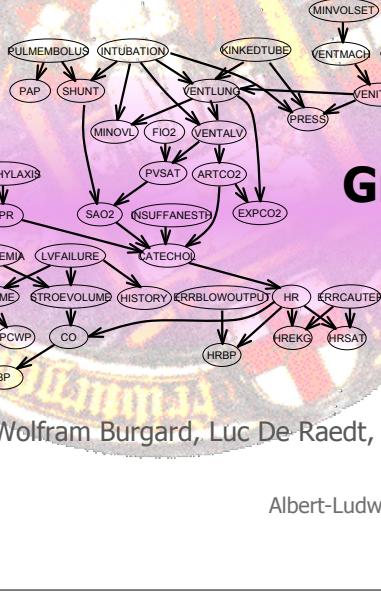


Based on J. A. Bilmes, "A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models", TR-97-021, U.C. Berkeley, April 1998; G. J. McLachlan, T. Krishnan, "The EM Algorithm and Extensions", John Wiley & Sons, Inc., 1997; D. Koller, course CS-228 handouts, Stanford University, 2001., N. Friedman & D. Koller's NIPS'99.

Advanced I

WS 06/07



Advanced
I
WS 06/07

Learning With Bayesian Networks

Fixed structure	Fixed variables	Hidden variables
Easiest problem counting	Selection of arcs New domain with no domain expert Data mining	
Numerical, nonlinear optimization, Multiple calls to BNs, Difficult for large networks	Encompasses to difficult subproblem, „Only“ Structural EM is known	Scientific discovery

Parameter Estimation

Structure learning ?

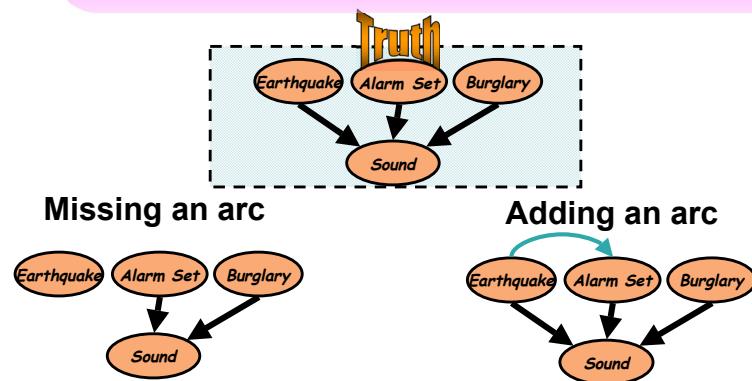
fully observed

Partially

Bayesian Networks - Learning

Advanced
I
WS 06/07

Why Struggle for Accurate Structure?



- Cannot be compensated for by fitting parameters
- Wrong assumptions about domain structure
- Increases the number of parameters to be estimated
- Wrong assumptions about domain structure

Bayesian Networks - Learning

Advanced
I
WS 06/07

Unknown Structure, (In)complete Data

E, B, A
<Y,N,N>
<Y,N,Y>
<N,N,Y>
<N,Y,Y>
.
<N,Y,Y>

E	B	P(A E,B)	
e	b	?	?
e	-	?	?
-	b	?	?
-	-	?	?

E, B, A
<Y,?,N>
<Y,N,?>
<N,N,Y>
<N,Y,Y>
.
<?,Y,Y>

| E | B | P(A E,B) | |
|---|---|------------|-----|
| e | b | .9 | .1 |
| e | - | .7 | .3 |
| - | b | .8 | .2 |
| - | - | .99 | .01 |

Learning algorithm

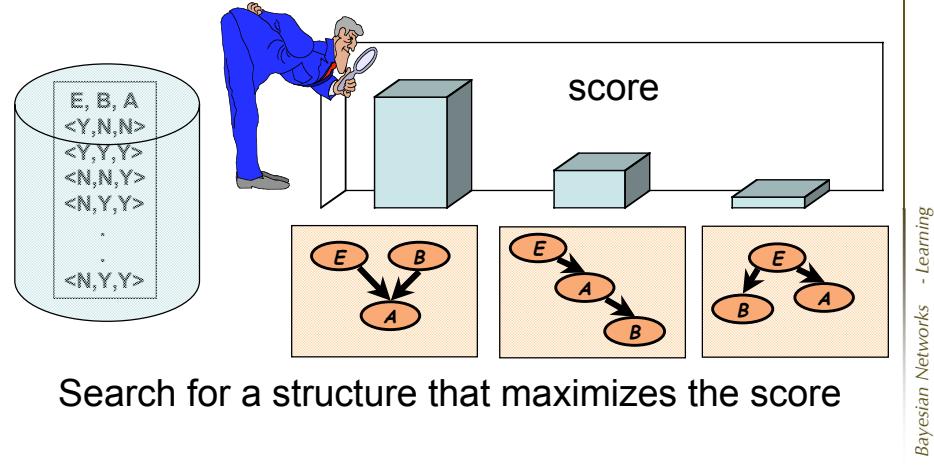
| | | | |
|---|---|------------|-----|
| E | B | P(A E,B) | |
| e | b | .9 | .1 |
| e | - | .7 | .3 |
| - | b | .8 | .2 |
| - | - | .99 | .01 |

- Network structure is not specified
- Data contains missing values
 - Need to consider assignments to missing values

Bayesian Networks - Learning

Score-based Learning

Define scoring function that evaluates how well a structure matches the data



Structure Search as Optimization

Input:

- Training data
- Scoring function
- Set of possible structures

Output:

- A network that maximizes the score

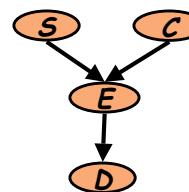
Heuristic Search

- Define a search space:
 - search states are possible structures
 - operators make small changes to structure
- Traverse space looking for high-scoring structures
- Search techniques:
 - Greedy hill-climbing
 - Best first search
 - Simulated Annealing
 - ...

Theorem: Finding maximal scoring structure with at most k parents per node is NP-hard for $k > 1$

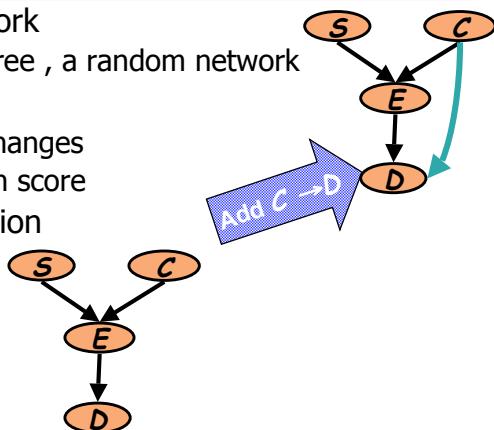
Typically: Local Search

- Start with a given network
 - empty network, best tree, a random network
- At each iteration
 - Evaluate all possible changes
 - Apply change based on score
- Stop when no modification improves score



Typically: Local Search

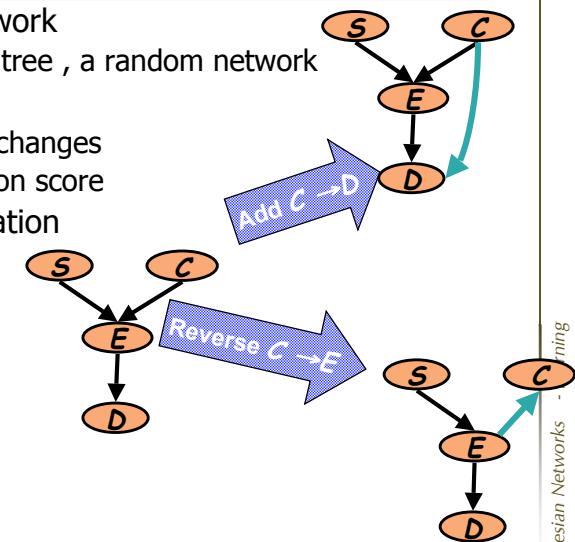
- Start with a given network
 - empty network, best tree , a random network
- At each iteration
 - Evaluate all possible changes
 - Apply change based on score
- Stop when no modification improves score



Bayesian Networks - Learning

Typically: Local Search

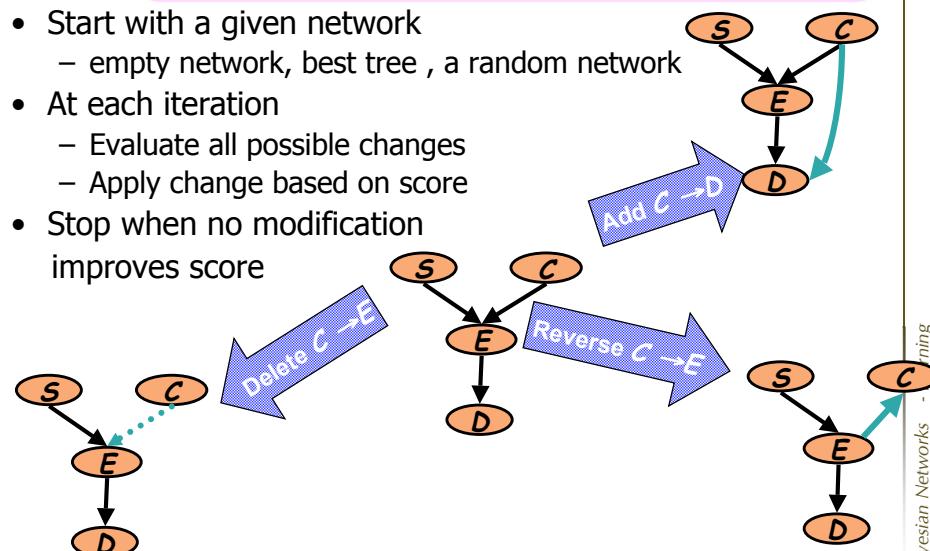
- Start with a given network
 - empty network, best tree , a random network
- At each iteration
 - Evaluate all possible changes
 - Apply change based on score
- Stop when no modification improves score



Bayesian Networks - Learning

Typically: Local Search

- Start with a given network
 - empty network, best tree , a random network
- At each iteration
 - Evaluate all possible changes
 - Apply change based on score
- Stop when no modification improves score

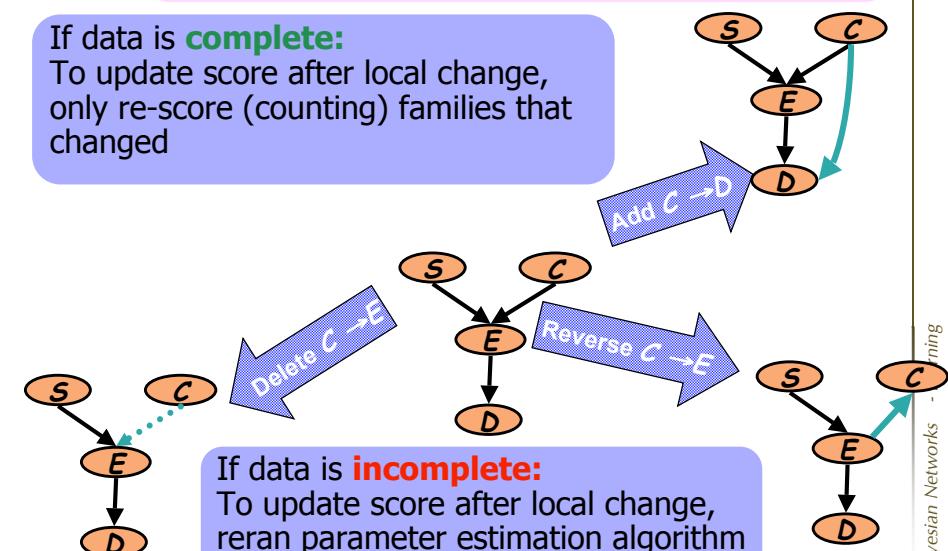


Bayesian Networks - Learning

Typically: Local Search

If data is **complete**:

To update score after local change,
only re-score (counting) families that
changed



Bayesian Networks - Learning

Local Search in Practice

- Local search can get stuck in:
 - **Local Maxima:**
 - All one-edge changes reduce the score
 - **Plateaux:**
 - Some one-edge changes leave the score unchanged
- Standard heuristics can escape both
 - Random restarts
 - TABU search
 - Simulated annealing

Local Search in Practice

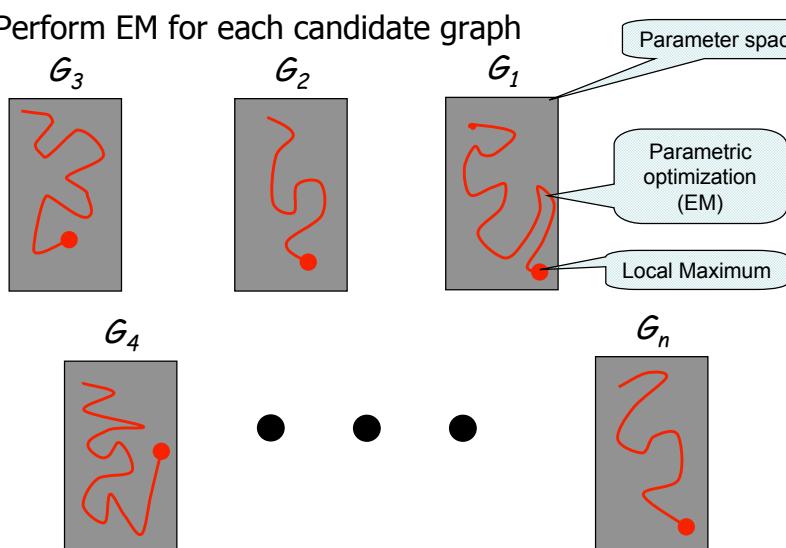
- Using LL as score, adding arcs always helps
 - Max score attained by fully connected network
 - Overfitting: A bad idea...
- Minimum Description Length:
 - Learning \Leftrightarrow data compression

$$MDL(BN | D) = \underbrace{-\log P(D | \Theta, G)}_{DL(\text{Data|model})} + \underbrace{\frac{\log N}{2} |\Theta|}_{DL(\text{Model})}$$

- Other: BIC (Bayesian Information Criterion), Bayesian score (BDe)

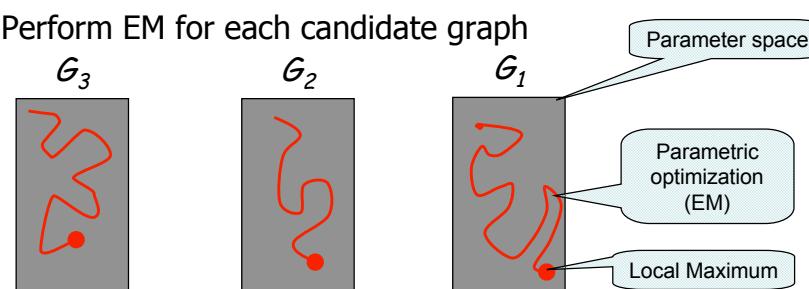
Local Search in Practice

- Perform EM for each candidate graph



Local Search in Practice

- Perform EM for each candidate graph



- ♦ Computationally expensive:
- Parameter optimization via EM — non-trivial
 - Need to perform EM for all candidate structures
 - Spend time even on poor candidates
- ⇒ In practice, considers only a few candidates

Structural EM

[Friedman et al. 98]

Recall, in complete data we had
 –Decomposition \Rightarrow efficient search

Idea:

- Instead of optimizing the real score...
- Find **decomposable** alternative score
- Such that maximizing new score
 \Rightarrow improvement in real score

Structural EM

[Friedman et al. 98]

Idea:

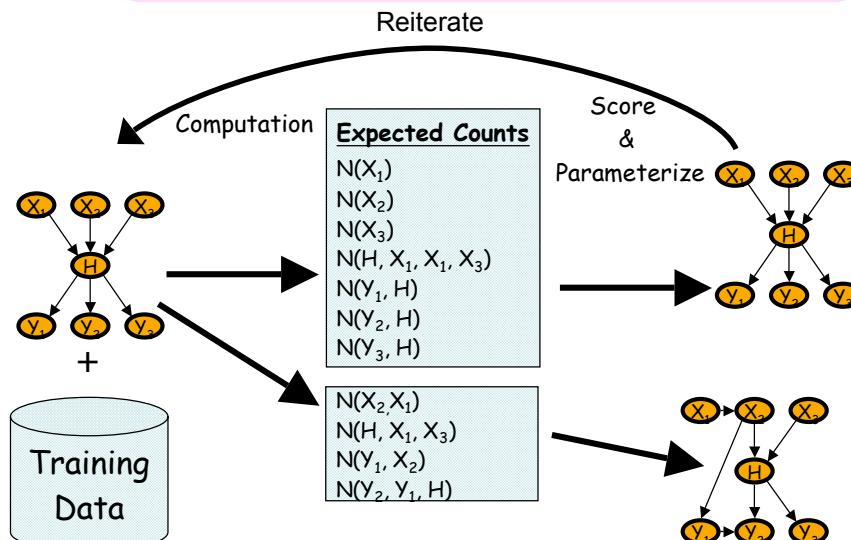
- Use current model to help evaluate new structures

Outline:

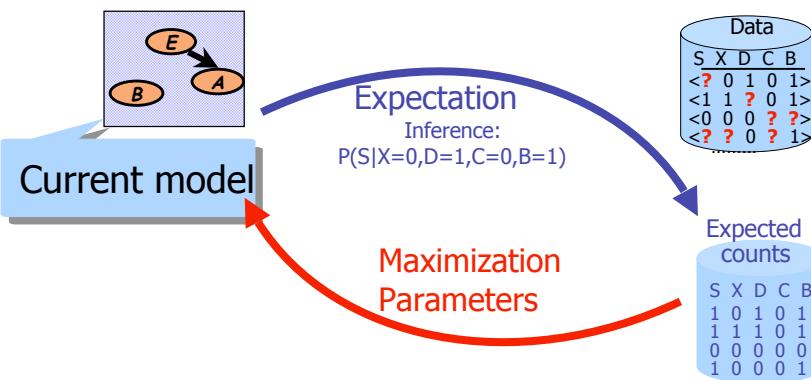
- Perform search in (Structure, Parameters) space
- At each iteration, use current model for finding either:
 - Better scoring parameters: “parametric” EM step or
 - Better scoring structure: “structural” EM step

Structural EM

[Friedman et al. 98]

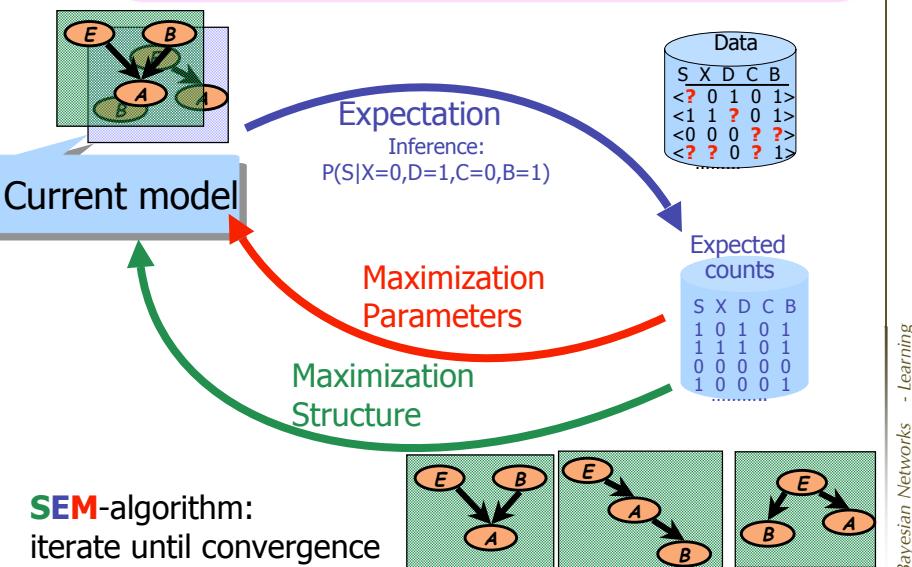


Structure Learning: incomplete data



EM-algorithm:
 iterate until convergence

Structure Learning: incomplete data



Structure Learning: Summary

- Expert **knowledge** + learning from **data**
- Structure learning involves parameter estimation (e.g. EM)
- Optimization w/ **score** functions
 - likelihood + complexity penalty = MDL
- **Local traversing** of space of possible structures:
 - add, reverse, delete (single) arcs
- Speed-up: **Structural EM**
 - Score candidates w.r.t. current best model