

4G10 BMI: notes for lecture 2

linear dynamical systems

Yashar Ahmadian

October 16, 2023

The models we discussed for dimensionality in Lecture 2 ignored temporal inter-dependencies that exist in the neural data. This is certainly a problematic assumption from the point of view of generative or causal modeling, considering that neural activity results from the *dynamics* of recurrently connected neural circuits; but it is also problematic as it leads to sub-optimal inferences that can significantly hinder the proper operation of a BMI. To see this, consider a (hypothetical) population of neurons in your motor cortex that represents the momentary speed of your hand along a given direction. As we tend to perform smooth movements with hand speed changing continuously in time, the latent state \mathbf{z}_t that we infer at time t should be constrained not only by the neural activity \mathbf{x}_t observed at time t , but also by other temporally adjacent neural observations (e.g. \mathbf{x}_{t-1} and \mathbf{x}_{t+1}), because of their statistical dependence (in the language of information theory, not only \mathbf{x}_t but also neural activity at other times contain information about \mathbf{z}_t). In other words, latent state inference ought to involve some form of ‘filtering’ or ‘smoothing’ of the observations.

The GPFA model discussed in Lecture 3 does take into account temporal inter-dependencies. However, the specific form of the Gaussian process prior on $\mathbf{z}_{1:T}$ in that specific model: (a) does not account for possible cross-correlations/interactions between different components of \mathbf{z} , and (b) hides the dynamical interactions responsible for creating the temporal dependencies. The model class we discuss below explicitly incorporates, in the generative model, the assumption of temporal continuity and smoothness in the underlying latent state trajectories, and moreover does so by explicit formulation in terms of dynamical systems. Bayesian inference on such a generative model then automatically uses “surrounding observations” to improve the inference of the latent states.

1 Gaussian linear dynamical systems

We start by describing the generative model and then, in the next section, discuss inference based on that model. As in the models in previous lectures, the generative model is specified by its two components: the prior on the latent trajectories, and the observation model relating the latent variables to the observed ones.

The generative prior on the latent trajectory. A simple form of prior over latent states that captures temporal continuity is the prior induced by linear stochastic dynamics of the form:

$$\mathbf{z}_{t+1} = A\mathbf{z}_t + B\boldsymbol{\epsilon}_{t+1}, \quad \boldsymbol{\epsilon}_{t+1} \sim \mathcal{N}(0, I_M), \quad (1)$$

with the initial condition distribution

$$\mathbf{z}_1 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0). \quad (2)$$

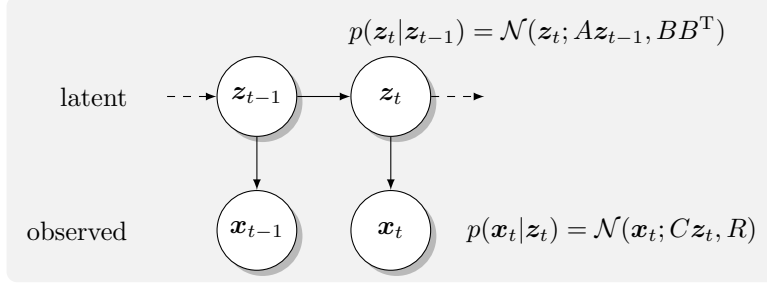


Figure 1: Graphical model associated with the LDS model.

Here A is the so-called “state matrix” (which describes the interaction between different components of \mathbf{z}_t), B is an “input matrix”, and ϵ_{t+1} is Gaussian white¹ noise. In other words, $p(\mathbf{z}_{t+1}|\mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t+1}; A\mathbf{z}_t, BB^T)$. Note that under the above model, $\mathbf{z}_{1:T}$ form a Markov chain.²

The system described above is mathematically a stochastic linear dynamical system; however, in the context of generative models, the term “**linear dynamical system**” (**LDS**) refers to the full generative model obtained when we combine the above prior with the FA-like Gaussian observation model, that we will now introduce.

The generative observation model. Similar to GPFA, we use a linear Gaussian model with independent observation noise, in which conditioned on \mathbf{z}_t , \mathbf{x}_t is independent of \mathbf{z} at other times. The dependence of \mathbf{x}_t on \mathbf{z}_t is given by

$$\mathbf{x}_t | \mathbf{z}_t \sim \mathcal{N}(C\mathbf{z}_t, R) \quad (3)$$

where R is a diagonal matrix.

The full generative model is described graphically in [Figure 1](#). Note that this model thus constitutes a hidden Markov model (HMM).³

1.1 Inference

In LDSs, two different inference problems arise. The first is the “filtering” problem of inferring the latent state \mathbf{z}_t based only on the observations *up until time t*. The second is called “smoothing”, and consists of inferring \mathbf{z}_t based on *all* observations, including future ones. While filtering is naturally more relevant to online BMI decoding applications, smoothing is an important building block for learning the parameters of the generative model.

Kalman filtering The filtering distribution $p(\mathbf{z}_t | \mathbf{x}_{0:t})$ is also known as the Kalman filter, which you will have come across in other 3F/4F modules. To derive the Kalman filter, we proceed by induction: we assume

¹Both spatially and temporally; *i.e.*, the noise vectors at different times are independent.

²Furthermore, since Gaussian variables remain Gaussian under linear transformations, it can be checked that the joint distribution of the sequence $\mathbf{z}_{1:T}$ is multivariate Gaussian; with zero mean and a $TM \times TM$ covariance matrix. This observation can also be extended to continuous time Gaussian LDS. Thus, mathematically speaking, the Gaussian LDS described above, or their continuous time generalizations, constitute a subclass of Gaussian processes, but of a form that does not fall in the class of GP’s employed in GPFA. For various reasons, LDS are not typically formulated in the language used for discussing GP’s, namely in terms of a covariance kernel.

³But prior exposition to or experience with HMM’s is not required – these notes are self-contained.

that $p(\mathbf{z}_t|\mathbf{x}_{0:t})$ is known and equal to $\mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_t, \Sigma_t)$, and go on showing that $p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t+1})$ is also Gaussian with moments $\boldsymbol{\mu}_{t+1}$ and Σ_{t+1} which we can obtain as a function of $\boldsymbol{\mu}_t$, Σ_t and \mathbf{x}_{t+1} :

$$p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t+1}) = p(\mathbf{z}_{t+1}|\mathbf{x}_{t+1}, \mathbf{x}_{0:t}) \propto p(\mathbf{x}_{t+1}|\mathbf{z}_{t+1}, \mathbf{x}_{0:t}) p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t}) \quad (\text{Bayes' theorem}) \quad (4)$$

filtering distrib

$$= p(\mathbf{x}_{t+1}|\mathbf{z}_{t+1}) p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t}) \quad (\text{conditional independence of observations}) \quad (5)$$

$$= p(\mathbf{x}_{t+1}|\mathbf{z}_{t+1}) \int p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{x}_{0:t}) p(\mathbf{z}_t|\mathbf{x}_{0:t}) d\mathbf{z}_t \quad (\text{sum and product rules}) \quad (6)$$

$$= p(\mathbf{x}_{t+1}|\mathbf{z}_{t+1}) \int p(\mathbf{z}_{t+1}|\mathbf{z}_t) p(\mathbf{z}_t|\mathbf{x}_{0:t}) d\mathbf{z}_t \quad (\text{conditional independence: Markov property}) \quad (7)$$

Note that the above is valid for any HMM. For our Gaussian LDS model, substituting the Gaussian distributions for the observation and transition distributions, and using the hypothesis of our induction, we obtain

$$p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t+1}) \propto \mathcal{N}(\mathbf{x}_{t+1}; C\mathbf{z}_{t+1}, R) \int \mathcal{N}(\mathbf{z}_{t+1}; A\mathbf{z}_t, BB^T) \mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_t, \Sigma_t) d\mathbf{z}_t. \quad (8)$$

Luckily, we do not have to explicitly evaluate the integral in the last equation: first note that (as in our earlier calculation of the marginal likelihood in pPCA) this integral is nothing but an expression for the marginal density of $\mathbf{z}_{t+1} = A\mathbf{z}_t + B\boldsymbol{\epsilon}_{t+1}$, with \mathbf{z}_t distributed according to $\mathbf{z}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)$ (by the assumption of our induction) and $\boldsymbol{\epsilon}_{t+1} \sim \mathcal{N}(\mathbf{0}, I_M)$. As \mathbf{z}_{t+1} is a linear combination of two Gaussian random variables, it is itself Gaussian (*i.e.*, its marginal distribution is Gaussian). In particular, it suffices to evaluate the mean and covariance matrix of \mathbf{z}_{t+1} in order to determine its marginal (Gaussian) distribution. Since $\boldsymbol{\epsilon}_{t+1}$ has zero mean, the mean of \mathbf{z}_{t+1} is $A\boldsymbol{\mu}_t$, and since $\boldsymbol{\epsilon}_{t+1}$ is independent of \mathbf{z}_t , the covariance matrix of \mathbf{z}_{t+1} is the sum of the covariance matrices of the two terms $A\mathbf{z}_t$ and $B\boldsymbol{\epsilon}_{t+1}$, which is $A\Sigma_t A^T + BB^T$. Thus, we can rewrite Equation (8) as

$$p(\mathbf{z}_{t+1}|\mathbf{x}_{0:t+1}) \propto \mathcal{N}(\mathbf{x}_{t+1}; C\mathbf{z}_{t+1}, R) \mathcal{N}(\mathbf{z}_{t+1}; A\boldsymbol{\mu}_t, A\Sigma_t A^T + BB^T). \quad (9)$$

It is easy to see that the right hand side is again a Gaussian density over \mathbf{z}_{t+1} (because both factors are exponentials of quadratic expressions in \mathbf{z}_{t+1}). The mean, $\boldsymbol{\mu}_{t+1}$, and covariance, Σ_{t+1} , can then be obtained using the same type of calculation as those that led to the expression for mean and covariance of the posterior distribution in pPCA (lecture 2). The result is

$$\Sigma_{t+1}^{-1} = C^T R^{-1} C + (A\Sigma_t A^T + BB^T)^{-1} \boldsymbol{\Sigma} \quad (10)$$

$$\boldsymbol{\mu}_{t+1} = \Sigma_{t+1} [(A\Sigma_t A^T + BB^T)^{-1} A\boldsymbol{\mu}_t + C^T R^{-1} \mathbf{x}_{t+1}] \quad (11)$$

Given that $p(\mathbf{z}_0)$ is Gaussian, inductive reasoning allows us to conclude that the filtering distribution is indeed Gaussian at all times. Finally, Equation (10) can be simplified using the Woodbury identity $(A + UBU^T)^{-1} = A^{-1} - A^{-1}U(B^{-1} + U^T A^{-1}U)^{-1}U^T A^{-1}$, leading to a more interpretable form:

$$\boldsymbol{\mu}_{t+1} = \underbrace{A\boldsymbol{\mu}_t}_{\text{prediction}} + \underbrace{\Sigma_{t+1} C^T R^{-1}}_{\text{Kalman gain}} \underbrace{(\mathbf{x}_{t+1} - CA\boldsymbol{\mu}_t)}_{\text{prediction error}} \quad (12)$$

Without the “correction” term $\boldsymbol{\mu}_t$ would evolve according to a deterministic (noise-free) version of the latent linear dynamics. This is indeed how the marginal means of \mathbf{z}_t evolve under the *prior* distribution on the latent trajectory. However, *a posteriori* (*i.e.*, after having observed $\mathbf{x}_{1:T}$) this prior prediction has to be corrected depending on the new observation \mathbf{x}_{t+1} , and according to Equation (12), the correction is proportional to the prediction error $\mathbf{x}_{t+1} - CA\boldsymbol{\mu}_t$.

$$\Sigma_{t+1}^{-1} = \left(\Sigma_t^{-1} + C^T R^{-1} C \right)^{-1} = \Sigma_t - \Sigma_t C^T (R + C \Sigma_t C^T)^{-1} C \Sigma_t$$

$$\Sigma_{t+1} = \Sigma_t - K_t C \Sigma_t = (I - K_t C) \Sigma_t$$

$$\begin{aligned} \boldsymbol{\mu}_{t+1} &= (I - K_t C) \Sigma_t \left[\Sigma_t^{-1} A \boldsymbol{\mu}_t + C^T R^{-1} \mathbf{x}_{t+1} \right] \\ &= (I - K_t C) A \boldsymbol{\mu}_t + (I - K_t C) \Sigma_t C^T R^{-1} \mathbf{x}_{t+1} \end{aligned}$$

$$= \left[I - \underbrace{\Sigma_t C^T (R + C \Sigma_t C^T)^{-1} C}_{K_t} \right] \Sigma_t$$

$$\begin{aligned}
&= A\mu_t + \Sigma_{t+1} C^T R^{-1} x_{t+1} - K_t C A \mu_t \\
&\therefore \mu_{t+1} = A\mu_t + K_t (x_{t+1} - C A \mu_t)
\end{aligned}$$

$K_t = \Sigma_{t+1} C^T R^{-1}$

$$\begin{aligned}
&K_t \cancel{=} \\
&= I - \Sigma_{t+1} \Sigma^{-1} \\
&= I - (\Sigma_t^{-1} + C^T R^{-1} C)^{-1} \Sigma_t^{-1} \\
&= (\Sigma_t^{-1} C^T R^{-1} C)^{-1} C^T R^{-1} C \\
&= \Sigma_{t+1} C^T R^{-1} \cancel{C}
\end{aligned}$$

Kalman smoothing Given an entire time series of T observations, the smoothing (marginal) distribution $p(z_t | x_{0:T})$ can be computed using an efficient algorithm known as Kalman smoothing (aka the Rauch-Tung-Striebel algorithm), which consists of a forward pass of Kalman filtering, as above, followed by a backwards pass. Indeed, we first note that the smoothing distribution at time step T is also the Gaussian $p(z_T | x_{0:T})$ obtained via filtering. From there, we can work backwards in much the same way as we worked forwards in the filtering pass. Suppose $p(z_{t+1} | x_{0:T})$ is known to be a Gaussian with mean $\tilde{\mu}_{t+1}$ and covariance $\tilde{\Sigma}_{t+1}$. To obtain $\tilde{\mu}_t$ and $\tilde{\Sigma}_t$ (and convince ourselves that $p(z_t | x_{0:T})$ is indeed Gaussian), we marginalise out z_{t+1} :

$$\begin{aligned}
p(z_t | x_{0:T}) &= \int dz_{t+1} \underbrace{p(z_{t+1} | x_{0:T})}_{\mathcal{N}(\tilde{\mu}_{t+1}, \tilde{\Sigma}_{t+1})} p(z_t | z_{t+1}, x_{0:T}) \\
&= \int p(z_t, z_{t+1} | x_{0:T}) \cdot dz_{t+1}
\end{aligned}$$

to find z_t . (13)

Moreover, we can use the conditional independence assumptions of the LDS model to split:

$$\begin{aligned}
p(z_t | z_{t+1}, x_{0:T}) &\propto p(z_t, z_{t+1}, x_{0:T}) \quad (\text{proportional as a function of } \cancel{z_t}, \cancel{x_{0:T}}?) \quad (14) \\
&= p(z_t | x_{0:t}) p(z_{t+1} | z_t) p(x_{t+1:T} | z_{t+1}) \quad (\text{conditional independence}) \quad (15)
\end{aligned}$$

$$\propto p(z_t | x_{0:t}) p(z_{t+1} | z_t) \quad (\text{just dropping the last term, constant in } z_t) \quad (16)$$

$$\stackrel{\text{K filter}}{=} \mathcal{N}(z_t; \mu_t, \Sigma_t) \mathcal{N}(z_{t+1}; A z_t, B B^T). \quad (17)$$

The product of the two normal densities in the last line is a (unnormalized) Gaussian density in z_t with mean⁴ $\mu_t - \Sigma_t A^T (A \Sigma_t A^T + B B^T)^{-1} A (\tilde{\mu}_{t+1} - \mu_t)$ and covariance $(\Sigma_t^{-1} + A^T (B B^T)^{-1} A)^{-1}$. Defining $P_t = A \Sigma_t A^T + B B^T$, we can now proceed with the marginalization in Equation (13), yielding

$$\begin{cases} \tilde{\mu}_t = \mu_t - \Sigma_t A^T P^{-1} A (\tilde{\mu}_{t+1} - \mu_t) \\ \tilde{\Sigma}_t = (\Sigma_t^{-1} + A^T (B B^T)^{-1} A)^{-1} + \Sigma_t A^T P^{-1} A \tilde{\Sigma}_{t+1} A^T P^{-1} A \Sigma_t \end{cases} \quad (18)$$

$$\quad (19)$$

This backward update step can again be massaged into a more interpretable form:

$$\boxed{
\begin{aligned}
G_t &\equiv \Sigma_t A^T (A \Sigma_t A^T + B B^T)^{-1} \quad \mu_t, \Sigma_t \text{ are from} \quad (20) \\
\tilde{\mu}_t &= \mu_t + G_t (\tilde{\mu}_{t+1} - A \mu_t) \quad (21) \\
\tilde{\Sigma}_t &= \Sigma_t + G_t (\tilde{\Sigma}_{t+1} - P_t) G_t^T \quad (22)
\end{aligned}
}$$

1.2 Parameter estimation

Here we will give a brief and incomplete summary of parameter estimation in LDS. Similar to pPCA and FA, the parameter estimation in LDS can be done by maximum likelihood, i.e. by maximizing $p(\{x_{1:T}^{(k)}\}_{k=1}^K | \theta)$, or equivalently the log-likelihood

$$\log p(\{x_{1:T}^{(k)}\}_{k=1}^K | \theta) = \sum_{k=1}^K \log p(x_{1:T}^{(k)} | \theta) \quad (23)$$

Here $\{x_{1:T}^{(k)}\}_{k=1}^K$ denotes a dataset of observed trajectories in K trials, and θ is the vector of parameters of the LDS model, namely a concatenation of $A, B, C, \text{diag}(R), \mu_0$, and Σ_0 . In writing the equation above, we have assumed that the observed neural activity in different trials are statistically independent, and therefore their log-probabilities sum over trials.

⁴The Woodbury identity is again used here.

The difficulty is that $p(\mathbf{x}_{1:T}^{(k)}|\boldsymbol{\theta})$ (which is actually a Gaussian distribution) has a complicated expression in terms of the parameters (more generally, *i.e.*, when we deviate from assumptions of linearity or Gaussianity, even calculating the density $p(\mathbf{x}_{1:T}^{(k)}|\boldsymbol{\theta})$, aka the marginal likelihood, is often intractable as it requires integration over the high-dimensional latent variable trajectories).

One classic approach for obtaining the maximum likelihood parameter estimates is the Expectation-Maximization (EM) algorithm. More modern gradient based (or second order) optimization methods are also popular and can be more efficient. Finally, one can go beyond parameter estimation and calculate an approximation to the Bayesian posterior distribution over the parameters, which also encodes their posterior uncertainty. A popular method for doing that is variational inference.

We will not cover these methods in this lecture. Here, we will only briefly sketch aspects of the EM algorithm for LDS; more details can be found in Section 13.3.2 of the [Pattern Recognition and Machine Learning](#) book by C.M. Bishop.

The expectation maximization (EM) algorithm. The EM algorithm is an iterative algorithm that alternates between the inference of the latent variables (E step) given the current estimate of the parameters, and updating of parameter estimates (M step) using the posterior distribution of the latent variables obtained in the previous E step.

In general, the E step consists of calculating the joint posterior distribution of the latent variables $p(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$. However, due to the HMM structure and Gaussianity of the posterior in our LDS, and also because of the structure of the sufficient statistics for the model parameters, all that we need from this posterior are the single-time-point posterior expectations $\mathbb{E}[\mathbf{z}_t]$, and posterior expectations of the pairwise products of latents over the same and consecutive time points, namely $\mathbb{E}[\mathbf{z}_t \mathbf{z}_t^T]$ and $\mathbb{E}[\mathbf{z}_{t+1} \mathbf{z}_t^T]$; here and below we denote posterior averages (over \mathbf{z}) by $\mathbb{E}[\cdot]$, keeping the conditioning on $\mathbf{x}_{1:T}$ implicit. The former are the smoothing posterior expectations $\tilde{\boldsymbol{\mu}}$ that, as described in the previous section, are calculated using the forward-backward algorithm. It turns out that a small addition to the forward-backward algorithm also yields the needed pairwise expectations (see Eq. 13.104 of the Bishop book).

In the M-step, we optimise the parameters given the above posterior expectations of the latents (which in turn were calculated based on the optimised parameters in the previous M-step). More precisely, the parameters are optimized to maximise the posterior expectation of the *joint* log-probability density, namely by solving the optimization problem⁵

$$\max_{\boldsymbol{\theta}} \sum_{k=1}^K \mathbb{E}[\log p(\mathbf{z}_{1:T}^{(k)}, \mathbf{x}_{1:T}^{(k)}|\boldsymbol{\theta})]. \quad (24)$$

Here, we will only review the M-step update of A (constituting a subset of model parameters). In particular, this will clarify why we (only in this case) need the pairwise posterior expectations. Based on the HMM structure (see [Figure 1](#)), let us first write out the components of the joint probabilities (for trial k) more explicitly as

$$\log p(\mathbf{z}_{1:T}^{(k)}, \mathbf{x}_{1:T}^{(k)}|\boldsymbol{\theta}) = \log p(\mathbf{z}_1^{(k)}|\boldsymbol{\theta}) + \sum_{t=2}^T \log p(\mathbf{z}_t^{(k)}|\mathbf{z}_{t-1}^{(k)}, \boldsymbol{\theta}) + \sum_{t=1}^T \log p(\mathbf{x}_t^{(k)}|\mathbf{z}_t^{(k)}, \boldsymbol{\theta}), \quad (25)$$

$$= \log p(\mathbf{z}_1^{(k)}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) + \sum_{t=2}^T \log p(\mathbf{z}_t^{(k)}|\mathbf{z}_{t-1}^{(k)}, A, B) + \sum_{t=1}^T \log p(\mathbf{x}_t^{(k)}|\mathbf{z}_t^{(k)}, C, R). \quad (26)$$

⁵Note how this is different from the direct maximization of model likelihood or log-likelihood, namely $\log p(\mathbf{x}_{1:T}|\boldsymbol{\theta})$; however, the fundamental theorem underlying and justifying the EM shows that the EM iterations converge at a local maximum of the model (“marginal”) likelihood $\log p(\mathbf{x}_{1:T}|\boldsymbol{\theta})$.

Since only the second sum above depends on A , we can focus on that. The summands of this sum are given by (as shown in previous section – see e.g., Equation (8))

$$\log \mathcal{N}(\mathbf{z}_t^{(k)}; A\mathbf{z}_{t-1}^{(k)}, BB^T) = -\frac{1}{2}(\mathbf{z}_t^{(k)} - A\mathbf{z}_{t-1}^{(k)})^T (BB^T)^{-1} (\mathbf{z}_t^{(k)} - A\mathbf{z}_{t-1}^{(k)}) + \text{const.} \quad (27)$$

where const. means independent of A . We see that this is quadratic in the \mathbf{z} 's and moreover only involves products of \mathbf{z} 's over the same or consecutive time steps. In particular, we see that the posterior expectation of this term, which we need to maximise in A , only involves the pairwise posterior expectations $\mathbb{E}[\mathbf{z}_t \mathbf{z}_t^T]$ and $\mathbb{E}[\mathbf{z}_t \mathbf{z}_{t-1}^T]$. Taking gradients with respect to A yields⁶

$$\sum_{k=1}^K \sum_{t=2}^T (BB^T)^{-1} \mathbb{E} \left[(\mathbf{z}_t^{(k)} - A\mathbf{z}_{t-1}^{(k)}) \mathbf{z}_{t-1}^{(k)T} \right] \quad (28)$$

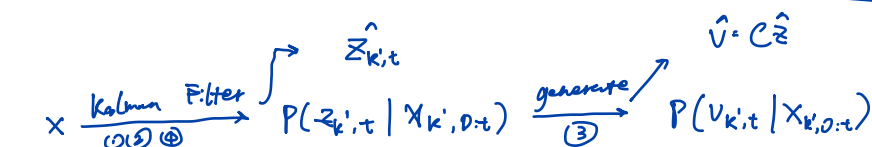
Setting this equal to zero (the necessary condition for local maximum), after multiplying by BB^T (and dividing by $K(T-1)$), yields the update rule

$$A^{\text{new}} = \overline{\mathbb{E}[\mathbf{z}_t \mathbf{z}_{t-1}^T]} \overline{\mathbb{E}[\mathbf{z}_{t-1} \mathbf{z}_{t-1}^T]}^{-1} \quad (29)$$

where $\overline{\mathbb{E}[\cdot]}$ denotes expectation over the posterior as well as averaging over trials and time steps $t = 2$ to T (and we have hidden the trial superscripts). Finally, note how this update generalizes the least squares solution derived in the lecture slides, in the case where the \mathbf{z} 's are fully observed, namely

$$A = \left(\overline{\mathbf{z}_t \mathbf{z}_{t-1}^T} \right) \left(\overline{\mathbf{z}_{t-1} \mathbf{z}_{t-1}^T} \right)^{-1} \quad (30)$$

where the bar denotes trial and time averaging; thus when the \mathbf{z} 's are latent, the formula is the same except trial/time-averaging has to be supplemented by a posterior average, because the \mathbf{z} 's were not directly observed.



maximize $\log P(\hat{\mathbf{z}}_{k,1:T}, \mathbf{X}_{k,1:T})$

eg4: $\mathbf{X}_{k,t} = \mathbf{D} \mathbf{z}_{k,t} + \mathbf{\xi}_{k,t} \sim N(0, S)$

$$\begin{aligned} \log P(\hat{\mathbf{z}}_{k,1:T}, \mathbf{X}_{k,1:T}) &= \log P(\mathbf{X}_{k,1:T} | \hat{\mathbf{z}}_{k,1:T}) P(\hat{\mathbf{z}}_{k,1:T}) \\ &= \frac{1}{K} \sum_k \left[\sum_{t=1}^T \log P(\mathbf{X}_{k,t} | \hat{\mathbf{z}}_{k,t}) + \log P(\hat{\mathbf{z}}_{k,1:T}) \right] \quad \text{(independent from D, S)} \\ &= \frac{1}{K} \sum_k \sum_{t=1}^T \left[-\frac{1}{2} \log 2\pi |S| - \frac{1}{2} (\mathbf{X}_{k,t} - \mathbf{D} \hat{\mathbf{z}}_{k,t})^T S^{-1} (\mathbf{X}_{k,t} - \mathbf{D} \hat{\mathbf{z}}_{k,t}) \right] \\ \frac{\partial}{\partial \mathbf{D}} &= \frac{1}{K} \sum_k \sum_{t=1}^T (\mathbf{X}_{k,t} \mathbf{z}_{k,t}^T - \mathbf{D} \mathbf{z}_{k,t} \mathbf{z}_{k,t}^T) = 0 \Rightarrow \mathbf{D} = \left(\frac{1}{K} \sum_k \sum_{t=1}^T \mathbf{X}_{k,t} \mathbf{z}_{k,t}^T \right) \left(\frac{1}{K} \sum_k \sum_{t=1}^T \mathbf{z}_{k,t} \mathbf{z}_{k,t}^T \right)^{-1} \end{aligned}$$

⁶To derive this, it is more straightforward to write Equation (27) explicitly in component form, as sums over matrix/vector components, and take a partial derivative w.r.t. A_{ij} and then at the end re-express the result in matrix/vector form.

$$\frac{\partial}{\partial S} = \sum \sum -\frac{1}{2}$$