

Readings:

K&F: 6.1, 6.2, 6.3, 14.1, 14.2, 14.3, 14.4,

Kalman Filters Gaussian MNs

Graphical Models – 10708
Carlos Guestrin
Carnegie Mellon University
December 1st, 2008

1

Multivariate Gaussian

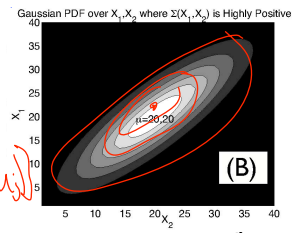
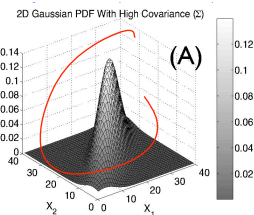
$$p(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

Mean vector:

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} \quad \mu_i = E[X_i]$$

Covariance matrix:

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22}^2 & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33}^2 \end{pmatrix} \quad \begin{aligned} &\sigma_{32} = \sigma_{23} \\ &\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)] \end{aligned}$$



Conditioning a Gaussian

Joint Gaussian:

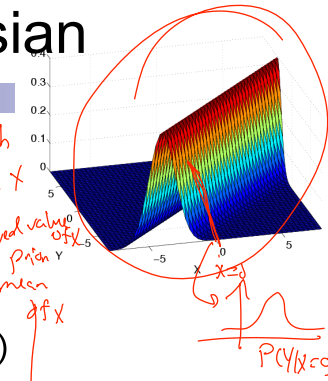
$$\square p(X, Y) \sim N(\mu; \Sigma)$$

Conditional linear Gaussian:

$$\square p(Y|X) \sim N(\mu_{Y|X}, \sigma_{Y|X}^2)$$

$$\mu_{Y|X=x} = \mu_Y + \frac{\sigma_{YX}}{\sigma_X^2}(x - \mu_X)$$

$$\sigma_{Y|X}^2 = \sigma_Y^2 - \frac{\sigma_{YX}^2}{\sigma_X^2}$$



for each fixed x
 observed value of x
 prior mean of x
 $P(Y|x=0)$
 posterior variance doesn't depend on observed value!!
 $\sigma_{Y|X}^2 \leq \sigma_Y^2$ ($\sigma_{Y|X}^2 = \sigma_Y^2$ iff $Y \perp X$)
 observations always decrease variance

3

Gaussian is a "Linear Model"

Conditional linear Gaussian:

$$\square p(Y|X) \sim N(\beta_0 + \beta X; \sigma^2)$$

$$\mu_{Y|X} = \mu_Y + \frac{\sigma_{YX}}{\sigma_X^2}(x - \mu_X)$$

$$\sigma_{Y|X}^2 = \sigma_Y^2 - \frac{\sigma_{YX}^2}{\sigma_X^2}$$

$$\begin{aligned} \mu_{Y|X} &= \underbrace{\mu_Y - \frac{\sigma_{YX}}{\sigma_X^2} \mu_X}_{\beta_0} + \underbrace{\frac{\sigma_{YX}}{\sigma_X^2}}_{\beta} x \\ &= \beta_0 + \beta x \end{aligned}$$

$$\text{equivalently: } Y = \beta_0 + \beta X + \varepsilon \quad \leftarrow \begin{array}{l} \text{white} \\ \text{Noise} \\ N(0, \sigma_{Y|X}^2) \end{array}$$

4

Conditioning a Gaussian

Joint Gaussian:

$$p(X,Y) \sim N(\mu; \Sigma)$$

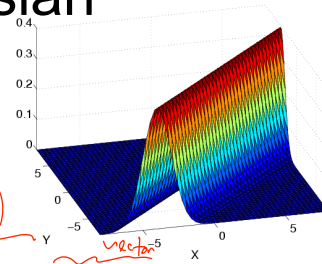
Conditional linear Gaussian:

$$p(Y|X) \sim N(\mu_{Y|X}, \Sigma_{YY|X})$$

$$\mu_{Y|X} = \mu_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (x - \mu_x)$$

$$\Sigma_{YY|X} = \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$$

Covariance of the posterior



5

Conditional Linear Gaussian (CLG) – general case

Conditional linear Gaussian:

$$p(Y|X) \sim N(\beta_0 + BX; \Sigma_{YY|X})$$

$$\mu_{Y|X} = \mu_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (x - \mu_x) = \mu_Y - \Sigma_{YX} \Sigma_{XX}^{-1} \mu_x + \Sigma_{YX} \Sigma_{XX}^{-1} x$$

$$\Sigma_{YY|X} = \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} = \beta_0 + Bx$$

$$Y = \beta_0 + BX + \epsilon \leftarrow \begin{matrix} \text{white Noise} \\ N(\vec{0}, \Sigma_{YY|X}) \end{matrix}$$

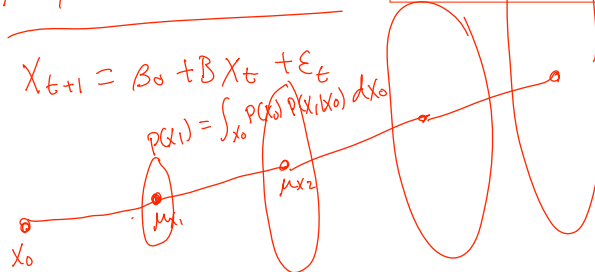
6

Understanding a linear Gaussian – the 2d case

- Variance increases over time (motion noise adds up)
- Object doesn't necessarily move in a straight line

$$Y = \beta_0 + B X + \varepsilon$$

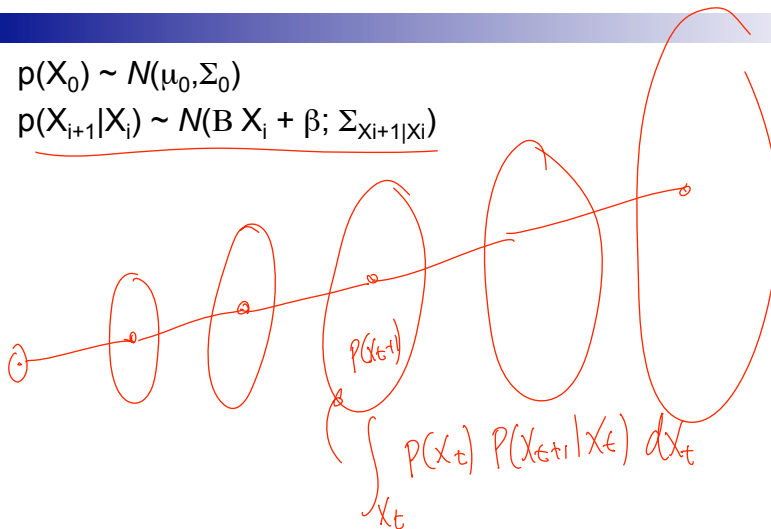
$$X_{t+1} = \beta_0 + B X_t + \varepsilon_t$$



7

Tracking with a Gaussian 1

- $p(X_0) \sim N(\mu_0, \Sigma_0)$
- $p(X_{i+1}|X_i) \sim N(B X_i + \beta; \Sigma_{X_{i+1}|X_i})$



8

Tracking with Gaussians 2 –

Making observations

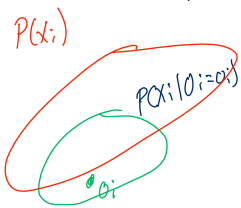
Why Gaussians
1. easy
reasonable many apps
2. central limit theorem
3. maxent...

- We have $p(X_i) \leftarrow$ prior
- Detector observes $O_i = o_i \leftarrow$ observation
- Want to compute $p(X_i | O_i = o_i) \leftarrow$ posterior
- Use Bayes rule: $p(X_i | O_i = o_i) \propto p(X_i) p(O_i = o_i | X_i)$

e.g. camera tracking
 $W \rightarrow$ transforms from
3d position X_i to 2d
camera obs. O_i

- Require a CLG observation model

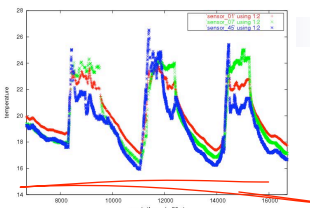
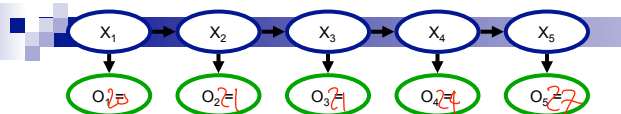
$$p(O_i | X_i) \sim N(W X_i + v; \Sigma_{O_i | X_i})$$



\Rightarrow intuitively
if true location is X_i
 $O_i = v + W X_i + \epsilon \leftarrow N(0, \Sigma_{\epsilon} | X_i)$
simplest case: unshifted, unbiased, $v = 0$
 $W = I$
 \Rightarrow

9

Operations in Kalman filter



- Compute $p(X_t | O_{1:t} = o_{1:t})$
- Start with $p(X_0)$
- At each time step t :
 - **Condition on observation** (posterior) $p(X_t | o_{1:t}) \propto p(X_t | o_{1:t-1}) p(o_t | X_t)$ (likelihood)
 - **Prediction** (Multiply transition model) $p(X_{t+1}, X_t | o_{1:t}) = p(X_{t+1} | X_t) p(X_t | o_{1:t})$
 - **Roll-up** (marginalize previous time step) $p(X_{t+1} | o_{1:t}) = \int_{X_t} p(X_{t+1}, x_t | o_{1:t}) dx_t$
- I'll describe one implementation of KF, there are others
 - Information filter

10

Exponential family representation of Gaussian: Canonical Form $e^{f(x)+c} \propto e^{f(x)}$

$$p(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

$\mathcal{L} \geq 0$
positive semi-definite

$$\begin{aligned} &\propto \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mu^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu^T \Sigma^{-1} \mu \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mu^T \Sigma^{-1} \mathbf{x} \right\} \quad \Sigma^{-1} = \Lambda, \quad \mu^T \Sigma^{-1} = \eta^T \\ &= \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Lambda \mathbf{x} + \eta^T \mathbf{x} \right\} \\ &= \exp \left\{ -\frac{1}{2} \sum_{ij} \lambda_{ij} x_i x_j + \sum_i \eta_i x_i \right\} = \exp \left\{ \sum_{ij} \lambda_{ij} f_{ij}(x_i, x_j) + \sum_i \eta_i f_i(x_i) \right\} \\ &\quad \rightarrow \text{features} \quad \begin{cases} f_i(x) = x_i \\ f_{ij}(x) = x_i x_j \end{cases} \end{aligned}$$

log linear model

11

Canonical form

$$\begin{aligned} p(X_1, \dots, X_n) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\} \\ &= K \exp \left\{ \eta^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \Lambda \mathbf{x} \right\} \end{aligned}$$

- Standard form and canonical forms are related:

$$\mu = \Lambda^{-1} \eta$$

$$\Sigma = \Lambda^{-1}$$

- Conditioning is easy in canonical form
- Marginalization easy in standard form

12

Conditioning in canonical form

$$p(X_t | o_{1:t}) \propto p(X_t | o_{1:t-1})p(o_t | X_t)$$

position prior observation

■ First multiply: $p(A, B) = p(A)p(B | A)$

$$p(A) : \eta_1, \Lambda_1$$

$$p(B | A) : \eta_2, \Lambda_2$$

$$p(A, B) : \eta_3 = \eta_1 + \eta_2, \Lambda_3 = \Lambda_1 + \Lambda_2$$

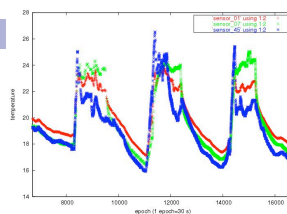
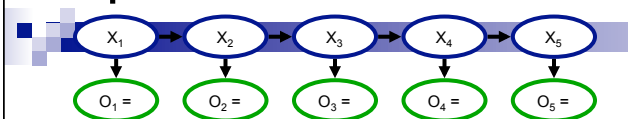
$$\Lambda_3 = \begin{pmatrix} \Lambda_{AA} & \Lambda_{AB} \\ \Lambda_{BA} & \Lambda_{BB} \end{pmatrix} \quad \eta_3 = \begin{pmatrix} \eta_A \\ \eta_B \end{pmatrix}$$

■ Then, condition on value $B = y$ $p(A | B = y)$

$$\begin{aligned} \eta_{A|B=y} &= \eta_A - \Lambda_{AB} \cdot y \\ \Lambda_{AA|B=y} &= \Lambda_{AA} \end{aligned}$$

13

Operations in Kalman filter



■ Compute $p(X_t | O_{1:t} = o_{1:t})$

■ Start with $p(X_0)$

■ At each time step t :

□ **Condition** on observation

$$p(X_t | o_{1:t}) \propto p(X_t | o_{1:t-1})p(o_t | X_t)$$

□ **Prediction** (Multiply transition model)

$$p(X_{t+1}, X_t | o_{1:t}) = p(X_{t+1} | X_t)p(X_t | o_{1:t})$$

□ **Roll-up** (marginalize previous time step)

$$p(X_{t+1} | o_{1:t}) = \int_{X_t} p(X_{t+1}, x_t | o_{1:t}) dx_t$$

matrix addition / selection of submatrices

14

Prediction & roll-up in canonical form

$$p(X_{t+1} | o_{1:t}) = \int_{X_t} p(X_{t+1} | x_t) p(x_t | o_{1:t}) dx_t$$

transition model (CLG) *posterior in current time step* (Gauss)

■ First multiply: $p(A, B) = p(A)p(B | A)$ $\eta = \begin{pmatrix} \eta_A \\ \eta_B \end{pmatrix}$

same as before

■ Then, marginalize X_t : $p(A) = \int_B p(A, b) db$ $\Lambda = \begin{pmatrix} \Lambda_{AA} & \Lambda_{AB} \\ \Lambda_{BA} & \Lambda_{BB} \end{pmatrix}$

marginal

$$\eta_A^m = \eta_A - \Lambda_{AB} \Lambda_{BB}^{-1} \eta_B$$

$$\Lambda_{AA}^m = \Lambda_{AA} - \Lambda_{AB} \Lambda_{BB}^{-1} \Lambda_{BA}$$

marginal

$$p(A) = N(\eta_A^m, \Lambda_{AA}^m)$$

Can also do EM for kalman filter

where does $p(X_{t+1} | x_t)$ come from?
 $p(o_t | x_t)$ learn from data!!
 $p(x_{t+1}, x_t) \propto p(x_{t+1} | x_t) p(x_t)$
 ratio matrix sustains

15

What if observations are not CLG?

- Often observations are not CLG

- CLG if $O_i = B X_i + \beta_o + \varepsilon$

- Consider a motion detector

- $O_i = 1$ if person is likely to be in the region

detector: in room / not in room

- Posterior is not Gaussian



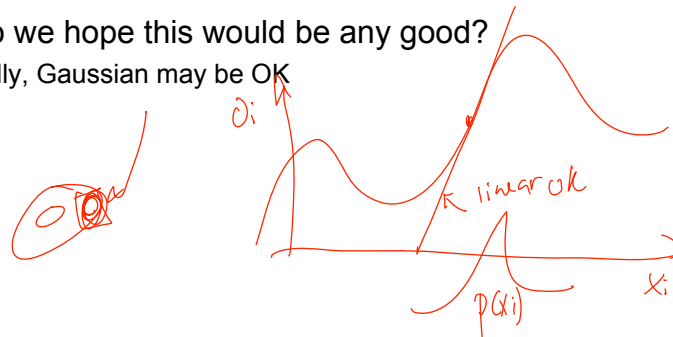
16

Linearization: incorporating non-linear evidence

- $p(O_i|X_i)$ not CLG, but... *not for observation model*
- Find a Gaussian approximation of $p(X_i, O_i) = p(X_i) p(O_i|X_i)$ *but for joint*
- Instantiate evidence $O_i = o_i$ and obtain a Gaussian for $p(X_i|O_i = o_i)$ *approximate*

- Why do we hope this would be any good?

- Locally, Gaussian may be OK



17

Linearization as integration

- Gaussian approximation of $p(X_i, O_i) = p(X_i) p(O_i|X_i)$ *already know $E[X_i]$ $E[X_i^2]$*

- Need to compute moments *mean* *variance*

$$\square E[O_i] = \int O_i p(O_i|x_i) p(x_i) dx_i do_i$$

$$\square E[O_i^2] = \int O_i^2 p(O_i|x_i) p(x_i) dx_i do_i$$

$$\square E[O_i X_i] = \int O_i x_i p(O_i|x_i) p(x_i) dx_i do_i$$

$f(O_i, x_i)$ Gaussian

- Note: Integral is product of a Gaussian with an arbitrary function *plug directly here*

18

Linearization as numerical integration

- **Product of a Gaussian with arbitrary function**

- Effective numerical integration with **Gaussian quadrature** method

- Approximate integral as **weighted sum over integration points**
- Gaussian quadrature defines location of points and weights

- Exact if arbitrary function is **polynomial of bounded degree**

- **Number of integration points exponential** in number of dimensions d

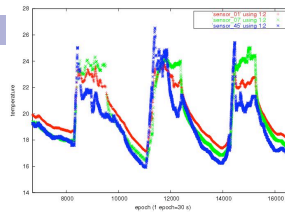
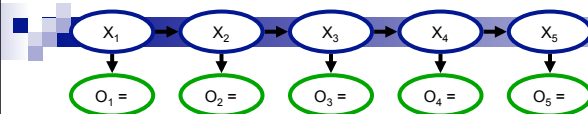
- **Exact monomials** requires exponentially fewer points

- For **$2d+1$ points**, this method is equivalent to effective **Unscented Kalman filter**
- **Generalizes to many more points**

can do this even if $p(o_i|x_i)$ is a black box
extended Kalman filter
requires derivative of $p(o_i|x_i)$

19

Operations in non-linear Kalman filter



- Compute $p(X_t | O_{1:t} = o_{1:t})$

- Start with $p(X_0)$

- At each time step t :

- **Condition** on observation (use **numerical integration**)

$$p(X_t | o_{1:t}) \propto p(X_t | o_{1:t-1})p(o_t | X_t)$$

- **Prediction** (Multiply transition model, use **numerical integration**)

$$p(X_{t+1}, X_t | o_{1:t}) = p(X_{t+1} | X_t)p(X_t | o_{1:t})$$

- **Roll-up** (marginalize previous time step)

$$p(X_{t+1} | o_{1:t}) = \int_{X_t} p(X_{t+1}, x_t | o_{1:t}) dx_t$$

20

Canonical form & Markov Nets

$$p(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

$$= K \exp \left\{ \eta^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \Lambda \mathbf{x} \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \sum_{i,j} \lambda_{ij} x_i x_j + \sum_i \eta_i x_i \right\}$$

$$= \exp \left\{ \sum_{i,j} \lambda_{ij} f_{ij}(x_i, x_j) + \sum_i \eta_i f_i(x_i) \right\}$$

edge
node
features
features

MN:

graph
structure

precision matrix $\Lambda \equiv \Sigma^{-1}$
 defines graph structure
 $\lambda_{ij} = 0$ no edge

21

What you need to know about Gaussians, Kalman Filters, Gaussian MNs

■ Kalman filter

- ☐ Probably most used BN
- ☐ Assumes Gaussian distributions
- ☐ Equivalent to linear system
- ☐ Simple matrix operations for computations

■ Non-linear Kalman filter

- ☐ Usually, observation or motion model not CLG
- ☐ Use numerical integration to find Gaussian approximation

■ Gaussian Markov Nets

- ☐ Sparsity in precision matrix equivalent to graph structure

■ Continuous and discrete (hybrid) model

- ☐ Much harder, but doable and interesting (see book)

22