

CS 2710 Foundations of AI

Lecture 16

Bayesian belief networks

Milos Hauskrecht

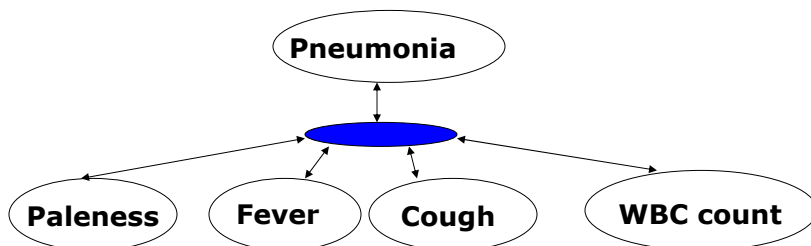
milos@cs.pitt.edu

5329 Sennott Square

CS 2710 Foundations of AI

Uncertainty

To make diagnostic inference possible we need to represent knowledge (axioms) that relate symptoms and diagnosis



Problem: disease/symptoms relations are not deterministic

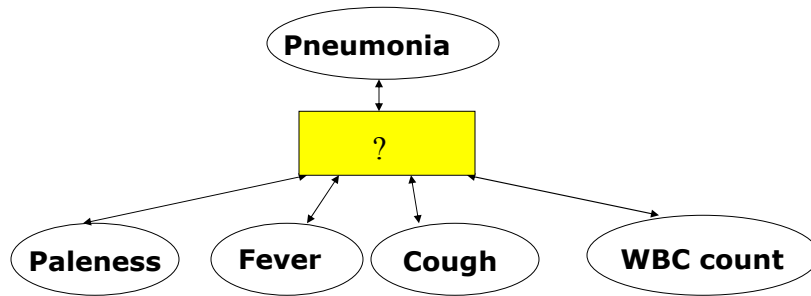
- They are **uncertain (or stochastic)** and vary from patient to patient

CS 2710 Foundations of AI

Modeling the uncertainty.

Key challenges:

- How to represent the relations in the presence of uncertainty?
- How to manipulate such knowledge to make inferences?
 - **Humans can reason with uncertainty.**



CS 2710 Foundations of AI

Methods for representing uncertainty

Probability theory

- A well defined theory for modeling and reasoning in the presence of uncertainty
- A natural choice to replace certainty factors

Facts (propositional statements)

- Are represented via **random variables** with two or more values

Example: *Pneumonia* is a random variable

values: *True* and *False*

- Each value can be achieved **with some probability:**

$$P(Pneumonia = True) = 0.001$$

$$P(WBCcount = high) = 0.005$$

CS 2710 Foundations of AI

Modeling uncertainty with probabilities

Probabilistic extension of propositional logic.

- **Propositions:**

- statements about the world
- Represented by the assignment of values to **random variables**

- **Random variables:**

- ! – **Boolean** *Pneumonia* is either *True, False*
Random variable Values
- ! – **Multi-valued** *Pain* is one of {*Nopain, Mild, Moderate, Severe*}
Random variable Values
- **Continuous** *HeartRate* is a value in $< 0 ; 250 >$
Random variable Values

Probabilities

Unconditional probabilities (prior probabilities)

$$P(\text{Pneumonia}) = 0.001 \quad \text{or} \quad P(\text{Pneumonia} = \text{True}) = 0.001$$

$$P(\text{Pneumonia} = \text{False}) = 0.999$$

$$P(\text{WBCcount} = \text{high}) = 0.005$$

Probability distribution

- Defines probabilities **for all possible value assignments to a random variable**
- Values are mutually exclusive

$$P(\text{Pneumonia} = \text{True}) = 0.001$$

$$P(\text{Pneumonia} = \text{False}) = 0.999$$

<i>Pneumonia</i>	P(Pneumonia)
<i>True</i>	0.001
<i>False</i>	0.999

Probability distribution

Defines probability for **all possible value assignments**

Example 1:

$$P(\text{Pneumonia} = \text{True}) = 0.001$$

$$P(\text{Pneumonia} = \text{False}) = 0.999$$

<i>Pneumonia</i>	P(<i>Pneumonia</i>)
<i>True</i>	0.001
<i>False</i>	0.999

$$P(\text{Pneumonia} = \text{True}) + P(\text{Pneumonia} = \text{False}) = 1$$

Probabilities sum to 1 !!!

Example 2:

$$P(\text{WBCcount} = \text{high}) = 0.005$$

$$P(\text{WBCcount} = \text{normal}) = 0.993$$

$$P(\text{WBCcount} = \text{low}) = 0.002$$

<i>WBCcount</i>	P(<i>WBCcount</i>)
<i>high</i>	0.005
<i>normal</i>	0.993
<i>low</i>	0.002

Joint probability distribution

Joint probability distribution (for a set variables)

- Defines probabilities for **all possible assignments of values to variables in the set**

Example: variables *Pneumonia* and *WBCcount*

$$\mathbf{P}(\text{pneumonia}, \text{WBCcount})$$

Is represented by 2×3 matrix

		<i>WBCcount</i>		
		<i>high</i>	<i>normal</i>	<i>low</i>
<i>Pneumonia</i>	<i>True</i>	0.0008	0.0001	0.0001
	<i>False</i>	0.0042	0.9929	0.0019

Joint probabilities

Marginalization

- reduces the dimension of the joint distribution
- Sums variables out

$P(\text{pneumonia}, \text{WBCcount})$ 2×3 matrix

		WBCcount			
		high	normal	low	
Pneumonia	True	0.0008	0.0001	0.0001	$P(\text{Pneumonia})$ 0.001 0.999
	False	0.0042	0.9929	0.0019	
		0.005	0.993	0.002	

$P(\text{WBCcount})$

Marginalization (here summing of columns or rows)

Full joint distribution

- **the joint distribution for all variables in the problem**
 - It defines the complete probability model for the problem

Example: pneumonia diagnosis

Variables: *Pneumonia*, *Fever*, *Paleness*, *WBCcount*, *Cough*

Full joint defines the probability for all possible assignments of values to *Pneumonia*, *Fever*, *Paleness*, *WBCcount*, *Cough*

$P(\text{Pneumonia}=T, \text{WBCcount}= \text{High}, \text{Fever}=T, \text{Cough}=T, \text{Paleness}=T)$

$P(\text{Pneumonia}=T, \text{WBCcount}= \text{High}, \text{Fever}=T, \text{Cough}=T, \text{Paleness}=F)$

$P(\text{Pneumonia}=T, \text{WBCcount}= \text{High}, \text{Fever}=T, \text{Cough}=F, \text{Paleness}=T)$

... etc

Conditional probabilities

Conditional probability distribution

- Defines probabilities for all possible assignments, given a fixed assignment to some other variable values

$$P(\text{Pneumonia} = \text{true} \mid \text{WBCcount} = \text{high})$$

$\mathbf{P}(\text{Pneumonia} \mid \text{WBCcount})$ 3 element vector of 2 elements

		<i>WBCcount</i>		
		<i>high</i>	<i>normal</i>	<i>low</i>
<i>Pneumonia</i>	<i>True</i>	0.08	0.0001	0.0001
	<i>False</i>	0.92	0.9999	0.9999
		1.0	1.0	1.0

$$P(\text{Pneumonia} = \text{true} \mid \text{WBCcount} = \text{high})$$

$$+ P(\text{Pneumonia} = \text{false} \mid \text{WBCcount} = \text{high})$$

CS 2710 Foundations of AI

Conditional probabilities

Conditional probability

- Is defined in terms of the joint probability:

$$P(A \mid B) = \frac{P(A, B)}{P(B)} \text{ s.t. } P(B) \neq 0$$

- Example:**

$$P(\text{pneumonia} = \text{true} \mid \text{WBCcount} = \text{high}) = \frac{P(\text{pneumonia} = \text{true}, \text{WBCcount} = \text{high})}{P(\text{WBCcount} = \text{high})}$$

$$P(\text{pneumonia} = \text{false} \mid \text{WBCcount} = \text{high}) = \frac{P(\text{pneumonia} = \text{false}, \text{WBCcount} = \text{high})}{P(\text{WBCcount} = \text{high})}$$

CS 2710 Foundations of AI

Conditional probabilities

- **Conditional probability distribution.**

$$P(A | B) = \frac{P(A, B)}{P(B)} \text{ s.t. } P(B) \neq 0$$

- **Product rule.** Joint probability can be expressed in terms of conditional probabilities

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

- **Chain rule.** Any joint probability can be expressed as a product of conditionals

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_n | X_1, \dots, X_{n-1})P(X_1, \dots, X_{n-1}) \\ &= P(X_n | X_1, \dots, X_{n-1})P(X_{n-1} | X_1, \dots, X_{n-2})P(X_1, \dots, X_{n-2}) \\ &= \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1}) \end{aligned}$$

Bayes rule

Conditional probability.

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

Bayes rule:

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

When is it useful?

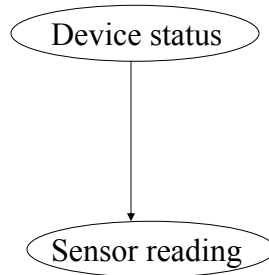
- When we are interested in computing the diagnostic query from the causal probability

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

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

Bayes Rule in 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(Device status)

normal	malfunctioning
0.9	0.1

P(Sensor reading | Device status)

Device\Sensor	high	low
normal	0.1	0.9
malfunctioning	0.6	0.4

Bayes Rule in a simple diagnostic inference.

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

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

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

- Note that typically the opposite conditional probabilities are given to us: they are much easier to estimate
- **Solution:** apply **Bayes rule** to reverse the conditioning variables

Probabilistic inference

Various inference tasks:

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

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

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

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

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

$$\mathbf{P}(Fever)$$

$$\mathbf{P}(Fever, ChestPain)$$

Inference

Any query can be computed from the full joint distribution !!!

- **Joint over a subset of variables** is obtained through marginalization

$$P(A = a, C = c) = \sum_i \sum_j P(A = a, B = b_i, C = c, D = d_j)$$

- **Conditional probability over set of variables**, given other variables' values is obtained through marginalization and definition of conditionals

$$\begin{aligned} P(D = d \mid A = a, C = c) &= \frac{P(A = a, C = c, D = d)}{P(A = a, C = c)} \\ &= \frac{\sum_i P(A = a, B = b_i, C = c, D = d)}{\sum_i \sum_j P(A = a, B = b_i, C = c, D = d_j)} \end{aligned}$$

Inference.

Any query can be computed from the full joint distribution !!!

- Any joint probability can be expressed as a product of conditionals via the **chain rule**.

$$\begin{aligned}P(X_1, X_2, \dots, X_n) &= P(X_n | X_1, \dots, X_{n-1})P(X_1, \dots, X_{n-1}) \\&= P(X_n | X_1, \dots, X_{n-1})P(X_{n-1} | X_1, \dots, X_{n-2})P(X_1, \dots, X_{n-2}) \\&= \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})\end{aligned}$$

- Sometimes it is easier to define the distribution in terms of conditional probabilities:
 - E.g. $\mathbf{P}(\text{Fever} | \text{Pneumonia} = T)$
 $\mathbf{P}(\text{Fever} | \text{Pneumonia} = F)$

Modeling uncertainty with probabilities

- Defining the **full joint distribution** makes it possible to represent and reason with uncertainty in a uniform way
- We are able to handle an arbitrary inference problem

Problems:

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

Medical diagnosis example.

- **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 marginal of $P(\text{Pneumonia}=T)$ 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

Modeling uncertainty with probabilities

- **Knowledge based system era (70s – early 80's)**

- **Extensional non-probabilistic models**
- Solve the space, time and acquisition bottlenecks in probability-based models
- froze the development and advancement of KB systems and contributed to the slow-down of AI in 80s in general

- Breakthrough (late 80s, beginning of 90s)

- **Bayesian belief networks**

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

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**

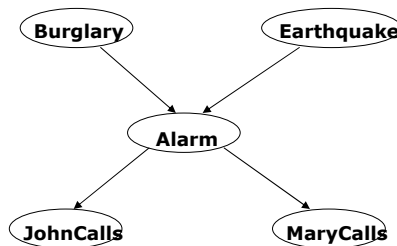
$$P(A, B | C) = P(A | C)P(B | C)$$

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

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

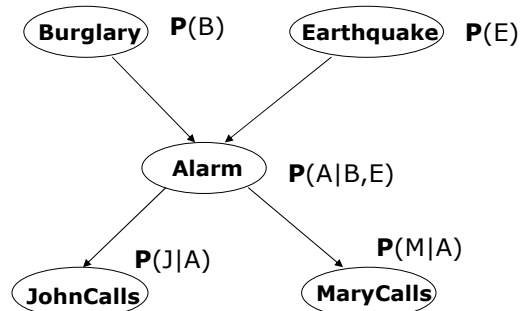


Bayesian belief network.

1. Directed acyclic graph

- **Nodes** = random variables
Burglary, Earthquake, Alarm, Mary calls and John calls
- **Links** = direct (causal) dependencies between variables.

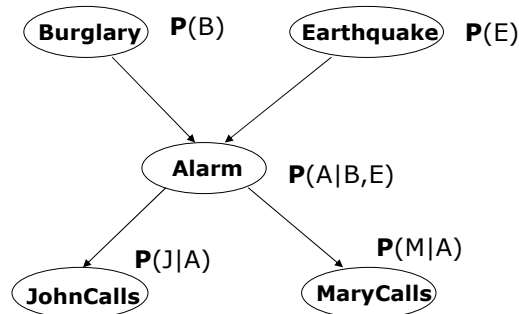
The chance of Alarm is influenced by Earthquake, The chance of John calling is affected by the Alarm



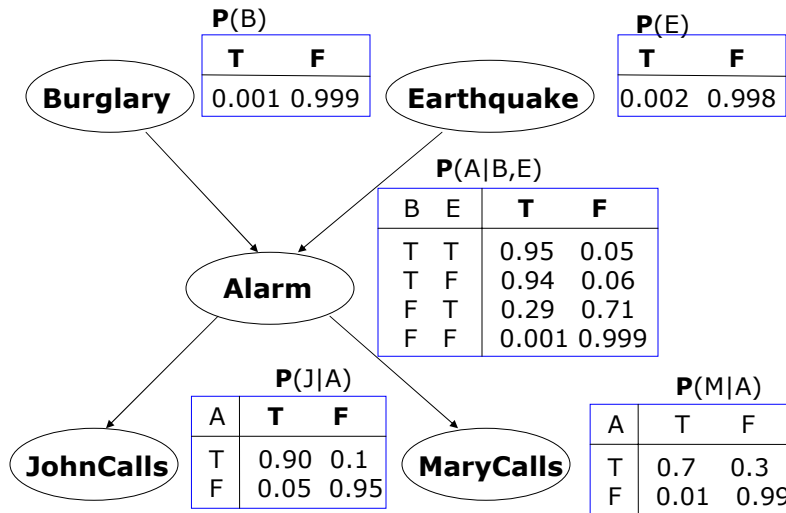
Bayesian belief network.

2. Local conditional distributions

- relate variables and their parents



Bayesian belief network.



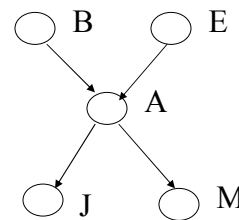
CS 2710 Foundations of AI

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

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

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

CS 2710 Foundations of AI

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:

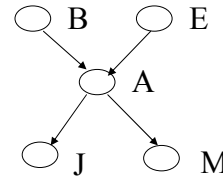
Assume the following assignment of values to random variables

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

Then its probability is:

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

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



Bayesian belief networks (BBNs)

Bayesian belief networks

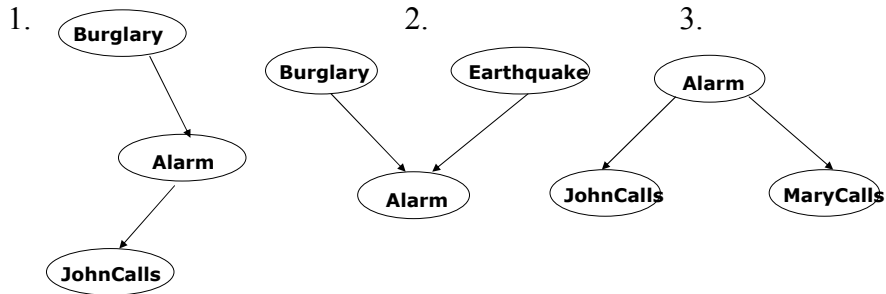
- Represent the full joint distribution over the variables more compactly using the product of local conditionals.
- **But how did we get to local parameterizations?**

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 !!!**

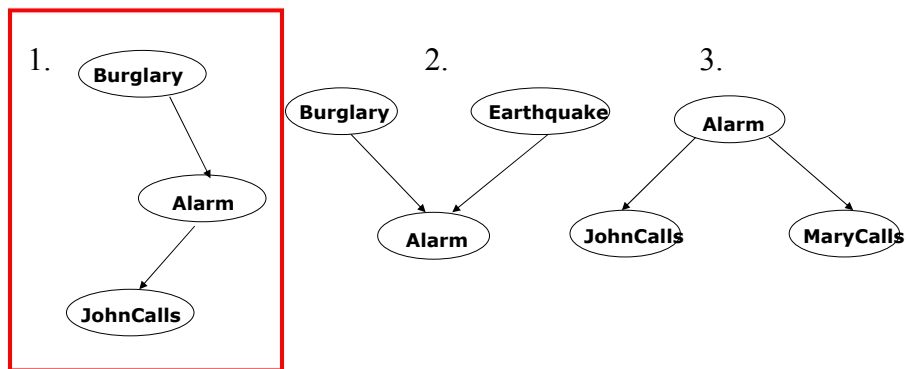
Independences in BBNs

3 basic independence structures:



CS 2710 Foundations of AI

Independences in BBNs



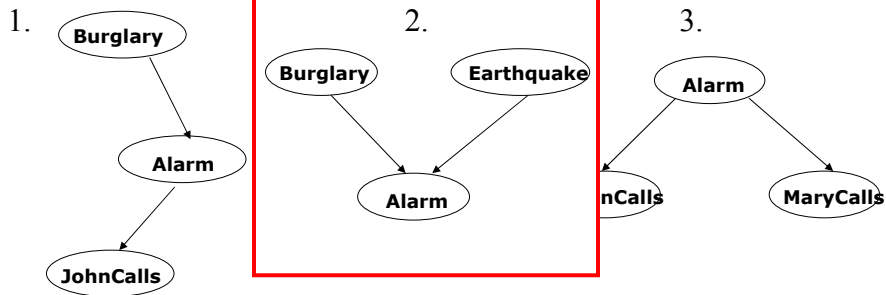
1. JohnCalls is **independent** of Burglary given Alarm

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

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

CS 2710 Foundations of AI

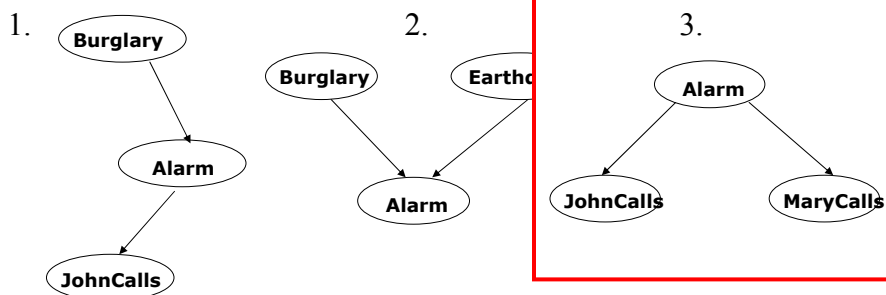
Independences in BBNs



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 BBN

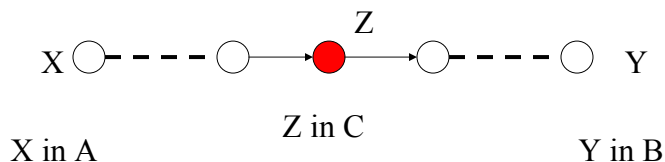
- BBN distribution models many conditional independence relations among distant variables and sets of variables
- These are defined in terms of the graphical criterion called d-separation
- **D-separation and independence**
 - Let X, Y and Z be three sets of nodes
 - If X and Y are d-separated by Z, 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 2710 Foundations of AI

Undirected path blocking

A is d-separated from B given C if every undirected path between them is **blocked**

- **1. Path blocking with a linear substructure**

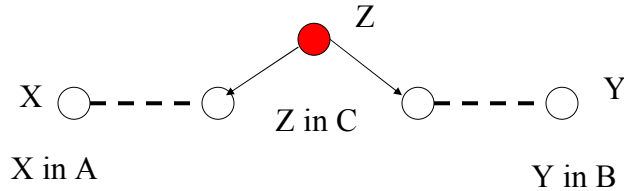


CS 2710 Foundations of AI

Undirected path blocking

A is d-separated from B given C if every undirected path between them is **blocked**

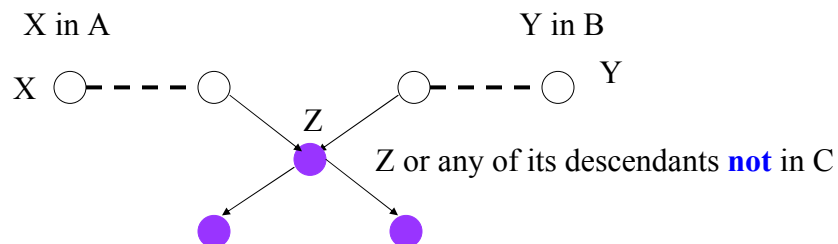
- 2. Path blocking with the wedge substructure



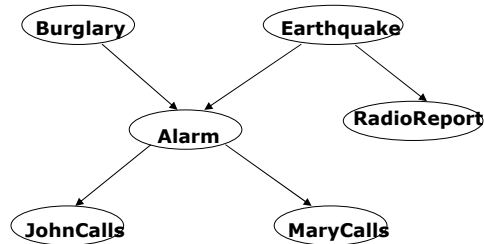
Undirected path blocking

A is d-separated from B given C if every undirected path between them is **blocked**

- 3. Path blocking with the vee substructure

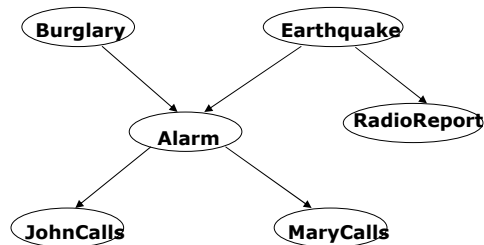


Independences in BBNs



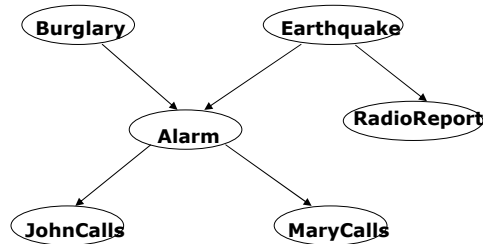
- Earthquake and Burglary are independent given MaryCalls ?

Independences in BBNs



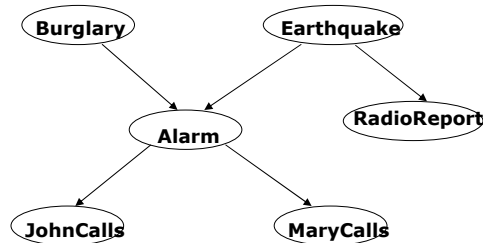
- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) ?

Independences in BBNs



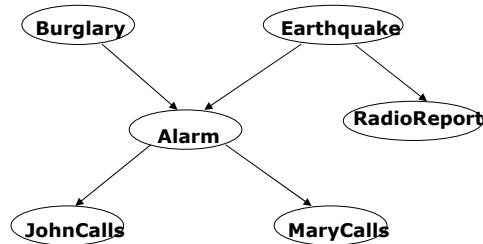
- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **F**
- Burglary and RadioReport are independent given Earthquake **?**

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 **?**

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**

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?**

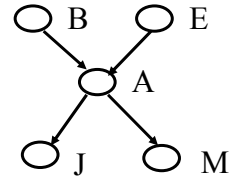
$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

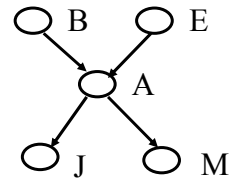
$$P(B=T, E=T, A=T, J=T, M=F) =$$



Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

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

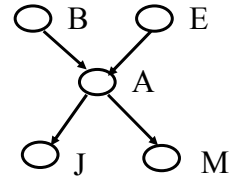


$$= P(J=T \mid B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F)$$

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

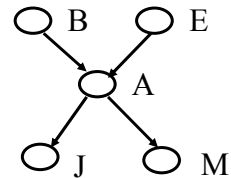
$$= \underline{P(J=T \mid A=T)} \underline{P(B=T, E=T, A=T, M=F)}$$

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

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

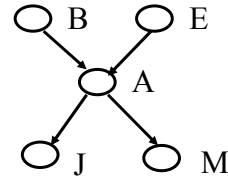
$$P(M=F \mid B=T, E=T, A=T) P(B=T, E=T, A=T)$$

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

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

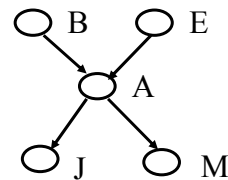
$$\underline{P(M=F \mid A=T)} \underline{P(B=T, E=T, A=T)}$$

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

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

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

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

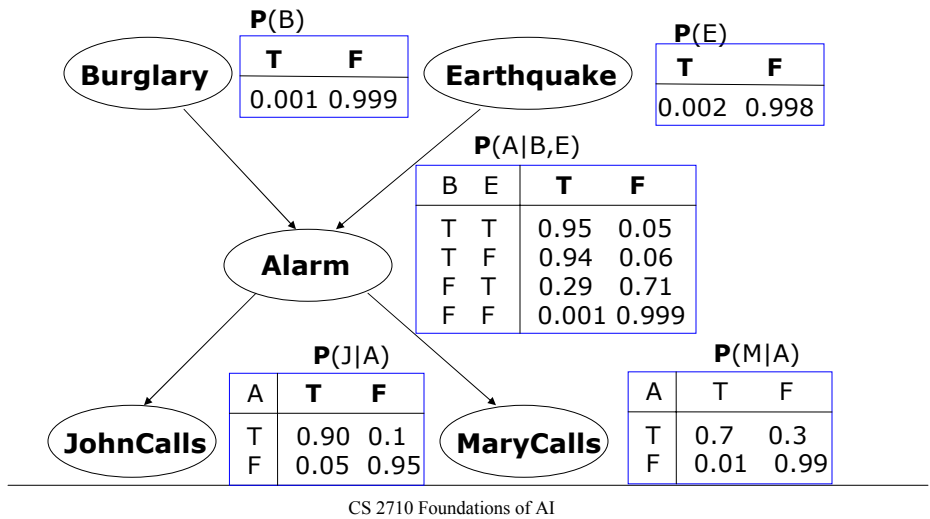
$$\underline{P(A=T \mid B=T, E=T)} \underline{P(B=T, E=T)}$$

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

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

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:

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

- What did we save?**

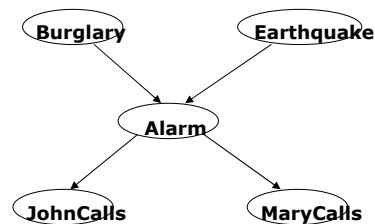
Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

$$2^5 - 1 = 31$$



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

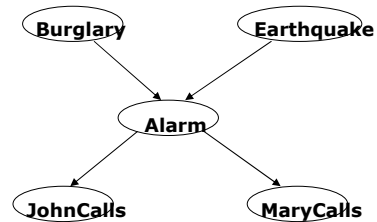
of parameters of the full joint:

$$2^5 = 32$$

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

