

CS 3750 Advanced Machine Learning

Lecture 7

Bayesian belief networks

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

CS 3750 Machine Learning

Modeling large probability distributions

- **Let** $\mathbf{X} = (X_1, X_2, \dots, X_n)$ be a set of features (represented as random variables), each with discrete set of values
- **Full joint distribution:** $P(\mathbf{X})$
 - joint distribution over all random variables defining the domain
 - Sufficient to do any probabilistic inference

Problems:

- **Space complexity.** To store full joint distribution requires to remember $O(u^n)$ numbers.
 n – number of random variables, u – number of values
- **Inference complexity.** To compute some queries requires $O(u^n)$ steps.

CS 3750 Machine Learning

Pneumonia example.

Random variables:

- **Pneumonia** (2 values: T,F),
- **Fever** (2: T,F),
- **Cough** (2: T,F),
- **WBCcount** (3: high, normal, low),
- **Paleness** (2: T,F)

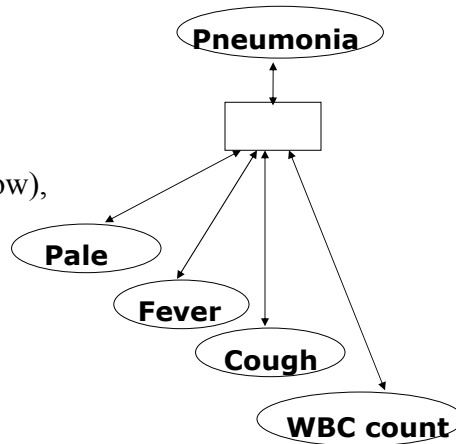
Joint distribution

- **Space complexity:**

Number of value assignments:

$$2*2*2*3*2=48$$

We need to define at least 47 probabilities.



Pneumonia example. Inferences

Any probabilistic query can be computed once the full joint is known

- **Time complexity** of computations
 - Assume we need to compute the probability of Pneumonia=T from the full joint
 - Sum over $2*2*3*2=24$ combinations

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

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

How to compute:

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

Bayesian belief networks (BBNs)

Bayesian belief networks.

- Represent the full joint distribution more compactly with smaller number of parameters.
- Take advantage of conditional and marginal independences among components in the distribution

- **A and B are independent**

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

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

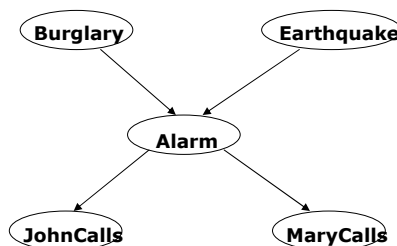
$$P(A, B | C) = P(A | C)P(B | C)$$

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

Alarm system example.

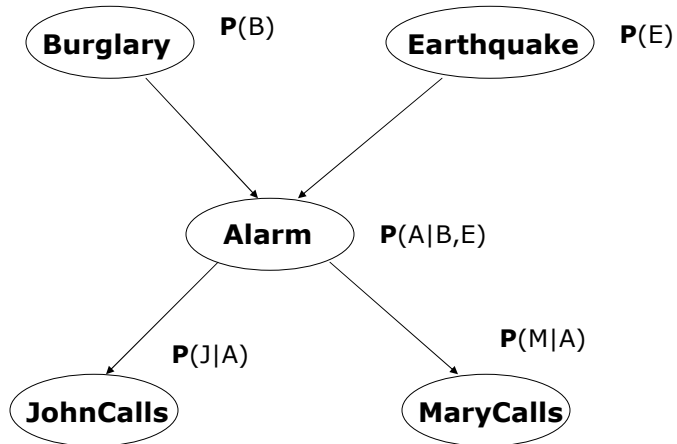
- Assume your house has an **alarm system** against **burglary**. You live in the seismically active area and the alarm system can get occasionally set off by an **earthquake**. You have two neighbors, **Mary** and **John**, who do not know each other. If they hear the alarm they call you, but this is not guaranteed.
- We want to represent the probability distribution of events:
 - Burglary, Earthquake, Alarm, Mary calls and John calls

Causal relations



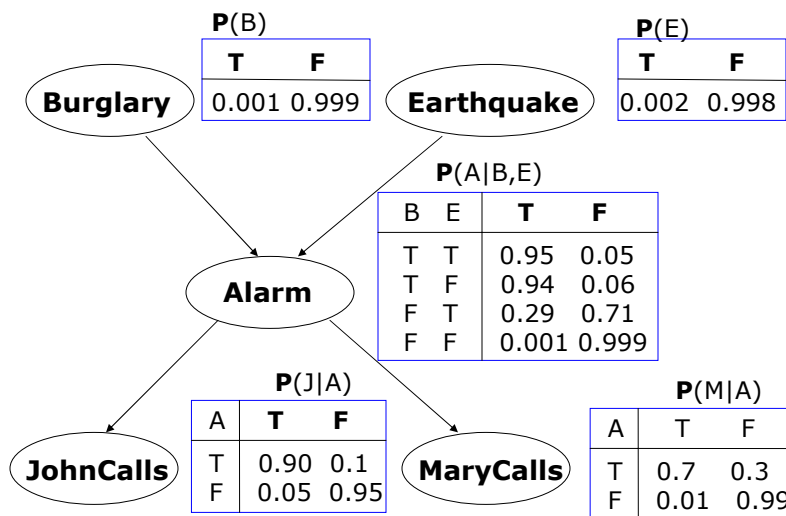
Bayesian belief network.

1. Graph reflecting direct (causal) dependencies between variables
2. Local conditional distributions relating variables and their parents



CS 3750 Machine Learning

Bayesian belief network.



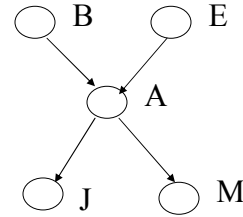
CS 3750 Machine Learning

Bayesian belief networks (general)

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

- **Directed acyclic graph**

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



- **Parameters**

- Local conditional probability distributions for every variable-parent configuration

$$P(X_i \mid pa(X_i))$$

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 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**

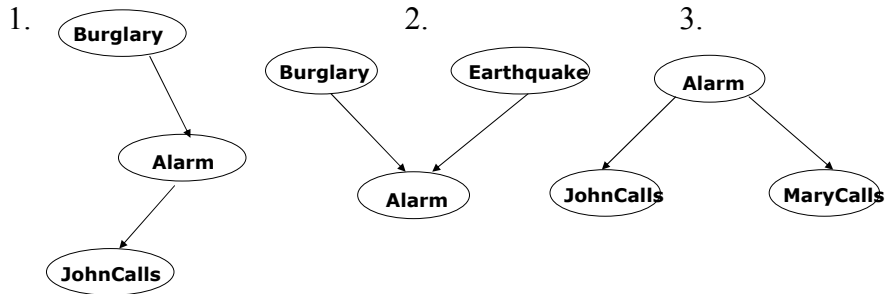
$$P(A \mid C, B) = P(A \mid C)$$

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

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

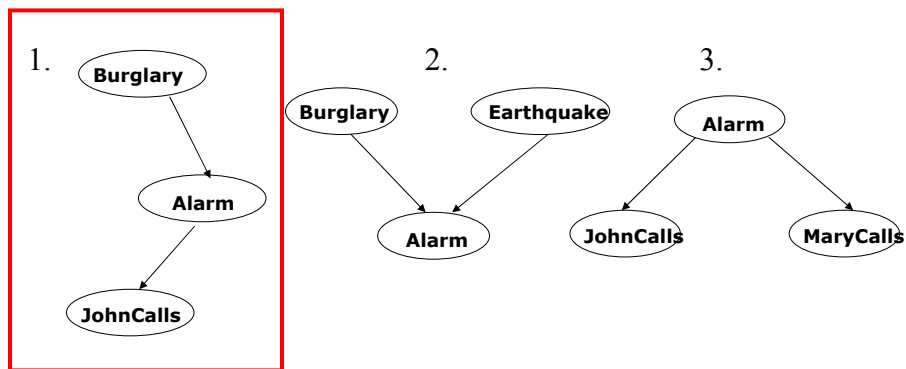
Independences in BBNs

3 basic independence structures:



CS 3750 Machine Learning

Independences in BBNs



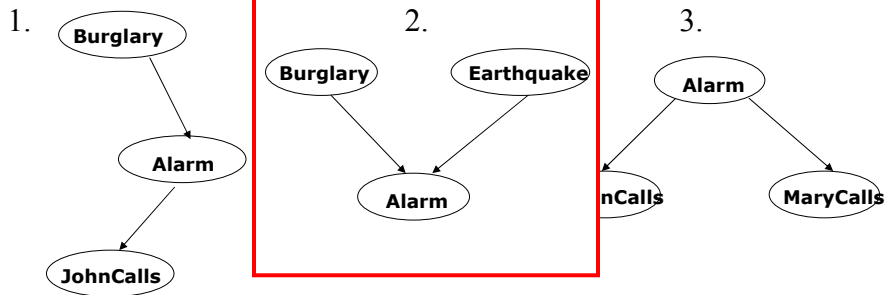
1. JohnCalls is **independent** of Burglary given Alarm

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

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

CS 3750 Machine Learning

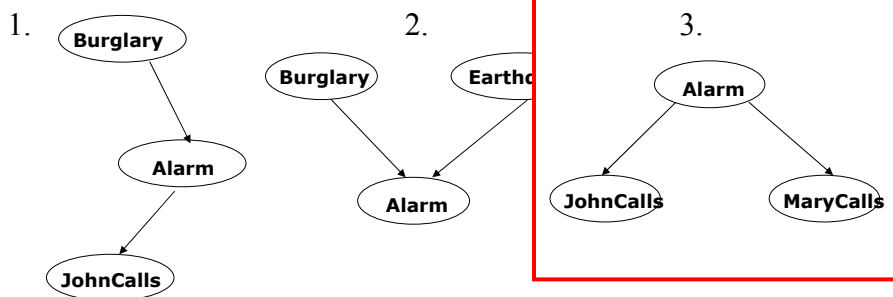
Independences in BBNs



2. Burglary **is independent** of Earthquake (not knowing Alarm)
Burglary and Earthquake **become dependent** given Alarm !!

$$P(B, E) = P(B)P(E)$$

Independences in BBNs



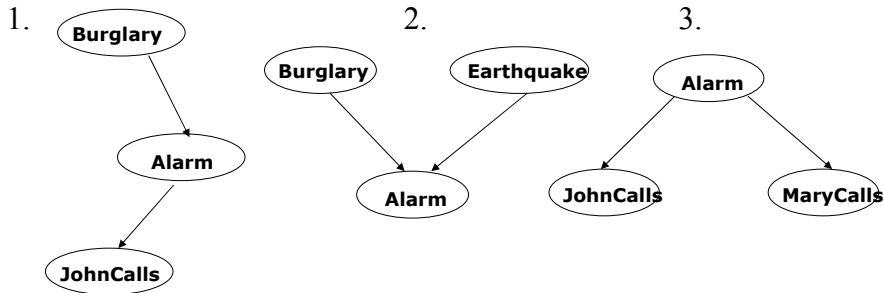
3. MaryCalls **is independent** of JohnCalls given Alarm

$$P(J \mid A, M) = P(J \mid A)$$

$$P(J, M \mid A) = P(J \mid A)P(M \mid A)$$

Independences in BBNs

3 basic independence structures encoded in BBNs



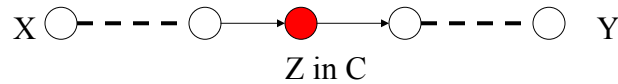
1. JohnCalls **is independent** of Burglary given Alarm
2. Burglary **is independent** of Earthquake (not knowing Alarm)
Burglary and Earthquake **become dependent** given Alarm !!
3. MaryCalls **is independent** of JohnCalls given Alarm

Independences in BBN

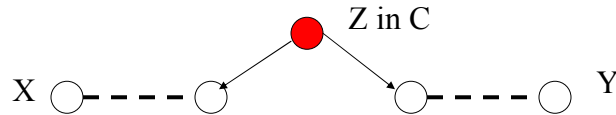
- BBN distribution models many conditional independence relations relating distant variables and sets
- These are defined in terms of the graphical criterion called d-separation
- **D-separation in the graph**
 - Let X, Y and Z be three sets of nodes
 - If X and Y are d-separated by Z then X and Y are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

Undirected path blocking

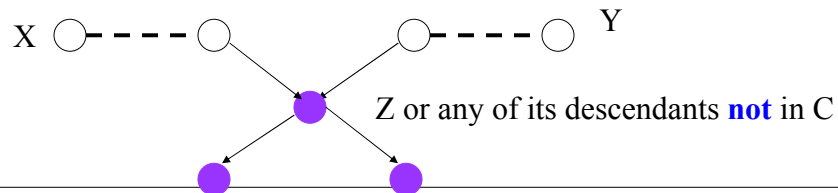
- 1. With linear substructure



- 2. With wedge substructure

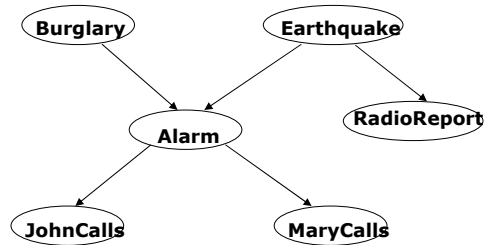


- 3. With vee substructure



CS 3750 Machine Learning

Independences in BBNs



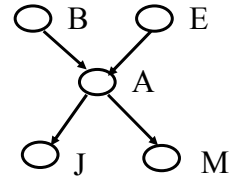
- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **F**
- Burglary and RadioReport are independent given Earthquake **T**
- Burglary and RadioReport are independent given MaryCalls **F**

CS 3750 Machine Learning

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

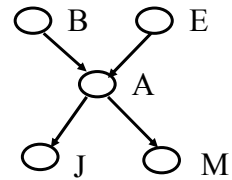
$$P(B=T, E=T, A=T, J=T, M=F) =$$



Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

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

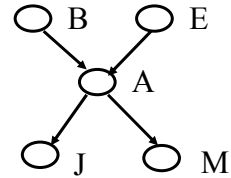


$$= P(J=T \mid B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F)$$

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

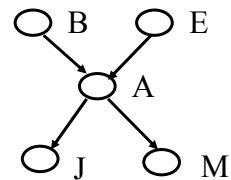
$$= \underline{P(J=T \mid A=T)} P(B=T, E=T, A=T, M=F)$$

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

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

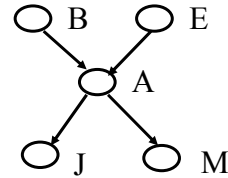
$$P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T)$$

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

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

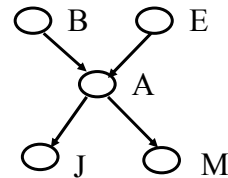


$$\begin{aligned}
 P(B=T, E=T, A=T, J=T, M=F) &= \\
 &= P(J=T \mid B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F) \\
 &= \underline{P(J=T \mid A=T)} P(B=T, E=T, A=T, M=F) \\
 &\quad P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T) \\
 &\quad \underline{P(M=F \mid A=T)} P(B=T, E=T, A=T) \\
 &\quad \underline{P(A=T \mid B=T, E=T)} P(B=T, E=T) \\
 &\quad P(B=T) P(E=T)
 \end{aligned}$$

CS 3750 Machine Learning

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



$$\begin{aligned}
 P(B=T, E=T, A=T, J=T, M=F) &= \\
 &= P(J=T \mid B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F) \\
 &= \underline{P(J=T \mid A=T)} P(B=T, E=T, A=T, M=F) \\
 &\quad P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T) \\
 &\quad \underline{P(M=F \mid A=T)} P(B=T, E=T, A=T) \\
 &\quad \underline{P(A=T \mid B=T, E=T)} P(B=T, E=T) \\
 &\quad P(B=T) P(E=T) \\
 &= P(J=T \mid A=T) P(M=F \mid A=T) P(A=T \mid B=T, E=T) P(B=T) P(E=T)
 \end{aligned}$$

CS 3750 Machine Learning

Parameter complexity problem

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

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

Parameters:

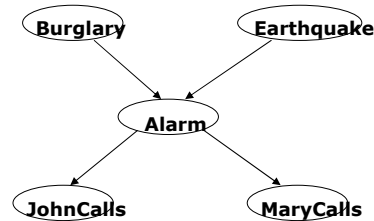
full joint: $2^5 = 32$

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 3750 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 3750 Machine Learning

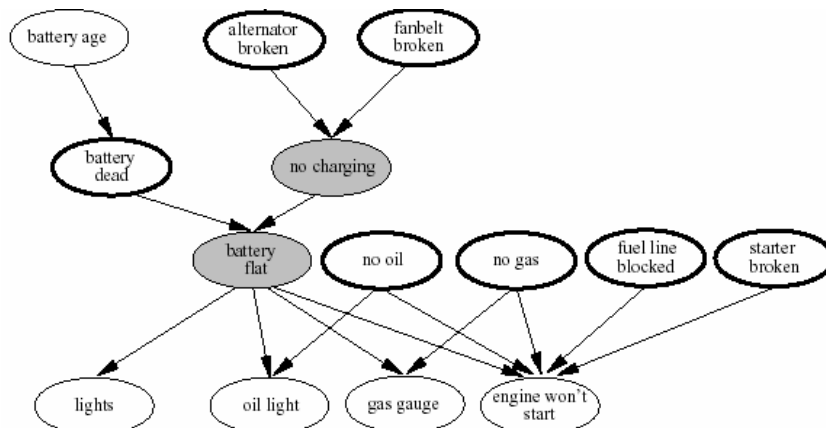
BBNs built in practice

- 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 3750 Machine Learning

Diagnosis of car engine

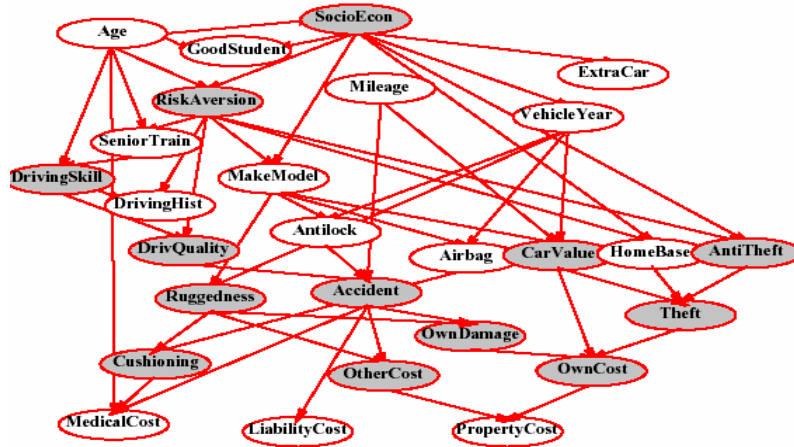
- Diagnose the engine start problem



CS 3750 Machine Learning

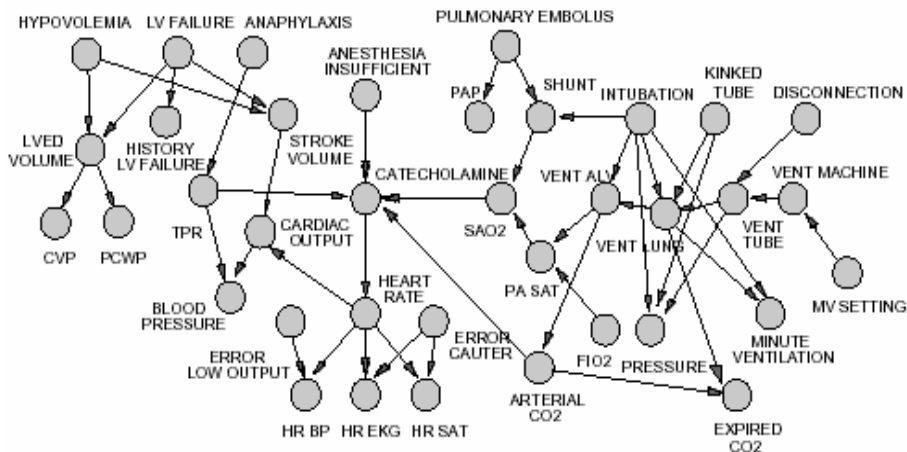
Car insurance example

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



CS 3750 Machine Learning

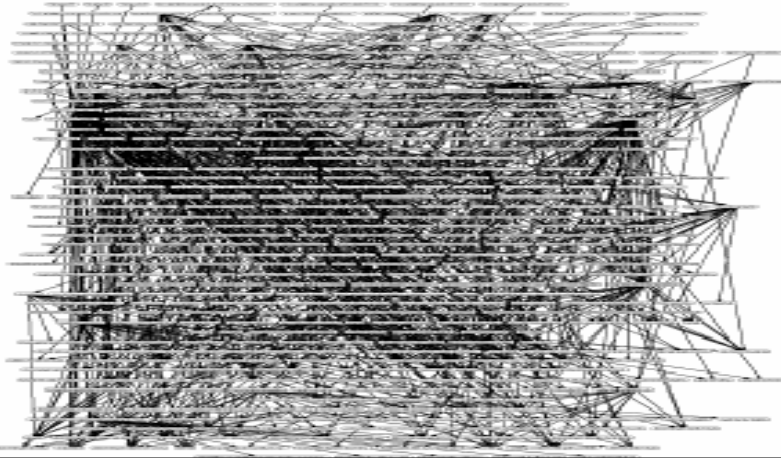
(ICU) Alarm network



CS 3750 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

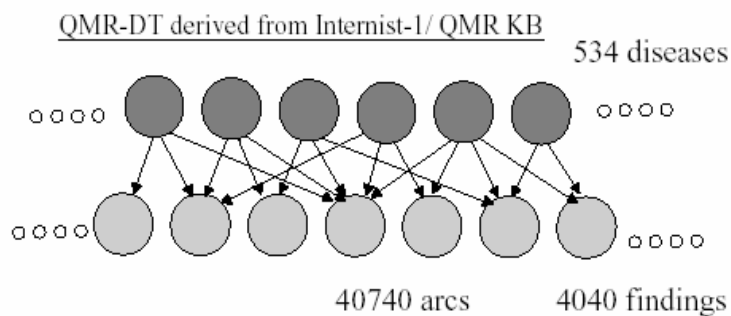


CS 3750 Machine Learning

QMR-DT

- Medical diagnosis in internal medicine

Bipartite network of disease/findings relations



CS 3750 Machine Learning

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)**

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

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

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

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

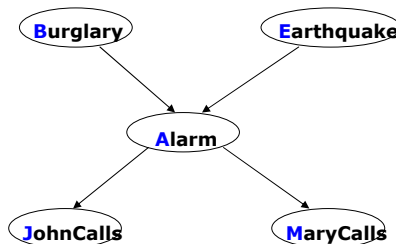
$$P(\text{Alarm})$$

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

CS 3750 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 3750 Machine Learning

Inference in Bayesian networks

Computing: $P(J = T)$

Approach 1. Blind approach.

- Sum out all uninstantiated 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: $16 * 4 = \mathbf{64}$

CS 3750 Machine Learning

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: $1 + 2 * (1) + 2 * (1 + 2 * (1)) = \mathbf{9}$

Number of products: $2 * (2 + 2 * (1) + 2 * (2 * (1))) = \mathbf{16}$

CS 3750 Machine Learning

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

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

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

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

- 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:

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

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

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

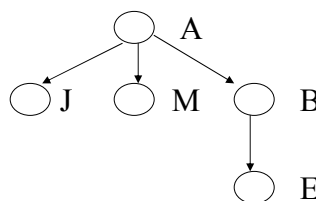
- **General technique: Variable elimination**

Inference in Bayesian networks

- General idea of variable elimination

$$\begin{aligned}
 P(\text{True}) &= 1 = \\
 &= \sum_{a \in T, F} \underbrace{\left[\sum_{j \in T, F} P(J=j \mid A=a) \right]}_{f_J(a)} \underbrace{\left[\sum_{m \in T, F} P(M=m \mid A=a) \right]}_{f_M(a)} \underbrace{\left[\sum_{b \in T, F} P(B=b) \left[\sum_{e \in T, F} P(A=a \mid B=b, E=e) P(E=e) \right] \right]}_{f_E(a,b)} \\
 &\hspace{15em} \underbrace{\hspace{10em}}_{f_B(a)}
 \end{aligned}$$

Variable order:



Results cashed in the tree structure

Inferences in Bayesian network

- **Exact inference algorithms:**
 - 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