

Review

6/30

* Probabilities

$P(X)$ unconditional

$P(Y|X)$ conditional

$P(X, Y)$ joint

* Rules

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

$$P(X|Y) = \frac{P(Y|X) P(X)}{P(Y)} \quad \text{Bayes rule}$$

$$P(X) = \sum_y P(X, Y=y) \quad \text{marginalization}$$

* Conditionalized versions

$$P(A, B, C, \dots | E) = P(A|E) P(B|A, E) P(C|A, B, E) \dots$$

$$P(X|Y, E) = P(Y|X, E) P(X|E) / P(Y|E)$$

$$P(X|E) = \sum_y P(X, Y=y | E)$$

* Independence

$$\left. \begin{aligned} P(X, Y) &= P(X) P(Y) \\ P(X|Y) &= P(X) \\ P(Y|X) &= P(Y) \end{aligned} \right\} \begin{array}{l} \text{each of these implies} \\ \text{the others} \end{array}$$

* Conditional independence

$$P(X, Y | E) = P(X | E) P(Y | E)$$

$$P(X | Y, E) = P(X | E)$$

$$P(Y | X, E) = P(Y | E)$$

each of these
implies the
others

Probabilistic inference

- Today: do probabilities capture patterns of common sense reasoning?

Examples - reasoning about:

- 1) multiple explanations of a single event
- 2) multiple events - common explanation
- 3) chain of intervening events

* Binary random variables

B = burglary?

E = earthquake?

A = alarm?

* Joint distribution.

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

product rule

* Pomain knowledge - what we're assuming to be true

$$- P(B=1) = 0.001 \rightarrow P(B=0) = 1 - P(B=1) = 0.999$$

$$\left. \begin{array}{l} - P(E=1|B=0) = 0.002 \\ P(E=1|B=1) = 0.002 \end{array} \right\} P(E=1|B) = P(E=1) = 0.002$$

Assumption: B & E are independent random variables

B	E	$P(A=1 B,E)$	$P(A=0 B,E)$
0	0	0.001	$1 - 0.001$
1	0	0.94	'
0	1	0.29	'
1	1	0.95	$1 - 0.95$

1) reasoning about multiple explanations.

$$\text{compare: } P(B=1) = 0.001$$

$$P(B=1|A=1) = ?$$

$$P(B=1|A=1, E=1) = ?$$

Bayes rule

$$P(B=1|A=1) = \frac{P(A=1|B=1) P(B=1)}{P(A=1)} \quad \leftarrow 0.001$$

Term in denominator :

$$P(A=1) = \sum_{\substack{b \in \{0,1\} \\ e \in \{0,1\}}} P(B=b, E=e, A=1) \quad \text{marginalization}$$

$$= \sum_{b,e} P(B=b) P(E=e | B=b) P(A=1 | E=e, B=b) \quad \text{product rule}$$

$$= \sum_{b,e} P(B=b) P(E=e) P(A=1 | B=b, E=e) \quad \text{independence}$$

$$\begin{aligned} &= P(B=0) P(E=0) P(A=1 | B=0, E=0) \\ &\quad + P(B=1) P(E=0) P(A=1 | B=1, E=0) \\ &\quad + P(B=0) P(E=1) P(A=1 | B=0, E=1) \\ &\quad + P(B=1) P(E=1) P(A=1 | B=1, E=1) \end{aligned}$$

$$= 0.00252$$

Term in numerator

$$P(A=1 | B=1) = \sum_{e \in \{0,1\}} P(A=1, E=e | B=1) \quad \text{conditionalized marginalization}$$

$$= \sum_e P(A=1 | E=e, B=1) P(E=e | B=1) \quad \text{conditionalized product rule}$$

$$= \sum_e P(A=1 | E=e, B=1) P(E=e) \quad \text{independence}$$

$$= P(A=1 | E=0, B=1) P(E=0) + P(A=1 | E=1, B=1) P(E=1)$$

$$= 0.94002$$

$$\therefore P(B=1 | A=1) = \frac{P(A=1 | B=1) P(B=1)}{P(A=1)}$$

$\downarrow 0.94002$ $\curvearrowright 0.001$
 $\uparrow 0.00252$

$$= 0.37$$

Compare to $P(B=1) = 0.001$

What about $P(B=1 | A=1, E=1)$?

* Conditionalized Bayes rule

$$P(B=1 | A=1, E=1) = \frac{P(A=1 | B=1, E=1) P(B=1 | E=1)}{P(A=1 | E=1)}$$

$$= \frac{0.95 \times 0.001}{P(A=1 | E=1) \leftarrow 0.29} = 0.0033$$

Term in denominator

$$P(A=1 | E=1) = \sum_b P(A=1, B=b | E=1) \quad \text{conditionalized marginalization}$$

= ... conditional product rule

= ... independence

$$= \sum P(A=1 | B=b, E=1) P(B=b) = 0.29$$

Summary

$$P(B=1) = 0.001$$

$$P(B=1 | A=1) = 0.37 \quad \uparrow \quad \text{non-monotonic}$$

$$P(B=1 | A=1, E=1) = 0.0033 \quad \downarrow$$

\Rightarrow earthquake "explains away" the alarm, decreasing our belief in burglary.

Arises from multiple (causal) explanations of observed event.

2) Multiple events with a common explanation.

Two more binary random variables.

J = John calls?

M = Mary calls?

Assumptions:

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

$$P(M|A) = P(M|A, B, E, J)$$

} conditional independence

* Joint distribution

$$= P(B) P(E|B) P(A|B, E) P(J|A, B, E) P(M|A, B, E, J) \quad \text{Product rule}$$

$$= P(B) P(E) P(A|B, E) P(J|A) P(M|A)$$

conditional independence

* Conditional probabilities (domain knowledge)

$$P(J=1 | A=0) = 0.05$$

$$P(J=1 | A=1) = 0.9$$

$$P(M=1 | A=0) = 0.01$$

$$P(M=1 | A=1) = 0.7$$

(compare $P(A=1) = 0.00252$ (from previous example)

$$P(A=1 | J=1) = ? \quad 0.0435 \quad \uparrow$$

$$P(A=1 | J=1, M=0) = ?$$

* Bayes rule

$$P(A=1 | J=1) = \frac{P(J=1) \cancel{P(A=1)} P(J=1 | A=1) P(A=1)}{P(J=1)}$$

$0.9 \quad 0.00252$
 0.0521

Term in denominator.:

$$P(J=1) = \sum_a P(A=a, J=1) \quad \text{marginalization}$$

$$= \sum_a P(A=a) P(J=1 | A=a) \quad \text{product rule}$$

$$= P(A=1) P(J=1 | A=1) + P(A=0) P(J=1 | A=0)$$

$$= (0.00252) (0.9) + (1 - 0.00252) (0.05)$$

$$= 0.0521$$

$$\therefore P(A=1 | J=1) = 0.0435$$

$$P(A=1 | J=1, M=0) = ?$$

* Bayes rule with multiple pieces of evidence

$$P(A=1 | J=1, M=0) = \frac{P(J=1, M=0 | A=1) P(A=1)}{P(J=1, M=0)}$$

$$= \frac{0.9 \quad 1-0.7 \quad 0.00252}{P(J=1, M=0) \quad 0.05}$$



(b/c J & M are conditionally ind (c.i.) given A)

Term in denominator

$$P(J=1, M=0) = \sum_a P(A=a, J=1, M=0) \text{ marginalization}$$

$$= \sum_a P(A=a) P(J=1 | A=a) P(M=0 | A=a, J=1) \text{ product rule}$$

$$= \sum_a P(A=a) P(J=1 | A=a) P(M=0 | A=a) \text{ conditional ind.}$$

$$= 0.05$$

$$P(A=1 | J=1, M=0) = \frac{(0.9)(1-0.7)(0.00252)}{0.05}$$

$$= 0.0136$$

$$P(A=1) = 0.00252$$

$$P(A=1 | J=1) = 0.0435 \quad \uparrow \quad \text{reproduces common sense.}$$

$$P(A=1 | J=1, M=0) = 0.0136 \quad \downarrow \quad \text{non-monotonic}$$

3) Reasoning about intervening events.

Compare $P(A=1) = 0.00252$

$$P(A=1 | J=1) = 0.0435 \quad \uparrow$$

$$P(A=1 | J=1, B=1) = ?$$

$$P(A=1 | J=1, B=1) = \frac{P(J=1 | A=1, B=1) P(A=1 | B=1)}{P(J=1 | B=1)}$$

conditionalized Bayes rule

$$= \frac{P(J=1 | A=1) P(A=1 | B=1)}{P(J=1 | B=1)} \quad \text{conditional independence}$$

$$= \frac{(0.9) (0.94002)}{(\quad)} \quad \text{from example \#1}$$

Denominator:

$$P(J=1 | B=1) = \sum_a P(A=a, J=1 | B=1) \quad \text{conditionalized marginalization}$$

$$= \sum_a P(A=a | B=1) P(J=1 | A=a, B=1) \quad \text{conditionalized product rule}$$

$$= \sum_a P(A=a | B=1) P(J=1 | A=a) \text{ conditional independence}$$

= plug & chug?

$$= 0.849$$

$$\begin{aligned} \therefore P(A=1 | J=1, B=1) &= \frac{(0.9)(0.94002)}{(0.849)} \\ &= \underline{\underline{0.9965}} \uparrow \uparrow \end{aligned}$$

Motivation

- Joint distribution $P(X_1=x_1, X_2=x_2, \dots, X_n=x_n)$ involves $O(2^n)$ numbers for n binary random variables.
- More compact representations
- more efficient algorithms for inference?

Alarm example

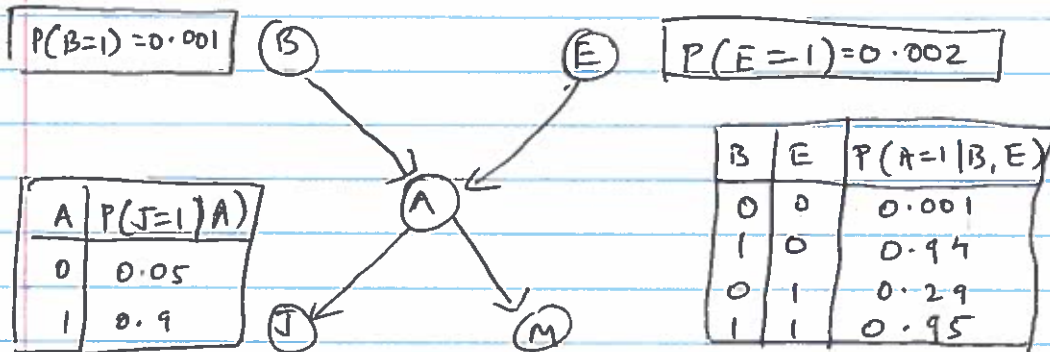
Binary variables $\{0, 1\}$

B	burglary	* Joint distribution
E	earthquake	
A	alarm	$P(B, E, A, J, M) = P(B)P(E B)$
J	John calls	$P(A B, E) P(J B, E, A) P(M B, E, A)$
M	Mary calls	product rule

$$= P(B) P(E) P(A|B,E) P(J|A) P(M|A)$$

conditional independence

* Directed acyclic graph (DAG)



* Conditional probability tables (CPTs)

A	$P(M=1 A)$
0	0.01
1	0.7

* Joint probabilities

$$P(B=1, E=0, A=1, J=1, M=1)$$

$$= P(B=1) P(E=0) P(A=1|B=1, E=0) P(J=1|A=1) P(M=1|A=1)$$

$$= (0.001) (1 - 0.002) (0.94) (0.9) (0.7) \dots$$

* Any query can be answered from joint distribution

$$\text{Ex: } P(B=1, E=0 \mid M=1)$$

$$\text{From product rule: } P(B=1, E=0 \mid M=1) = \frac{P(B=1, E=0, M=1)}{P(M=1)}$$

From marginalization

$$\text{numerator: } P(B=1, E=0, M=1) = \sum_{a, j} P(B=1, E=0, A=a, J=j, M=1)$$

$$\text{denominator: } P(M=1) = \sum_{b, e, a, j} P(B=b, E=e, A=a, J=j, M=1)$$

More efficient algorithms? Yes

Exploit structure of DAG.

Belief Network (BN)

* A BN is a DAG in which

- (i) nodes represent random variables
- (ii) edges represent conditional dependencies
- (iii) CPTs describe how each node depends on its parents.

* Conditional independence

Generally true in any domain that

$$\begin{aligned} P(x_1, x_2, \dots, x_n) &= P(x_1) P(x_2 | x_1) \dots P(x_n | x_1, \dots, x_{n-1}) \quad \text{product rule} \\ &= \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1}) \end{aligned}$$

In a given domain, suppose that

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

parents(x_i) is some subset of $\{x_1, x_2, \dots, x_{i-1}\}$.

BIG IDEA: represent dependencies by a graph!

* Constructing a BN:

1) choose random variables

2) ^{choose} Ordering

3) while there are variables left:

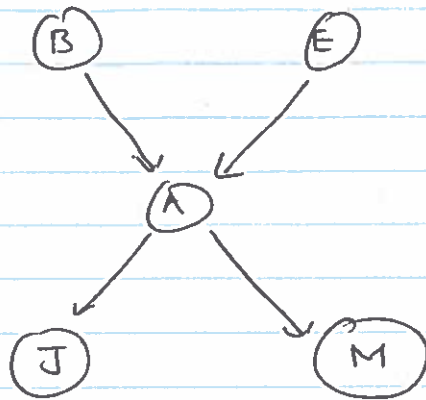
(a) add node x_i to the BN

(b) set parents of x_i to the minimal subset
Satisfying (*).

(c) define CPT $P(x_i | \text{pa}(x_i))$
 \uparrow parents of x_i

Ex:

$\{B, E, A, J, M\}$



* advantages

- ~~complete model~~, compact, consistent representation of joint probs

Ex: for binary variables, if $k = \max \# \text{parent. of any node (in-degree of graph)}$,

then $O(n2^k)$ numbers will be needed to write out CPTs,

versus $O(2^n)$ to represent joint distribution.

if $k \ll n$, huge savings!

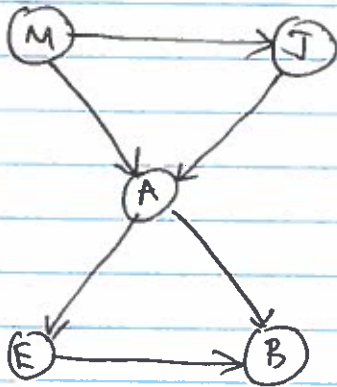
- clean separation of qualitative vs quantitative knowledge.

DAG encodes conditional independencies
CPTs encode numerical influences

* Node ordering

- Best order is to add "root causes", then the variables they influence, and so on.

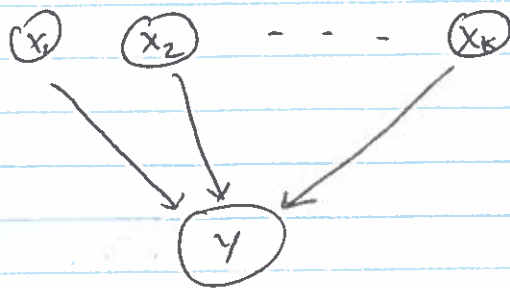
- Ex: "wrong order" $\{M, J, A, E, B\}$.
What BN do we get?



from misordered graph.

- (conditional) independences in world not obvious
- more numbers in CPT to specify same joint distribution.
- less natural, more difficult to assess or learn CPTs from data.

* Representing CPTs



for simplicity,
consider $x_i \in \{0, 1\}$
 $y \in \{0, 1\}$

How to represent $P(y=1 | x_1, x_2, \dots, x_k)$?

Possible answers:

i) lookup table $O(2^k)$ can store arbitrary CPT

2 ^k rows	x_1	x_2	...	x_k	$P(y=1 x_1, x_2, \dots, x_k)$
	0	0	- - -	0	0.6
	1	0	- - -	0	0.2
	⋮	⋮	⋮	⋮	⋮
	1	1	1 1 1	1	0.5

What if k is too large