

CS 2750 Machine Learning

Lecture 13

Bayesian belief networks

Milos Hauskrecht

milos@cs.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

Attributes:

- modeled by random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$ with:
 - **Continuous values**
 - **Discrete values**

E.g. *blood pressure* with numerical values

or *chest pain* with discrete values

[no-pain, mild, moderate, strong]

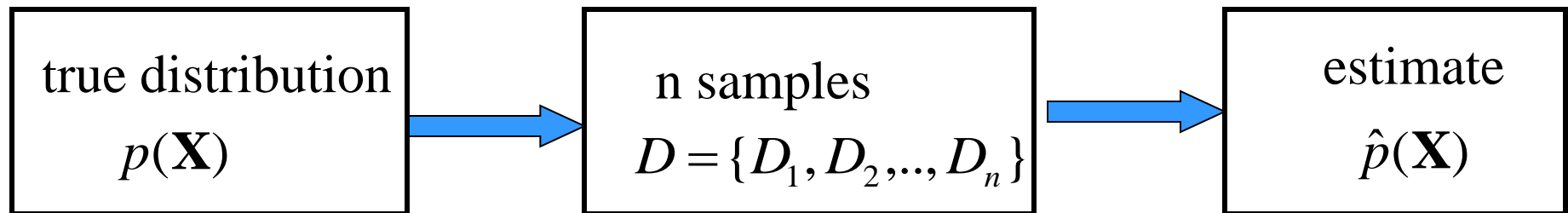
Underlying true probability distribution:

$$p(\mathbf{X})$$

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 **(identical) distribution** (fixed $p(\mathbf{X})$)

How to learn complex distributions

How to learn complex multivariate distributions $\hat{p}(\mathbf{X})$ with large number of variables?

One solution:

- **Decompose the distribution using conditional independence relations**
- **Decompose the parameter estimation problem to a set of smaller parameter estimation tasks**

Decomposition of distributions under conditional independence assumption is the main idea behind **Bayesian belief networks**

Example

Problem description:

- **Disease:** pneumonia
- **Patient symptoms (findings, lab tests):**
 - Fever, Cough, Paleness, WBC (white blood cells) count, Chest pain, etc.

Representation of a patient case:

- Symptoms and disease are represented as random variables

Our objectives:

- **Describe a multivariate distribution representing the relations between symptoms and disease**
- **Design of inference and learning procedures for the multivariate model**

Modeling uncertainty with probabilities

- **Full joint distribution:**

- Assume $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$ are all random variables that define the domain
- Full joint: $P(\mathbf{X})$ or $P(X_1, X_2, \dots, X_d)$

Full joint it is sufficient to do any type of probabilistic inference:

- Computation of joint probabilities for sets of variables

$$P(X_1, X_2, X_3) \quad P(X_1, X_{10})$$

- Computation of conditional probabilities

$$P(X_1 \mid X_2 = \text{True}, X_3 = \text{False})$$

Marginalization

Joint probability distribution (for a set variables)

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

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

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

$P(\text{WBCcount})$


Marginalization (summing of rows, or columns)
- summing out variables

Variable independence

- The joint distribution over a subset of variables can be always computed from the joint distribution through marginalization
- Not the other way around !!!
 - **Only exception:** when variables are independent

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

$P(pneumonia, WBCcount)$		$WBCcount$			$P(Pneumonia)$
		<i>high</i>	<i>normal</i>	<i>low</i>	
<i>Pneumonia</i>	<i>True</i>	0.0008	0.0001	0.0001	0.001
	<i>False</i>	0.0042	0.9929	0.0019	0.999
		0.005	0.993	0.002	

$P(WBCcount)$ 

Conditional probability

Conditional probability :

- Probability of A given B

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

- Conditional probability is defined in terms of joint probabilities
- Joint probabilities can be expressed in terms of conditional probabilities

$$P(A, B) = P(A | B)P(B) \quad \text{(product rule)}$$

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1}) \quad \text{(chain rule)}$$

- Conditional probability – is useful for **various probabilistic inferences**

$$P(Pneumonia = True | Fever = True, WBCcount = high, Cough = True)$$

Inference: joint distribution

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 a 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: Chain rule

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

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

- It is often easier to define the distribution in terms of conditional probabilities:
 - E.g. $\mathbf{P}(\textit{Fever} \mid \textit{Pneumonia} = T)$
 $\mathbf{P}(\textit{Fever} \mid \textit{Pneumonia} = F)$

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

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

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

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**

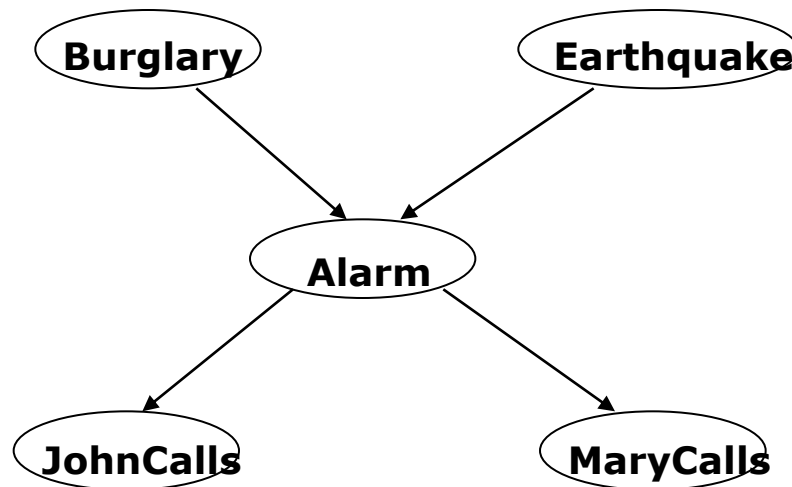
$$P(A, B | C) = P(A | C)P(B | C)$$

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

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

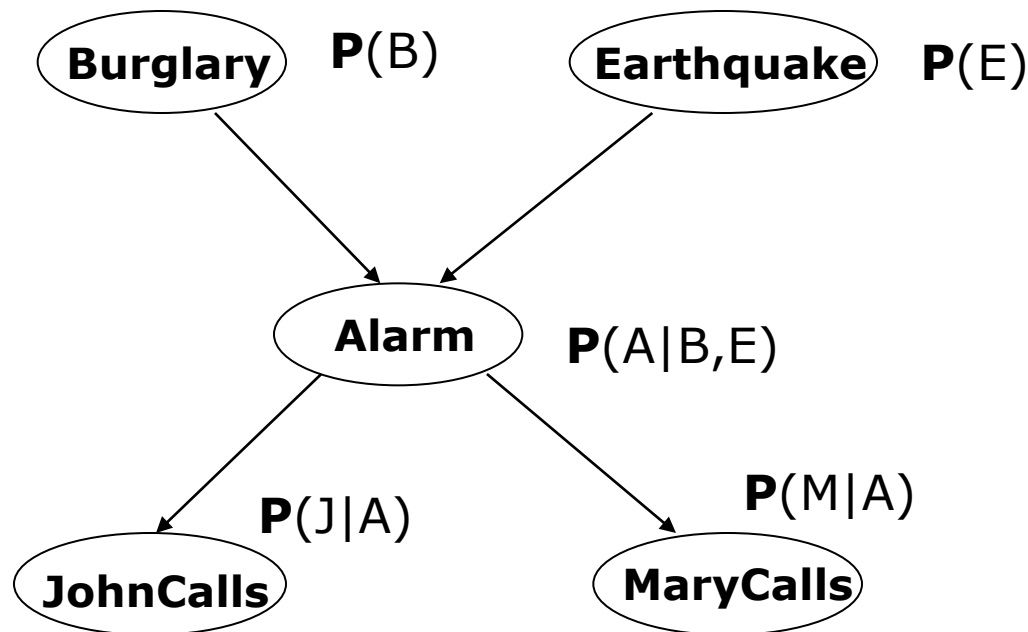


Bayesian belief network

1. Directed acyclic graph

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

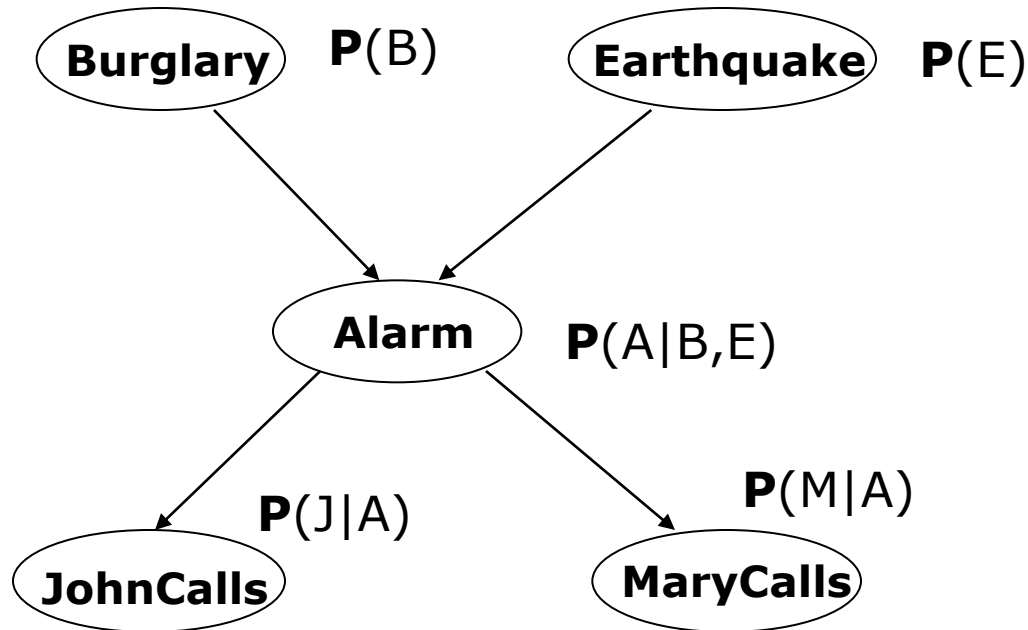
The chance of Alarm being is influenced by Earthquake,
The chance of John calling is affected by the Alarm



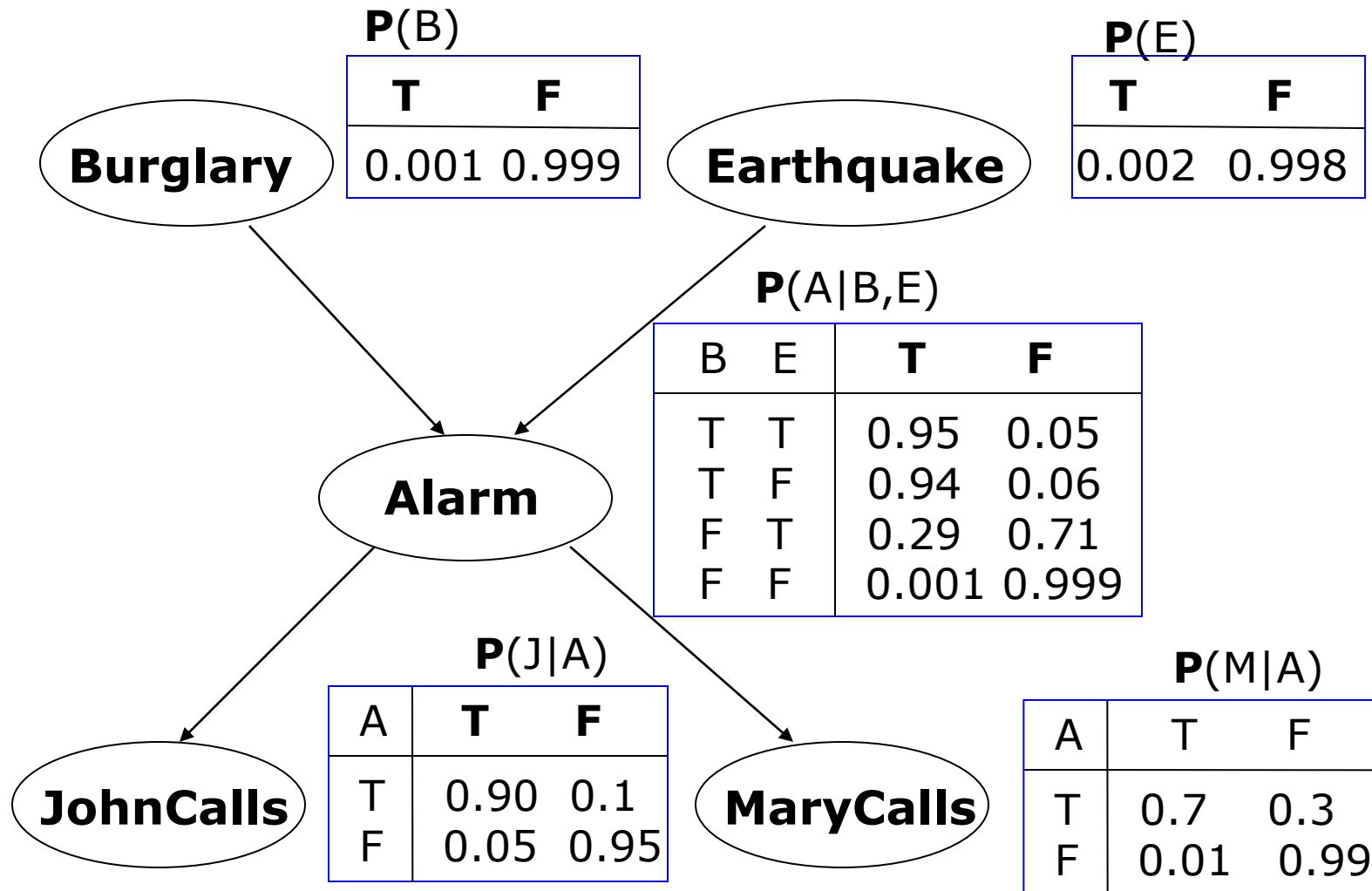
Bayesian belief network

2. Local conditional distributions

- relate variables and their parents



Bayesian belief network



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:

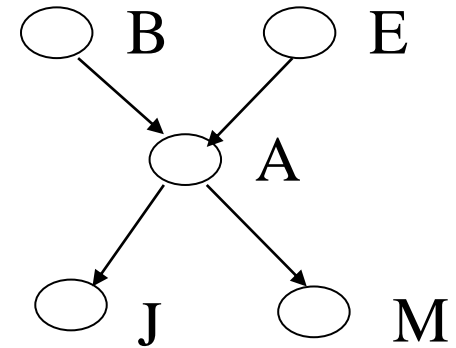
Assume the following assignment of values to random variables

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

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**

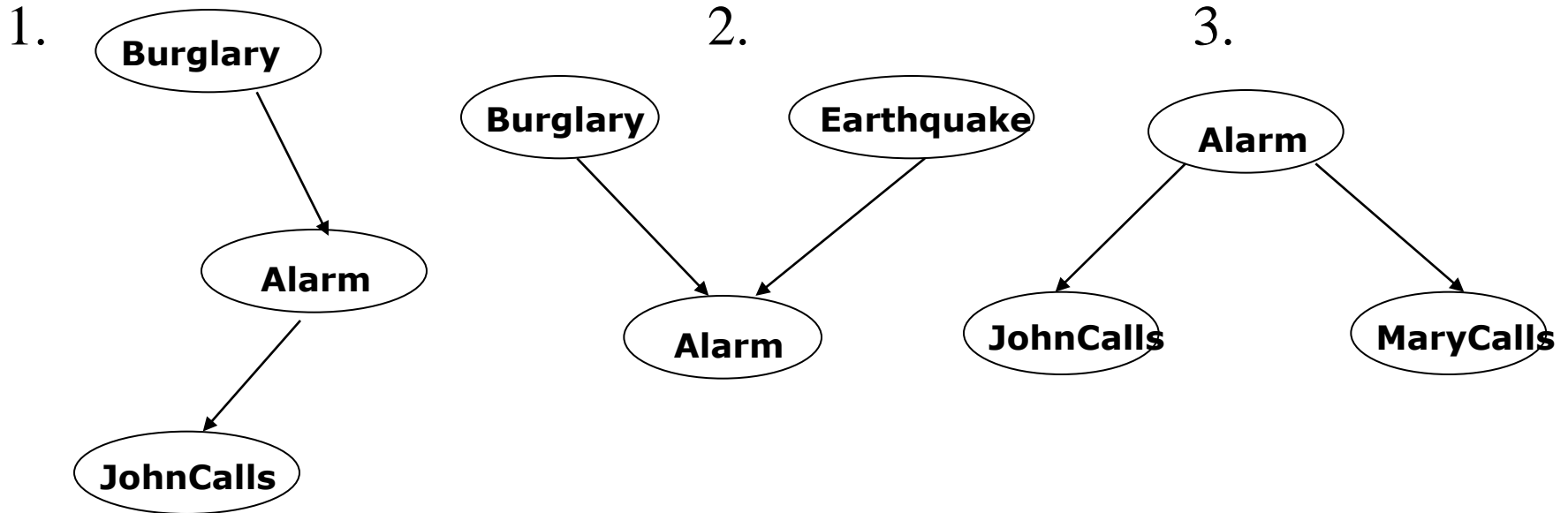
$$P(A | C, B) = P(A | C)$$

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

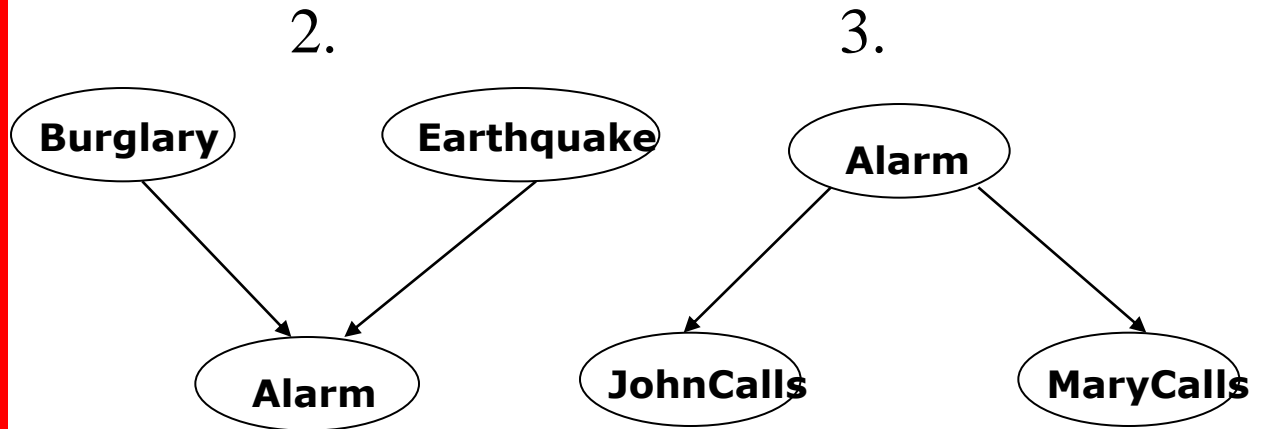
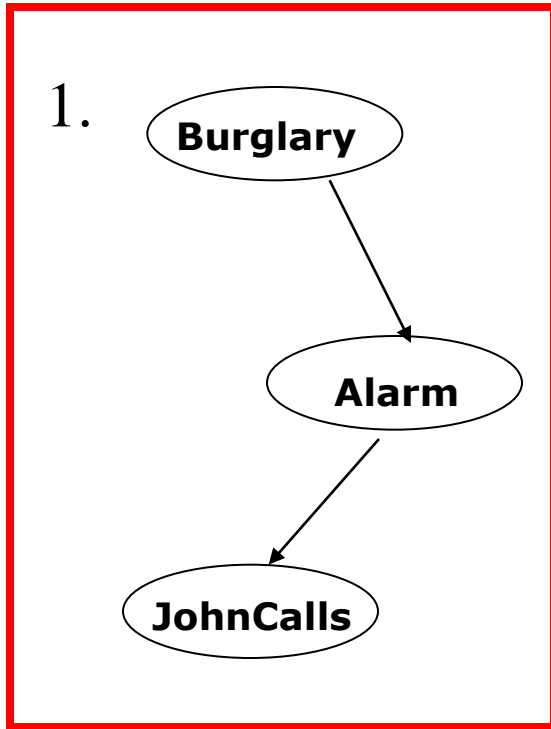
- **The graph structure implies the decomposition !!!**

Independences in BBNs

3 basic independence structures:



Independences in BBNs

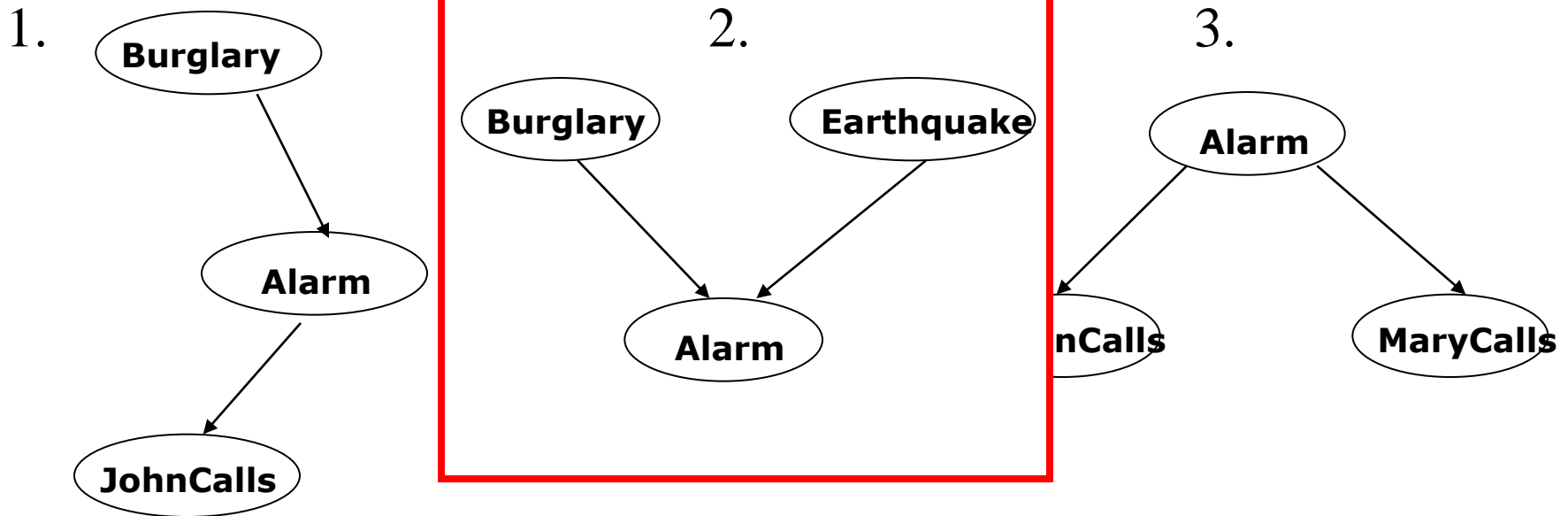


1. JohnCalls **is independent** of Burglary given Alarm

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

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

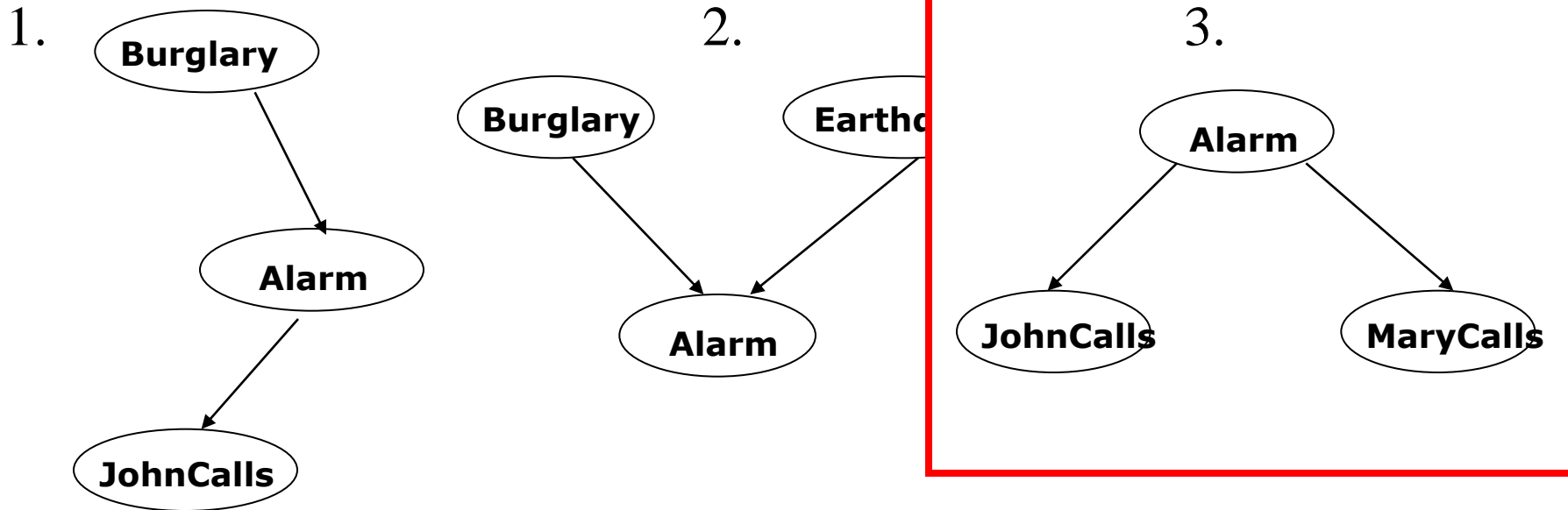
Independences in BBNs



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

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

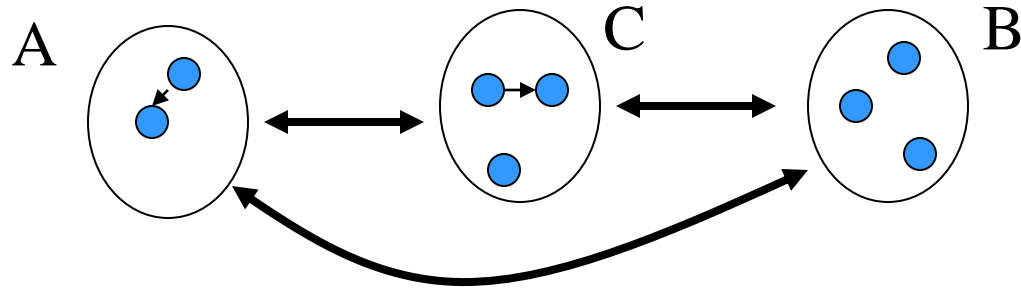
$$P(J, M \mid A) = P(J \mid A)P(M \mid A)$$

Independence 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

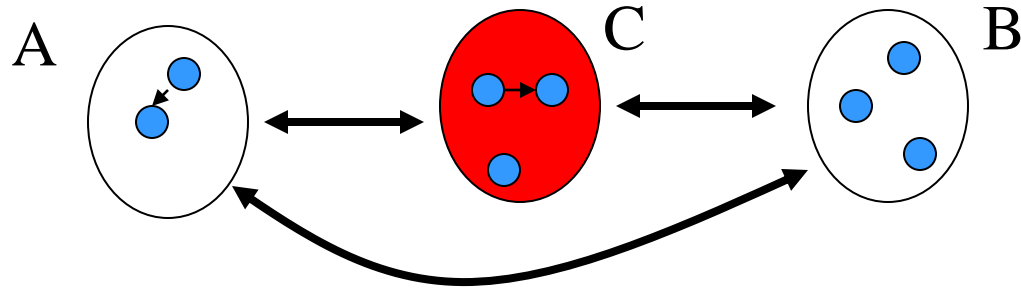
Undirected path blocking

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



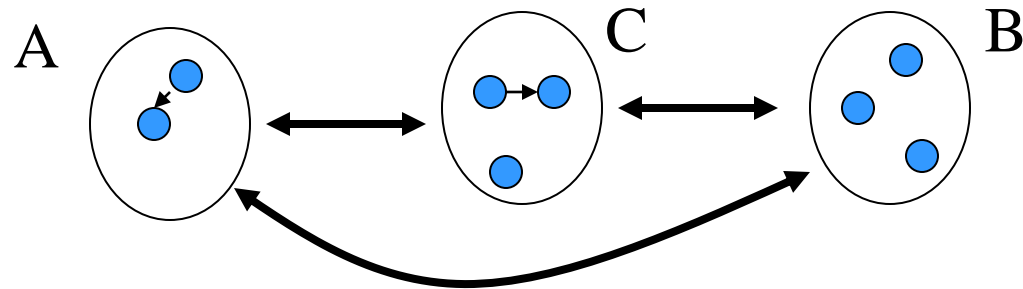
Undirected path blocking

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

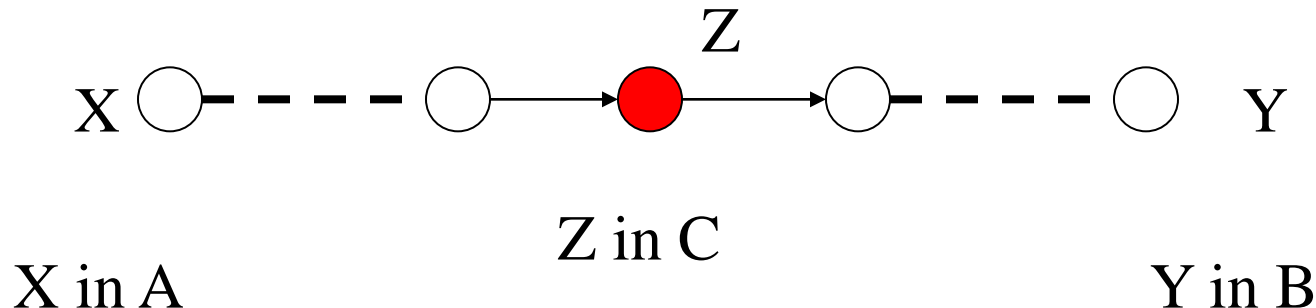


Undirected path blocking

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



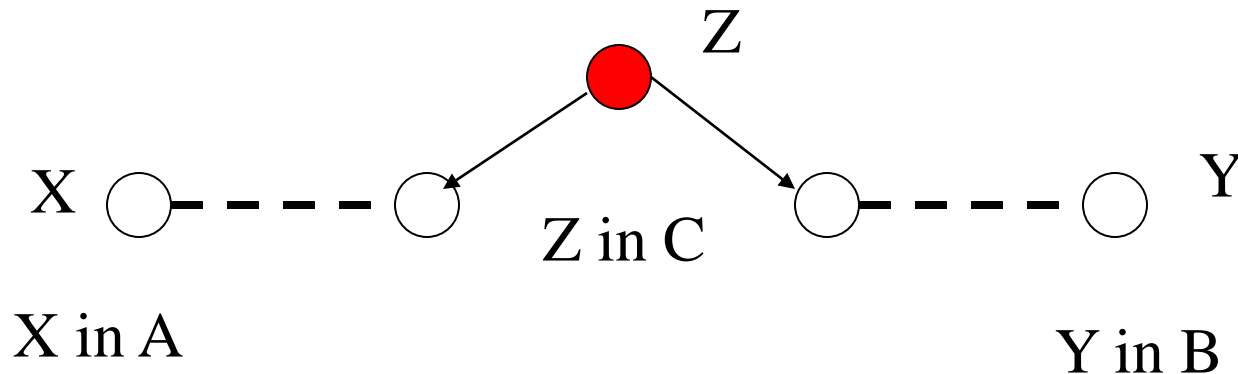
- 1. Path blocking with a linear substructure



Undirected path blocking

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

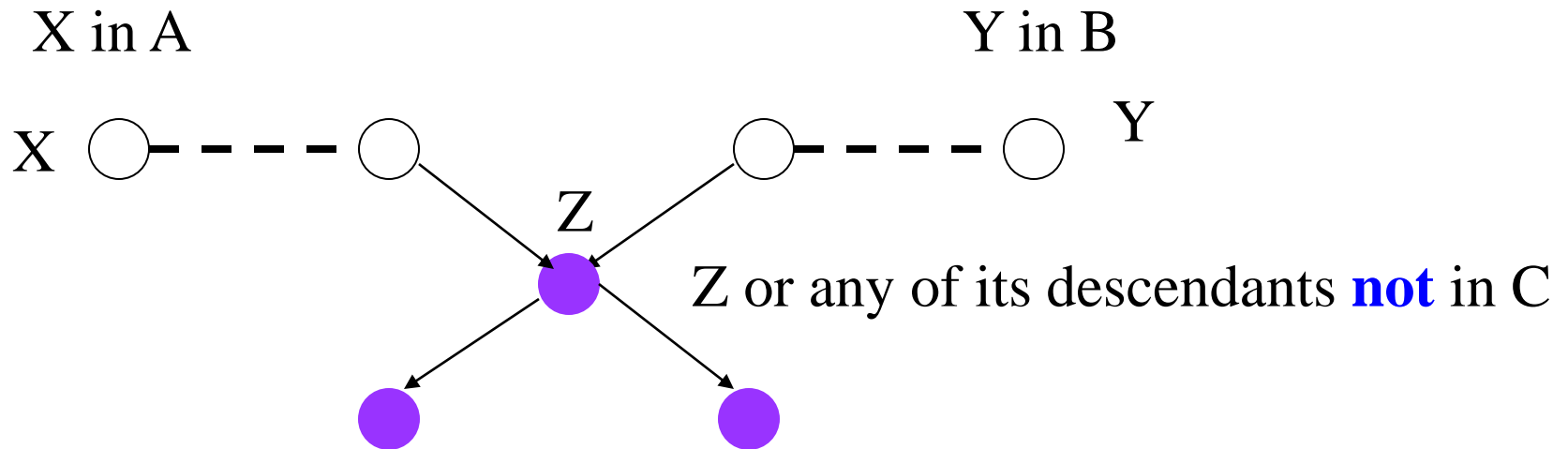
- **2. Path blocking with the wedge substructure**



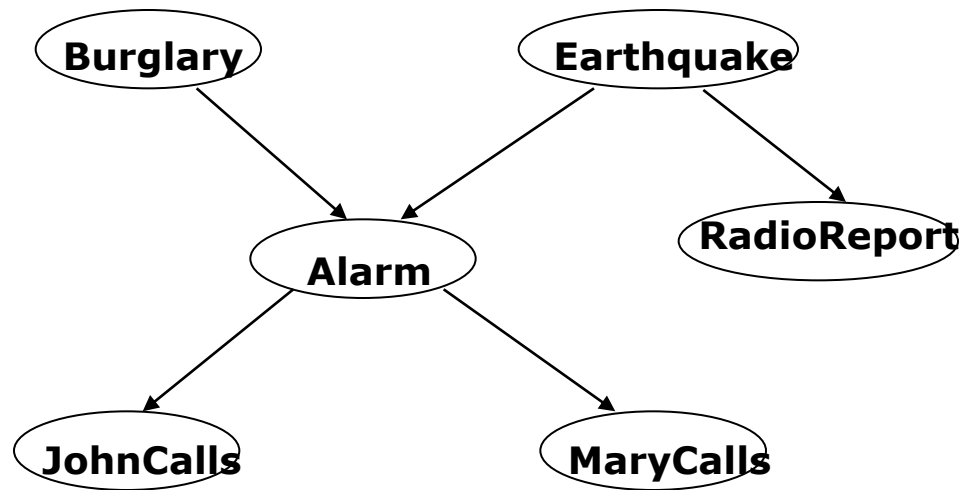
Undirected path blocking

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

- 3. Path blocking with the vee substructure

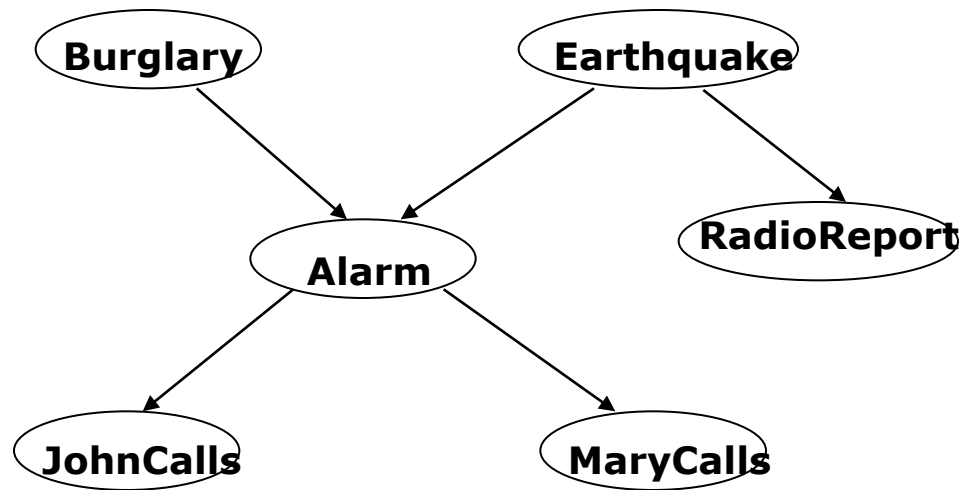


Independences in BBNs



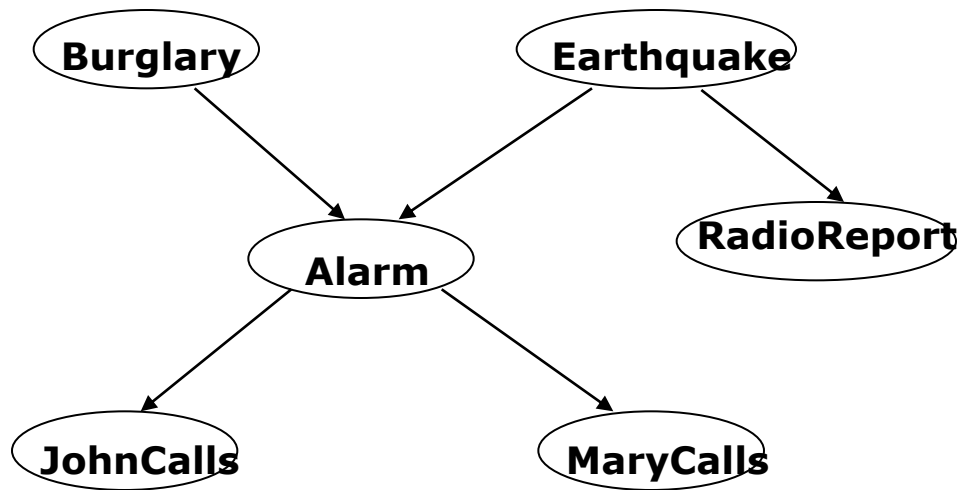
- Earthquake and Burglary are independent given MaryCalls ?

Independences in BBNs



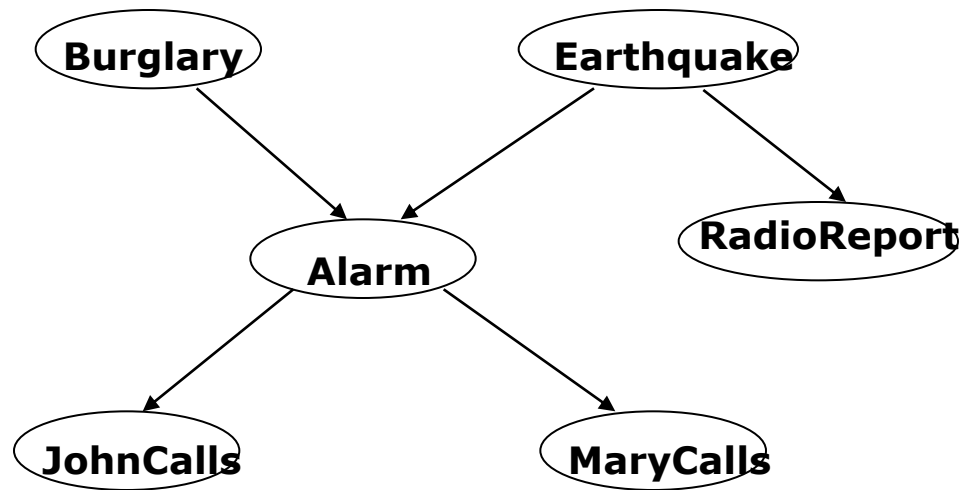
- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **?**

Independences in BBNs



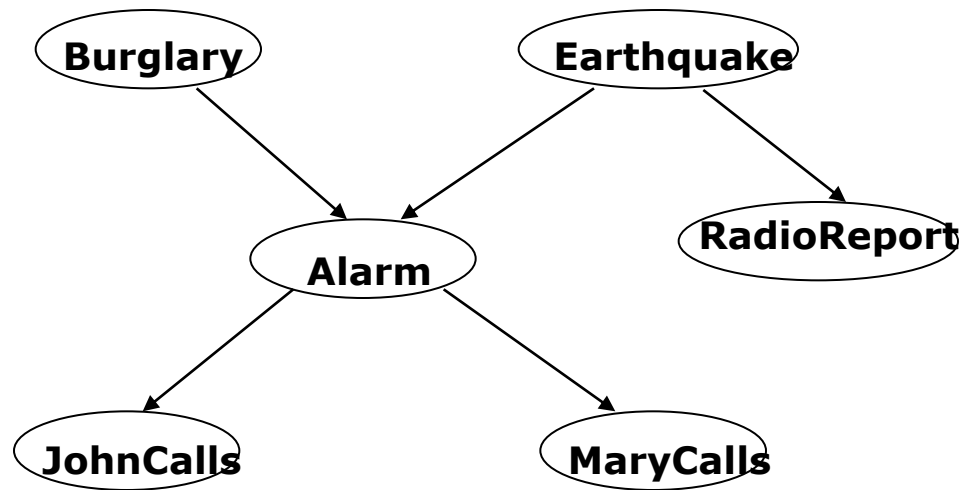
- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **F**
- Burglary and RadioReport are independent given Earthquake **?**

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

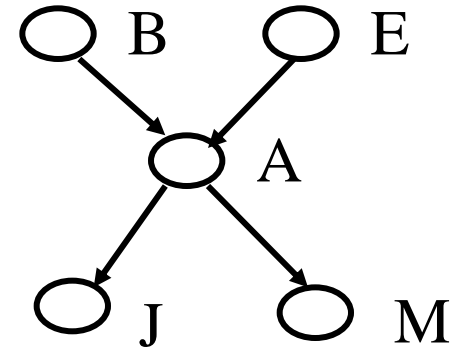


- Earthquake and Burglary are independent given MaryCalls **F**
- Burglary and MaryCalls are independent (not knowing Alarm) **F**
- Burglary and RadioReport are independent given Earthquake **T**
- Burglary and RadioReport are independent given MaryCalls **F**

Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

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



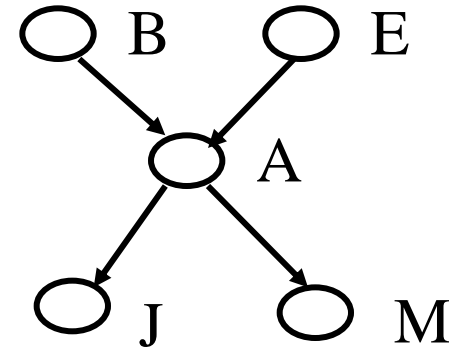
Full joint distribution in BBNs

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:

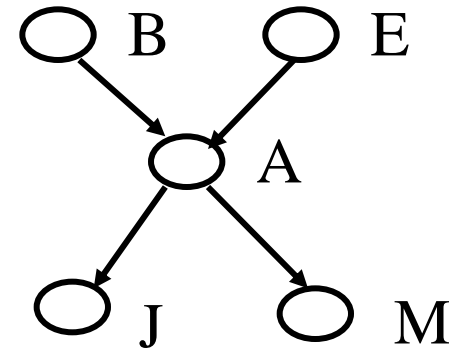
$$P(B = T, E = T, A = T, J = T, M = F) =$$

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

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

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

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



Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

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

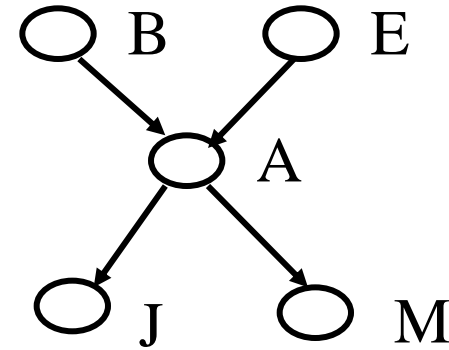
$$= P(J = T \mid B = T, E = T, A = T, M = F)P(B = T, E = T, A = T, M = F)$$

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

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

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

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



Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:

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

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

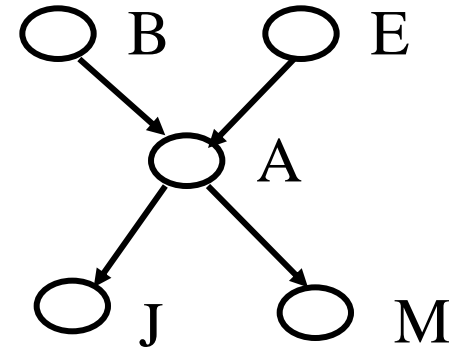
$$= \underline{P(J = T \mid A = T)}P(B = T, E = T, A = T, M = F)$$

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

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

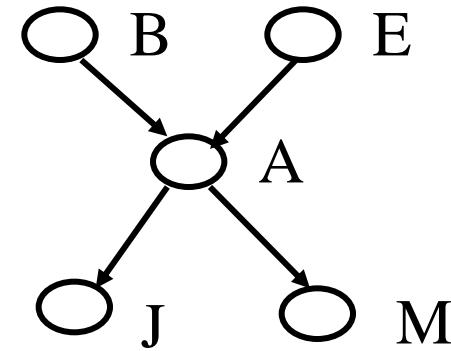
$$\underline{P(A = T \mid B = T, E = T)}P(B = T, E = T)$$

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



Full joint distribution in BBNs

Rewrite the full joint probability using the product rule:



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

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

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

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

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

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

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

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

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?**

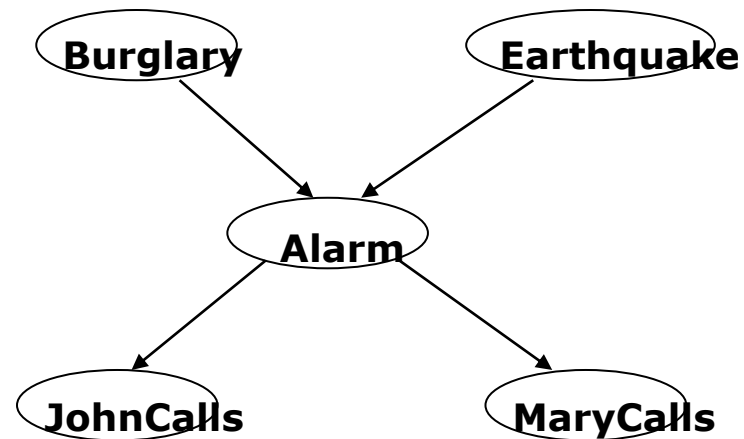
Alarm example: 5 binary (True, False) variables

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:

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

- What did we save?**

Alarm example: 5 binary (True, False) variables

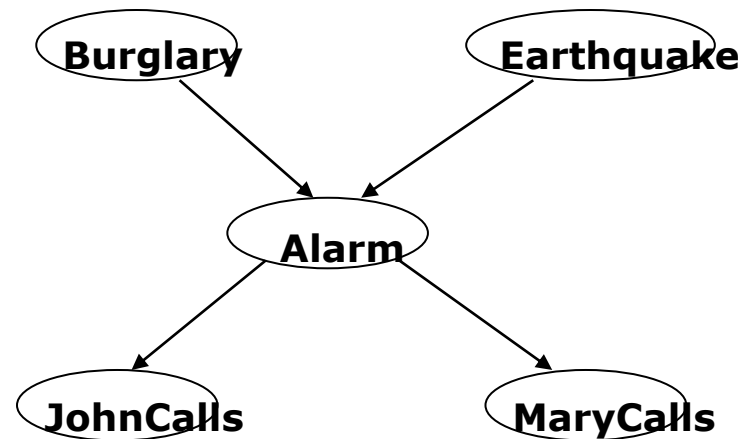
of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

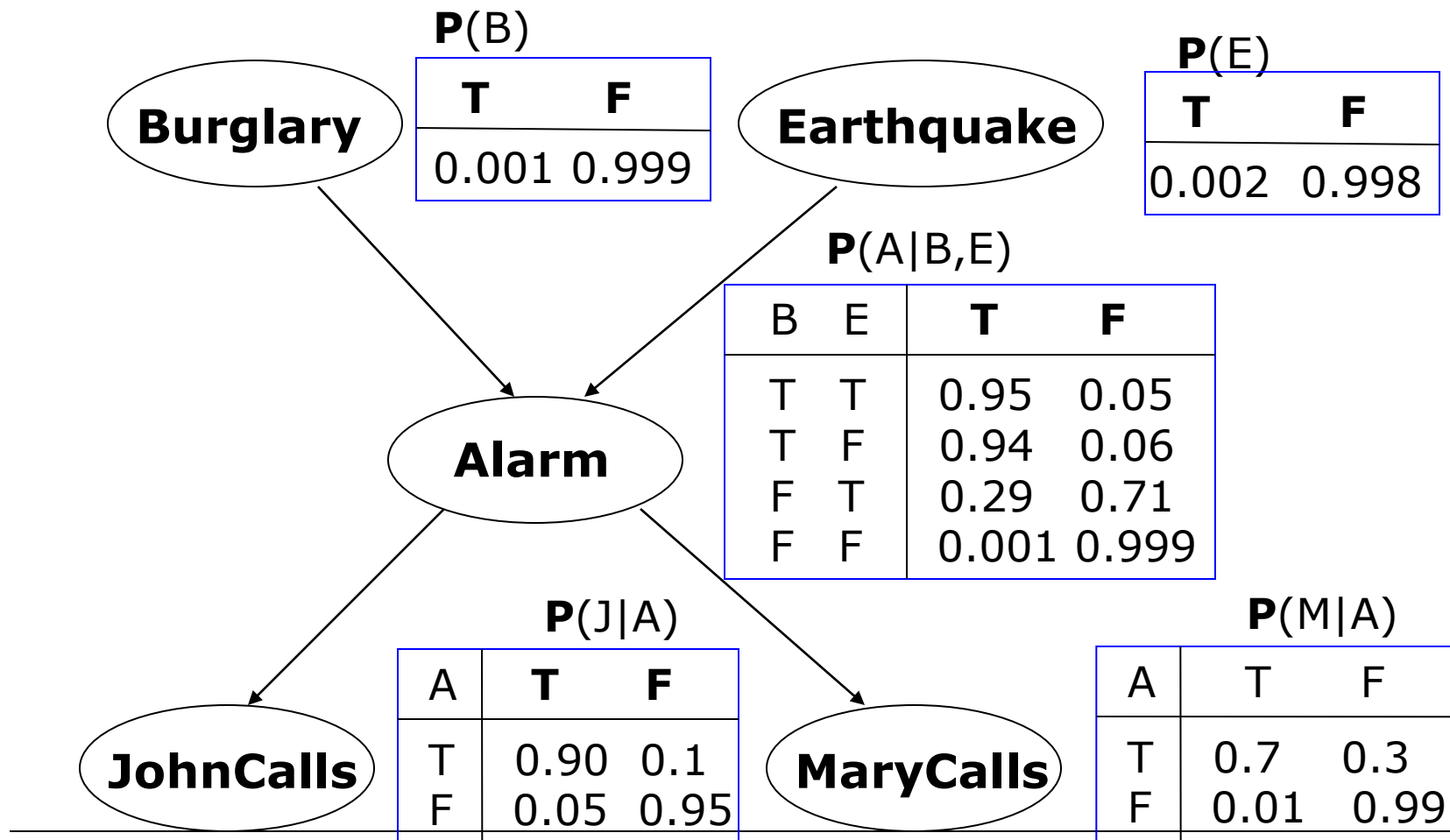
$$2^5 - 1 = 31$$

of parameters of the BBN: ?



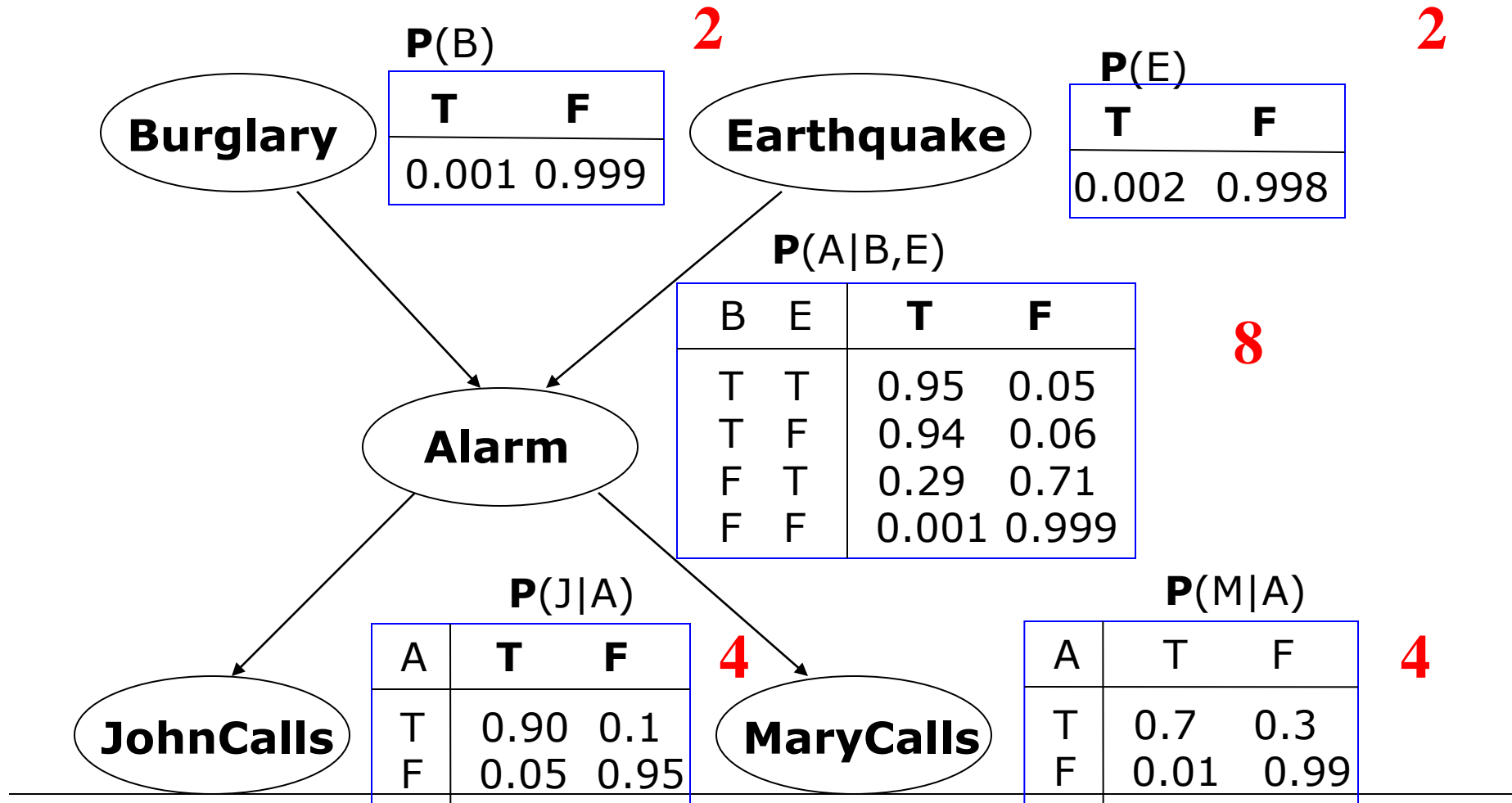
Bayesian belief network.

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



Bayesian belief network.

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



Parameter complexity problem

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

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

- What did we save?**

Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

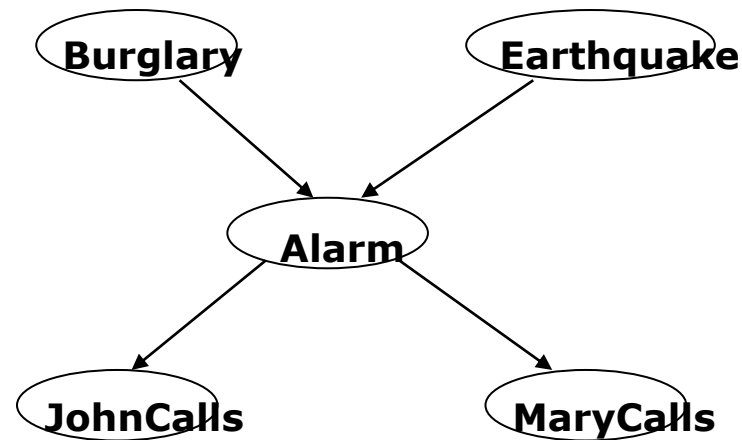
$$2^5 - 1 = 31$$

of parameters of the BBN:

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

One parameter in every conditional is for free:

?



Parameter complexity problem

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

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

- What did we save?**

Alarm example: 5 binary (True, False) variables

of parameters of the full joint:

$$2^5 = 32$$

One parameter is for free:

$$2^5 - 1 = 31$$

of parameters of the BBN:

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

One parameter in every conditional is for free:

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

