# EM for BNs

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 24th, 2008

1

---

# Thus far, fully supervised learning

- We have assumed fully supervised learning:

$$\langle A{=}t, F{=}f, S{=}t, H{=}t, N{=}f \rangle$$

$$\vdots$$

- Many real problems have missing data:

$$\langle A{=}t, F{=}f, S{=}?, H{=}t, N{=}? \rangle$$

2

---

1

## The general learning problem with missing data

- Marginal likelihood – **x** is observed, **z** is missing:

$$
\begin{aligned}
\ell(\mathcal{D}:\theta) &= \log \prod_{j=1}^{m} P(x^{(j)} \mid \theta) \\
&= \sum_{j=1}^{m} \log P(x^{(j)} \mid \theta) \\
&= \sum_{j=1}^{m} \log \sum_{z} P(z, x^{(j)} \mid \theta)
\end{aligned}
$$

# E-step

- **x** is observed, **z** is missing
- Compute probability of missing data given current choice of $\theta$
  - Q(**z**|**x**$^{(j)}$) for each **x**$^{(j)}$
    - e.g., probability computed during classification step
    - corresponds to "classification step" in K-means

$$
Q^{(t+1)}(z \mid x^{(j)}) = P(z \mid x^{(j)}, \theta^{(t)})
$$

inference in a BN

2

# Jensen's inequality

$$\ell(\mathcal{D} : \theta) = \sum_{j=1}^{m} \log \sum_{z} P(z, x^{(j)} \mid \theta)$$

- **Theorem**: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\log 2 = \log(1+1) \geq \log 1 + \log 1 = 0$$

# Applying Jensen's inequality

EM: optimizes lower bound on $\ell(D:\theta)$

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\mathcal{D} : \theta^{(t)}) = \sum_{j=1}^{m} \log \sum_{z} Q^{(t+1)}(z \mid x^{(j)}) \frac{P(z, x^{(j)} \mid \theta^{(t)})}{Q^{(t+1)}(z \mid x^{(j)})}$$

$"P(z)"$        $"f(z)"$

Vald

$$\ell(D:\theta^{(t)}) \geq \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z|x^{(j)}) \log \frac{P(z, x^{(j)}|\theta^{(t)})}{Q^{(t+1)}(z|x^{(j)})}$$

wrt $\theta, Q$

$$= \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z|x^{(j)}) \log P(z, x^{(j)}|\theta^{(t)}) - \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z|x^{(j)}) \log Q^{(t+1)}(z|x^{(j)})$$

weighted log-likelihood of fully observed data

$m \hat{H}_{Q^{(t+1)}}(Z|X)$

3

# The M-step maximizes lower bound on weighted data

*Constant*

- Lower bound from Jensen's:

$$\ell(\mathcal{D} : \theta^{(t)}) \geq \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z \mid x^{(j)}) \log P(z, x^{(j)} \mid \theta^{(t)}) + \overset{\frown}{H}(Q^{(t+1)})$$

*fix Q :*
*optimize θ*

*weighted LL completed data*

- Corresponds to weighted dataset:
  - $<\underline{\mathbf{x}^{(1)}, \mathbf{z}=1}>$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}^{(1)})$  *.8*
  - $<\mathbf{x}^{(1)}, \mathbf{z}=2>$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}^{(1)})$  *.15*
  - $<\mathbf{x}^{(1)}, \mathbf{z}=3>$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}^{(1)})$  *.05*
  - $<\underline{\mathbf{x}^{(2)}, \mathbf{z}=1}>$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}^{(2)})$  *.1*
  - $<\mathbf{x}^{(2)}, \mathbf{z}=2>$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}^{(2)})$  *.8*
  - $<\mathbf{x}^{(2)}, \mathbf{z}=3>$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}^{(2)})$  *.1*
  - …

# The M-step

*discovering hidden variables?*
*→ see book*

$$\ell(\mathcal{D} : \theta^{(t)}) \geq \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z \mid x^{(j)}) \log P(z, x^{(j)} \mid \theta^{(t)}) + H(Q^{(t+1)})$$

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{z} Q^{(t+1)}(z \mid x^{(j)}) \log P(z, x^{(j)} \mid \theta)$$

*$E_{Q^{(t+1)}}\left[ count(X=t, Z=f) \right] = \sum_{j=1}^{m} \mathbb{1}(x^{(j)}=t) \, Q^{(t+1)}(z=f|x^{(j)})$*

- Use expected counts instead of counts:
  - If learning requires Count(**x**,**z**)
  - Use $E_{Q(t+1)}$[Count(**x**,**z**)]

# Convergence of EM

- Define potential function F($\theta$,Q):

$$\ell(\mathcal{D} : \theta^{(t)}) \geq F(\theta, Q) = \sum_{j=1}^{m} \sum_{z} Q(z \mid x^{(j)}) \log \frac{P(z, x^{(j)} \mid \theta)}{Q(z \mid x^{(j)})}$$

*Initialization can matter a lot...*

- EM corresponds to coordinate ascent on F
  - □ Thus, maximizes lower bound on marginal log likelihood
  - □ As seen in Machine Learning class last semester

*fixing Q → max over $\theta$*
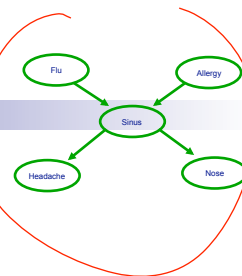*fixing $\theta$ → max over Q*

# Data likelihood for BNs

- Given structure, log likelihood of fully observed data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$

$$= \sum_{j=1}^{m} \log P(f^{(j)}) P(a^{(j)}) P(s^{(j)} \mid f^{(j)}, a^{(j)}) \ldots$$
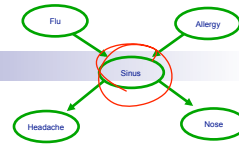
↑ *decomposes*
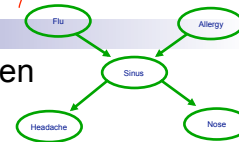
5

# Marginal likelihood



■ What if S is hidden?

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$

$$= \sum_{j=1}^{m} \log \sum_{S} P(f^{(j)}) \, P(a^{(j)}) \, P(S \mid a^{(j)}, f^{(j)}) \, P(h^{(j)} \mid S) \, P(n^{(j)} \mid S)$$

no longer decomposes!!

# Log likelihood for BNs with hidden data

$C, \{x_1, \ldots x_n\}$    $H = \{x_1, \ldots x_n\}/O$



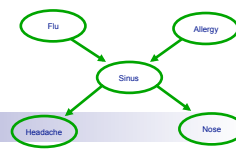■ Marginal likelihood – **O** is observed, **H** is hidden

$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^{m} \log P(\mathbf{o}^{(j)} \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{h}} P(\mathbf{h}, \mathbf{o}^{(j)} \mid \theta)$$

What if I just want to compute this?

else, e.g., Variable elimination

6

# E-step for BNs

- E-step computes probability of hidden vars **h** given **o**

$$Q^{(t+1)}(\mathbf{h} \mid \mathbf{o}^{(j)}) \leftarrow P(\mathbf{h} \mid \mathbf{o}^{(j)}, \theta^{(t)})$$

compute with:
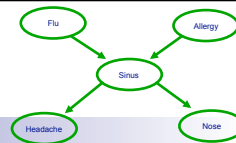
- Corresponds to inference in BN

Naively, must represent joint over H|O$^{(j)}$
if there are K hidden vars, then joint is huge!!

# The M-step for BNs

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{h}} Q^{(t+1)}(\mathbf{h} \mid \mathbf{o}^{(j)}) \log P(\mathbf{h}, \mathbf{o}^{(j)} \mid \theta)$$

- Use expected counts instead of counts:
  - □ If learning requires Count(**h**,**o**)
  - □ Use $E_{Q(t+1)}$[Count(**h**,**o**)]

# M-step for each CPT



- **M-step decomposes per CPT**
  - □ Standard MLE:
  $$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\text{Count}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\text{Count}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

  - □ M-step uses expected counts:
  $$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\text{ExCount}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\text{ExCount}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

if var & parents hidden

$\text{ExCount}[X_i = x_i, Pa_{X_i} = z]$
$= \sum_{j=1}^{m} Q^{(t+1)}(x_i, z \mid O^{(j)})$

hidden

e.g., $P(H \mid S = f) \overset{MLE}{=} \dfrac{\text{Count}(H = t, S = f)}{\text{Count}(S = f)}$   S is hidden

$\text{ExCount}(H = t, S = f) = \sum_{j=1}^{m} \mathbb{I}(H^{(j)} = t) \, Q^{(t+1)}(S = f \mid h^{(j)}, f^{(j)}, a^{(j)}, n^{(j)})$

---

# Computing expected counts



$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\text{ExCount}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\text{ExCount}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

- **M-step requires expected counts:**
  - □ Observe **O=o**
  - □ For a set of vars **A**, must compute ExCount(**A=a**)
  - □ Some of **A** in example *j* will be observed
    - denote by $\mathbf{A_O} = \mathbf{a_O}^{(j)}$    $A_O \subseteq O$
  - □ Some of **A** will be hidden
    - denote by $\mathbf{A_H}$    $A_H \cap O = \emptyset$
- **Use inference (E-step computes expected counts):**
  - □ ExCount$^{(t+1)}$(**A$_O$ = a$_O$, A$_H$ = a$_H$**) $= \sum_{j=1}^{m} \mathbb{I}(A_O^{(j)} = a_O) \, Q^{(t+1)}(A_H = a_H \mid O^{(j)})$

requires inference
for each
datapoint

all observed

if $A_O = \emptyset$,   $= \sum_{j=1}^{m} Q^{(t+1)}(A = a \mid O^{(j)})$

# Data need not be hidden in the same way

Flu → Sinus ← Allergy; Sinus → Headache, Sinus → Nose

- When data is fully observed
  - A data point is $\langle F=t, A=f, S=t \rangle$

  *ugly notation*
  *same everything*

- When data is partially observed $\langle F=t, A=f, S=f \rangle$
  - A data point is $\langle F=t, A=f, S=? \rangle$
  $\langle F=f, A=f, S=? \rangle$

- But unobserved variables can be different for different data points
  - e.g., $\langle F=t, A=t, S=? \rangle$ $\langle F=?, A=?, S=t \rangle$
  $\langle F=t, A=?, S=? \rangle$

- Same framework, just change definition of expected counts
  - Observed vars in point $j$, $O_j$ ← changes for each $j$
  - Consider set of vars **A** $\quad A_{O_j} = A \wedge O_j \qquad A_{H_j} = A / A_{O_j}$
  - ExCount$^{(t+1)}$(**A = a**) $= \sum_{j=1}^{m} \mathbb{1}(A_{O_j} = a_{O_j}) \, Q^{(t+1)}(A_{H_j} = a_{H_j} \mid O_j^{(j)})$

---

# Poster printing → *facilities usually don't print on weekends*

- Poster session:
  - ~~Friday~~ *Monday* Dec 1st, 3-6pm in the NSH Atrium.
  - There will be a popular vote for best poster. Invite your friends!
  - please be ready to set up your poster at 2:45pm sharp.
- We will provide posterboards, easels and pins.
  - The posterboards are 30x40 inches
  - We don't have a specific poster format for you to use.
    - You can either bring a big poster or a print a set of regular sized pages and pin them together.
- Unfortunately, we don't have a budget to pay for printing. If you are an SCS student, SCS has a poster printer you can use which prints on a 36" wide roll of paper.
  - If you are a student outside SCS, you will need to check with your department to see if there are printing facilities for big posters (I don't know what is offered outside SCS),  or print a set of regular sized pages.
- We are looking forward to a great poster session!

9

# EM for BNs & identifiability: a superficial discussion

- **What happens if a leaf is never observed?**

$X$

*y is observed but*
*X is hidden*
*e.g., clustering*
*"can learn"*

$$\text{Ecount}(X=x, Y=y)$$
$$= \sum_{j=1}^{m} \mathbb{I}(Y^{(j)}=y) \, Q(X=x \mid Y=y^{(j)})$$

*Y is hidden but X is observed?*

$Y$

$$\text{Ex Count}(X=x, Y=y)$$
$$= \sum_{j=1}^{m} \mathbb{I}(X^{(j)}=x) \, Q(Y=y \mid X=x^{(j)})$$

*Can't learn!!*

$$\ell(D; \theta) = \sum_{j=1}^{m} \log \sum_{y} P(y, x^{(j)} \mid \theta)$$
$$= \sum_{j=1}^{m} \log \sum_{y} P(y \mid x^{(j)}, \theta) \, P(x^{(j)} \mid \theta)$$

*independently of $\theta$*
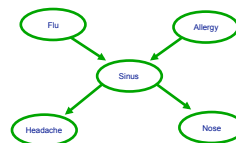
---

# Learning structure with missing data
[K&F 18.4]

- Known BN structure: Use expected counts, learning algorithm doesn't change
- Unknown BN structure:
  - □ Can use expected counts and score model as when we talked about structure learning
  - □ But, very slow...
    - e.g., greedy algorithm would need to redo inference for every edge we test…
- (Much Faster) **Structure-EM**: Iterate:
  - □ compute expected counts
  - □ do a some structure search (e.g., many greedy steps)
  - □ repeat
- **Theorem**: Converges to local optima of marginal log-likelihood
  - □ details in the book

Flu    Allergy

Sinus

Headache    Nose

# What you need to know about learning BNs with missing data

- EM for Bayes Nets
- E-step: inference computes expected counts
  - Only need expected counts over $X_i$ and $\mathbf{Pa}_{xi}$
- M-step: expected counts used to estimate parameters
- Which variables are hidden can change per datapoint
  - Also, use labeled and unlabeled data $\rightarrow$ some data points are complete, some include hidden variables
- Structure-EM:
  - iterate between computing expected counts & many structure search steps

# MNs & CRFs with missing data

- MNs with missing data
  - Models P(**X**), part of **X** hidden
  - Use EM to optimize
  - Same ideas as BN

- CRFs with missing data
  - Models P(**Y**|**X**)
  - What's hidden?
    - Part of **Y**:  $Y_O , Y_H$    $\ell(D:\theta) = \sum_{j=1}^{m} \log \sum_{y_H} P(y_O^{(i)}, y_H | x^{(i)}, \theta)$
    
      can use EM
    - All of **Y**:  $\ell(D:\theta) = \sum_{j=1}^{m} \log \sum_{y} P(y | x^{(i)}, \theta)$   $\Rightarrow$  Can't use EM
    
      $1$
    - Part of **X**:  Can't use EM, because  need to average over  hidden $X$, but  have no model of $P(X)$

11

# Kalman Filters
# Gaussian BNs

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 24th, 2008

23

---

# Adventures of our BN hero

- Compact representation for probability distributions    1. Naïve Bayes
- Fast inference
- Fast learning
- Approximate inference

2 and 3.
Hidden Markov models (HMMs)
Kalman Filters

- But… Who are the most popular kids?

*is an HMM with Gaussian "CPTs"*
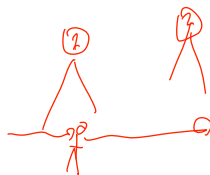
24

# The Kalman Filter

- An HMM with Gaussian distributions
- Has been around for at least 60 years
- Possibly the most used graphical model ever
- It's what
  - does your cruise control
  - tracks missiles
  - controls robots
  - …
- And it's so simple…
  - Possibly explaining why it's so used
- Many interesting models build on it…
  - An example of a Gaussian BN (more on this later)

25

# Example of KF – SLAT
## Simultaneous Localization and Tracking

[Funiak, Guestrin, Paskin, Sukthankar '06]

- Place some cameras around an environment, don't know where they are
- Could measure all locations, but requires lots of grad. student (Stano) time
- Intuition:
  - A person walks around
  - If camera 1 sees person, then camera 2 sees person, learn about relative positions of cameras
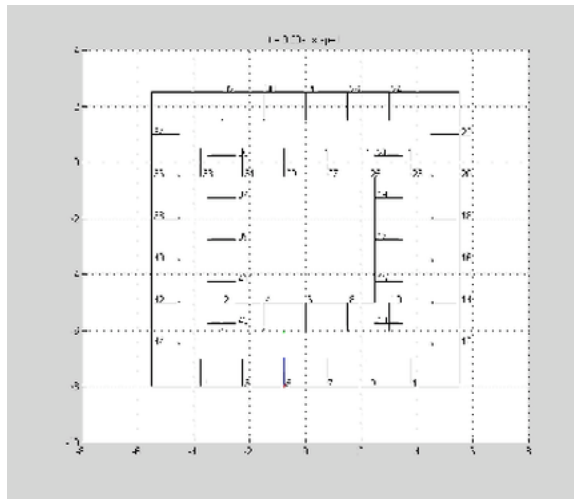
26

13

# Example of KF – SLAT
## Simultaneous Localization and Tracking

[Funiak, Guestrin, Paskin, Sukthankar '06]



27

# Multivariate Gaussian

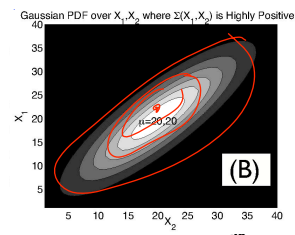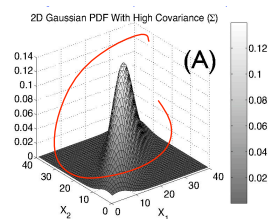$$p(X_1, \ldots, X_n) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

Mean vector:
$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$$

Covariance matrix:
$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 \end{pmatrix}$$

$$\sigma_{32} = \sigma_{23}$$



2D Gaussian PDF With High Covariance (Σ)

(A)

Gaussian PDF over X₁,X₂ where Σ(X₁,X₂) is Highly Positive

μ=20,20

(B)

14

# Conditioning a Gaussian

- **Joint Gaussian:**
  - ☐ $p(X,Y) \sim N(\mu; \Sigma)$
- **Conditional linear Gaussian:**
  - ☐ $p(Y|X) \sim N(\mu_{Y|X}; \sigma^2_{Y|X})$ ← *gaussian*

*observed value of X*
*prior mean of X*

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

↑ *prior*

$$\sigma^2_{Y|X} = \sigma^2_Y - \frac{\sigma^2_{YX}}{\sigma^2_X}$$

*posterior variance*  ↑ *prior variance*  *≥0*

← *Posterior variance doesn't depend on observed value!!*

$P(Y|X=0)$

$\sigma^2_{Y|X} \leq \sigma^2_Y$  $(\sigma^2_{Y|X} = \sigma^2_Y$ iff $Y \perp X)$

*observations always decrease variance*