

CS 1675 Introduction to Machine Learning  
Lecture 16

Bayesian belief networks:  
learning and inference

Milos Hauskrecht

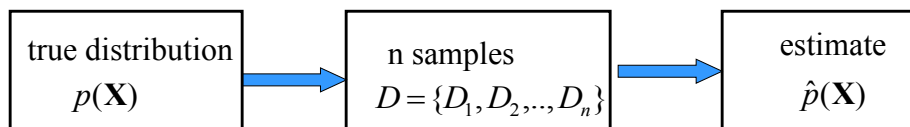
[milos@pitt.edu](mailto:milos@pitt.edu)

5329 Sennott Square

Density estimation

**Data:**  $D = \{D_1, D_2, \dots, D_n\}$   
 $D_i = \mathbf{x}_i$  a vector of attribute values

**Objective:** try to estimate the underlying true probability distribution over variables  $\mathbf{X}$ ,  $p(\mathbf{X})$ , using examples in  $D$



**Standard (iid) assumptions:** Samples

- are **independent** of each other
- come from the same **(i)dentical (d)istribution** (fixed  $p(\mathbf{X})$ )

## Modeling complex distributions

**Question:** How to model and learn complex multivariate distributions  $\hat{p}(\mathbf{X})$  with a large number of variables?

**Example: modeling of disease – symptoms relations**

- **Disease:** pneumonia
- **Patient symptoms (findings, lab tests):**
  - Fever, Cough, Paleness, WBC (white blood cells) count, Chest pain, etc.
- **Model of the full joint distribution:**  
 $P(\text{Pneumonia}, \text{Fever}, \text{Cough}, \text{Paleness}, \text{WBC}, \text{Chest pain})$

One probability per assignment of values to variables:

$P(\text{Pneumonia} = T, \text{Fever} = T, \text{Cough} = T, \text{WBC} = \text{High}, \text{Chest pain} = T)$

---

## Bayesian belief networks (BBNs)

**Bayesian belief networks** (late 80s, beginning of 90s)

- Give solutions to the space, acquisition bottlenecks
- Partial solutions for time complexities

**Key features:**

- 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
- **X and Y are independent**  $P(X, Y) = P(X)P(Y)$
- **X and Y are conditionally independent given Z**

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

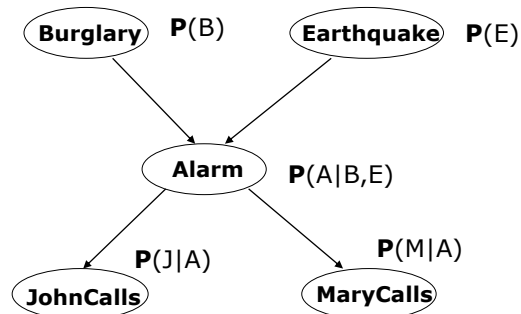
$$P(X | Y, Z) = P(X | Z)$$

---

## Bayesian belief network

### 1. Directed acyclic graph

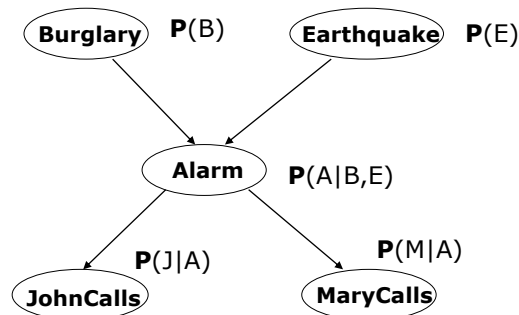
- **Nodes** = random variables  
Burglary, Earthquake, Alarm, Mary calls and John calls
- **Links** = direct (causal) dependencies between variables.  
The chance of Alarm being is influenced by Earthquake,  
The chance of John calling is affected by the Alarm



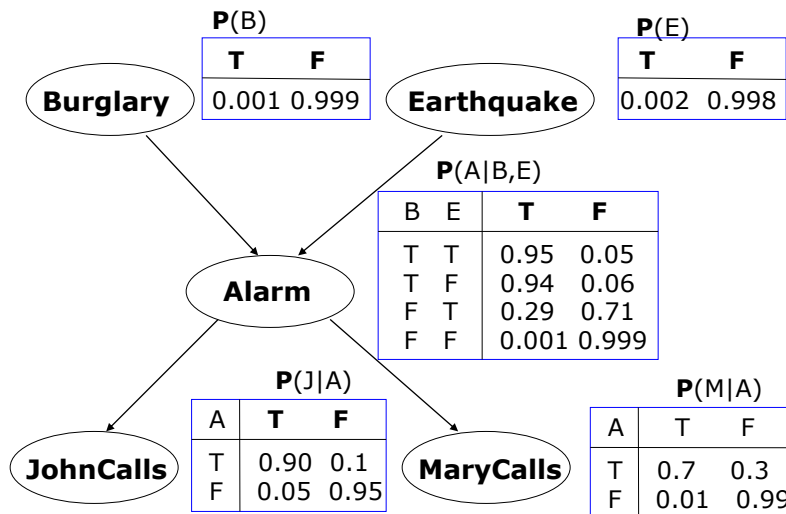
## Bayesian belief network

### 2. Local conditional distributions

- relating variables and their parents



## Bayesian belief network



## Full joint distribution in BBNs

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

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

### Example:

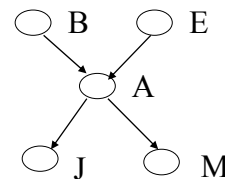
Assume the following assignment of values to random variables

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

Then its probability is:

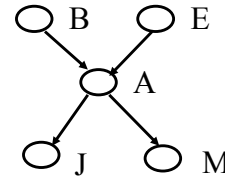
$$P(B = T, E = T, A = T, J = T, M = F) =$$

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



## Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

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

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

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

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

## 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: binary (True, False) variables

# of parameters of the full joint:

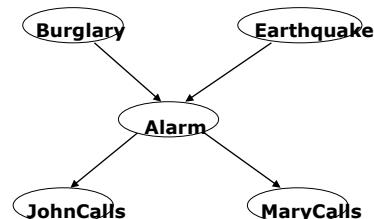
$$2^5 = 32$$

One parameter is for free:

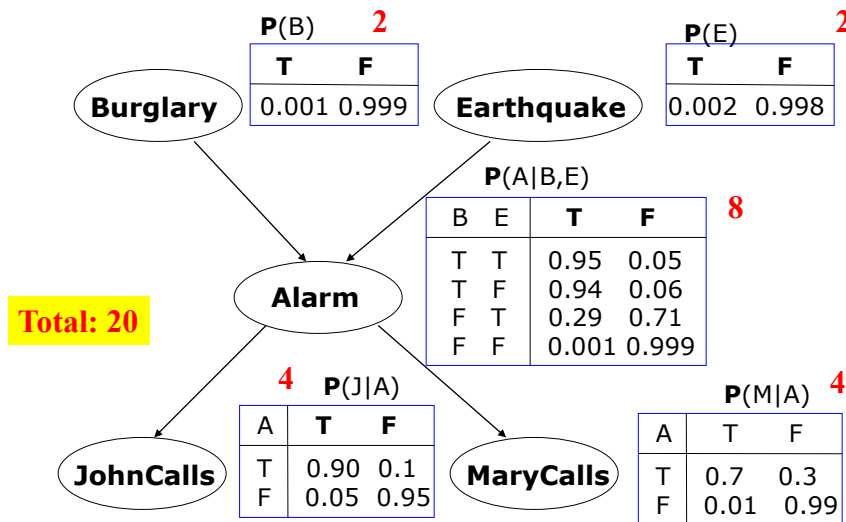
$$2^5 - 1 = 31$$

# of parameters of the BBN:

?



## Bayesian belief network: parameters count



## Parameter complexity problem

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

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

- 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:

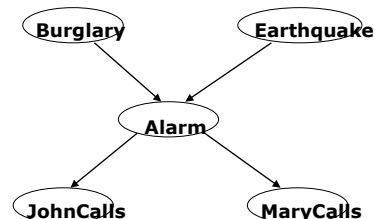
$$2^5 - 1 = 31$$

# of parameters of the BBN:

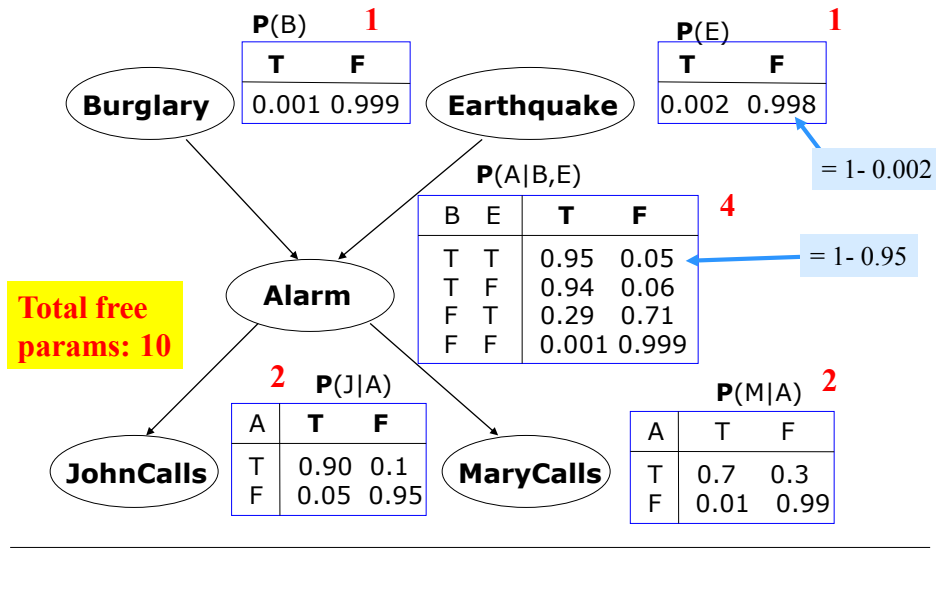
$$2^3 + 2(2^2) + 2(2) = 20$$

One parameter in every conditional is for free:

?



## Bayesian belief network: free parameters



## Parameter complexity problem

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

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

- 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:

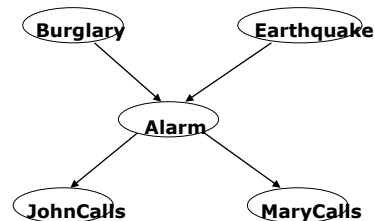
$$2^5 - 1 = 31$$

# of parameters of the BBN:

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

One parameter in every conditional is for free:

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

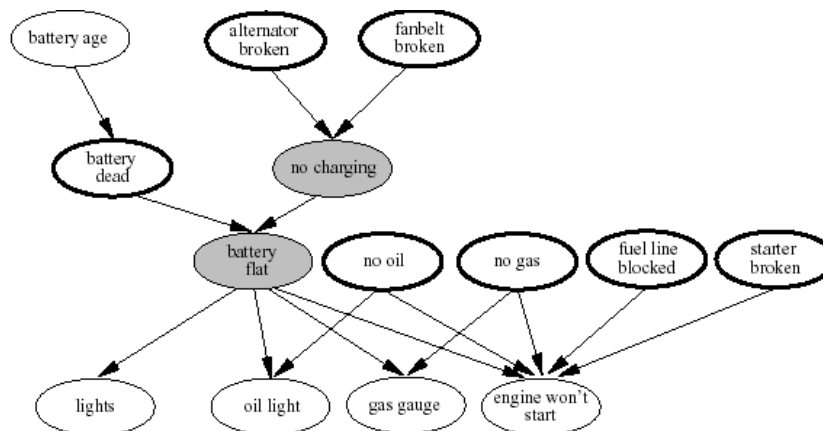


## BBNs examples

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

## Diagnosis of car engine

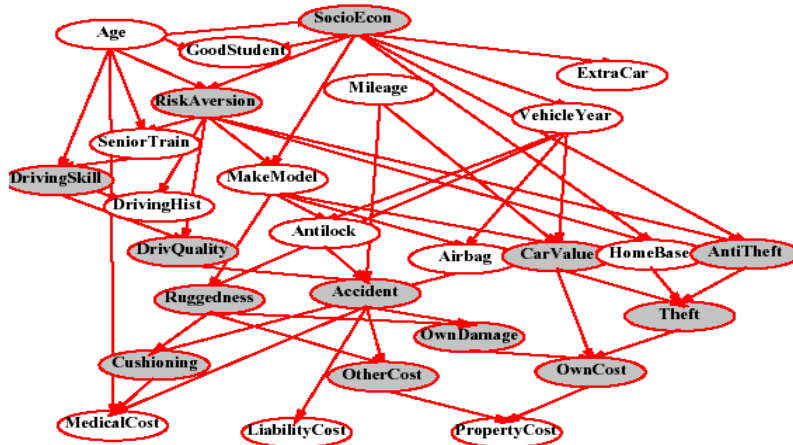
- Diagnose the engine start problem



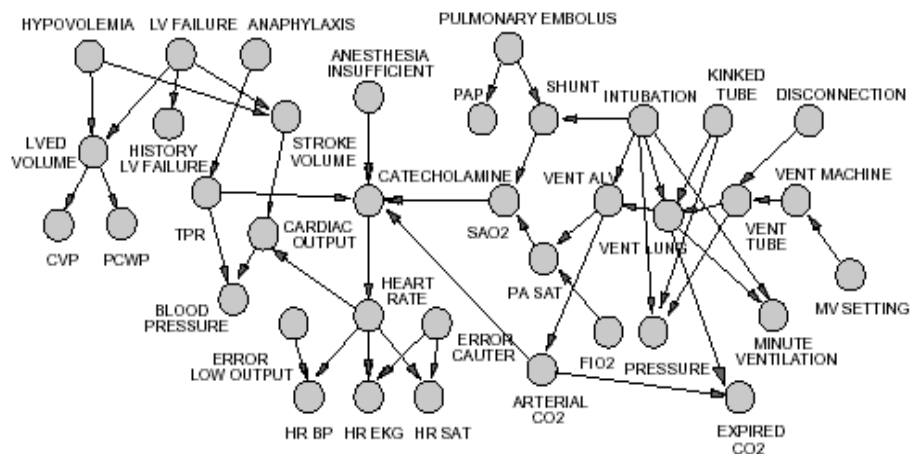


## Car insurance example

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

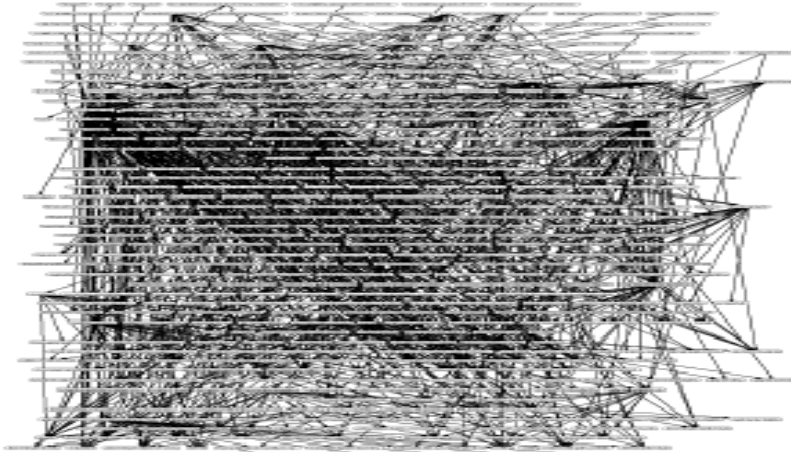


## (ICU) Alarm network



## CPCS

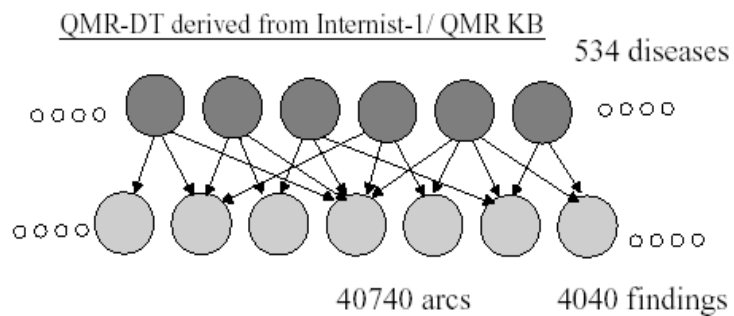
- Computer-based Patient Case Simulation system (CPCS-PM) developed by Parker and Miller (at University of Pittsburgh)
- 422 nodes and 867 arcs



## QMR-DT

- Medical diagnosis in internal medicine

Bipartite network of disease/findings relations



## Naïve Bayes model

A special (simple) Bayesian belief network

- used as a generative classifier model

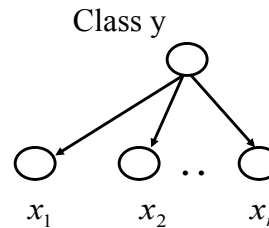
- Model of  $p(\mathbf{x}, y) = p(\mathbf{x} | y) p(y)$

- Class variable  $y$

$$p(y)$$

- Attributes are independent given  $y$

$$p(\mathbf{x} | y = i) = \prod_{j=1}^d p(x_j | y = i)$$



**Learning:**

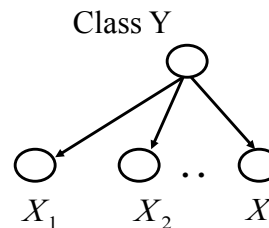
- Parameterize models of  $p(y)$  and all  $p(x_j | y=i)$
- ML estimates of the parameters

## Naïve Bayes model

A special (simple) Bayesian belief network

- used as a generative classifier model

Model of  $p(\mathbf{x}, y) = p(\mathbf{x} | y) p(y)$



**Classification:** given  $\mathbf{x}$  select the class

- Select the class with the maximum posterior
- Calculation of a posterior is an example of BBN inference

$$p(y = i | \mathbf{x}) = \frac{p(y = i) p(\mathbf{x} | y = i)}{\sum_{u=1}^k p(y = u) p(\mathbf{x} | y = u)} = \frac{p(y = i) \prod_{j=1}^d p(x_j | y = i)}{\sum_{u=1}^k p(y = u) \prod_{j=1}^d p(x_j | y = u)}$$

**Remember:** we can calculate the probabilities from the full joint

## Learning of BBN

### Learning.

- **Learning of parameters of conditional probabilities**
- **Learning of the network structure**

### Variables:

- **Observable** – values present in every data sample
- **Hidden** – they values are never observed in data
- **Missing values** – values sometimes present, sometimes not

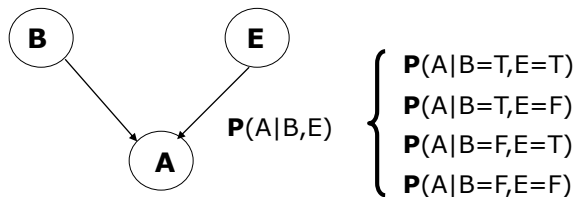
### Next:

- Learning of the parameters of BBN
  - Values for all variables are observable
- 

## Estimation of parameters of BBN

- **Idea:** decompose the estimation problem for the full joint over a large number of variables to a set of smaller estimation problems corresponding to local parent-variable conditionals.
- **Example:** Assume A,E,B are binary with *True*, *False* values

**Learning of  $P(A|B,E)$  = 4 estimation problems**



- **Assumption that enables the decomposition:** parameters of conditional distributions are independent
-

## Estimates of parameters of BBN

- Two assumptions that permit the decomposition:

- **Sample independence**

$$P(D | \Theta, \xi) = \prod_{u=1}^N P(D_u | \Theta, \xi)$$

- **Parameter independence**

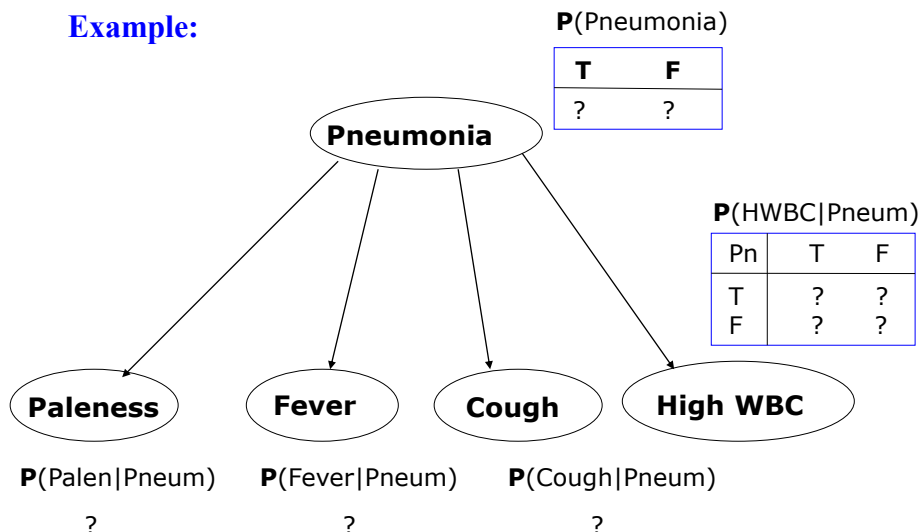
$$p(\Theta | D, \xi) = \prod_{i=1}^n \prod_{j=1}^{q_i} p(\theta_{ij} | D, \xi)$$

# of nodes  
# of parents' values

Parameters of **each conditional** (one for every assignment of values to parent variables) can be learned independently

## Learning of BBN parameters. Example.

**Example:**

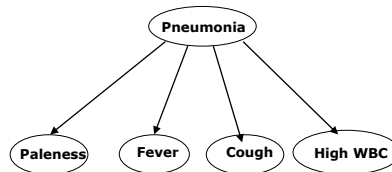


## Learning of BBN parameters. Example.

Data D (different patient cases):

Pal Fev Cou HWB Pneu

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



## Estimates of parameters of BBN

- Much like multiple **coin toss or roll of a dice** problems.
- A “smaller” learning problem corresponds to the learning of exactly one conditional distribution
- **Example:**  $P(\text{Fever} | \text{Pneumonia} = T)$
- **Problem:** How to pick the data to learn?

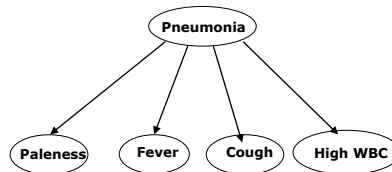
## Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 1:** Select data points with Pneumonia=T

Pal Fev Cou HWB Pneu

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



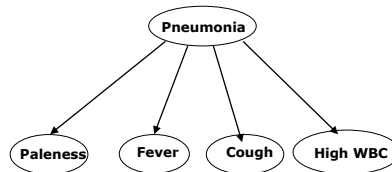
## Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 1:** Ignore the rest

Pal Fev Cou HWB Pneu

F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	T	T	T	T
F	T	F	T	T



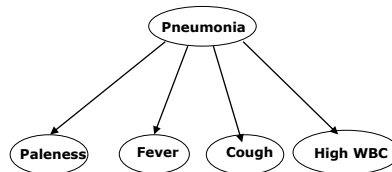
## Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 2:** Select values of the random variable defining the distribution of Fever

Pal Fev Cou HWB Pneu

F	<b>F</b>	T	T	T
F	<b>F</b>	T	F	T
F	<b>T</b>	T	T	T
T	<b>T</b>	T	T	T
F	<b>T</b>	F	T	T



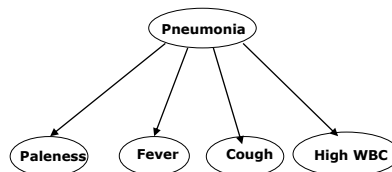
## Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 2:** Ignore the rest

Fev

**F**  
**F**  
**T**  
**T**  
**T**





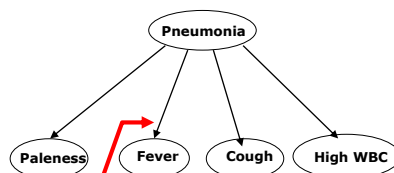
## Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} \mid \text{Pneumonia} = T)$

**Step 3a: Learning the ML estimate**

**Fev**

**F**  
**F**  
**T**  
**T**  
**T**



$P(\text{Fever} \mid \text{Pneumonia} = T)$

T	F
0.6	0.4

## Learning of BBN parameters. Bayesian learning.

**Learn:**  $P(\text{Fever} \mid \text{Pneumonia} = T)$

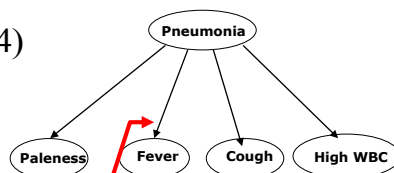
**Step 3b: Learning the Bayesian posterior**

**Assume the prior**

$$\theta_{\text{Fever} \mid \text{Pneumonia} \neq T} \sim \text{Beta}(3, 4)$$

**Fev**

**F**  
**F**  
**T**  
**T**  
**T**



**Posterior:**

$$\theta_{\text{Fever} \mid \text{Pneumonia} \neq T} \sim \text{Beta}(6, 6)$$

$$\theta_{\text{Fever} \mid \text{Pneumonia} \neq T}^{\text{MAP}} = \frac{6-1}{6+6-2} = 0.5$$

**MAP estimates**

T	F
0.5	0.5

## Estimates of parameters of BBN

Much like multiple **coin toss or roll of a dice** problems.

- A “smaller” learning problem corresponds to the learning of exactly one conditional distribution

**Example:**

$$\mathbf{P}(Fever | Pneumonia = T)$$

**Problem:** How to pick the data to learn?

**Answer:**

1. Select data points with Pneumonia=T  
(ignore the rest)
  2. Focus on (select) only values of the random variable defining the distribution (Fever)
  3. Learn the parameters of the local conditionals the same way as we learned the parameters of a biased coin or a die
- 

## Probabilistic inferences

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
  - Simplifies the representation and learning of a model
  - Can be used for the different **inference tasks** ....
-

## Bayes theorem

### Conditional/joint probability relations.

$$P(A|B) = \frac{P(A,B)}{P(B)} \quad \xrightarrow{\quad} \quad P(A,B) = P(B|A)P(A)$$

### Bayes theorem (switches conditioning events) :

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

### When is it useful?

- When we are interested in computing the diagnostic query from the causal probability

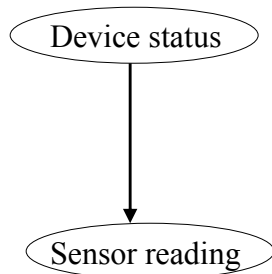
$$P(\text{cause} | \text{effect}) = \frac{P(\text{effect} | \text{cause})P(\text{cause})}{P(\text{effect})}$$

- **Reason:** It is often easier to assess causal probability
  - E.g. Probability of pneumonia causing fever  
vs. probability of pneumonia given fever

## Example: a simple diagnostic inference

- **Device** (equipment) operating *normally* or *malfunctioning*.
  - Operation of the device sensed indirectly via a sensor
- **Sensor reading** is either *High* or *Low*

**BBN**



**P(Device status)**

<b>normal</b>	<b>malfunctioning</b>
0.9	0.1

**P(Sensor reading | Device status)**

Status\Sensor	<b>High</b>	<b>Low</b>
normal	0.1	0.9
malfunc	0.6	0.4

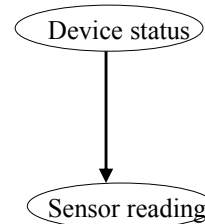
## Example: a simple diagnostic inference

- **Diagnostic inference:** compute the probability of device operating normally or malfunctioning **given a sensor reading**

$P(\text{Device status} \mid \text{Sensor reading} = \text{high}) =$

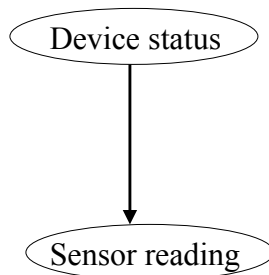
$$= \begin{pmatrix} P(\text{Device status} = \text{normal} \mid \text{Sensor reading} = \text{high}) \\ P(\text{Device status} = \text{malfunc} \mid \text{Sensor reading} = \text{high}) \end{pmatrix}$$

- Note we have the opposite conditional probabilities: they are much easier to estimate
- **Solution:** apply **Bayes theorem** to reverse the conditioning variables



## Example: a simple diagnostic inference

- **Device** (equipment) operating *normally* or *malfunctioning*.
  - Operation of the device sensed indirectly via a sensor
- **Sensor reading** is either *High* or *Low*



$P(\text{Device status})$

normal	malfunctioning
0.9	0.1

$P(\text{Sensor reading} \mid \text{Device status})$

Status\Sensor	High	Low
normal	0.1	0.9
malfunc	0.6	0.4

$P(\text{Device status} \mid \text{Sensor reading} = \text{high}) = ?$

## Bayes theorem

Assume a variable A with multiple values  $a_1, a_2, \dots, a_k$

Bayes theorem can be rewritten as:

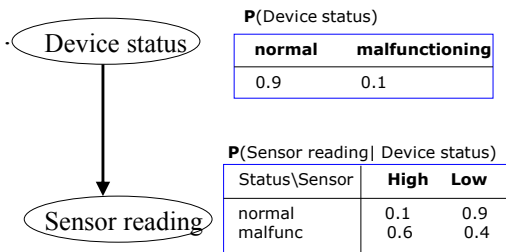
$$P(A = a_j | B = b) = \frac{P(B = b | A = a_j)P(A = a_j)}{P(B = b)}$$

$$= \frac{P(B = b | A = a_j)P(A = a_j)}{\sum_{j=1}^k P(B = b | A = a_j)P(A = a_j)}$$

Used in practice when we want to compute:

$P(A | B = b)$  for all values of  $a_1, a_2, \dots, a_k$

## Example: a simple diagnostic inference

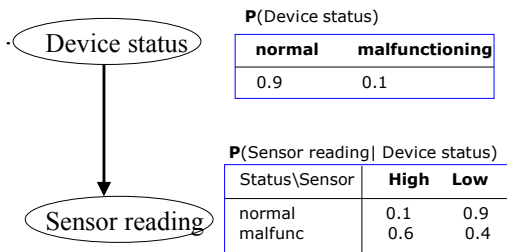


$P(\text{Device status} | \text{Sensor reading} = \text{high}) = ?$

$P(\text{Device status} = \text{norm} | \text{Sensor reading} = \text{high}) =$

$$= \frac{P(\text{Sensor reading} = \text{high}, \text{Device status} = \text{norm})}{P(\text{Sensor reading} = \text{high})}$$

## Example: a simple diagnostic inference



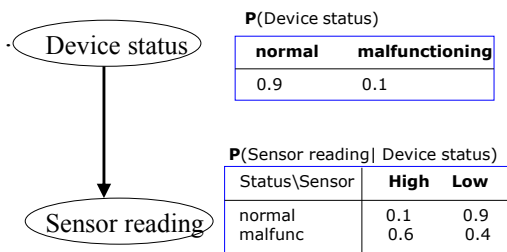
$$P(\text{Device status} \mid \text{Sensor reading} = \text{high}) = ?$$

$$P(\text{Device status} = \text{norm} \mid \text{Sensor reading} = \text{high}) =$$

$$= \frac{P(\text{Sensor reading} = \text{high}, \text{Device status} = \text{norm})}{P(\text{Sensor reading} = \text{high})}$$

$$= \frac{P(\text{Sensor reading} = \text{high} \mid \text{Device status} = \text{norm})P(\text{Device status} = \text{norm})}{P(\text{Sensor reading} = \text{high})}$$

## Example: a simple diagnostic inference



$$P(\text{Device status} \mid \text{Sensor reading} = \text{high}) = ?$$

$$P(\text{Device status} = \text{norm} \mid \text{Sensor reading} = \text{high}) =$$

$$= \frac{P(\text{Sensor reading} = \text{high} \mid \text{Device status} = \text{norm})P(\text{Device status} = \text{norm})}{P(\text{Sensor reading} = \text{high})}$$

$$+ P(\text{Sensor reading} = \text{high} \mid \text{Device status} = \text{malf}) P(\text{Device status} = \text{malf})$$

## Inference in Bayesian networks

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
- Simplifies the representation and learning of a model
- 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)**

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

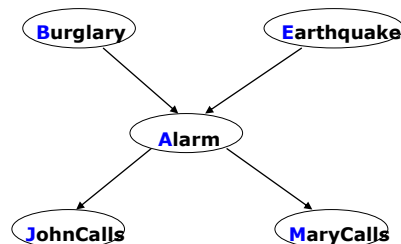
- **Other probabilistic queries** (queries on joint distributions).

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

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

## 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



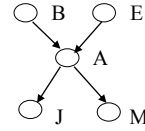
- Assume we want to compute:  $P(J = T)$
-

## 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



$$P(J = T) =$$

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

$$= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T \mid A = a) P(M = m \mid A = a) P(A = a \mid B = b, E = e) P(B = b) P(E = e)$$

**Computational cost:**

Number of additions: ?

Number of products: ?

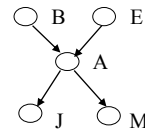
---

## 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



$$P(J = T) =$$

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

$$= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T \mid A = a) P(M = m \mid A = a) P(A = a \mid B = b, E = e) P(B = b) P(E = e)$$

**Computational cost:**

Number of additions: 15

Number of products: ?

---

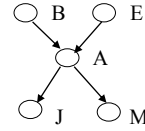


## Inference in Bayesian networks

**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



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

**Computational cost:**

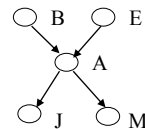
Number of additions: 15

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)



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

**Computational cost:**

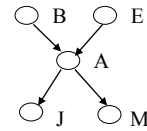
Number of additions:  $1 + 2 * [1 + 1 + 2 * 1] = ?$

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)



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

1

**Computational cost:**

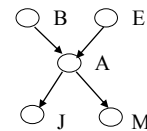
Number of additions: ?

---

## 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)



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

2\*1

**Computational cost:**

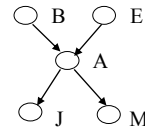
Number of additions: ?

---

## 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)



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

→ →  
 2\*2\*1

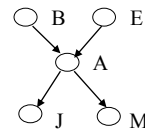
**Computational cost:**

Number of additions: ?

## 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)



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

→ → → →  
 2\*1      2\*2\*1

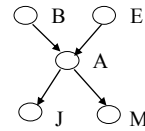
**Computational cost:**

Number of additions: ?

## 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)



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

Diagram illustrating the computational cost of the naive approach (Approach 1). Red arrows indicate the number of additions required for each summation step:

- Sum over  $a \in T, F$ : 1 addition
- Sum over  $m \in T, F$ :  $2 \times 1$  additions
- Sum over  $b \in T, F$ :  $2 \times 1$  additions
- Sum over  $e \in T, F$ :  $2 \times 2 \times 1$  additions

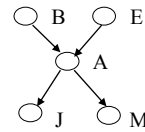
**Computational cost:**

Number of additions: ?

## 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)



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

Diagram illustrating the computational cost of the naive approach (Approach 1). Red arrows indicate the number of additions required for each summation step:

- Sum over  $a \in T, F$ : 1 addition
- Sum over  $m \in T, F$ :  $2 \times 1$  additions
- Sum over  $b \in T, F$ :  $2 \times 1$  additions
- Sum over  $e \in T, F$ :  $2 \times 2 \times 1$  additions

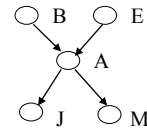
**Computational cost:**

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

## 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)



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

↓  
1

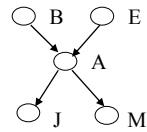
**Computational cost:**

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)



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

↓  
2\*2 \* 2\*1

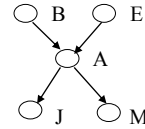
**Computational cost:**

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)



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

Diagram illustrating the computational cost of the naive approach. Blue arrows show the flow of computation from the innermost sum to the outermost. Red arrows indicate the cost of each summation step:

- Innermost sum:  $2 \times 2$
- Second sum:  $2 \times 2 \times 1$
- Third sum:  $2 \times 2 \times 2 \times 1$

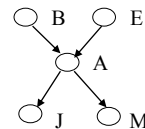
**Computational cost:**

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)



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

Diagram illustrating the computational cost of the naive approach. Blue arrows show the flow of computation from the innermost sum to the outermost. Red arrows indicate the cost of each summation step:

- Innermost sum:  $2 \times 2$
- Second sum:  $2 \times 2 \times 1$
- Third sum:  $2 \times 2 \times 2 \times 1$

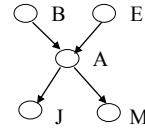
**Computational cost:**

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

## 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)



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

#### Computational cost:

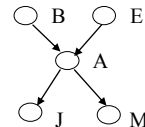
Number of additions:  $1+2*[1+1+2*1]=9$

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

## Variable elimination

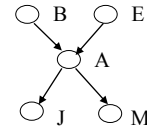
### Variable elimination:

- Similar idea but interleave sum and products one variable at the time during inference
- E.g. Query  $P(J=T)$  requires to eliminate A,B,E,M and this can be done in different order



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

## Variable elimination



**Assume order:** M, E, B, A to calculate  $P(J=T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J=T | A=a) P(M=m | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \left[ \sum_{m \in T, F} P(M=m | A=a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \quad 1 \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \left[ \sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \tau_1(A=a, B=b) \\
 &= \sum_{a \in T, F} P(J=T | A=a) \left[ \sum_{b \in T, F} P(B=b) \tau_1(A=a, B=b) \right] \\
 &= \sum_{a \in T, F} P(J=T | A=a) \tau_2(A=a) = \boxed{P(J=T)}
 \end{aligned}$$

## Inference in Bayesian network

- **Exact inference algorithms:**
  - **Variable elimination**
  - Recursive decomposition (Cooper, Darwiche)
  - Symbolic inference (D'Ambrosio)
  - Belief propagation 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