

Probabilistic Graphical Models

Lecture 12 – Dynamical Models

CS/CNS/EE 155
Andreas Krause

Announcements

- Homework 3 out tonight
 - Start early!!
- Project milestones due today
 - Please email to TAs

Parameter learning for log-linear models

- Feature functions $\phi_i(C_i)$ defined over cliques
- Log linear model over undirected graph G
 - Feature functions $\phi_1(C_1), \dots, \phi_k(C_k)$
 - Domains C_i can overlap
- Joint distribution

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp\left(\underbrace{\sum_i w_i^T \phi_i(C_i)}\right)$$

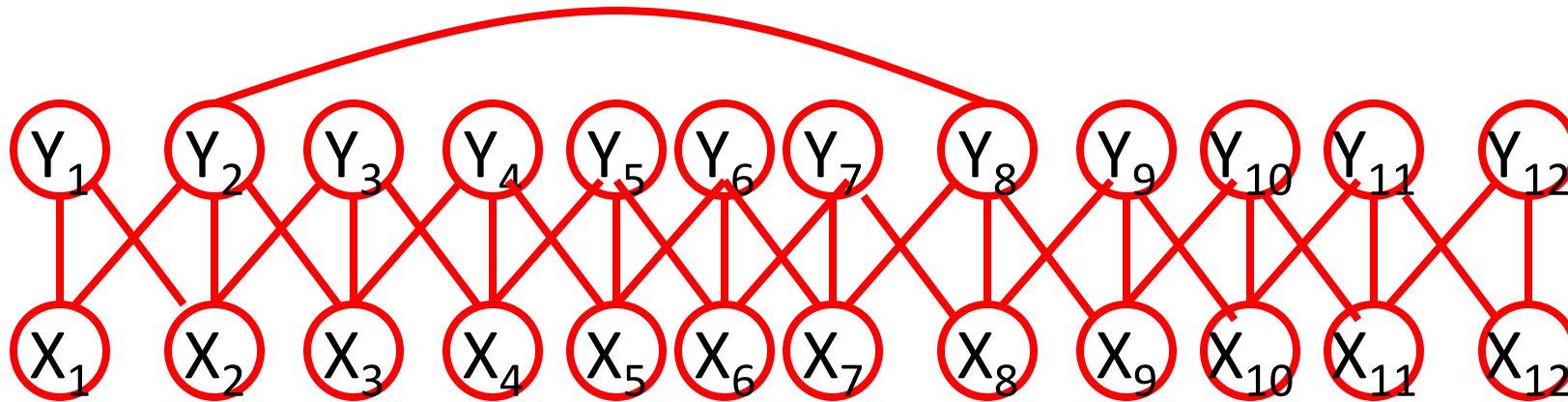
- How do we get weights w_i ?

Log-linear conditional random field

- Define log-linear model over outputs Y
 - No assumptions about inputs X
- Feature functions $\phi_i(C_i, x)$ defined over cliques and inputs
 $C_i \subseteq Y$
- Joint distribution

$$P(Y_1, \dots, Y_n \mid x) = \frac{1}{Z(x)} \exp\left(\sum_i w_i^T \phi_i(C_i, x)\right)$$

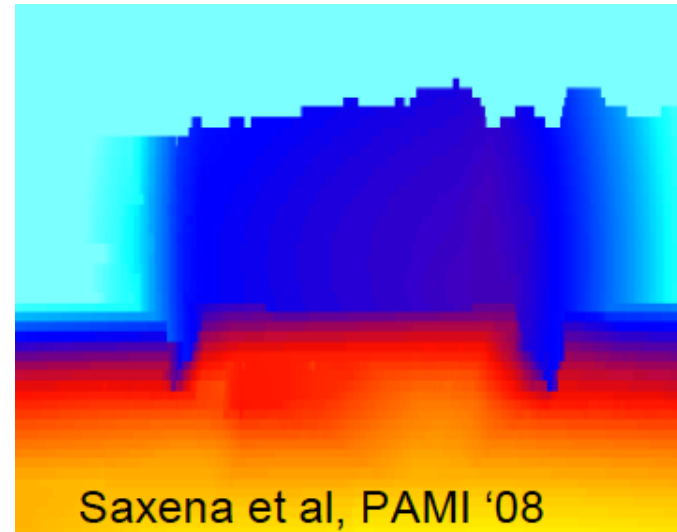
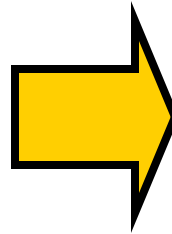
Example: CRFs in NLP



Mrs. Greene spoke today in New York. Green chairs the finance committee

- Classify into Person, Location or Other

Example: CRFs in vision



Gradient of conditional log-likelihood

- Partial derivative

$$\frac{\partial \log P(\mathcal{D}_Y \mid w, \mathcal{D}_X)}{\partial w_i} = \sum_j \left[\phi_i(\mathbf{c}_i^{(j)}, x^{(j)}) + \underbrace{\sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w, \underline{x^{(j)}}) \phi_i(\mathbf{c}_i, \underline{x^{(j)}})}_{\text{Req. Inference}} \right]$$

- Requires one inference per feature and per data point
Can be very expensive
Can do "pseudo" likelihood est. (approximate)
- Can optimize using conjugate gradient

Exponential Family Distributions

- Distributions for log-linear models

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp\left(\sum_i w_i^T \phi_i(C_i)\right)$$

- More generally: **Exponential family distributions**

$$P(x) = h(x) \exp(w^T \phi(x) - A(w))$$

- $h(x)$: Base measure \leftarrow often constant
- w : natural parameters
- $\phi(x)$: Sufficient statistics
- $A(w)$: log-partition function $\leftarrow A(w) = \log \mathcal{Z}(w)$
- Hereby x can be continuous (defined over any set)

Examples

- Exp. Family:

$$P(x) = h(x) \exp(w^T \phi(x) - A(w))$$

$h(x)$: Base measure

w : natural parameters

$\phi(x)$: Sufficient statistics

$A(w)$: log-partition function

- Gaussian distribution

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$h(x) = \frac{1}{\sqrt{2\pi}}$$

$$-\frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}$$

$$\phi(x) = \left[-\frac{x^2}{2}, x\right], \quad A(w) = \frac{\mu^2}{2\sigma^2} - \log \sigma$$

$$w = \left[\frac{1}{\sigma^2}, \frac{\mu}{\sigma^2}\right]$$

- Other examples: Multinomial, Poisson, Exponential, Gamma, Weibull, chi-square, Dirichlet, Geometric, ...

Moments and gradients

$$P(x) = h(x) \exp(w^T \phi(x) - A(w))$$

- Correspondence between moments and log-partition function (just like in log-linear models)

$$\frac{\partial A(w)}{\partial w_i} = \int p(x | w) \phi_i(x) dx = \mathbb{E}[\phi_i | w]$$

$$\frac{\partial^2 A(w)}{\partial w_i \partial w_j} = \text{Cov}(\phi_i, \phi_j | w)$$

- Can compute moments from derivatives, and derivatives from moments!
- MLE \Leftrightarrow moment matching

Conjugate priors in Exponential Family

$$P(x | w) = h(x) \exp(w^T \phi(x) - A(w))$$

Any exponential family likelihood has a conjugate prior

$$P(w|\alpha, \beta) = \exp(\underline{\alpha}^T w - \beta A(w) - \underline{B(\alpha, \beta)})$$

$$P(x|w) P(w|\alpha, \beta) \propto \exp(w^T(\phi(x) + \alpha) - (\beta + 1) A(w))$$

Exponential family graphical models

- So far, only defined graphical models over discrete variables.
- Can define GMs over continuous distributions!
- For exponential family distributions:

$$p(X_1, \dots, X_n) = \prod_i h_i(C_i) \exp\left(\sum_i w_i^T \phi_i(C_i) - A(w)\right)$$
$$\exp A(w) = \int \int \dots \int \prod_i h_i(C_i) \exp\left(\sum_i w_i^T \phi_i(C_i) - A(w)\right) dx_1 \dots dx_n$$

- Can do much of what we discussed (VE, JT, parameter learning, etc.) for such exponential family models
- Important example: **Gaussian Networks**

Multivariate Gaussian distribution

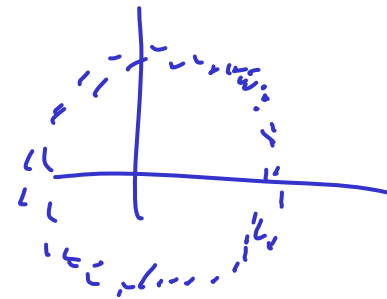
$$\mathcal{N}(x; \Sigma, \mu) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & & & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$

$$\Sigma^T = \Sigma$$

$$\text{pos. def. } x^T \Sigma x \geq 0 \quad \forall x$$

- Joint distribution over n random variables $P(X_1, \dots, X_n)$
- $\sigma_{jk} = E[(X_j - \mu_j)(X_k - \mu_k)] = \text{Cov}(X_j, X_k)$
- X_j and X_k independent $\Leftrightarrow \sigma_{jk} = 0$



Marginalization

- Suppose $(X_1, \dots, X_n) \sim N(\mu, \Sigma)$

- What is $P(X_1)$??

$$P(X_1) = \int \int \int \dots P(X_1, \dots, X_n) dX_2, \dots, dX_n$$

$$P(X_1) = \mathcal{N}(X_1; \mu_1, \sigma_1^2)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & & & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$

Σ_{AA} μ_A

- More generally: Let $A = \{i_1, \dots, i_k\} \subseteq \{1, \dots, N\}$

- Write $X_A = (X_{i_1}, \dots, X_{i_k})$

- X_A $\sim N(\mu_A, \Sigma_{AA})$

$$\Sigma_{AA} = \begin{pmatrix} \sigma_{i_1 i_1}^2 & \sigma_{i_1 i_2} & \dots & \sigma_{i_1 i_k} \\ \vdots & & & \vdots \\ \sigma_{i_k i_1} & \sigma_{i_k i_2} & \dots & \sigma_{i_k i_k}^2 \end{pmatrix} \quad \mu_A = \begin{pmatrix} \mu_{i_1} \\ \mu_{i_2} \\ \vdots \\ \mu_{i_k} \end{pmatrix}$$

Conditioning

- Suppose $(X_1, \dots, X_n) \sim N(\mu, \Sigma)$
- Decompose as (X_A, X_B)
- What is $P(X_A | X_B)$??

$$\Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

- $P(X_A = x_A | X_B = x_B) = N(x_A; \mu_{A|B}, \Sigma_{A|B})$ where

$$\mu_{A|B} = \mu_A + \underbrace{\Sigma_{AB} \Sigma_{BB}^{-1}}_w (x_B - \mu_B)$$

$$\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

- Computable using linear algebra!

Does not depend on x_B

Conditional linear Gaussians

$$\mu_{A|B} = \mu_A + \underbrace{\Sigma_{AB}\Sigma_{BB}^{-1}}_W (x_B - \mu_B)$$

$$\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$$

$$x_{A|B} = \mu_A + Wx_B - W\mu_B + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \Sigma_{A|B})$$
$$= \underline{Wx_B + b} + \epsilon \quad \leftarrow \text{noise}$$

$$W = \Sigma_{AB}\Sigma_{BB}^{-1}$$

$$b = \mu_A - W\mu_B$$

Canonical Representation

$$p(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$
$$\propto \exp(\eta^T x - \frac{1}{2} x^T \underline{\Lambda} x) = \exp \left(\sum_i \eta_i x_i - \frac{1}{2} \sum_{i,j} x_i x_j \Lambda_{i,j} \right)$$

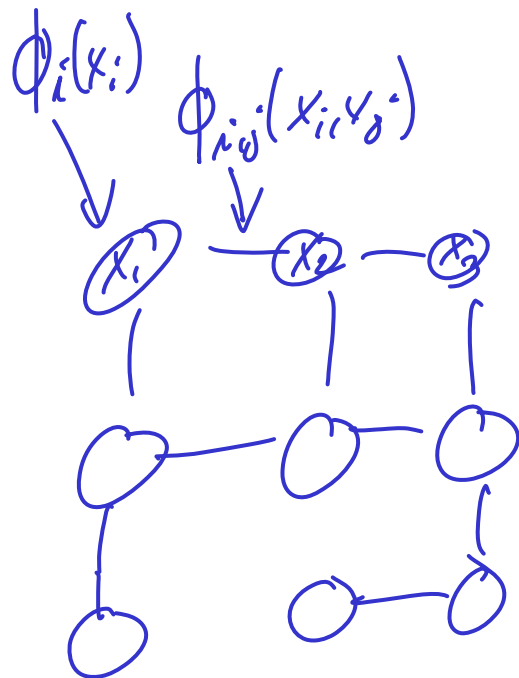
- Multivariate Gaussians in exponential family!
- Standard vs canonical form:

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

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

Gaussian Networks

$$p(X_1, \dots, X_n) \propto \exp\left(-\frac{1}{2} \sum_{i,j} \lambda_{i,j} x_i x_j + \sum_i \eta_i x_i\right)$$



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

$$\phi_{ij}(x_i, x_j) = -\frac{1}{2} x_i x_j \quad \phi_i(x_i) = x_i$$

No edge between X_i, X_j

$$\Leftrightarrow \lambda_{ij} = 0$$

Zeros in precision matrix Λ indicate missing edges in log-linear model!

Inference in Gaussian Networks

- Can compute marginal distributions in $O(n^3)$!
- For large numbers n of variables, still intractable
- If Gaussian Network has low treewidth, can use variable elimination / JT inference!
- Need to be able to multiply and marginalize factors!

$$g = \int_{x_i} \prod_j f_j$$

Multiplying factors in Gaussians

$$P(X_A) = \mathcal{N}(x_A; \Lambda_1, \eta_1) \quad \begin{array}{l} |A| = k, |B| = m \\ \Lambda_1 \in \mathbb{R}^{k \times k} \end{array}$$

$$P(X_B | X_A) \propto \mathcal{N}(x_B, x_A; \Lambda_2, \eta_2) \quad \leftarrow \begin{array}{l} \text{CLG representation} \\ \Lambda_2 \in \mathbb{R}^{m \times m} \end{array}$$

$$P(X_A, X_B) = P(X_A) P(X_B | X_A)$$

$$\propto \exp(x_A^T \Lambda_1 x_A + \eta_1^T x_A) \exp\left((x_A, x_B)^T \Lambda_2 (x_A, x_B) + \eta_2^T \begin{pmatrix} x_A \\ x_B \end{pmatrix}\right)$$

$$= \mathcal{N}(x_A, x_B; \Lambda, \eta)$$

$$\Lambda = \Lambda_1 + \Lambda_2 = \begin{pmatrix} \boxed{\Lambda_1} & 0 \end{pmatrix} + \begin{pmatrix} \Lambda_2 \end{pmatrix}$$

$$\eta = \eta_1 + \eta_2$$

Conditioning in canonical form

- Joint distribution $(X_A, X_B) \sim N(\eta_{AB}, \Lambda_{AB})$

$$\begin{array}{cc} \Lambda_{AA} & \Lambda_{AB} \\ \Lambda_{BA} & \Lambda_{BB} \end{array}$$

- Conditioning: $P(X_A \mid X_B = x_B) = N(x_A; \eta_{A|B=x_B}, \Lambda_{A|B=x_B})$

$$\eta_{A|B=x_B} = \eta_A - \Lambda_{AB}x_B$$

$$\underline{\Lambda_{A|B=x_B}} = \underline{\Lambda_{AA}}$$

Marginalizing in canonical form

- Recall conversion formulas

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

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

- Marginal distribution

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

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

Standard vs. canonical form

Standard form

Canonical form

Marginalization

$$\underline{\mu_A^m = \mu_A}$$

$$\underline{\Sigma_{AA}^m = \Sigma_{AA}}$$

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

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

Conditioning

$$\mu_{A|B=x_B} = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (x_B - \mu_B)$$

$$\Sigma_{A|B=x_B} = \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

$$\eta_{A|B=x_B} = \eta_A - \Lambda_{AB} x_B$$

$$\Lambda_{A|B=x_B} = \underline{\Lambda_{AA}}$$

- In standard form, marginalization is easy
- In canonical form, conditioning is easy!

Variable elimination

$$P(X_1, \dots, X_n) = P(X_1) P(X_2 | X_1) \dots P(X_n | X_{n-1})$$

$\underbrace{\quad}_{\Lambda_1} \quad \underbrace{\quad}_{\Lambda_2} \quad \dots \quad \underbrace{\quad}_{\Lambda_n}$

$$\textcircled{X_1} \rightarrow \textcircled{X_2} \rightarrow \textcircled{X_3} \rightarrow \textcircled{X_4} = \mathcal{N}(\cdot; \Lambda, \eta)$$

$$\Lambda = \begin{pmatrix} \Lambda_1 & & & \\ & \Lambda_2 & & \\ & & \Lambda_3 & \\ & & & \ddots \\ & & & & \Lambda_n \end{pmatrix} \quad \text{block diagonal}$$

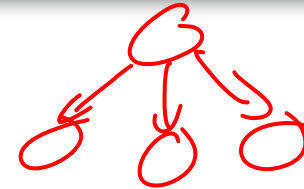
- In Gaussian Markov Networks, Variable elimination = Gaussian elimination (fast for low bandwidth = low treewidth matrices)

Dynamical models

HMMs / Kalman Filters

- Most famous Graphical models:

- Naïve Bayes model
- Hidden Markov model
- Kalman Filter



- Hidden Markov models

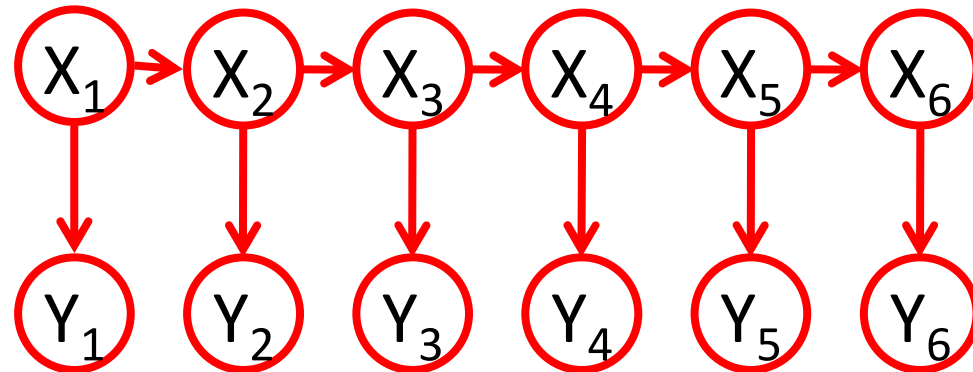
- Speech recognition
- Sequence analysis in comp. bio

- Kalman Filters control

- Cruise control in cars
- GPS navigation devices
- Tracking missiles..

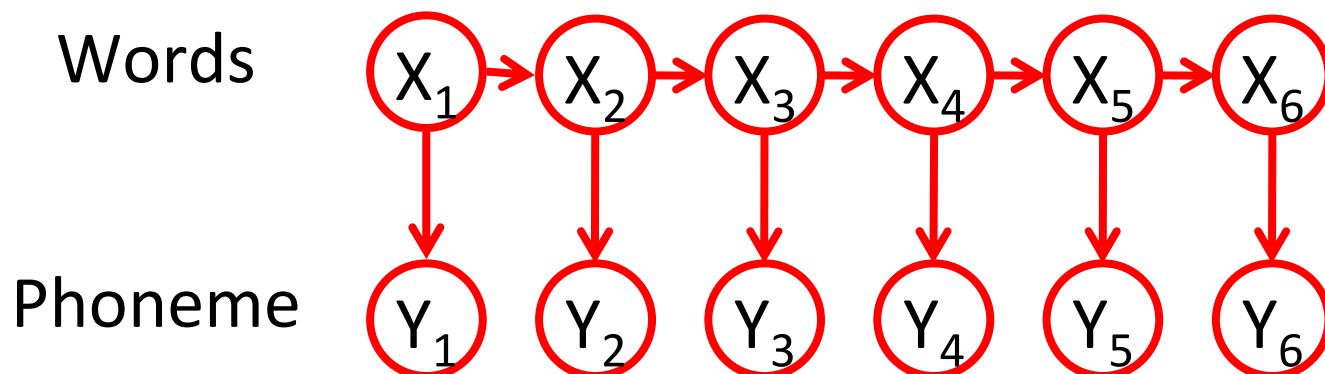
- Very simple models but very powerful!!

HMMs / Kalman Filters



- X_1, \dots, X_T : Unobserved (hidden) variables
- Y_1, \dots, Y_T : Observations
- HMMs: X_i Multinomial, Y_i arbitrary
- Kalman Filters: X_i, Y_i Gaussian distributions
 - Non-linear KF: X_i Gaussian, Y_i arbitrary

HMMs for speech recognition



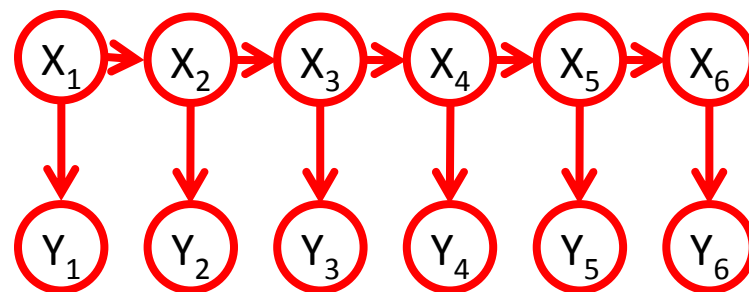
“He ate the cookies on the couch”

- Infer spoken words from audio signals

Hidden Markov Models

- Inference:

- In principle, can use VE, JT etc.
- New variables X_t, Y_t at each time step \rightarrow need to rerun

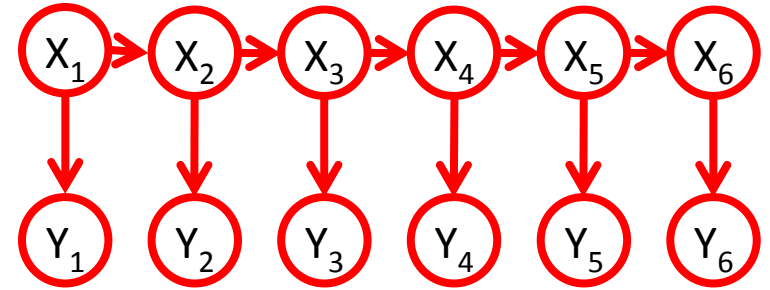


- Bayesian Filtering:

- Suppose we already have computed $P(X_t \mid y_{1,\dots,t})$
- Want to efficiently compute $P(X_{t+1} \mid y_{1,\dots,t+1})$

Bayesian filtering

- Start with $P(X_1)$
- At time t
 - Assume we have $P(X_t | y_{1..t-1})$
 - Condition: $P(X_t | y_{1..t})$



$$P(X_t | y_{1..t}) \propto P(X_t | y_{1..t-1}) \underbrace{P(Y_t | X_t, y_{1..t-1})}_{\text{cond. ind. } P(Y_t | X_t)}$$

- Prediction: $P(X_{t+1}, X_t | y_{1..t})$

$$P(X_{t+1}, X_t | y_{1..t}) = P(X_t | y_{1..t}) \cdot \underbrace{P(X_{t+1} | X_t, y_{1..t})}_{= P(X_{t+1} | X_t)}$$

"Roll up"

- Marginalization: $P(X_{t+1} | y_{1..t})$

$$P(X_{t+1} | y_{1..t}) = \sum_{X_t} P(X_{t+1}, X_t | y_{1..t})$$

Parameter learning in HMMs

- Assume we have labels for hidden variables
- Assume stationarity
 - $P(X_{t+1} | X_t)$ is same over all time steps
 - $P(Y_t | X_t)$ is same over all time steps
 - Violates parameter independence (\rightarrow parameter “sharing”)

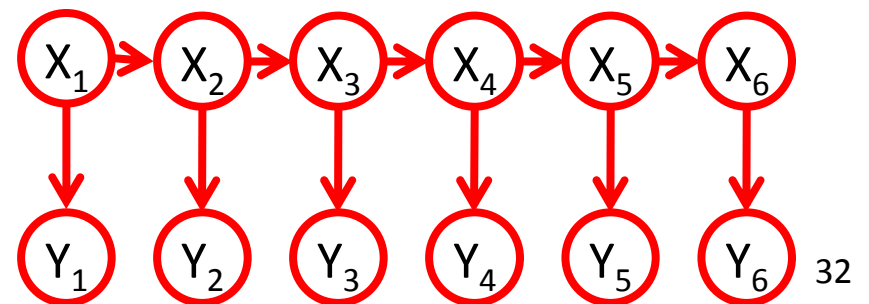
- Example: compute parameters for $P(X_{t+1}=x | X_t=x')$

$$\begin{aligned}
 \log P(x_1, \dots, x_T, y_1, \dots, y_T | \theta) &= \log P(x_1) \prod_{t=2}^T P(x_t | x_{t-1}, \theta_1) \prod_{t=1}^T P(y_t | x_t, \theta_2) \\
 &= \log P(x_1 | \theta_3) + \sum_t \log P(x_t | x_{t-1}, \theta_1) + \sum_t \log P(y_t | x_t, \theta_2) \\
 \theta_{x_t=x | x_{t-1}=x'} &= \frac{\text{count}(x', x)}{T-1}
 \end{aligned}$$

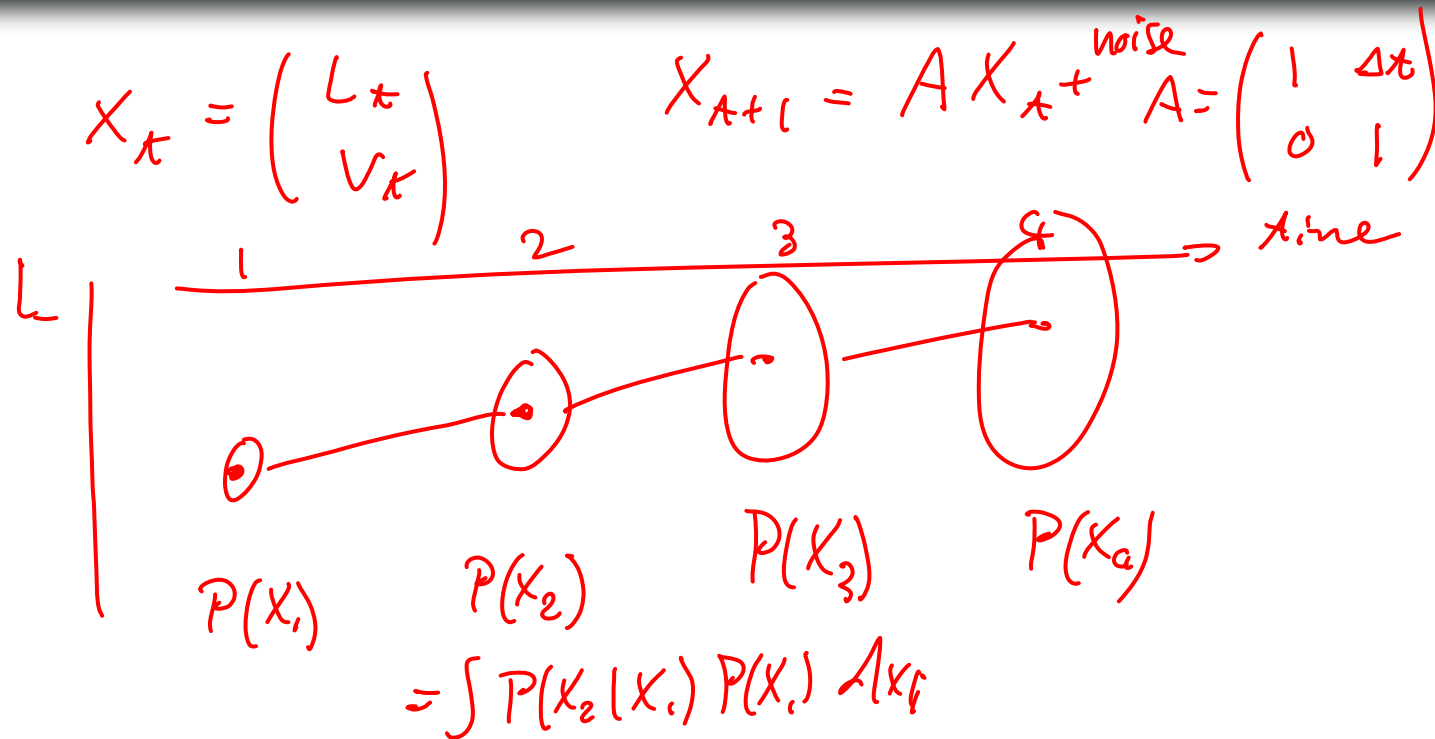
- What if we don't have labels for hidden vars?
 - \rightarrow Use EM (later this course)

Kalman Filters (Gaussian HMMs)

- X_1, \dots, X_T : Location of object being tracked
- Y_1, \dots, Y_T : Observations
- $P(X_1)$: Prior belief about location at time 1
- $P(X_{t+1} | X_t)$: “Motion model”
 - How do I expect my target to move in the environment?
 - Represented as CLG: $X_{t+1} = A X_t + N(0, \Sigma_M)$
- $P(Y_t | X_t)$: “Sensor model”
 - What do I observe if target is at location X_t ?
 - Represented as CLG: $Y_t = H X_t + N(0, \Sigma_O)$



Understanding Motion model

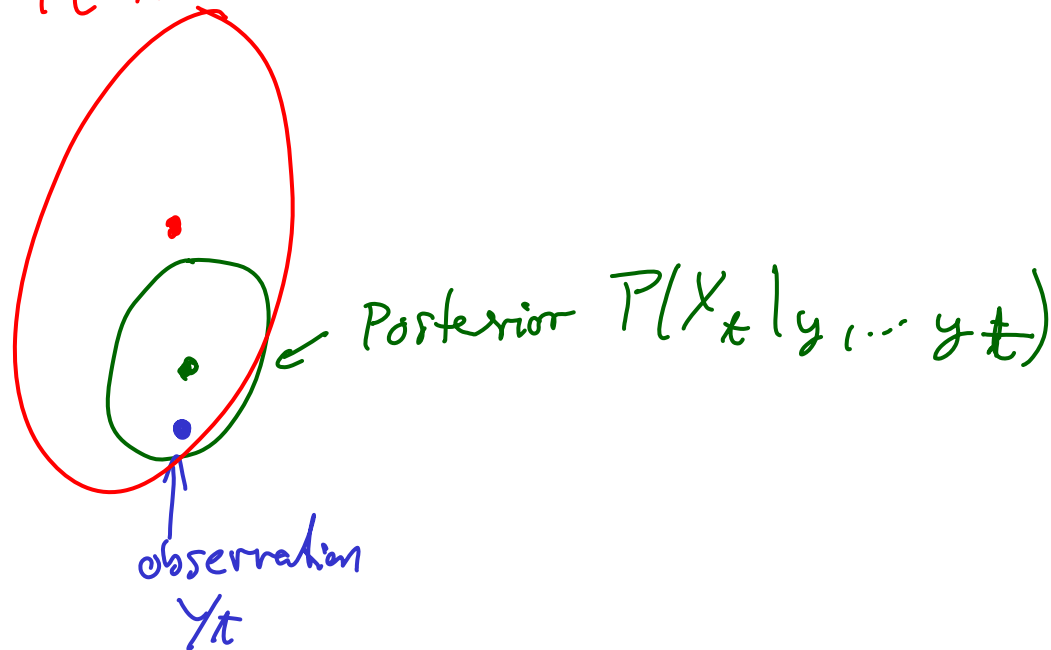


Understanding sensor model

$$X_t = \begin{pmatrix} L_t \\ V_t \end{pmatrix}$$

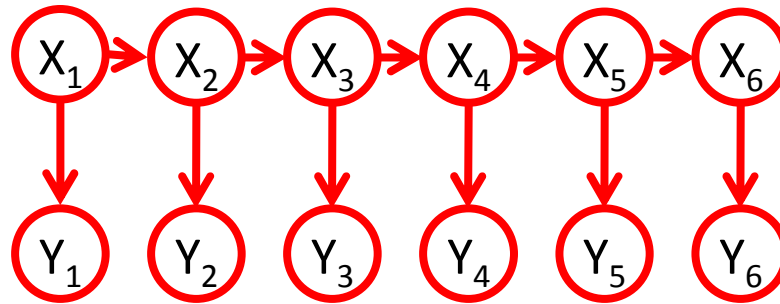
$$y_t = \underbrace{(1, 0)}_{\text{only obs. location}} X_t + \text{noise}$$

$$P(X_t | y_1, \dots, y_{t-1})$$



Bayesian Filtering for KFs

- Can use Gaussian elimination to perform inference in “unrolled” model



- Start with prior belief $P(X_1)$
- At every timestep have belief $P(X_t \mid y_{1:t-1})$
 - Condition on observation: $P(X_t \mid y_{1:t})$ $\leftarrow \mathcal{N}(x_t; \hat{x}_t, \hat{\Sigma}_t)$ "sensor model"
 - Predict (multiply motion model): $P(X_{t+1}, X_t \mid y_{1:t})$ \leftarrow Multiply Likelihood
 - "Roll-up" (marginalize prev. time): $P(X_{t+1} \mid y_{1:t})$ \leftarrow Multiply motion model

Implementation

- Current belief: $P(x_t | y_{1:t-1}) = \mathcal{N}(x_t; \eta_{x_t}, \Lambda_{x_t})$

- Multiply sensor and motion model

Motion model
 $P(x_{t+1} | x_t) \leftarrow$ Represented as CLG, canonical params Λ_M, η_M

$$\begin{aligned} P(x_{t+1}, x_t | y_{1:t}) &= P(x_t | y_{1:t-1}) P(x_{t+1} | x_t) \\ &= \mathcal{N}(\cdot; \eta_{x_t} + \eta_M, \Lambda_{x_t} + \Lambda_M) \end{aligned}$$

- Marginalize $P(x_{t+1} | y_{1:t}) = \mathcal{N}(\cdot; \eta_A^m, \Lambda_{AA}^m)$

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

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

$$A = x_{t+1} | y_{1:t-1}$$

$$B = x_t | y_{1:t-1}$$

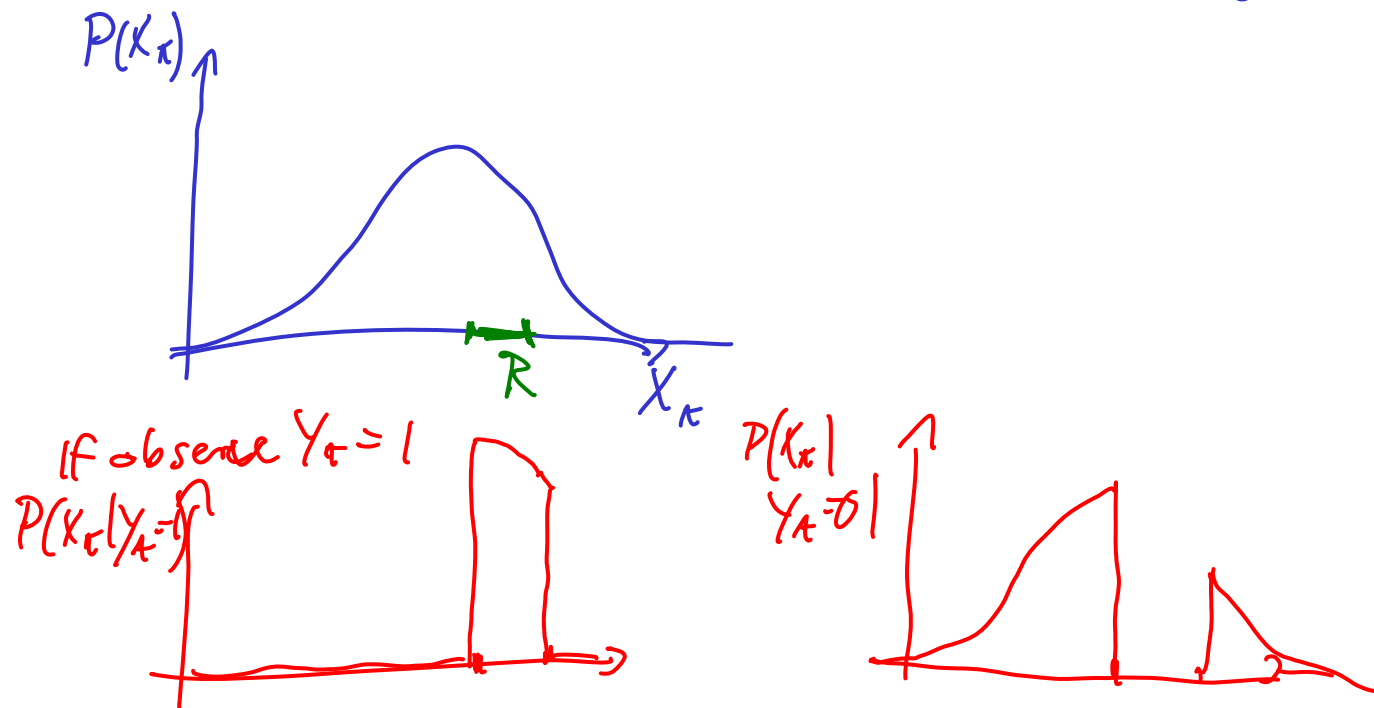
What if observations not “linear”?

- Linear observations:

- $Y_t = H X_t + \text{noise}$

- Nonlinear observations:

“Motion detector” : $Y_t = 1$ if $X_t \in \mathcal{R}$
 $= 0$ otherwise



Incorporating Non-gaussian observations

- Nonlinear observation $\Rightarrow P(Y_t | X_t)$ not Gaussian
- First approach: Approximate $P(Y_t | X_t)$ as CLG
 - Linearize $P(Y_t | X_t)$ around current estimate $E[X_t | y_{1..t-1}]$
 - Known as Extended Kalman Filter (EKF)
 - Can perform poorly if $P(Y_t | X_t)$ highly nonlinear
- Second approach: Approximate $P(Y_t, X_t)$ as Gaussian
 - Takes correlation in X_t into account
 - After obtaining approximation, condition on $Y_t=y_t$ (now a “linear” observation)

Finding Gaussian approximations

- Need to find Gaussian approximation of $P(X_t, Y_t)$
- How?
 - Gaussians in Exponential Family \rightarrow Moment matching!!

- $E[Y_t] = \int y P(y|x) dy = \int y \underbrace{P(y|x)}_{\text{nonlinear}} \underbrace{P(x)}_{\text{Gaussian}} dx dy$

- $E[Y_t^2] = \int y^2 P(y|x) P(x) dx dy$

- $E[X_t Y_t] = \int x y P(y|x) P(x) dx dy$

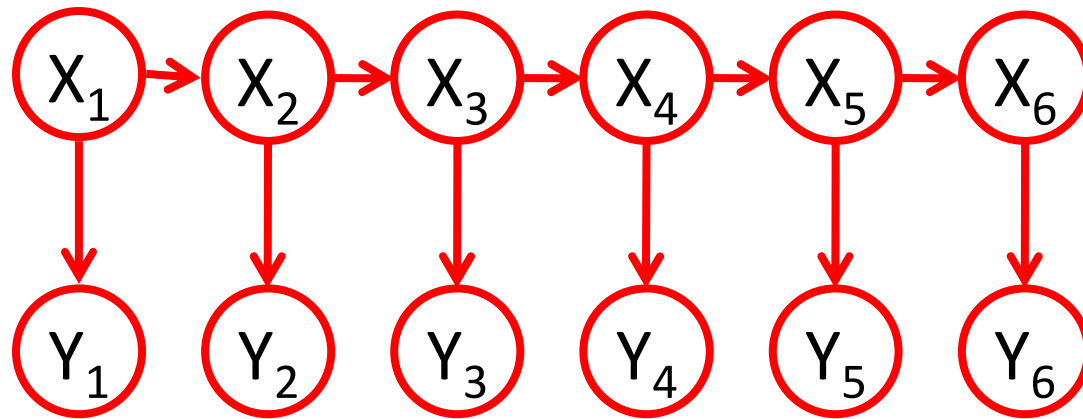
Linearization by integration

- Need to integrate product of Gaussian with arbitrary function
- Can do that by numerical integration
 - Approximate integral as weighted sum of evaluation points
$$\int f(x,y) p(x) dx dy \approx \sum_i w_i f(x^{(i)}, y^{(i)})$$

\uparrow How should we choose?
 - Gaussian quadrature defines locations and weights of points
 - For 1 dim: **Exact** for polynomials of degree D if choosing 2D points using Gaussian quadrature
 - For higher dimensions: Need exponentially many points to achieve exact evaluation for polynomials
- Application of this is known as “Unscented” Kalman Filter (UKF)

Factored dynamical models

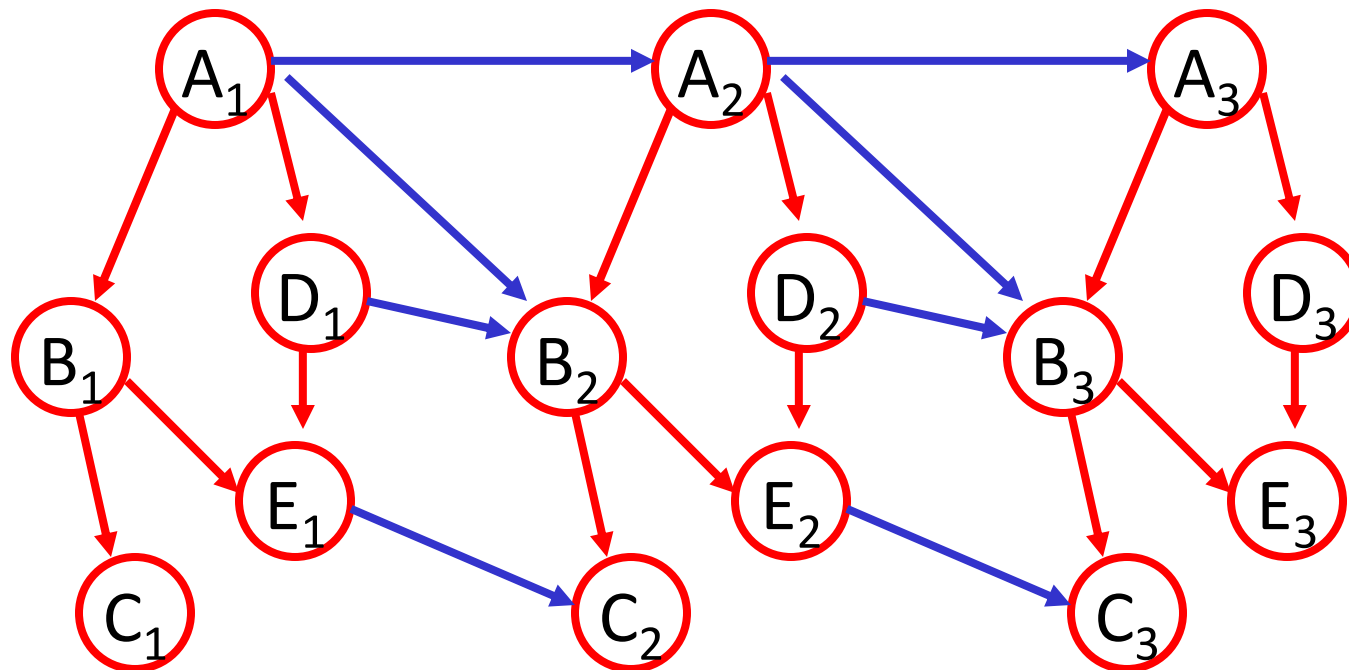
- So far: HMMs and Kalman filters



- What if we have more than one variable at each time step?
 - E.g., temperature at different locations, or road conditions in a road network?
 - ➔ Spatio-temporal models

Dynamic Bayesian Networks

- At every timestep have a Bayesian Network



- Variables at each time step t called a “slice” S_t
- “Temporal” edges connecting S_{t+1} with S_t

Tasks

- Read Koller & Friedman Chapters 6.2.3, 15.1