# Param. Learning (MLE)

# Structure Learning *for BN*

## The Good

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 1st, 2008

---

# Learning the CPTs



Data

$\mathbf{X}^{(1)}$

...

$\mathbf{X}^{(m)}$

For each discrete variable $X_i$    $P_{aX_i} = U$

$$P(X_i \mid P_{aX_i}) = P(X_i \mid U)$$

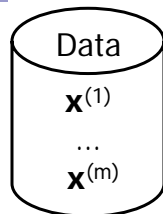$$\hat{P}_{MLE}(x_i \mid U) = \frac{Count(X_i = x_i, U = u)}{Count(U = u)}$$

Why??

MLE:    $P(X_i = x_i \mid X_j = x_j) = \dfrac{Count(X_i = x_i, X_j = x_j)}{Count(X_j = x_j)}$
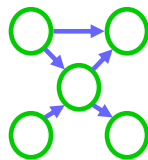
# Learning the CPTs

Data
$\mathbf{x}^{(1)}$
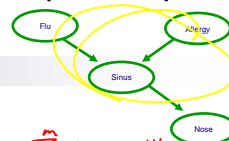…
$\mathbf{x}^{(m)}$

For each discrete variable $X_i$

MLE: $P(X_i = x_i \mid X_j = x_j) = \dfrac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$

**WHY??????????**

if only one var

then take derivative, set to $\phi$

all is good

# Maximum likelihood estimation (MLE) of BN parameters – example

$\log a \cdot b = \log a + \log b$

- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \log \prod_{j=1}^{m} P(x^{(j)} \mid \theta_G, G) = \sum_{j=1}^{m} \log P(x^{(j)} \mid \theta_G, G)$$

for the example

$\sum_{j=1}^{m} \log P(f^{(j)}, a^{(j)}, s^{(j)}, n^{(j)} \mid \theta_G, G) = \sum_{j=1}^{m} \log P(f^{(j)} \mid \theta_G, G) \cdot P(a^{(j)} \mid \theta_G, G) \, P(s^{(j)} \mid a^{(j)}, f^{(j)}, \theta_G, G)$

$P(n^{(j)} \mid s^{(j)}, \theta_G, G)$

$= \sum_{j=1}^{m} \left[ \log P(f^{(j)} \mid \theta_G, G) + \log P(a^{(j)} \mid \theta_G, G) + \log P(s^{(j)} \mid a^{(j)}, f^{(j)}, \theta_G, G) + \log P(n^{(j)} \mid s^{(j)}, \theta_G, G) \right.$

$= \sum_{j=1}^{m} \log P(f^{(j)} \mid \theta_F, G) + \sum_{j=1}^{m} \log P(a^{(j)} \mid \theta_A, G) + \sum_{j=1}^{m} \log P(s^{(j)} \mid a^{(j)}, f^{(j)}, \theta_{S|FA}, G) + \sum_{j=1}^{m} \log P(n^{(j)} \mid s^{(j)}, \theta_{N|S}, G)$

$P(F)$    $P(A)$    $P(S|FA)$    $P(N/S)$

one for each CPT

Broke up problem into independent subproblems:

# Maximum likelihood estimation (MLE) of BN parameters – General case

- Data: $\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(m)}$  $|\mathcal{D}|=m$
- Restriction: $\mathbf{x}^{(j)}[\mathbf{Pa}_{Xi}] \to$ assignment to $\mathbf{Pa}_{Xi}$ in $\mathbf{x}^{(j)}$
- Given structure, log likelihood of data:

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

$$= \max_{\theta_{X_1|Pa_{X_1}}} \max_{\theta_{X_2|Pa_{X_2}}} \ldots \max_{\theta_{X_n|Pa_{X_n}}} \sum_{i=1}^{n} \left[ \sum_{j=1}^{m} \log P\left(x_i^{(j)} \mid x^{(j)}[Pa_{X_i}], \theta_{X_i|Pa_{X_i}}\right) \right]$$

$$= \sum_{i=1}^{n} \left[ \max_{\theta_{X_i|Pa_{X_i}}} \sum_{j=1}^{m} \log P\left(x_i^{(j)} \mid x^{(j)}[Pa_{X_i}], \theta_{X_i|Pa_{X_i}}\right) \right]$$

indep. Max prob.

Sol. MLE: $\hat{P}(X_i = x_i \mid U = u) = \dfrac{Count(x_i, u)}{Count(u)}$

# Taking derivatives of MLE of BN parameters – General case

param sharing
decomposition a
little diff. e.g. HMM
$X_1 \to X_2 \to X_3 \to \ldots$
same CPT

$$\log P(\mathcal{D} \mid \theta_\mathcal{G}, \mathcal{G}) = \sum_{j=1}^{m} \sum_{i=1}^{n} \log P\left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)}\left[\mathbf{Pa}_{X_i}\right]\right)$$

$$\frac{\partial \log P(D|\theta_G, G)}{\partial \theta_{X_7|Pa_{X_7}}} = \frac{\partial}{\partial \theta_{X_7|Pa_{X_7}}} \left[ \sum_{j=1}^{m} \sum_{i=1}^{n} \log P\left(X_i = x_i^{(j)} \mid Pa_{X_i} = x^{(j)}[Pa_{X_i}], \theta_{X_i|Pa_{X_i}}\right) \right]$$

$$= \sum_{j=1}^{n} \frac{\partial}{\partial \theta_{X_7|Pa_{X_7}}} \log P\left(X_7 = x_7^{(j)} \mid Pa_{X_7} = x^{(j)}[Pa_{X_7}], \theta_{X_7|Pa_{X_7}}\right) = 0$$

Same as usual $\implies$  $\hat{P}(X_7^{=x_7} \mid Pa_{X_7} = U) \overset{MLE}{=} \dfrac{Count(X_7 = x_7, U = u)}{Count(U = u)}$

# General MLE for a CPT

- Take a CPT: P(X|**U**)
- Log likelihood term for this CPT


- Parameter $\theta_{X=x|U=u}$ :

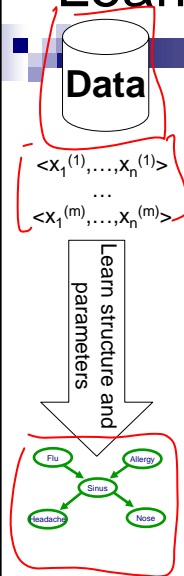MLE:   $P(X = x \mid U = u) = \theta_{X=x|U=u} = \dfrac{\text{Count}(X = x, U = u)}{\text{Count}(U = u)}$

# Where are we with learning BNs?

- Given structure, estimate parameters
  - Maximum likelihood estimation
  - Later Bayesian learning
- What about learning structure?

# Learning the structure of a BN

**Data**

$<x_1^{(1)},\ldots,x_n^{(1)}>$
$\ldots$
$<x_1^{(m)},\ldots,x_n^{(m)}>$

Learn structure and parameters

Flu   Allergy
Sinus
Headache   Nose

- **Constraint-based approach**
  - □ BN encodes conditional independencies
  - □ Test conditional independencies in data
  - □ Find an I-map
- **Score-based approach**
  - □ Finding a structure and parameters is a density estimation task
  - □ Evaluate model as we evaluated parameters
    - ■ Maximum likelihood
    - ■ Bayesian
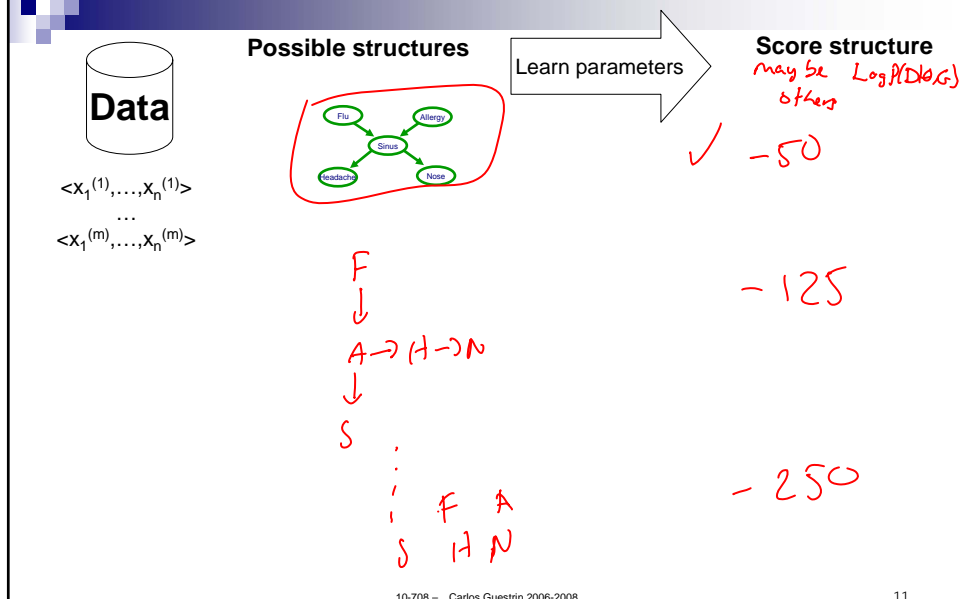    - ■ etc.

---

# Remember: Obtaining a P-map?

- Given the independence assertions that are true for *P*
  - □ Obtain skeleton
  - □ Obtain immoralities
- From skeleton and immoralities, obtain every (and any) BN structure from the equivalence class

- **Constraint-based approach**:
  - □ Use Learn PDAG algorithm
  - □ Key question: **Independence test**

# Score-based approach

**Possible structures**          Learn parameters          **Score structure**

may be  Log P(D|θ,G)
others

**Data**

$<x_1^{(1)},\ldots,x_n^{(1)}>$
$\ldots$
$<x_1^{(m)},\ldots,x_n^{(m)}>$

✓  −50

F
↓
A → H → N
↓
S
⋮
⋮      F  A
S   H  N

−125

−250

---

# Information-theoretic interpretation of maximum likelihood

- Given structure, log likelihood of data:

fixed

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^{m} \sum_{i=1}^{n} \log P\left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)}\left[\mathbf{Pa}_{X_i}\right]\right)$$

$$= \sum_{i=1}^{n} \left[ \sum_{j=1}^{m} \log P\left(x_i = x_i^{(j)} \mid Pa_{x_i} = x^{(j)}\left[Pa_{x_i}\right]\theta\right)\right]$$

$$= \sum_{i=1}^{n} \left[ \sum_{x_i} \sum_{\substack{u \in \\ Val(Pa_{x_i})}} Count(X_i = x_i, Pa_{x_i} = u) \log P(x_i = x_i \mid Pa_{x_i} = u)\theta\right]$$

$$= m \sum_{i=1}^{n} \left[ \sum_{x_i} \sum_{\substack{u \in \\ Val(Pa_{x_i})}} \hat{P}(x_i, Pa_{x_i} = u) \log P(X_i = x_i \mid Pa_{x_i} = u)\right]$$

$X_i = x_i, Pa_{x_i} = u$

$Count(X_i = x_i, Pa_{x_i} = u)$
number of times

MLE:
$\hat{P}(X_i = x_i, Pa_{x_i} = u) = \dfrac{Count(x_i, u)}{m}$

if MLE  $\hat{P}(X_i = x_i \mid Pa_N = u)$
$\theta \in$

# Information-theoretic interpretation of maximum likelihood 2

- Given structure, log likelihood of data:

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) \overset{MLE}{=} m \sum_i \sum_{x_i, \mathbf{Pa}_{x_i, \mathcal{G}}} \hat{P}(x_i, \mathbf{Pa}_{x_i, \mathcal{G}}) \log \hat{P}(x_i \mid \mathbf{Pa}_{x_i, \mathcal{G}})$$

$-H(X_i \mid Pa_{X_i})$

$$= - m \sum_i \hat{H}(X_i \mid Pa_{X_i})$$

$$= m \sum_i \hat{I}(X_i, Pa_{X_i}) - m \sum_i \hat{H}(X_i)$$

doesn't depend on structure

$H(A \mid B) = -\sum_a \sum_b P(a,b) \log P(a \mid b)$

$H(A \mid B) = H(A) \qquad A \perp B$

$H(A \mid B) = 0 \qquad A \, \& \, B$ perfectly correlated

$H(A \mid B) \geqslant 0$

$I(A, B) = H(A) - H(A \mid B)$

$I(X_i, Pa_{X_i}) = H(X_i) - H(X_i \mid Pa_{X_i})$

---

# Decomposable score

- Log data likelihood

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i}) - m \sum_i \hat{H}(X_i)$$

- Decomposable score:
  - Decomposes over families in BN (node and its parents)
  - Will lead to significant computational efficiency!!!
  - Score($G : D$) = $\sum_i$ FamScore($X_i | \mathbf{Pa}_{Xi} : D$)

for MLE    Fam Score $(X_i \mid Pa_{X_i} : D) = m \hat{I}(X_i, Pa_{X_i}) - m \hat{H}(X_i)$

# Announcements

- **Recitation tomorrow**
  - Don't miss it!

- **HW2**
  - Out today
  - Due in 2 weeks

- **Projects!!!** ☺
  - Proposals due Oct. 8th in class
  - Individually or groups of two
  - Details on course website
  - Project suggestions will be up soon!!!

*(handwritten notes:)*
① it is good related to research, but must be new

② it must have something to do w: Graphical models

---

# BN code release!!!!

- Pre-release of a C++ library for probabilistic inference and learning

- Features:
  - basic datastructures (random variables, processes, linear algebra)
  - distributions (Gaussian, multinomial, ...)
  - basic graph structures (directed, undirected)
  - graphical models (Bayesian network, MRF, junction trees)
  - inference algorithms (variable elimination, loopy belief propagation, filtering)
- Limited amount of learning (IPF, Chow Liu, order-based search)
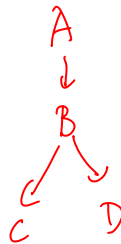
- Supported platforms:
  - Linux (tested on Ubuntu 8.04)
  - MacOS X (tested on 10.4/10.5)
  - limited Windows support

- Will be made available to the class early next week.

# How many trees are there?

*# parents = 1*

**Nonetheless – Efficient optimal algorithm finds best tree**

A
↓
B
↙ ↘
C  D

C

$A \rightarrow C \rightarrow D \rightarrow B$

$2^{\Theta(n \log n)}$

---
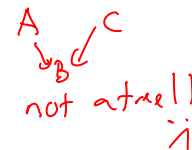
# Scoring a tree 1: I-equivalent trees

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i}) - m \sum_i \hat{H}(X_i)$$

$A \rightarrow B \rightarrow C$

Score: $m[I(B,A) + I(C,B) - H(A) - H(B) - H(C) + I(A, \emptyset)]$

$A \leftarrow B \rightarrow C$

$m[I(A,B) + I(C,B) - H(B) + H(A) - H(C)]$

I-equivalent
same score!!

A ↘ C
  B
not a tree!!

$A \leftarrow B \leftarrow C$

$m[I(A,B) + I(B,C) - H(A) - H(B) - H(C)]$

SAME Skeleton ⇒ Same Score!! (only for trees)

# Scoring a tree 2: similar trees

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i}) - m \sum_i \hat{H}(X_i)$$

$G_1$

$A \to B \to C$

$m[I(A,B) + I(B,C) - H(A) - H(B) - H(C)]$

$G_2$

$B \to A \to C$

$m[I(B,A) + I(A,C) - H(A) - H(B) - H(C)]$

$Score(tree) = m \sum_{ij} I(X_i, X_j) - m \sum_i H(X_i)$

is in skeleton

constant ignore

Score($G_1$) − Score($G_2$)
= $I(B,C) - I(A,C)$

edges Skeleton

same:
$A - B$
⟹ Score ← $I(A,B)$

different

$G_1$
$B - C$
$I(B,C)$

$G_2$
$A - C$
$I(A,C)$

---

# Chow-Liu tree learning algorithm 1

best tree

- For each pair of variables $X_i, X_j$
  □ Compute empirical distribution:
  $$\bar{P}(x_i, x_j) \overset{MLE}{\equiv} \frac{Count(x_i, x_j)}{m}$$
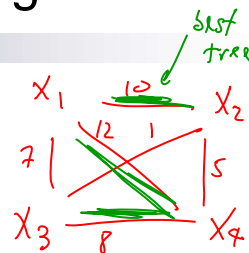  □ Compute mutual information:
  $$\bar{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$

- Define a graph
  □ Nodes $X_1, \ldots, X_n$   $w_{ij}$
  □ Edge (i,j) gets weight $\hat{I}(X_i, X_j)$

$X_1$  10  $X_2$
7  12  1  5
$X_3$  8  $X_4$

$\max_{trees} Score(tree)$
$= \sum_{ij} I(X_i, X_j)$
$= \sum_{ij} w_{ij}$

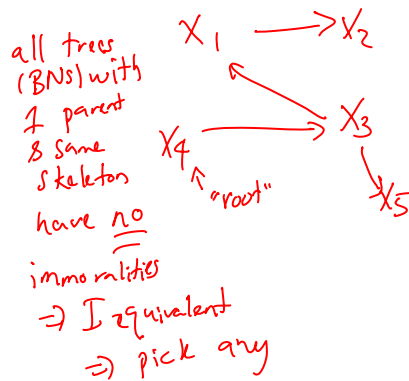find Maximum Spanning tree

best tree BN

# Chow-Liu tree learning algorithm 2

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \mathbf{Pa}_{x_i,\mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

- Optimal tree BN
  - Compute maximum weight spanning tree
  - Directions in BN: pick any node as <u>root</u>, breadth-first-search defines directions

*(handwritten annotations)*

all trees (BNs) with 1 parent & same skeleton have <u>no</u> immoralities
$\Rightarrow$ I equivalent
$\Rightarrow$ pick any

$X_1 \longrightarrow X_2$
$X_4 \longrightarrow X_3$
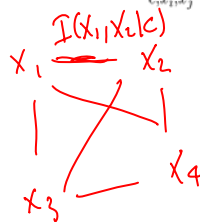"root"
$X_5$

using Chow-Liu OPTIMAL tree BN

---

# Can we extend Chow-Liu 1

- Tree augmented naïve Bayes (TAN)
  [Friedman et al. '97]
  - Naïve Bayes model overcounts, because correlation between features not considered
  - Same as Chow-Liu, but score edges with:

$$\hat{I}(X_i, X_j \mid C) = \sum_{c, x_i, x_j} \hat{P}(c, x_i, x_j) \log \frac{\hat{P}(x_i, x_j \mid c)}{\hat{P}(x_i \mid c)\hat{P}(x_j \mid c)}$$

*(handwritten annotations)*

$C \longrightarrow X_n$
$X_1 \quad X_2$

TAN can reduce double counting

$C$
$X_1 \rightarrow X_2 \rightarrow X_3 \quad X_4$

$I(X_1, X_2 \mid C)$
$X_1 \qquad X_2$
$X_3 \qquad X_4$

Maximum Spanning tree, $\Rightarrow$ optimal TAN

SCORE:
$I(C, \emptyset) + I(C, X_1) + I(X_2, C X_1) + I(X_3, C X_2) + I(X_4, C X_2)$
+ entropy