

Bayesian Networks Part 1

CS 760@UW-Madison



Goals for the lecture



you should understand the following concepts

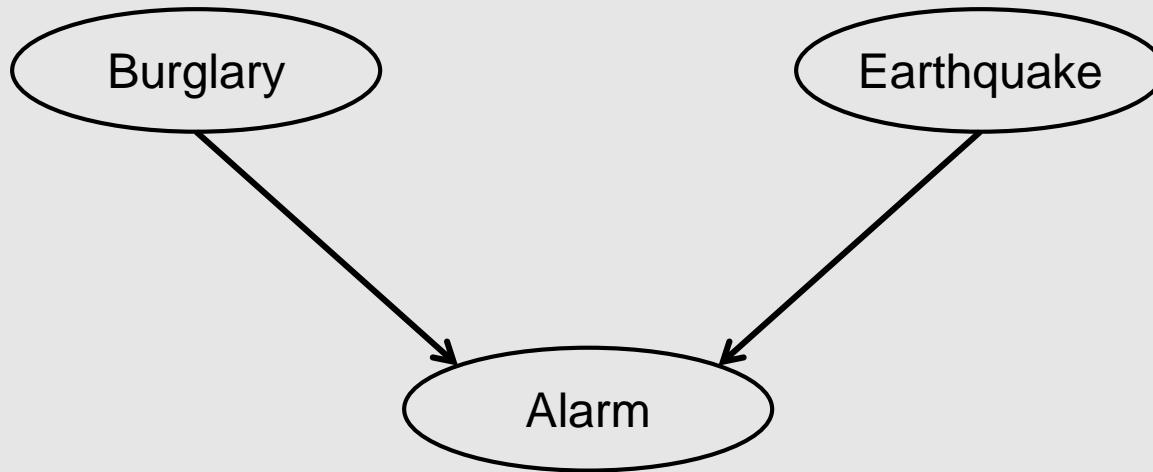
- the Bayesian network representation
- inference by enumeration
- Introduce the learning tasks for Bayes nets

Bayesian network example

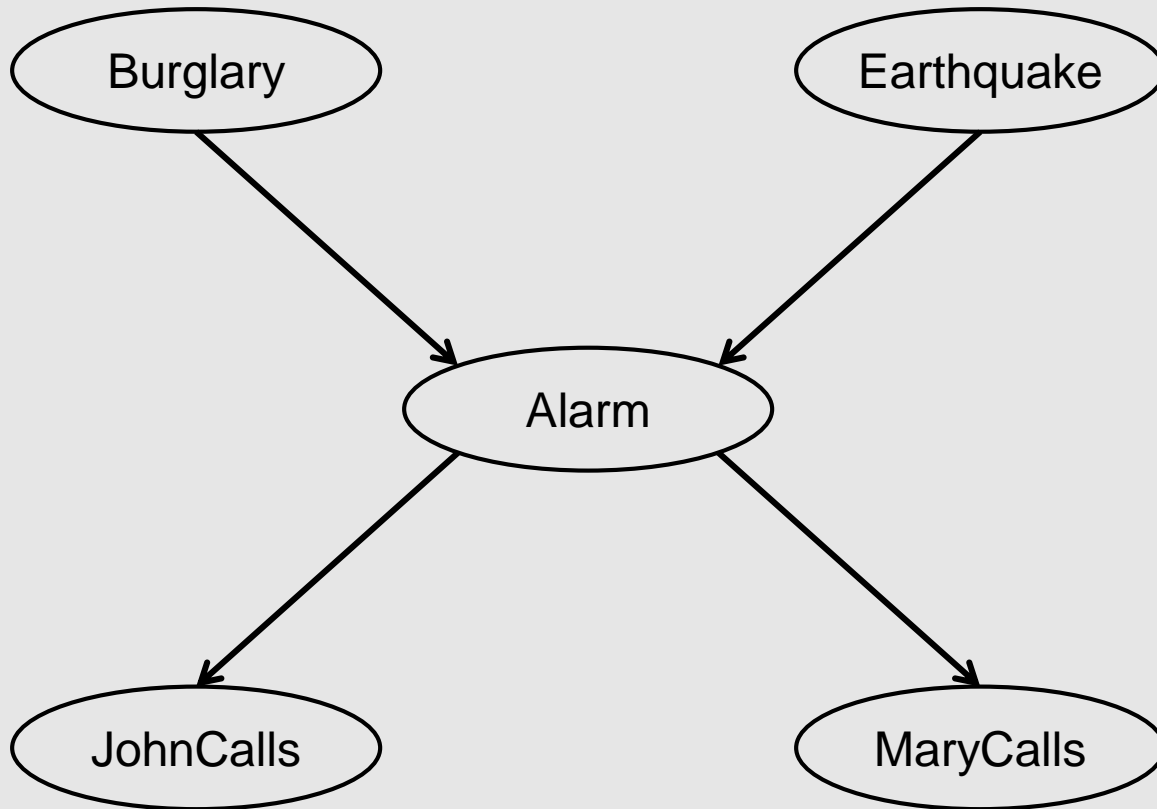


- Consider the following 5 binary random variables:
 - B = a burglary occurs at your house
 - E = an earthquake occurs at your house
 - A = the alarm goes off
 - J = John calls to report the alarm
 - M = Mary calls to report the alarm
- Suppose Burglary or Earthquake can trigger Alarm, and Alarm can trigger John's call or Mary's call
- Now we want to answer queries like **what is $P(B | M, J)$?**

Bayesian network example



Bayesian network example



Bayesian network example

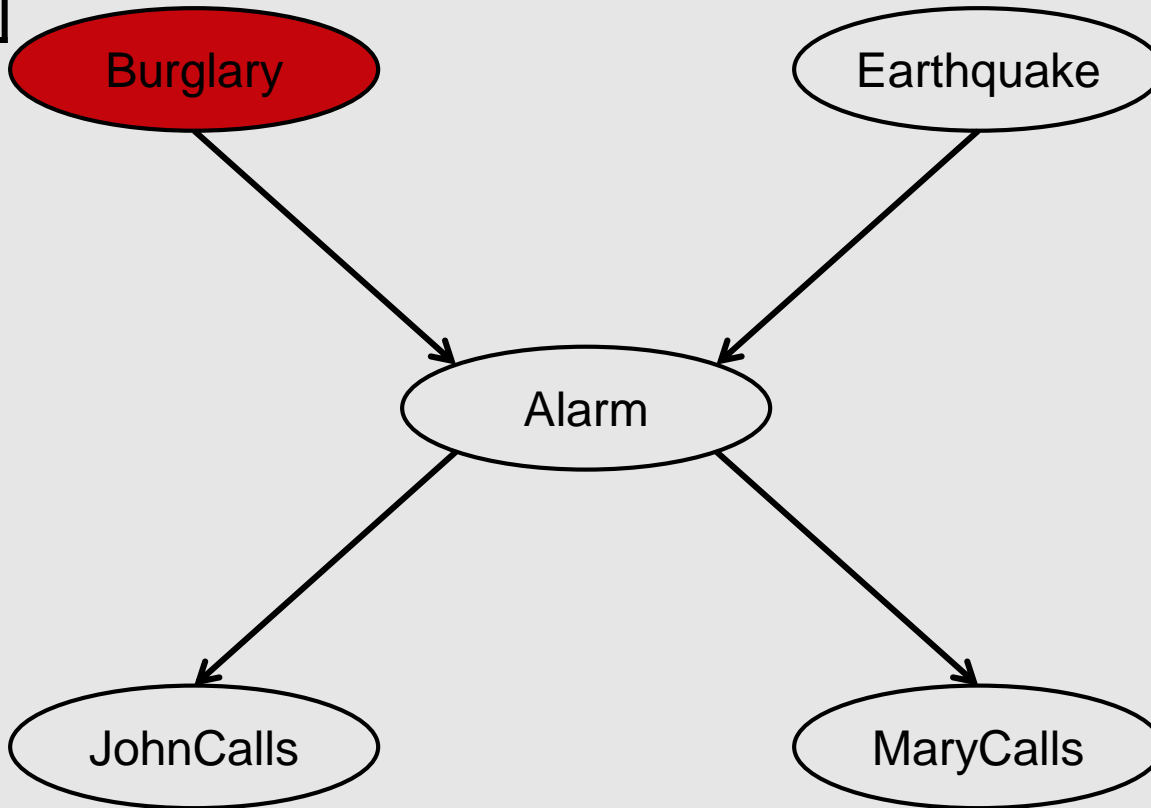


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



Bayesian network example

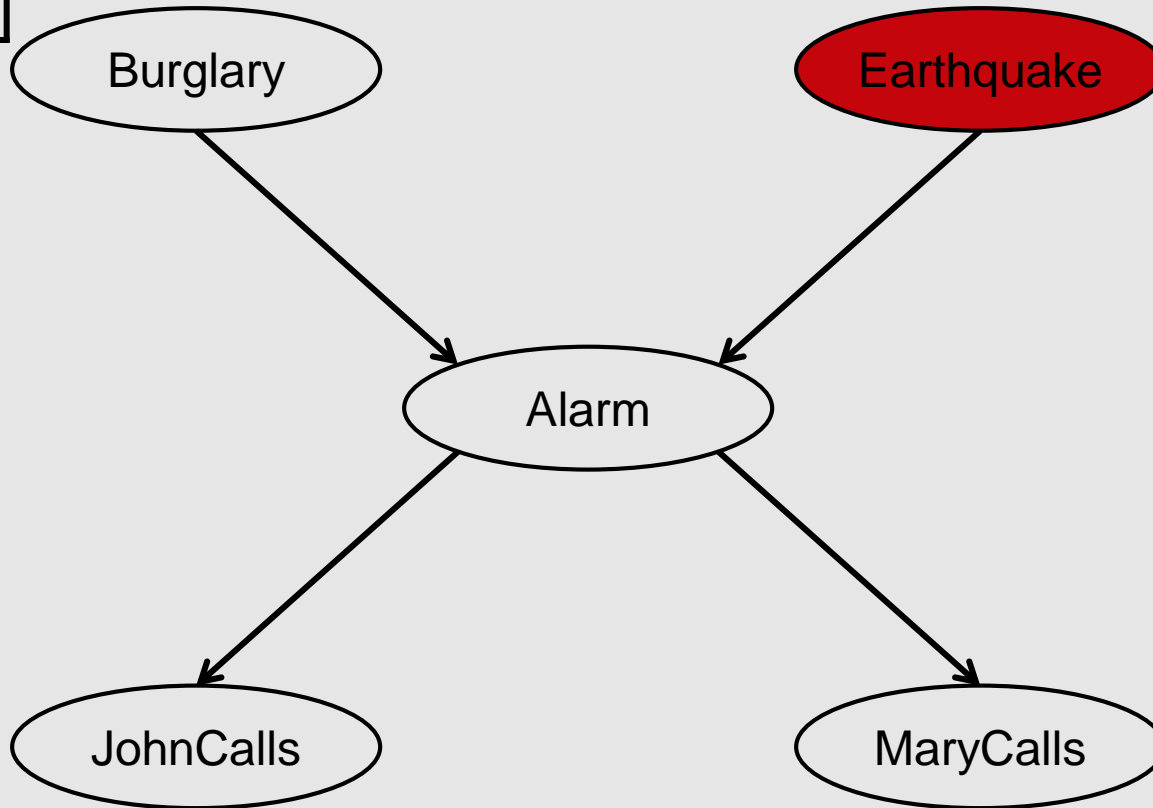


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



Bayesian network example

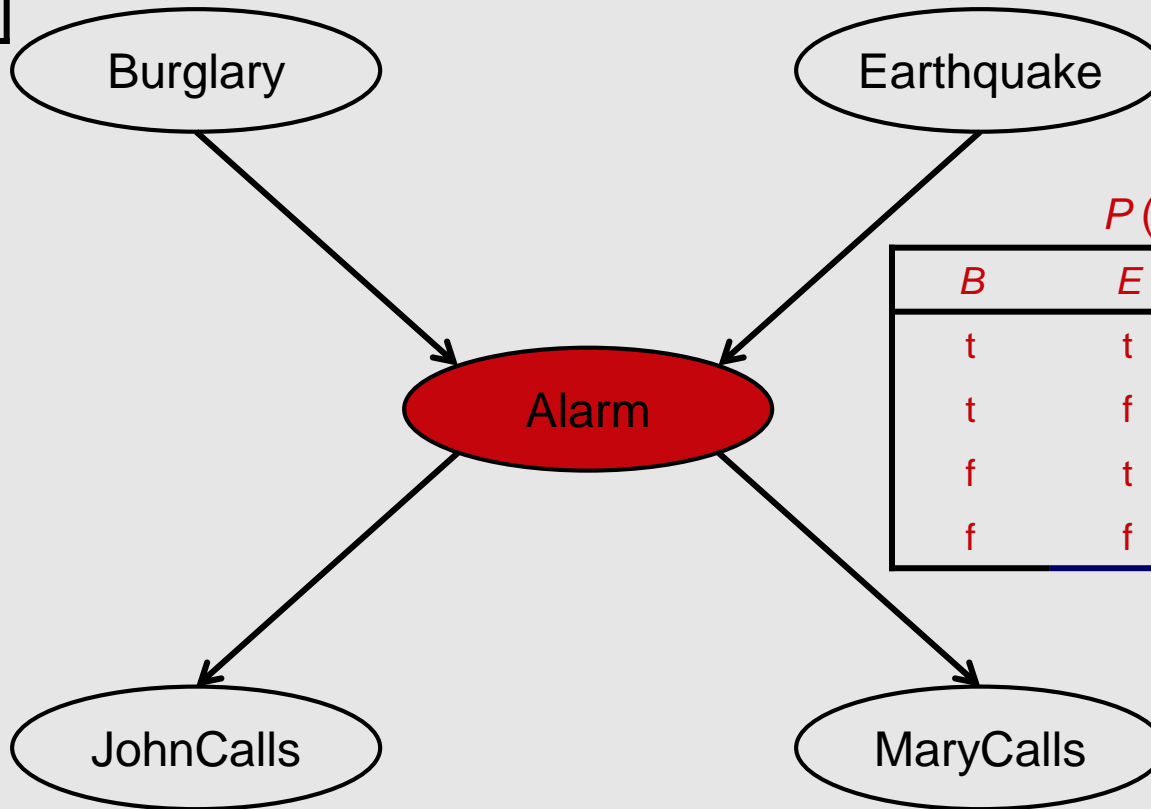


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



$P(A|B, E)$

<i>B</i>	<i>E</i>	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 network example

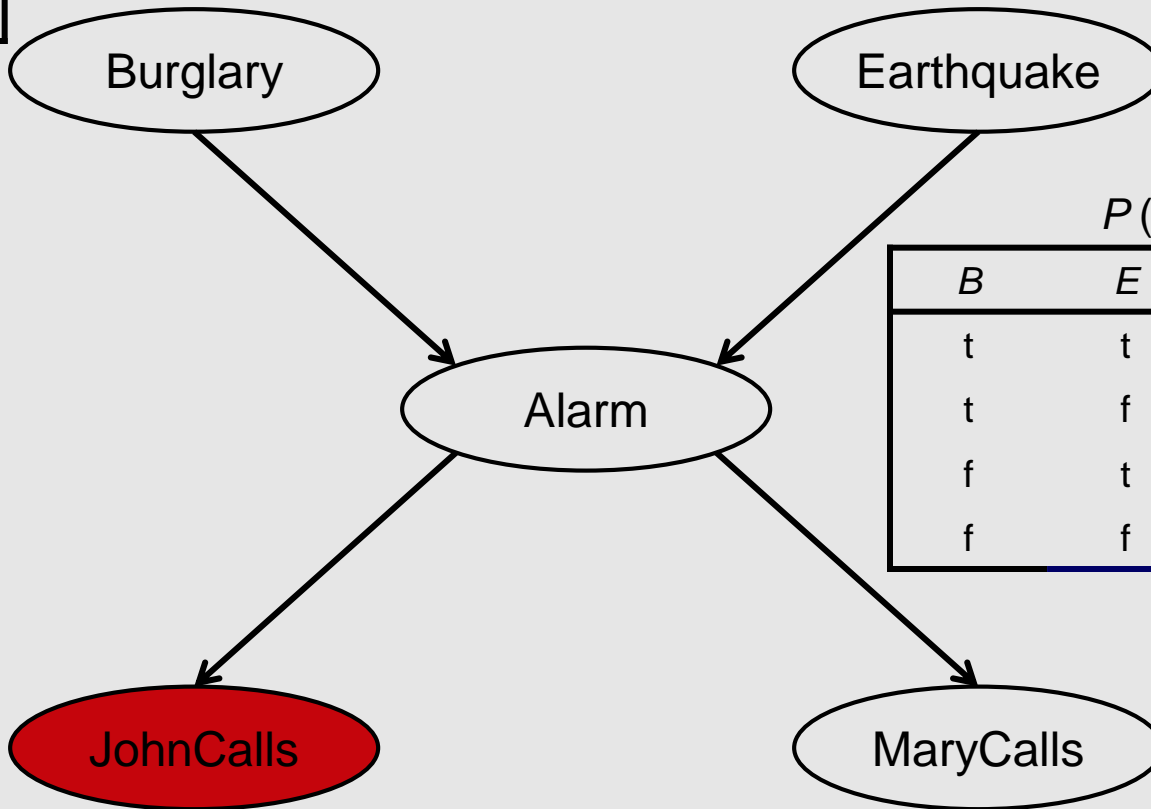


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



$P(A | B, E)$

<i>B</i>	<i>E</i>	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

$P(J | A)$

<i>A</i>	t	f
t	0.9	0.1
f	0.05	0.95

Bayesian network example

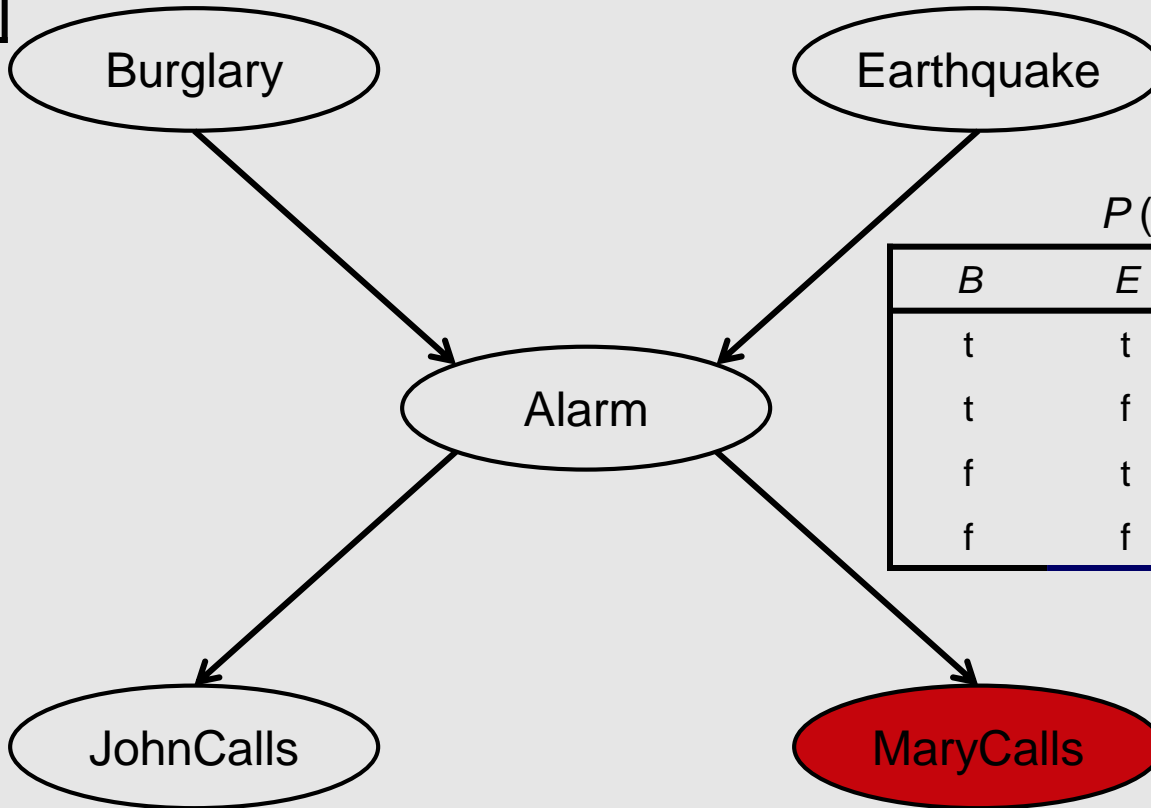


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



$P(A | B, E)$

<i>B</i>	<i>E</i>	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

$P(J | A)$

<i>A</i>	t	f
t	0.9	0.1
f	0.05	0.95

$P(M | A)$

<i>A</i>	t	f
t	0.7	0.3
f	0.01	0.99

Bayesian network example

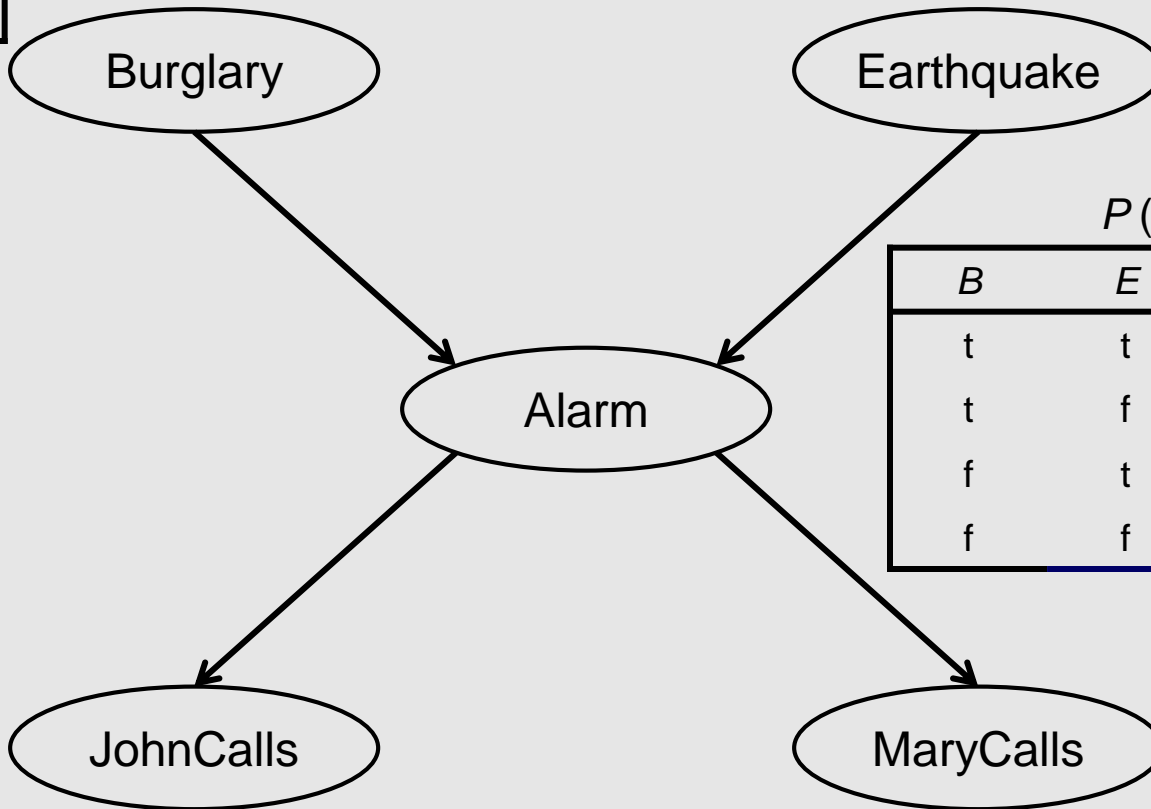


$P(B)$

t	f
0.001	0.999

$P(E)$

t	f
0.001	0.999



$P(A | B, E)$

<i>B</i>	<i>E</i>	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

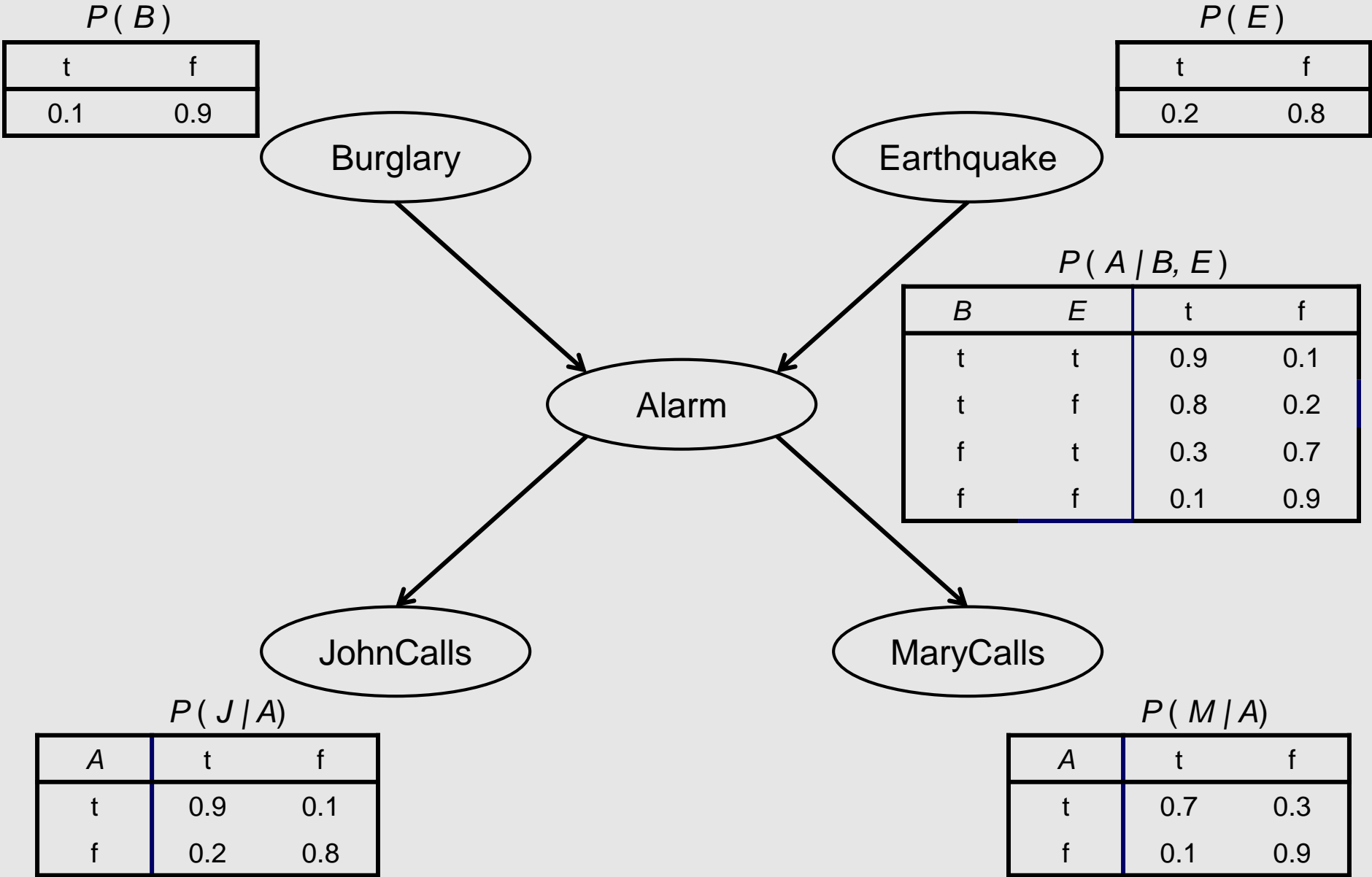
$P(J | A)$

<i>A</i>	t	f
t	0.9	0.1
f	0.05	0.95

$P(M | A)$

<i>A</i>	t	f
t	0.7	0.3
f	0.01	0.99

Bayesian network example (different parameters)





Bayesian networks

- a BN consists of a **Directed Acyclic Graph (DAG)** and a set of **conditional probability distributions**
- in the DAG
 - each node denotes random a variable
 - each edge from X to Y represents that X *directly influences* Y
 - (formally: each variable X is independent of its non-descendants given its parents)
- each node X has a *conditional probability distribution* (CPD) representing $P(X \mid \text{Parents}(X))$

Bayesian networks



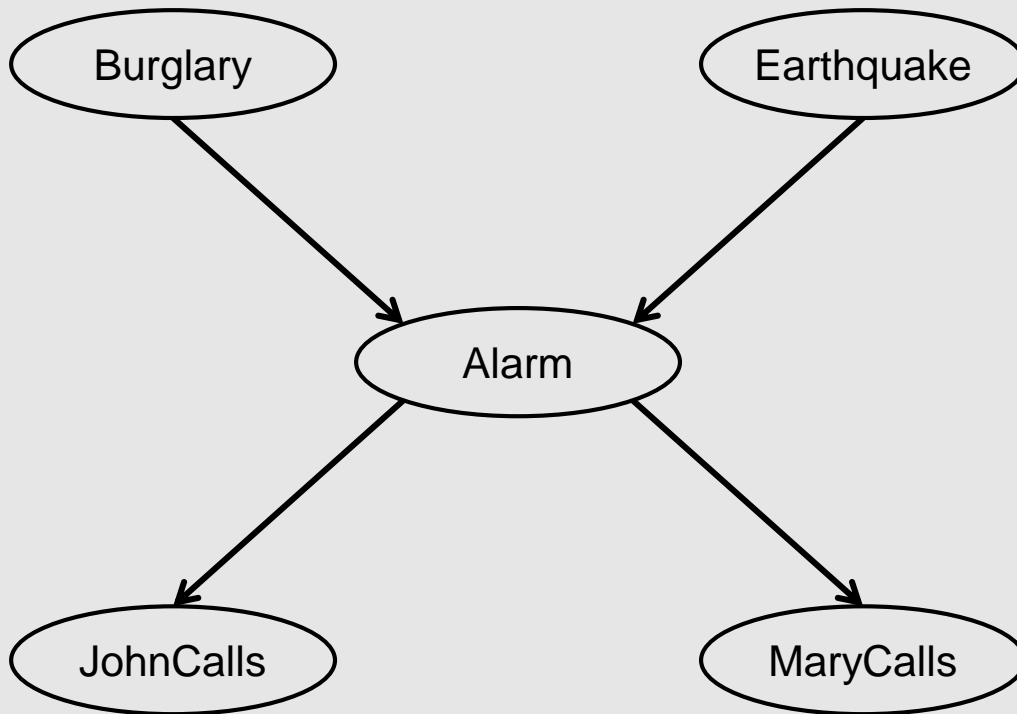
- using the chain rule, a joint probability distribution can always be expressed as

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

- a BN provides a compact representation of a joint probability distribution. It corresponds to the assumption:

$$P(X_1, \dots, X_n) = P(X_1) \prod_{i=2}^n P(X_i | \text{Parents}(X_i))$$

Bayesian networks



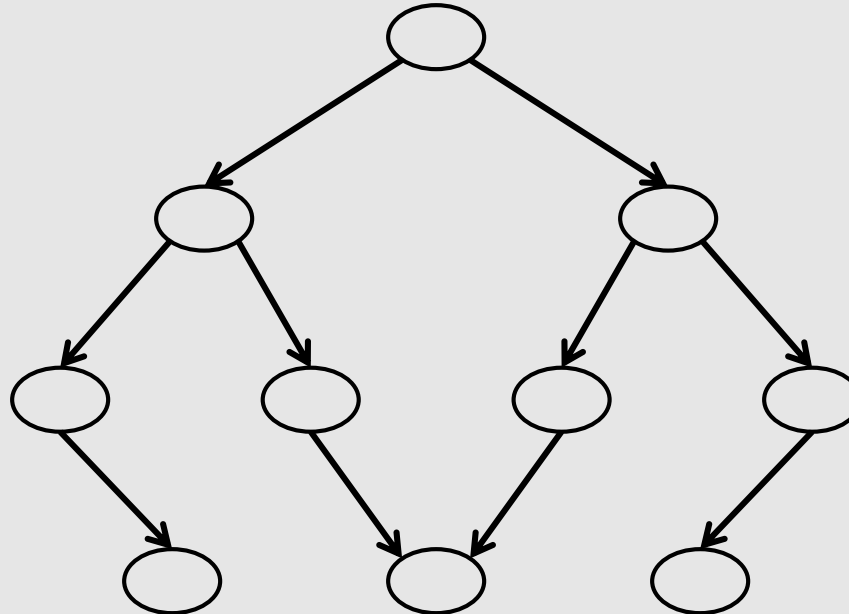
$$\begin{aligned} P(B, E, A, J, M) \\ &= P(B) \\ &\times P(E) \\ &\times P(A | B, E) \\ &\times P(J | A) \\ &\times P(M | A) \end{aligned}$$

- a standard representation of the joint distribution for the Alarm example has $2^5 = 32$ parameters
- the BN representation of this distribution has 20 parameters



Bayesian networks

- consider a case with 10 binary random variables
- How many parameters does a BN with the following graph structure have?



- How many parameters does the standard table representation of the joint distribution have?

Advantages of Bayesian network representation

- Captures independence and conditional independence where they exist
- Encodes the relevant portion of the full joint among variables where dependencies exist
- Uses a graphical representation which lends insight into the complexity of inference

Inference



The inference task in Bayesian networks



Given: values for some variables in the network (*evidence*), and a set of *query* variables

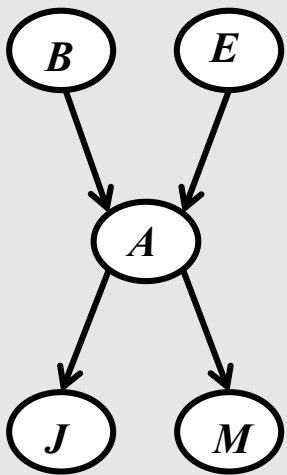
Do: compute the posterior distribution over the query variables

- variables that are neither evidence variables nor query variables are *hidden* variables
- the BN representation is flexible enough that any set can be the evidence variables and any set can be the query variables



Inference by enumeration

- let a denote $A=\text{true}$, and $\neg a$ denote $A=\text{false}$
- suppose we're given the query: $P(b | j, m)$
“probability the house is being burglarized given that John and Mary both called”
- from the graph structure we can first compute:



$$P(b, j, m) = \sum_{e, \neg e} \sum_{a, \neg a} P(b)P(E)P(A | b, E)P(j | A)P(m | A)$$

sum over possible values for E and A variables ($e, \neg e, a, \neg a$)



Inference by enumeration

$$P(b, j, m) = \sum_{e, \neg e} \sum_{a, \neg a} P(b)P(E)P(A | b, E)P(j | A)P(m | A)$$

$$= P(b) \sum_{e, \neg e} \sum_{a, \neg a} P(E)P(A | b, E)P(j | A)P(m | A)$$

$P(B)$
0.001

$P(E)$
0.001

B *E* *A* *J* *M*

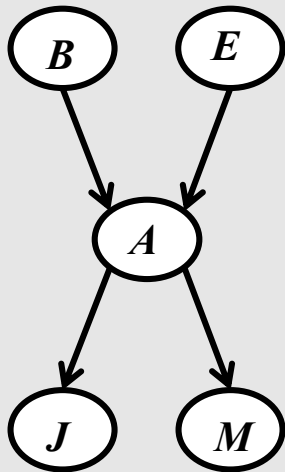
$$= 0.001 \times (0.001 \times 0.95 \times 0.9 \times 0.7 +$$

$$0.001 \times 0.05 \times 0.05 \times 0.01 +$$

$$0.999 \times 0.94 \times 0.9 \times 0.7 +$$

$$0.999 \times 0.06 \times 0.05 \times 0.01)$$

e, a
e, ¬a
¬e, a
¬e, ¬a



<i>B</i>	<i>E</i>	$P(A)$
t	t	0.95
t	f	0.94
f	t	0.29
f	f	0.001

<i>A</i>	$P(J)$
t	0.9
f	0.05

<i>A</i>	$P(M)$
t	0.7
f	0.01



Inference by enumeration

- now do equivalent calculation for $P(\neg b, j, m)$
- and determine $P(b | j, m)$

$$P(b | j, m) = \frac{P(b, j, m)}{P(j, m)} = \frac{P(b, j, m)}{P(b, j, m) + P(\neg b, j, m)}$$

Comments on BN inference



- *inference by enumeration* is an *exact* method (i.e. it computes the exact answer to a given query)
- it requires summing over a joint distribution whose size is exponential in the number of variables
- in many cases we can do exact inference efficiently in large networks
 - key insight: save computation by pushing sums inward
- in general, the Bayes net inference problem is NP-hard
- there are also methods for approximate inference – these get an answer which is “close”
- in general, the approximate inference problem is NP-hard also, but approximate methods work well for many real-world problems



Learning

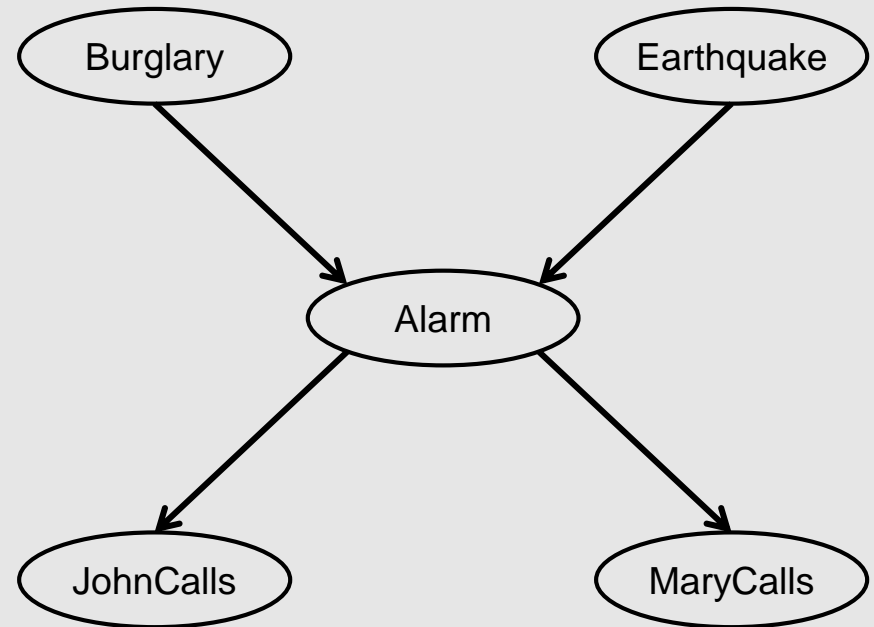


The parameter learning task



- Given: a set of training instances, the graph structure of a BN

B	E	A	J	M
f	f	f	t	f
f	t	f	f	f
f	f	t	f	t
		...		



- Do: infer the parameters of the CPDs

The structure learning task



- Given: a set of training instances

B	E	A	J	M
f	f	f	t	f
f	t	f	f	f
f	f	t	f	t
		...		

- Do: infer the graph structure (and perhaps the parameters of the CPDs too)

Parameter learning and MLE



- *maximum likelihood estimation* (MLE)
 - given a model structure (e.g. a Bayes net graph) G and a set of data D
 - set the model parameters θ to maximize $P(D | G, \theta)$
- i.e. make the data D look as likely as possible under the model $P(D | G, \theta)$

Maximum likelihood estimation review



consider trying to estimate the parameter θ (probability of heads) of a biased coin from a sequence of flips (1 stands for head)

$$\mathbf{x} = \{1, 1, 1, 0, 1, 0, 0, 1, 0, 1\}$$

the likelihood function for θ is given by:

$$\begin{aligned} L(\theta : x_1, \dots, x_n) &= \theta^{x_1} (1 - \theta)^{1-x_1} \dots \theta^{x_n} (1 - \theta)^{1-x_n} \\ &= \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i} \end{aligned}$$

What's MLE of the parameter?



THANK YOU

Some of the slides in these lectures have been adapted/borrowed from materials developed by Mark Craven, David Page, Jude Shavlik, Tom Mitchell, Nina Balcan, Elad Hazan, Tom Dietterich, and Pedro Domingos.

