

A Tutorial on Inference and Learning in Bayesian Networks

Irina Rish

IBM T.J.Watson Research Center

rish@us.ibm.com

<http://www.research.ibm.com/people/r/rish/>

Outline

- Motivation: learning probabilistic models from data
- Representation: Bayesian network models
- Probabilistic inference in Bayesian Networks
 - Exact inference
 - Approximate inference
- Learning Bayesian Networks
 - Learning parameters
 - Learning graph structure (model selection)
- Summary

Bayesian Networks

Structured, graphical representation of probabilistic relationships between several random variables

Explicit representation of conditional independencies

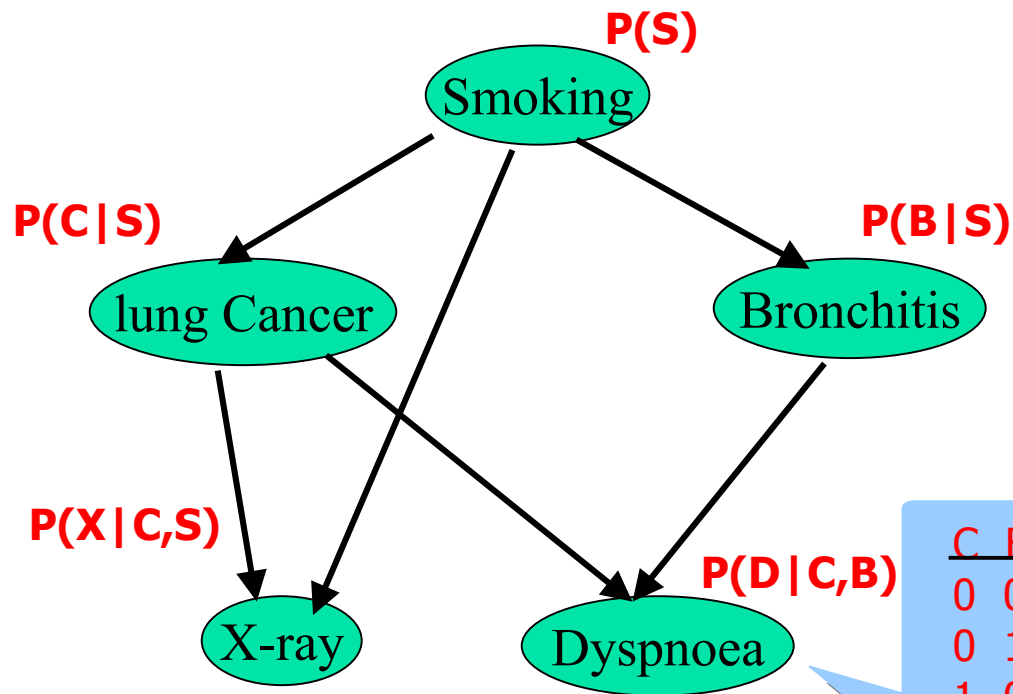
Missing arcs encode conditional independence

Efficient representation of joint PDF $P(X)$

Generative model (not just discriminative): allows arbitrary queries to be answered, e.g.

$$P(\text{lung cancer}=\text{yes} \mid \text{smoking}=\text{no}, \text{positive X-ray}=\text{yes}) = ?$$

Bayesian Network: $\text{BN} = (\mathbf{G}, \Theta)$



\mathbf{G} - directed acyclic graph (DAG)
 nodes – random variables
 edges – direct dependencies

Θ - set of parameters in all conditional probability distributions (CPDs)

CPD:

C	B	D=0	D=1
0	0	0.1	0.9
0	1	0.7	0.3
1	0	0.8	0.2
1	1	0.9	0.1

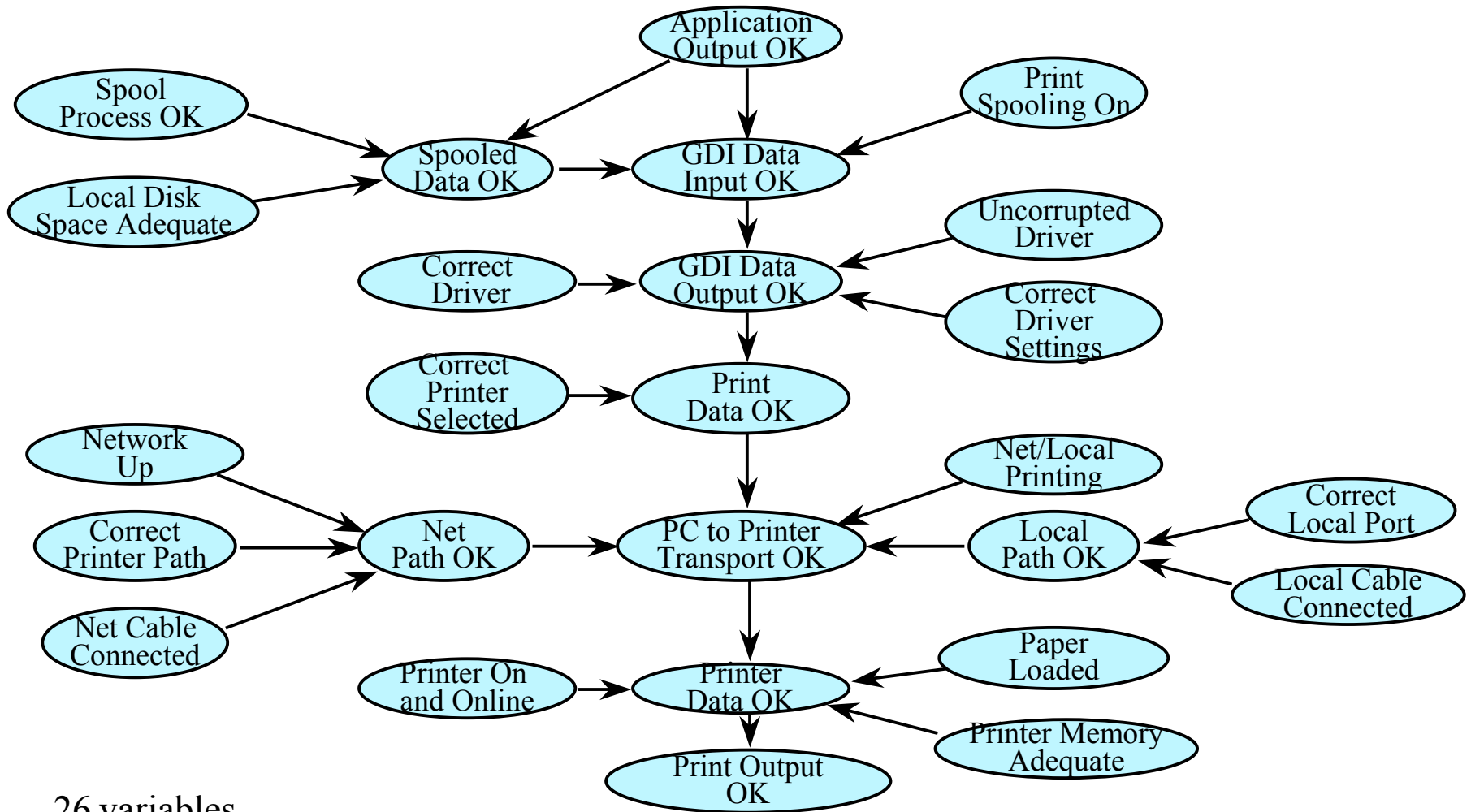
CPD of node X:
 $P(X | \text{parents}(X))$

Compact representation of joint distribution in a **product form** (chain rule):

$$P(S, C, B, X, D) = P(S) P(C|S) P(B|S) P(X|C,S) P(D|C,B)$$

$$1 + 2 + 2 + 4 + 4 = 13 \text{ parameters instead of } 2^5 = 32$$

Example: Printer Troubleshooting



26 variables

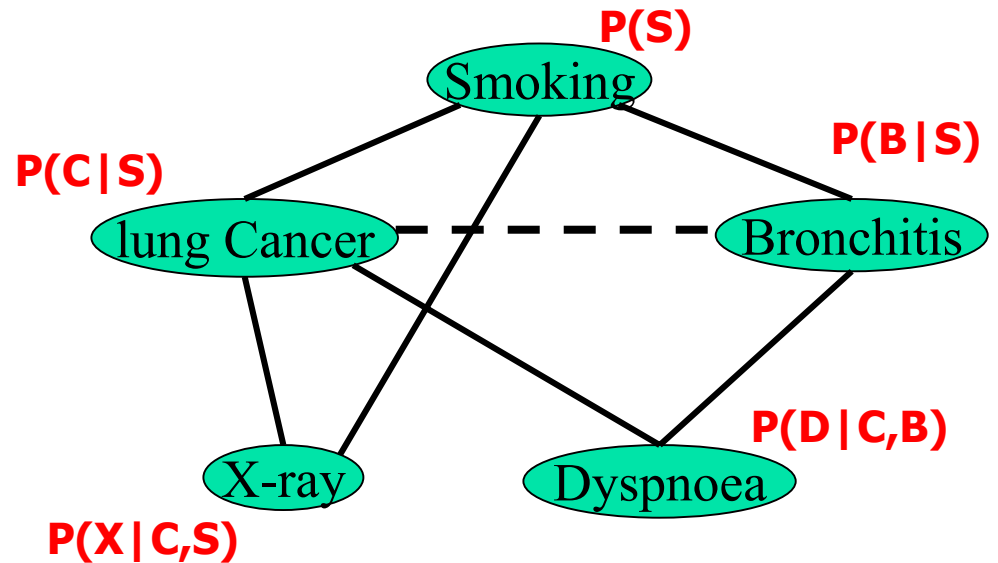
Instead of 2^{26} parameters we get

$$99 = 17x1 + 1x2^1 + 2x2^2 + 3x2^3 + 3x2^4$$

"Moral" graph of a BN

Moralization algorithm:

1. Connect ("marry") parents of each node.
2. Drop the directionality of the edges.



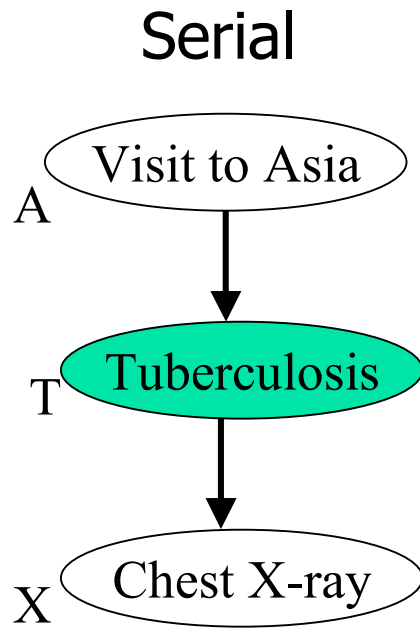
Resulting undirected graph is called the "moral" graph of BN

Interpretation:

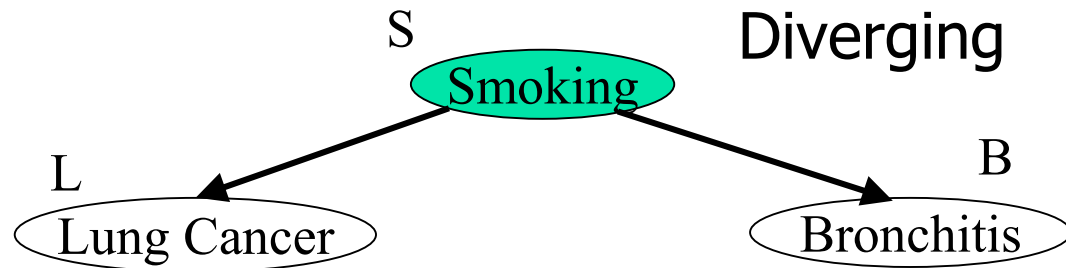
every pair of nodes that occur together in a CPD is connected by an edge in the moral graph.

CPD for X and its k parents (called "**family**") is represented by a clique of size **($k+1$)** in the moral graph, and contains $d^k (d-1)$ probability parameters where d is the number of values each variable can have (domain size).

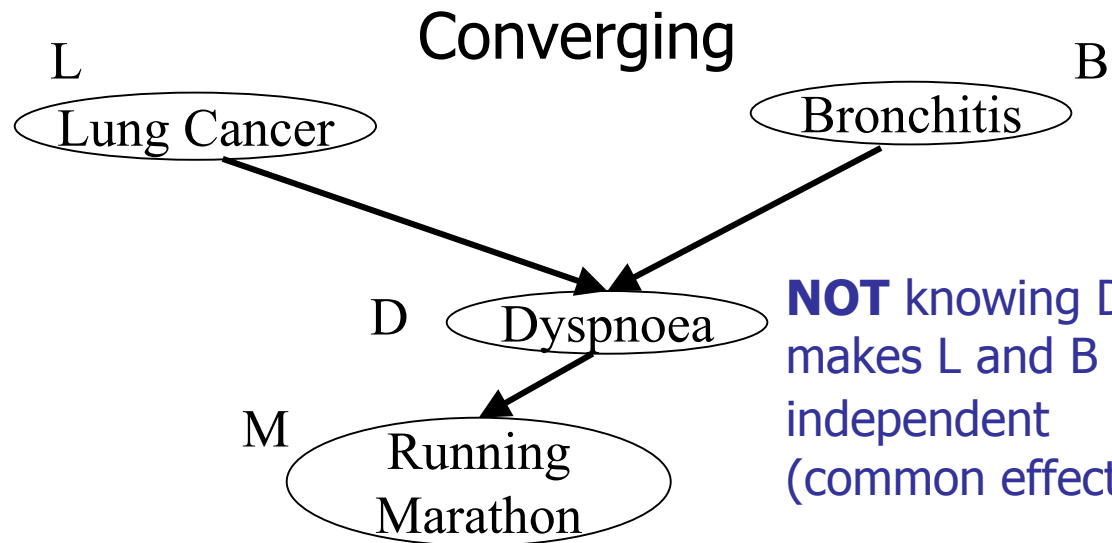
Conditional Independence in BNs: Three types of connections



Knowing T makes A and X independent (intermediate cause)



Knowing S makes L and B independent (common cause)



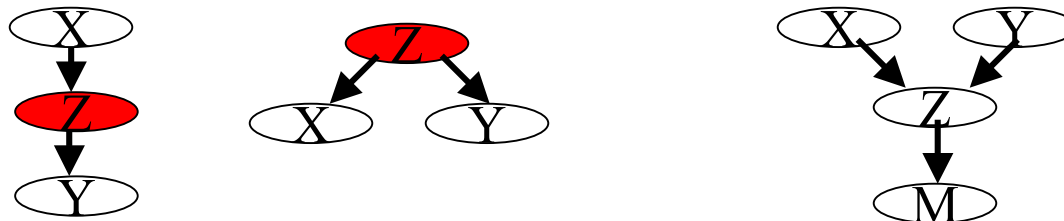
NOT knowing D or M makes L and B independent (common effect)

d-separation

Nodes X and Y are *d-separated* if on *any (undirected) path* between X and Y there is some variable Z such that is either

Z is in a *serial* or *diverging* connection and Z is *known*, or

Z is in a *converging* connection and neither Z nor any of Z 's descendants are known

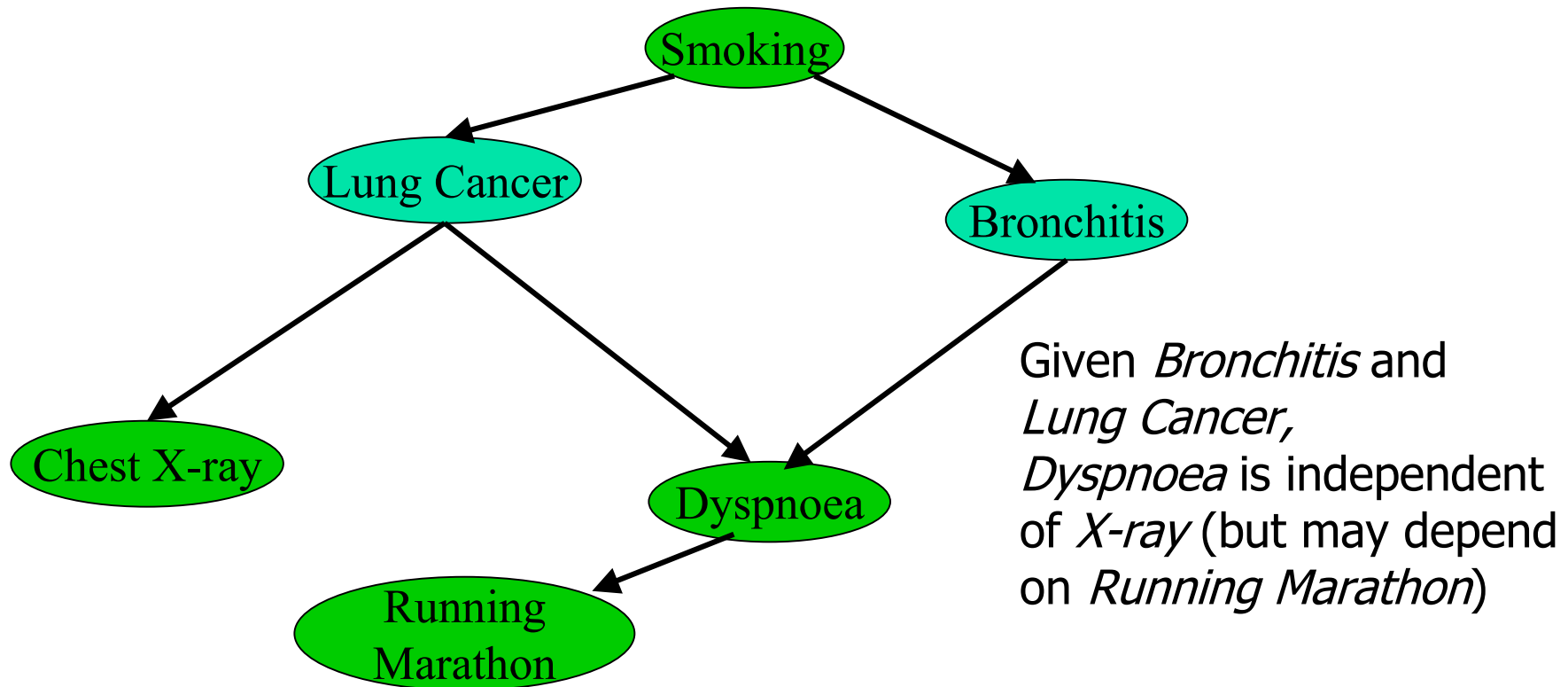


Nodes X and Y are called *d-connected* if they are not d-separated (there exists an undirected path between X and Y not d-separated by any node or a set of nodes)

If nodes X and Y are *d-separated* by Z , then X and Y are *conditionally independent* given Z (see Pearl, 1988)

Independence Relations in BN

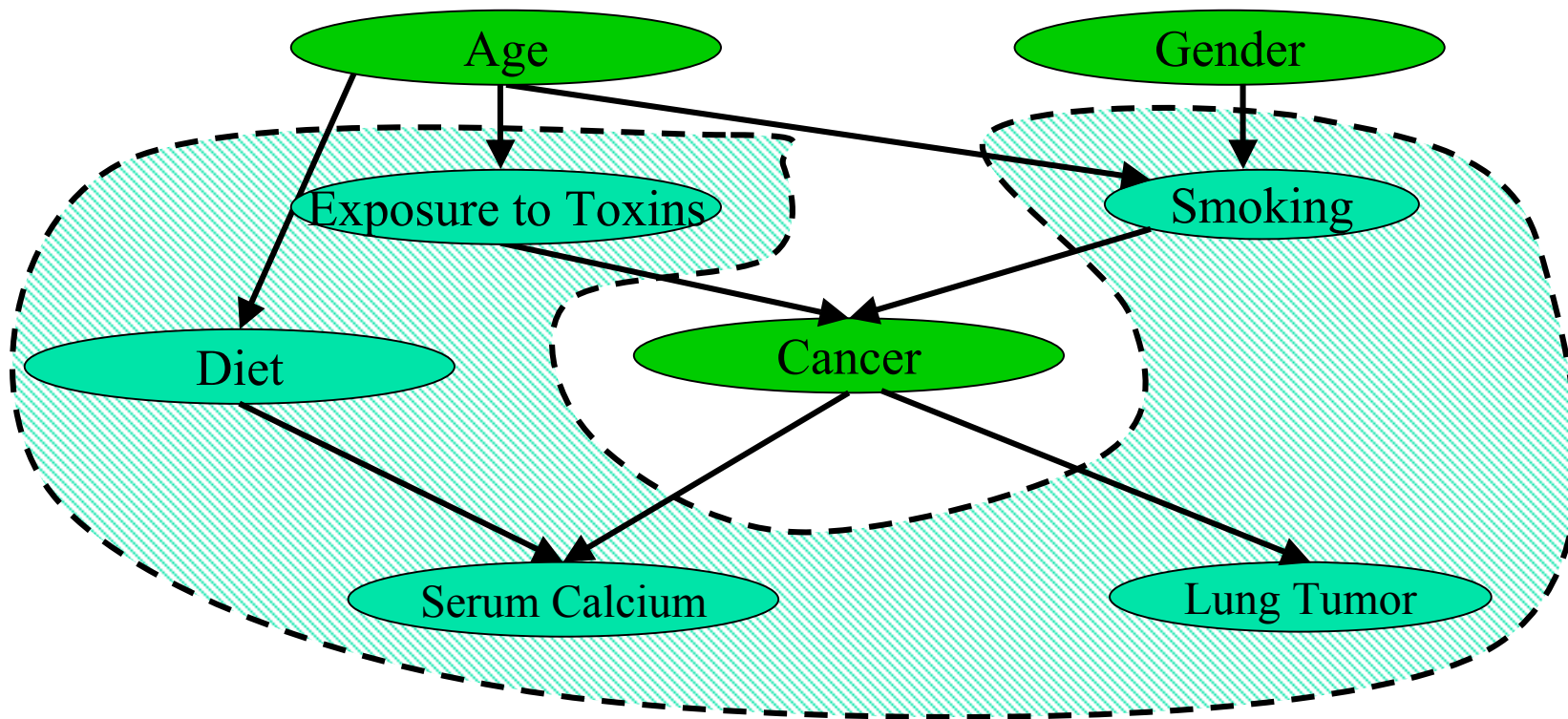
A variable (node) is conditionally independent of its non-descendants given its parents



Markov Blanket

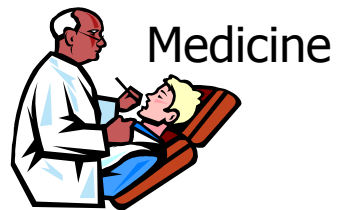
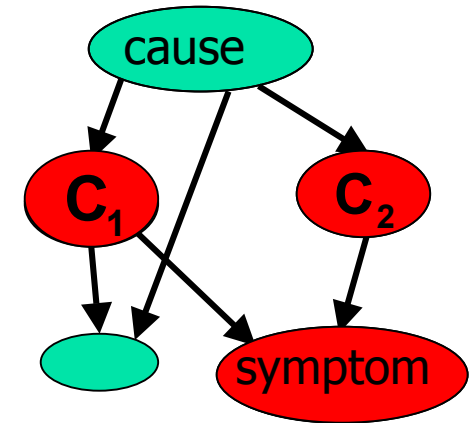
A node is conditionally independent of ALL other nodes given its *Markov blanket*, i.e. its *parents*, *children*, and “*spouses*” (parents of common children)

(Proof left as a homework problem ☺)

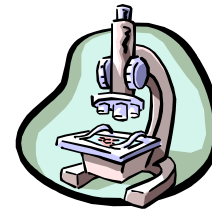


What are BNs useful for?

- Diagnosis: $P(\text{cause}|\text{symptom})=?$
- Prediction: $P(\text{symptom}|\text{cause})=?$
- Classification: $\max_{\text{class}} P(\text{class}|\text{data})$
- Decision-making (given a cost function)



Speech
recognition



Application Examples

APRI system developed at AT&T Bell Labs

learns & uses Bayesian networks from data to identify customers
liable to default on bill payments

NASA Vista system

predict failures in propulsion systems

considers time criticality & suggests highest utility action

dynamically decide what information to show

Application Examples

Office Assistant in MS Office 97/ MS Office 95

Extension of Answer wizard

uses naïve Bayesian networks

help based on past experience (keyboard/mouse use) and task user is doing currently

This is the “smiley face” you get in your MS Office applications

Microsoft Pregnancy and Child-Care

Available on MSN in Health section

Frequently occurring children’s symptoms are linked to expert modules that repeatedly ask parents relevant questions

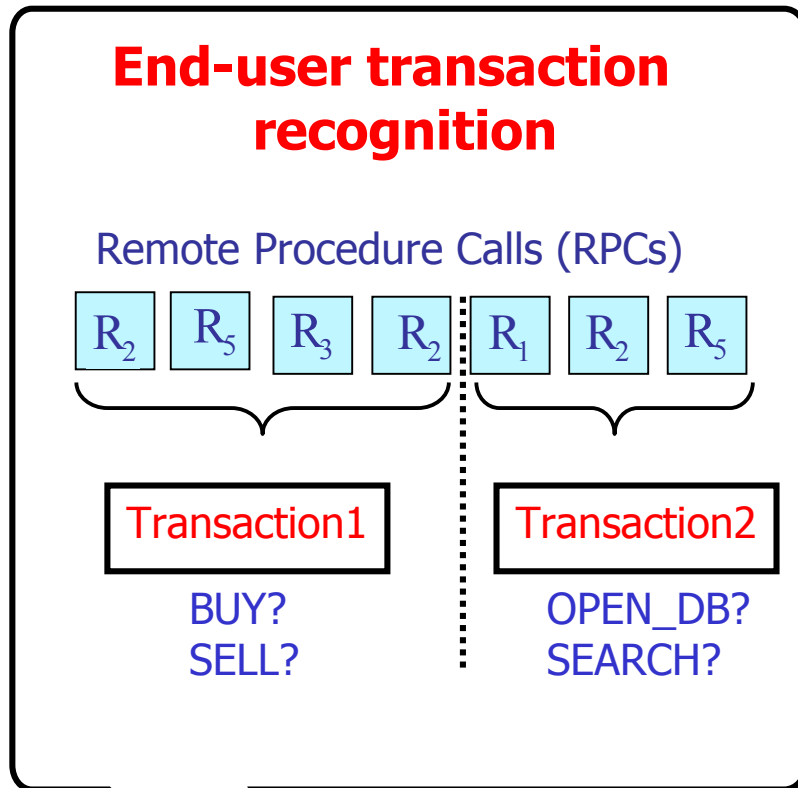
Asks next best question based on provided information

Presents articles that are deemed relevant based on information provided

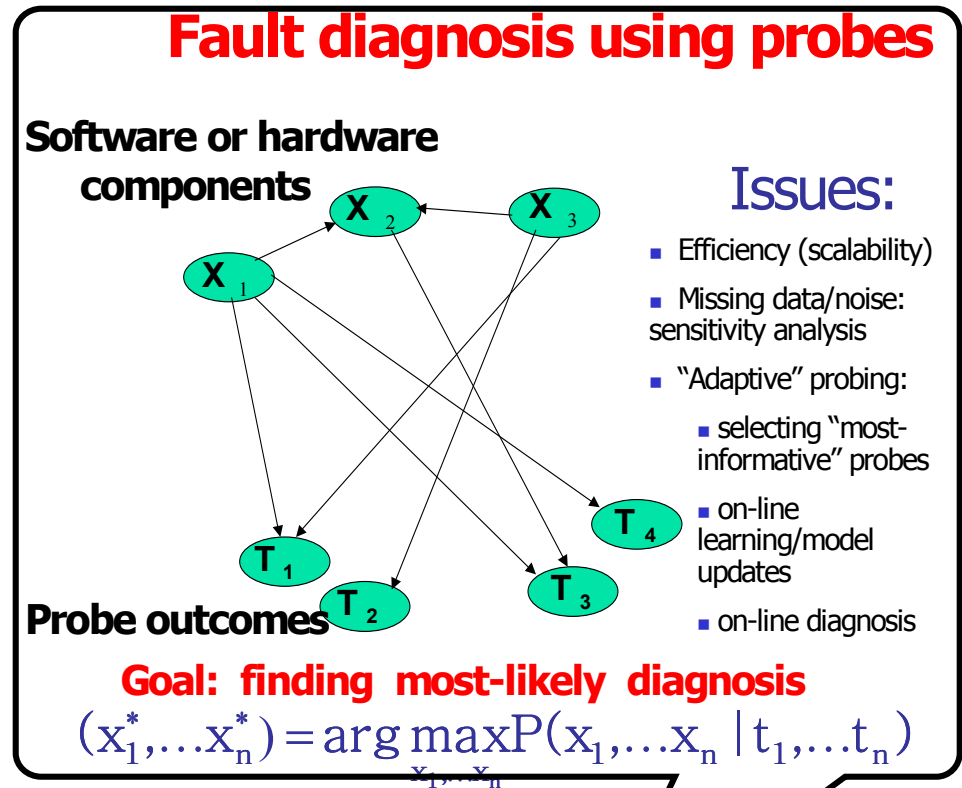
IBM's systems management applications

Machine Learning for Systems @ Watson

(Hellerstein, Jayram, Rish (2000))



(Rish, Brodie, Ma (2001))



Pattern discovery, classification, diagnosis and prediction

Probabilistic Inference Tasks

- Belief updating:

$$\mathbf{BEL}(X_i) = \mathbf{P}(X_i = x_i \mid \mathbf{evidence})$$

- Finding most probable explanation (MPE)

$$\bar{\mathbf{x}}^* = \mathbf{argmax}_{\bar{\mathbf{x}}} \mathbf{P}(\bar{\mathbf{x}}, \mathbf{e})$$

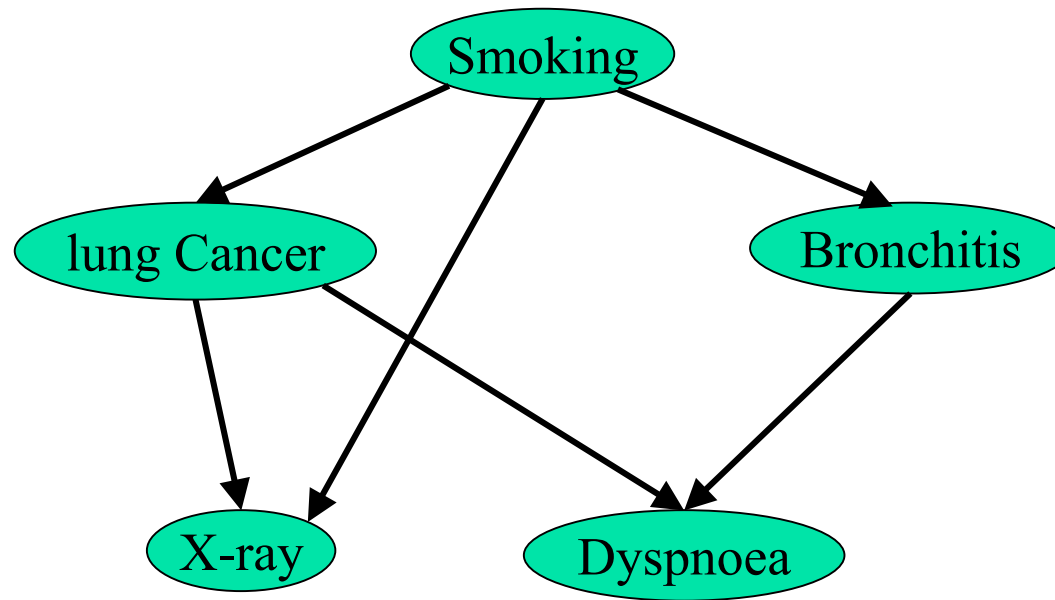
- Finding maximum a-posteriori hypothesis

$$(\mathbf{a}_1^*, \dots, \mathbf{a}_k^*) = \mathbf{argmax}_{\bar{\mathbf{a}}} \sum_{\bar{\mathbf{x}}/A} \mathbf{P}(\bar{\mathbf{x}}, \mathbf{e}) \quad \begin{array}{l} A \subseteq X: \\ \text{hypothesis variables} \end{array}$$

- Finding maximum-expected-utility (MEU) decision

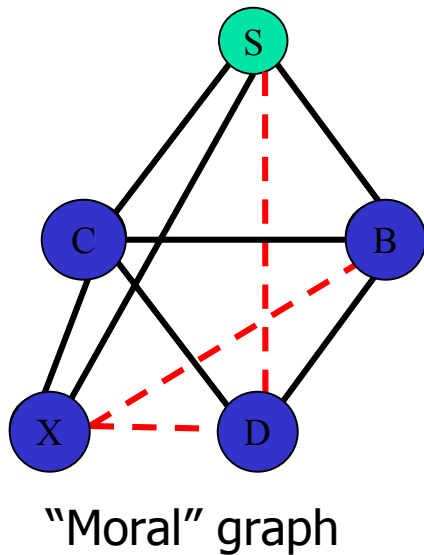
$$(\mathbf{d}_1^*, \dots, \mathbf{d}_k^*) = \mathbf{argmax}_{\bar{\mathbf{d}}} \sum_{\bar{\mathbf{x}}/D} \mathbf{P}(\bar{\mathbf{x}}, \mathbf{e}) \mathbf{U}(\bar{\mathbf{x}}) \quad \begin{array}{l} D \subseteq X: \text{ decision variables} \\ \mathbf{U}(\bar{\mathbf{x}}): \text{ utility function} \end{array}$$

Belief Updating Task: Example



$P(\text{smoking} \mid \text{dyspnoea}=\text{yes}) = ?$

Belief updating: find $P(X|\text{evidence})$



$$P(s|d=1) = \frac{P(s, d=1)}{P(d=1)} \propto P(s, d=1) =$$

$$\sum_{d=1, b, x, c} P(s) \underbrace{P(c|s)} P(b|s) \underbrace{P(x|c, s) P(d|c, b)} =$$

$$P(s) \sum_{d=1} \sum_b P(b|s) \sum_x \underbrace{\sum_c P(c|s) P(x|c, s) P(d|c, b)}$$

Variable Elimination

$$f(s, d, b, x)$$

$W^*=4$

Complexity: $O(n \exp(w^*))$

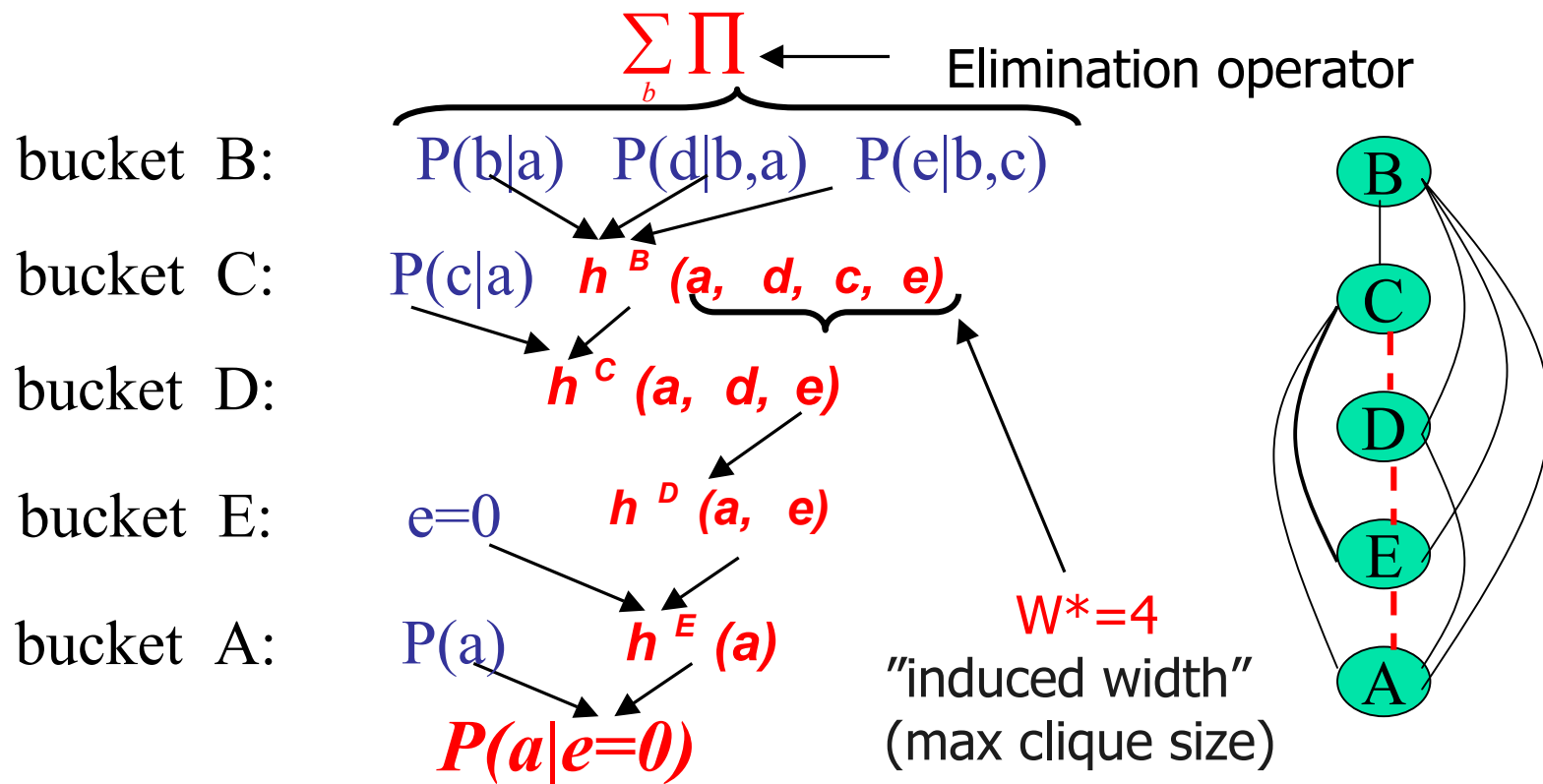
"induced width"
(max induced clique size)

Efficient inference: variable orderings, conditioning, approximations

Variable elimination algorithms

(also called "bucket-elimination")

Belief updating: **VE-algorithm** *elim-bel* (Dechter 1996)

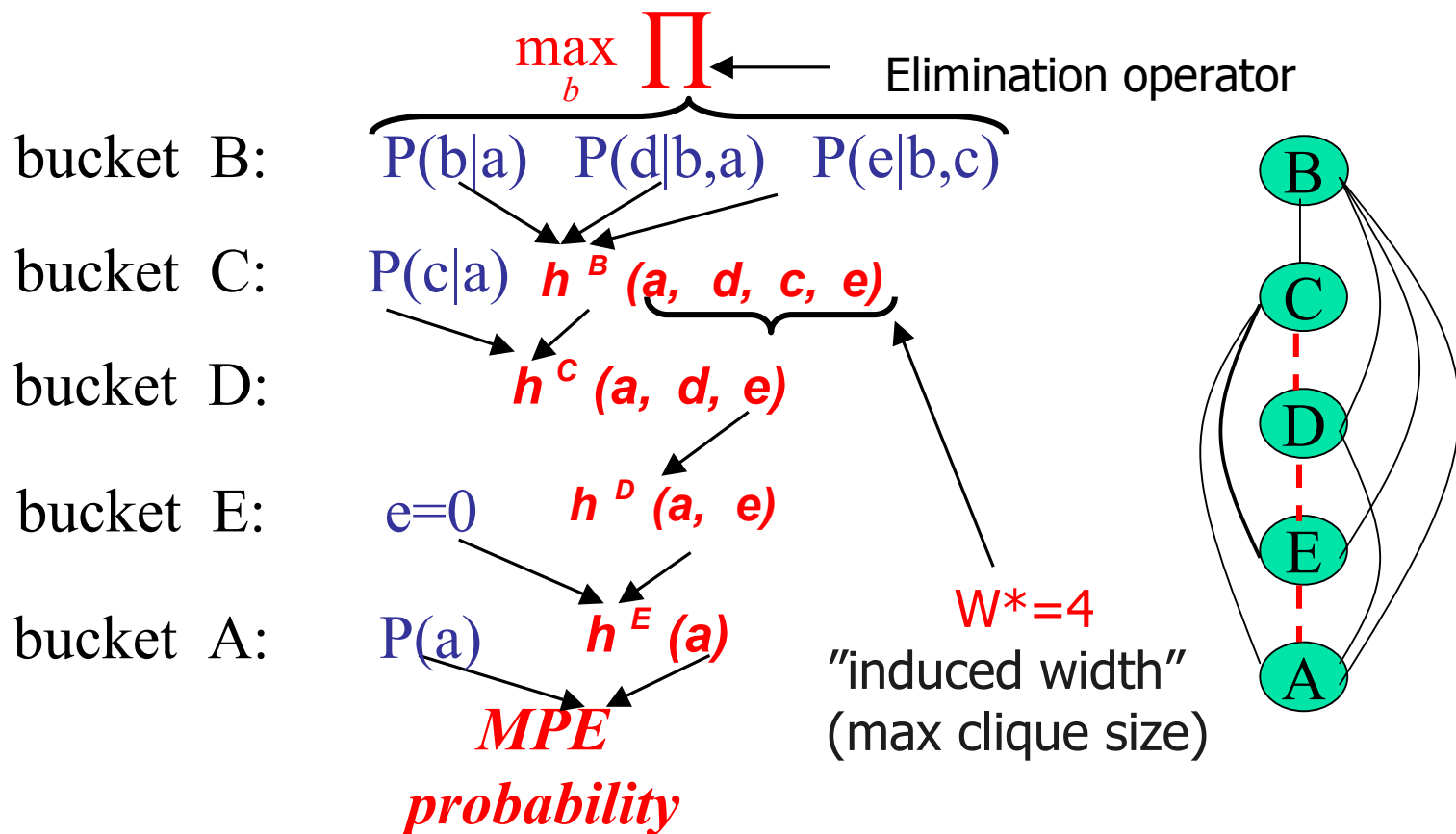


Finding $MPE = \max_{\bar{x}} P(\bar{x})$

VE-algorithm *elim-mpe* (Dechter 1996)

\sum is replaced by *max* :

$$MPE = \max_{a,e,d,c,b} P(a)P(c|a)P(b|a)P(d|a,b)P(e|b,c)$$



Generating the MPE-solution

5. $b' = \arg \max_b P(b | a') \times P(d' | b, a') \times P(e' | b, c')$

4. $c' = \arg \max_c P(c | a') \times h^B(a', d', c, e')$

3. $d' = \arg \max_d h^C(a', d, e')$

2. $e' = 0$

1. $a' = \arg \max_a P(a) \cdot h^E(a)$



B: $P(b|a) \quad P(d|b,a) \quad P(e|b,c)$

C: $P(c|a) \quad h^B(a, d, c, e)$

D: $h^C(a, d, e)$

E: $e=0 \quad h^D(a, e)$

A: $P(a) \quad h^E(a)$

Return (a', b', c', d', e')

Complexity of VE-inference: $O(n \exp(w_o^*))$

w_o^* – the induced width of moral graph along ordering o

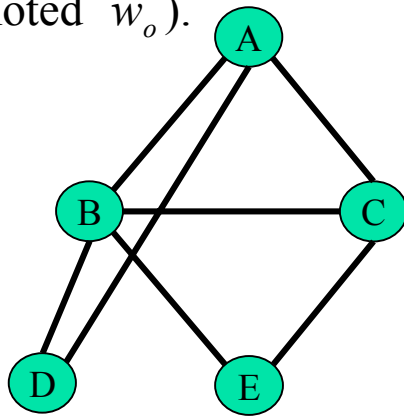
Meaning : $w_o^* + 1 =$ size of largest clique created during inference

The width $w_o(X)$ of a variable X in graph G along the ordering o is the number of nodes preceding X in the ordering and connected to X (*earlier neighbors*).

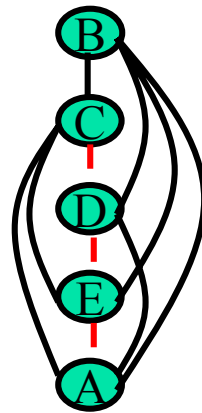
The width w_o of a graph is the maximum width $w_o(X)$ among all nodes.

The *induced graph* G' along the ordering o is obtained by recursively connecting earlier neighbors of each node, from last to the first in the ordering.

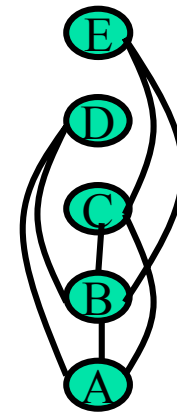
The width of the induced graph G' is called the *induced width* of the graph G (denoted w_o^*).



"Moral" graph



$$w_{o_1}^* = 4$$

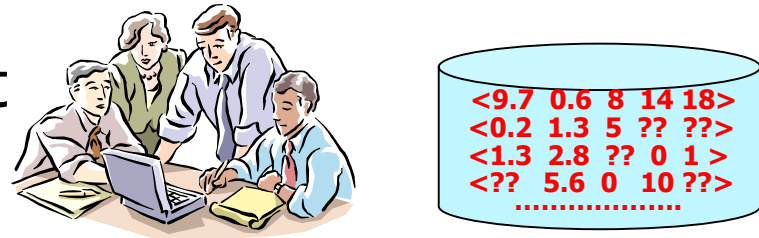


$$w_{o_2}^* = 2$$

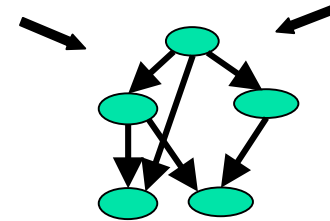
Ordering is important! But finding min- w^* ordering is NP-hard... Inference is also NP-hard in general case [Cooper].

Learning Bayesian Networks

- **Combining** domain expert knowledge with data



- Efficient **representation** and **inference**



- **Incremental** learning: $P(H) \nearrow$ or \searrow

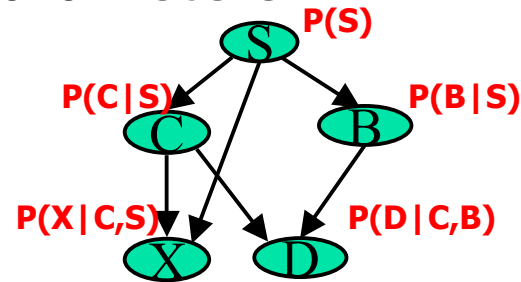
- Handling **missing data**: $\langle 1.3 \ 2.8 \ ?? \ 0 \ 1 \rangle$

- Learning **causal** relationships: $S \rightarrow C$

Learning tasks: four main cases

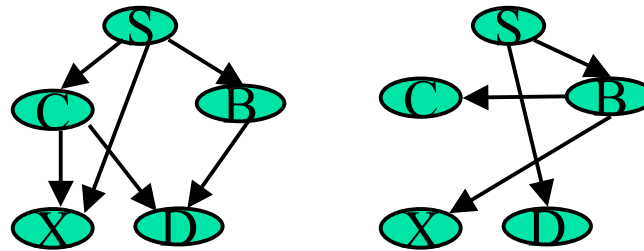
■ Known graph – learn parameters

- Complete data:
parameter estimation (ML, MAP)
- Incomplete data:
non-linear parametric
optimization (gradient descent, EM)



■ Unknown graph – learn graph and parameters

- Complete data:
optimization (search
in space of graphs)
- Incomplete data:
structural EM,
mixture models

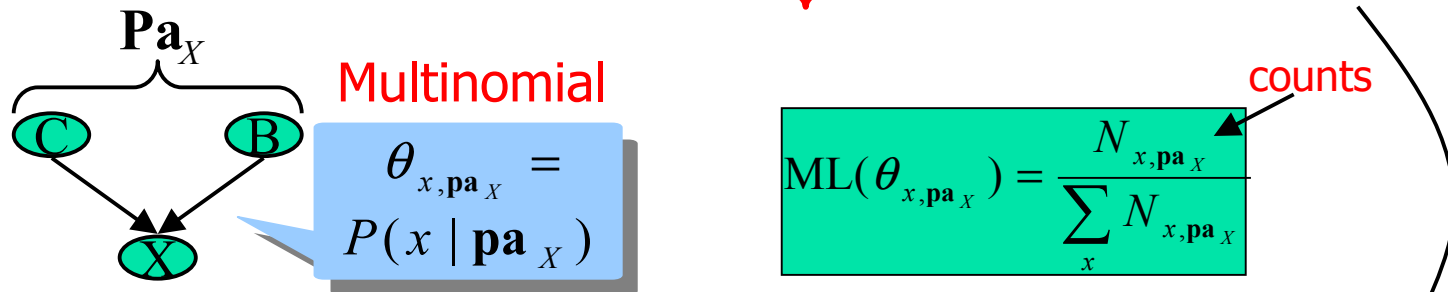


$$\hat{G} = \arg \max_G \text{Score}(G)$$

Learning Parameters: complete data

(overview)

- ML-estimate: $\max_{\Theta} \log P(D | \Theta)$ - decomposable!



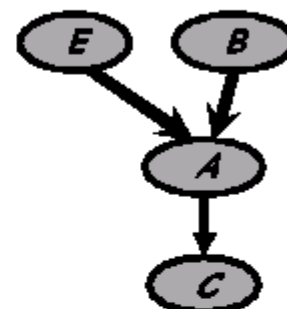
- MAP-estimate (Bayesian statistics) $\max_{\Theta} \log P(D | \Theta) P(\Theta)$

Conjugate priors - **Dirichlet** $\text{Dir}(\theta_{\text{pa}_X} | \alpha_{1, \text{pa}_X}, \dots, \alpha_{m, \text{pa}_X})$

$$\text{MAP}(\theta_{x, \text{pa}_X}) = \frac{N_{x, \text{pa}_X} + \alpha_{x, \text{pa}_X}}{\sum_x N_{x, \text{pa}_X} + \underbrace{\sum_x \alpha_{x, \text{pa}_X}}_{\text{Equivalent sample size (prior knowledge)}}$$

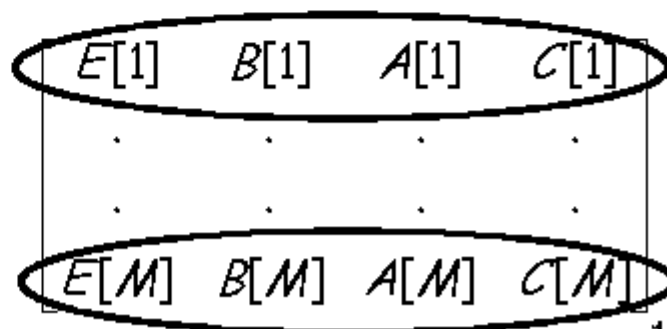
Likelihood Function

◆ By definition of network, we get



$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$

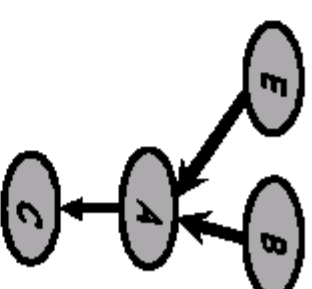
$$= \prod_m \left(\begin{array}{l} P(E[m] : \Theta) \\ P(B[m] : \Theta) \\ P(A[m] | B[m], E[m] : \Theta) \\ P(C[m] | A[m] : \Theta) \end{array} \right)$$



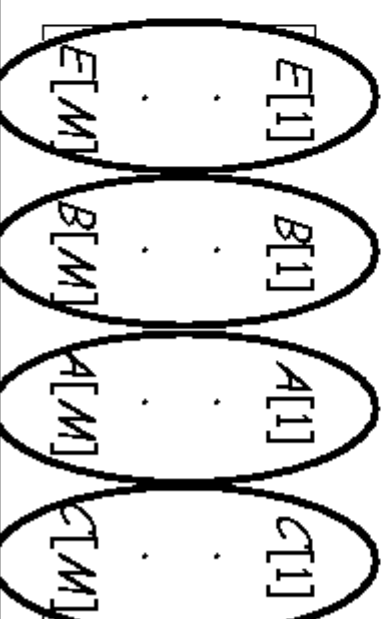
Likelihood Function

- ◆ Rewriting terms, we get

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$



$$\begin{aligned} & \prod_m P(E[m] : \Theta) \\ & \prod_m P(B[m] : \Theta) \\ & = \prod_m P(A[m] | B[m], E[m] : \Theta) \\ & \prod_m P(C[m] | A[m] : \Theta) \end{aligned}$$



General Bayesian Networks

Generalizing for any Bayesian network:

$$\begin{aligned} \mathcal{L}(\Theta : D) &= \prod_m P(x_1[m], \dots, x_n[m] : \Theta) \\ &= \prod_m \prod_i P(x_i[m] \mid Pa_i[m] : \Theta_i) \\ &= \prod_i \mathcal{L}_i(\Theta_i : D) \end{aligned}$$

Decomposition

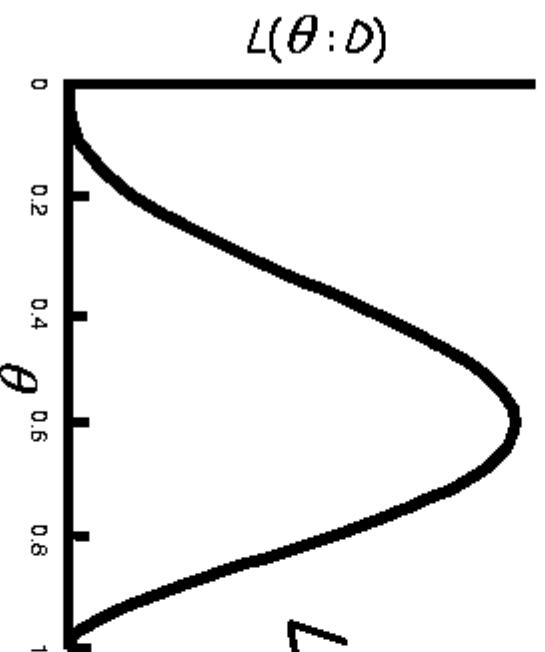
⇒ **Independent estimation problems**

Likelihood Function: Multinomials

$$L(\theta : D) = P(D | \theta) = \prod_m P(x[m] | \theta)$$

- ◆ The likelihood for the sequence H, T, T, H, H is

$$L(\theta : D) = \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta$$



Count of k^{th} outcome in D

Probability of k^{th} outcome

General case:
$$L(\theta : D) = \prod_{k=1}^K \theta_k^{N_k}$$

Dirichlet Priors

- ◆ Recall that the likelihood function is

$$L(\Theta : \mathcal{D}) = \prod_{k=1}^K \theta_k^{N_k}$$

- ◆ **Dirichlet** prior with hyperparameters $\alpha_1, \dots, \alpha_K$

$$P(\Theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

\Rightarrow the posterior has the same form, with

hyperparameters $\alpha_1 + N_1, \dots, \alpha_K + N_K$

$$P(\Theta | \mathcal{D}) \propto P(\Theta)P(\mathcal{D} | \Theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1} \prod_{k=1}^K \theta_k^{N_k} = \prod_{k=1}^K \theta_k^{\alpha_k + N_k - 1}$$

Learning Parameters: Summary

- ◆ Estimation relies on **sufficient statistics**
 - For multinomials: counts $\mathcal{N}(x_i, p_{q_j})$
 - Parameter estimation

$$\hat{\theta}_{x_i|p_{q_j}} = \frac{\mathcal{N}(x_i, p_{q_j})}{\mathcal{N}(p_{q_j})} \qquad \tilde{\theta}_{x_i|p_{q_j}} = \frac{\alpha(x_i, p_{q_j}) + \mathcal{N}(x_i, p_{q_j})}{\alpha(p_{q_j}) + \mathcal{N}(p_{q_j})}$$

MLE

Bayesian (Dirichlet)

- ◆ Both are asymptotically equivalent and consistent
- ◆ Both can be implemented in an on-line manner by accumulating sufficient statistics

Learning Parameters: incomplete data

Non-decomposable marginal likelihood (hidden nodes)

Initial parameters

Current model
 (G, Θ)

Expectation

Compute EXPECTED
Counts via inference in BN

Data

S	X	D	C	B
<?>	0	1	0	1
<1>	1	?>	0	1
<0>	0	0	?>	?>
<?>	?>	0	?>	1

Expected counts

Maximization

Update parameters
(ML, MAP)

EM-algorithm:
iterate until convergence

$$E_{P(X)} [N_{x, \text{pa}_x}] = \sum_{k=1}^N p(x, \text{pa}_x | y^k, \Theta, G)$$

Learning graph structure

$$\text{Find } \hat{G} = \arg \max_G \text{Score}(G)$$

NP-hard optimization

■ Heuristic search:

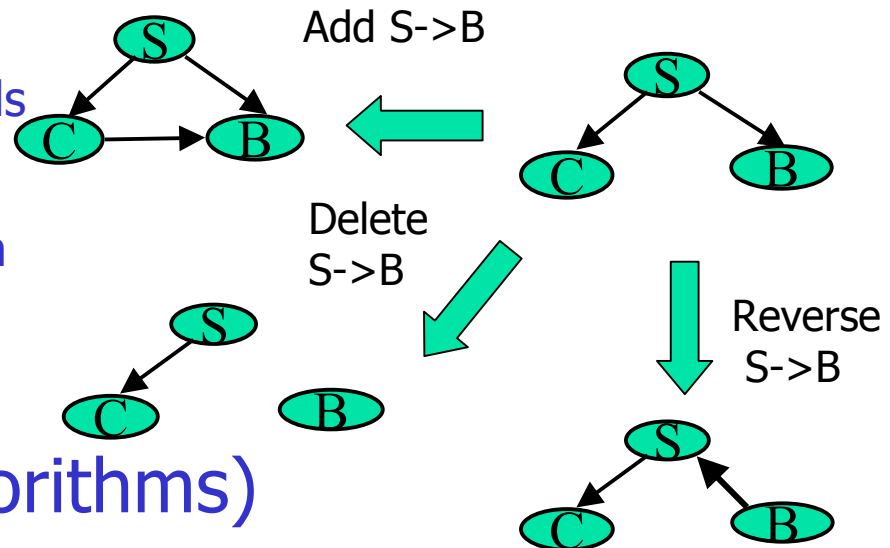
Complete data – local computations

Incomplete data (score non-decomposable): stochastic methods

Local greedy search; K2 algorithm

■ Constrained-based methods (PC/IC algorithms)

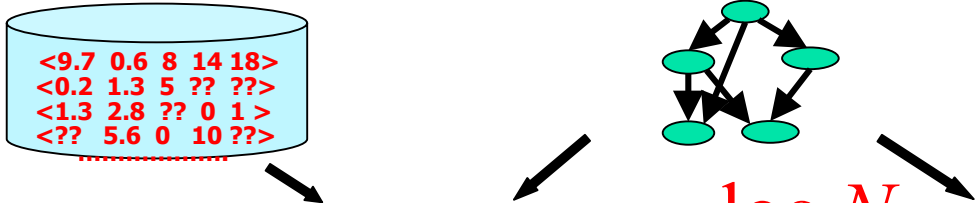
➤ Data impose independence relations (constraints) on graph structure



Scoring function:

Minimum Description Length (MDL)

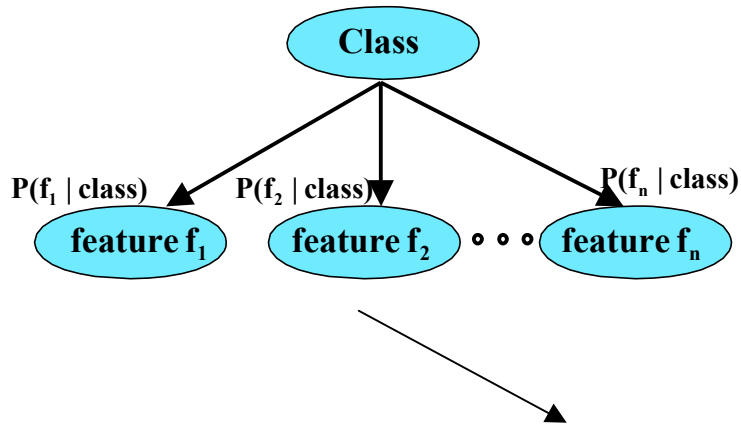
- Learning \Leftrightarrow data compression


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

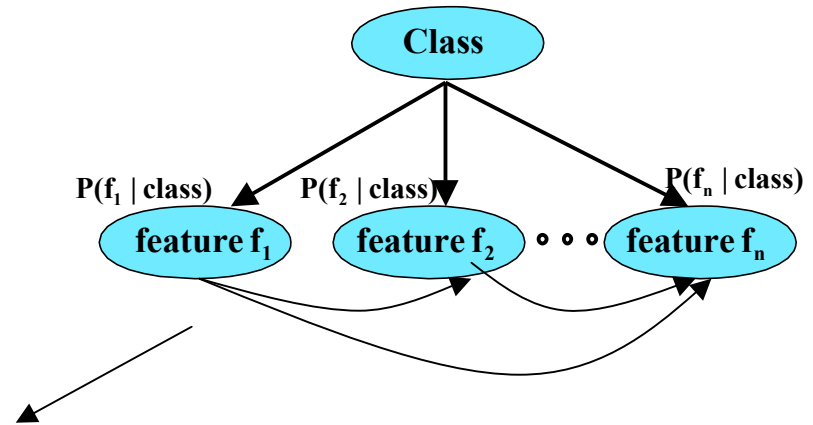
- Other: MDL = -BIC (Bayesian Information Criterion)
- Bayesian score (BDe) - asymptotically equivalent to MDL

Model selection trade-offs

Naïve Bayes – too simple
(less parameters, but bad model)



Unrestricted BN – too complex
(possible overfitting + complexity)

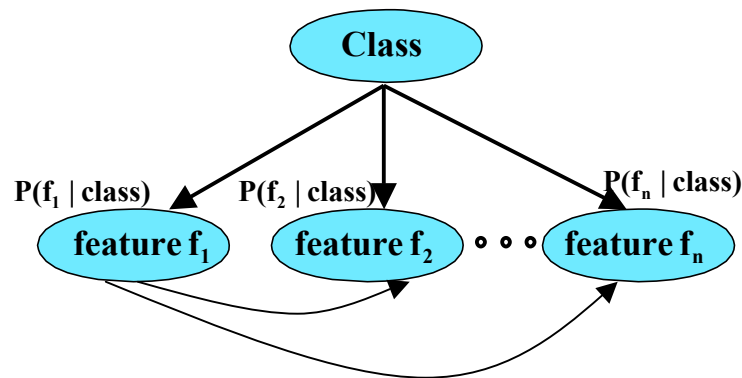


Various approximations between the two extremes

TAN:

tree-augmented Naïve Bayes
[Friedman et al. 1997]

Based on Chow-Liu Tree Method
(CL) for learning trees
[Chow-Liu, 1968]

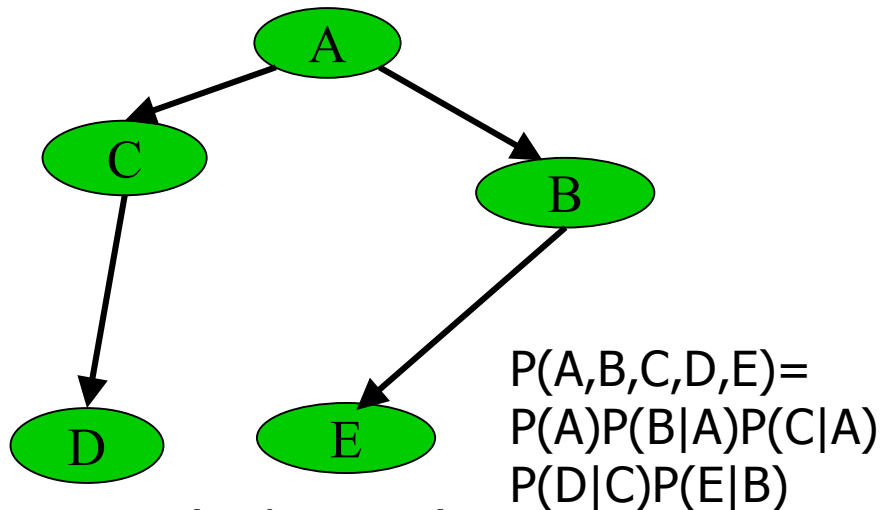


Tree-structured distributions

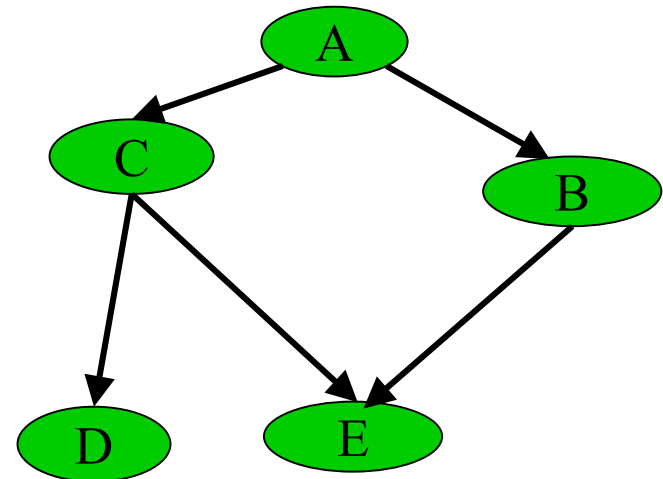
A joint probability distribution is tree-structured if it can be written as

$$P(\mathbf{X}) = \prod_{i=1}^n P(x_i | x_{j(i)})$$

where $x_{j(i)}$ is the parent of x_i in Bayesian network for $P(\mathbf{x})$ (a directed tree)



A tree (with root A)

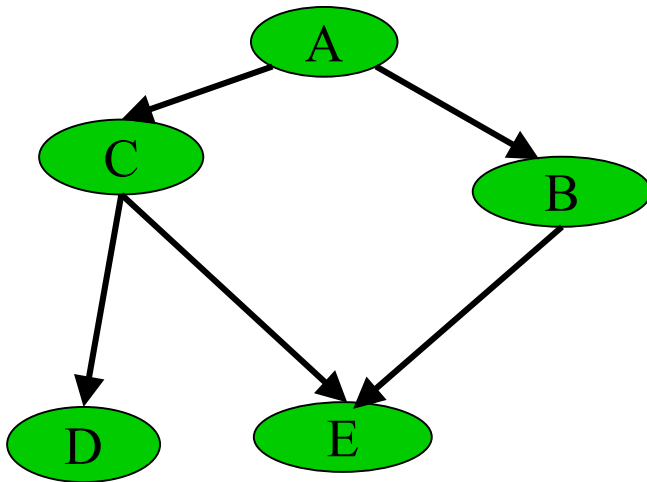


Not a tree – has an (undirected) cycle

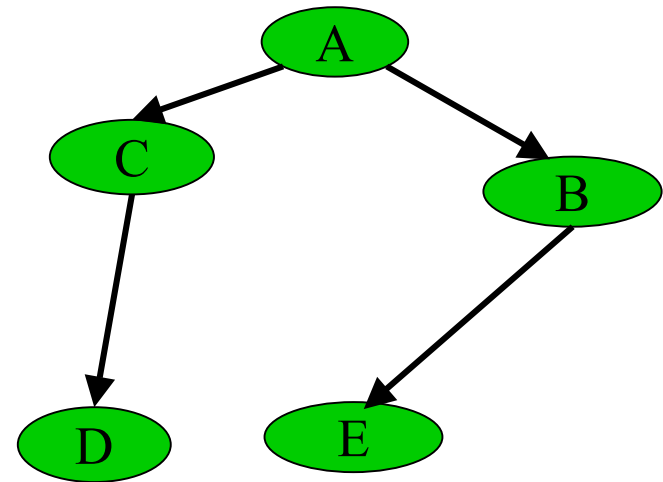
A tree requires only **$[(d-1) + d(d-1)(n-1)]$** parameters, where d is domain size
Moreover, inference in trees is $O(n)$ (linear) since their $w^*=1$

Approximations by trees

True distribution $P(X)$



Tree-approximation $P'(X)$



How good is approximation? Use cross-entropy (KL-divergence):

$$D(P, P') = P \sum_{\mathbf{x}} P(\mathbf{x}) \log \frac{P(\mathbf{x})}{P'(\mathbf{x})}$$

$D(P, P')$ is non-negative, and $D(P, P')=0$ if and only if P coincides with P' (on a set of measure 1)

How to find the best tree-approximation?

Optimal trees: Chow-Liu result

- **Lemma**

Given a joint PDF $P(x)$ and a fixed tree structure T , the best approximation $P'(x)$ (i.e., $P'(x)$ that minimizes $D(P, P')$) satisfies

$$P'(x_i | x_{j(i)}) = P(x_i | x_{j(i)}) \quad \text{for all } i = 1, \dots, n$$

Such $P'(x)$ is called the projection of $P(x)$ on T .

- **Theorem** [Chow and Liu, 1968]

Given a joint PDF $P(x)$, the KL-divergence $D(P, P')$ is minimized by projecting $P(x)$ on a *maximum-weight spanning tree (MSWT)* over nodes in X , where the weight on the edge (X_i, X_j) is defined by the mutual information measure

$$I(X_i; X_j) = \sum_{x_i, x_j} P(x_i, x_j) \log \frac{P(x_i, x_j)}{P(x_i)P(x_j)}$$

Note, that $I(X; Y) = 0$ when X and Y are independent and that $I(X; Y) = D(P(x, y), P(x)P(y))$

Proofs

Proof of Lemma :

$$\begin{aligned}
 D(P, P') &= -\sum_{\mathbf{x}} P(\mathbf{x}) \sum_{i=1}^n \log P'(x_i | x_{j(i)}) + \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x}) = -\sum_{\mathbf{x}} P(\mathbf{x}) \sum_{i=1}^n \log P'(x_i | x_{j(i)}) - H(X) = \\
 &= -\sum_{i=1}^n \sum_{\mathbf{x}} P(\mathbf{x}) \log P'(x_i | x_{j(i)}) - H(X) = -\sum_{i=1}^n \sum_{x_i, x_{j(i)}} P(x_i, x_{j(i)}) \log P'(x_i | x_{j(i)}) - H(X) = \quad (1)
 \end{aligned}$$

$$= -\sum_{i=1}^n \sum_{x_{j(i)}} P(x_{j(i)}) \sum_{x_i} P(x_i | x_{j(i)}) \log P'(x_i | x_{j(i)}) - H(X) \quad (2)$$

A known fact : given $P(x)$, the maximum of $\sum_{x_i} P(x) \log P'(x)$ is achieved by the choice $P'(x) = P(x)$.

Therefore, for any value of i and $x_{j(i)}$, the term $\sum_{x_i} P(x_i | x_{j(i)}) \log P'(x_i | x_{j(i)})$ is maximized by

choosing $P'(x_i | x_{j(i)}) = P(x_i | x_{j(i)})$ (and thus the total $D(P, P')$ is minimized), which proves the Lemma.

Proof of Theorem :

Replacing $P'(x_i | x_{j(i)}) = P(x_i | x_{j(i)})$ in the expression (1) yields

$$\begin{aligned}
 D(P, P') &= -\sum_{i=1}^n \sum_{x_i, x_{j(i)}} P(x_i, x_{j(i)}) \log [P(x_i x_{j(i)}) / P(x_{j(i)})] - H(X) = \\
 &= -\sum_{i=1}^n \sum_{x_i, x_{j(i)}} P(x_i, x_{j(i)}) \left[\log \frac{P(x_i x_{j(i)})}{P(x_{j(i)}) P(x_i)} + \log P(x_i) \right] - H(X) = \\
 &= -\sum_{i=1}^n I(X_i, X_{j(i)}) - \sum_{i=1}^n \sum_{x_i} P(x_i) \log P(x_i) - H(X).
 \end{aligned}$$

The last two terms are independent of the choice of the tree, and thus $D(P, P')$ is minimized by maximizing the sum of edge weights $I(X_i, X_{j(i)})$.

Chow-Liu algorithm

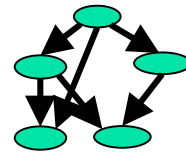
[As presented in Pearl, 1988]

1. From the given distribution $P(x)$ (or from data generated by $P(x)$), compute the joint distribution $P(x_i | x_j)$ for all $i \neq j$
2. Using the pairwise distributions from step 1, compute the mutual information $I(X_i; X_j)$ for each pair of nodes and assign it as the weight to the corresponding edge (X_i, X_j) .
3. Compute the maximum-weight spanning tree (MSWT):
 - a. Start from the empty tree over n variables
 - b. Insert the two largest-weight edges
 - c. Find the next largest-weight edge and add it to the tree if no cycle is formed; otherwise, discard the edge and repeat this step.
 - d. Repeat step (c) until $n-1$ edges have been selected (a tree is constructed).
4. Select an arbitrary root node, and direct the edges outwards from the root.
5. Tree approximation $P'(x)$ can be computed as a projection of $P(x)$ on the resulting directed tree (using the product-form of $P'(x)$).

Summary:

Learning and inference in BNs

- **Bayesian Networks** – graphical probabilistic models
- Efficient **representation** and **inference**
- Expert **knowledge** + learning from **data**
- Learning:
 - **parameters** (parameter estimation, EM)
 - **structure** (optimization w/ **score** functions – e.g., **MDL**)
 - **Complexity trade-off:**
 - **NB, BNs and trees**
- There is much more: causality, modeling time (DBNs, HMMs), approximate inference, on-line learning, active learning, etc.



Online/print resources on BNs

Conferences & Journals

UAI, ICML, AAI, AISTAT, KDD

MLJ, DM&KD, JAIR, IEEE KDD, IJAR, IEEE PAMI

Books and Papers

Bayesian Networks without Tears by Eugene Charniak. AI Magazine: Winter 1991.

Probabilistic Reasoning in Intelligent Systems by Judea Pearl. Morgan Kaufmann: 1998.

Probabilistic Reasoning in Expert Systems by Richard Neapolitan. Wiley: 1990.

CACM special issue on Real-world applications of BNs, March 1995

Online/Print Resources on BNs

AUAI online: www.auai.org. Links to:

Electronic proceedings for UAI conferences

Other sites with information on BNs and reasoning under uncertainty

Several tutorials and important articles

Research groups & companies working in this area

Other societies, mailing lists and conferences

Publicly available s/w for BNs

List of BN software maintained by Russell Almond at
bayes.stat.washington.edu/almond/belief.html

several free packages: generally research only

commercial packages: most powerful (& expensive) is
HUGIN; others include Netica and Dxpess

we are working on developing a Java based BN toolkit here at
Watson