

1

Absolute Independence

- ▶ Two random variables are **independent** when knowing one *tells us nothing* about the other
- ▶ Dice example: just because one die comes up 5, this tells us nothing about what the next die will be
 - ▶ The **conditional probability** of the second die being 5, given that the first one was, is the *same as before*
- ▶ To get the probability of two 5's, we multiply
 - ▶ 1/6 of the time, the 1st die comes up 5
 - ▶ Of those cases, 1/6 of the time, the 2nd die is 5 too
 - ▶ Thus, the joint probability is $(1/6 \times 1/6) = 1/36$

▶

2

Defining Independence

- ▶ **General Definition:** X and Y are **independent** whenever the following holds (all are equivalent):

$$1. P(X | Y) = P(X)$$

$$2. P(Y | X) = P(Y)$$

$$3. P(X, Y) = P(X)P(Y)$$

▶

3

Examples of Independence

- ▶ Given a fair coin, the result of the current flip is independent of all prior flips:
 - $P(H) = 0.5$
 - $P(H | T) = P(H) = 0.5$
 - $P(H | TTHHHTTTTHH) = P(H) = 0.5$
 - $P(H | HHHHHHHHHHHHHHHHHHHHHHHHH) = P(H) = 0.5$
- ▶ Examples abound in games of chance (cards, dice, roulette)
- ▶ Independence assumptions make many calculations and models much simpler (*if they correctly apply*)

▶

4

Examples of Non-Independence

- Many real-life variables are NOT independent!
- Often, knowing that the one variable holds makes the other *much more likely* to hold true:

$$P(\text{PassTest} | \text{StudyTest}) \geq P(\text{PassTest})$$

$$P(\text{Fire} | \text{Smoke}) \geq P(\text{Fire})$$

$$P(\text{TorontoWins} | \text{RedSoxError}) \geq P(\text{TorontoWins})$$

- Other times, evidence can make it *much less likely*:

$$P(\text{BucsWin} | \text{BradyRetires}) \leq P(\text{BucsWin})$$

$$P(\text{Innocent} | \text{SmokingGun}) \leq P(\text{Innocent})$$

$$P(\text{HealthCare} | \text{Poverty}) \leq P(\text{HealthCare})$$

5

Simplifying Joint Distributions Using Independence Relations

- Suppose we have a joint distribution on cavities and toothaches:

	toothache	\neg toothache
cavity	0.1	0.02
\neg cavity	0.08	0.8

- Now consider an independent weather variable with distribution:

sunny	\neg sunny
0.7	0.3

- We could combine the two into a single table:

	sunny		\neg sunny	
	toothache	\neg toothache	toothache	\neg toothache
cavity	0.07	0.014	0.03	0.006
\neg cavity	0.056	0.56	0.024	0.24

- Or, save space and simply multiply when we need to calculate one of the combinations, storing fewer values overall (4 vs. 7)

6

Conditional Independence

- Sometimes, we find that two variables are not independent when considered *by themselves*
- However, if we had some *other evidence* (i.e., knew the value of other variables), then **conditional upon** that new evidence, they *are in fact* independent
- Again, we can define this in three equivalent ways; we say that X and Y are **conditionally independent given** Z iff:

$$1. P(X, Y | Z) = P(X | Z)P(Y | Z)$$

$$2. P(X | Y, Z) = P(X | Z)$$

$$3. P(Y | X, Z) = P(Y | Z)$$

7

An Example

- A classic case is multiple effects with a *common cause*
- Consider disease diagnosis, with 3 binary variables: B (patient has bubonic plague), P (patient has pustules), and G (patient has swollen glands)
 - Generally, there are all sorts of dependencies between the two symptoms: i.e., having one can make the other more likely
- However, if we actually *know* if someone has plague or not, then the chance of having pustules *does not* depend on whether they have swollen glands (i.e., it gives us no additional information)
 - $P(P = \text{true} | B = \text{true}, G) = P(P = \text{true} | B = \text{true})$
 - $P(P = \text{true} | B = \text{false}, G) = P(P = \text{true} | B = \text{false})$
- Variables P and G are **conditionally independent given** B

8

Example, continued

- ▶ The full joint distribution on 3 binary variables (B, P, G) requires $(2 \times 2 \times 2) - 1 = 7$ values
- ▶ However, we can use the **chain rule**, along with conditional independence to reduce this number:

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

$$= P(P | B) P(G | B) P(B)$$

- ▶ Now we only need 2 values for each of the first 2 distributions, plus 1 for the second: **5 total values**
- ▶ As the number of variables increases in a model or system, this sort of savings can grow very large

Artificial Intelligence (CS 131)

9

9

Naïve Bayesian Models

- ▶ Plague example reflects a common assumption that effects (symptoms) are all completely independent of one another, given a common cause (disease)
- ▶ If we let D be a disease, and s_1, s_2, \dots, s_n be the symptoms, this naïve Bayes assumption gives us:

$$P(D, s_1, s_2, \dots, s_n) = P(D) P(s_1 | D) P(s_2 | D) \dots P(s_n | D)$$

- ▶ Thus, we can calculate probability of disease given symptoms:

$$P(D | s_1, \dots, s_n) = \frac{P(D, s_1, \dots, s_n)}{P(s_1, \dots, s_n)}$$

$$= \frac{P(D) P(s_1 | D) P(s_2 | D) \dots P(s_n | D)}{P(s_1, \dots, s_n)}$$

Artificial Intelligence (CS 131)

10

10

Updating Naïve Beliefs

- ▶ This basic **factored representation** allows us to *update* our beliefs easily given new evidence
- ▶ Our *current belief* in the joint probability of disease and symptoms is given by the formula:

$$P(D, s_1, \dots, s_n) = P(D) \prod_{i=1}^n P(s_i | D)$$

- ▶ When a new symptom (s_{n+1}) appears, we can then easily account for it as follows, with minimal need to re-calculate:

$$P(D, s_1, \dots, s_n, s_{n+1}) = P(D) \prod_{i=1}^{n+1} P(s_i | D)$$

$$= P(D, s_1, \dots, s_n) P(s_{n+1} | D)$$

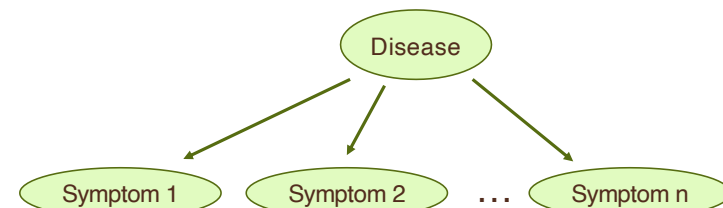
Artificial Intelligence (CS 131)

11

11

A Graphical Representation

- ▶ We could represent our assumption about connections between disease and symptoms using a **graphical model**
 - ▶ **Nodes** are the variables (disease, symptoms)
 - ▶ **Edges/arrows** show connections between them
 - ▶ No edge if things are **independent** of one another



Artificial Intelligence (CS 131)

12

12

Bayesian Networks (BNs)

► Generalizes the naïve disease example:

1. Set of **nodes**, one per variable.
2. Edges compose a **directed acyclic graph (DAG)**.
 - Each arrow-link, $A \rightarrow B$, between variables represents influence of variable A on variable B .
 - Cycles (loops) are forbidden, which both represents common sense, and makes certain algorithms (which we will see later) work properly
3. **Conditional probability tables (CPTs)** for each node, conditioned on its **parents** (if it has none, this is just a simple, non-conditional, one-variable distribution):

$$P(\text{Child} | \text{Parent}_1, \text{Parent}_2, \dots, \text{Parent}_n)$$

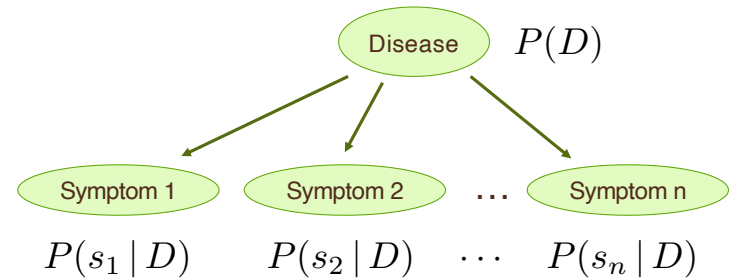


Artificial Intelligence (CS 131) 13

13

A Graphical Representation

► In disease example, with n possible symptoms we would have $(n + 1)$ CPTs, including one for disease itself:



Artificial Intelligence (CS 131) 14

14

Joint Probability in BNs

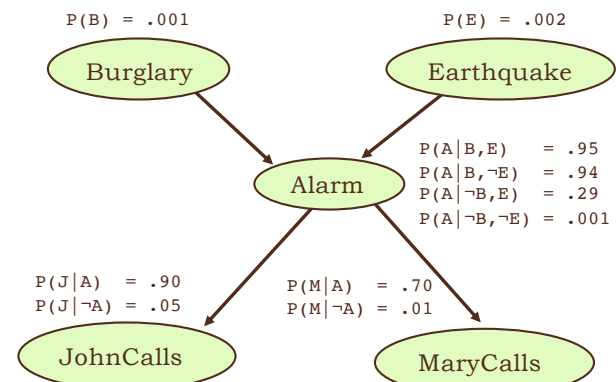
- Each BN is capable of representing a **full joint probability distribution** over its node variables, which means:
1. All probability questions of interest can be answered by looking at the numbers from the BN
 2. Memory cost of individual CPTs can be significantly smaller than that for the explicit full joint distribution
 3. Inference can be made much more simple



Artificial Intelligence (CS 131) 15

15

An Example: Earthquakes & Burglaries



Artificial Intelligence (CS 131)

16

16

Parameters of a BN

- ▶ In general, to specify a CPT in a Bayes Network, we need:
 - ▶ A CPT entry for every possible outcome of parent nodes
 - ▶ For nodes without parents, this is just the prior distribution
 - ▶ $n-1$ entries for each such, where n is the number of values the child node can have
- ▶ Thus, if $|X|$ is the number of values that a variable X has, the total number of parameters is:

$$\sum_{X \in BN} \left[(|X| - 1) \times \prod_{A \in Parents(X)} |A| \right]$$

Artificial Intelligence (CS 131) 17

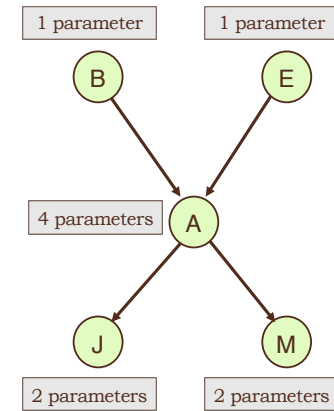
17

Parameters of a BN, cont'd.

- ▶ Thus for the earthquake case, where all the variables are binary, we require:

1. 1 value for $P(B)$ and $P(E)$
2. 4 values for $P(A | B, E)$
3. 2 values for $P(J | A)$ and $P(M | A)$

- ▶ The number of values needed gets smaller as the BN has more and more independence



Artificial Intelligence (CS 131) 18

18

General Principles for Building BN's

- ▶ To build a BN for a given domain, we need to:
 1. Figure out what *variables* we need (nodes)
 2. Figure out how they *depend upon* one another (edges)
 3. Give *probabilities* for these dependencies (CPTs)
- ▶ Often, identifying each of these is a matter of expert opinion or research

Artificial Intelligence (CS 131) 19

19

An Algorithm for Building BNs

Loop for $(i = 1 \dots n)$ variables:

1. Pick variable X_i to add to the graph.
2. Find *smallest* set of parents you can, choosing them so that the new variable is **conditionally independent** of all others previously added, given those parents:

$$P(X_i | X_1, X_2, \dots, X_{i-1}) = P(X_i | Parents(X_i))$$

3. Draw arrows from $Parents(X_i)$ to X_i
4. Specify the CPT: $P(X_i | Parents(X_i))$

Artificial Intelligence (CS 131) 20

20

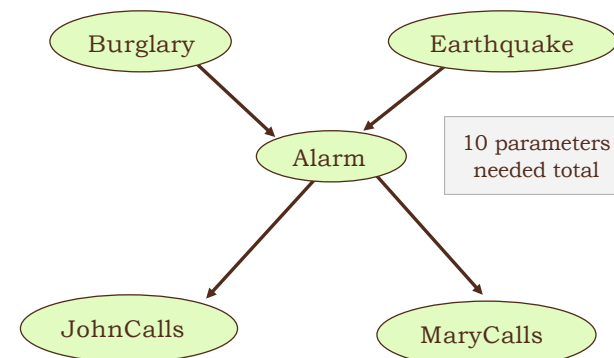
Properties of the Algorithm

- ▶ The graph is always **acyclic**
 - ▶ Arrows always go from existing prior nodes to new ones (and never vice-versa)
 - ▶ No looping back to nodes added before
- ▶ Result encodes all the probabilistic information we ever need to calculate any questions we might ask about the variables in the system
- ▶ *Compactness and sparseness* (number of connections) depends upon the order in which variables are considered when adding them to the BN, however

Artificial Intelligence (CS 131) 21

21

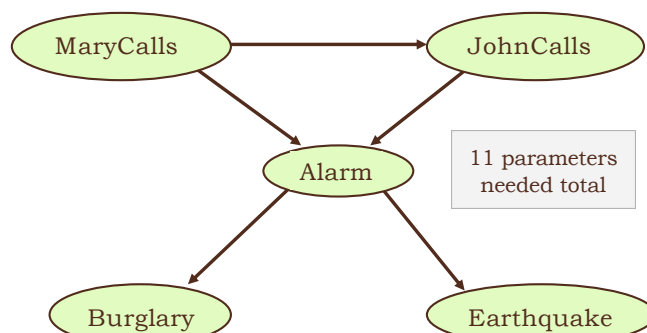
Original BN: Built Using Order (B, E, A, J, M)



Artificial Intelligence (CS 131) 22

22

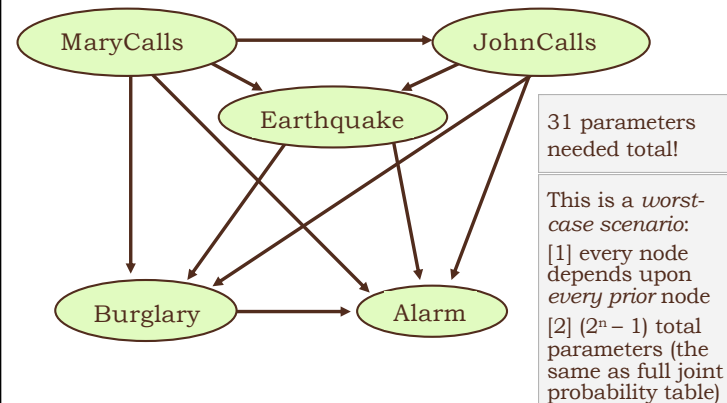
Another Order: (M, J, A, B, E)



Artificial Intelligence (CS 131) 23

23

Still Another Order: (M, J, E, B, A)



Artificial Intelligence (CS 131) 24

24