

# Artificial Intelligence



Lesson 11  
(From Russell & Norvig)

168

Ram Meshulam 2004

# Conditional probability

- Conditional or posterior probabilities  
e.g.,  $P(\text{cavity} \mid \text{toothache}) = 0.8$   
i.e., given that *toothache* is all I know
- Notation for conditional distributions:  
 $\mathbf{P}(\text{Cavity} \mid \text{Toothache}) = 2\text{-element vector of } 2\text{-element vectors}$
- If we know more, e.g., *cavity* is also given, then we have  
 $P(\text{cavity} \mid \text{toothache}, \text{cavity}) = 1$
- New evidence may be irrelevant, allowing simplification, e.g.,  
 $P(\text{cavity} \mid \text{toothache}, \text{sunny}) = P(\text{cavity} \mid \text{toothache}) = 0.8$
- This kind of inference, sanctioned by domain knowledge, is crucial

169

Ram Meshulam 2004

# Inference by enumeration

- Start with the joint probability distribution:

	<i>toothache</i>		$\neg$ <i>toothache</i>	
	<i>catch</i>	$\neg$ <i>catch</i>	<i>catch</i>	$\neg$ <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
$\neg$ <i>cavity</i>	.016	.064	.144	.576

- Can also compute conditional probabilities:

$$\begin{aligned}
 P(\neg \text{cavity} \mid \text{toothache}) &= \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})} \\
 &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} = 0.4
 \end{aligned}$$

170

Ram Meshulam 2004

# Independence

- $A$  and  $B$  are independent iff  
 $\mathbf{P}(A/B) = \mathbf{P}(A)$  or  $\mathbf{P}(B/A) = \mathbf{P}(B)$  or  $\mathbf{P}(A, B) = \mathbf{P}(A) \mathbf{P}(B)$



$$\begin{aligned}
 \mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) \\
 = \mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}) \mathbf{P}(\text{Weather})
 \end{aligned}$$

- 32 entries reduced to 12; for  $n$  independent biased coins,  $O(2^n) \rightarrow O(n)$
- Absolute independence powerful but rare
- Dentistry is a large field with hundreds of variables, none of which are independent. What to do?

171

Ram Meshulam 2004

## Conditional independence

- $\mathbf{P}(\text{Toothache}, \text{Cavity}, \text{Catch})$  has  $2^3 - 1 = 7$  independent entries
- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
  - (1)  $\mathbf{P}(\text{catch} \mid \text{toothache}, \text{cavity}) = \mathbf{P}(\text{catch} \mid \text{cavity})$
- The same independence holds if I haven't got a cavity:
  - (2)  $\mathbf{P}(\text{catch} \mid \text{toothache}, \neg \text{cavity}) = \mathbf{P}(\text{catch} \mid \neg \text{cavity})$
- *Catch* is **conditionally independent** of *Toothache* given *Cavity*:  
 $\mathbf{P}(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = \mathbf{P}(\text{Catch} \mid \text{Cavity})$
- Equivalent statements:  
 $\mathbf{P}(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = \mathbf{P}(\text{Toothache} \mid \text{Cavity})$   
 $\mathbf{P}(\text{Toothache}, \text{Catch} \mid \text{Cavity}) = \mathbf{P}(\text{Toothache} \mid \text{Cavity}) \mathbf{P}(\text{Catch} \mid \text{Cavity})$

172

Ram Meshulam 2004

## Bayesian networks

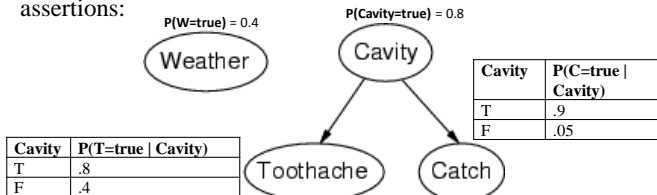
- A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions
- It describes how variables interact locally
- Local interactions chain together to give global, indirect interactions
- **Syntax:**
  - a set of nodes, one per variable
  - a directed, acyclic graph (link  $\approx$  "directly influences")
  - a conditional distribution for each node given its parents:  
 $\mathbf{P}(X_i \mid \text{Parents}(X_i))$ - conditional probability table (CPT)

173

Ram Meshulam 2004

## Example 1

- Topology of network encodes conditional independence assertions:



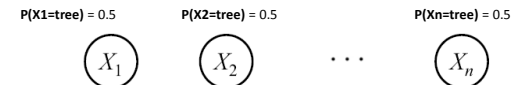
- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity*
- It is usually easy for a domain expert to decide what direct influences exist

174

Ram Meshulam 2004

## Example 2

- $N$  independent coin flips :



- No interactions between variables: absolute independence
- Does every Bayes Net can represent every full joint?
- No. For example, Only distributions whose variables are absolutely independent can be represented by a Bayes' net with no arcs.

175

Ram Meshulam 2004

## Calculation of Joint Probability

- Given its parents, each node is conditionally independent of everything except its descendants

- Thus,

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_n) = \prod_{i=1, \dots, n} P(x_i | \text{parents}(X_i))$$

→ full joint distribution table

- Every BN over a domain implicitly represents some joint distribution over that domain

176

Ram Meshulam 2004

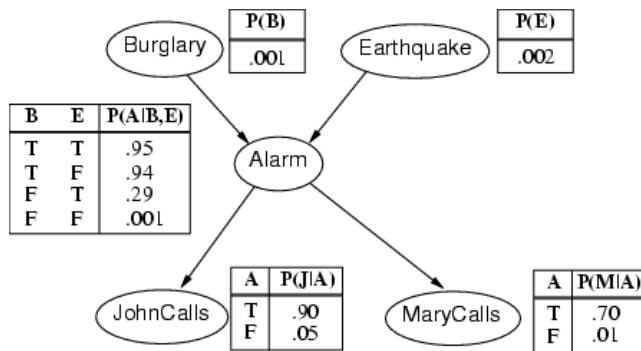
## Example 3

- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?
- Variables: *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*
- Network topology reflects "causal" knowledge:
  - A burglar can set the alarm off
  - An earthquake can set the alarm off
  - The alarm can cause Mary to call
  - The alarm can cause John to call

177

Ram Meshulam 2004

## Example contd.



178

Ram Meshulam 2004

## Answering queries

- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?
  - $P(b|j, \neg m) = P(b, j, \neg m) / P(j, \neg m)$
  - $P(b, j, \neg m) = P(b, e, a, j, \neg m) + P(b, \neg e, a, j, \neg m) + P(b, e, \neg a, j, \neg m) + P(b, \neg e, \neg a, j, \neg m) =$   
 $P(b)P(e)P(a|b, e)P(j|a)P(\neg m|\neg a) +$   
 $P(b)P(e)P(\neg a|b, e)P(j|\neg a)P(\neg m|\neg a) +$   
 $P(b)P(\neg e)P(a|b, \neg e)P(j|a)P(\neg m|a) +$   
 $P(b)P(\neg e)P(\neg a|b, \neg e)P(j|\neg a)P(\neg m|\neg a)$
  - Do the same to calculate  $P(\neg b, j, \neg m)$  and normalize
  - Worst case, for a network with  $n$  Boolean variables,  $O(n^2^n)$ .

179

Ram Meshulam 2004

## Laziness and Ignorance

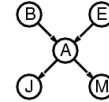
- The probabilities actually summarize a potentially infinite set of circumstances in which the alarm might fail to go off
  - high humidity
  - power failure
  - dead battery
  - cut wires
  - a dead mouse stuck inside the bell
- John or Mary might fail to call and report it
  - out to lunch
  - on vacation
  - temporarily deaf
  - passing helicopter

180

Ram Meshulam 2004

## Compactness

- A CPT for Boolean  $X_i$  with  $k$  Boolean parents has  $2^k$  rows for the combinations of parent values
- Each row requires one number  $p$  for  $X_i = \text{true}$  (the number for  $X_i = \text{false}$  is just  $1-p$ )
- If each variable has no more than  $k$  parents, the complete network requires  $O(n \cdot 2^k)$  numbers
- I.e., grows linearly with  $n$ , vs.  $O(2^n)$  for the full joint distribution
- For burglary net,  $1 + 1 + 4 + 2 + 2 = 10$  numbers (vs.  $2^5 - 1 = 31$ )
- We utilize the property of **locally structured system**

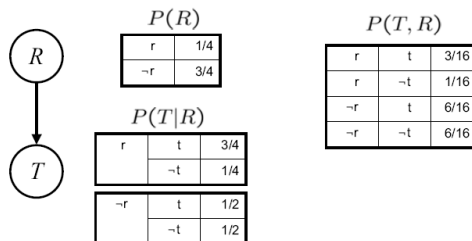


181

Ram Meshulam 2004

## Casualty?

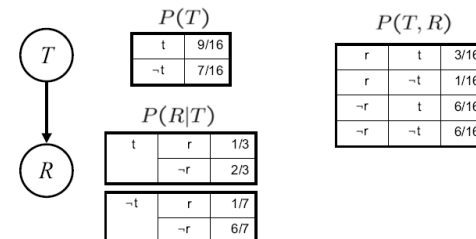
- Rain causes Traffic
- Let's build the joint:



182

Ram Meshulam 2004

## Reverse Casualty?



183

Ram Meshulam 2004

## Casualty?

- What do the arrows really mean?
- Topology may happen to encode causal structure
- Topology really encodes conditional independencies
- When Bayes' nets reflect the true causal patterns:
  - Often simpler (nodes have fewer parents)
  - Often easier to think about
  - Often easier to elicit from experts
- BNs need not actually be causal
  - Sometimes no causal net exists over the domain
  - E.g. consider the variables Traffic and RoofDrips
  - End up with arrows that reflect correlation, not causation

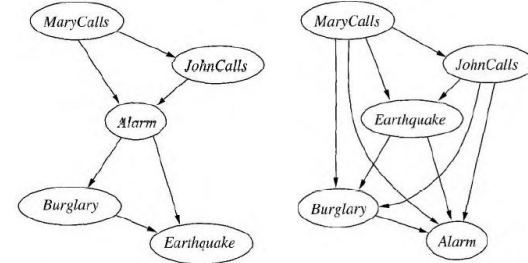
184

Ram Meshulam 2004

## Example 2, Again

Consider the following 2 orders for insertion:

- (a) MaryCalls, JohnCalls, Alarm, Burglary, Earthquake
  - Since,  $P(\text{Burglary}|\text{Alarm, JohnCalls, MaryCalls}) = P(\text{Burglary}|\text{Alarm})$
- (b) Mary Calls, JohnCalls, Earthquake, Burglary, Alarm.



185

Ram Meshulam 2004

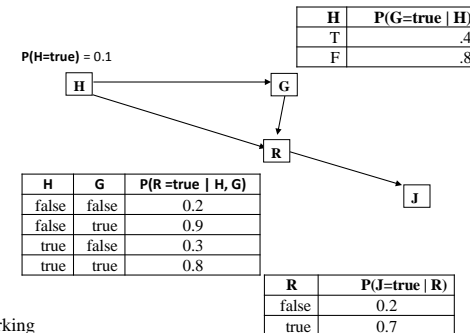
## Connection Types

Name	Diagram	X ind. Z?	X ind. Z, given Y?
Casual chain		Not necessarily	Yes
Common Cause		No	Yes
Common Effect		Yes	No

186

Ram Meshulam 2004

## Test Question



- H - Hardworking
- G - Good Grader
- R - Excellent Recommendation
- J - Landed a good Job

187

Ram Meshulam 2004

## What can be inferred?

- i:  $P(H,G) = P(H) \cdot P(G)$  ✗
- ii  $P(J|R,H) = P(J|R)$  ✓
- iii  $P(J) \neq P(J|H)$  ✗

Q: What is the value of  $P(H,G, \neg R, \neg J)$ ?

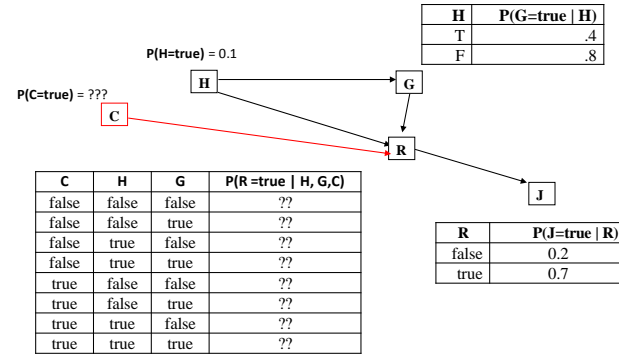
A:  $P(H,G, \neg R, \neg J) = P(H) \cdot P(G|H) \cdot P(\neg R|H,G) \cdot P(\neg J|H,G, \neg R) = P(H) \cdot P(G|H) \cdot P(\neg R|H,G) \cdot P(\neg J | \neg R) = 0.1 * 0.4 * 0.2 * 0.8 = 0.0064$

Q: What if we want to add another parameter, C= Has The Right Connections?

188

Ram Meshulam 2004

## Answer

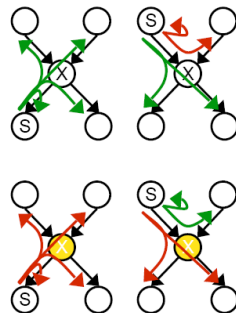


189

Ram Meshulam 2004

## Reachability (the Bayes Ball)

- Shade evidence nodes
- Start at source node
- Try to reach target by search
- States: node, along with previous arc
- Successor function:
  - Unobserved nodes:
    - To any child
    - To any parent if coming from a child
  - Observed nodes:
    - From parent to parent
- If you can't reach a node, it's conditionally independent of the start node

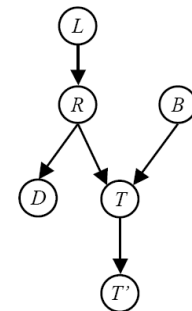


190

Ram Meshulam 2004

## Example

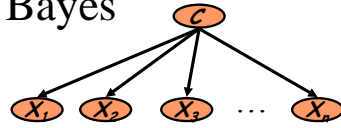
- L ind. T', given T? **Yes**
- L ind. B? **Yes**
- L ind. B, given T? **No**
- L ind. B, given T'?' **No**
- L ind. B, given T and R? **Yes**



191

Ram Meshulam 2004

## Naïve Bayes



- **Conditional Independence Assumption:** features are independent of each other given the class:

$$P(X_1, \dots, X_n | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_n | C)$$

- What can we model with naïve bayes?
- Any process where,
  - Each cause has lots of “independent” effects
  - Easy to estimate the CPT fro each effect
  - We want to reason about the probability of different causes given observed effects

192

Ram Meshulam 2004

## Naive Bayes Classifiers

Task: Classify a new instance  $D$  based on a tuple of attribute values into one of the classes  $c_j \in C$

$$D = \langle x_1, x_2, \dots, x_n \rangle$$

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | x_1, x_2, \dots, x_n)$$

$$= \operatorname{argmax}_{c \in C} \frac{P(x_1, x_2, \dots, x_n | c)P(c)}{P(x_1, x_2, \dots, x_n)}$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

193

CIS 391 - Intro to AI

## Summary

- Bayesian networks provide a natural representation for (causally induced) conditional independence
- Topology + CPTs = compact representation of joint distribution
- Generally easy for domain experts to construct

194

Ram Meshulam 2004