

Latent Variable Models for Sequential Data

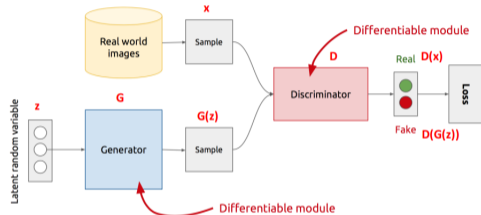
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

April 8, 2019

Recap: Deep Generative Models - GAN

- GAN: Generative Adversarial Network[†]
- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)



- Can be thought of as a two-player minimax game

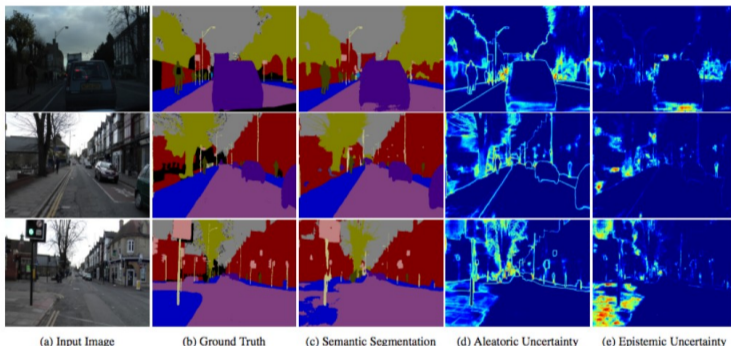
$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

- With the generator G fixed, the optimal discriminator $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$
- At the global minimum of the objective, $p_g = p_{data}$

[†] Generative Adversarial Nets (Goodfellow et al, 2013), Figure: <https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

Deep Learning and Uncertainty

- Estimating uncertainty is important

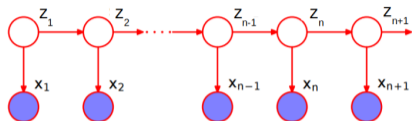


What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? [Kendall & Gal, NIPS, 2017]

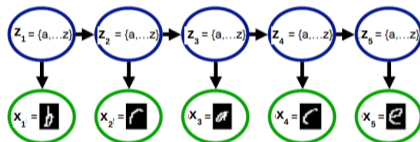
- Aleatoric uncertainty**, capturing inherent noise in the data; **Epistemic uncertainty**, capturing models lack of knowledge

Latent Variable Models for Sequential Data

- Task: Given a sequence of observations, infer the latent state of each observation



- An example: Recognizing a sequence of handwritten characters

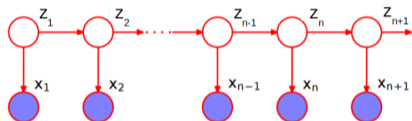


- In this example, the latent state z_n at step n is a **discrete value**
- Another example: Given a sequence of observed noisy 2D coordinates x_n of an object, infer its latent state z_n , e.g., actual coordinates, velocity, acceleration, etc. at each step $n = 1, 2, \dots$
 - In this example, the latent state z_n at step n is a **continuous vector**

Latent Variable Models for Sequential Data

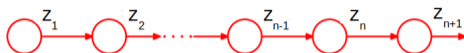
- Consider the following latent variable model for a sequence of observations $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$

$$\begin{aligned}\mathbf{x}_n | \mathbf{z}_n &\sim p(\mathbf{x}_n | \mathbf{z}_n) && \text{(i.i.d. draws of } \mathbf{x}_n \text{ given } \mathbf{z}_n) \\ \mathbf{z}_n | \mathbf{z}_{n-1} &\sim p(\mathbf{z}_n | \mathbf{z}_{n-1}) && \text{(first-order dependence b/w } \mathbf{z}_n \text{'s)}\end{aligned}$$



- $p(\mathbf{z}_n | \mathbf{z}_{n-1})$ is called state-transition model, $p(\mathbf{x}_n | \mathbf{z}_n)$ is called observation/emission model
 - Note: In some cases, the parameters defining these distributions may be known
- If latent states \mathbf{z}_n are discrete, we get a **Hidden Markov Model (HMM)**
- If latent states \mathbf{z}_n are continuous vectors, we get a **State-Space Model (SSM)**
- In both cases, observations \mathbf{x}_n can be anything (discrete/real)

State-Transition Model



- For discrete states case (HMM), $p(\mathbf{z}_n|\mathbf{z}_{n-1})$ will be a **discrete distribution**, e.g.,

$$p(\mathbf{z}_n|\mathbf{z}_{n-1} = \ell) = \text{multinoulli}(\boldsymbol{\pi}_\ell)$$

where $\boldsymbol{\pi}_\ell = [\pi_{\ell,1}, \dots, \pi_{\ell,K}]$ is $K \times 1$ a **transition prob. vector**, s.t. $p(\mathbf{z}_n = k|\mathbf{z}_{n-1} = \ell) = \pi_{\ell,k}$

- For HMM, $p(\mathbf{z}_n|\mathbf{z}_{n-1})$ is fully defined by a $K \times K$ **transition prob. matrix** $\Pi = [\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots, \boldsymbol{\pi}_K]$
- For continuous states (SSM), $p(\mathbf{z}_n|\mathbf{z}_{n-1})$ will be a **continuous distribution**, e.g., Gaussian

$$p(\mathbf{z}_n|\mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{A}\mathbf{z}_{n-1}, \mathbf{I}_K)$$

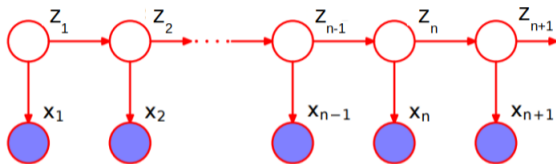
- Note: More powerful transition models usually employ nonlinear mappings between \mathbf{z}_{n-1} and \mathbf{z}_n
- For both HMM and SSM, there is also an **initial state distribution** $p(\mathbf{z}_1)$, e.g.,

$$p(\mathbf{z}_1) = \text{multinoulli}(\boldsymbol{\pi}_0) \quad (\text{for HMM})$$

$$p(\mathbf{z}_1) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K) \quad (\text{for SSM})$$

Observation/Emission Model

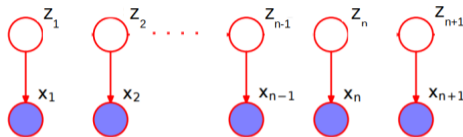
- The type of observation model distribution $p(\mathbf{x}_n|\mathbf{z}_n)$ depends on the type of data



- For discrete observations (e.g., words), $p(\mathbf{x}_n|\mathbf{z}_n)$ is a discrete distribution (e.g., multinoulli)
- For continuous observations (e.g., images, location of an object, etc.), $p(\mathbf{x}_n|\mathbf{z}_n)$ is a continuous distribution (e.g., Gaussian)
- Note: More powerful observation models usually employ nonlinear mappings between \mathbf{z}_n and \mathbf{x}_n

A Special Case

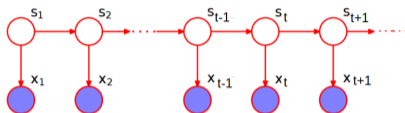
- What if we have i.i.d. latent states, i.e., $p(\mathbf{z}_n | \mathbf{z}_{n-1}) = p(\mathbf{z}_n)$?



- HMM becomes a **standard Mixture Model**. Reason: $p(\mathbf{z}_n | \mathbf{z}_{n-1} = \ell) = p(\mathbf{z}_n) = \text{multinoulli}(\boldsymbol{\pi})$
- SSM becomes **PPCA/factor analysis**. Reason: $p(\mathbf{z}_n | \mathbf{z}_{n-1}) = p(\mathbf{z}_n) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ or $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi})$
- Therefore, inference algorithms for HMM/SSM are often very similar to mixture models/PPCA
 - Only main difference is how the latent variables \mathbf{z}_n 's are inferred (because these are no longer i.i.d.)
 - E.g., if using EM, only E step needs to change. Given the expectations, the M step updates are derived similarly to how it's done in mixture models and PPCA (Bishop Chap 13 has EM for HMM and SSM)

State Space Models (SSM)

- Today we will mainly **focus on SSM** (when the latent variables are continuous vectors)



Using 's' instead of 'z'
to refer to states

Using 't' to denote the
'time-step'

- Most of the details of methods we will see apply to HMMs too (but \mathbf{s}_t will be discrete)
- In the most general form, the transition and observation models in an SSM can be expressed as

$$\mathbf{s}_t | \mathbf{s}_{t-1} = g_t(\mathbf{s}_{t-1}) + \epsilon_t \quad (\text{must be a cont. dist. over } \mathbf{s}_t)$$

$$\mathbf{x}_t | \mathbf{s}_t = h_t(\mathbf{s}_t) + \delta_t \quad (\text{can be any dist. over } \mathbf{x}_t)$$

- Here g_t and h_t are functions (can be linear/nonlinear)
- Assuming zero-mean Gaussian noise $\epsilon_t \sim \mathcal{N}(0, \mathbf{Q}_t)$, $\delta_t \sim \mathcal{N}(0, \mathbf{R}_t)$, we get a **Gaussian SSM**

$$\mathbf{s}_t | \mathbf{s}_{t-1} \sim \mathcal{N}(\mathbf{s}_t | g_t(\mathbf{s}_{t-1}), \mathbf{Q}_t)$$

$$\mathbf{x}_t | \mathbf{s}_t \sim \mathcal{N}(\mathbf{x}_t | h_t(\mathbf{s}_t), \mathbf{R}_t)$$

- Note: If $g_t, h_t, \mathbf{Q}_t, \mathbf{R}_t$ are independent of t then the model is called **stationary**

State Space Models (SSM)

- A simple example of a state-space model

$$\begin{aligned} \mathbf{s}_t | \mathbf{s}_{t-1} &= \mathbf{s}_{t-1} + \epsilon_t \\ \mathbf{x}_t | \mathbf{s}_t &= \mathbf{s}_t + \delta_t \end{aligned} \quad (\text{assumes } \mathbf{x}_t \text{ and } \mathbf{s}_t \text{ to be of same size})$$

- Another simple but more general example (latent states and observations of diff. dimensions)

$$\begin{aligned} \mathbf{s}_t | \mathbf{s}_{t-1} &= \mathbf{A}_t \mathbf{s}_{t-1} + \epsilon_t && (\mathbf{A}_t \text{ is } K \times K) \\ \mathbf{x}_t | \mathbf{s}_t &= \mathbf{B}_t \mathbf{s}_t + \delta_t && (\mathbf{B}_t \text{ is } D \times K) \end{aligned}$$

- The above can also be written as follows

$$\begin{aligned} \mathbf{s}_t | \mathbf{s}_{t-1} &\sim \mathcal{N}(\mathbf{s}_t | \mathbf{A}_t \mathbf{s}_{t-1}, \mathbf{Q}_t) \\ \mathbf{x}_t | \mathbf{s}_t &\sim \mathcal{N}(\mathbf{x}_t | \mathbf{B}_t \mathbf{s}_t, \mathbf{R}_t) \end{aligned}$$

- This is a [Linear Gaussian SSM](#); also called [Linear Dynamical System \(LDS\)](#)
- Note: $\mathbf{A}_t, \mathbf{B}_t, \mathbf{Q}_t, \mathbf{R}_t$ may be known (fixed) or may be required to be learned

Linear Gaussian SSM (LDS): An Example

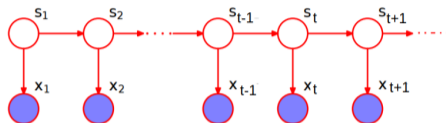
- Consider the linear Gaussian SSM: $\mathbf{s}_t | \mathbf{s}_{t-1} = \mathbf{A}_t \mathbf{s}_{t-1} + \epsilon_t$ and $\mathbf{x}_t | \mathbf{s}_t = \mathbf{B}_t \mathbf{s}_t + \delta_t$
- Suppose $\mathbf{x}_t \in \mathbb{R}^2$ denotes the (noisy) observed 2D location of an object
- Suppose $\mathbf{s}_t \in \mathbb{R}^6$ denotes its “state” vector $\mathbf{s}_t = [pos_1, vel_1, accel_1, pos_2, vel_2, accel_2]$
- Assuming a pre-defined \mathbf{A}_t , \mathbf{B}_t , a possible linear Gaussian SSM to model this data will be

$$\mathbf{s}_t = \mathbf{A}_t \mathbf{s}_{t-1} + \epsilon_t$$
$$\mathbf{A}_t = \begin{bmatrix} 1 & \Delta t & \frac{1}{2}(\Delta t)^2 & 0 & 0 & 0 \\ 0 & 1 & \Delta t & 0 & 0 & 0 \\ 0 & 0 & e^{-\alpha \Delta t} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & \Delta t & \frac{1}{2}(\Delta t)^2 \\ 0 & 0 & 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 0 & 0 & e^{-\alpha \Delta t} \end{bmatrix}$$

$$\mathbf{x}_t = \mathbf{B}_t \mathbf{s}_t + \delta_t$$
$$\mathbf{B}_t = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

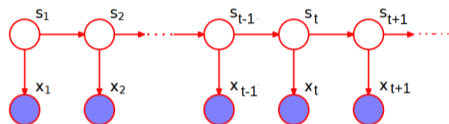
Typical Inference Tasks in Gaussian SSM

- One of the key tasks: Given sequence $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$, infer the latent states $\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots$



- This is usually solved in one of the following two ways
 - Infer the distribution $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$ given the past observations: “Filtering Problem”
 - Infer the distribution $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ given all (past/future) observations: “Smoothing Problem”
- Other tasks we may be interested in
 - Predicting future state(s) given observations seen thus far: $p(\mathbf{s}_{t+h} | \mathbf{x}_1, \dots, \mathbf{x}_t)$ for $h \geq 1$
 - Predict next observation(s) given observations seen thus far: $p(\mathbf{x}_{t+h} | \mathbf{x}_1, \dots, \mathbf{x}_t)$ for $h \geq 1$
- Today, we’ll mainly focus on the filtering problem (solved using the [Kalman Filtering](#) algorithm)

Kalman Filtering



- Recall that $\mathbf{s}_t | \mathbf{s}_{t-1} \sim \mathcal{N}(\mathbf{s}_t | \mathbf{A}_t \mathbf{s}_{t-1}, \mathbf{Q}_t)$ and $\mathbf{x}_t | \mathbf{s}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{B}_t \mathbf{s}_t, \mathbf{R}_t)$
- Let's assume a stationary SSM, i.e., $\mathbf{A}_t = \mathbf{A}$, $\mathbf{B}_t = \mathbf{B}$, $\mathbf{Q}_t = \mathbf{Q}$, and $\mathbf{R}_t = \mathbf{R}$
- Kalman Filtering gives an exact way to infer $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$ in a linear Gaussian SSM
 - Note: The “exactness” assumes we are given \mathbf{A} , \mathbf{B} , \mathbf{Q} , \mathbf{R} are known (or have estimated these)
- Using Bayes rule, our target will be

$$p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t) \propto p(\mathbf{x}_t | \mathbf{s}_t) p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1})$$

- The “prior” above is: $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1}) = \int p(\mathbf{s}_t | \mathbf{s}_{t-1}) p(\mathbf{s}_{t-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1}) d\mathbf{s}_{t-1}$

Kalman Filtering

- Thus the Kalman Filtering problem computes the following

$$p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t) \propto \underbrace{p(\mathbf{x}_t | \mathbf{s}_t)}_{\mathcal{N}(\mathbf{x}_t | \mathbf{B}\mathbf{s}_t, \mathbf{R})} \int \underbrace{p(\mathbf{s}_t | \mathbf{s}_{t-1})}_{\mathcal{N}(\mathbf{s}_t | \mathbf{A}\mathbf{s}_{t-1}, \mathbf{Q})} p(\mathbf{s}_{t-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1}) d\mathbf{s}_{t-1}$$

- Note that the LHS is the **posterior on \mathbf{s}_t** , the RHS consists of a **posterior on \mathbf{s}_{t-1}**
- This suggests a simple **“forward algorithm”** to recursively compute $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$
 - For Kalman smoothing problem $p(\mathbf{z}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$, a similar recursive **“forward-backward” algorithm** exists (the backup slides contain an illustration for the same)
- In this Linear Gaussian SSM, $p(\mathbf{s}_{t-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1})$ would be a Gaussian, say $\mathcal{N}(\mathbf{s}_{t-1} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$
 - Reason: Starting with $p(\mathbf{s}_0) = \mathcal{N}(\mathbf{s}_0 | \mathbf{0}, \mathbf{I}_K)$, the posterior over \mathbf{s}_t will be Gaussian at each step t
- Also, using Gaussian's properties, we know that

$$\int \mathcal{N}(\mathbf{s}_t | \mathbf{A}\mathbf{s}_{t-1}, \mathbf{Q}) \mathcal{N}(\mathbf{s}_{t-1} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{s}_{t-1} = \mathcal{N}(\mathbf{s}_t | \mathbf{A}\boldsymbol{\mu}, \mathbf{Q} + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

Kalman Filtering

- We can now compute the desired posterior

$$p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t) \propto \mathcal{N}(\mathbf{x}_t | \mathbf{B}\mathbf{s}_t, \mathbf{R}) \times \mathcal{N}(\mathbf{s}_t | \mathbf{A}\boldsymbol{\mu}, \mathbf{Q} + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

- This again is a Gaussian (Gaussian likelihood and Gaussian prior), given by

$$p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t) = \mathcal{N}(\mathbf{s}_t | \boldsymbol{\mu}', \boldsymbol{\Sigma}')$$

where the Gaussian posterior's covariance matrix and mean vector are given by

$$\begin{aligned}\boldsymbol{\Sigma}' &= [(\mathbf{Q} + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)^{-1} + \mathbf{B}^\top \mathbf{R}^{-1} \mathbf{B}]^{-1} \\ \boldsymbol{\mu}' &= \boldsymbol{\Sigma}' [\mathbf{B}^\top \mathbf{R}^{-1} \mathbf{x}_t + (\mathbf{Q} + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)^{-1} \mathbf{A}\boldsymbol{\mu}]\end{aligned}$$

- Thus we get closed form expressions for the parameters $(\boldsymbol{\Sigma}', \boldsymbol{\mu}')$ of $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$ in terms of the parameters $(\boldsymbol{\Sigma}, \boldsymbol{\mu})$ of $p(\mathbf{s}_{t-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1})$

Kalman Filtering: Predicting Future Observations

- We saw how to compute $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$ which was a Gaussian $\mathcal{N}(\mathbf{s}_t | \boldsymbol{\mu}', \boldsymbol{\Sigma}')$
- Often we are also interested in predicting the future observations

$$\begin{aligned} p(\mathbf{x}_{t+1} | \mathbf{x}_1, \dots, \mathbf{x}_t) &= \int p(\mathbf{x}_{t+1} | \mathbf{s}_{t+1}) p(\mathbf{s}_{t+1} | \mathbf{x}_1, \dots, \mathbf{x}_t) d\mathbf{s}_{t+1} \\ &= \int \underbrace{p(\mathbf{x}_{t+1} | \mathbf{s}_{t+1})}_{\mathcal{N}(\mathbf{x}_{t+1} | \mathbf{B}\mathbf{s}_{t+1}, \mathbf{R})} \int \underbrace{p(\mathbf{s}_{t+1} | \mathbf{s}_t)}_{\mathcal{N}(\mathbf{s}_{t+1} | \mathbf{A}\mathbf{s}_t, \mathbf{Q})} \underbrace{p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)}_{\mathcal{N}(\mathbf{s}_t | \boldsymbol{\mu}', \boldsymbol{\Sigma}')} d\mathbf{s}_t d\mathbf{s}_{t+1} \end{aligned}$$

- This requires two integrals but the final result is again a Gaussian (expression not shown here)

Kalman Filtering: Some Notes

- Note that we assumed the LDS parameters $\mathbf{A}_t, \mathbf{B}_t, \mathbf{Q}_t, \mathbf{R}_t$ are known

$$\mathbf{s}_t | \mathbf{s}_{t-1} \sim \mathcal{N}(\mathbf{s}_t | \mathbf{A}_t \mathbf{s}_{t-1}, \mathbf{Q}_t)$$

$$\mathbf{x}_t | \mathbf{s}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{B}_t \mathbf{s}_t, \mathbf{R}_t)$$

- Usually these aren't known (unless we have some domain knowledge about the underlying system)
- We can use iterative methods to estimate these parameters
 - Basically, we can alternate between inferring the states and inferring the parameters
- This can be done using approximate inference methods such as EM, MCMC, or VB

Other Extensions of SSM/LDS

- Nonlinear dynamical systems: Assume state-transition and observation models to be nonlinear

$$\begin{aligned}\mathbf{s}_t | \mathbf{s}_{t-1} &= g(\mathbf{s}_{t-1}) + \epsilon_t \\ \mathbf{x}_t | \mathbf{s}_t &= h(\mathbf{s}_t) + \delta_t\end{aligned}$$

- The functions g and h can be nonlinear functions, modeled using deep neural nets, or GP. Another way is to model these as linear approximations of nonlinear functions (Extended Kalman Filter)
- **Switching LDS/SSM**: Assumes data to be generated from a **mixture of M LDS/SSM**
 - For each observation \mathbf{x}_t , first draw a cluster id $c_t \in \{1, \dots, M\}$ from a multinoulli
 - Suppose $c_t = m$. Now generate the observation \mathbf{x}_t using the the m -th LDS/SSM

$$\begin{aligned}\mathbf{s}_t | \mathbf{s}_{t-1}, c_t = m &\sim \mathcal{N}(\mathbf{s}_t | \mathbf{A}^{(m)} \mathbf{s}_{t-1}, \mathbf{Q}^{(m)}) \\ \mathbf{x}_t | \mathbf{s}_t, c_t = m &\sim \mathcal{N}(\mathbf{x}_t | \mathbf{B}^{(m)} \mathbf{s}_t, \mathbf{R}^{(m)})\end{aligned}$$

- It's a hybrid LDS – the “state” consists of two latent variables c_t, \mathbf{z}_t (discrete and continuous)

Summary

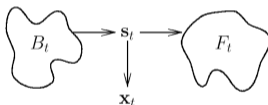
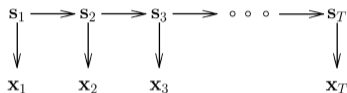
- SSM/LDS allows modeling non i.i.d. sequential data
- Gaussian assumption on transition/observation models helps inference considerably
- These basic models have been extended to more sophisticated models, e.g.,
 - Non-Gaussian LDS
 - Deep LDS
- Inference for HMM is also based on similar principles (e.g., forward and forward-backward algorithm), except that the latent variables are discrete
- The general principle (time-evolving latent variables) can be applied in a wide range of probabilistic models to enable them handle dynamic/time-evolving data
 - E.g., in LDA, we can make the topic assignments of adjacent words follow a Markov relationship (results in an HMM-LDA type model)

Backup Slides: Kalman Smoothing

Kalman Smoothing in SSMs

Goal: Infer $p(\mathbf{s}_t | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ given all the observations (both past and future)

Note that each state variable \mathbf{s}_t separates the graph into three independent parts



$$B_t = \{\mathbf{x}_1 \dots \mathbf{x}_{t-1}, \mathbf{s}_1 \dots \mathbf{s}_{t-1}\}$$

$$F_t = \{\mathbf{x}_{t+1} \dots \mathbf{x}_T, \mathbf{s}_{t+1} \dots \mathbf{s}_T\}$$

$$p(B_t, \mathbf{s}_t, \mathbf{x}_t, F_t) = p(B_t, \mathbf{s}_t) p(\mathbf{x}_t | \mathbf{s}_t) p(F_t | \mathbf{s}_t)$$

Kalman Smoothing in SSMs

- Goal: marginal probability $p(\mathbf{s}_t | \mathbf{x}_1, \dots, \mathbf{x}_T)$ of each state (i.e., smoothing)
- Let's look at the joint probability first:

$$\begin{aligned} p(\mathbf{s}_t, \mathbf{x}_1.. \mathbf{x}_T) &= \int_{\mathbf{s}_1.. \mathbf{s}_{t-1}} \int_{\mathbf{s}_{t+1}.. \mathbf{s}_T} p(B_t, \mathbf{s}_t, \mathbf{x}_t, F_t) \\ &= \left(\int_{\mathbf{s}_1.. \mathbf{s}_{t-1}} p(B_t, \mathbf{s}_t) \right) p(\mathbf{x}_t | \mathbf{s}_t) \left(\int_{\mathbf{s}_{t+1}.. \mathbf{s}_T} p(F_t | \mathbf{s}_t) \right) \\ &= p(B_t^x, \mathbf{s}_t) p(\mathbf{x}_t | \mathbf{s}_t) p(F_t^x | \mathbf{s}_t) \end{aligned}$$

$$B_t^x = \{\mathbf{x}_1.. \mathbf{x}_{t-1}\}$$

$$F_t^x = \{\mathbf{x}_{t+1}.. \mathbf{x}_T\}$$

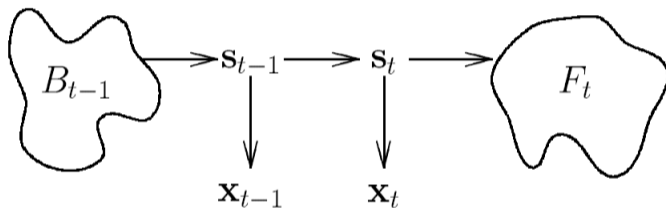
$$\alpha_t(\mathbf{s}_t) = p(B_t^x, \mathbf{s}_t) p(\mathbf{x}_t | \mathbf{s}_t) = p(B_t^x, \mathbf{x}_t, \mathbf{s}_t)$$

$$\beta_t(\mathbf{s}_t) = p(F_t^x | \mathbf{s}_t)$$

$$p(\mathbf{s}_t, \mathbf{x}_1.. \mathbf{x}_T) = \alpha_t(\mathbf{s}_t) \beta_t(\mathbf{s}_t)$$

- From the joint, we can compute $p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{\mathbf{s}_t} p(\mathbf{s}_t, \mathbf{x}_1, \dots, \mathbf{x}_T)$, and $p(\mathbf{s}_t | \mathbf{x}_1, \dots, \mathbf{x}_T)$ using Bayes rule

Estimation via Forward-Backward Recursion



Denote $B_t = B_{t-1} \cup \{s_{t-1}, x_{t-1}\}$ and $F_{t-1} = \{s_t, x_t\} \cup F_t$

Estimation via Forward-Backward Recursion

Denote $B_t = B_{t-1} \cup \{\mathbf{s}_{t-1}, \mathbf{x}_{t-1}\}$ and $F_{t-1} = \{\mathbf{s}_t, \mathbf{x}_t\} \cup F_t$

Can compute α and β recursively

$$\begin{aligned}\alpha_t(\mathbf{s}_t) &= p(\mathbf{x}_t|\mathbf{s}_t)p(B_t^x, \mathbf{s}_t) = p(\mathbf{x}_t|\mathbf{s}_t) \int_{\mathbf{z}} p(B_{t-1}^x, \mathbf{s}_{t-1} = \mathbf{z}, \mathbf{x}_{t-1}, \mathbf{s}_t) \\ &= p(\mathbf{x}_t|\mathbf{s}_t) \int_{\mathbf{z}} p(B_{t-1}^x, \mathbf{s}_{t-1} = \mathbf{z})p(\mathbf{x}_{t-1}|\mathbf{s}_{t-1} = \mathbf{z})p(\mathbf{s}_t|\mathbf{s}_{t-1} = \mathbf{z}) \\ &= p(\mathbf{x}_t|\mathbf{s}_t) \int_{\mathbf{z}} p(\mathbf{s}_t|\mathbf{s}_{t-1} = \mathbf{z})\alpha_{t-1}(\mathbf{z})\end{aligned}$$

Forward recursion for α

$$\begin{aligned}\beta_{t-1}(\mathbf{s}_{t-1}) &= p(F_{t-1}^x|\mathbf{s}_{t-1}) = \int_{\mathbf{z}} p(\mathbf{s}_t = \mathbf{z}, \mathbf{x}_t, F_t^x|\mathbf{s}_{t-1}) \\ &= \int_{\mathbf{z}} p(\mathbf{s}_t = \mathbf{z}|\mathbf{s}_{t-1})p(\mathbf{x}_t|\mathbf{s}_t = \mathbf{z})p(F_t^x|\mathbf{s}_t = \mathbf{z}) \\ &= \int_{\mathbf{z}} p(\mathbf{s}_t = \mathbf{z}|\mathbf{s}_{t-1})p(\mathbf{x}_t|\mathbf{s}_t = \mathbf{z})\beta_t(\mathbf{z})\end{aligned}$$

Backward recursion for β

Initialize as $\alpha_1(\mathbf{s}_1) = p(\mathbf{s}_1)p(\mathbf{x}_1|\mathbf{s}_1)$ and $\beta_T(\mathbf{s}_T) = 1$