

**CS 2750 Machine Learning  
Lecture 15**

**Bayesian belief networks.  
Inference.**

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

**Midterm exam**

**Wednesday, March 17, 2004**

- **In class**
- **Closed book**
- **Material covered before Spring break**
- **Last year midterm will be posted on the web**

---

CS 2750 Machine Learning

## Project proposals

**Due: Wednesday, March 24, 2004**

- **1-2 pages long**

### Proposal

- **Written proposal:**
  1. Outline of a learning problem, type of data you have available. Why is the problem important?
  2. Learning methods you plan to try and implement for the problem. References to previous work.
  3. How do you plan to test, compare learning approaches
  4. Schedule of work (approximate timeline of work)
- **A 3-slide PPT presentation summarizing points 1-4**

## Modeling uncertainty with probabilities

- **Full joint distribution:** joint distribution over all random variables defining the domain
  - it is sufficient to represent the complete domain and to do any type of probabilistic inferences

### Problems:

- **Space complexity.** To store full joint distribution requires to remember  $O(d^n)$  numbers.
  - $n$  – number of random variables,  $d$  – number of values
- **Inference complexity.** To compute some queries requires  $O(d^n)$  steps.
- **Acquisition problem.** Who is going to define all of the probability entries?

## Pneumonia example. Complexities.

- **Space complexity.**
  - Pneumonia (2 values: T,F), Fever (2: T,F), Cough (2: T,F), WBCcount (3: high, normal, low), paleness (2: T,F)
  - Number of assignments:  $2*2*2*3*2=48$
  - We need to define at least 47 probabilities.
- **Time complexity.**
  - Assume we need to compute the probability of Pneumonia=T from the full joint

$$P(\text{Pneumonia} = T) = \sum_{i \in T, F} \sum_{j \in T, F} \sum_{k=h, n, l} \sum_{u \in T, F} P(\text{Fever} = i, \text{Cough} = j, \text{WBCcount} = k, \text{Pale} = u)$$

- Sum over  $2*2*3*2=24$  combinations

---

CS 2750 Machine Learning

## Bayesian belief networks (BBNs)

### Bayesian belief networks.

- Represent the full joint distribution over the variables more compactly with 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**

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

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

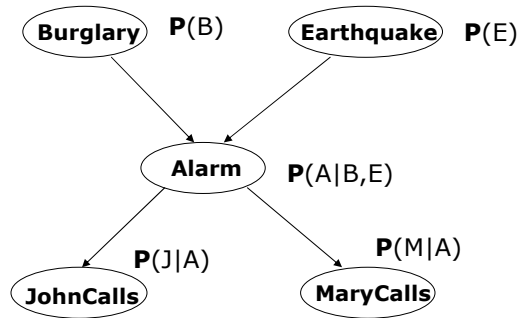
---

CS 2750 Machine Learning

## Bayesian belief network.

### 1. Directed acyclic graph

- **Nodes** = random variables
- **Links** = missing links encode independences.

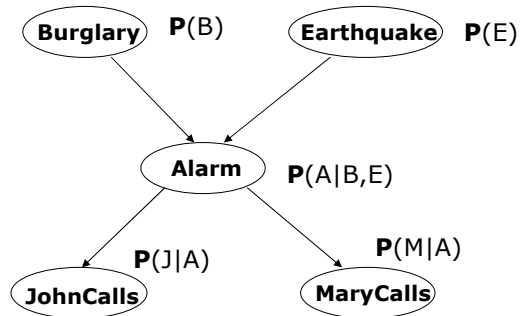


CS 2750 Machine Learning

## Bayesian belief network.

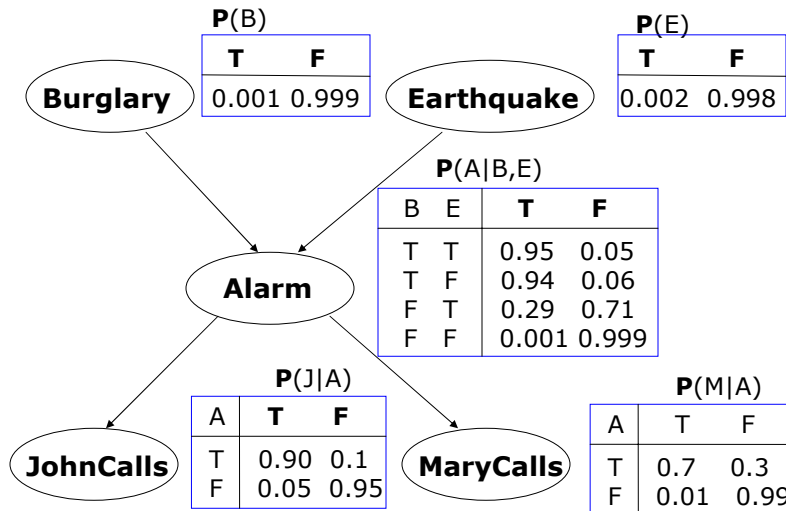
### 2. Local conditional distributions

- relate variables and their parents



CS 2750 Machine Learning

## Bayesian belief network.



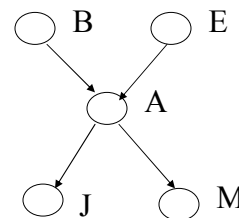
CS 2750 Machine Learning

## 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

CS 2750 Machine Learning

## 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:

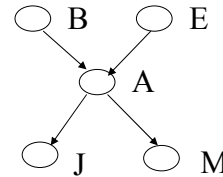
Assume the following assignment of values to random variables

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

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**

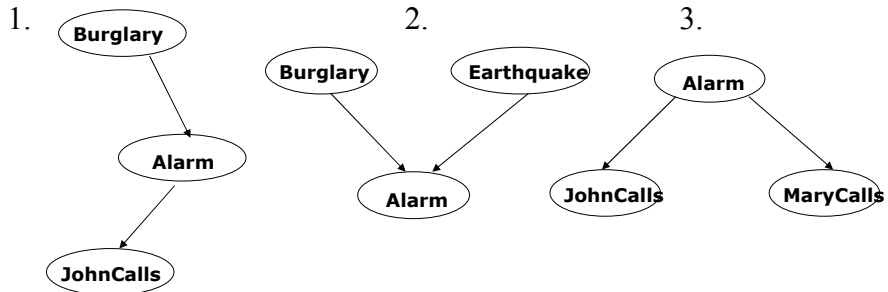
$$P(A \mid C, B) = P(A \mid C)$$

$$P(A, B \mid C) = P(A \mid C)P(B \mid C)$$

- **The graph structure implies the decomposition !!!**

## Independences in BBNs

3 basic independence structures:

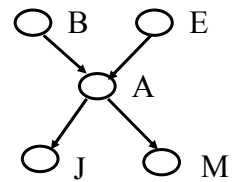


CS 2750 Machine Learning

## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

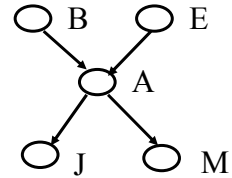
$$P(B=T, E=T, A=T, J=T, M=F) =$$



CS 2750 Machine Learning

## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



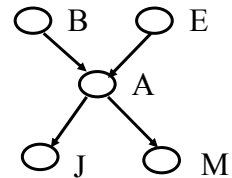
$$P(B=T, E=T, A=T, J=T, M=F) =$$

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

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

## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

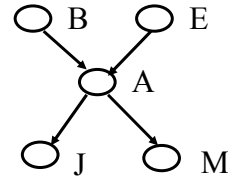
$$P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T)$$

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



## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

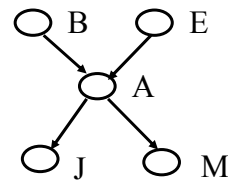
$$P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T)$$

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

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

## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

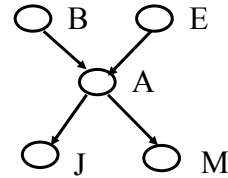
$$\underline{P(M=F \mid A=T)} P(B=T, E=T, A=T)$$

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

$$P(B=T) P(E=T)$$

## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

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

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

$$P(B=T) P(E=T)$$

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

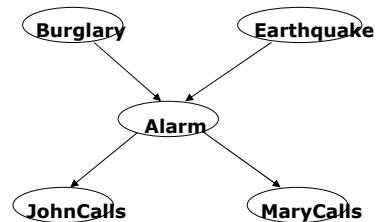
CS 2750 Machine Learning

## Parameter complexity problem

- In the BBN the full joint distribution is expressed as a product of conditionals (of smaller) complexity

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

Parameters:  
full joint: ?  
BBN: ?



CS 2750 Machine Learning

## Parameter complexity problem

- In the BBN the full joint distribution is expressed as a product of conditionals (of smaller) complexity

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

**Parameters:**

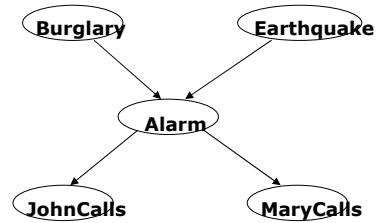
**full joint:**  $2^5 = 32$

**BBN:**  $2^3 + 2(2^2) + 2(2) = 20$

**Parameters to be defined:**

**full joint:**  $2^5 - 1 = 31$

**BBN:**  $2^2 + 2(2) + 2(1) = 10$



CS 2750 Machine Learning

## Model acquisition problem

**The structure of the BBN** typically reflects causal relations

- BBNs are also sometime referred to as **causal networks**
- Causal structure is very intuitive in many applications domain and it is relatively easy to obtain from the domain expert

**Probability parameters of BBN** correspond to conditional distributions relating a random variable and its parents only

- Their complexity much smaller than the full joint
- Easier to come up (estimate) the probabilities from expert or automatically by learning from data

CS 2750 Machine Learning

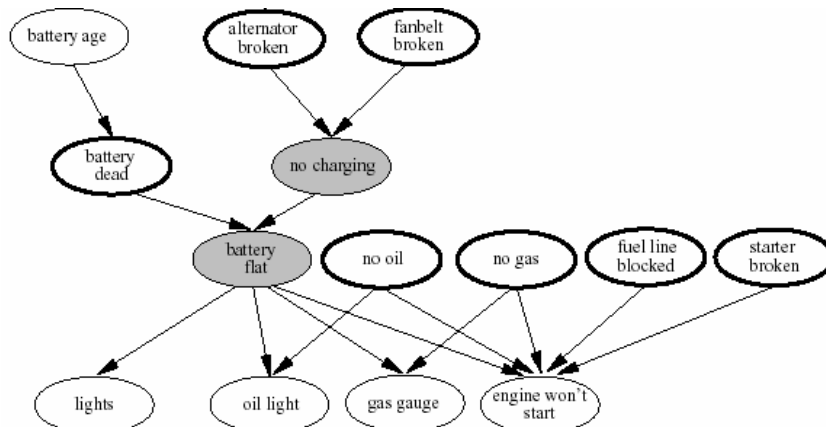
## BBNs built in practice

- **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
  - Insurance, credit applications

CS 2750 Machine Learning

## Diagnosis of car engine

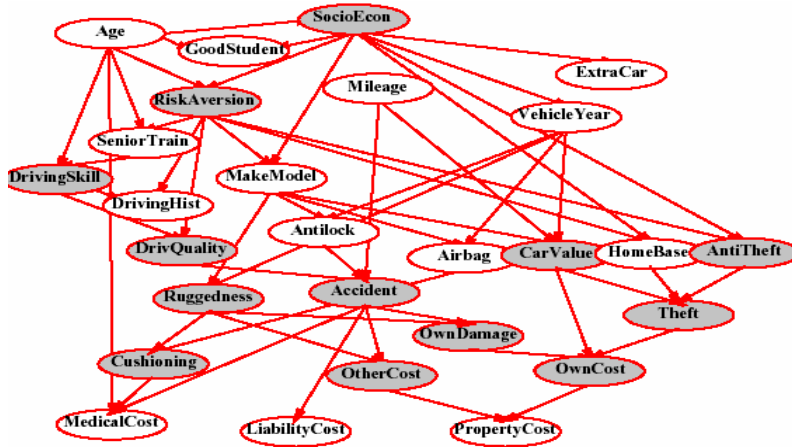
- Diagnose the engine start problem



CS 2750 Machine Learning

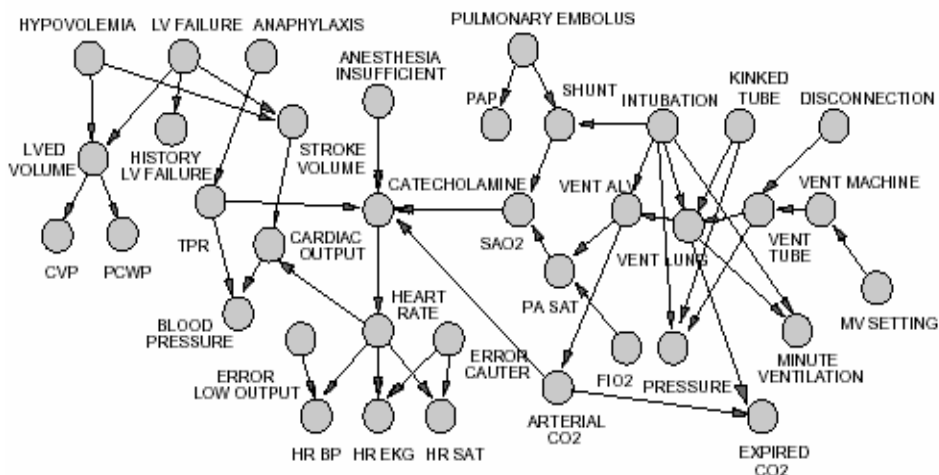
## Car insurance example

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



CS 2750 Machine Learning

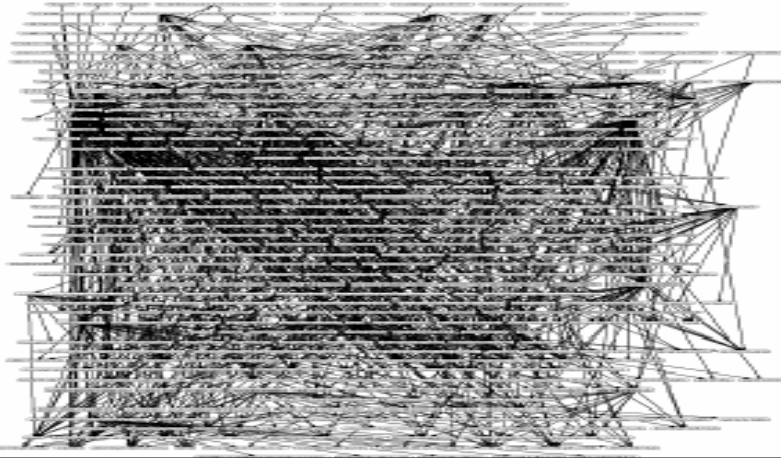
## (ICU) Alarm network



CS 2750 Machine Learning

## CPCS

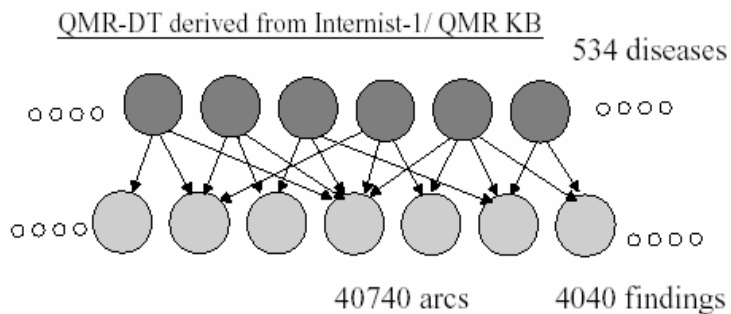
- Computer-based Patient Case Simulation system (CPCS-PM) developed by Parker and Miller (at University of Pittsburgh)
- 422 nodes and 867 arcs



CS 2750 Machine Learning

## QMR-DT

- **Medical diagnosis in internal medicine**
  - Bipartite network of disease/findings relations
  - Derived from the Internist-1/QMR knowledge base



CS 2750 Machine Learning

## Inference in Bayesian networks

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
  - Smaller number of parameters
- But we are interested in solving various **inference tasks**:
  - **Diagnostic task. (from effect to cause)**

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

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

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

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

$$P(\text{Alarm})$$

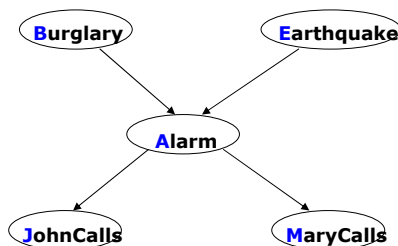
- **Question:** Can we take advantage of independences to construct special algorithms and speedup the inference?

---

CS 2750 Machine Learning

## 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:  $P(J = T)$

---

CS 2750 Machine Learning

## 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:  $16 * 4 = \mathbf{64}$

## 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: ?

Number of products: ?



## 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:  $1 + 2*(1) + 2*(1 + 2*(1)) = 9$

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

## 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:  $P(B = T, J = T)$

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

- 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:  $P(B = T, J = T)$

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

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

- 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

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

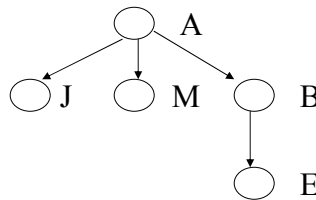
- **General technique: Variable elimination**

## 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

## Inferences in Bayesian network

- **Exact inference algorithms:**
  - Symbolic inference (D<sup>2</sup> Ambrosio)
  - 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