = (\hat{\mathbf{C}}^{-1})'$, the transformation can be optimised by maximising the following auxiliary function

$$\begin{aligned} \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = & \frac{1}{2} \sum_{j,m} \gamma_m(j) \log \left(\frac{\text{abs}(|\hat{\mathbf{B}}'|)^2}{\text{abs}(|\text{diag}(\hat{\mathbf{B}}'_{[k]} \mathbf{W}_m \hat{\mathbf{B}}'_{[k]})|) \text{abs}(|\text{diag}(\hat{\mathbf{B}}'_{[p-k]} \hat{\mathbf{\Sigma}} \hat{\mathbf{B}}'_{[p-k]})|)} \right) \\ & - \frac{p}{2} \sum_{j,m} \gamma_m(j) (\log(2\pi) + 1) \end{aligned} \quad (129)$$

where $\gamma_m(j)$ denotes the posterior of the mixture component m given the j th observation and the current model set. The sufficient statistics \mathbf{W}_m and $\hat{\mathbf{\Sigma}}$ are defined as follows

$$\mathbf{W}_m = \frac{\sum_j \gamma_m(j) (\mathbf{o}_j - \hat{\boldsymbol{\mu}}_m) (\mathbf{o}_j - \hat{\boldsymbol{\mu}}_m)'}{\sum_j \gamma_m(j)} \quad (130)$$

$$\hat{\mathbf{\Sigma}} = \frac{\sum_{j,m} \gamma_m(j) (\hat{\boldsymbol{\mu}}_m - \hat{\boldsymbol{\mu}}) (\hat{\boldsymbol{\mu}}_m - \hat{\boldsymbol{\mu}})'}{\sum_{j,m} \gamma_m(j)} + \frac{\sum_{j,m} \gamma_m(j) (\mathbf{o}_j - \hat{\boldsymbol{\mu}}_m) (\mathbf{o}_j - \hat{\boldsymbol{\mu}}_m)'}{\sum_{j,m} \gamma_m(j)} \quad (131)$$

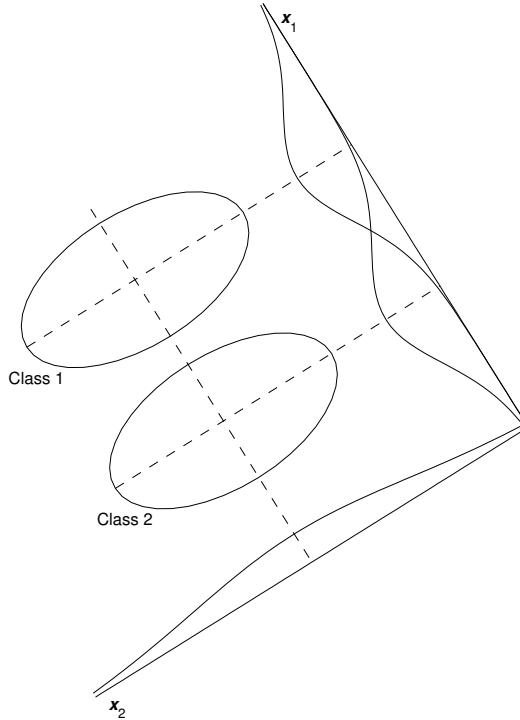


Figure 15: Linear discriminant analysis in case of two classes for which the original axes orientation is not optimal.

where

$$\hat{\boldsymbol{\mu}}_m = \frac{\sum_j \gamma_m(j) \boldsymbol{o}_j}{\sum_j \gamma_m(j)} \quad (132)$$

$$\hat{\boldsymbol{\mu}} = \frac{\sum_{j,m} \gamma_m(j) \boldsymbol{o}_j}{\sum_{j,m} \gamma_m(j)} \quad (133)$$

There exists a simple iterative scheme to optimise the columns of the matrix $\hat{\boldsymbol{B}}$ which was presented in [11]. The ML solutions for the state density mean vector and covariance matrix are

$$\hat{\boldsymbol{\mu}}_m^{(x)} = \hat{\boldsymbol{B}}'_{[k]} \hat{\boldsymbol{\mu}}_m \quad (134)$$

$$\hat{\boldsymbol{\Sigma}}_m^{(x)} = \text{diag}(\hat{\boldsymbol{B}}'_{[k]} \boldsymbol{W}_m \hat{\boldsymbol{B}}_{[k]}) \quad (135)$$

and the new observation noise parameters are simply

$$\hat{\boldsymbol{\mu}}^{(o)} = \hat{\boldsymbol{\mu}}_{[p-k]} \quad (136)$$

$$\hat{\boldsymbol{\Sigma}}^{(o)} = \hat{\boldsymbol{\Sigma}}_{[p-k]} \quad (137)$$

The benefit of the LDA and HLDA is that the likelihood is calculated only in the useful subspace and therefore, storage of the observation noise parameters is not needed.

6.1 Linear Discriminant Analysis Hidden Markov Models

Finding a better subspace in HMM framework using linear discriminant analysis and HLDA was studied in [12, 21, 29]. In [12], the close relationship between semi-tied full covariance matrices

(STC) and HLDA was presented. STC is based on class specific transformations which are applied to the state dependent output covariance matrices of a HMM. HMM with HLDA can be viewed as a STC model when the state-space and observation space have the same dimensionality, $k = p$, and there exists only one transform class.

A linear discriminant analysis HMM can be expressed with the following generative equations

$$\mathbf{x}_t \sim \mathcal{M}^{hmm} \quad (138)$$

$$\mathbf{o}_t = \mathbf{C} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{v} \end{bmatrix}, \quad \mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (139)$$

where a single Gaussian observation noise is used. The generative model above can also be modified to represent other interesting models. If the observation noise is replaced by another state vector distributed according to a HMM, the equations define a two stream HMM with a linear mixing process. Depending on the form of the observation matrix, the model can represent an independent stream HMM or a factorial HMM [18]. As noted before, many interesting models based on linear discriminant analysis observation process are omitted in this paper.

The optimisation of the parameters of a HMM with HLDA observation process is the same as in the static case presented above. Only the observation indices must be replaced by time indices and the component posterior replaced by the joint likelihood of the current HMM state and the current mixture component given the observation sequence and the current model parameters. The posteriors can be obtained using the standard forward-backward recursions by replacing the posteriors of the observation given the HMM state accordingly.

6.2 Dynamical Linear Discriminant Analysis

Dynamical Linear Discriminant Analysis (DLDA) is based on continuous hidden k dimensional state variable vectors, \mathbf{x}_t , which evolve according to linear first-order Markov dynamics. A p dimensional observation vector, \mathbf{o}_t is generated from the current state by a linear observation process. DLDA can be described by the following two equations

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \quad (140)$$

$$\mathbf{o}_t = \mathbf{C} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{v} \end{bmatrix}, \quad \mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (141)$$

where the relationship between the state vectors and the observation vectors is deterministic. Again, only a simple observation noise distribution is used. Defining $\mathbf{B} = (\mathbf{C}^{-1})'$ the likelihood of two consecutive observations can be represented as

$$p(\mathbf{o}_t | \mathbf{o}_{t-1}) = \text{abs}|\mathbf{B}'| \mathcal{N}(\mathbf{B}'_{[k]} \mathbf{o}_t; \mathbf{A}\mathbf{B}'_{[k]} \mathbf{o}_{t-1} + \boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \mathcal{N}(\mathbf{B}'_{[p-k]} \mathbf{o}_t; \boldsymbol{\mu}^{(o)}, \boldsymbol{\Sigma}^{(o)}) \quad (142)$$

where $\mathbf{B}_{[k]}$ denotes the first k columns of the matrix \mathbf{B} and an absolute value of the determinant of the transformation matrix has to be taken to guarantee non-negative likelihoods. The transposition in the definition of \mathbf{B} is needed to stick to our notation in the optimisation below. The likelihood of an observation sequence, $\mathbf{O} = \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T$, is simply

$$p(\mathbf{O}) = \prod_{t=2}^T p(\mathbf{o}_t | \mathbf{o}_{t-1}) \quad (143)$$

and the auxiliary function to be maximised is therefore defined by

$$\mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) = \sum_{t=2}^T \log p(\mathbf{o}_t | \mathbf{o}_{t-1}) \quad (144)$$

Again, in the likelihood calculations the observation noise term can be ignored since it does not carry any class discriminating information.

6.2.1 Parameter Optimisation

If diagonal state transition and observation noise covariance matrices are used, the parameter optimisation can be done easily row by row. It is shown in Appendix G that the column vectors $1 \leq i \leq k$ of the matrix $\hat{\mathbf{B}} = (\hat{\mathbf{C}}^{-1})'$ are given by

$$\hat{\mathbf{b}}_i = \mathbf{k}_i \mathbf{G}_i^{-1} (\mathbf{k}_i + \delta \mathbf{c}_i) \quad (145)$$

where \mathbf{c}_i is a vector of the cofactors of the matrix $\hat{\mathbf{B}}$ corresponding to the elements in the vector $\hat{\mathbf{b}}_i$, δ satisfies a usual quadratic expression in Eq. 289 and

$$\mathbf{G}_j = \frac{1}{\sigma_j^{(x)2}} \sum_{t=2}^T (\mathbf{o}_t \mathbf{o}'_t - a_{jj} \mathbf{o}_t \mathbf{o}'_{t-1} - a_{jj} \mathbf{o}_{t-1} \mathbf{o}_t + a_{jj}^2 \mathbf{o}_{t-1} \mathbf{o}'_{t-1}) \quad (146)$$

$$\mathbf{k}_j = \frac{\mu_j^{(x)2}}{\sigma_j^{(x)2}} \sum_{t=2}^T (\mathbf{o}_t - a_{jj} \mathbf{o}_{t-1}) \quad (147)$$

where a_{jj} is the j th diagonal element of the state transition matrix \mathbf{A} . The rest $p - k$ columns of the matrix $\hat{\mathbf{B}}$ can be obtained by just replacing the statistics as follows

$$\mathbf{G}_j = \frac{1}{\sigma_{(j-k)}^{(o)2}} \sum_{t=2}^T \mathbf{o}_t \mathbf{o}'_t \quad (148)$$

$$\mathbf{k}_j = \frac{\mu_{(j-k)}^{(o)2}}{\sigma_{(j-k)}^{(o)2}} \sum_{t=2}^T \mathbf{o}_t \quad (149)$$

The state evolution parameters of DLDA can be obtained as follows

$$\hat{\mathbf{A}} = \left(\text{diag} \left(\hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_t \mathbf{o}'_{t-1} \right) \hat{\mathbf{B}}_{[k]} - \frac{1}{T-1} \hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_t \right) \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right)' \hat{\mathbf{B}}_{[k]} \right) \right. \\ \left. \left(\text{diag} \left(\hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_{t-1} \mathbf{o}'_{t-1} \right) \hat{\mathbf{B}}_{[k]} - \frac{1}{T-1} \hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right) \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right)' \hat{\mathbf{B}}_{[k]} \right) \right)^{-1} \quad (150)$$

$$\hat{\boldsymbol{\mu}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1}) \quad (151)$$

$$\hat{\boldsymbol{\Sigma}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \text{diag} \left(\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \mathbf{o}'_t \hat{\mathbf{B}}_{[k]} - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \left[\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \mathbf{o}'_{t-1} \hat{\mathbf{B}}_{[k]} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \right]' \right) \quad (152)$$

and finally the observation noise parameters can be re-estimated as follows

$$\hat{\boldsymbol{\mu}}^{(o)} = \frac{1}{T-1} \hat{\mathbf{B}}'_{[p-k]} \sum_{t=2}^T \mathbf{o}_t \quad (153)$$

$$\hat{\boldsymbol{\Sigma}}^{(o)} = \frac{1}{T-1} \sum_{t=2}^T \text{diag} \left(\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t \mathbf{o}'_t \hat{\mathbf{B}}_{[p-k]} - \hat{\boldsymbol{\mu}}^{(o)} \hat{\boldsymbol{\mu}}^{(o)'} \right) \quad (154)$$

7 Implementation Issues

In this section, some implementation issues are discussed. Firstly, there exists a nice result from matrix algebra that can be exploited in the inversion of the large matrices in case of models based on factor analysis observation process. Secondly, some issues in variance flooring and parameter tying are presented.

7.1 Efficient E Step Calculation

A p by p matrix of the form $\mathbf{R} + \mathbf{STU}$ has to be inverted in the E step of factor analysis based models reviewed in this paper; namely, $\mathbf{K} = \mathbf{C}'(\mathbf{C}\mathbf{C}' + \mathbf{\Sigma}^{(o)})^{-1}$ in case of factor analysis, $\mathbf{K}_m = \mathbf{\Sigma}^{(x)}\mathbf{C}'(\mathbf{C}\mathbf{\Sigma}^{(x)}\mathbf{C}' + \mathbf{\Sigma}^{(o)})^{-1}$ in case of mixtures of factors, $\mathbf{K}_t = \mathbf{\Sigma}_t^{t-1}\mathbf{C}'(\mathbf{C}\mathbf{\Sigma}_t^{t-1}\mathbf{C}' + \mathbf{\Sigma}^{(o)})^{-1}$ in case of linear dynamical systems and $\mathbf{K}_{jmn} = \mathbf{\Sigma}_{jn}^{(x)}\mathbf{C}'_j(\mathbf{C}_j\mathbf{\Sigma}_{jn}^{(x)}\mathbf{C}'_j + \mathbf{\Sigma}_{jm}^{(o)})^{-1}$ in case of factor analysed HMMs. Since mostly diagonal covariance matrices are used, it is possible to use the following matrix inversion lemma to implement the inversion efficiently [23]

$$(\mathbf{R} + \mathbf{STU})^{-1} = \mathbf{R}^{-1} - \mathbf{R}^{-1}\mathbf{S}(\mathbf{T}^{-1} + \mathbf{UR}^{-1}\mathbf{S})^{-1}\mathbf{UR}^{-1} \quad (155)$$

where \mathbf{R} is replaced by a diagonal output noise covariance matrix in all the previous cases and is therefore easy to invert, and the k by k matrix $(\mathbf{T}^{-1} + \mathbf{UR}^{-1}\mathbf{S})$ can be inverted more efficiently since $k < p$ in case of subspace modelling.

Similar result for the determinants of the matrices of the form $\mathbf{R} + \mathbf{STU}$ can be employed in the likelihood calculations. Since the matrices corresponding to \mathbf{R} and \mathbf{T} are often diagonal apart from LDSs where $\mathbf{\Sigma}_t^{t-1}$ is unfortunately not diagonal but is still in the subspace, the determinant

$$|\mathbf{R} + \mathbf{STU}| = |\mathbf{R}||\mathbf{T}||\mathbf{T}^{-1} + \mathbf{UR}^{-1}\mathbf{S}| \quad (156)$$

can be obtained efficiently. The inverse $(\mathbf{T}^{-1} + \mathbf{UR}^{-1}\mathbf{S})^{-1}$ in the subspace can be computed using Cholesky decomposition since the matrices are symmetric positive definite and the determinant $|\mathbf{T}^{-1} + \mathbf{UR}^{-1}\mathbf{S}|$ can be obtained as a by-product.

7.2 Variance Flooring for FAHMMs

The number of Gaussian components in a large system based on FAHMMs is often huge. To guarantee that the covariance matrices are invertible after re-estimation, flooring of the variance components has to be considered. In a standard diagonal covariance HMM the global variance of the training data is first computed. If the i th component of a re-estimated variance, $\sigma_{jmi}^{(o)2}$, falls below a certain fraction of the corresponding global variance component, σ_{fi}^2 , it is floored to σ_{fi}^2 [42]. The floor is often chosen to be one hundredth of the global variance.

In case of factor analysed HMMs, there are two mixture distributions, one in state level and the other in observation level. In the experiments carried out so far, flooring was found to be an issue only with very few observation noise variance elements. The full covariance matrices $\mathbf{C}_j\mathbf{\Sigma}_{jn}^{(x)}\mathbf{C}'_j + \mathbf{\Sigma}_{jm}^{(o)}$ have to be invertible for every combination of observation and state-space noise components. Since the matrices $\mathbf{C}_j\mathbf{\Sigma}_{jn}^{(x)}\mathbf{C}'_j$ are generally singular, even one zero variance element in a single observation noise component can make the full covariance matrix singular. This was found to happen very seldom with a context dependent large vocabulary speech recognition system when there was not enough training examples assigned to a particular state. The standard variance flooring scheme of HMM systems on the observation noise variances was found to be working without any dramatic influence on the training and recognition likelihood scores.

7.3 Parameter Tying

As the number of model parameters is increased, more data is required to obtain reliable estimates in the re-estimation process. As discussed earlier, tying of the model parameters becomes important as the complexity of the models is increased. The parameter tying has only been seen in this

paper in the re-estimation of the mixture of linear dynamical systems and the factor analysers with mixtures of observation processes. It is though possible to arbitrarily share any parameters in different levels of every model presented in this paper.

In general, the parameter tying requires simply accumulation of the sufficient statistics in the E step of the EM algorithm and the M step equations remain as usual. For example, if the observation matrices of single FAHMM over all the states are to be tied, the statistics \mathbf{G}_l and \mathbf{k}_l have to be summed over all the states $j \in (1, N_s)$ before updating the parameters in the usual way. Also, parameters over several models can be tied similarly. For example, entire state parameters within several models can be tied as is often done in the tied state hidden Markov models in continuous speech recognition.

8 Conclusion

Several interesting linear Gaussian models have been presented in this paper. The models fall into category of general state-space models with assumptions on linearity and distributions which are either Gaussians or mixtures of Gaussians. The state vectors can be regarded as generated by a single distribution in case of static models or either linear first-order Gauss-Markov random process or hidden Markov (HMM) model in case of dynamic models.

This paper complements the attempt to unify the set of linear Gaussian models presented in [38] although both these papers have still not covered all the possible forms. Some variants of linear discriminant analysis (LDA) [12] have been omitted as well as all the possible combinations of LDA with different mixture assumptions. Instead, factor analysis models and factor analysed hidden Markov models were presented in very generic fashion combining the way of representing mixtures of factors such as used in independent factor analysis [1], and mixtures of factor analysis observation processes with different parameter sharing schemes such as used in mixtures of factor analysers [15] and shared factor analysis [22].

Factor analysed hidden Markov model (FAHMM) provides a very generic extension to currently very popular time-series model, the hidden Markov model. FAHMM is based on state vectors generated by a HMM and an observation process that performs shared factor analysis using the state vector elements as factors. FAHMMs provide a flexible framework for time-series modelling since all the advantages of HMMs are present such as efficient segmentation using Viterbi algorithm but, in addition, a variety of possible mixing and tying schemes of the model parameters can be employed. An important benefit of FAHMMs in comparison to the standard diagonal covariance HMMs is the modelling of spatial correlation which is often present due to non-ideal feature extraction algorithms.

Several models suitable in segment model (SM) [34] framework have also been presented. It has been noted that the major drawback in the segment models using linear dynamical systems (LDSs) as the models for feature vector dynamics is the unimodal nature of the observations. It is well known that many feature vectors such as in speech recognition have multi-modal distributions due to the source and environment variability. In this paper, a form of mixture of LDSs where the parameters of several independent LDSs other than the observation noise parameters are tied so representing crudely a LDS with a mixture of Gaussians as the observation noise distribution. Other forms of LDSs with mixture model noise sources do not have tractable training schemes and some sort of approximate methods have to be considered such as particle filtering [8]. There are still many interesting alternatives to extend the class of models for feature dynamics in SM framework. It will be a worthwhile to investigate these models since all the promises of SMs have not yet been utilised.

Acknowledgements

A-V.I. Rosti is funded by the Cambridge European Trust, EPSRC, Finnish Cultural Foundation, Nokia Foundation, Research Scholarship Foundation of Tampere, Tampere Chamber of Commerce

and Industry, and Tampere Graduate School in Information Science and Engineering.

References

- [1] H. Attias. Independent factor analysis. *Neural Computation*, 11(4):803–851, 1999.
- [2] L.E. Baum and T. Petrie. Statistical inference for probabilistic functions of finite state Markov chains. *Annals of Mathematical Statistics*, 37:1554–1563, 1966.
- [3] L.E. Baum, T. Petrie, G. Soules, and N. Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Annals of Mathematical Statistics*, 41(1):164–171, 1970.
- [4] C.S. Blackburn. *Articulatory Methods for Speech Production and Recognition*. PhD thesis, University of Cambridge, 1996.
- [5] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39:1–38, 1977.
- [6] V. Digalakis. *Segment-Based Stochastic Models of Spectral Dynamics for Continuous Speech Recognition*. PhD thesis, Boston University, 1992.
- [7] V. Digalakis, J.R. Rohlicek, and M. Ostendorf. ML estimation of a stochastic linear system with the EM algorithm and its application to speech recognition. *IEEE Transactions on Speech and Audio Processing*, 1(4):431–442, 1993.
- [8] A. Doucet, S. Godsill, and C. Andrieu. On sequential Monte Carlo sampling methods for Bayesian filtering. *Statistics and Computing*, 10(3):197–208, 2000.
- [9] A. Doucet, S.J. Godsill, and M. West. Monte carlo filtering and smoothing with application to time-varying spectral estimation. In *Proceedings International Conference on Acoustics, Speech and Signal Processing*, pages 701–704, 2000.
- [10] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, 2 edition, 1990.
- [11] M.J.F. Gales. Semi-tied covariance matrices for hidden Markov models. *IEEE Transactions on Speech and Audio Processing*, 7(3):272–281, 1999.
- [12] M.J.F. Gales. Maximum likelihood multiple subspace projections for hidden Markov models. *IEEE Transactions on Speech and Audio Processing*, 10(2):37–47, 2002.
- [13] M.J.F. Gales, K.M. Knill, and S.J. Young. State-based Gaussian selection in large vocabulary continuous speech recognition using HMMs. *IEEE Transactions on Speech and Audio Processing*, 7(2):152–161, 1999.
- [14] Z. Ghahramani. Learning dynamic Bayesian networks. In C.L. Giles and M. Gori, editors, *Adaptive Processing of Sequences and Data Structures*, volume 1387 of *Lecture Notes in Computer Science*, pages 168–197. Springer, 1998.
- [15] Z. Ghahramani and G. Hinton. The EM algorithm for mixtures of factor analyzers. Technical Report CRG-TR-96-1, Department of Computer Science, University of Toronto, 1996.
- [16] Z. Ghahramani and G. Hinton. Parameter estimation for linear dynamical systems. Technical Report CRG-TR-96-2, Department of Computer Science, University of Toronto, 1996.
- [17] Z. Ghahramani and G.E. Hinton. Variational learning for switching state-space models. *Neural Computation*, 12(4):963–996, 1998.

- [18] Z. Ghahramani and M.I. Jordan. Factorial hidden Markov models. *Machine Learning*, 29:245–273, 1997.
- [19] S. Godsill, A. Doucet, and M. West. Maximum a posteriori sequence estimation using Monte Carlo particle filters. *Annals of the Institute of Statistical Mathematics*, 52, 2001.
- [20] R.A. Gopinath. Constrained maximum likelihood modeling with gaussian distributions. In *Proceedings Broadcast News Transcription and Understanding Workshop*, 1998.
- [21] R.A. Gopinath. Maximum likelihood modeling with Gaussian distributions for classification. In *Proceedings International Conference on Acoustics, Speech, and Signal Processing*, volume 2, pages 661–664, 1998.
- [22] R.A. Gopinath, B. Ramabhadran, and S. Dharanipragada. Factor analysis invariant to linear transformations of data. In *Proceedings International Conference on Speech and Language Processing*, pages 397–400, 1998.
- [23] D.A. Harville. *Matrix Algebra from a Statistician’s Perspective*. Springer, 1997.
- [24] M. Isard and A. Blake. Contour tracking by stochastic propagation of conditional density. In *Proceedings European Conference on Computer Vision*, volume 1, pages 343–356, 1996.
- [25] M. Isard and A. Blake. Condensation – conditional density propagation for visual tracking. *International Journal of Computer Vision*, 29(1):5–28, 1998.
- [26] R.A. Johnson and D.W. Wichern. *Applied Multivariate Statistical Analysis*. Prentice Hall, 4 edition, 1998.
- [27] R.E. Kalman. A new approach to linear filtering and prediction problems. *Transactions of the American Society of Mechanical Engineering, Series D, Journal of Basic Engineering*, 82:35–45, 1960.
- [28] R.E. Kalman and R.S. Bucy. New results in linear filtering and prediction theory. *Transactions of the American Society of Mechanical Engineering, Series D, Journal of Basic Engineering*, 83:95–108, 1961.
- [29] N. Kumar. *Investigation of Silicon-Auditory Models and Generalization of Linear Discriminant Analysis for Improved Speech Recognition*. PhD thesis, Johns Hopkins University, 1997.
- [30] C.J. Leggetter and P.C. Woodland. Maximum likelihood linear regression for speaker adaptation of continuous density HMMs. *Computer Speech and Language*, 9(2):171–185, 1995.
- [31] L.A. Liporace. Maximum likelihood estimation for multivariate observations of Markov sources. *IEEE Transactions on Information Theory*, IT-28(5):729–734, 1982.
- [32] T.P. Minka. From hidden Markov models to linear dynamical systems. Technical Report #531, Vision and Modeling Group, MIT Media Laboratory, 1999. Available at <http://vismod.www.media.mit.edu/~tpminka/papers/learning.html>.
- [33] K. Murphy. Learning switching Kalman filter models. Technical Report 98-10, Compaq Cambridge Research Lab., 1998. Available at <http://www.cs.berkeley.edu/~murphyk/publ.html>.
- [34] M. Ostendorf, V. Digalakis, and O. Kimball. From HMM’s to segment models: A unified view of stochastic modeling for speech recognition. *IEEE Transactions on Speech and Audio Processing*, 4(5):360–378, 1996.
- [35] L.R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–285, 1989.

- [36] H.E. Rauch. Solutions to the linear smoothing problem. *IEEE Transactions on Automatic Control*, 8:371–372, 1963.
- [37] H.E. Rauch, F. Tung, and C.T. Striebel. Maximum likelihood estimates of linear dynamic systems. *American Institute of Aeronautics and Astronautics Journal*, 3(8):1445–1450, 1965.
- [38] S. Roweis and Z. Ghahramani. A unifying review of linear Gaussian models. *Neural Computation*, 11(2):305–345, 1999.
- [39] D.B. Rubin and D.T. Thayer. EM algorithms for ML factor analysis. *Psychometrika*, 47(1):69–76, 1982.
- [40] L. Saul and M. Rahim. Maximum likelihood and minimum classification error factor analysis for automatic speech recognition. *IEEE Transactions on Speech and Audio Processing*, 8(2):115–125, 1999.
- [41] A.J. Viterbi. Error bounds for convolutional codes and an asymptotically optimal decoding algorithm. *IEEE Transactions on Information Theory*, IT-13:260–269, 1967.
- [42] S. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland. *The HTK Book (for HTK Version 3.0)*. Cambridge University, 2000.
- [43] R.S. Zemel. *A Minimum Description Length Framework for Unsupervised Learning*. PhD thesis, University of Toronto, 1994.

A Useful Results from Matrix Algebra

This appendix reviews some useful results from the matrix algebra for the derivations in the following appendices. First, inversion formulae for partitioned matrices are presented and then conditional density of two multivariate Gaussians is derived.

A.1 Inverting Partitioned Matrices

Let \mathbf{A} represent an arbitrary $(n + m)$ by $(n + m)$ matrix with the following partitioning

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \quad (157)$$

where \mathbf{A}_1 is an n by n matrix, \mathbf{A}_2 an n by m matrix, \mathbf{A}_3 an m by n matrix and \mathbf{A}_4 an m by m matrix. Suppose that \mathbf{A}_1 is non-singular. The matrix \mathbf{A} is singular if and only if the n by n matrix

$$(\mathbf{A}|\mathbf{A}_4) = \mathbf{A}_1 - \mathbf{A}_2\mathbf{A}_4^{-1}\mathbf{A}_3 \quad (158)$$

is non-singular [23]. The matrix $(\mathbf{A}|\mathbf{A}_4)$ is called the Schur complement of \mathbf{A}_4 in \mathbf{A} and the following two identities apply

$$\begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{A}|\mathbf{A}_4)^{-1} & -(\mathbf{A}|\mathbf{A}_4)^{-1}\mathbf{A}_2\mathbf{A}_4^{-1} \\ -\mathbf{A}_4^{-1}\mathbf{A}_3(\mathbf{A}|\mathbf{A}_4)^{-1} & \mathbf{A}_4^{-1} + \mathbf{A}_4^{-1}\mathbf{A}_3(\mathbf{A}|\mathbf{A}_4)^{-1}\mathbf{A}_2\mathbf{A}_4^{-1} \end{bmatrix} \quad (159)$$

$$= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_4^{-1} \end{bmatrix} + \begin{bmatrix} \mathbf{I}_n & \\ & -\mathbf{A}_4^{-1}\mathbf{A}_3 \end{bmatrix} (\mathbf{A}|\mathbf{A}_4)^{-1} \begin{bmatrix} \mathbf{I}_n \\ -\mathbf{A}_2\mathbf{A}_4^{-1} \end{bmatrix}' \quad (160)$$

A.2 Conditioning Multivariate Gaussians

Let \mathbf{x} and \mathbf{y} be p and k dimensional Gaussian distributed random vectors with mean vectors, $\boldsymbol{\mu}_x$ and $\boldsymbol{\mu}_y$, and covariance matrices, $\boldsymbol{\Sigma}_x$ and $\boldsymbol{\Sigma}_y$, respectively. Suppose also that \mathbf{x} and \mathbf{y} are also jointly Gaussian with $\boldsymbol{\Sigma}_{yx}$ and $\boldsymbol{\Sigma}_{xy}$ as their cross covariance matrices. The joint distribution can be represented as follows

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix}\right) \quad (161)$$

It should be noted that $\boldsymbol{\Sigma}_{yx} = \boldsymbol{\Sigma}_{xy}'$.

The posterior distribution of \mathbf{x} given \mathbf{y} is obtained by the definition, $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}, \mathbf{y})/p(\mathbf{y})$,

$$p(\mathbf{x}|\mathbf{y}) = (2\pi)^{-\frac{(p+k)}{2}} \left| \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix} \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix}' \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix} \right\} \\ \left/ \left((2\pi)^{-\frac{k}{2}} |\boldsymbol{\Sigma}_y|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu}_y)' \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \right\} \right) \right. \quad (162)$$

$$= (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_y^{-1}\boldsymbol{\Sigma}_{yx}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \gamma \right\} \quad (163)$$

where the last determinant is obtained using the following identity for partitioned matrices [23]

$$\left| \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix} \right| = |\boldsymbol{\Sigma}_y| |\boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_y^{-1}\boldsymbol{\Sigma}_{yx}| \quad (164)$$

and the definition of the scalar γ is

$$\gamma = \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix}' \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix} - (\mathbf{y} - \boldsymbol{\mu}_y)' \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \quad (165)$$

The partitioned joint covariance matrix, denoted by Σ below, can be inverted by using the identity in Eq. 160. Obviously, the first term in the right hand side of the identity cancels the last term in Eq. 165 and therefore it simplifies to the following form

$$\gamma = \left((\mathbf{x} - \boldsymbol{\mu}_x)' - (\mathbf{y} - \boldsymbol{\mu}_y)' \Sigma_y^{-1} \Sigma_{yx} \right) (\Sigma | \Sigma_y)^{-1} \left((\mathbf{x} - \boldsymbol{\mu}_x) - \Sigma_{xy} \Sigma_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \right) \quad (166)$$

which is the Mahalanobis distance between the vectors \mathbf{x} and $\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \Sigma_{xy} \Sigma_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y)$ with covariance matrix $\Sigma_{x|y} = (\Sigma | \Sigma_y) = \Sigma_x - \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx}$.

Indeed, by substituting γ back into Eq. 163 it is obvious that the posterior distribution of \mathbf{x} given \mathbf{y} is Gaussian distributed with mean vector $\boldsymbol{\mu}_{x|y}$ and covariance matrix $\Sigma_{x|y}$; i.e.,

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_x + \Sigma_{xy} \Sigma_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y), \Sigma_x - \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx}) \quad (167)$$

B Derivation of the EM Algorithm for Factor Analysis

The expectation maximisation algorithm for factor analysis is presented in this appendix. The derivation of this standard algorithm was first done in [39]. Here the derivation is presented to keep the paper self contained and using generic matrix manipulation. In the E step, the sufficient statistics required in the M step are obtained. The M step is very generic in that only few modifications are needed in the EM algorithms for the other models with factor analysis observation process.

B.1 E Step for Factor Analysis

The state posteriors have to be obtained in the E step of the EM algorithm. It is easy to see from the factor analysis observation equation in Eq. 11 that the likelihood of the observation vector, \mathbf{o}_j , can be obtained as follows

$$p(\mathbf{o}_j) = \mathcal{N}(\mathbf{o}_j; \boldsymbol{\mu}^{(o)}, \mathbf{C}\mathbf{C}' + \Sigma^{(o)}) \quad (168)$$

and since the observation and the state vector are also jointly Gaussian with the following joint density

$$p(\mathbf{o}_j, \mathbf{x}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{o}_j \\ \mathbf{x} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}^{(o)} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{C}\mathbf{C}' + \Sigma^{(o)} & \mathbf{C} \\ \mathbf{C}' & \mathbf{I} \end{bmatrix}\right) \quad (169)$$

and the state posteriors can be obtained using Eq. 167

$$p(\mathbf{x}|\mathbf{o}_j) = \mathcal{N}(\mathbf{x}; \mathbf{K}(\mathbf{o}_j - \boldsymbol{\mu}^{(o)}), \mathbf{I} - \mathbf{K}\mathbf{C}) \quad (170)$$

where $\mathbf{K} = \mathbf{C}'(\mathbf{C}\mathbf{C}' + \Sigma^{(o)})^{-1}$.

The required first and second-order sufficient statistics can be now represented as follows

$$\hat{\mathbf{x}}(j) = E\{\mathbf{x}|\mathbf{o}_j, \mathcal{M}\} = \mathbf{K}(\mathbf{o}_j - \boldsymbol{\mu}^{(o)}) \quad (171)$$

$$\hat{\mathbf{R}}(j) = E\{\mathbf{x}\mathbf{x}'|\mathbf{o}_j, \mathcal{M}\} = \mathbf{I} - \mathbf{K}\mathbf{C} + \hat{\mathbf{x}}(j)\hat{\mathbf{x}}'(j) \quad (172)$$

where \mathcal{M} is added to show explicitly that the old model parameters are used.

B.2 M Step for Factor Analysis

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{C}}, \hat{\boldsymbol{\mu}}^{(o)}, \hat{\Sigma}^{(o)})$, for the observation process the following auxiliary function has to be maximised

$$\begin{aligned} \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = & \quad (173) \\ & -\frac{1}{2} \sum_{j=1}^N \left(\log |\hat{\Sigma}^{(o)}| + E\left\{ (\mathbf{o}_j - \hat{\mathbf{C}}\mathbf{x} - \hat{\boldsymbol{\mu}}^{(o)})' \Sigma^{(o)-1} (\mathbf{o}_j - \hat{\mathbf{C}}\mathbf{x} - \hat{\boldsymbol{\mu}}^{(o)}) \middle| \mathcal{O}, \mathcal{M} \right\} \right) \end{aligned}$$

Differentiating Eq. 173 with respect to $\hat{\boldsymbol{\mu}}^{(o)}$ and equating the resulting expression to zero

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}^{(o)}} = \hat{\boldsymbol{\Sigma}}^{(o)-1} \sum_{j=1}^N (\mathbf{o}_j - \hat{\mathbf{C}} \hat{\mathbf{x}}(j) - \hat{\boldsymbol{\mu}}^{(o)}) = \mathbf{0} \quad (174)$$

$$\hat{\boldsymbol{\mu}}^{(o)} = \frac{1}{N} \sum_{j=1}^N (\mathbf{o}_j - \hat{\mathbf{C}} \hat{\mathbf{x}}(j)) \quad (175)$$

Differentiating Eq. 173 with respect to $\hat{\mathbf{C}}$, substituting the above ML estimate of the observation noise mean and setting the resulting equation to zero, the new observation matrix can be solved using only the sufficient statistics from the E step as follows

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\mathbf{C}}} = \hat{\boldsymbol{\Sigma}}^{(o)-1} \sum_{j=1}^N (\mathbf{o}_j \hat{\mathbf{x}}'(j) - \hat{\boldsymbol{\mu}}^{(o)} \hat{\mathbf{x}}'(j) - \hat{\mathbf{C}} \hat{\mathbf{R}}(j)) = \mathbf{0} \quad (176)$$

$$\hat{\mathbf{C}} = \left(\sum_{j=1}^N \mathbf{o}_j \hat{\mathbf{x}}'(j) - \frac{1}{N} \sum_{j=1}^N \mathbf{o}_j \sum_{j=1}^N \hat{\mathbf{x}}'(j) \right) \left(\sum_{j=1}^N \hat{\mathbf{R}}(j) - \frac{1}{N} \sum_{j=1}^N \hat{\mathbf{x}}(j) \sum_{j=1}^N \hat{\mathbf{x}}'(j) \right)^{-1} \quad (177)$$

When re-estimating the model parameters, the observation matrix has to be updated before the observation noise mean vector or they can be re-estimated simultaneously using matrix notation. By redefining the following sums of the sufficient statistics

$$\begin{aligned} \boldsymbol{\Gamma}_1 &= \sum_{j=1}^N \mathbf{o}_j \hat{\mathbf{x}}'(j) \quad , \boldsymbol{\Gamma}_2 = \sum_{j=1}^N \hat{\mathbf{R}}(j) \\ \boldsymbol{\zeta}_1 &= \sum_{j=1}^N \mathbf{o}_j \quad , \boldsymbol{\zeta}_2 = \sum_{j=1}^N \hat{\mathbf{x}}(j) \end{aligned} \quad (178)$$

it can be checked that the parameters $\hat{\mathbf{C}}$ and $\hat{\boldsymbol{\mu}}^{(o)}$ can be maximised simultaneously as follows

$$\begin{bmatrix} \hat{\mathbf{C}} & \hat{\boldsymbol{\mu}}^{(o)} \end{bmatrix} = \left(\sum_{j=1}^N \begin{bmatrix} \mathbf{o}_j \hat{\mathbf{x}}'(j) & \mathbf{o}_j \end{bmatrix} \right) \left(\sum_{j=1}^N \begin{bmatrix} \hat{\mathbf{R}}(j) & \hat{\mathbf{x}}(j) \\ \hat{\mathbf{x}}'(j) & 1 \end{bmatrix} \right)^{-1} = \begin{bmatrix} \boldsymbol{\Gamma}_1 & \boldsymbol{\zeta}_1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Gamma}_2 & \boldsymbol{\zeta}_2 \\ \boldsymbol{\zeta}_2' & N \end{bmatrix}^{-1} \quad (179)$$

Now, the Schur's complement, $(\boldsymbol{\Gamma}|N)$, in Eq. 160 for the matrix to be inverted is

$$(\boldsymbol{\Gamma}|N) = \boldsymbol{\Gamma}_2 - \frac{1}{N} \boldsymbol{\zeta}_2 \boldsymbol{\zeta}_2' \quad (180)$$

and the first element of the matrix product in Eq. 179 is

$$\boldsymbol{\Gamma}_1 (\boldsymbol{\Gamma}|N)^{-1} - \frac{1}{N} \boldsymbol{\zeta}_1 \boldsymbol{\zeta}_2' (\boldsymbol{\Gamma}|N)^{-1} = \left(\boldsymbol{\Gamma}_1 - \frac{1}{N} \boldsymbol{\zeta}_1 \boldsymbol{\zeta}_2' \right) \left(\boldsymbol{\Gamma}_2 - \frac{1}{N} \boldsymbol{\zeta}_2 \boldsymbol{\zeta}_2' \right)^{-1} \quad (181)$$

which is exactly the ML estimate for $\hat{\mathbf{C}}$. The second element of the matrix product in Eq. 179 is

$$\begin{aligned} & -\frac{1}{N} \boldsymbol{\Gamma}_1 (\boldsymbol{\Gamma}|N)^{-1} \boldsymbol{\zeta}_2 + \frac{1}{N} \boldsymbol{\zeta}_1 + \frac{1}{N^2} \boldsymbol{\zeta}_1 \boldsymbol{\zeta}_2' (\boldsymbol{\Gamma}|N)^{-1} \boldsymbol{\zeta}_2 \\ & = \frac{1}{N} \left(\boldsymbol{\zeta}_1 - \left(\boldsymbol{\Gamma}_1 (\boldsymbol{\Gamma}|N)^{-1} - \frac{1}{N} \boldsymbol{\zeta}_1 \boldsymbol{\zeta}_2' (\boldsymbol{\Gamma}|N)^{-1} \right) \boldsymbol{\zeta}_2 \right) \end{aligned} \quad (182)$$

which is exactly the ML estimate for $\hat{\boldsymbol{\mu}}^{(o)}$.

To find the new observation noise covariance matrix, it is easier to differentiate Eq. 173 with respect to the inverse of $\hat{\boldsymbol{\Sigma}}^{(o)}$ as follows

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\Sigma}}^{(o)-1}} = \frac{1}{2} \sum_{j=1}^N \left(\hat{\boldsymbol{\Sigma}}^{(o)} - E \left\{ (\mathbf{o}_j - \hat{\mathbf{C}} \mathbf{x} - \hat{\boldsymbol{\mu}}^{(o)}) (\mathbf{o}_j - \hat{\mathbf{C}} \mathbf{x} - \hat{\boldsymbol{\mu}}^{(o)})' \middle| \mathcal{O}, \mathcal{M} \right\} \right) = \mathbf{0} \quad (183)$$

$$\hat{\boldsymbol{\Sigma}}^{(o)} = \frac{1}{N} \sum_{j=1}^N \text{diag} \left(\mathbf{o}_j \mathbf{o}_j' - \begin{bmatrix} \hat{\mathbf{C}} & \hat{\boldsymbol{\mu}}^{(o)} \end{bmatrix} \begin{bmatrix} \mathbf{o}_j \hat{\mathbf{x}}'(j) & \mathbf{o}_j \end{bmatrix}' \right) \quad (184)$$

where the ML estimates of the observation matrix and noise mean vector must be used. The last form is obtained by noting that the terms inside the expectation in Eq. 183 can be reorganised as follows

$$\begin{aligned} & \mathbf{o}_j \mathbf{o}_j' - \\ & \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] \left[\mathbf{o}_j \hat{\mathbf{x}}'(j) \ \mathbf{o}_j \right]' - \left[\mathbf{o}_j \hat{\mathbf{x}}'(j) \ \mathbf{o}_j \right] \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right]' + \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right] \begin{bmatrix} \hat{\mathbf{R}}^{(j)} & \hat{\mathbf{x}}^{(j)} \\ \hat{\mathbf{x}}'^{(j)} & 1 \end{bmatrix} \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right]' \end{aligned} \quad (185)$$

and substituting the first term of Eq. 179 in place of $\left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}^{(o)} \right]$.

C Derivation of the EM Algorithm for Mixtures of Factors Factor Analyser

The expectation maximisation algorithm for factor analysis with mixtures of factors is presented in this appendix. The E step is a bit more complicated than in general factor analysis. The M step is very generic and can be applied to the factor analysed hidden Markov models as well.

C.1 E Step for Mixtures of Factors

The joint posterior of the observation and the state vectors given the mixture component can be written as

$$p(\mathbf{o}_j, \mathbf{x}|n) = \mathcal{N}\left(\begin{bmatrix} \mathbf{o}_j \\ \mathbf{x} \end{bmatrix}; \begin{bmatrix} \mathbf{C}\boldsymbol{\mu}_n^{(x)} + \boldsymbol{\mu}^{(o)} \\ \boldsymbol{\mu}_n^{(x)} \end{bmatrix}, \begin{bmatrix} \mathbf{C}\boldsymbol{\Sigma}_n^{(x)}\mathbf{C}' + \boldsymbol{\Sigma}^{(o)} & \mathbf{C}\boldsymbol{\Sigma}_n^{(x)} \\ \boldsymbol{\Sigma}_n^{(x)}\mathbf{C}' & \boldsymbol{\Sigma}_n^{(x)} \end{bmatrix}\right) \quad (186)$$

Using the conditioning of multivariate Gaussians presented in Appendix A.2 the posterior of the state vector given the observation and the mixture component is

$$p(\mathbf{x}|\mathbf{o}_j, n) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_n^{(x)} + \mathbf{K}_n(\mathbf{o}_j - \mathbf{C}\boldsymbol{\mu}_n^{(x)} - \boldsymbol{\mu}^{(o)}), \boldsymbol{\Sigma}_n^{(x)} - \mathbf{K}_n\mathbf{C}\boldsymbol{\Sigma}_n^{(x)}) \quad (187)$$

where $\mathbf{K}_n = \boldsymbol{\Sigma}_n^{(x)}\mathbf{C}'(\mathbf{C}\boldsymbol{\Sigma}_n^{(x)}\mathbf{C}' + \boldsymbol{\Sigma}^{(o)-1})$ and the sufficient statistics needed for the E step are

$$\hat{\mathbf{x}}_{jn} = E\{\mathbf{x}|\mathbf{o}_j, n\} = \boldsymbol{\mu}_n^{(x)} + \mathbf{K}_n(\mathbf{o}_j - \mathbf{C}\boldsymbol{\mu}_n^{(x)} - \boldsymbol{\mu}^{(o)}) \quad (188)$$

$$\hat{\mathbf{R}}_{jn} = E\{\mathbf{x}\mathbf{x}'|\mathbf{o}_j, n\} = \boldsymbol{\Sigma}_n^{(x)} - \mathbf{K}_n\mathbf{C}\boldsymbol{\Sigma}_n^{(x)} + \hat{\mathbf{x}}_{jn}\hat{\mathbf{x}}_{jn}' \quad (189)$$

In case of fully tied observation process the posteriors of the state vector given the observation has to be obtained. These are easily obtained using the mixture component posteriors and the above state posteriors as follows

$$\hat{\mathbf{x}}_j = E\{\mathbf{x}|\mathbf{o}_j\} = \sum_{n=1}^{M^{(x)}} P(n|\mathbf{o}_j)E\{\mathbf{x}|\mathbf{o}_j, n\} = \sum_{n=1}^{M^{(x)}} \gamma_{jn}\hat{\mathbf{x}}_{jn} \quad (190)$$

$$\hat{\mathbf{R}}_j = E\{\mathbf{x}\mathbf{x}'|\mathbf{o}_j\} = \sum_{n=1}^{M^{(x)}} P(n|\mathbf{o}_j)E\{\mathbf{x}\mathbf{x}'|\mathbf{o}_j, n\} = \sum_{n=1}^{M^{(x)}} \gamma_{jn}\hat{\mathbf{R}}_{jn} \quad (191)$$

C.2 M Step for the Mixture Component Priors

To find the new mixture component priors, $\hat{c}_n^{(x)}$, the following auxiliary function has to be maximised

$$\mathcal{Q}_\omega(\mathcal{M}, \hat{\mathcal{M}}) = \sum_{j=1}^N \sum_{n=1}^{M^{(x)}} \gamma_{jn} \log \hat{c}_n^{(x)} \quad (192)$$

which is not possible by just setting the first derivative to zero since the logarithm is monotonically increasing function. The constraint $\sum_{n=1}^{M^{(x)}} \hat{c}_n^{(x)} = 1$ can be used to find the maximum using Lagrange method. By introducing a Lagrange multiplier λ , the maximisation of Eq. 192 is equivalent to maximising the following function

$$\mathcal{L}(\hat{c}_n^{(x)}, \lambda) = \sum_{j=1}^N \sum_{n=1}^{M^{(x)}} \gamma_{jn} \log \hat{c}_n^{(x)} - \lambda \left(\sum_{n=1}^{M^{(x)}} \hat{c}_n^{(x)} - 1 \right) \quad (193)$$

Differentiating the above equation with respect to $\hat{c}_n^{(x)}$ and λ and setting the derivatives to zeroes lead to the following pair of equations

$$\sum_{j=1}^N \gamma_{jn} \frac{1}{\hat{c}_n^{(x)}} - \lambda = 0 \quad (194)$$

$$\sum_{n=1}^{M^{(x)}} \hat{c}_n^{(x)} - 1 = 0 \quad (195)$$

By eliminating λ , the new component priors can be found as follows

$$\hat{c}_n^{(x)} = \frac{1}{N} \sum_{j=1}^N \gamma_{jn} \quad (196)$$

C.3 M Step for Mixtures of Factors

To find the new state parameters, $\hat{\mathcal{M}} = (\hat{\boldsymbol{\mu}}_n^{(x)}, \hat{\boldsymbol{\Sigma}}_n^{(x)})$, the following auxiliary function has to be maximised

$$\mathcal{Q}_s(\mathcal{M}, \hat{\mathcal{M}}) = -\frac{1}{2} \sum_{j=1}^N \sum_{n=1}^{M^{(x)}} \gamma_{jn} \left(\log |\hat{\boldsymbol{\Sigma}}_n^{(x)}| + E \left\{ (\mathbf{x} - \hat{\boldsymbol{\mu}}_n^{(x)})' \hat{\boldsymbol{\Sigma}}_n^{(x)-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_n^{(x)}) \mid \mathbf{o}_j, n, \mathcal{M} \right\} \right) \quad (197)$$

Differentiating Eq. 197 with respect to the mixture component mean vector, $\hat{\boldsymbol{\mu}}_n^{(x)}$, and setting the resulting equation to zero, the new mean vectors can be solved using the sufficient statistics from the E step as follows

$$\frac{\partial \mathcal{Q}_s(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}_n^{(x)}} = \hat{\boldsymbol{\Sigma}}_n^{(x)-1} \sum_{j=1}^N \gamma_{jn} (\hat{\boldsymbol{\mu}}_n^{(x)} - \hat{\mathbf{x}}_{jn}) = \mathbf{0} \quad (198)$$

$$\hat{\boldsymbol{\mu}}_n^{(x)} = \frac{\sum_{j=1}^N \gamma_{jn} \hat{\mathbf{x}}_{jn}}{\sum_{j=1}^N \gamma_{jn}} \quad (199)$$

To find the new state noise covariance matrix, it is easier to differentiate Eq. 197 with respect to the inverse of $\hat{\boldsymbol{\Sigma}}_n^{(x)}$ as follows

$$\frac{\partial \mathcal{Q}_s(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\Sigma}}_n^{(x)-1}} = \frac{1}{2} \sum_{j=1}^N \gamma_{jn} \left(\hat{\boldsymbol{\Sigma}}_n^{(x)} - (\hat{\mathbf{R}}_{jn} - \hat{\mathbf{x}}_{jn} \hat{\boldsymbol{\mu}}_n^{(x)'} - \hat{\boldsymbol{\mu}}_n^{(x)} \hat{\mathbf{x}}_{jn}' + \hat{\boldsymbol{\mu}}_n^{(x)} \hat{\boldsymbol{\mu}}_n^{(x)'}) \right) = \mathbf{0} \quad (200)$$

$$\hat{\boldsymbol{\Sigma}}_n^{(x)} = \frac{\sum_{j=1}^N \gamma_{jn} \hat{\mathbf{R}}_{jn}}{\sum_{j=1}^N \gamma_{jn}} - \hat{\boldsymbol{\mu}}_n^{(x)} \hat{\boldsymbol{\mu}}_n^{(x)'} \quad (201)$$

D Derivation of the EM Algorithm for Shared Factor Analysis

The expectation maximisation algorithm for shared factor analysis is presented in this appendix. It was first introduced in [22] but the derivation is in a simple consistent form. The E step is very similar to the factor analysis but the statistics have to be estimated for each component and the component posteriors have to be obtained as well. The M step is again very generic and can be applied to the mixture of linear dynamical systems as well as factor analysed hidden Markov models.

D.1 E Step for Shared Factor Analysis

In addition to the state posteriors, the component posteriors have to be estimated as follows

$$\gamma_{jm} = P_\omega(m|\mathbf{o}_j, \mathcal{M}) = \frac{c_m^{(o)} p(\mathbf{o}_j|m, \mathcal{M})}{\sum_{l=1}^{M^{(o)}} c_l^{(o)} p(\mathbf{o}_j|l, \mathcal{M})} \quad (202)$$

where the component posterior of the observation is simply

$$p(\mathbf{o}_j|m, \mathcal{M}) = \mathcal{N}(\mathbf{o}_j; \boldsymbol{\mu}_m^{(o)}, \mathbf{C}\mathbf{C}' + \boldsymbol{\Sigma}_m^{(o)}) \quad (203)$$

The state posteriors can be obtained like in the factor analysis by defining the state posteriors given the observation and the mixture component as follows

$$\hat{\mathbf{x}}_{jm} = \mathbf{K}_m(\mathbf{o}_j - \boldsymbol{\mu}_m) \quad (204)$$

$$\hat{\mathbf{R}}_{jm} = \mathbf{I} - \mathbf{K}_m\mathbf{C} + \hat{\mathbf{x}}_{jm}\hat{\mathbf{x}}_{jm}' \quad (205)$$

where $\mathbf{K}_m = \mathbf{C}'(\mathbf{C}\mathbf{C}' + \boldsymbol{\Sigma}_m^{(o)})^{-1}$.

D.2 M Step for Shared Factor Analysis

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{C}}, \hat{\boldsymbol{\mu}}_m^{(o)}, \hat{\boldsymbol{\Sigma}}_m^{(o)})$, for the observation process the following auxiliary function has to be maximised

$$\begin{aligned} \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = & \quad (206) \\ & -\frac{1}{2} \sum_{j=1}^N \sum_{m=1}^{M^{(o)}} \gamma_{jm} \left(\log |\hat{\boldsymbol{\Sigma}}_m^{(o)}| + E \left\{ (\mathbf{o}_j - \hat{\mathbf{C}}\mathbf{x} - \hat{\boldsymbol{\mu}}_m^{(o)})' \hat{\boldsymbol{\Sigma}}_m^{(o)-1} (\mathbf{o}_j - \hat{\mathbf{C}}\mathbf{x} - \hat{\boldsymbol{\mu}}_m^{(o)}) \middle| \mathcal{O}, \mathcal{M} \right\} \right) \end{aligned}$$

The M step in case of mixture models is not as elegant as before since the optimisation of the observation matrix cannot be done without using the old model parameters. In case of full covariance matrices the optimisation is a bit difficult. Fortunately, if diagonal covariance matrices were used, the auxiliary function can be rewritten ignoring all terms independent of the observation matrix and defining $\hat{\mathbf{B}} = \hat{\mathbf{C}}'$ as follows

$$\mathcal{Q}'_o(\mathcal{M}, \hat{\mathcal{M}}) = -\frac{1}{2} \sum_{l=1}^p \left(\hat{\mathbf{b}}_l' \mathbf{G}_l \hat{\mathbf{b}}_l - 2\hat{\mathbf{b}}_l' \mathbf{k}_l \right) \quad (207)$$

where $\hat{\mathbf{b}}_l$ is the l th column of the matrix $\hat{\mathbf{B}}$ and

$$\mathbf{G}_l = \sum_{m=1}^{M^{(o)}} \frac{1}{\sigma_{ml}^{(o)2}} \sum_{j=1}^N \gamma_{jm} \hat{\mathbf{R}}_{jm} \quad (208)$$

$$\mathbf{k}_l = \sum_{m=1}^{M^{(o)}} \frac{1}{\sigma_{ml}^{(o)2}} \sum_{j=1}^N \gamma_{jm} (o_{jl} - \mu_{ml}^{(o)}) \hat{\mathbf{x}}_{jm} \quad (209)$$

where $\sigma_{ml}^{(o)2}$ is the l th diagonal element of the observation noise covariance matrix $\Sigma_m^{(o)}$, o_{jl} is the l th element of the observation vector \mathbf{o}_j and $\mu_{ml}^{(o)}$ is the l th element of the observation noise mean vector $\boldsymbol{\mu}_m^{(o)}$.

Differentiating the auxiliary function in Eq. 207 with respect to $\hat{\mathbf{b}}_l$ and equating the resulting expression to zero, the vectors $\hat{\mathbf{b}}_l$ can be re-estimated as follows

$$\frac{\partial \mathcal{Q}'_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\mathbf{b}}_l} = \mathbf{k}_l - \mathbf{G}_l \hat{\mathbf{b}}_l = \mathbf{0} \quad (210)$$

$$\hat{\mathbf{b}}_l = \mathbf{G}_l^{-1} \mathbf{k}_l \quad (211)$$

which are the row vectors of the new observation matrix $\hat{\mathbf{C}}$.

Differentiating Eq. 206 with respect to the component observation mean vectors and equating the resulting expression to zero, the new estimates for the means can be obtained as follows

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}_m^{(o)}} = \hat{\Sigma}_m^{(o)-1} \sum_{j=1}^N \gamma_{jm} (\mathbf{o}_j - \hat{\mathbf{C}} \hat{\mathbf{x}}_{jm} - \hat{\boldsymbol{\mu}}_m^{(o)}) = \mathbf{0} \quad (212)$$

$$\hat{\boldsymbol{\mu}}_m^{(o)} = \frac{\sum_{j=1}^N \gamma_{jm} (\mathbf{o}_j - \hat{\mathbf{C}} \hat{\mathbf{x}}_{jm})}{\sum_{j=1}^N \gamma_{jm}} \quad (213)$$

To find the new component observation noise covariance matrices, it is easier to differentiate Eq. 206 with respect to the inverse of $\Sigma_m^{(o)}$ as follows

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\Sigma}_m^{(o)-1}} = \quad (214)$$

$$\frac{1}{2} \sum_{j=1}^N \gamma_{jm} \left(\hat{\Sigma}_m^{(o)} - E \left\{ (\mathbf{o}_j - \hat{\mathbf{C}} \mathbf{x} - \hat{\boldsymbol{\mu}}_m^{(o)}) (\mathbf{o}_j - \hat{\mathbf{C}} \mathbf{x} - \hat{\boldsymbol{\mu}}_m^{(o)})' \middle| \mathcal{O}, \mathcal{M} \right\} \right) = \mathbf{0}$$

$$\begin{aligned} \hat{\Sigma}_m^{(o)} = & \frac{1}{N} \sum_{j=1}^N \gamma_{jm} \text{diag} \left(\mathbf{o}_j \mathbf{o}_j' - \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right] \left[\mathbf{o}_j \hat{\mathbf{x}}_{jm}' \mathbf{o}_j \right]' \right. \\ & \left. - \left[\mathbf{o}_j \hat{\mathbf{x}}_{jm}' \mathbf{o}_j \right] \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right]' + \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right] \left[\begin{array}{c} \hat{\mathbf{R}}_{jm} \hat{\mathbf{x}}_{jm} \\ \hat{\mathbf{x}}_{jm}' \\ 1 \end{array} \right] \left[\hat{\mathbf{C}} \hat{\boldsymbol{\mu}}_m^{(o)} \right]' \right) \end{aligned} \quad (215)$$

Unfortunately, this expression cannot be further simplified due to the form of the observation matrix re-estimation formula.

E Derivation of the EM Algorithm for Linear Dynamical System

Derivation of the Kalman filter and smoother recursions as well as the M step of the EM algorithm for a linear dynamical system are presented in this appendix. The Kalman filter and smoother recursions are derived using the generic forward-backward algorithm and conditioning multivariate Gaussians as described in App. A.2. The M step is rather straightforward maximisation of quadratic matrix forms. The derivation presented here is very similar to [16] but more attention is paid to the matrix notation used in this paper.

E.1 Kalman Filter

Following the derivation in [32] so called generic scaled forward-backward algorithm is used. Let the forward variable, $\alpha_{\mathbf{x}}(t)$, be defined as follows

$$\alpha_{\mathbf{x}}(t) = p(\mathbf{x}_t, \mathbf{o}_1, \dots, \mathbf{o}_t) = p(\mathbf{o}_1, \dots, \mathbf{o}_t)p(\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_t) = \left(\prod_{\tau=1}^t \kappa_{\tau} \right) \hat{\alpha}_{\mathbf{x}}(t) \quad (216)$$

where $\kappa_t = p(\mathbf{o}_t|\mathbf{o}_1, \dots, \mathbf{o}_{t-1})$ are the scaling factors and $\hat{\alpha}_{\mathbf{x}}(t) = p(\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_t)$ are the scaled forward variables. Both the scaling factors and forward variables are Gaussian distributed due to the linear Gaussian model assumption. Let $\mathbf{x}^{(\tau)}(t)$ denote $E\{\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_{\tau}\}$ and $\Sigma^{(\tau)}(t)$ denote $E\{(\mathbf{x}_t - \mathbf{x}^{(\tau)}(t))(\mathbf{x}_t - \mathbf{x}^{(\tau)}(t))'|\mathbf{o}_1, \dots, \mathbf{o}_{\tau}\}$ then the scaled forward variable is distributed as

$$\hat{\alpha}_{\mathbf{x}}(t) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t)}(t), \Sigma^{(t)}(t)) \quad (217)$$

Using the first-order Markov property and the state conditional independence of the observations in Eqs. 51 and 52 the forward variable, $\alpha_{\mathbf{x}}(t)$, can be rewritten in the following recursive form

$$\alpha_{\mathbf{x}}(t) = p(\mathbf{x}_t, \mathbf{o}_1, \dots, \mathbf{o}_t) = p(\mathbf{o}_t|\mathbf{x}_t)p(\mathbf{x}_t, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}) \quad (218)$$

$$= p(\mathbf{o}_t|\mathbf{x}_t) \int p(\mathbf{x}_t, \mathbf{x}_{t-1} = \mathbf{z}, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}) d\mathbf{z} \quad (219)$$

$$= p(\mathbf{o}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1} = \mathbf{z})p(\mathbf{x}_{t-1} = \mathbf{z}, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}) d\mathbf{z} \quad (220)$$

$$= p(\mathbf{o}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1} = \mathbf{z}) \left(\prod_{\tau=1}^{t-1} \kappa_{\tau} \right) \hat{\alpha}_{\mathbf{x}}(t-1) d\mathbf{z} \quad (221)$$

$$\kappa_t \hat{\alpha}_{\mathbf{x}}(t) = p(\mathbf{o}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1} = \mathbf{z}) \hat{\alpha}_{\mathbf{x}}(t-1) d\mathbf{z} \quad (222)$$

which starts off with the initial value $\hat{\alpha}_{\mathbf{x}}(1) = p(\mathbf{x}_1|\mathbf{o}_1)$. All the terms in Eq. 222 can be expressed with Gaussians as follows

$$\kappa_t \hat{\alpha}_{\mathbf{x}}(t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}\mathbf{x}_t + \boldsymbol{\mu}^{(o)}, \Sigma^{(o)}) \int \mathcal{N}(\mathbf{x}_t; \mathbf{A}\mathbf{z} + \boldsymbol{\mu}^{(x)}, \Sigma^{(x)}) \mathcal{N}(\mathbf{z}; \mathbf{x}^{(t-1)}(t-1), \Sigma^{(t-1)}(t-1)) d\mathbf{z} \quad (223)$$

where the both terms inside the integral are also jointly Gaussian as follows

$$p(\mathbf{z}, \mathbf{x}_t) = \mathcal{N} \left(\begin{bmatrix} \mathbf{z} \\ \mathbf{x}_t \end{bmatrix}; \begin{bmatrix} \mathbf{x}^{(t-1)}(t-1) \\ \mathbf{A}\mathbf{x}^{(t-1)}(t-1) + \boldsymbol{\mu}^{(x)} \end{bmatrix}, \begin{bmatrix} \Sigma^{(t-1)}(t-1) & \Sigma^{(t-1)}(t-1)\mathbf{A}' \\ \mathbf{A}\Sigma^{(t-1)}(t-1) & \mathbf{A}\Sigma^{(t-1)}(t-1)\mathbf{A}' + \Sigma^{(x)} \end{bmatrix} \right) \quad (224)$$

and since the integration is carried out with respect to \mathbf{z} the term on the right hand side reduces to the probability of \mathbf{x}_t and by defining $\mathbf{x}^{(t-1)}(t) = \mathbf{A}\mathbf{x}^{(t-1)}(t-1) + \boldsymbol{\mu}^{(x)}$ and $\Sigma^{(t-1)}(t) = \mathbf{A}\Sigma^{(t-1)}(t-1)\mathbf{A}' + \Sigma^{(x)}$, Eq. 223 becomes

$$\kappa_t \hat{\alpha}_{\mathbf{x}}(t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}\mathbf{x}_t + \boldsymbol{\mu}^{(o)}, \Sigma^{(o)}) \mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t-1)}(t), \Sigma^{(t-1)}(t)) \quad (225)$$

$$= p(\mathbf{o}_t|\mathbf{x}_t, \mathbf{o}_1, \dots, \mathbf{o}_{t-1})p(\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_{t-1}) \quad (226)$$

which is not yet in the required form but can be transformed to it by using the Bayes' formula as follows

$$p(\mathbf{o}_t|\mathbf{x}_t, \mathbf{o}_1, \dots, \mathbf{o}_{t-1})p(\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_{t-1}) = p(\mathbf{o}_t|\mathbf{o}_1, \dots, \mathbf{o}_{t-1})p(\mathbf{x}_t|\mathbf{o}_1, \dots, \mathbf{o}_t) \quad (227)$$

and noticing that \mathbf{o}_t and \mathbf{x}_t are also jointly Gaussian

$$p(\mathbf{o}_t, \mathbf{x}_t) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_t \\ \mathbf{o}_t \end{bmatrix}; \begin{bmatrix} \mathbf{x}^{(t-1)}(t) \\ \mathbf{C}\mathbf{x}^{(t-1)}(t) + \boldsymbol{\mu}^{(o)} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}^{(t-1)}(t) & \boldsymbol{\Sigma}^{(t-1)}(t)\mathbf{C}' \\ \mathbf{C}\boldsymbol{\Sigma}^{(t-1)}(t) & \mathbf{C}\boldsymbol{\Sigma}^{(t-1)}(t)\mathbf{C}' + \boldsymbol{\Sigma}^{(o)} \end{bmatrix}\right) \quad (228)$$

The latter form in Eq. 227 can be easily obtained by applying the conditioning of two multivariate Gaussians in Eq. 167. Therefore Eq. 225 becomes

$$\begin{aligned} \kappa_t \hat{\alpha}_{\mathbf{x}}(t) &= \mathcal{N}(\mathbf{o}_t; \mathbf{C}\mathbf{x}^{(t-1)}(t) + \boldsymbol{\mu}^{(o)}, \mathbf{C}\boldsymbol{\Sigma}^{(t-1)}(t)\mathbf{C}' + \boldsymbol{\Sigma}^{(o)}) \\ &\mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t-1)}(t) + \mathbf{K}(t)(\mathbf{o}_t - \mathbf{C}\mathbf{x}^{(t-1)}(t) - \boldsymbol{\mu}^{(o)}), \boldsymbol{\Sigma}^{(t-1)}(t) - \mathbf{K}(t)\mathbf{C}\boldsymbol{\Sigma}^{(t-1)}(t)) \end{aligned} \quad (229)$$

where $\mathbf{K}(t) = \boldsymbol{\Sigma}^{(t-1)}(t)\mathbf{C}'(\mathbf{C}\boldsymbol{\Sigma}^{(t-1)}(t)\mathbf{C}' + \boldsymbol{\Sigma}^{(o)})^{-1}$ and which is now the product of the probability of the current observation given the history up to the time instant t and the scaled forward variable as defined earlier.

E.2 Kalman Smoother

The backward variable is defined as usual

$$\begin{aligned} \beta_{\mathbf{x}}(t) &= \\ p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{x}_t) &= p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{o}_1, \dots, \mathbf{o}_t) \frac{p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{x}_t)}{p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{o}_1, \dots, \mathbf{o}_t)} \left(\prod_{\tau=t+1}^T \kappa_{\tau} \right) \hat{\beta}_{\mathbf{x}}(t) \end{aligned} \quad (230)$$

where $\kappa_t = p(\mathbf{o}_t | \mathbf{o}_1, \dots, \mathbf{o}_{t-1})$ are the same scaling factors as the ones used with the forward variables. Using again the first-order Markov property and the state conditional independence of the observations in Eqs. 51 and 52 the backward variable, $\beta_{\mathbf{x}}(t-1)$, can be rewritten in the following recursive form

$$\beta_{\mathbf{x}}(t-1) = \int p(\mathbf{x}_t = \mathbf{z}, \mathbf{o}_t, \dots, \mathbf{o}_T | \mathbf{x}_{t-1}) d\mathbf{z} \quad (231)$$

$$= \int p(\mathbf{x}_t = \mathbf{z} | \mathbf{x}_{t-1}) p(\mathbf{o}_t | \mathbf{x}_t = \mathbf{z}) p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{x}_t = \mathbf{z}) d\mathbf{z} \quad (232)$$

$$= \int p(\mathbf{x}_t = \mathbf{z} | \mathbf{x}_{t-1}) p(\mathbf{o}_t | \mathbf{x}_t = \mathbf{z}) \left(\prod_{\tau=t+1}^T \kappa_{\tau} \right) \hat{\beta}_{\mathbf{z}}(t) d\mathbf{z} \quad (233)$$

$$\hat{\beta}_{\mathbf{x}}(t-1) = \frac{1}{\kappa_t} \int p(\mathbf{x}_t = \mathbf{z} | \mathbf{x}_{t-1}) p(\mathbf{o}_t | \mathbf{x}_t = \mathbf{z}) \hat{\beta}_{\mathbf{z}}(t) d\mathbf{z} \quad (234)$$

The Kalman smoother estimates, $\hat{\mathbf{x}}(t)$ and $\hat{\boldsymbol{\Sigma}}(t)$, are the mean vector and covariance matrix of the state vector, \mathbf{x}_t , given the entire observation sequence, \mathbf{O} . Therefore it can be represented as the product of the scaled forward and backward variables as follows

$$\hat{\alpha}_{\mathbf{x}}(t) \hat{\beta}_{\mathbf{x}}(t) = p(\mathbf{x}_t | \mathbf{o}_1, \dots, \mathbf{o}_t) \frac{p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{x}_t)}{p(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | \mathbf{o}_1, \dots, \mathbf{o}_t)} \quad (235)$$

$$= p(\mathbf{x}_t | \mathbf{O}) = \mathcal{N}(\mathbf{x}_t; \hat{\mathbf{x}}(t), \hat{\boldsymbol{\Sigma}}(t)) \quad (236)$$

The backward recursion can be derived by substituting 53, 54, 217, 225, 234 and 236 into the

product of the scaled forward and backward variables and doing some algebra as follows

$$\hat{\alpha}_{\mathbf{x}}(t)\hat{\beta}_{\mathbf{x}}(t) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t)}(t), \boldsymbol{\Sigma}^{(t)}(t)) \int \mathcal{N}(\mathbf{z}; \mathbf{A}\mathbf{x}_t + \boldsymbol{\mu}^{(x)}, \boldsymbol{\Sigma}^{(x)}) \mathcal{N}(\mathbf{o}_{t+1}; \mathbf{C}\mathbf{z}, \boldsymbol{\Sigma}^{(o)}) \quad (237)$$

$$\frac{\mathcal{N}(\mathbf{z}; \hat{\mathbf{x}}(t+1), \hat{\boldsymbol{\Sigma}}(t+1))}{\kappa_{t+1}\hat{\alpha}_{\mathbf{z}}(t+1)} d\mathbf{z} \\ = \int \mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t)}(t) + \mathbf{J}(t)(\mathbf{z} - \mathbf{A}\mathbf{x}^{(t)}(t) - \boldsymbol{\mu}^{(x)}), \boldsymbol{\Sigma}^{(t)}(t) - \mathbf{J}(t)\mathbf{A}\boldsymbol{\Sigma}^{(t)}(t)) \quad (238)$$

$$\mathcal{N}(\mathbf{z}; \hat{\mathbf{x}}(t+1), \hat{\boldsymbol{\Sigma}}(t+1)) d\mathbf{z} \\ = \mathcal{N}(\mathbf{x}_t; \mathbf{x}^{(t)}(t) + \mathbf{J}(t)(\hat{\mathbf{x}}(t+1) - \mathbf{x}^{(t)}(t+1)), \quad (239)$$

$$\boldsymbol{\Sigma}^{(t)}(t) + \mathbf{J}(t)(\hat{\boldsymbol{\Sigma}}(t+1) - \boldsymbol{\Sigma}^{(t)}(t+1))\mathbf{J}'(t) \quad (240)$$

where $\mathbf{J}(t) = \boldsymbol{\Sigma}^{(t)}(t)\mathbf{A}'(\boldsymbol{\Sigma}^{(t)}(t+1))^{-1}$ and Kalman smoother recursions result obviously.

The cross covariance of two consecutive state vectors is also required. It can be obtained by the scaled forward and backward variables as follows

$$p(\mathbf{x}_{t-1}, \mathbf{x}_t | \mathcal{O}) = \frac{1}{\kappa_t} \hat{\alpha}_{\mathbf{x}}(t-1) p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{o}_t | \mathbf{x}_t) \hat{\beta}_{\mathbf{x}_t}(t) \quad (241)$$

Since the current state and the previous state are jointly Gaussian it is easy to obtain the cross covariance matrix $\hat{\boldsymbol{\Sigma}}^{(t-1)}(t) = \hat{\boldsymbol{\Sigma}}(t)\mathbf{J}'(t-1)$.

E.3 M Step

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{C}}, \hat{\boldsymbol{\mu}}^{(o)}, \hat{\boldsymbol{\Sigma}}^{(o)})$, for the observation process the following auxiliary function has to be maximised

$$\mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = \quad (242) \\ -\frac{1}{2} \sum_{t=1}^T \left(\log |\hat{\boldsymbol{\Sigma}}^{(o)}| + E \left\{ (\mathbf{o}_t - \hat{\mathbf{C}}\mathbf{x}_t - \hat{\boldsymbol{\mu}}^{(o)})' \hat{\boldsymbol{\Sigma}}^{(o)-1} (\mathbf{o}_t - \hat{\mathbf{C}}\mathbf{x}_t - \hat{\boldsymbol{\mu}}^{(o)}) \middle| \mathcal{O}, \mathcal{M} \right\} \right)$$

which is exactly the same as Eq. 173 in the M step for the factor analysis except the index is time, t , rather than a set index, j . So, the derivation follows directly from Appendix B.2 by replacing N with T and j with t .

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{A}}, \hat{\boldsymbol{\mu}}^{(x)}, \hat{\boldsymbol{\Sigma}}^{(x)})$, for the state evolution process the following auxiliary function has to be maximised

$$\mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}}) = \quad (243) \\ -\frac{1}{2} \sum_{t=2}^T \left(\log |\hat{\boldsymbol{\Sigma}}^{(x)}| + E \left\{ (\mathbf{x}_t - \hat{\mathbf{A}}\mathbf{x}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)})' \hat{\boldsymbol{\Sigma}}^{(x)-1} (\mathbf{x}_t - \hat{\mathbf{A}}\mathbf{x}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)}) \middle| \mathcal{O}, \mathcal{M} \right\} \right)$$

To find the new state evolution noise mean vector the auxiliary function in Eq. 243 can be differentiated with respect to $\hat{\boldsymbol{\mu}}^{(x)}$ as follows

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}^{(x)}} = \hat{\boldsymbol{\Sigma}}^{(x)-1} \sum_{t=2}^T (\hat{\mathbf{x}}(t) - \hat{\mathbf{A}}\hat{\mathbf{x}}(t-1) - \hat{\boldsymbol{\mu}}^{(x)}) = \mathbf{0} \quad (244)$$

$$\hat{\boldsymbol{\mu}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T (\hat{\mathbf{x}}(t) - \hat{\mathbf{A}}\hat{\mathbf{x}}(t-1)) \quad (245)$$

Differentiating Eq. 243 with respect to $\hat{\mathbf{A}}$ and setting the resulting equation to zero the new state transition matrix can be solved using only the sufficient statistics from the E step as follows

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\mathbf{A}}} = \hat{\Sigma}^{(x)-1} \sum_{t=2}^T (\hat{\mathbf{R}}^{(t-1)}(t) - \hat{\boldsymbol{\mu}}^{(x)} \hat{\mathbf{x}}'(t-1) - \hat{\mathbf{A}} \hat{\mathbf{R}}(t-1)) = \mathbf{0} \quad (246)$$

$$\begin{aligned} \hat{\mathbf{A}} &= \left(\sum_{t=2}^T \hat{\mathbf{R}}^{(t-1)}(t) - \frac{1}{T-1} \sum_{t=2}^T \hat{\mathbf{x}}(t) \sum_{t=2}^T \hat{\mathbf{x}}'(t-1) \right) \\ &\quad \left(\sum_{t=2}^T \hat{\mathbf{R}}(t-1) - \frac{1}{T-1} \sum_{t=2}^T \hat{\mathbf{x}}(t-1) \sum_{t=2}^T \hat{\mathbf{x}}'(t-1) \right)^{-1} \end{aligned} \quad (247)$$

When re-estimating the model parameters, the state evolution matrix has to be updated before the state evolution noise mean vector or they can be re-estimated simultaneously using matrix notation. By defining the following sums of the sufficient statistics

$$\begin{aligned} \boldsymbol{\Gamma}_3 &= \sum_{t=2}^T \hat{\mathbf{R}}^{(t-1)}(t) \quad , \boldsymbol{\Gamma}_4 = \sum_{t=2}^T \hat{\mathbf{R}}(t-1) \\ \boldsymbol{\zeta}_3 &= \sum_{t=2}^T \hat{\mathbf{x}}(t) \quad , \boldsymbol{\zeta}_4 = \sum_{t=2}^T \hat{\mathbf{x}}(t-1) \end{aligned} \quad (248)$$

it can be checked that the parameters $\hat{\mathbf{A}}$ and $\hat{\boldsymbol{\mu}}^{(x)}$ can be estimated simultaneously as follows

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{A}} & \hat{\boldsymbol{\mu}}^{(x)} \end{bmatrix} &= \\ \left(\sum_{t=2}^T \begin{bmatrix} \hat{\mathbf{R}}^{(t-1)}(t) & \hat{\mathbf{x}}(t) \end{bmatrix} \right) \left(\sum_{t=2}^T \begin{bmatrix} \hat{\mathbf{R}}(t-1) & \hat{\mathbf{x}}(t-1) \\ \hat{\mathbf{x}}'(t-1) & 1 \end{bmatrix} \right)^{-1} &= \begin{bmatrix} \boldsymbol{\Gamma}_3 & \boldsymbol{\zeta}_3 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Gamma}_4 & \boldsymbol{\zeta}_4 \\ \boldsymbol{\zeta}_4' & T-1 \end{bmatrix}^{-1} \end{aligned} \quad (249)$$

Now, the Schur's complement, $(\boldsymbol{\Gamma}|T-1)$, in Eq. 159 for the matrix to be inverted is

$$(\boldsymbol{\Gamma}|T-1) = \boldsymbol{\Gamma}_4 - \frac{1}{T-1} \boldsymbol{\zeta}_4 \boldsymbol{\zeta}_4' \quad (250)$$

and the first element of the matrix product in Eq. 249 is

$$\boldsymbol{\Gamma}_3 (\boldsymbol{\Gamma}|T-1)^{-1} - \frac{1}{T-1} \boldsymbol{\zeta}_3 \boldsymbol{\zeta}_4' (\boldsymbol{\Gamma}|T-1)^{-1} = \left(\boldsymbol{\Gamma}_3 - \frac{1}{T-1} \boldsymbol{\zeta}_3 \boldsymbol{\zeta}_4' \right) \left(\boldsymbol{\Gamma}_4 - \frac{1}{T} \boldsymbol{\zeta}_4 \boldsymbol{\zeta}_4' \right)^{-1} \quad (251)$$

which is exactly the ML estimate for $\hat{\mathbf{A}}$. The second element of the matrix product in Eq. 249 is

$$\begin{aligned} -\frac{1}{T-1} \boldsymbol{\Gamma}_3 (\boldsymbol{\Gamma}|T-1)^{-1} \boldsymbol{\zeta}_4 + \frac{1}{T-1} \boldsymbol{\zeta}_3 + \frac{1}{(T-1)^2} \boldsymbol{\zeta}_3 \boldsymbol{\zeta}_4' (\boldsymbol{\Gamma}|T-1)^{-1} \boldsymbol{\zeta}_4 &= \\ \frac{1}{T-1} \left(\boldsymbol{\zeta}_3 - (\boldsymbol{\Gamma}_3 (\boldsymbol{\Gamma}|T-1)^{-1} - \frac{1}{T-1} \boldsymbol{\zeta}_3 \boldsymbol{\zeta}_4' (\boldsymbol{\Gamma}|T-1)^{-1}) \boldsymbol{\zeta}_4 \right) \end{aligned} \quad (252)$$

which is exactly the ML estimate for $\hat{\boldsymbol{\mu}}^{(x)}$.

To find the new state evolution noise covariance matrix, the auxiliary function in Eq. 243 can be differentiated with respect to the inverse of $\hat{\Sigma}^{(x)}$ as follows

$$\begin{aligned} \frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\Sigma}^{(x)-1}} &= \\ \frac{1}{2} \sum_{t=2}^T \left(\hat{\Sigma}^{(x)} - E \left\{ (\mathbf{x}_t - \hat{\mathbf{A}} \mathbf{x}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)}) (\mathbf{x}_t - \hat{\mathbf{A}} \mathbf{x}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)})' \middle| \mathcal{O}, \mathcal{M} \right\} \right) &= \mathbf{0} \end{aligned} \quad (253)$$

$$\hat{\Sigma}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \left(\hat{\mathbf{R}}(t) - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \left[\hat{\mathbf{R}}^{(t-1)}(t) \hat{\mathbf{x}}(t) \right]' \right) \quad (254)$$

where the last form is obtained by noting that the terms inside the expectation in Eq. 254 can be reorganised as follows

$$\begin{aligned} & \hat{\mathbf{R}}(t) - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \left[\hat{\mathbf{R}}^{(t-1)}(t) \hat{\mathbf{x}}(t) \right]' \\ & - \left[\hat{\mathbf{R}}^{(t-1)}(t) \hat{\mathbf{x}}(t) \right] \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right]' + \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \left[\begin{array}{c} \hat{\mathbf{R}}(t-1) \hat{\mathbf{x}}(t-1) \\ \hat{\mathbf{x}}'(t-1) \quad 1 \end{array} \right] \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right]' \end{aligned} \quad (255)$$

and substituting the first term in Eq. 249 into the place of $\left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(o)} \right]$.

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\boldsymbol{\mu}}^{(i)}, \hat{\Sigma}^{(i)})$, for the initial state distribution the following auxiliary function has to be maximised

$$\mathcal{Q}_i(\mathcal{M}, \hat{\mathcal{M}}) = -\frac{1}{2} \log |\hat{\Sigma}^{(i)}| - \frac{1}{2} E \left\{ (\mathbf{x}_1 - \hat{\boldsymbol{\mu}}^{(i)})' \hat{\Sigma}^{(i)-1} (\mathbf{x}_1 - \hat{\boldsymbol{\mu}}^{(i)}) \mid \mathcal{O}, \mathcal{M} \right\} \quad (256)$$

Differentiating 256 with respect to the initial state mean vector $\hat{\boldsymbol{\mu}}^{(i)}$ yields

$$\frac{\partial \mathcal{Q}_i(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}^{(i)}} = \hat{\Sigma}^{(i)-1} (\hat{\mathbf{x}}(1) - \hat{\boldsymbol{\mu}}^{(i)}) = \mathbf{0} \quad (257)$$

$$\hat{\boldsymbol{\mu}}^{(i)} = \hat{\mathbf{x}}(1) \quad (258)$$

and with respect to the inverse of the initial state covariance matrix $\hat{\Sigma}^{(i)}$ yields

$$\frac{\partial \mathcal{Q}_i(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\Sigma}^{(i)-1}} = \frac{1}{2} \hat{\Sigma}^{(i)} - \frac{1}{2} (\hat{\mathbf{R}}(1) - \hat{\mathbf{x}}(1) \hat{\boldsymbol{\mu}}^{(i)'} - \hat{\boldsymbol{\mu}}^{(i)} \hat{\mathbf{x}}(1)' + \hat{\boldsymbol{\mu}}^{(i)} \hat{\boldsymbol{\mu}}^{(i)'}) = \mathbf{0} \quad (259)$$

$$\hat{\Sigma}^{(i)} = \hat{\mathbf{R}}(1) - \hat{\boldsymbol{\mu}}^{(i)} \hat{\boldsymbol{\mu}}^{(i)'} \quad (260)$$

F Derivation of the EM Algorithm for FAHMM

The expectation maximisation algorithm for factor analysed hidden Markov models is presented in this appendix. The E step is very similar to the hidden Markov models except for the different observation posteriors given the current state. The M step is also a combination of the hidden Markov model, the mixtures of factors and shared factor analysis provided the correct statistics are used.

F.1 E Step

The state posterior $p(\mathbf{x}_t | j, n)$ and the observation posterior $p(\mathbf{o}_t | \mathbf{x}_t, j, m)$ are defined as the following two Gaussians

$$p(\mathbf{x}_t | j, n) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{jn}^{(x)}, \boldsymbol{\Sigma}_{jn}^{(x)}) \quad (261)$$

$$p(\mathbf{o}_t | \mathbf{x}_t, j, m) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}_j \mathbf{x}_t + \boldsymbol{\mu}_{jm}^{(o)}, \boldsymbol{\Sigma}_{jm}^{(o)}) \quad (262)$$

where $\boldsymbol{\mu}_{jn}^{(x)}$ and $\boldsymbol{\Sigma}_{jn}^{(x)}$ are the mean vector and covariance matrix associated with the state j and mixture component n . \mathbf{C}_j is a p by k observation matrix and $\boldsymbol{\mu}_{jm}^{(o)}$, $\boldsymbol{\Sigma}_{jm}^{(o)}$ are the observation noise mean vector and covariance matrix. Since the state vectors and observations are also jointly Gaussian given the state and the mixture components it is easy to obtain the posterior distribution of an observation \mathbf{o}_t using again the conditioning of multivariate Gaussians as follows

$$p(\mathbf{o}_t | j, m, n) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}_j \boldsymbol{\mu}_{jn}^{(x)} + \boldsymbol{\mu}_{jm}^{(o)}, \mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' + \boldsymbol{\Sigma}_{jm}^{(o)}) \quad (263)$$

It should be noted that various independence assumptions have been applied. Now, the likelihood of the observation given only the state can be marginalised by using the mixture priors and summing over the mixture components as follows

$$b_j(\mathbf{o}_t) = \sum_{m=1}^{M^{(o)}} c_{jm}^{(o)} \sum_{n=1}^{M^{(x)}} c_{jn}^{(x)} p(\mathbf{o}_t | j, m, n) \quad (264)$$

The state posteriors, $\gamma_j(t)$ and $\xi_{ij}(t)$, can be obtained using the forward-backward algorithm in the usual way. To find the joint likelihood of the state and the state and observation mixture components, $\gamma_{jmn}(t)$, the forward and backward variables are used as usual

$$\gamma_{jmn}(t) = \frac{1}{p(\mathbf{O})} c_{jm}^{(o)} c_{jn}^{(x)} b_{jmn}(\mathbf{o}_t) \sum_{i=1}^{N_s} a_{ij} \alpha_i(t-1) \beta_j(t) \quad (265)$$

where $b_{jmn}(\mathbf{o}_t) = p(\mathbf{o}_t | j, m, n)$.

The first and second-order sufficient statistics $\hat{\mathbf{x}}_{jmn}(t) = E\{\mathbf{x}_t | j, m, n, \mathbf{O}\}$ and $\hat{\mathbf{R}}_{jmn}(t) = E\{\mathbf{x}_t \mathbf{x}_t' | j, m, n, \mathbf{O}\}$ for the M step can be obtained by using the independence assumption of HMMs as follows

$$p(\mathbf{x}_t | j, m, n, \mathbf{O}) = \frac{p(\mathbf{o}_t, \mathbf{x}_t | j, m, n)}{p(\mathbf{o}_t | j, m, n)} = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{jn}^{(x)} + \mathbf{K}_{jmn}(\mathbf{o}_t - \mathbf{C}_j \boldsymbol{\mu}_{jn}^{(x)} - \boldsymbol{\mu}_{jm}^{(o)}), \boldsymbol{\Sigma}_{jn}^{(x)} - \mathbf{K}_{jmn} \mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)}) \quad (266)$$

where $\mathbf{K}_{jmn} = \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' (\mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} \mathbf{C}_j' + \boldsymbol{\Sigma}_{jm}^{(o)})^{-1}$. Again, the joint Gaussianity of the observations and the state vectors were used as well as the conditioning of multivariate Gaussians as described in Appendix A.2. The statistics are then

$$\hat{\mathbf{x}}_{jmn}(t) = \boldsymbol{\mu}_{jn}^{(x)} + \mathbf{K}_{jmn}(\mathbf{o}_t - \mathbf{C}_j \boldsymbol{\mu}_{jn}^{(x)} - \boldsymbol{\mu}_{jm}^{(o)}) \quad (267)$$

$$\hat{\mathbf{R}}_{jmn}(t) = \boldsymbol{\Sigma}_{jn}^{(x)} - \mathbf{K}_{jmn} \mathbf{C}_j \boldsymbol{\Sigma}_{jn}^{(x)} + \hat{\mathbf{x}}_{jmn}(t) \hat{\mathbf{x}}_{jmn}'(t) \quad (268)$$

F.2 M Step

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\pi}_j, \hat{a}_{ij})$, for the HMM state transition matrix the following auxiliary function has to be maximised

$$\mathcal{Q}_s(\mathcal{M}, \hat{\mathcal{M}}) = \sum_{j=1}^{N_s} P_{q_1}(j | \mathbf{O}, \mathcal{M}) \log \hat{\pi}_j + \sum_{t=2}^T \sum_{i=1}^{N_s} \sum_{j=1}^{N_s} P_{q_{t-1}, q_t}(i, j | \mathbf{O}, \mathcal{M}) \log \hat{a}_{ij} \quad (269)$$

It is not possible to maximise the above expression by setting its first derivative with respect to the parameters to zero but Lagrange multipliers can be used since $\hat{\pi}_j$ and \hat{a}_{ij} are subject to constraints $\sum_{j=1}^{N_s} \hat{\pi}_j = 1$ and $\forall i : \sum_{j=1}^{N_s} \hat{a}_{ij} = 1$, respectively. The derivation follows closely to the derivation of the new mixture priors in a mixture of factor analysers presented in Appendix C.2. Defining $\gamma_j(t) = P_{q_t}(j | \mathbf{O}, \mathcal{M})$ and $\xi_{ij}(t) = P_{q_{t-1}, q_t}(i, j | \mathbf{O}, \mathcal{M})$ the ML estimates for $\hat{\pi}_j$ and \hat{a}_{ij} are

$$\hat{\pi}_j = \frac{\gamma_j(1)}{\sum_{i=1}^{N_s} \gamma_i(1)} \quad (270)$$

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T \xi_{ij}(t)}{\sum_{t=2}^T \gamma_i(t-1)} \quad (271)$$

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\boldsymbol{\mu}}_{jn}^{(x)}, \hat{\boldsymbol{\Sigma}}_{jn}^{(x)})$, for the HMM state vector distributions the following auxiliary function has to be maximised

$$\begin{aligned} \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}}) = & \quad (272) \\ & -\frac{1}{2} \sum_{t=1}^T \sum_{j=1}^{N_s} \sum_{m=1}^{M^{(o)}} \sum_{n=1}^{M^{(x)}} \gamma_{jmn}(t) \left(\log |\hat{\boldsymbol{\Sigma}}_{jn}^{(x)}| + E \left\{ (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jn}^{(x)})' \hat{\boldsymbol{\Sigma}}_{jn}^{(x)-1} (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jn}^{(x)}) \middle| \mathbf{O}, j, m, n, \mathcal{M} \right\} \right) \end{aligned}$$

where $\gamma_{jmn}(t) = P_{q_t, \omega^o, \omega^x}(j, m, n | \mathbf{O}, \mathcal{M})$.

Differentiating the auxiliary function with respect to the state mean vectors yields

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}_{jn}^{(x)}} = \hat{\boldsymbol{\Sigma}}_{jn}^{(x)-1} \sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) (\hat{\boldsymbol{\mu}}_{jn}^{(x)} - \hat{\mathbf{x}}_{jmn}(t)) = \mathbf{0} \quad (273)$$

$$\hat{\boldsymbol{\mu}}_{jn}^{(x)} = \frac{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) \hat{\mathbf{x}}_{jmn}(t)}{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t)} \quad (274)$$

and with respect to the inverse covariance matrices yields

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\Sigma}}_{jn}^{(x)-1}} = \quad (275)$$

$$\frac{1}{2} \sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) \left(\hat{\boldsymbol{\Sigma}}_{jn}^{(x)} - (\hat{\mathbf{R}}_{jmn}(t) - \hat{\mathbf{x}}_{jmn}(t) \hat{\boldsymbol{\mu}}_{jn}^{(x)'} - \hat{\boldsymbol{\mu}}_{jn}^{(x)} \hat{\mathbf{x}}_{jmn}'(t) + \hat{\boldsymbol{\mu}}_{jn}^{(x)} \hat{\boldsymbol{\mu}}_{jn}^{(x)'}) \right) = \mathbf{0}$$

$$\hat{\boldsymbol{\Sigma}}_{jn}^{(x)} = \frac{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t) \hat{\mathbf{R}}_{jmn}(t)}{\sum_{t=1}^T \sum_{m=1}^{M^{(o)}} \gamma_{jmn}(t)} - \hat{\boldsymbol{\mu}}_{jn}^{(x)} \hat{\boldsymbol{\mu}}_{jn}^{(x)'} \quad (276)$$

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{C}}_j, \hat{\boldsymbol{\mu}}_{jm}^{(o)}, \hat{\boldsymbol{\Sigma}}_{jm}^{(o)})$, for the observation process, exactly the same auxiliary function as in case of shared factor analysis in Eq. 206 has to be maximised provided the sufficient statistics are replaced accordingly and the state-space component is marginalised by summing over n .

G Optimisation of Dynamical Linear Discriminant Analysis

The algorithm for the optimisation of the parameters of a dynamical linear discriminant analysis is presented in this appendix. The algorithm differs from the expectation maximisation algorithm presented in previous chapters in that there are no latent variables present.

G.1 Observation Matrix Optimisation

To find the new transformation matrix, $\hat{\mathbf{C}}$, the following auxiliary function has to be maximised

$$\begin{aligned} \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = & \quad (277) \\ & \frac{1}{2} \sum_{t=2}^T \left(\log(|\hat{\mathbf{B}}'|^2) - (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \mathbf{A} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \boldsymbol{\mu}^{(x)})' \boldsymbol{\Sigma}^{(x)-1} (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \mathbf{A} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \boldsymbol{\mu}^{(x)}) \right. \\ & \left. - (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t - \boldsymbol{\mu}^{(o)})' \boldsymbol{\Sigma}^{(o)-1} (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t - \boldsymbol{\mu}^{(o)}) \right) \end{aligned}$$

where $\hat{\mathbf{B}} = (\hat{\mathbf{C}}^{-1})'$. The optimisation can now be done one column of $\hat{\mathbf{B}}$ at a time easily if diagonal transition matrix is used. The auxiliary function can be then rewritten as follows

$$\mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}}) = \tag{278}$$

$$(T-1) \log(\hat{\mathbf{b}}'_i \mathbf{c}_i) - \frac{1}{2} \sum_{j=1}^k (\hat{\mathbf{b}}'_j \mathbf{G}_j^{(x)} \hat{\mathbf{b}}_j - 2\hat{\mathbf{b}}'_j \mathbf{k}_j^{(x)}) - \frac{1}{2} \sum_{j=k+1}^p (\hat{\mathbf{b}}'_j \mathbf{G}_j^{(o)} \hat{\mathbf{b}}_j - 2\hat{\mathbf{b}}'_j \mathbf{k}_j^{(o)})$$

where \mathbf{c}_i is a vector of the cofactors of the matrix $\hat{\mathbf{B}}$ corresponding to the elements in the vector $\hat{\mathbf{b}}_i$ and

$$\mathbf{G}_j^{(x)} = \frac{1}{\sigma_j^{(x)2}} \sum_{t=2}^T (\mathbf{o}_t \mathbf{o}'_t - a_{jj} \mathbf{o}_t \mathbf{o}'_{t-1} - a_{jj} \mathbf{o}_{t-1} \mathbf{o}_t + a_{jj}^2 \mathbf{o}_{t-1} \mathbf{o}'_{t-1}) \tag{279}$$

$$\mathbf{k}_j^{(x)} = \frac{\mu_j^{(x)2}}{\sigma_j^{(x)2}} \sum_{t=2}^T (\mathbf{o}_t - a_{jj} \mathbf{o}_{t-1}) \tag{280}$$

$$\mathbf{G}_j^{(o)} = \frac{1}{\sigma_{(j-k)}^{(o)2}} \sum_{t=2}^T \mathbf{o}_t \mathbf{o}'_t \tag{281}$$

$$\mathbf{k}_j^{(o)} = \frac{\mu_{(j-k)}^{(o)2}}{\sigma_{(j-k)}^{(o)2}} \sum_{t=2}^T \mathbf{o}_t \tag{282}$$

where $\mu_j^{(x)2}$ and $\mu_j^{(o)2}$ are the j th elements of the observation and state evolution noise mean vectors $\boldsymbol{\mu}^{(x)}$ and $\boldsymbol{\mu}^{(o)}$, respectively, $\sigma_j^{(x)2}$, $\sigma_j^{(o)2}$ and a_{jj} are the j th diagonal elements of the observation and state evolution noise covariance, and the state transition matrices, respectively.

For the rows $1 \leq i \leq k$, differentiating the auxiliary function with respect to $\hat{\mathbf{b}}_i$ yields

$$\frac{\partial \mathcal{Q}_o(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\mathbf{b}}_i} = (T-1) \frac{\mathbf{c}_i}{\hat{\mathbf{b}}'_i \mathbf{c}_i} - \mathbf{G}_i^{(x)} \hat{\mathbf{b}}_i + \mathbf{k}_i^{(x)} = \mathbf{0} \tag{283}$$

and after rearranging yields

$$\hat{\mathbf{b}}_i \hat{\mathbf{b}}'_i \mathbf{c}_i = \mathbf{G}_i^{(x)-1} \mathbf{k}_i^{(x)} \hat{\mathbf{b}}'_i \mathbf{c}_i + (T-1) \mathbf{G}_i^{(x)-1} \mathbf{c}_i \tag{284}$$

Since $\hat{\mathbf{b}}'_i \mathbf{c}_i$ is a scalar and considering the direction of the vector $\hat{\mathbf{b}}_i$, it can be represented as follows

$$\hat{\mathbf{b}}_i = \delta \mathbf{G}_i^{(x)-1} (\lambda \mathbf{k}_i^{(x)} + \mathbf{c}_i) \tag{285}$$

and the values for δ and λ have to be found. Substituting this expression into Eq. 284, multiplying by $\mathbf{G}_i^{(x)}$ from the left and rearranging yields

$$\delta(\lambda\delta - 1)(\lambda \mathbf{k}_i^{(x)'} + \mathbf{c}'_i) \mathbf{G}_i^{(x)-1} \mathbf{c}_i \mathbf{k}_i^{(x)} = ((T-1) - \delta^2 (\lambda \mathbf{k}_i^{(x)'} + \mathbf{c}'_i) \mathbf{G}_i^{(x)-1} \mathbf{c}_i) \mathbf{c}_i \tag{286}$$

which holds always if and only if

$$\lambda\delta = 1 \tag{287}$$

$$(T-1) = \delta^2 (\lambda \mathbf{k}_i^{(x)'} + \mathbf{c}'_i) \mathbf{G}_i^{(x)-1} \mathbf{c}_i \tag{288}$$

and these can be expressed as one simple quadratic equation

$$\delta^2 \mathbf{c}'_i \mathbf{G}_i^{(x)-1} \mathbf{c}_i + \delta \mathbf{k}_i^{(x)'} \mathbf{G}_i^{(x)-1} \mathbf{c}_i - (T-1) = 0 \tag{289}$$

It can be shown that both roots of this equation are maxima. Now that knowing δ and λ , the vector $\hat{\mathbf{b}}_i$ can be expressed as

$$\hat{\mathbf{b}}_i = \mathbf{k}_i^{(x)} \mathbf{G}_i^{(x)-1} (\mathbf{k}_i^{(x)} + \delta \mathbf{c}_i) \quad (290)$$

and the auxiliary function for the i th row ignoring all the terms independent of δ can be represented as follows

$$\mathcal{Q}_o^{(i)}(\mathcal{M}, \hat{\mathcal{M}}) = (T-1) \log(|\delta \epsilon_1 + \epsilon_2|) - \frac{1}{2} \delta^2 \epsilon_1 \quad (291)$$

where $\epsilon_1 = \mathbf{c}_i' \mathbf{G}_i^{(x)-1} \mathbf{c}_i$ and $\epsilon_2 = \mathbf{k}_i^{(x)'} \mathbf{G}_i^{(x)-1} \mathbf{c}_i$. Since it is not possible to ensure that $\epsilon_2 > 0$, the value of δ is selected that maximises $\mathcal{Q}_o^{(i)}(\mathcal{M}, \hat{\mathcal{M}})$.

For the rows $k+1 \leq i \leq p$, the optimisation is exactly the same as above except for using the values $\mathbf{G}_i^{(o)}$ and $\mathbf{k}_i^{(o)}$ instead of $\mathbf{G}_i^{(x)}$ and $\mathbf{k}_i^{(x)}$.

G.2 Observation Noise Parameter Optimisation

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\boldsymbol{\Sigma}}^{(o)}, \hat{\boldsymbol{\mu}}^{(o)})$, for the observation noise, the following auxiliary function has to be maximised

$$\mathcal{Q}_v(\mathcal{M}, \hat{\mathcal{M}}) = -\frac{1}{2} \sum_{t=2}^T \left(\log |\hat{\boldsymbol{\Sigma}}^{(o)}| + (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t - \hat{\boldsymbol{\mu}}^{(o)})' \hat{\boldsymbol{\Sigma}}^{(o)-1} (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t - \hat{\boldsymbol{\mu}}^{(o)}) \right) \quad (292)$$

To find the new observation noise mean vector the auxiliary function in Eq. 292 has to be differentiated with respect to $\hat{\boldsymbol{\mu}}^{(o)}$ as follows

$$\frac{\partial \mathcal{Q}_v(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}^{(o)}} = \hat{\boldsymbol{\Sigma}}^{(o)-1} \sum_{t=2}^T (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t - \hat{\boldsymbol{\mu}}^{(o)}) = \mathbf{0} \quad (293)$$

$$\hat{\boldsymbol{\mu}}^{(o)} = \frac{1}{T-1} \hat{\mathbf{B}}'_{[p-k]} \sum_{t=2}^T \mathbf{o}_t \quad (294)$$

To find the new observation noise covariance matrix, the auxiliary function in Eq. 292 has to be differentiated with respect to the inverse of $\hat{\boldsymbol{\Sigma}}^{(o)}$ as follows

$$\frac{\partial \mathcal{Q}_v(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\Sigma}}^{(o)-1}} = \quad (295)$$

$$\frac{1}{2} \sum_{t=2}^T \left(\hat{\boldsymbol{\Sigma}}^{(o)} - (\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t \mathbf{o}_t' \hat{\mathbf{B}}_{[p-k]} - \hat{\boldsymbol{\mu}}^{(o)} \mathbf{o}_t' \hat{\mathbf{B}}_{[p-k]} - \hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t \hat{\boldsymbol{\mu}}^{(o)'} + \hat{\boldsymbol{\mu}}^{(o)} \hat{\boldsymbol{\mu}}^{(o)'}) \right) = \mathbf{0}$$

$$\hat{\boldsymbol{\Sigma}}^{(o)} = \frac{1}{T-1} \sum_{t=2}^T \text{diag}(\hat{\mathbf{B}}'_{[p-k]} \mathbf{o}_t \mathbf{o}_t' \hat{\mathbf{B}}_{[p-k]} - \boldsymbol{\mu}^{(o)} \boldsymbol{\mu}^{(o)'}) \quad (296)$$

G.3 State Evolution Process Parameter Optimisation

To find the new parameters, $\hat{\mathcal{M}} = (\hat{\mathbf{A}}, \hat{\boldsymbol{\Sigma}}^{(x)}, \hat{\boldsymbol{\mu}}^{(x)})$, for the state evolution process the following auxiliary function has to be maximised

$$\mathcal{Q}_s(\mathcal{M}, \hat{\mathcal{M}}) = \quad (297)$$

$$-\frac{1}{2} \sum_{t=2}^T \left(\log |\hat{\boldsymbol{\Sigma}}^{(x)}| + (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)})' \hat{\boldsymbol{\Sigma}}^{(o)-1} (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)}) \right)$$

To find the new state evolution noise mean vector, the auxiliary function in Eq. 297 can be differentiated with respect to $\hat{\boldsymbol{\mu}}^{(x)}$ as follows

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\mu}}^{(x)}} = \hat{\boldsymbol{\Sigma}}^{(x)-1} \sum_{t=2}^T (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)}) = \mathbf{0} \quad (298)$$

$$\hat{\boldsymbol{\mu}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1}) \quad (299)$$

Differentiating Eq. 297 with respect to $\hat{\mathbf{A}}$ and setting the resulting equation to zero, the new state transition matrix can be solved remembering the diagonality assumption

$$\frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\mathbf{A}}} = \hat{\boldsymbol{\Sigma}}^{(x)-1} \sum_{t=2}^T (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \mathbf{o}'_{t-1} \hat{\mathbf{B}}_{[k]} - \hat{\boldsymbol{\mu}}^{(x)} \mathbf{o}'_{t-1} \hat{\mathbf{B}}_{[k]} - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} \mathbf{o}'_{t-1} \hat{\mathbf{B}}_{[k]}) = \mathbf{0} \quad (300)$$

$$\begin{aligned} \hat{\mathbf{A}} = & \left(\text{diag} \left(\hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_t \mathbf{o}'_{t-1} \right) \hat{\mathbf{B}}_{[k]} - \frac{1}{T-1} \hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_t \right) \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right)' \hat{\mathbf{B}}_{[k]} \right) \right) \\ & \left(\text{diag} \left(\hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_{t-1} \mathbf{o}'_{t-1} \right) \hat{\mathbf{B}}_{[k]} - \frac{1}{T-1} \hat{\mathbf{B}}'_{[k]} \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right) \left(\sum_{t=2}^T \mathbf{o}_{t-1} \right)' \hat{\mathbf{B}}_{[k]} \right) \right)^{-1} \end{aligned} \quad (301)$$

which follows like the state transition matrix in case of linear dynamical systems in Appendix E.3.

To find the new state evolution noise covariance matrix, the auxiliary function in Eq. 297 can be differentiated with respect to the inverse of $\hat{\boldsymbol{\Sigma}}^{(x)}$ as follows

$$\begin{aligned} \frac{\partial \mathcal{Q}_x(\mathcal{M}, \hat{\mathcal{M}})}{\partial \hat{\boldsymbol{\Sigma}}^{(x)-1}} = & \quad (302) \\ \frac{1}{2} \sum_{t=2}^T & \left(\hat{\boldsymbol{\Sigma}}^{(x)} - (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)}) (\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t - \hat{\mathbf{A}} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_{t-1} - \hat{\boldsymbol{\mu}}^{(x)})' \right) = \mathbf{0} \end{aligned}$$

$$\hat{\boldsymbol{\Sigma}}^{(x)} = \frac{1}{T-1} \sum_{t=2}^T \text{diag} \left(\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \mathbf{o}'_t \hat{\mathbf{B}}_{[k]} - \left[\hat{\mathbf{A}} \hat{\boldsymbol{\mu}}^{(x)} \right] \left[\hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \mathbf{o}'_{t-1} \hat{\mathbf{B}}_{[k]} \hat{\mathbf{B}}'_{[k]} \mathbf{o}_t \right]' \right) \quad (303)$$

which of course holds only if diagonal covariance matrices are used.