

**CS 3710 Advanced Topics in AI
Lecture 3**

**Probabilistic graphical
models**

Milos Hauskrecht

milos@cs.pitt.edu

5329 Sennott Square

CS 3710 Probabilistic graphical models

Modeling uncertainty with probabilities

- **Representing large multivariate distributions directly and exhaustively is hopeless:**
 - The number of parameters is exponential in the number of random variables
 - Inference can be exponential in the number of variables
- Breakthrough (late 80s, beginning of 90s)
 - **Bayesian belief networks**
 - Give solutions to the space, acquisition bottlenecks
 - Partial solutions for time complexities

CS 3710 Probabilistic graphical models

Graphical models

Aim: alleviate the representational and computational bottlenecks

Idea: Take advantage of the structure, more specifically, **independences and conditional independences** that hold among random variables

Two classes of models:

- **Bayesian belief networks**
 - Modeling asymmetric (causal) effects and dependencies
- **Markov random fields**
 - Modeling symmetric effects and dependencies among random variables
 - Used often to model spatial dependences (image analysis)

Bayesian belief networks (BBNs)

Bayesian belief networks.

- Represent the full joint distribution over the variables more compactly using a **smaller number of parameters**.
- Take advantage of **conditional and marginal independences** among random variables

- **A and B are independent**

$$P(A, B) = P(A)P(B)$$

- **A and B are conditionally independent given C**

$$P(A, B | C) = P(A | C)P(B | C)$$

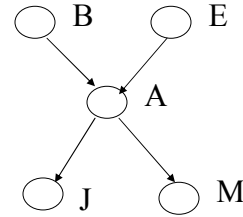
$$P(A | C, B) = P(A | C)$$

Bayesian belief networks (general)

Two components: $B = (S, \Theta_S)$

- Directed acyclic graph**

- Nodes correspond to random variables
- (Missing) links encode independences



- Parameters**

- Local conditional probability distributions for every variable-parent configuration

$$P(X_i | pa(X_i))$$

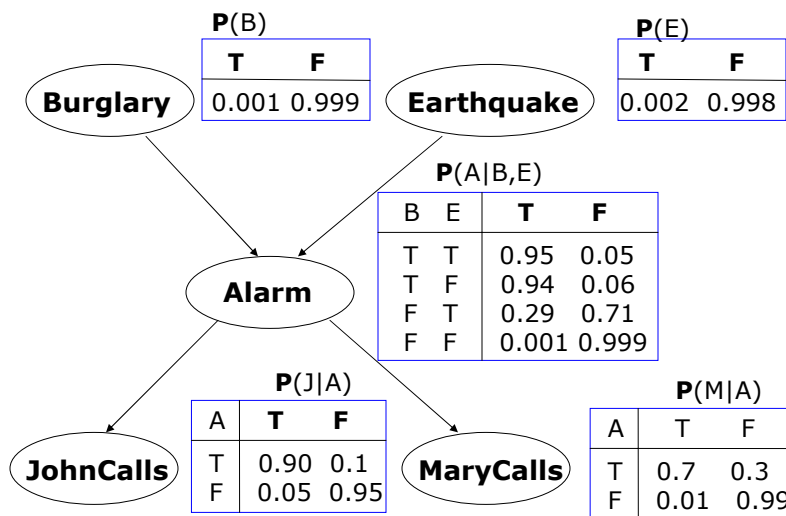
Where:

$pa(X_i)$ - stand for parents of X_i

$P(A|B,E)$

B	E	T	F
T	T	0.95	0.05
T	F	0.94	0.06
F	T	0.29	0.71
F	F	0.001	0.999

Bayesian belief network.



Full joint distribution in BBNs

Full joint distribution is defined in terms of local conditional distributions (obtained via the chain rule):

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

Example:

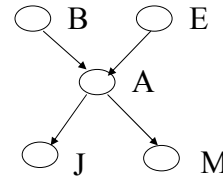
Assume the following assignment of values to random variables

$$B=T, E=T, A=T, J=T, M=F$$

Then its probability is:

$$P(B=T, E=T, A=T, J=T, M=F) =$$

$$P(B=T)P(E=T)P(A=T \mid B=T, E=T)P(J=T \mid A=T)P(M=F \mid A=T)$$



Bayesian belief networks (BBNs)

Bayesian belief networks

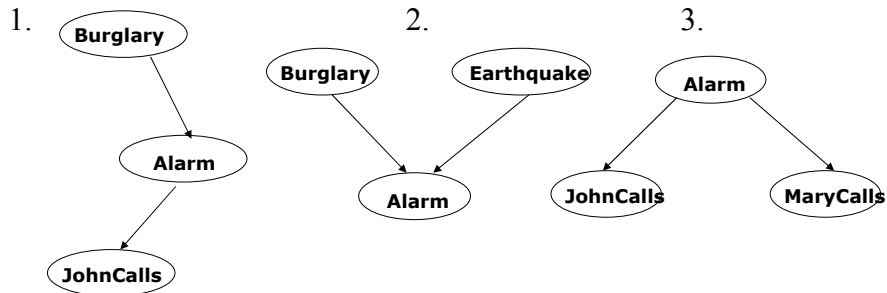
- Represent the full joint distribution over the variables more compactly using the product of local conditionals.
- **But how did we get to local parameterizations?**

Answer:

- **Graphical structure** encodes **conditional and marginal independences** among random variables
- **A and B are independent** $P(A, B) = P(A)P(B)$
- **A and B are conditionally independent given C**
$$P(A \mid C, B) = P(A \mid C)$$
$$P(A, B \mid C) = P(A \mid C)P(B \mid C)$$
- **The graph structure implies the decomposition !!!**

Independences in BBNs

3 basic independence structures:



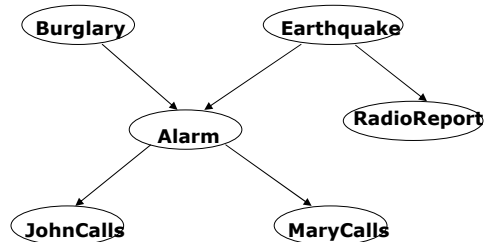
CS 3710 Probabilistic graphical models

Independences in BBN

- BBN distribution models many conditional independence relations among distant variables and sets of variables
- These are defined in terms of the graphical criterion called d-separation
- **D-separation and independence**
 - Let X, Y and Z be three sets of nodes
 - If X and Y are d-separated by Z , then X and Y are conditionally independent given Z
- **D-separation :**
 - A is d-separated from B given C if every undirected path between them is **blocked with C**
- **Path blocking**
 - 3 cases that expand on three basic independence structures

CS 3710 Probabilistic graphical models

Independences in BBNs



- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **F**
- Burglary and RadioReport are independent given Earthquake **T**
- Burglary and RadioReport are independent given MaryCalls **F**

CS 3710 Probabilistic graphical models

Bayesian belief networks (BBNs)

Bayesian belief networks

- Represents the full joint distribution over the variables more compactly using the product of local conditionals.
- **So how did we get to local parameterizations?**

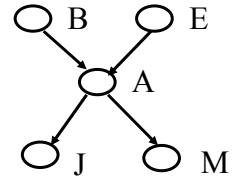
$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

- **The decomposition is implied by the set of independences encoded in the belief network.**

CS 3710 Probabilistic graphical models

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



$$P(B=T, E=T, A=T, J=T, M=F) =$$

$$= P(J=T | B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F)$$

$$= \underline{P(J=T | A=T)} \underline{P(B=T, E=T, A=T, M=F)}$$

$$P(M=F | B=T, E=T, A=T) P(B=T, E=T, A=T)$$

$$\underline{P(M=F | A=T)} \underline{P(B=T, E=T, A=T)}$$

$$\underline{P(A=T | B=T, E=T)} \underline{P(B=T, E=T)}$$

$$\underline{P(B=T)} \underline{P(E=T)}$$

$$= P(J=T | A=T) P(M=F | A=T) P(A=T | B=T, E=T) P(B=T) P(E=T)$$

CS 3710 Probabilistic graphical models

Parameter complexity problem

- In the BBN the **full joint distribution** is defined as:

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i | pa(X_i))$$

- What did we save?**

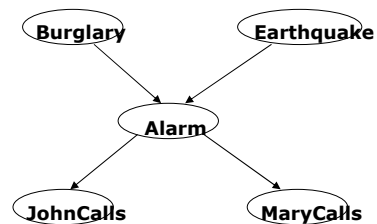
Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

$$2^5 - 1 = 31$$



CS 3710 Probabilistic graphical models

Parameter complexity problem

- In the BBN the **full joint distribution** is defined as:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} P(X_i \mid pa(X_i))$$

- What did we save?

Alarm example: 5 binary (True, False) variables

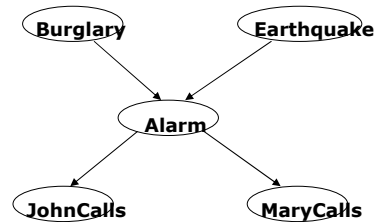
of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

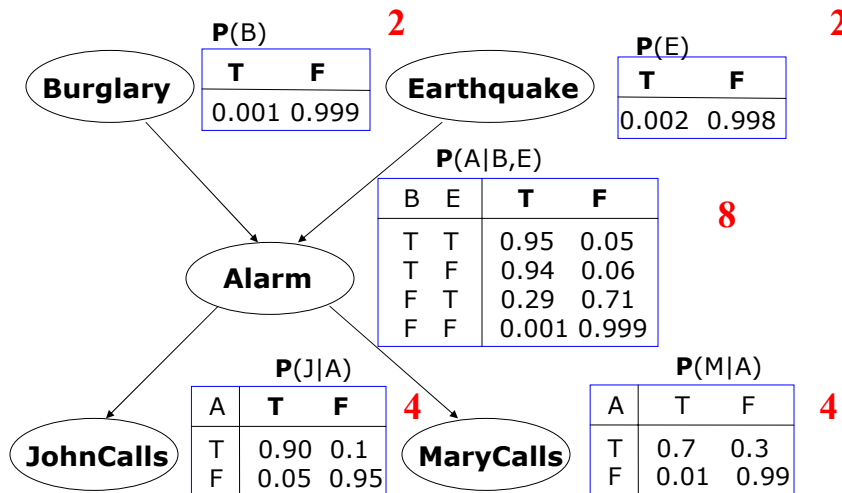
$$2^5 - 1 = 31$$

of parameters of the BBN: ?



Bayesian belief network.

- In the BBN the **full joint distribution** is expressed using a set of local conditional distributions



Parameter complexity problem

- In the BBN the **full joint distribution** is defined as:

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

- What did we save?**

Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

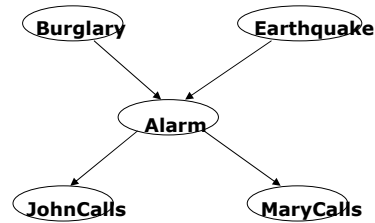
$$2^5 - 1 = 31$$

of parameters of the BBN:

$$2^3 + 2(2^2) + 2(2) = 20$$

One parameter in every conditional is for free:

?



Parameter complexity problem

- In the BBN the **full joint distribution** is defined as:

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

- What did we save?**

Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

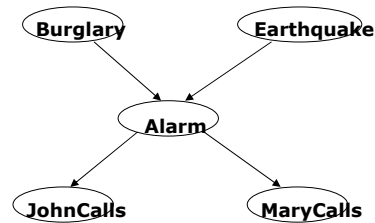
$$2^5 - 1 = 31$$

of parameters of the BBN:

$$2^3 + 2(2^2) + 2(2) = 20$$

One parameter in every conditional is for free:

$$2^2 + 2(2) + 2(1) = 10$$



Model acquisition problem

The structure of the BBN

- typically reflects causal relations
(BBNs are also sometime referred to as **causal networks**)
- Causal structure is intuitive in many applications domain and it is relatively easy to define to the domain expert

Probability parameters of BBN

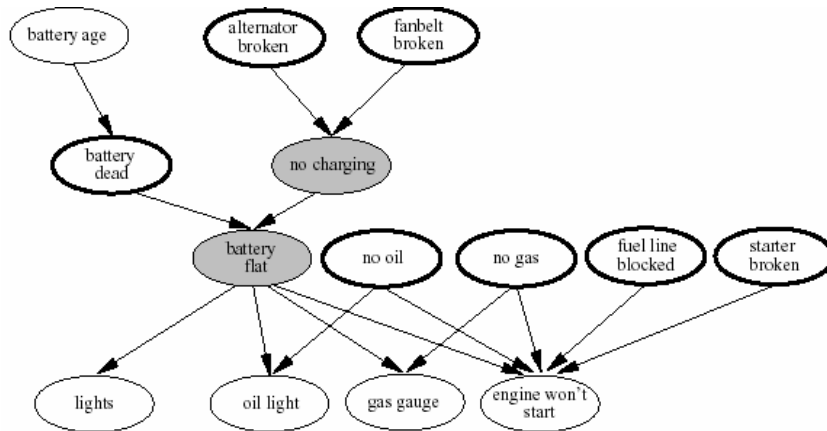
- are conditional distributions relating random variables and their parents
- Complexity is much smaller than the full joint
- It is much easier to obtain such probabilities from the expert or learn them automatically from data

BBNs built in practice

- **In various areas:**
 - Intelligent user interfaces (Microsoft)
 - Troubleshooting, diagnosis of a technical device
 - Medical diagnosis:
 - Pathfinder (Intellipath)
 - CPSC
 - Munin
 - QMR-DT
 - Collaborative filtering
 - Military applications
 - Business and finance
 - Insurance, credit applications

Diagnosis of car engine

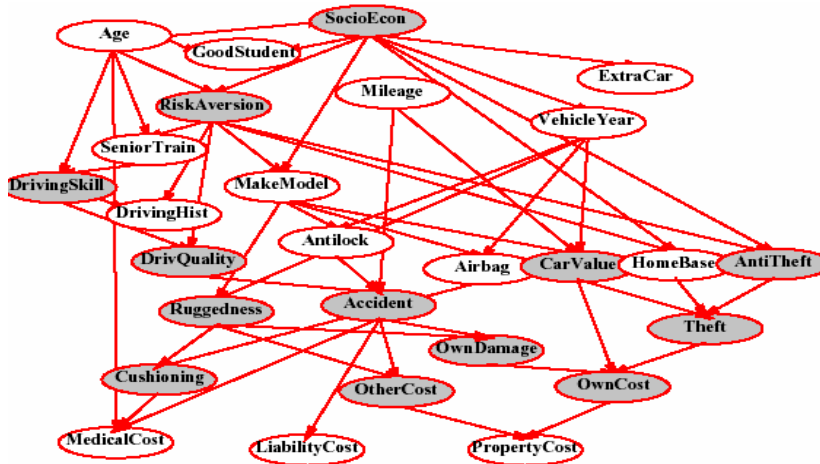
- Diagnose the engine start problem



CS 3710 Probabilistic graphical models

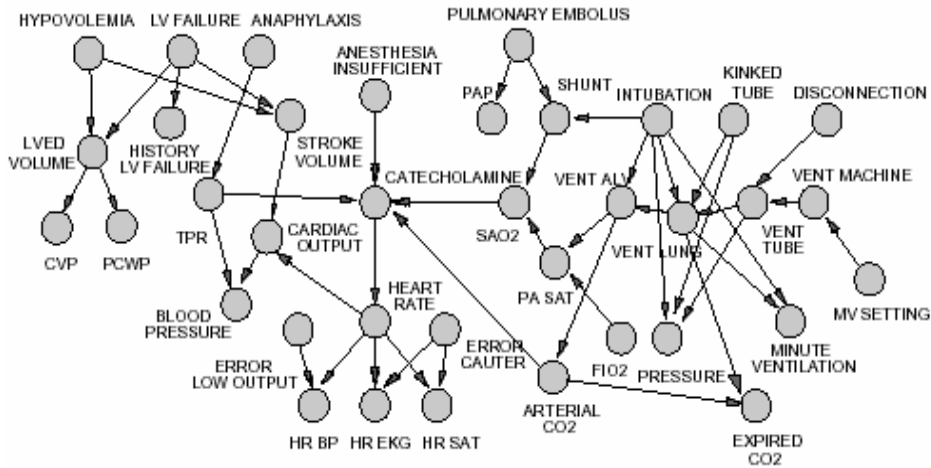
Car insurance example

- Predict claim costs (medical, liability) based on application data



CS 3710 Probabilistic graphical models

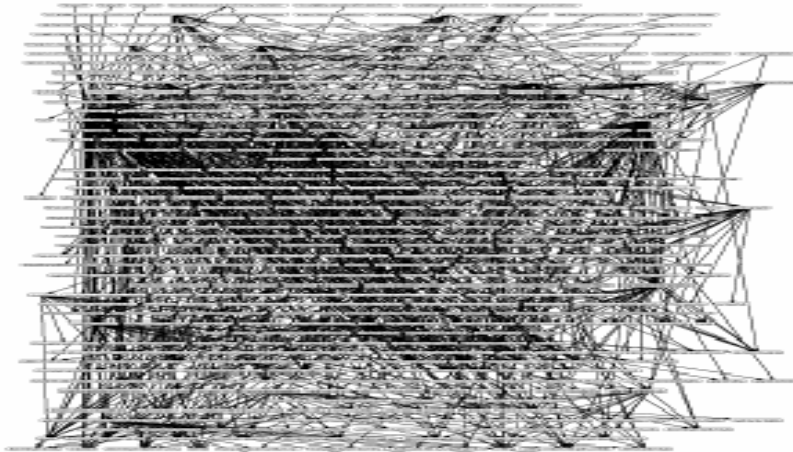
(ICU) Alarm network



CS 3710 Probabilistic graphical models

CPCS

- Computer-based Patient Case Simulation system (CPCS-PM) developed by Parker and Miller (University of Pittsburgh)
- 422 nodes and 867 arcs

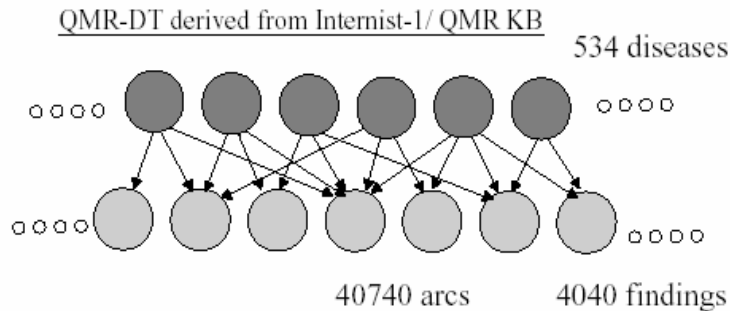


CS 3710 Probabilistic graphical models

QMR-DT

- **Medical diagnosis in internal medicine**

Bipartite network of disease/findings relations



CS 3710 Probabilistic graphical models

Inference in Bayesian networks

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
- Simplifies the acquisition of a probabilistic model
- But we are interested in solving various **inference tasks**:

- **Diagnostic task. (from effect to cause)**

$$\mathbf{P}(\textit{Burglary} \mid \textit{JohnCalls} = T)$$

- **Prediction task. (from cause to effect)**

$$\mathbf{P}(\textit{JohnCalls} \mid \textit{Burglary} = T)$$

- **Other probabilistic queries** (queries on joint distributions).

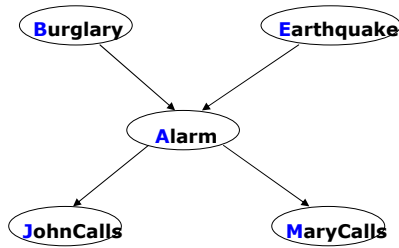
$$\mathbf{P}(\textit{Alarm})$$

- **Main issue:** Can we take advantage of independences to construct special algorithms and speeding up the inference?

CS 3710 Probabilistic graphical models

Inference in Bayesian network

- **Bad news:**
 - Exact inference problem in BBNs is NP-hard (Cooper)
 - Approximate inference is NP-hard (Dagum, Luby)
- **But** very often we can achieve significant improvements
- Assume our Alarm network



- Assume we want to compute: $P(J = T)$

Inference in Bayesian networks

Computing: $P(J = T)$

Approach 1. Blind approach.

- Sum out all un-instantiated variables from the full joint,
- express the joint distribution as a product of conditionals

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(B = b, E = e, A = a, J = T, M = m) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e)
 \end{aligned}$$

Computational cost:

Number of additions: ?

Number of products: ?

Inference in Bayesian networks

Computing: $P(J = T)$

Approach 1. Blind approach.

- Sum out all un-instantiated variables from the full joint,
- express the joint distribution as a product of conditionals

$$\begin{aligned} P(J = T) &= \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(B = b, E = e, A = a, J = T, M = m) \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \end{aligned}$$

Computational cost:

Number of additions: 15

Number of products: ?

Inference in Bayesian networks

Computing: $P(J = T)$

Approach 1. Blind approach.

- Sum out all un-instantiated variables from the full joint,
- express the joint distribution as a product of conditionals

$$\begin{aligned} P(J = T) &= \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(B = b, E = e, A = a, J = T, M = m) \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \end{aligned}$$

Computational cost:

Number of additions: 15

Number of products: $16 * 4 = 64$

Inference in Bayesian networks

Approach 2. Interleave sums and products

- Combines sums and product in a smart way (multiplications by constants can be taken out of the sum)

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

Computational cost:

Number of additions: $1 + 2 * [1 + 1 + 2 * 1] = ?$

Number of products: $2 * [2 + 2 * (1 + 2 * 1)] = ?$

Inference in Bayesian networks

Approach 2. Interleave sums and products

- Combines sums and product in a smart way (multiplications by constants can be taken out of the sum)

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

Computational cost:

Number of additions: $1 + 2 * [1 + 1 + 2 * 1] = 9$

Number of products: $2 * [2 + 2 * (1 + 2 * 1)] = ?$

Inference in Bayesian networks

Approach 2. Interleave sums and products

- Combines sums and product in a smart way (multiplications by constants can be taken out of the sum)

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in \mathcal{T}, F} \sum_{e \in \mathcal{T}, F} \sum_{a \in \mathcal{T}, F} \sum_{m \in \mathcal{T}, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in \mathcal{T}, F} \sum_{a \in \mathcal{T}, F} \sum_{m \in \mathcal{T}, F} P(J = T | A = a) P(M = m | A = a) P(B = b) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &= \sum_{a \in \mathcal{T}, F} P(J = T | A = a) \left[\sum_{m \in \mathcal{T}, F} P(M = m | A = a) \right] \left[\sum_{b \in \mathcal{T}, F} P(B = b) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

Computational cost:

Number of additions: $1 + 2 * [1 + 1 + 2 * 1] = 9$

Number of products: $2 * [2 + 2 * (1 + 2 * 1)] = 16$

Inference in Bayesian networks

- The smart interleaving of sums and products can help us to speed up the computation of joint probability queries
- What if we want to compute: $P(B = T, J = T)$

$$\begin{aligned}
 P(B = T, J = T) &= \\
 &= \sum_{a \in \mathcal{T}, F} P(J = T | A = a) \left[\sum_{m \in \mathcal{T}, F} P(M = m | A = a) \right] \left[P(B = T) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = T, E = e) P(E = e) \right] \right] \\
 P(J = T) &= \quad \updownarrow \quad \updownarrow \quad \updownarrow \quad \updownarrow \quad \updownarrow \\
 &= \sum_{a \in \mathcal{T}, F} P(J = T | A = a) \left[\sum_{m \in \mathcal{T}, F} P(M = m | A = a) \right] \left[\sum_{b \in \mathcal{T}, F} P(B = b) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

- A lot of shared computation
 - Smart caching of results can save the time for more queries

Inference in Bayesian networks

- The smart interleaving of sums and products can help us to speed up the computation of joint probability queries
- What if we want to compute: $P(B = T, J = T)$

$$P(B = T, J = T) = \sum_{a \in \mathcal{T}, F} P(J = T | A = a) \left[\sum_{m \in \mathcal{T}, F} P(M = m | A = a) \right] \left[P(B = T) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = T, E = e) P(E = e) \right] \right]$$

$$P(J = T) = \sum_{a \in \mathcal{T}, F} P(J = T | A = a) \left[\sum_{m \in \mathcal{T}, F} P(M = m | A = a) \right] \left[\sum_{b \in \mathcal{T}, F} P(B = b) \left[\sum_{e \in \mathcal{T}, F} P(A = a | B = b, E = e) P(E = e) \right] \right]$$

- A lot of shared computation
 - Smart caching of results can save the time if more queries

Inference in Bayesian networks

- When caching of results becomes handy?
- What if we want to compute a diagnostic query:

$$P(B = T | J = T) = \frac{P(B = T, J = T)}{P(J = T)}$$

- Exactly probabilities we have just compared !!
- There are other queries when caching and ordering of sums and products can be shared and saves computation

$$\mathbf{P}(B | J = T) = \frac{\mathbf{P}(B, J = T)}{P(J = T)} = \alpha \mathbf{P}(B, J = T)$$

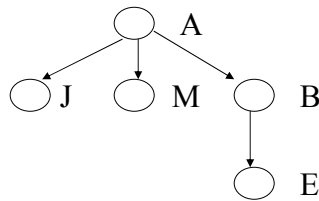
- General technique: **Variable elimination**

Inference in Bayesian networks

- General idea of variable elimination

$$\begin{aligned}
 P(\text{True}) &= 1 = \\
 &= \sum_{a \in T, F} \underbrace{\left[\sum_{j \in T, F} P(J=j | A=a) \right]}_{f_J(a)} \underbrace{\left[\sum_{m \in T, F} P(M=m | A=a) \right]}_{f_M(a)} \underbrace{\left[\sum_{b \in T, F} P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \right]}_{f_E(a, b)} \\
 &\hspace{15em} \underbrace{\hspace{15em}}_{f_B(a)}
 \end{aligned}$$

Variable order:



Results cached in
the tree structure

Complexity:
treewidth of the graph

Inference in Bayesian network

- **Exact inference algorithms:**
 - Variable elimination
 - Symbolic inference (D'Ambrosio)
 - Recursive decomposition (Cooper)
 - Message passing algorithm (Pearl)
 - Clustering and joint tree approach (Lauritzen, Spiegelhalter)
 - Arc reversal (Olmsted, Schachter)
- **Approximate inference algorithms:**
 - **Monte Carlo methods:**
 - Forward sampling, Likelihood sampling
 - Variational methods

Markov random fields

- **Probabilistic models with symmetric dependences.**
 - Typically models of spatially varying quantities

$$P(x) \propto \prod_{c \in cl(x)} f_c(x_c)$$

$f_c(x_c)$ - A potential function (defined over factors)

$$P(x) = \frac{1}{Z} \exp\left(-\sum_{c \in cl(x)} \phi_c(x_c)\right)$$

- Gibbs (Boltzman) distribution

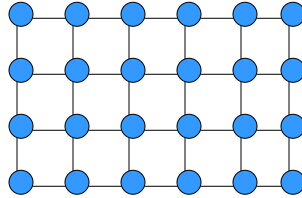
$$Z = \sum_{x \in \{x\}} \exp\left(-\sum_{c \in cl(x)} \phi_c(x_c)\right) \quad \text{- A partition function}$$

Markov random fields

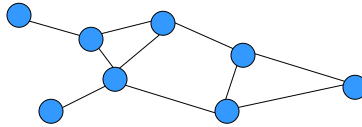
- **Interactions induced by the factorized form are captured by an undirected network (also called independence graph)**
- $G = (S, E)$
 - $S=1, 2, \dots, N$ correspond to random variables
 - $(i, j) \in E \Leftrightarrow \exists c : \{i, j\} \subset c$
or x_i and x_j appear within the same factor c
- **Consequence:**
 - factors c correspond to cliques of the graph

Markov random fields

- regular lattice (Ising model)



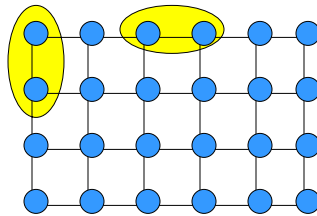
- Arbitrary graph



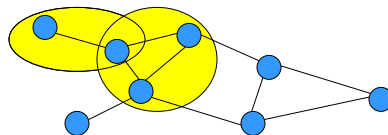
CS 3710 Probabilistic graphical models

Markov random fields

- regular lattice (Ising model)



- Arbitrary graph

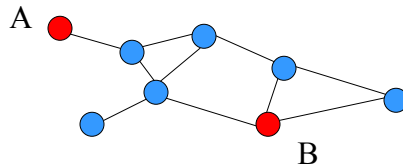


CS 3710 Probabilistic graphical models

Markov random fields

- **Pairwise Markov property**

- Two nodes in the network that are not directly connected can be made independent given all other nodes



$$P(x_A, x_B | x_r) = \frac{P(x_A, x_B, x_r)}{P(x_r)} \propto \exp\left(-\sum_{c:c \cap A \neq \{\}} \phi_c(x_c) - \sum_{c:c \cap B \neq \{\}} \phi_c(x_c)\right)$$
$$\propto \exp\left(-\sum_{c:c \cap A \neq \{\}} \phi_c(x_c)\right) = P(x_A | x_r)$$

Markov random fields

- **Pairwise Markov property**

- Two nodes in the network that are not directly connected can be made independent given all other nodes

- **Local Markov property**

- A set of nodes (variables) can be made independent from the rest of nodes variables given its immediate neighbors

- **Global Markov property**

- A vertex set A is independent of the vertex set B (A and B are disjoint) given set C if all chains in between elements in A and B intersect C