

Linear Dynamical Systems and State Space Models

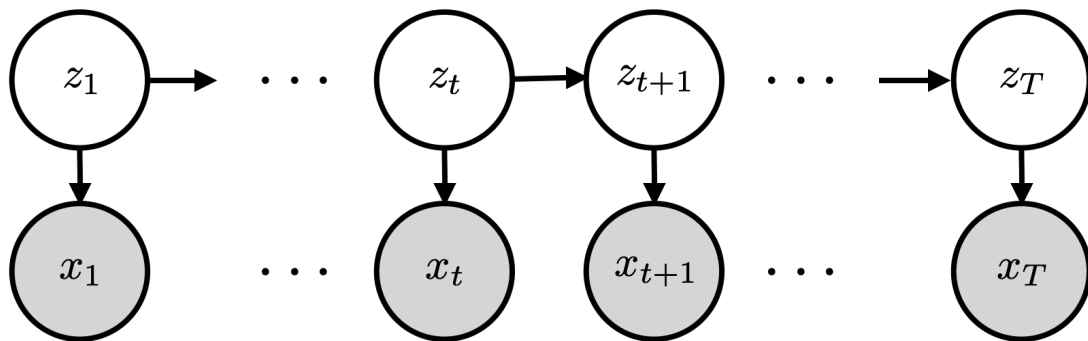
STATS 305C: Applied Statistics

Scott Linderman

May 11, 2022

Hidden Markov Models

Hidden Markov Models (HMMs) assume a particular factorization of the joint distribution on latent states (z_t) and observations (\mathbf{x}_t). The graphical model looks like this:



This graphical model says that the joint distribution factors as,

$$p(z_{1:T}, \mathbf{x}_{1:T}) = p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t | z_t). \quad (1)$$

We call this an HMM because $p(z_1) \prod_{t=2}^T p(z_t | z_{t-1})$ is a Markov chain.

Hidden Markov Models II

We are interested in questions like:

- ▶ What are the *predictive distributions* of $p(z_{t+1} \mid \mathbf{x}_{1:t})$?
- ▶ What is the *posterior marginal* distribution $p(z_t \mid \mathbf{x}_{1:T})$?
- ▶ What is the *posterior pairwise marginal* distribution $p(z_t, z_{t+1} \mid \mathbf{x}_{1:T})$?
- ▶ What is the *posterior mode* $\mathbf{z}_{1:T}^* = \arg \max p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T})$?
- ▶ How can we *sample the posterior* $p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T})$ of an HMM?
- ▶ What is the *marginal likelihood* $p(\mathbf{x}_{1:T})$?
- ▶ How can we *learn the parameters* of an HMM?

Question: Why might these sound like hard problems?

State space models

Note that nothing above assumes that z_t is a discrete random variable!

HMM's are a special case of more general **state space models** with discrete states.

State space models assume the same graphical model but allow for arbitrary types of latent states.

For example, suppose that $\mathbf{z}_t \in \mathbb{R}^D$ are continuous valued latent states and that,

$$p(\mathbf{z}_{1:T}) = p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t \mid \mathbf{z}_{t-1}) \quad (2)$$

$$= \mathcal{N}(\mathbf{z}_1 \mid \mathbf{b}_1, \mathbf{Q}_1) \prod_{t=2}^T \mathcal{N}(\mathbf{z}_t \mid \mathbf{A}\mathbf{z}_{t-1} + \mathbf{b}, \mathbf{Q}) \quad (3)$$

This is called a Gaussian **linear dynamical system** (LDS).

Stability of Gaussian linear dynamical systems

Question: What is the asymptotic mean of a Gaussian LDS, $\lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{z}_t]$?

Question: When is a Gaussian LDS stable? I.e. when is the asymptotic mean finite?

$$\mu_{\infty} = \lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{z}_t]$$

$$\mu_{\infty} = A\mu_{\infty} + b$$

$$\mu_{\infty} = (I - A)^{-1} b$$

$$|\text{eig}(A)| < 1$$

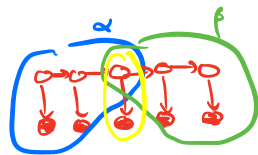
EE263: LDS

Message passing in HMMs

In the HMM with discrete states, we showed how to compute posterior marginal distributions using message passing,

$$p(z_t | \mathbf{x}_{1:T}) \propto \sum_{z_1} \cdots \sum_{z_{t-1}} \sum_{z_{t+1}} \cdots \sum_{z_T} p(z_{1:T}, \mathbf{x}_{1:T}) \quad (4)$$

$$= \alpha_t(z_t) p(\mathbf{x}_t | z_t) \beta_t(z_t) \quad (5)$$



where the *forward* and *backward* messages are defined recursively

$$\alpha_t(z_t) = \sum_{z_{t-1}} p(z_t | z_{t-1}) p(\mathbf{x}_{t-1} | z_{t-1}) \alpha_{t-1}(z_{t-1}) \quad (6)$$

$$\beta_t(z_t) = \sum_{z_{t+1}} p(z_{t+1} | z_t) p(\mathbf{x}_{t+1} | z_{t+1}) \beta_{t+1}(z_{t+1}) \quad (7)$$

The initial conditions are $\alpha_1(z_1) = p(z_1)$ and $\beta_T(z_T) = 1$.

What do the forward messages tell us?

The forward messages are equivalent to,

$$\alpha_t(z_t) = \sum_{z_1} \cdots \sum_{z_{t-1}} p(z_{1:t}, \mathbf{x}_{1:t-1}) \quad (8)$$

$$p(z_t, \mathbf{x}_{1:t-1}). \quad (9)$$

The normalized message is the *predictive distribution*,

$$\frac{\alpha_t(z_t)}{\sum_{z'_t} \alpha_t(z'_t)} = \frac{p(z_t, \mathbf{x}_{1:t-1})}{\sum_{z'_t} p(z'_t, \mathbf{x}_{1:t-1})} = \frac{p(z_t, \mathbf{x}_{1:t-1})}{p(\mathbf{x}_{1:t-1})} = p(z_t \mid \mathbf{x}_{1:t-1}), \quad (10)$$

The final normalizing constant yields the marginal likelihood, $\sum_{z_T} \alpha_T(z_T) = p(\mathbf{x}_{1:T})$.

Message passing in state space models

The same recursive algorithm applies (in theory) to any state space model with the same factorization, but the sums are replaced with integrals,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:T}) \propto \int d\mathbf{z}_1 \cdots \int d\mathbf{z}_{t-1} \int d\mathbf{z}_{t+1} \cdots \int d\mathbf{z}_T p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) \quad (11)$$

$$= \alpha_t(\mathbf{z}_t) p(\mathbf{x}_t \mid \mathbf{z}_t) \beta_t(\mathbf{z}_t) \quad (12)$$

where the *forward and backward messages* are defined recursively

$$\alpha_t(\mathbf{z}_t) = \int p(\mathbf{z}_t \mid \mathbf{z}_{t-1}) p(\mathbf{x}_{t-1} \mid \mathbf{z}_{t-1}) \alpha_{t-1}(\mathbf{z}_{t-1}) d\mathbf{z}_{t-1} \quad (13)$$

$$\beta_t(\mathbf{z}_t) = \int p(\mathbf{z}_{t+1} \mid \mathbf{z}_t) p(\mathbf{x}_{t+1} \mid \mathbf{z}_{t+1}) \beta_{t+1}(\mathbf{z}_{t+1}) d\mathbf{z}_{t+1} \quad (14)$$

The initial conditions are $\alpha_1(\mathbf{z}_1) = p(\mathbf{z}_1)$ and $\beta_T(\mathbf{z}_T) \propto 1$.

Forward pass in a linear dynamical system

Consider an linear dynamical system (LDS) with Gaussian emissions,

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{z}_t) \quad (15)$$

$$= \mathcal{N}(\mathbf{z}_1 | \mathbf{b}_1, \mathbf{Q}_1) \prod_{t=2}^T \mathcal{N}(\mathbf{z}_t | \mathbf{A}\mathbf{z}_{t-1} + \mathbf{b}, \mathbf{Q}) \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t | \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \quad (16)$$

Let's derive the forward message $\alpha_{t+1}(\mathbf{z}_{t+1})$. Assume $\alpha_t(\mathbf{z}_t) \propto \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$. [Induction]

$$\alpha_{t+1}(\mathbf{z}_{t+1}) = \int p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_t | \mathbf{z}_t) \alpha_t(\mathbf{z}_t) d\mathbf{z}_t \quad (17)$$

$$= \int \underbrace{\mathcal{N}(\mathbf{z}_{t+1} | \mathbf{A}\mathbf{z}_t + \mathbf{b}, \mathbf{Q})}_{\text{predict}} \underbrace{\mathcal{N}(\mathbf{x}_t | \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})}_{\text{update}} d\mathbf{z}_t \quad (18)$$

The update step

The first step is the **update step**, where we **condition on** the emission \mathbf{x}_t ,

Exercise: Expand the densities, collect terms, and complete the square to compute the mean $\boldsymbol{\mu}_{t|t}$ and covariance $\boldsymbol{\Sigma}_{t|t}$ after the update step,

$$\mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \propto \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}). \quad (19)$$

$$\propto \exp \left\{ -\frac{1}{2} (\mathbf{x}_t - \mathbf{C}\mathbf{z}_t - \mathbf{d})^\top \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{C}\mathbf{z}_t - \mathbf{d}) - \frac{1}{2} (\mathbf{z}_t - \boldsymbol{\mu}_{t|t-1})^\top \boldsymbol{\Sigma}_{t|t-1}^{-1} (\mathbf{z}_t - \boldsymbol{\mu}_{t|t-1}) \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{z}_t^\top \mathbf{J}_{t|t} \mathbf{z}_t + \mathbf{h}_{t|t}^\top \mathbf{z}_t \right\}$$

$$\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1} + \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top (\mathbf{C} \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top + \mathbf{R})^{-1} \mathbf{C} \boldsymbol{\Sigma}_{t|t-1} (\mathbf{x}_t - \mathbf{C} \boldsymbol{\mu}_{t|t-1} - \mathbf{d})$$

$$\mathbf{J}_{t|t} = \mathbf{C}^\top \mathbf{R}^{-1} \mathbf{C} + \boldsymbol{\Sigma}_{t|t-1}^{-1}$$

$$\mathbf{h}_{t|t} = \mathbf{J}_{t|t}^{-1} \mathbf{h}_{t|t}$$

$$= (\mathbf{C}^\top \mathbf{R}^{-1} \mathbf{C} + \boldsymbol{\Sigma}_{t|t-1}^{-1})^{-1} [\mathbf{C}^\top \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{d}) + \boldsymbol{\Sigma}_{t|t-1}^{-1} \boldsymbol{\mu}_{t|t-1}]$$

$$\mathbf{h}_{t|t} = \mathbf{C}^\top \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{d}) + \boldsymbol{\Sigma}_{t|t-1}^{-1} \boldsymbol{\mu}_{t|t-1}$$

The update step II

Write the joint distribution,

$$p(\mathbf{z}_t, \mathbf{x}_t \mid \mathbf{x}_{1:t-1}) = \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (20)$$

$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_{t|t-1} \\ \mathbf{C}\boldsymbol{\mu}_{t|t-1} + \mathbf{d} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{t|t-1} & \boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top \\ \mathbf{C}\boldsymbol{\Sigma}_{t|t-1} & \mathbf{C}\boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top + \mathbf{R} \end{bmatrix}\right) \quad (21)$$

Since $(\mathbf{z}_t, \mathbf{x}_t)$ are jointly Gaussian, \mathbf{z}_t must be conditionally Gaussian as well,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) = \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}). \quad (22)$$

Exercise: Now use the **Schur complement** from Week 1 to solve for $\boldsymbol{\mu}_{t|t}$ and $\boldsymbol{\Sigma}_{t|t}$

matrix inversion lemma (2x)

$$\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1} + \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top (\mathbf{C} \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top + \mathbf{R})^{-1} (\mathbf{x}_t - \mathbf{C} \boldsymbol{\mu}_{t|t-1} - \mathbf{d})$$

$$\boldsymbol{\Sigma}_{t|t} = \boldsymbol{\Sigma}_{t|t-1} - \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top (\mathbf{C} \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top + \mathbf{R})^{-1} \mathbf{C} \boldsymbol{\Sigma}_{t|t-1}$$

$$\boldsymbol{\mu}_{t|t} = (\mathbf{C}^\top \mathbf{R}^{-1} \mathbf{C} + \boldsymbol{\Sigma}_{t|t-1}^{-1})^{-1} [\mathbf{C}^\top \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{d}) + \boldsymbol{\Sigma}_{t|t-1}^{-1} \boldsymbol{\mu}_{t|t-1}]$$

The update step III

Exercise: Write $\mu_{t|t}$ and $\Sigma_{t|t}$ in terms of the **Kalman gain**,

$$K_t = \Sigma_{t|t-1} \mathbf{C}^\top (\mathbf{C} \Sigma_{t|t-1} \mathbf{C}^\top + \mathbf{R})^{-1} \quad \hat{x}_t = \mathbf{C} \mu_{t|t-1} + d \quad (23)$$

What is the Kalman gain doing?

$$\mu_{t|t} = \mu_{t|t-1} + \Sigma_{t|t-1} \mathbf{C}^\top (\mathbf{C} \Sigma_{t|t-1} \mathbf{C}^\top + \mathbf{R})^{-1} (x_t - \mathbf{C} \mu_{t|t-1} - d)$$

$$= \underbrace{\mu_{t|t-1}}_{\text{predicted mean}} + K_t \underbrace{(x_t - \hat{x}_t)}_{\text{residual}}$$

The predict step

The predict step yields $p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) = \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$. To complete the forward pass, we need the **predict step**,

$$\alpha_{t+1}(\mathbf{z}_{t+1}) = \int p(\mathbf{z}_{t+1} \mid \mathbf{z}_t) p(\mathbf{x}_t \mid \mathbf{z}_t) \alpha_t(\mathbf{z}_t) d\mathbf{z}_t \quad (24)$$

$$= \int \mathcal{N}(\mathbf{z}_{t+1} \mid \mathbf{A}\mathbf{z}_t + \mathbf{b}, \mathbf{Q}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) d\mathbf{z}_t \quad (25)$$

$$= \mathcal{N}(\mathbf{z}_{t+1} \mid \boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t}) \quad (26)$$

Exercise: Solve for the mean $\boldsymbol{\mu}_{t+1|t}$ and covariance $\boldsymbol{\Sigma}_{t+1|t}$ after the predict step.

$$\boldsymbol{\mu}_{t+1|t} = \mathbf{A}\boldsymbol{\mu}_{t|t} + \mathbf{b}$$

$$\boldsymbol{\Sigma}_{t+1|t} = \mathbf{Q} + \mathbf{A}\boldsymbol{\Sigma}_{t|t}\mathbf{A}^T$$

Completing the recursions

That wraps up the recursions! All that's left is the base case, which comes from the initial state distribution,

$$\mu_{1|0} = \mathbf{b}_1 \quad \text{and} \quad \Sigma_{1|0} = \mathbf{Q}_1. \quad (27)$$

$$\phi(x_1) = \mathcal{N}(x_1 | \mathbf{b}_1, \mathbf{Q}_1)$$

Computing the marginal likelihood

Like in the discrete HMM, we can compute the marginal likelihood along the forward pass.

$$p(\mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{x}_{1:t-1}) \quad (28)$$

$$= \prod_{t=1}^T \int p(\mathbf{x}_t \mid \mathbf{z}_t) p(\mathbf{z}_t \mid \mathbf{x}_{1:t-1}) d\mathbf{z}_t \quad (29)$$

$$= \prod_{t=1}^T \int \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) d\mathbf{z}_t \quad (30)$$

Exercise: Obtain a closed form expression for the integrals.

$$\rightarrow \prod_t \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\boldsymbol{\mu}_{t|t-1} + \mathbf{d}, \mathbf{R} + \mathbf{C}\boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^T)$$

Computing the smoothing distributions

- ▶ The forward pass gives us the filtering distributions $p(\mathbf{z}_t \mid \mathbf{x}_{1:t})$. How can we compute the smoothing distributions, $p(\mathbf{z}_t \mid \mathbf{x}_{1:T})$?
- ▶ In the discrete HMM we essentially ran the *same algorithm in reverse* to get the backward messages, starting from $\beta_T(\mathbf{z}_T) \propto 1$.
- ▶ We can do the same sort of thing here, but it's a bit funky because we need to start with an improper Gaussian distribution $\beta_T(\mathbf{z}_T) \propto \mathcal{N}(\mathbf{0}, \infty I)$.
- ▶ Instead, we'll derive an alternative way of computing the smoothing distributions.

Bayesian Smoothing

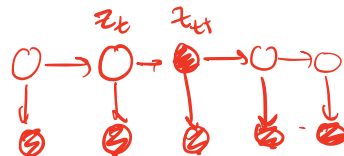
Note: \mathbf{z}_t is conditionally independent of $\mathbf{x}_{t+1:T}$ given \mathbf{z}_{t+1} , so

$$p(\mathbf{z}_t | \mathbf{x}_{1:T})$$

$$\neq p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:T}) = p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:t})$$

$$= \frac{p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{x}_{1:t})}{p(\mathbf{z}_{t+1} | \mathbf{x}_{1:t})}$$

$$= \frac{p(\mathbf{z}_t | \mathbf{x}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t)}{p(\mathbf{z}_{t+1} | \mathbf{x}_{1:t})}$$



(31)

(32)

(33)

Question: what rules did we apply in each of these steps?

Bayesian Smoothing II

Now we can write the joint distribution as,

$$p(\mathbf{z}_t, \mathbf{z}_{t+1} \mid \mathbf{x}_{1:T}) = p(\mathbf{z}_t \mid \mathbf{z}_{t+1}, \mathbf{x}_{1:T}) p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T}) \quad (34)$$

$$= \frac{p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) p(\mathbf{z}_{t+1} \mid \mathbf{z}_t) p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T})}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})}. \quad (35)$$

Marginalizing over \mathbf{z}_{t+1} gives us,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:T}) = p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) \int \frac{\overset{\text{model}}{p(\mathbf{z}_{t+1} \mid \mathbf{z}_t)} \overset{\text{Smoothing}}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T})}}{\underset{\text{filtering}}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})}} d\mathbf{z}_{t+1} \quad (36)$$

Question: Can we compute each of these terms?

The Rauch-Tung-Striebel Smoother, aka Kalman Smoother

These recursions apply to any Markovian state space model. For the special case of a Gaussian linear dynamical system, the smoothing distributions are again Gaussians,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}_t \mid \underline{\mu}_{t|T}, \underline{\Sigma}_{t|T}) \quad (37)$$

where

$$\underline{\mu}_{t|T} = \mu_{t|t} + \mathbf{G}_t(\mu_{t+1|T} - \mu_{t+1|t}) \quad (38)$$

$$\underline{\Sigma}_{t|T} = \Sigma_{t|t} + \mathbf{G}_t(\Sigma_{t+1|T} - \Sigma_{t+1|t})\mathbf{G}_t^\top \quad (39)$$

$$\mathbf{G}_t \triangleq \Sigma_{t|t}\mathbf{A}^\top \Sigma_{t+1|t}^{-1}. \quad (40)$$

This is called the **Rauch-Tung-Striebel (RTS) smoother** or the **Kalman smoother**.

$$\mu_{t|t} \longrightarrow \mu_{T|T}$$

Kalman smoothing in information form

So far we've worked with the *mean parameters* $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, but working with *natural parameters* \mathbf{J} and \mathbf{h} offers another perspective.

Let's go back to the basics,

$$p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T}) \propto p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) \quad (41)$$

$$= p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t \mid \mathbf{z}_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{z}_t) \quad (42)$$

$$= \mathcal{N}(\mathbf{z}_1 \mid \mathbf{b}_1, \mathbf{Q}_1) \prod_{t=2}^T \mathcal{N}(\mathbf{z}_t \mid \mathbf{A}\mathbf{z}_{t-1} + \mathbf{b}, \mathbf{Q}) \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \quad (43)$$

Kalman smoothing in information form II

Expand the Gaussian densities,

$$p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T}) \propto \exp \left\{ -\frac{1}{2} (\mathbf{z}_1 - \mathbf{b}_1)^\top \mathbf{Q}_1^{-1} (\mathbf{z}_1 - \mathbf{b}_1) \right. \quad (44)$$

$$\left. -\frac{1}{2} \sum_{t=2}^T (\mathbf{z}_t - \mathbf{A}\mathbf{z}_{t-1} - \mathbf{b})^\top \mathbf{Q}^{-1} (\mathbf{z}_t - \mathbf{A}\mathbf{z}_{t-1} - \mathbf{b}) \right. \quad (45)$$

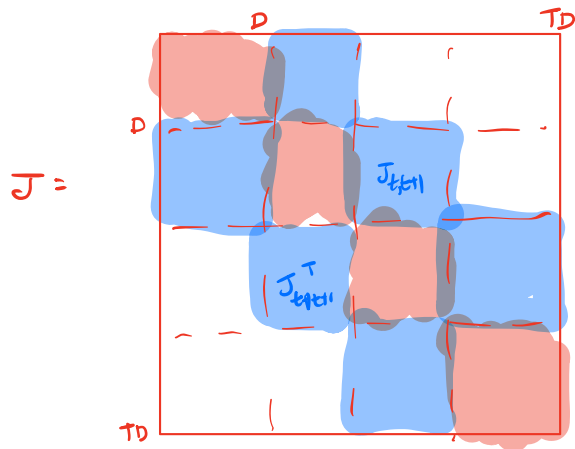
$$\left. -\frac{1}{2} \sum_{t=1}^T (\mathbf{x}_t - \mathbf{C}\mathbf{z}_t - \mathbf{d})^\top \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{C}\mathbf{z}_t - \mathbf{d}) \right\} \quad (46)$$

This is a giant quadratic expression in $\mathbf{z}_{1:T}$; i.e. a multivariate normal distribution on \mathbb{R}^{TD} .

We can write it in terms of its natural parameters $\mathbf{J} \in \mathbb{R}^{TD \times TD}$ and $\mathbf{h} \in \mathbb{R}^{TD}$

Kalman smoothing in information form III

Question: Which entries in J are nonzero?



block tridiagonal
 $T(D^2) + (T-1)D^2$

Duality between message passing and sparse linear algebra

Recall that to get mean from the natural parameters we have,

$$p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}_{1:T} \mid \mathbf{J}^{-1}\mathbf{h}, \mathbf{J}^{-1}). \quad (47)$$

In other words, the posterior mean is the solution of a linear system $\mathbf{J}^{-1}\mathbf{h}$.

Typically, this would cost $O((TD)^3)$, but since \mathbf{J} is block-tridiagonal (or more generally, banded), we can compute it in only $O(TD^3)$ time.

The algorithm for solving this sparse linear system is essentially the same as the message passing algorithm we derived today.

Message passing in nonlinear dynamical systems

Question: What would you do if you were given a nonlinear model,

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t | f(\mathbf{z}_{t-1}), \mathbf{Q})?$$

$$q(\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_{t-1} | \mu_{t-1|t-2}, \Sigma_{t-1|t-2})$$

$$f(\mathbf{z}_{t-1}) \approx f(\mu_{t-1|t-2}) + \nabla f(\mu_{t-1|t-2})(\mathbf{z}_{t-1} - \mu_{t-1|t-2}) + \mathcal{O}(\|\mathbf{z}_{t-1} - \mu_{t-1|t-2}\|^2)$$

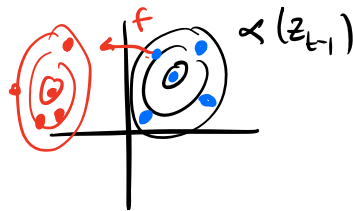
\Rightarrow extended kalman filter

c.f. unscented kalman filter

- neural network

- \rightarrow Deep Kalman filter

- \rightarrow Structured VAEs



Sequential Monte Carlo

Recall that the forward messages are proportional to the predictive distributions $p(\mathbf{z}_t \mid \mathbf{x}_{1:t-1})$. We can view the forward recursions as an expectation,

$$\alpha_t(\mathbf{z}_t) = \int p(\mathbf{z}_t \mid \mathbf{z}_{t-1}) p(\mathbf{x}_{t-1} \mid \mathbf{z}_{t-1}) \alpha_{t-1}(\mathbf{z}_{t-1}) d\mathbf{z}_{t-1} \quad (48)$$

$$\propto \mathbb{E}_{\mathbf{z}_{t-1} \sim p(\mathbf{z}_{t-1} \mid \mathbf{x}_{1:t-2})} [p(\mathbf{z}_t \mid \mathbf{z}_{t-1}) p(\mathbf{x}_{t-1} \mid \mathbf{z}_{t-1})] \quad (49)$$

One natural idea is to approximate this expectation with Monte Carlo,

$$\hat{\alpha}_t(\mathbf{z}_t) \approx \frac{1}{S} \sum_{s=1}^S \left[w_{t-1}^{(s)} p(\mathbf{z}_t \mid \mathbf{z}_{t-1}^{(s)}) \right] \quad (50)$$

where we have defined the **weights** $w_{t-1}^{(s)} \triangleq p(\mathbf{x}_{t-1} \mid \mathbf{z}_{t-1}^{(s)})$.

How do we sample $\mathbf{z}_{t-1}^{(s)} \stackrel{\text{iid}}{\sim} p(\mathbf{z}_{t-1} \mid \mathbf{x}_{1:t-2})$? Let's sample the normalized $\hat{\alpha}_{t-1}(\mathbf{z}_{t-1})$ instead!

Sequential Monte Carlo II

The normalizing constant is,

$$\int \hat{\alpha}_{t-1}(\mathbf{z}_{t-1}) d\mathbf{z}_{t-1} = \frac{1}{S} \sum_{s=1}^S w_{t-2}^{(s)} \int p(\mathbf{z}_{t-1} | \mathbf{z}_{t-2}^{(s)}) d\mathbf{z}_{t-1} = \frac{1}{S} \sum_{s=1}^S w_{t-2}^{(s)}. \quad (51)$$

Use this to define the *normalized forward message* (i.e. the Monte Carlo estimate of the predictive distribution) is,

$$\bar{\alpha}_{t-1}(\mathbf{z}_{t-1}) \triangleq \frac{\hat{\alpha}_{t-1}(\mathbf{z}_{t-1})}{\int \hat{\alpha}_{t-1}(\mathbf{z}'_{t-1}) d\mathbf{z}'_{t-1}} = \sum_{s=1}^S \bar{w}_{t-2}^{(s)} p(\mathbf{z}_{t-1} | \mathbf{z}_{t-2}^{(s)}) \quad (52)$$

where $\bar{w}_{t-2}^{(s)} = \frac{w_{t-2}^{(s)}}{\sum_{s'} w_{t-2}^{(s')}}$ is the normalized weight of sample $\mathbf{z}_{t-2}^{(s)}$.

The normalized forward message is just a mixture distribution with weights $\bar{w}_{t-2}^{(s)}$!

Putting it all together

Combining the above, we have the following algorithm for the forward pass:

1. Let $\bar{\alpha}_1(\mathbf{z}_1) = p(\mathbf{z}_1)$
2. For $t = 1, \dots, T$:
 - a. Sample $\mathbf{z}_t^{(s)} \stackrel{\text{iid}}{\sim} \bar{\alpha}_t(\mathbf{z}_t)$ for $s = 1, \dots, S$
 - b. Compute weights $w_t^{(s)} = p(\mathbf{x}_t | \mathbf{z}_t^{(s)})$ and normalize $\bar{w}_t^{(s)} = w_t^{(s)} / \sum_{s'} w_t^{(s')}$.
 - c. Compute normalized forward message $\bar{\alpha}_{t+1}(\mathbf{z}_{t+1}) = \sum_{s=1}^S \bar{w}_t^{(s)} p(\mathbf{z}_{t+1} | \mathbf{z}_t^{(s)})$.

This is called **sequential Monte Carlo** (SMC) using the model dynamics as the proposal.

Note that Step 2a can **resample** the same $\mathbf{z}_{t-1}^{(s)}$ multiple times according to its weight.

Question: How can you approximate the marginal likelihood $p(\mathbf{x}_{1:T})$ using the weights? *Hint: look back to Slide 7.*

Generalizations

- Instead of sampling $\bar{\alpha}_t(\mathbf{z}_t)$, we could have sampled with a **proposal distribution** $r(\mathbf{z}_t | \mathbf{z}_{t-1}^{(s)})$ instead and corrected for it by defining the weights to be,

$$w_t^{(s)} = \frac{p(\mathbf{z}_t | \mathbf{z}_{t-1}^{(s)}) p(\mathbf{x}_t | \mathbf{z}_t)}{r(\mathbf{z}_t | \mathbf{z}_{t-1}^{(s)})} \quad (53)$$

Moreover, the proposal distribution can “look ahead” to future data \mathbf{x}_t .

References I