

State Space Models

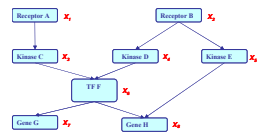
Probabilistic Graphical Models (10-708)

Lecture 13, part II
Nov 5th, 2007

Eric Xing

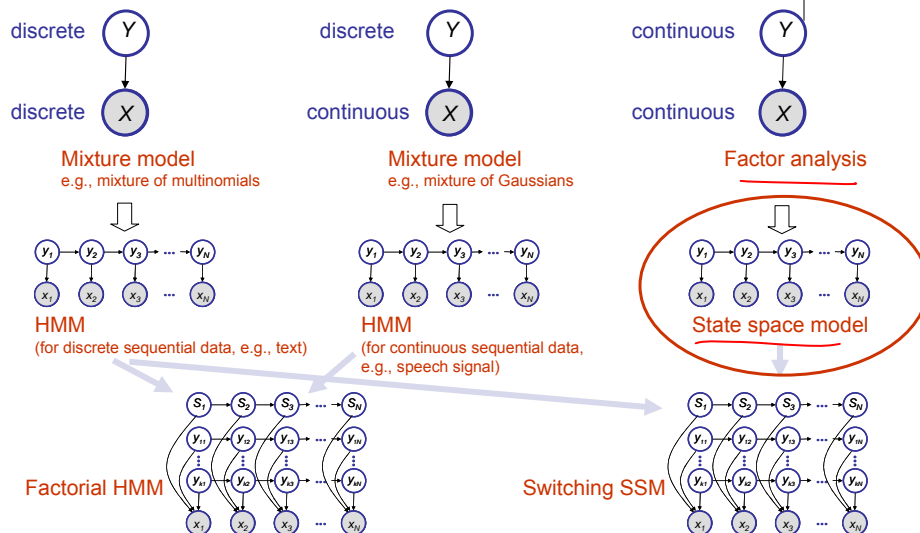
Reading: J-Chap. 15, K&F chapter 19.1 -19.3

2



1

A road map to more complex dynamic models

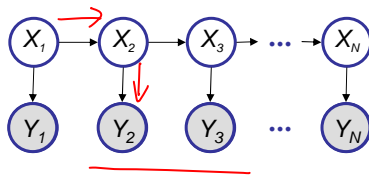


Eric Xing

2

State space models (SSM):

- A sequential FA or a continuous state HMM



$$\underline{\mathbf{x}}_t = A \underline{\mathbf{x}}_{t-1} + G \underline{\mathbf{w}}_t$$

$$\underline{\mathbf{y}}_t = C \underline{\mathbf{x}}_t + \underline{\mathbf{v}}_t$$

$$\underline{\mathbf{w}}_t \sim \mathcal{N}(0; \underline{Q}), \quad \underline{\mathbf{v}}_t \sim \mathcal{N}(0; \underline{R})$$

$$\underline{\mathbf{x}}_0 \sim \mathcal{N}(0; \underline{\Sigma}_0),$$

This is a linear dynamic system.

- In general,

$$\underline{\mathbf{x}}_t = f(\underline{\mathbf{x}}_{t-1}) + G \underline{\mathbf{w}}_t$$

$$\underline{\mathbf{y}}_t = g(\underline{\mathbf{x}}_{t-1}) + \underline{\mathbf{v}}_t$$

where f is an (arbitrary) dynamic model, and g is an (arbitrary) observation model

Eric Xing

3

LDS for 2D tracking

- Dynamics: new position = old position + $\Delta \times$ velocity + noise
(constant velocity model, Gaussian noise)

$$\begin{aligned} \text{positions} \quad \left\{ \begin{pmatrix} x_t^1 \\ x_t^2 \end{pmatrix} \right. &= \begin{pmatrix} 1 & 0 & \Delta & 0 \\ 0 & 1 & 0 & \Delta \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_{t-1}^1 \\ x_{t-1}^2 \\ \dot{x}_{t-1}^1 \\ \dot{x}_{t-1}^2 \end{pmatrix} + \text{noise} \\ \text{Velocity} \quad \left\{ \begin{pmatrix} \dot{x}_t^1 \\ \dot{x}_t^2 \end{pmatrix} \right. & \end{aligned}$$

$\underline{x}_t^1 = \underline{x}_{t-1}^1 + \Delta \underline{\dot{x}}_{t-1}^1$

- Observation: project out first two components (we observe Cartesian position of object - linear!)

$$\underline{\mathbf{y}}_t = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ \dot{x}_t^1 \\ \dot{x}_t^2 \end{pmatrix} + \text{noise}$$

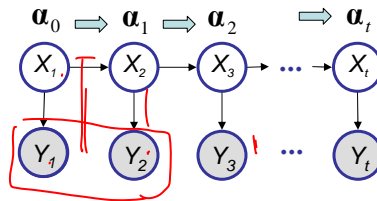
Eric Xing

4

The inference problem 1

- Filtering \rightarrow given y_1, \dots, y_t , estimate x_t : $P(x_t | y_{1:t}) = \alpha(x_t)$
 - The **Kalman filter** is a way to perform exact **online inference** (sequential Bayesian updating) in an LDS.
 - It is the Gaussian analog of the **forward algorithm** for HMMs:

$$p(X_t = i | y_{1:t}) = \alpha_t^i \underbrace{p(y_t | X_t = i)} \sum_j \underbrace{p(X_t = i | X_{t-1} = j)} \alpha_{t-1}^j$$

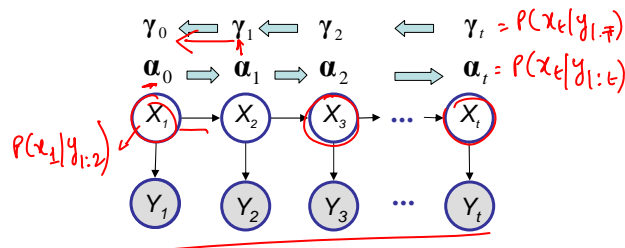


Eric Xing

5

The inference problem 2

- Smoothing \rightarrow given y_1, \dots, y_T , estimate x_t ($t < T$): $P(x_t | y_{1:T}) = \gamma_t(x)$
 - The Rauch-Tung-Strievel smoother is a way to perform exact off-line inference in an LDS. It is the Gaussian analog of the forwards-backwards (alpha-gamma) algorithm:

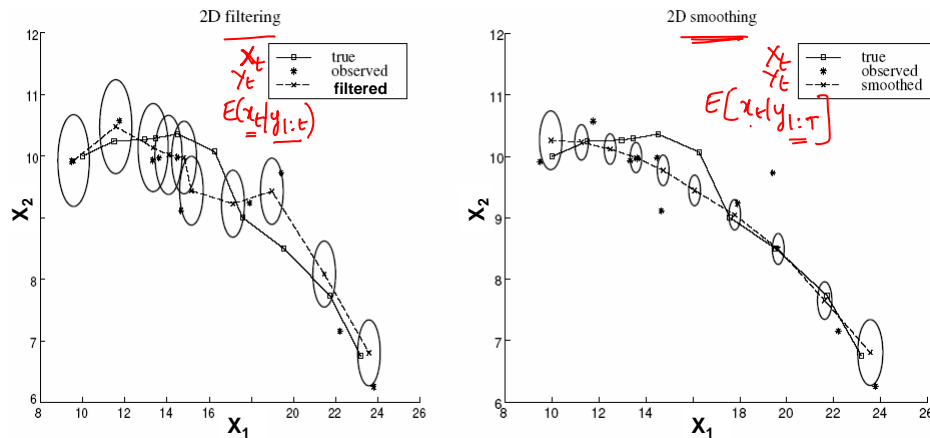


$$p(X_t = i | y_{1:T}) = \gamma_t^i \sum_j \alpha_t^j \underbrace{P(X_{t+1}^j | X_t^i)} \gamma_{t+1}^j$$

Eric Xing

6

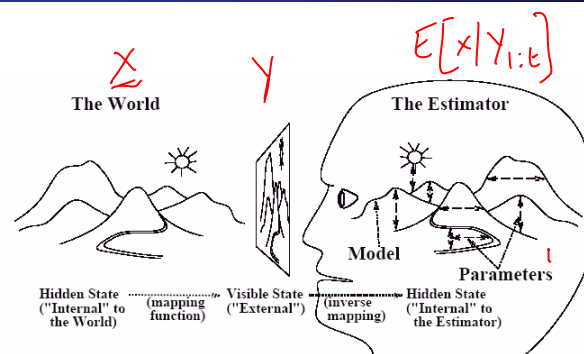
2D tracking



Eric Xing

7

Kalman filtering in the brain?



Eric Xing

8

Kalman filtering derivation

- Since all CPDs are linear Gaussian, the system defines a large multivariate Gaussian.
 - Hence all marginals are Gaussian.
 - Hence we can represent the belief state $p(X_t | y_{1:t})$ as a Gaussian
 - mean $\hat{x}_{t|t} \equiv E(X_t | y_{1:t}, y_t)$ $\hat{x}_{t|t} := \alpha E(x | y_{1:t})$
 - covariance $P_{t|t} \equiv E(X_t X_t^T | y_{1:t}, y_t)$ $P(x_t | y_{1:t})$ $P(x_t | y_{1:t})$
 - Hence, instead of marginalization for message passing, we will directly estimate the means and covariances of the required marginals
 - It is common to work with the inverse covariance (precision) matrix $P_{t|t}^{-1}$; this is called information form.

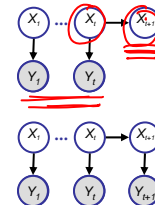
Eric Xing

9

Kalman filtering derivation

- Kalman filtering is a recursive procedure to update the belief state: $p(x_t | y_{1:t}) \rightarrow p(x_{t+1} | y_{1:t+1})$

- Predict step: compute $p(X_{t+1} | y_{1:t})$ from prior belief $p(X_t | y_{1:t})$ and dynamical model $p(X_{t+1} | X_t)$ --- time update
- Update step: compute new belief $p(X_{t+1} | y_{1:t+1})$ from prediction $p(X_{t+1} | y_{1:t})$, observation y_{t+1} and observation model $p(y_{t+1} | X_{t+1})$ --- measurement update



$$p(\underline{x_{t+1}} | \underline{y_{1:t+1}}) = \sum_{x_t} p(x_{t+1}, x_t | y_{1:t}) = \sum_{x_t} p(x_{t+1} | x_t) p(x_t | y_{1:t})$$

$$= \sum_{x_t} p(x_{t+1} | x_t) \cdot p(x_t | y_{1:t})$$

$$\underline{p(x_{t+1} | y_{1:t+1})} = p(x_{t+1} | y_{1:t}, y_{t+1})$$

$$\propto p(y_{t+1} | y_{1:t}, x_{t+1}) p(x_{t+1} | y_{1:t})$$

$$\propto p(y_{t+1} | x_{t+1}) p(x_{t+1} | y_{1:t})$$

Eric Xing

10

Predict step

- Dynamical Model: $\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{G}\mathbf{w}_t$, $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}; \mathbf{Q})$

- One step ahead prediction of state:

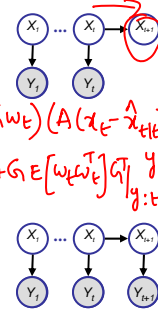
$$E[\mathbf{x}_{t+1} | \mathbf{y}_{1:t}] = \hat{\mathbf{x}}_{t+1|t} = E[\mathbf{A}\mathbf{x}_t + \mathbf{G}\mathbf{w}_t | \mathbf{y}_{1:t}] = \mathbf{A}\hat{\mathbf{x}}_{t|t} + \mathbf{0}$$

$$E[(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t})(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t})^T | \mathbf{y}_{1:t}] = E[(\mathbf{A}(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t}) + \mathbf{G}\mathbf{w}_t)(\mathbf{A}(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t}) + \mathbf{G}\mathbf{w}_t)^T | \mathbf{y}_{1:t}]$$

$$= \mathbf{A} E[(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t})(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t})^T | \mathbf{y}_{1:t}] \mathbf{A}^T + \mathbf{G} E[\mathbf{w}_t \mathbf{w}_t^T] \mathbf{G}^T$$

$$= \mathbf{A} \mathbf{P}_{t|t} \mathbf{A}^T + \mathbf{G} \mathbf{Q} \mathbf{G}^T$$

$$P_{t+1|t} = \mathbf{A} \mathbf{P}_{t|t} \mathbf{A}^T + \mathbf{G} \mathbf{Q} \mathbf{G}^T$$



- Observation model: $\mathbf{y}_{t+1} = \mathbf{C}\mathbf{x}_{t+1} + \mathbf{v}_{t+1}$, $\mathbf{v}_{t+1} \sim \mathcal{N}(\mathbf{0}; \mathbf{R})$

- One step ahead prediction of observation:

$$\hat{\mathbf{y}}_{t+1|t} = \mathbf{C} \hat{\mathbf{x}}_{t+1|t}$$

$$E[(\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t})(\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t})^T | \mathbf{y}_{1:t}] = E[(\mathbf{C}(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t}) + \mathbf{v}_{t+1})(\mathbf{C}(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t}) + \mathbf{v}_{t+1})^T | \mathbf{y}_{1:t}]$$

$$= \mathbf{C} E[(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t})(\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t})^T | \mathbf{y}_{1:t}] \mathbf{C}^T + E[\mathbf{v}_{t+1} \mathbf{v}_{t+1}^T]$$

$$= \mathbf{C} P_{t+1|t} \mathbf{C}^T + \mathbf{R}$$

Eric Xing

Update step

- Summarizing results from previous slide, we have

$$\underline{p(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} | \mathbf{y}_{1:t})} \sim \mathcal{N}(\underline{\mathbf{m}}_{t+1}, \underline{\mathbf{V}}_{t+1}), \text{ where}$$

$$\underline{\mathbf{m}}_{t+1} = \begin{pmatrix} \hat{\mathbf{x}}_{t+1|t} \\ \hat{\mathbf{y}}_{t+1|t} \end{pmatrix}, \quad \underline{\mathbf{V}}_{t+1} = \begin{pmatrix} \mathbf{P}_{t+1|t} & \mathbf{P}_{t+1|t} \mathbf{C}^T \\ \mathbf{C} \mathbf{P}_{t+1|t} & \mathbf{C} \mathbf{P}_{t+1|t} \mathbf{C}^T + \mathbf{R} \end{pmatrix}$$

- Remember the formulas for conditional Gaussian distributions:

$$p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \middle| \mu, \Sigma\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2^m, \mathbf{V}_2^m) \quad p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{m}_2^m = \mu_2$$

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{x}_2 - \mu_2)$$

$$\mathbf{V}_2^m = \Sigma_{22}$$

$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

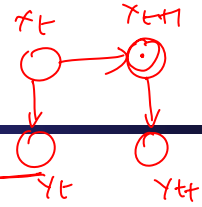
$$p(\mathbf{x}_{t+1} | \mathbf{y}_{t+1}, \mathbf{y}_{1:t})$$

$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1})$$

Eric Xing

12

Kalman Filter



- Measurement updates:

$$\hat{\mathbf{x}}_{t+1|t+1} = \hat{\mathbf{x}}_{t+1|t} + K_{t+1}(\mathbf{y}_{t+1} - C\hat{\mathbf{x}}_{t+1|t})$$

$$P_{t+1|t+1} = P_{t+1|t} - K_{t+1}CP_{t+1|t}$$

- where K_{t+1} is the *Kalman gain matrix*

$$K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1}$$

Example of KF in 1D

- Consider noisy observations of a 1D particle doing a random walk:

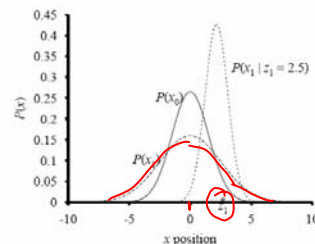
$$\mathbf{x}_{t|t-1} = \mathbf{x}_{t-1} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(0, \sigma_x) \quad \mathbf{z}_t = \mathbf{x}_t + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(0, \sigma_z)$$

- KF equations: $P_{t+1|t} = AP_{t|t}A^T + GQG^T = \sigma_t + \sigma_x$, $\hat{\mathbf{x}}_{t+1|t} = A\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t}$

$$K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1} = (\sigma_t + \sigma_x)(\sigma_t + \sigma_x + \sigma_z)$$

$$\hat{\mathbf{x}}_{t+1|t+1} = \hat{\mathbf{x}}_{t+1|t} + K_{t+1}(\mathbf{z}_{t+1} - C\hat{\mathbf{x}}_{t+1|t}) = \frac{(\sigma_t + \sigma_x)\mathbf{z}_{t+1} + \sigma_z\hat{\mathbf{x}}_{t|t}}{\sigma_t + \sigma_x + \sigma_z}$$

$$P_{t+1|t+1} = P_{t+1|t} - K_{t+1}CP_{t+1|t} = \frac{(\sigma_t + \sigma_x)\sigma_z}{\sigma_t + \sigma_x + \sigma_z}$$



KF intuition

- The KF update of the mean is

$$\hat{x}_{t+1|t+1} = \hat{x}_{t+1|t} + K_{t+1}(z_{t+1} - C\hat{x}_{t+1|t}) = \frac{(\sigma_t + \sigma_x)z_{t+1} + \sigma_z\hat{x}_{t|t}}{\sigma_t + \sigma_x + \sigma_z}$$

- the term $(z_{t+1} - C\hat{x}_{t+1|t})$ is called the *innovation*
- New belief is convex combination of updates from prior and observation, weighted by Kalman Gain matrix:

$$K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1} = (\sigma_t + \sigma_x)(\sigma_t + \sigma_x + \sigma_z)^{-1}$$

- If the observation is unreliable, σ_z (i.e., R) is large so K_{t+1} is small, so we pay more attention to the prediction.
- If the old prior is unreliable (large σ_t) or the process is very unpredictable (large σ_x), we pay more attention to the observation.

Eric Xing

15

KF, RLS and LMS

- The KF update of the mean is

$$\hat{x}_{t+1|t+1} = A\hat{x}_{t|t} + K_{t+1}(y_{t+1} - C\hat{x}_{t+1|t})$$

$$y_t = Cx_t + v_t$$

- Consider the special case where the hidden state is a constant, $x_t = \theta$, but the "observation matrix" C is a time-varying vector, $C = x_t^T$.
 - Hence the observation model at each time slide, $y_t = x_t^T \theta + v_t$, is a linear regression

- We can estimate recursively using the Kalman filter:

$$\hat{\theta}_{t+1} = \hat{\theta}_t + P_{t+1}R^{-1}(y_{t+1} - x_t^T \hat{\theta}_t)x_t$$

This is called the recursive least squares (RLS) algorithm.

- We can approximate $P_{t+1}R^{-1} \approx \eta_{t+1}$ by a scalar constant. This is called the least mean squares (LMS) algorithm.
- We can adapt η_t online using stochastic approximation theory.

Eric Xing

16

Complexity of one KF step

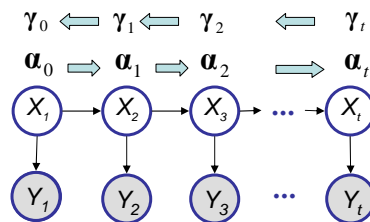
- Let $X_t \in \mathbb{R}^{N_x}$ and $Y_t \in \mathbb{R}^{N_y}$,
- Computing $\underline{P_{t+1|t}} = \underline{A P_t A^T + G Q G^T}$ takes $\underline{O(N_x^3)}$ time, assuming dense P and dense A .
- Computing $\underline{K_{t+1}} = \underline{P_{t+1|t} C^T (C P_{t+1|t} C^T + R)^{-1}}$ takes $\underline{O(N_y^3)}$ time.
- So overall time is, in general, $\underline{\max \{N_x^3, N_y^3\}}$

Eric Xing

17

The inference problem 2

- Smoothing \rightarrow given y_1, \dots, y_T , estimate x_t ($t < T$) $p(x_t | y_{1:T})$
 - The Rauch-Tung-Striebel smoother is a way to perform exact off-line inference in an LDS. It is the Gaussian analog of the forwards-backwards (alpha-gamma) algorithm:



$$p(X_t = i | y_{1:T}) = \gamma_t^i \sum_j \alpha_i^j P(X_{t+1}^j | X_t^i) \gamma_{t+1}^j$$

Eric Xing

18

RTS smoother derivation

$$p(x_t | y_{1:T}) = \gamma(x_t) \rho(x_{t+1} | x_t)$$

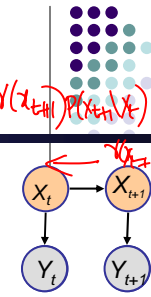
- Smoothing \rightarrow given y_1, \dots, y_T , estimate $P(x_t | y_{1:T})$ ($t < T$)

- Step 1: joint distribution of \underline{x}_t and \underline{x}_{t+1} conditioned on $\underline{y}_{1:t}$

- Use $\underline{x}_{t+1} = A\underline{x}_t + Gw_t; w_t \sim \mathcal{N}(0; Q);$

$$E[\underline{x}_t | y_{1:t}] = \hat{x}_{t|t} \quad E[\underline{x}_{t+1} | y_{1:t}] = \hat{x}_{t+1|t}$$

$$\begin{aligned} & E[(x_t - \hat{x}_{t|t})(x_{t+1} - \hat{x}_{t+1|t})^T | y_{1:t}] \\ &= E[(x_t - \hat{x}_{t|t})(Ax_t + Gw_t - A\hat{x}_{t|t})^T | y_{1:t}] \\ &= E[(x_t - \hat{x}_{t|t})(A(x_t - \hat{x}_{t|t}) + Gw_t)^T | y_{1:t}] \\ &= P_{t|t} A^T \\ & E[(x_{t+1} - \hat{x}_{t+1|t})(x_{t+1} - \hat{x}_{t+1|t})^T | y_{1:t}] = A P_{t|t} A^T + G Q G^T \end{aligned}$$



Eric Xing

19

RTS smoother derivation

- Following the results from previous slide, we need to derive $p(\underline{x}_{t+1}, \underline{x}_t | y_{1:t}) \sim \mathcal{N}(m, V)$, where

$$m = \begin{pmatrix} \hat{x}_{t|t} \\ \hat{x}_{t+1|t} \end{pmatrix}, \quad V = \begin{pmatrix} P_{t|t} & P_{t|t} A^T \\ A P_{t|t} & P_{t+1|t} \end{pmatrix}$$

$$\begin{aligned} & p(x_t | y_{1:t}) \\ & \uparrow \\ & p(x_t | x_{t+1}, y_{1:t}) \end{aligned}$$

- all the quantities here are available after a forward KF pass
- Remember the formulas for conditional Gaussian distributions:

$$p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \mu, \Sigma\right) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right),$$

$$\begin{aligned} p(x_2) &= \mathcal{N}(x_2 | m_2^m, V_2^m) \\ m_2^m &= \mu_2 \\ V_2^m &= \Sigma_{22} \end{aligned}$$

$$\begin{aligned} p(x_1 | x_2) &= \mathcal{N}(x_1 | m_{12}, V_{12}) \\ m_{12} &= \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2) \\ V_{12} &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \end{aligned}$$

$$p(x_t | x_{t+1}, y_{1:t})$$

$$E[x_t | x_{t+1}, y_{0:t}] = \hat{x}_{t|t} + L_t (x_{t+1} - \hat{x}_{t+1|t})$$

$$\text{Var}[x_t | x_{t+1}, y_{0:t}] = P_{t|t} - L_t P_{t+1|t} L_t^T$$

$$L_t = P_{t|t} A^T P_{t+1|t}^{-1}$$

Eric Xing

20

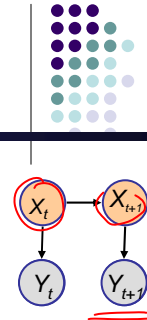
RTS smoother derivation

$$\begin{aligned} E[x_t | x_{t+1}, y_{0:t}] &= \hat{x}_{t|t} + L_t(x_{t+1} - \hat{x}_{t+1|t}) \\ \text{Var}[x_t | x_{t+1}, y_{0:t}] &= P_{t|t} - L_t P_{t+1|t} L_t^T \end{aligned}$$

- Step 2: compute $\hat{x}_{t|T} = E[x_t | y_{0:T}]$ using results above

- Use $E[x_t | x_{t+1}, y_{0:T}] = E[x_t | x_{t+1}, y_{0:t}]$
- Use $E[X|Z] = E[E[X|Y,Z]|Z]$

$$\begin{aligned} \hat{x}_{t|T} E[x_t | y_{0:T}] &= E[E[x_t | x_{t+1}, y_{0:T}] | y_{0:T}] \\ &= E[E[x_t | x_{t+1}, y_{0:t}] | y_{0:T}] \\ &= E[\hat{x}_{t|t} + L_t(x_{t+1} - \hat{x}_{t+1|t}) | y_{0:T}] \\ &= \hat{x}_{t|t} + L_t(\hat{x}_{t+1|T} - \hat{x}_{t+1|t}) \\ \underbrace{P(x_t | x_{t+1}, y_{1:t})}_{\text{red}} &\rightarrow P(x_t | y_{1:T}) \end{aligned}$$



Eric Xing

21

RTS derivation

- Repeat the same process for Variance
 - Refer to Jordan chapter 15

- The RTS smoother results:

$$\hat{x}_{t|T} = \hat{x}_{t|t} + L_t(\hat{x}_{t+1|T} - \hat{x}_{t+1|t})$$

$$P_{t|T} = P_{t|t} + L_t(P_{t+1|T} - P_{t+1|t})L_t^T$$

Eric Xing

22

Learning SSMs

$$\begin{aligned}x_{t+1} &= \boxed{A}x_t + \boxed{B}w_t \\ y_t &= \boxed{C}x_t + \boxed{D}v_t\end{aligned}$$

- Complete log likelihood

$$\begin{aligned}\ell_c(\theta, D) &= \sum_n \log p(x_n, y_n) = \sum_n \log p(x_1) + \sum_n \sum_t \log p(x_{n,t} | x_{n,t-1}) + \sum_n \sum_t \log p(y_{n,t} | x_{n,t}) \\ &= f_1(x_1; \Sigma_0) + f_2(\langle x_t x_{t-1}^T \rangle, \langle x_t x_t^T \rangle, \langle x_t \rangle : \forall t; A, Q, G) + f_3(\langle x_t x_t^T \rangle, \langle x_t \rangle : \forall t; C, R)\end{aligned}$$

- EM

- E-step: compute $\langle \underline{x_t x_{t-1}^T} \rangle, \langle \underline{x_t x_t^T} \rangle, \langle \underline{x_t} \rangle | y_1, \dots, y_T$

these quantities can be inferred via KF and RTS filters, etc.,

$$e, g, \langle x_t x_t^T \rangle \equiv \text{var}(x_t, x_t^T) + E(x_t)^2 = P_{t|T} + \hat{x}_{t|T}^2$$

- M-step: MLE using

$$\begin{aligned}\ell_c(\theta, D) &= f_1(x_1; \Sigma_0) + f_2(\langle x_t x_{t-1}^T \rangle, \langle x_t x_t^T \rangle, \langle x_t \rangle : \forall t; A, Q, G) + f_3(\langle x_t x_t^T \rangle, \langle x_t \rangle : \forall t; C, R) \\ &\text{c.f., M-step in factor analysis}\end{aligned}$$

Eric Xing

23

Nonlinear systems

- In robotics and other problems, the motion model and the observation model are often nonlinear:

$$x_t = f(x_{t-1}) + w_t, \quad y_t = g(x_t) + v_t$$

- An optimal closed form solution to the filtering problem is no longer possible.
- The nonlinear functions f and g are sometimes represented by neural networks (multi-layer perceptrons or radial basis function networks).
- The parameters of f and g may be learned offline using EM, where we do gradient descent (back propagation) in the M step, c.f. learning a MRF/CRF with hidden nodes.
- Or we may learn the parameters online by adding them to the state space: $x_t' = (\underline{x_t}, \theta)$. This makes the problem even more nonlinear.

Eric Xing

24

Extended Kalman Filter (EKF)

- The basic idea of the EKF is to linearize f and g using a second order Taylor expansion, and then apply the standard KF.
- i.e., we approximate a stationary nonlinear system with a non-stationary linear system.

$$x_t = f(\hat{x}_{t|t-1}) + A_{\hat{x}_{t|t-1}}(x_{t-1} - \hat{x}_{t-1|t-1}) + w_t$$

$$y_t = g(\hat{x}_{t|t-1}) + C_{\hat{x}_{t|t-1}}(x_t - \hat{x}_{t|t-1}) + v_t$$

where $\hat{x}_{t|t-1} = f(\hat{x}_{t-1|t-1})$ and $A_{\hat{x}} \stackrel{\text{def}}{=} \frac{\partial f}{\partial x} \Big|_{\hat{x}}$ and $C_{\hat{x}} \stackrel{\text{def}}{=} \frac{\partial g}{\partial x} \Big|_{\hat{x}}$

- The noise covariance (Q and R) is not changed, i.e., the additional error due to linearization is not modeled.

Eric Xing

25

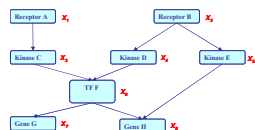
Complex Graphical Models

Probabilistic Graphical Models (10-708)

Lecture 14, Nov 5th, 2007

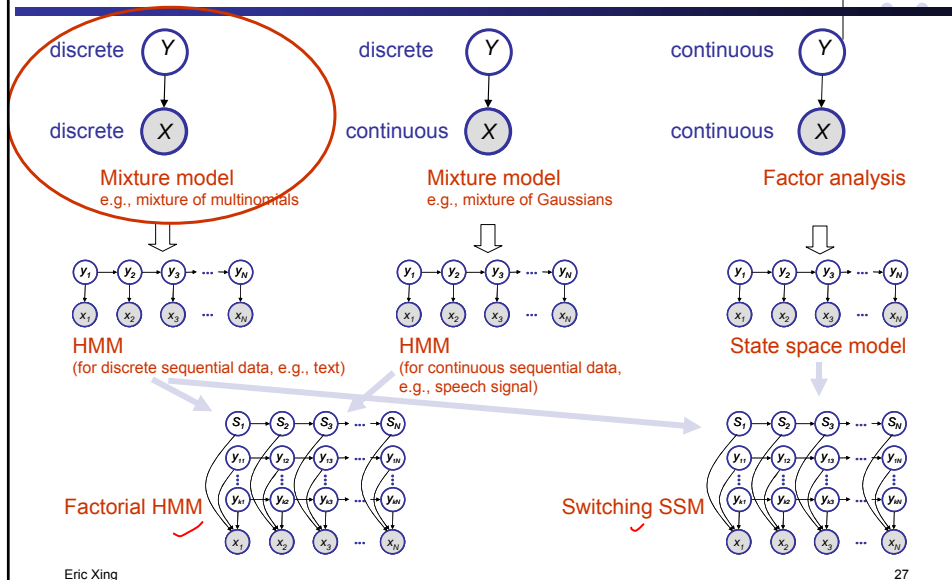
Eric Xing

Reading: K&F chapter 20.1 - 20.3



26

A road map to more complex dynamic models



Eric Xing

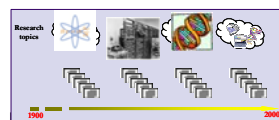
27

NLP and Data Mining



We want:

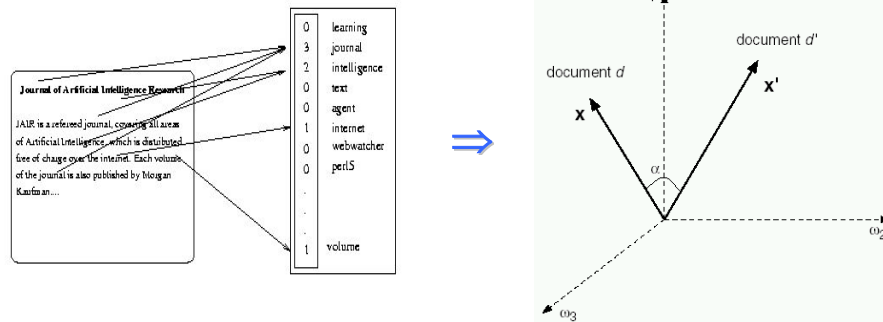
- Semantic-based search
- infer topics and categorize documents
- Multimedia inference
- Automatic translation
- Predict how topics evolve
- ...



28

The Vector Space Model

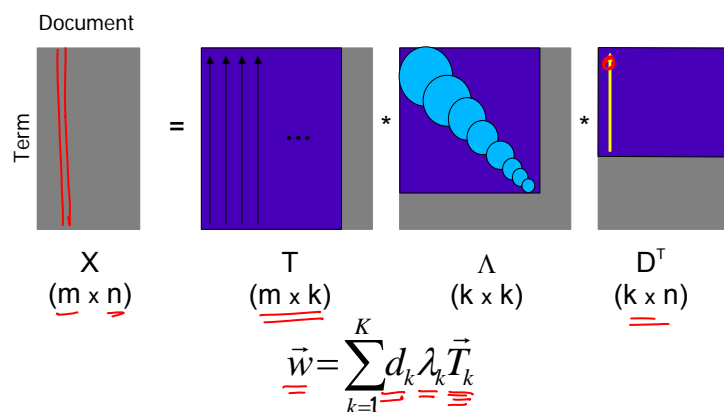
- Represent each document by a high-dimensional vector in the space of words



Eric Xing

29

Latent Semantic Indexing

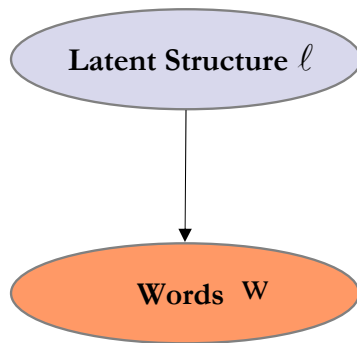


- LSA does not define a properly normalized probability distribution of observed and latent entities
 - Does not support probabilistic reasoning under uncertainty and data fusion

Eric Xing

30

Latent Semantic Structure



Distribution over words

$$\underline{P(\mathbf{w})} = \sum_{\ell} \underline{P(\mathbf{w}, \ell)}$$

Inferring latent structure

$$\underline{P(\ell | \mathbf{w})} = \frac{P(\mathbf{w} | \ell)P(\ell)}{P(\mathbf{w})}$$

Prediction

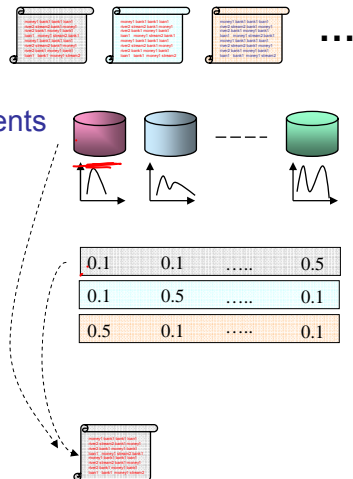
$$\underline{P(w_{n+1} | \mathbf{w})} = \dots$$

Eric Xing

31

Admixture Models

- Objects are **bags** of elements
- Mixtures are **distributions** over elements
- Objects have **mixing** vector θ
 - Represents each mixtures' contributions
- Object is **generated** as follows:
 - Pick a mixture component from θ
 - Pick an element from that component



Eric Xing

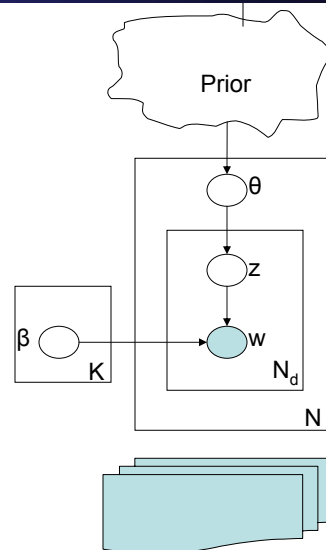
32

Topic Models = Admixture Models

Generating a document

- Draw θ from the prior
- For each word n
- Draw z_n from $\text{multinomial}(\theta)$
 - Draw $w_n | z_n, \{\beta_{1:k}\}$ from $\text{multinomial}(\beta_{z_n})$

Which prior to use?



Eric Xing

33

Choice of Prior

- Dirichlet (LDA) (Blei et al. 2003)
 - Conjugate prior means efficient inference
 - Can **only** capture variations in each topic's intensity **independently**
- Logistic Normal (CTM=LoNTAM) (Blei & Lafferty 2005, Ahmed & Xing 2006)
 - Capture the intuition that some topics are highly correlated and can rise up in intensity together
 - **Not** a conjugate prior implies **hard** inference

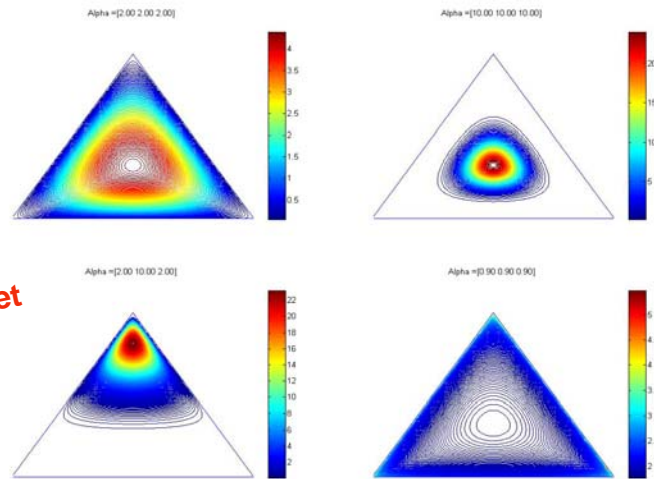
Eric Xing

34

Logistic Normal Vs. Dirichlet



Dirichlet



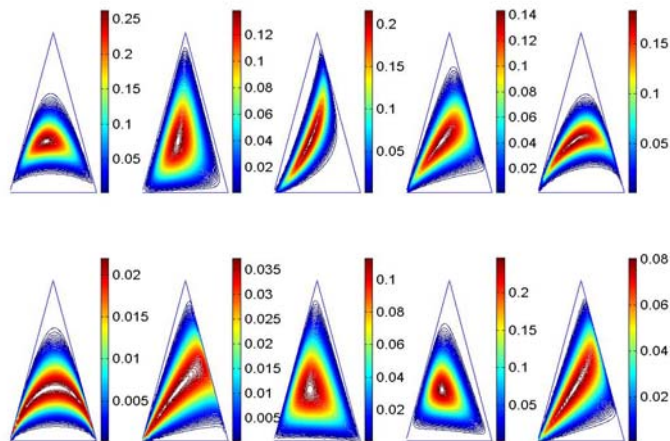
Eric Xing

35

Logistic Normal Vs. Dirichlet



Logistic Normal

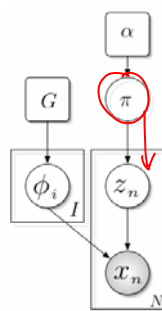


Eric Xing

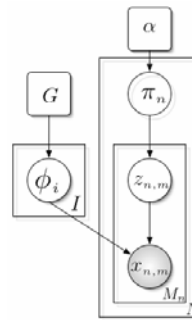
36

Mixed Membership Model (M³)

- Mixture versus admixture



A Bayesian mixture model



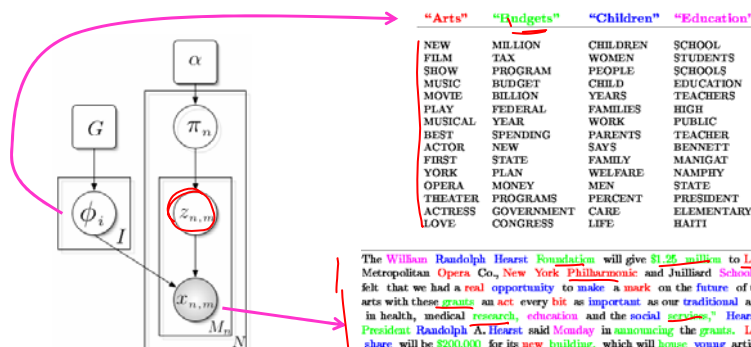
A Bayesian admixture model:
Mixed membership model

Eric Xing

37

Latent Dirichlet Allocation: M³ in text mining

- A document is a bag of words each generated from a randomly selected topic



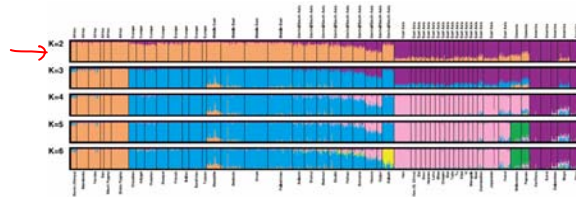
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants, an act every bit as important as our traditional areas of support in health, medical research, education and the social sciences," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Eric Xing

38

-

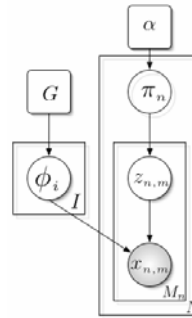


Noah A. Rosenberg,^{1*} Jonathan K. Pritchard,² James L. Weber,³
Howard M. Cann,⁴ Kenneth K. Kidd,⁵ Lev A. Zhivotovskiy,⁶
Marcus W. Feldman⁷

- Eric Xing

39

-
- Figure 1: A graphical model for the multi-view problem. The model consists of three main components: a generative model G , a prior distribution π , and a likelihood function. G takes input ϕ_i and produces output x_n . π takes input α and produces output z_n . The likelihood function takes input z_n and produces output x_n . The output x_n is shared between G and the likelihood function.



$$p(\mathcal{D}) = \sum_{(z_{n,m})} \int \cdots \int \left(\prod_n \left(\prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | \mathcal{G}) d\pi_1 \cdots d\pi_N d\phi$$

- $$p(\pi_n | \mathcal{D}) \sim \sum_{\{z_{n,m}\}} \int \left(\prod_n \left(\prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | \mathcal{G}) d\pi_{-i} d\phi$$

Eric Xing

40

Approaches to inference



- Exact inference algorithms
 - The elimination algorithm
 - The junction tree algorithms
- Approximate inference techniques
 - Monte Carlo algorithms: ✓
 - Stochastic simulation / sampling methods
 - Markov chain Monte Carlo methods
 - Variational algorithms: ✓
 - Belief propagation ←
 - Variational inference