# Bayesian Networks – (Structure) Learning

Machine Learning – CSE546

Carlos Guestrin

University of Washington
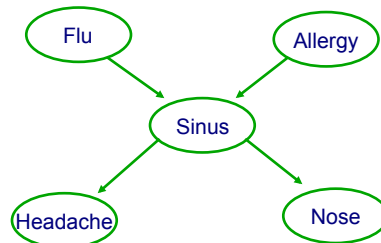
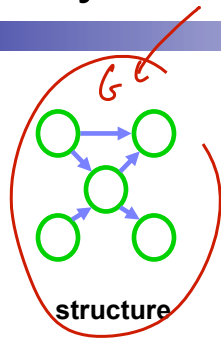November 25, 2013
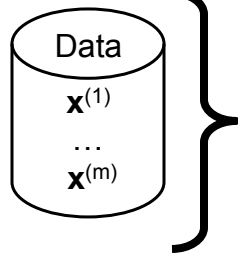
1

---

# Review

- Bayesian Networks
  - □ Compact representation for probability distributions
  - □ Exponential reduction in number of parameters
- Fast probabilistic inference
  - □ As shown in demo examples
  - □ Compute P(X|e)
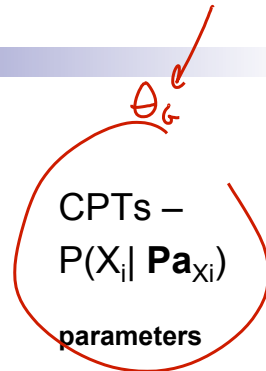- Today
  - □ Learn BN structure

2

# Learning Bayes nets

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

$G$

**structure**

$\theta_G$

CPTs –
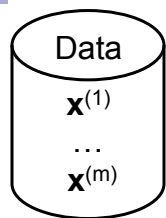$P(X_i | \mathbf{Pa}_{Xi})$

**parameters**

+

max likelihood approach
structure, & params

$$P(D | G, \theta_G)$$

3

---

# Learning the CPTs

$|Y|$ is # of assignments,
e.g. $Pa_{x_i} = \{F, A, H\}$, binary
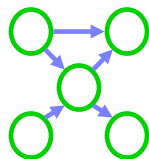$|Y| = |F, A, H| = 2^3$, $\widehat{Count}(Pa_{x_i} = c)$

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

given $G$

For each discrete variable $X_i$

$\hat{P}(S = t | A = t) \overset{MLE}{=} \dfrac{Count(S = t, A = t)}{Count(A = t)}$

$\hat{P}(X_i = x_i | Pa_{x_i} = w) \overset{MLE}{=} \dfrac{Count(X_i = x_i, Pa_{x_i} = w)}{Count(Pa_{x_i} = w)}$

Small (huge) subtlety: $Count(Pa_{x_i} = w) = 0$ ?

Smoothing / AKA regularization / AKA Bayesian learning

$\widehat{Count}(Y = y) = Count(Y = y) + \alpha \dfrac{1}{|Y|}$, for $\alpha > 0$ usually $\alpha \approx 1$

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

4

2

# Information-theoretic interpretation of maximum likelihood 1

*n variables, m data points*

Flu, Allergy, Sinus, Headache, Nose

- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) \overset{iid}{=} \log \prod_{j=1}^{m} P(x_1^{(j)}, \ldots, x_n^{(j)} \mid \theta_{\mathcal{G}}, \mathcal{G})$$

$$= \log \prod_{j=1}^{m} \prod_{i=1}^{n} P(x_i^{(j)} \mid Pa_{X_i, \mathcal{G}} = u_i^{(j)}, \theta_{\mathcal{G}})$$

$$= \sum_{j=1}^{m} \sum_{i=1}^{n} \log P(x_i^{(j)} \mid Pa_{X_i, \mathcal{G}} = u_i^{(j)}, \theta_{\mathcal{G}}) \quad \leftarrow \begin{array}{c} \text{max} \\ \mathcal{G} \end{array}$$

$$x^{(j)} \leftarrow (F=t, A=f, S=t, \ldots)$$
$$\underbrace{\qquad}_{u_i^{(j)}} \underbrace{\qquad}_{x_i^{(j)}}$$

5

---

# Information-theoretic interpretation of maximum likelihood 2

$P(y) = \frac{count(y=y)}{m}$

Flu, Allergy, Sinus, Headache, Nose

- Given structure, log likelihood of data:

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

$u_i^{(j)}$

$$= \sum_{i=1}^{n} \sum_{j=1}^{m} \log P(x_i^{(j)} \mid Pa_{X_i, \mathcal{G}} = u_i^{(j)})$$

$$= \sum_{i=1}^{n} \sum_{x_i} \sum_{u_i} count(X_i = x_i, Pa_{X_i, \mathcal{G}} = u_i)$$
$$\log P(x_i \mid Pa_{X_i, \mathcal{G}} = u_i)$$

$$= m \sum_{i=1}^{n} \sum_{x_i} \sum_{u_i} \hat{P}(x_i \mid Pa_{X_i, \mathcal{G}} = u_i) \log \hat{P}(x_i \mid Pa_{X_i, \mathcal{G}} = u_i)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{-H(X_i \mid Pa_{X_i, \mathcal{G}})}$$

e.g. $\sum_{j=1}^{m} \log P(h^{(j)} \mid s^{(j)})$

$= count(H=t, S=t) \cdot \log P(H=t \mid S=t)$
$+$
$count(H=f, S=f) \cdot \log P(H=f \mid S=f)$
$+$
$count(H=f, S=t) \cdot \log P(H=f \mid S=t)$
$+$
$count(H=t, S=f) \cdot \log P(H=f \mid S=f)$

6

3

# Information-theoretic interpretation of maximum likelihood 3

$I(A,B)$
$= I(B,A)$

Flu | Allergy | Sinus | Headache | Nose

- Given structure, log likelihood of data:

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

$$\underset{G}{\max} \; = -m \sum_{i=1}^n \hat{H}(X_i \mid Pa_{x_i, G}) \equiv \underset{G}{\min} \, m \sum_{i=1}^n \hat{H}(X_i \mid Pa_{x_i, G})$$

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

$\equiv \underset{G}{\max} \; m \sum_{i=1}^n I(X_i, Pa_{x_i, G}) \; - \; m \sum_{i=1}^n \hat{H}(X_i)$

How uncertain $X_i$ is given parents $\Rightarrow$ minimize this over $G$

Mutual Information — information theoretic measure of dependency

Constant w.r.t. $G$

Mutual information
$I(A,B)$
$= H(A) - H(A|B)$

©Carlos Guestrin 2005-2013    7

---

# Decomposable score

OK: $A \to B \searrow C$  |  NOT: $A \to B$, $C$

- Log data likelihood

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

families          Constant

- Decomposable score:
  - Decomposes over families in BN (node and its parents)
  - Will lead to significant computational efficiency!!!
  - $\underset{G}{\max}$ Score($G : D$) = $\sum_{i=1}^n$ FamScore($X_i | \mathbf{Pa}_{X_i} : D$)
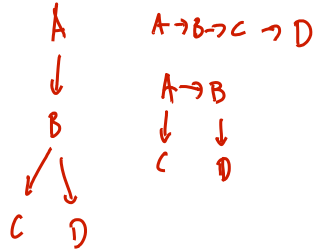
e.g. $= \sum_{i=1}^n I(X_i, Pa_{x_i, G})$

also get this decomposition for other losses

©Carlos Guestrin 2005-2013    8

4

# How many trees are there?

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

$A \to B \to C \to D$

$A \to B$
$\quad \downarrow \quad \downarrow$
$\quad C \quad D$
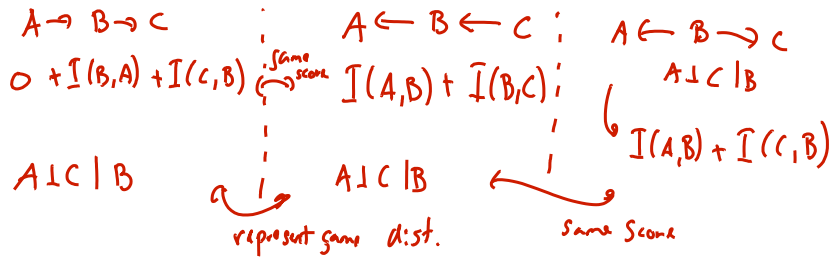
HHM is a tree

Naive Bayes is a tree

How many trees are there?
For n variables: $O(n^{\log n})$

exhaustive search is impossible

9

---

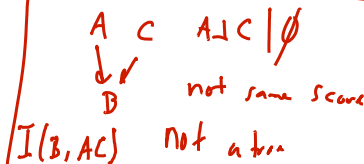# Scoring a tree 1: equivalent trees

$I(A,D) = 0$

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i, \mathcal{G}}) - m \sum_i \hat{H}(X_i) = \max_{\mathcal{G}} \sum_{i=1}^{n} \hat{I}(X_i, \mathbf{Pa}_{\mathcal{G}})$$

$A \to B \to C$
$0 + I(B,A) + I(C,B) \xleftrightarrow{\text{same score}}$
$A \perp C \mid B$
represent same dist.

$A \leftarrow B \leftarrow C$
$I(A,B) + \hat{I}(B,C)$
$A \perp C \mid B$

$A \leftarrow B \to C$
$A \perp C \mid B$
$I(A,B) + \hat{I}(C,B)$
Same Score

Every graph with same independence assumptions has same Score

$A \quad C \qquad A \perp C \mid \emptyset$
$\quad \searrow \swarrow$
$\quad B \qquad$ not same score
$I(B, AC) \quad$ not a tree

10

5

# Scoring a tree 2: similar trees

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

A
B
C

Score
$I(A,B) + I(B,C)$

A
B  C

$I(A,B) + I(A,C)$

$$\max_{G} \text{ Score of tree} \equiv \text{Score}(T) = \sum_{(i,j) \in T} I(X_i, X_j) = \text{Sum over edges of Score of edge}$$

11

---

# Chow-Liu tree learning algorithm 1

- For each pair of variables $X_i, X_j$
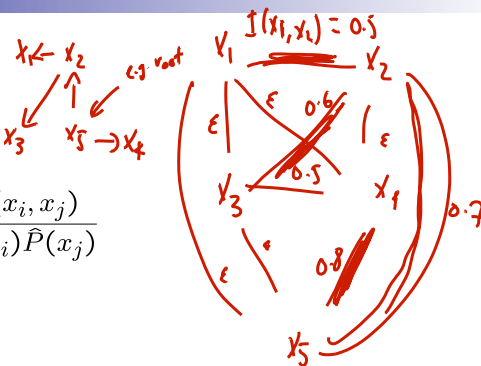  - Compute empirical distribution:
  $$\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$$
  - Compute mutual information:
  $$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$
- Define a graph
  - Nodes $X_1, \dots, X_n$
  - Edge (i,j) gets weight $\hat{I}(X_i, X_j)$

$X_1 \leftarrow X_2$   e.g. root   $X_1 \quad I(x_1, x_2) = 0.5 \quad X_2$

$X_3 \quad X_5 \rightarrow X_4$

$\varepsilon \quad 0.6 \quad \varepsilon$

$0.5$

$X_3 \quad X_1 \quad 0.7$

$\varepsilon \quad 0.8$

$X_5$

D   Maximum Spanning tree:
Find tree with maximum weight on edges
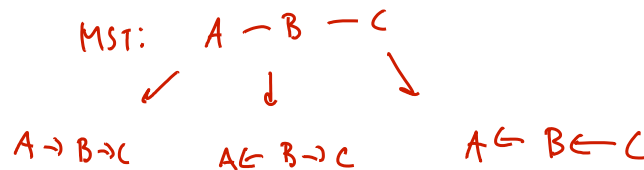$\Big\}$ $\longmapsto$ Complexity about $O(E \log E)$

12

6

# 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 root, breadth-first-search defines directions

*t all root give same score*

*MST:*  $A \frown B - C$

$A \rightarrow B \rightarrow C$    $A \leftarrow B \rightarrow C$    $A \leftarrow B \leftarrow C$

# Structure learning for general graphs

- In a tree, a node only has one parent

- **Theorem**:
  - The problem of learning a BN structure with at most *d* parents is NP-hard for any (fixed) *d>1*

  *for d >1*

- Most structure learning approaches use heuristics
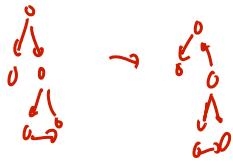  - (Quickly) Describe the two simplest heuristic

# Learn BN structure using local search

**Starting from Chow-Liu tree**

**Local search,** possible moves:
- Add edge
- Delete edge
- Invert edge

**Score using BIC**

Penalize for dense graphs

Push away from fully connected graph

Converge to local optima

©Carlos Guestrin 2005-2013                    15

---

# Learn Graphical Model Structure using LASSO

- Graph structure is about selecting parents:

$P(X_i | Pa_{X_i, G}) \leftarrow$ "logistic regression"

$Pa_{X_i, G} \subset \{X_1, \dots X_{i-1}, X_{i+1}, \dots X_n\}$

- If no independence assumptions, then CPTs depend on all parents:

$P(H | FASN)$

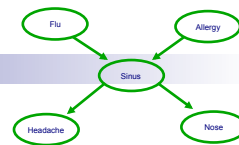- With independence assumptions, depend on key variables:

$P(H | FASN) = P(H|S) \leftarrow$ Sparse conditional model where other dependencies are zero

- One approach for structure learning, sparse logistic regression!

LR for each variable: $P(X_1 | X_2 \dots X_n) \leftarrow$ Sparse LR
$\parallel$
caviat: this approach not appropriate (add edges from all non-zero parents to $X_i$)
for BNs, but used in other graphical models, like Markov Networks, undirected

©Carlos Guestrin 2005-2013                    16

8

# What you need to know about learning BN structures

- Decomposable scores
  - Maximum likelihood
  - Information theoretic interpretation
- Best tree (Chow-Liu)
- Beyond tree-like models is NP-hard
- Use heuristics, such as:
  - Local search
  - LASSO

17

# Learning Theory

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 27, 2013

18

# What now…

- We have explored **many** ways of learning from data
- But…
  - How good is our classifier, really?
  - How much data do I need to make it "good enough"?

# A simple setting…

- Classification
  - N data points *iid*
  - **Finite** number of possible hypothesis (e.g., dec. trees of depth d)
- A learner finds a hypothesis $h$ that is **consistent** with training data
  - Gets zero error in training – $\text{error}_{train}(h) = 0$
- What is the probability that $h$ has more than $\varepsilon$ true error?
  - $\text{error}_{true}(h) \geq \varepsilon$      For some    $\varepsilon > 0$

# How likely is a bad hypothesis to get *N* data points right?

- Hypothesis *h* that is **consistent** with training data → got *N* i.i.d. points right   $\varepsilon > 0$
  - □ h "bad" if it gets all this data right, but has high true error
- Prob. *h* with error$_{true}$(h) ≥ ε  gets one data point right
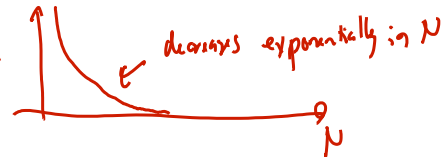
  less than   $1-\varepsilon$ | if error $\varepsilon = 0.25$
  75% points are
  correct $= 1-\varepsilon$

- Prob. *h* with error$_{true}$(h) ≥ ε  gets *N* data points right
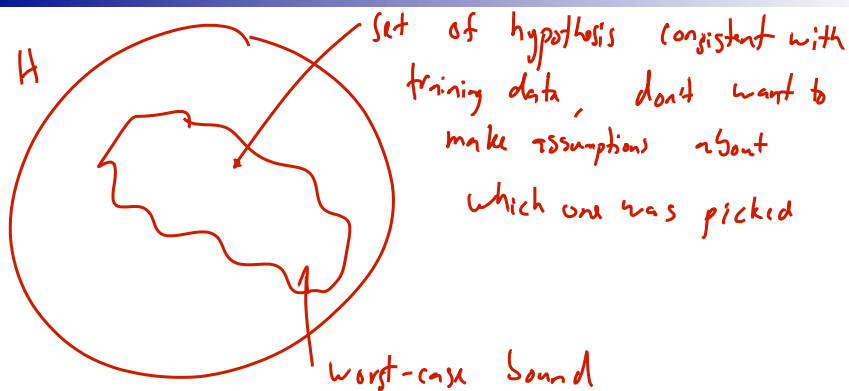
  less than   $(1-\varepsilon)^N$     Prob bad h wins     decreases exponentially in N

21

---

# But there are many possible hypothesis that are consistent with training data

H

Set of hypothesis consistent with
training data, don't want to
make assumptions about
which one was picked

worst-case bound

22

11

# How likely is learner to pick a bad hypothesis

- Prob. $h$ with error$_{true}$(h) $\geq \epsilon$ gets $N$ data points right

  less than $(1-\epsilon)^N$

  $\rightarrow h_1, \dots h_k$

- There are $k$ hypothesis consistent with data
  - How likely is learner to pick a bad one?    Some bad, Some good

$$P\left( \exists h \text{ Consistent with data}^{train}, \text{error}_{true}(h) \geq \epsilon \right) \quad \exists \text{ deal with worst case}$$

$$= P\left( \text{error}_{true}(h_1) \geq \epsilon \text{ OR error}_{true}(h_2) \geq \epsilon \text{ OR} \dots \text{ OR error}_{true}(h_k) \geq \epsilon \right)$$

$$\text{Bound?}$$
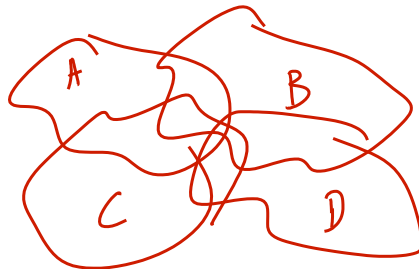
23

# Union bound

- $P(A \text{ or } B \text{ or } C \text{ or } D \text{ or } \dots) \leq P(A) + P(B) + P(C) + P(D) \dots$

24

# How likely is learner to pick a bad hypothesis

- Prob. a particular $h$ with error$_{true}$(h) $\geq \epsilon$ gets $N$ data points right *less than* $(1-\epsilon)^N$

- There are $k$ hypothesis consistent with data
  - How likely is it that learner will pick a bad one out of these $k$ choices?

$$P(\exists h \text{ consistent with train data, } error_{true}(h) \geq \epsilon) \leq k(1-\epsilon)^N$$

$$\leq |H|(1-\epsilon)^N$$

*what's $k$?*

$$K \leq |H|$$

*total # hypothesis*

*(vrzy loose)*

©Carlos Guestrin 2005-2013

25

# Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem**: Hypothesis space $H$ finite, dataset $D$ with $N$ i.i.d. samples, $0 < \epsilon < 1$ : for any learned hypothesis $h$ that is consistent on the training data:

$$P(error_{true}(h) \geq \epsilon) \leq |H|e^{-N\epsilon}$$

*prob picking bad*

*Because exponentially*

*decreases exponentially*

*N*

$$\leq |H|(1-\epsilon)^N \leq |H|(e^{-\epsilon})^N = |H|e^{-\epsilon N}$$

*for $0 \leq \epsilon \leq 1$*

$$1-\epsilon \lesssim e^{-\epsilon}$$

©Carlos Guestrin 2005-2013

26

13