

Knowledge Representation and Reasoning

Introduction to Probabilistic Reasoning System

Nur Ulfa Maulidevi
(ulfa@informatika.org)

KK IF – Teknik Informatika – STEI

ITB

Inteligensi Artifisial
(Artificial Intelligence)

Probabilistic Reasoning System (PRS)

Agent in real world need to handle uncertainty

Logical Agents handle uncertainty by disjunction → cannot tell us how likely the different conditions are

Probability theory provides a quantitative way of encoding likelihood

Given state(s) e , what is the probability that x happens → $P(x|e)$

Joint Probability Distribution

Bayesian/ Belief Network

Joint Probability Distribution

- Random variables

- Function: discrete domain $\rightarrow [0, 1]$
- Sums to 1 over the domain
 - Raining is a propositional random variable
 - $\text{Raining}(\text{true}) = 0.2$
 - $P(\text{Raining} = \text{true}) = 0.2$
 - $\text{Raining}(\text{false}) = 0.8$
 - $P(\text{Raining} = \text{false}) = 0.8$

- Joint distribution

- Probability assignment to all combinations of values of random variables

Inference using 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

- For any proposition ϕ , sum the atomic events where it is true: $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
- $P(\textit{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$
- $P(\textit{cavity}) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2$
- $P(\textit{cavity} \cup \textit{toothache}) = ?$
- $P(\neg \textit{cavity} \mid \textit{toothache}) = P(\neg \textit{cavity} \cap \textit{toothache}) / P(\textit{toothache})$
 $= (0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064)$
 $= 0.4$

Inference using 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

- If you have n binary propositional variables \rightarrow requires 2^n numbers to build Joint Probability Distribution
- \rightarrow Bayesian Network (We want to exploit independences in the domain)

Bayes' Rule:

$$\begin{aligned} P(A \mid B) &= P(A \cap B) / P(B) \\ &= P(B \mid A) P(A) / P(B) \end{aligned}$$

Conditioning:

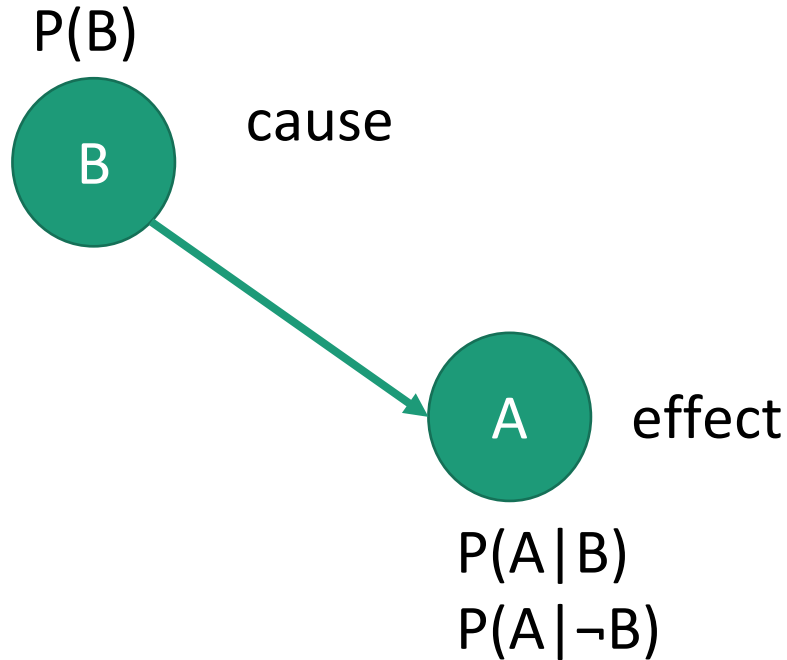
$$\begin{aligned} P(A) &= P(A \mid B) P(B) + P(A \mid \neg B) P(\neg B) \\ &= P(A \cap B) + P(A \cap \neg B) \end{aligned}$$

Structure of Bayesian Network

Nodes (variable)

Directed arc

Numerical
Parameters



Should be
Directed Acyclic Graph (DAG)

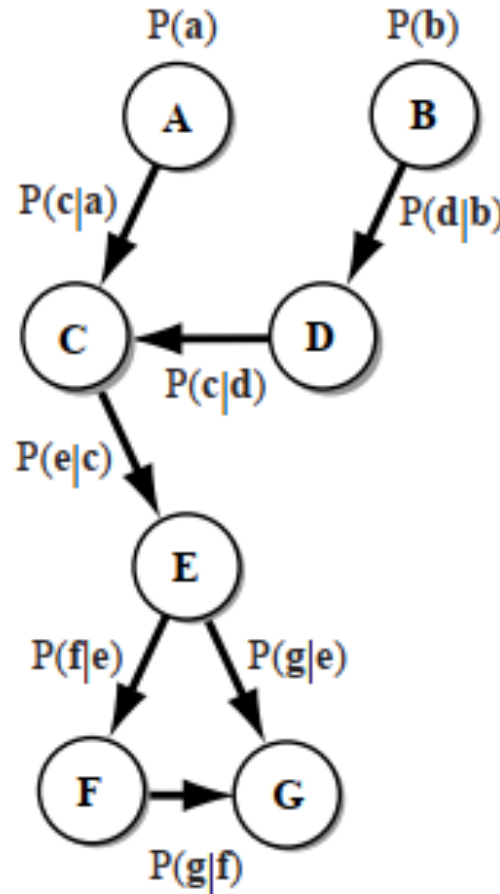
- Prior probability
- Probability conditional tables

Structure of Bayesian Network

Representation of causal dependencies graphically (Hart et al., 2001)

BN have capability probabilistic reasoning like full joint probability distribution. It can answer any question about the domain.

-- How we exploit Independence?



Bayesian belief nets

Causal networks

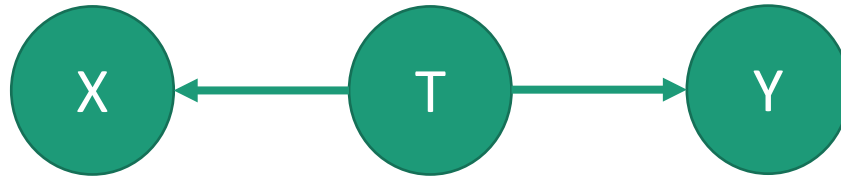
Belief nets

Independence

- A and B are independent iff
 - $P(A \cap B) = P(A) \cdot P(B)$
 - $P(A \mid B) = P(A)$
 - $P(B \mid A) = P(B)$
- Independence is essential for efficient probabilistic reasoning
- A and B are conditionally independent given C iff
 - $P(A \mid B, C) = P(A \mid C)$
 - $P(B \mid A, C) = P(B \mid C)$
 - $P(A \cap B \mid C) = P(A \mid C) \cdot P(B \mid C)$

Example of Independence

- X is late (X)
- Traffic Jam (T)
- Y is late (Y)
- None of these propositions are independent of one other
- X and Y are conditionally independent given T

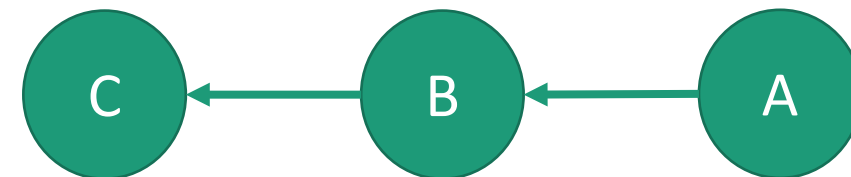
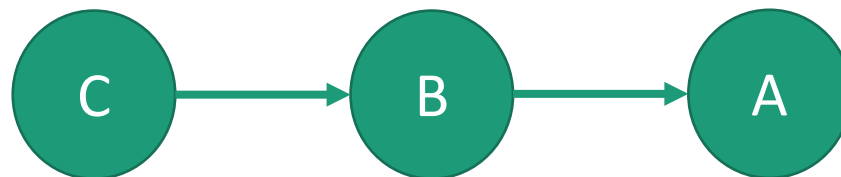


Types of Connections in Bayesian Network

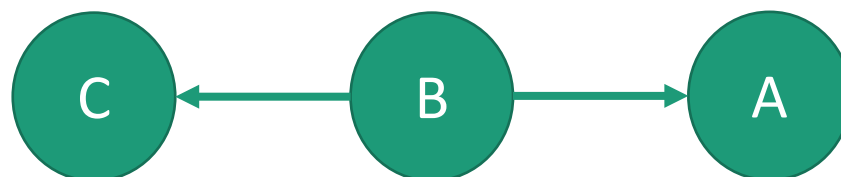
Forward

Backward

Serial

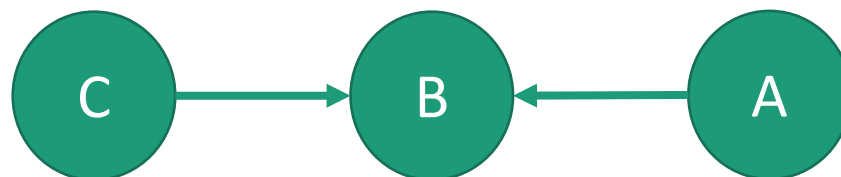


Diverging



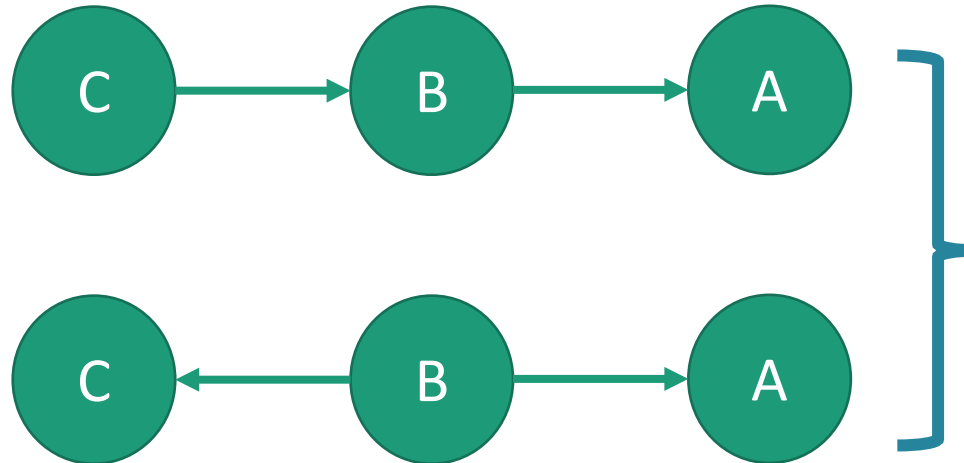
Why? To exploit Independence

Converging



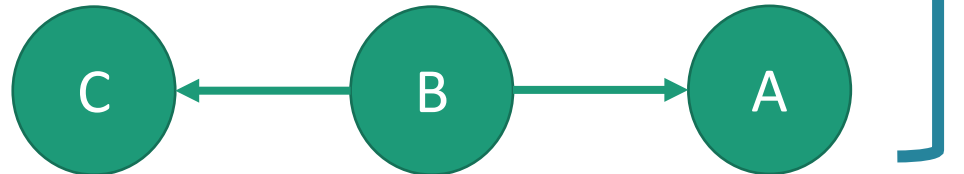
Independence in Connection

Serial



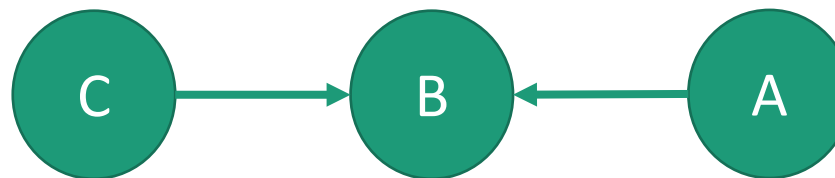
Knowing C will tell us about A, but if we know B, knowing C will tell us nothing about A (C and A conditionally independent or d-separated)

Diverging



Knowing C will tell us nothing about A without knowing B, but if we see evidence about B, C and A becomes dependent

Converging

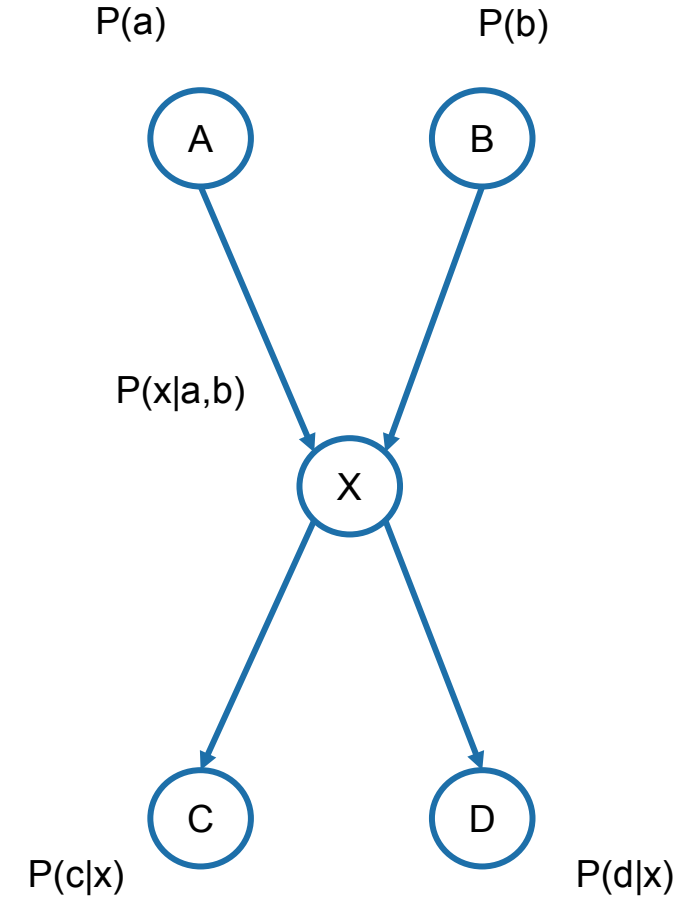


Reduction Number of Probabilities

In a domain with N binary propositional variables (2 possibilities value), one needs 2^N numbers to specify the joint probability distribution. $N=5$: need 32 probabilities

Independence and conditional independence relationships among variables can greatly reduce the number of probabilities that need to be specified in order to define the full joint distribution (Russel & Norvig, 2013)

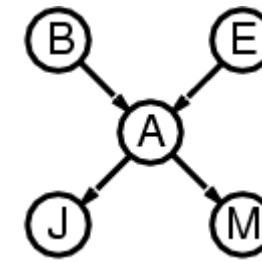
For 5 binary variables with causal networks: need $2+2+8+4+4=20$ probabilities (or 10 with complements).



Bayesian Network as Joint Probability Distribution (chain rule)

The full joint distribution is defined as the product of the local conditional distributions:

$$\mathbf{P}(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

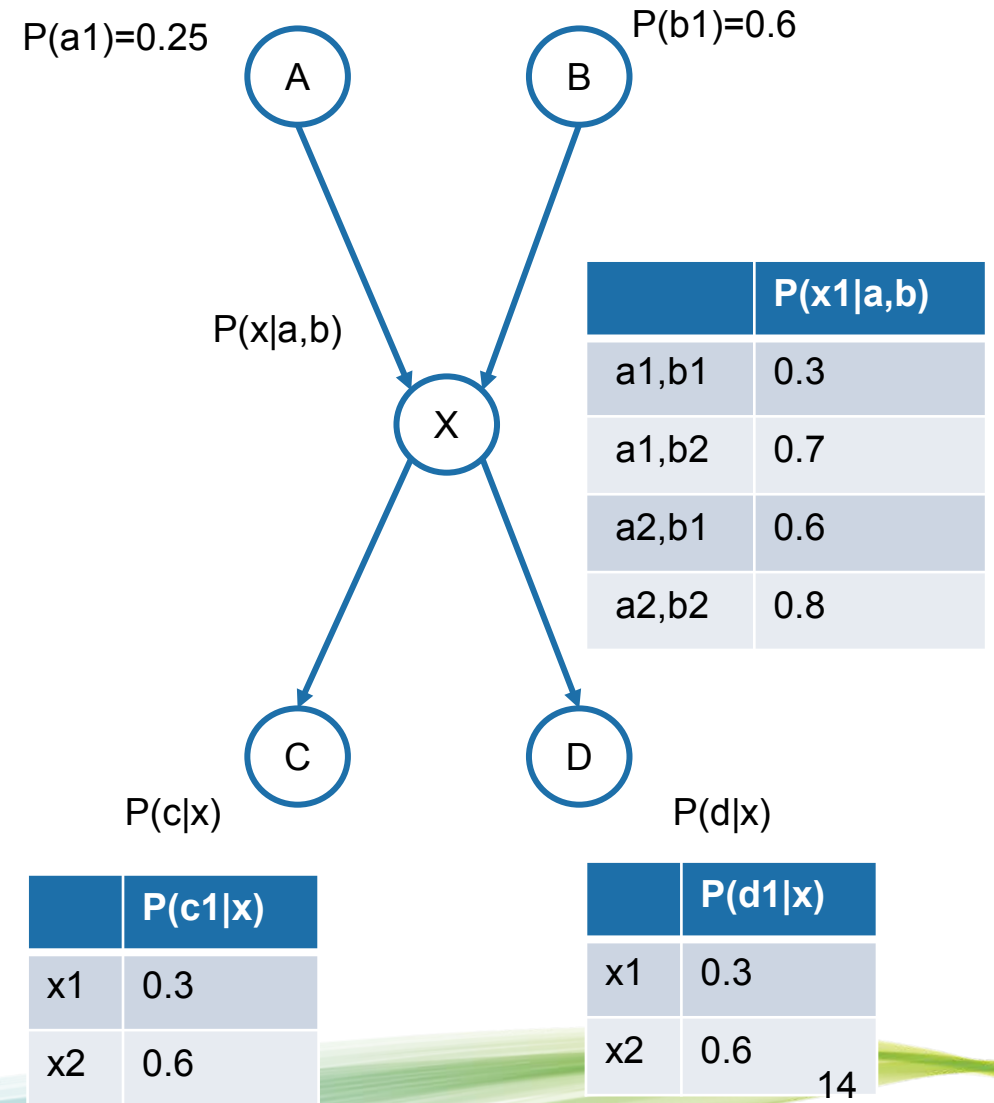


e.g., $\mathbf{P}(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$= \mathbf{P}(j \mid a) \mathbf{P}(m \mid a) \mathbf{P}(a \mid \neg b, \neg e) \mathbf{P}(\neg b) \mathbf{P}(\neg e)$$

Example

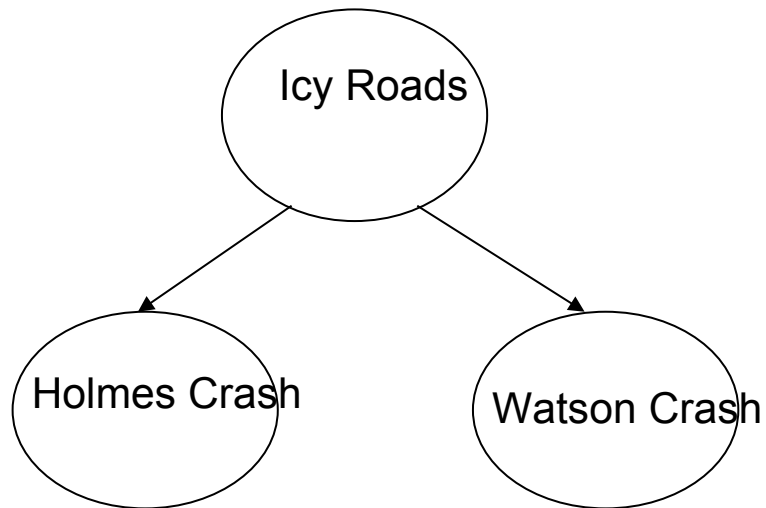
- We can determine the value of any entry in the joint probability.
- $P(a2,b1,x2,c2,d1)$
 $=P(a2)P(b1)P(x2|a2,b1)P(c2|x2)P(d1|x2)$
 $=0.75*0.6*0.4*0.4*0.6$
 $=0.0432$



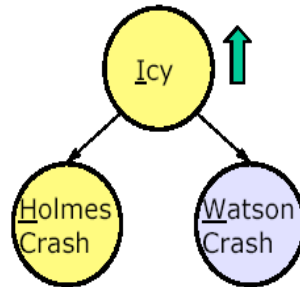
Example of Bayesian Network

Icy Roads

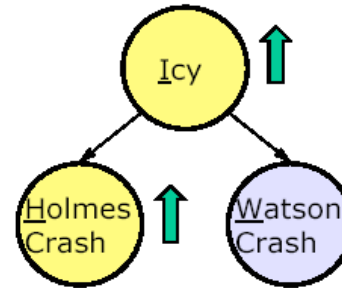
Inspector Smith is waiting for Holmes and Watson, who are driving (separately) to meet him. It is winter. His secretary tells him that Watson has had an accident. He says, "It must be that the roads are icy. I bet that Holmes will have an accident too. I should go to lunch." But, his secretary says, "No, the roads are not icy, look at the window." So, he says, "I guess I better wait for Holmes."



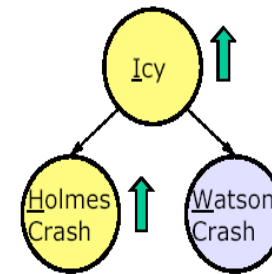
"Causal" Component



"Causal" Component

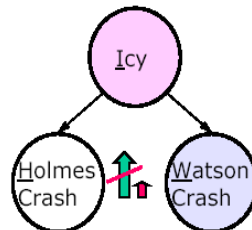


"Causal" Component



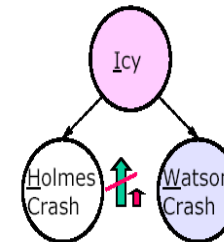
H and W are dependent,

"Causal" Component



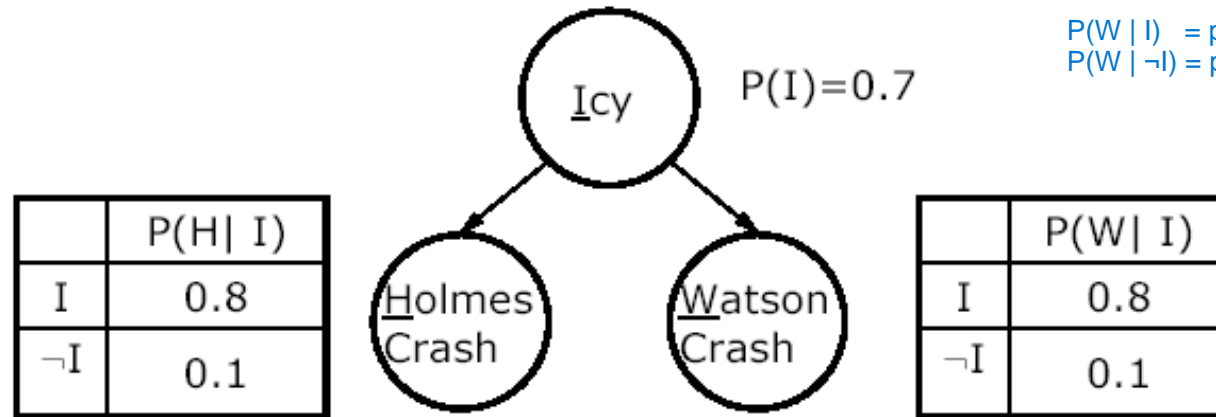
H and W are dependent,

"Causal" Component



H and W are dependent, but conditionally independent given I

Icy Roads with Numbers

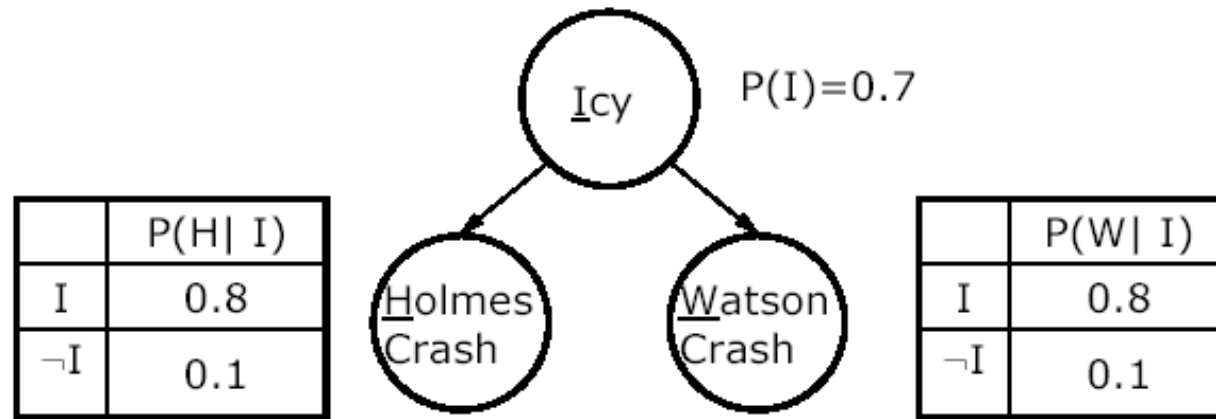


$P(W|I)$ = prob watson kecelakaan JIKA icy roads
 $P(W|\neg I)$ = prob watson kecelakaan JIKA TIDAK icy roads

Probability that Watson Crashes:

$$\begin{aligned} P(W) &= P(W|I) P(I) + P(W|\neg I) P(\neg I) \\ &= 0.8 \cdot 0.7 + 0.1 \cdot 0.3 \\ &= 0.56 + 0.03 \\ &= 0.59 \end{aligned}$$

Icy Roads with Numbers

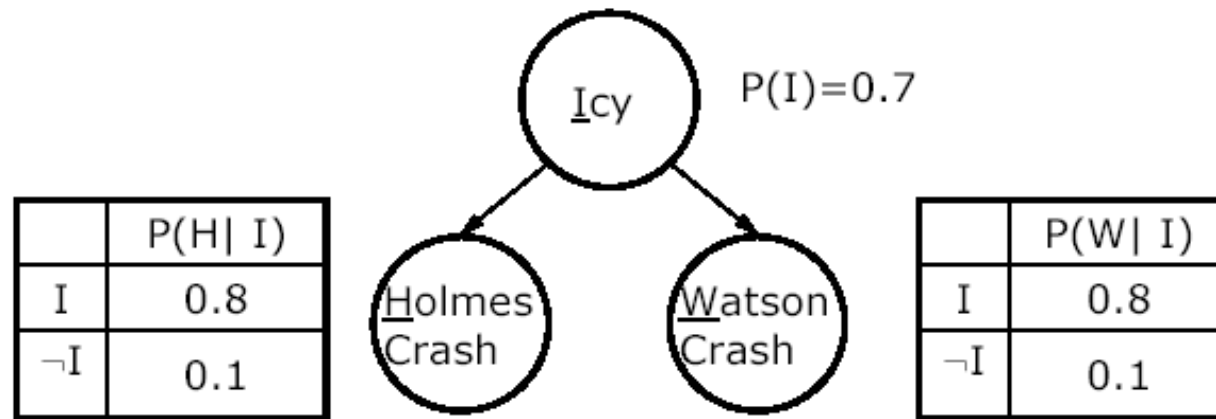


Probability of Icy given Watson (Bayes' Rule):

$$\begin{aligned} P(I | W) &= P(W | I) P(I) / P(W) \\ &= 0.8 \cdot 0.7 / 0.59 \\ &= 0.95 \end{aligned}$$

We started with $P(I) = 0.7$; knowing that Watson crashed raised the probability to 0.95

Icy Roads with Numbers



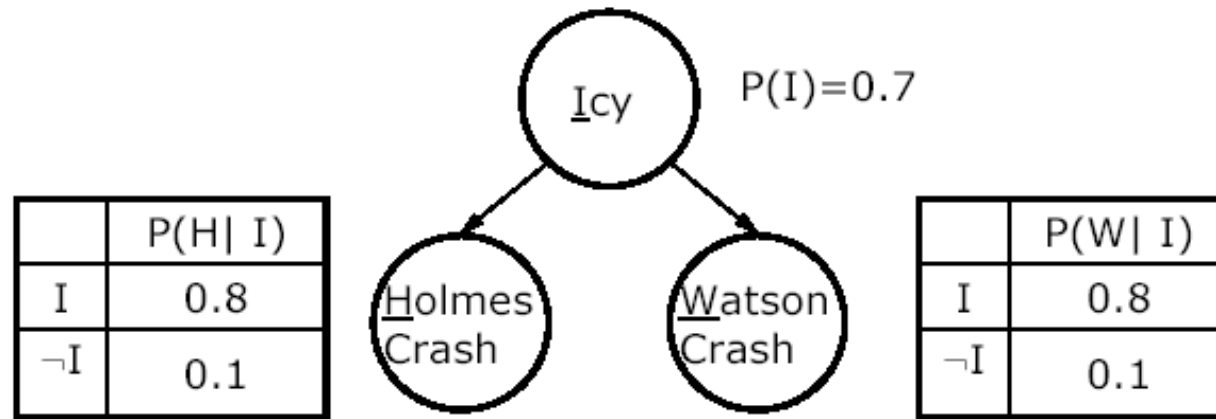
Probability of Holmes given Watson :

$$\begin{aligned} P(H|W) &= P(H|W, I)P(I|W) + P(H|W, \neg I) P(\neg I| W) \\ &= P(H|I)P(I|W) + P(H|\neg I) P(\neg I| W) \\ &= 0.8 \cdot 0.95 + 0.1 \cdot 0.05 \\ &= 0.765 \end{aligned}$$

We started with $P(H) = 0.59$; knowing that

Watson crashed raised the probability to 0.765

Icy Roads with Numbers



Probability of Holmes given Icy and Watson :
 $P(H|W, \neg I) = P(H|\neg I) = 0.1$

H and W are d-separated given I, so H and W are conditionally independent given I

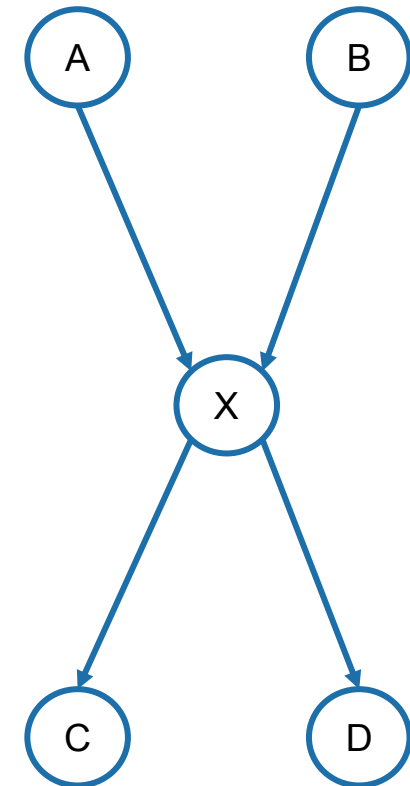
Classification using Bayesian Network

Masayu Leylia Khodra
(masayu@informatika.org)

KK IF – Teknik Informatika – STEI ITB

Inteligensi Artifisial
(*Artificial Intelligence*)

- X represents the fish : x1=salmon and x2=sea bass.
- X is influenced by A and B.
- A represents time of year: a1 = winter, a2 = spring, a3 = summer and a4 = autumn. Probability distribution on A is uniform.
- B represents geographical area where the fish was caught: b1 = north Atlantic and b2 = south Atlantic. The probabilities that any fish came from those areas are 0.6 and 0.4.
- C represents lightness with c1 = light, c2 = medium and c3 = dark
- D represents thickness with d1 = wide and d2 = thin.



The probability that the fish was caught in the summer in the north Atlantic and is a sea bass that is dark and thin.



The probability that the fish was caught in the summer (**a3**) in the north Atlantic (**b1**) and is a sea bass (**x2**) that is dark (**c3**) and thin (**d2**).



$P(a3, b1, x2, c3, d2)$

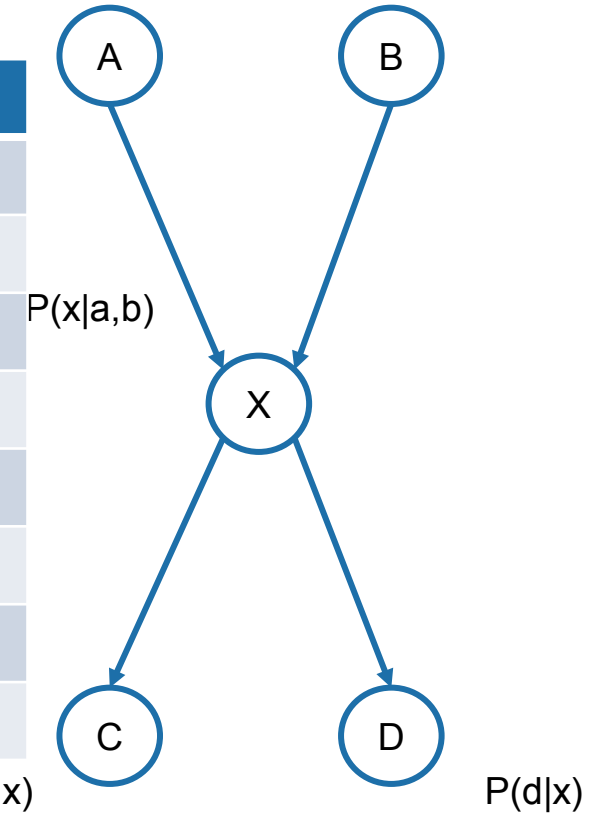
Inference: Example

$$P(a_3, b_1, x_2, c_3, d_2) = P(a_3)P(b_1)P(x_2|a_3, b_1)P(c_3|x_2)P(d_2|x_2) = 0.012$$

P(a)			
P(a1)	P(a2)	P(a3)	P(a4)
0.25	0.25	0.25	0.25

P(b)
P(b1)
0.6

	P(x1 a,b)
a1,b1	0.3
a1,b2	0.7
a2,b1	0.6
a2,b2	0.8
a3,b1	0.4
a3,b2	0.1
a4,b1	0.2
a4,b2	0.3



	P(c1 x)	P(c2 x)	P(c3 x)
x1	0.6	0.2	0.2
x2	0.2	0.3	0.5

	P(d1 x)
x1	0.3
x2	0.6

Classification

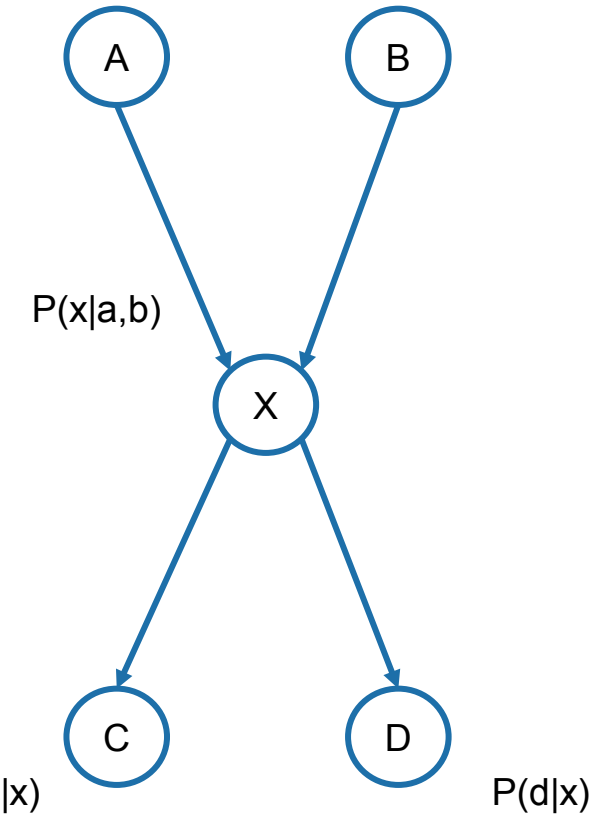
Classify the fish that is light (c1) and caught in the south Atlantic (b2), but we do not know what time of year the fish was caught nor its thickness.

Maximum a posterior probability:
 $P(x_1|c_1, b_2)$ vs $P(x_2|c_1, b_2)$

P(a)			
P(a1)	P(a2)	P(a3)	P(a4)
0.25	0.25	0.25	0.25

P(b)
P(b1)
0.6

	P(x1 a,b)
a1,b1	0.3
a1,b2	0.7
a2,b1	0.6
a2,b2	0.8
a3,b1	0.4
a3,b2	0.1
a4,b1	0.2
a4,b2	0.3



	P(c1 x)	P(c2 x)	P(c3 x)
x1	0.6	0.2	0.2
x2	0.2	0.3	0.5

	P(d1 x)
x1	0.3
x2	0.6

Classification (2)

Q: query

e: evidence of all variables

$$P(Q|e) = P(q,e)/P(e) = \alpha P(Q,e)$$

$$P(x_1|c_1,b_2) = P(x_1,c_1,b_2)/P(c_1,b_2)$$

$$= \alpha \sum P(x_1, \mathbf{a}, b_2, c_1, \mathbf{d})$$

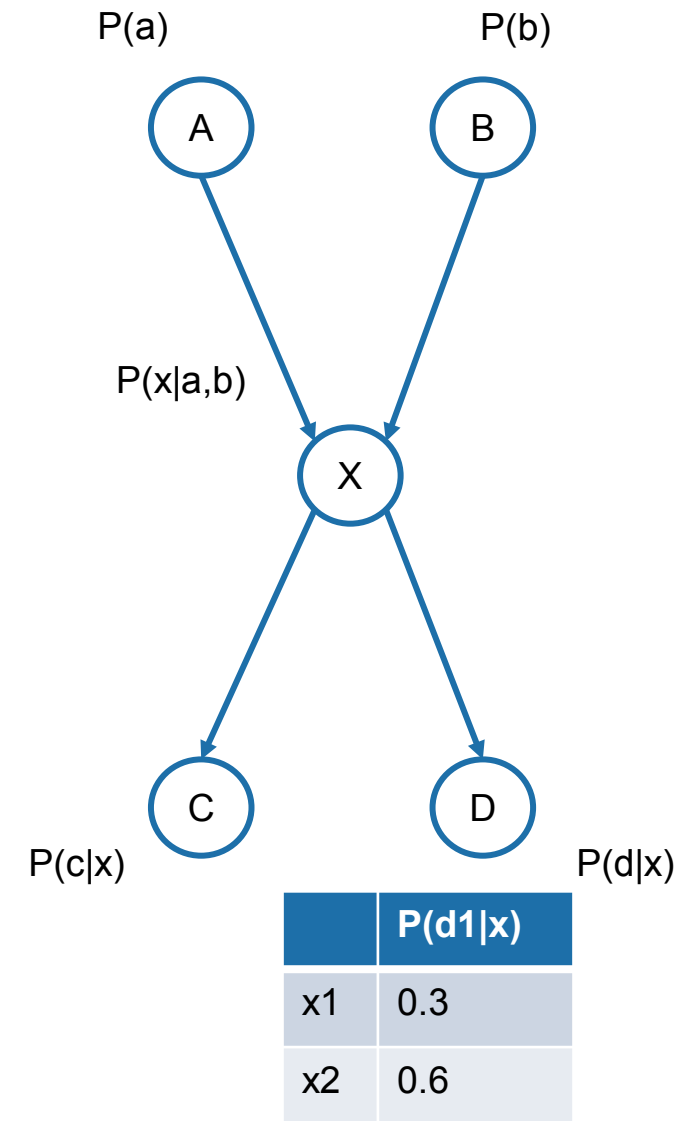
$$= \alpha \sum P(\mathbf{a}).P(b_2).P(x_1|\mathbf{a},b_2).P(c_1|x_1).P(\mathbf{d}|x_1)$$

$$= \alpha P(b_2).P(c_1|x_1) \sum P(\mathbf{a}).P(x_1|\mathbf{a},b_2).P(\mathbf{d}|x_1)$$

$$= \alpha P(b_2).P(c_1|x_1) [\sum P(\mathbf{a}).P(x_1|\mathbf{a},b_2)] [\sum P(\mathbf{d}|x_1)] = \alpha$$

$$\begin{aligned} &P(\mathbf{a}_1).P(x_1|\mathbf{a}_1,b_2) + \\ &P(\mathbf{a}_2).P(x_1|\mathbf{a}_2,b_2) + \\ &P(\mathbf{a}_3).P(x_1|\mathbf{a}_3,b_2) + \\ &P(\mathbf{a}_4).P(x_1|\mathbf{a}_4,b_2) \end{aligned}$$

$$P(d_1|x_1) + P(d_2|x_1) = 1.0$$



Classification (3)

$$P(x_1|c_1,b_2)=P(x_1,c_1,b_2)/P(c_1,b_2) \\ =\alpha P(b_2).P(c_1|x_1).[\sum P(\mathbf{a}).P(x_1|\mathbf{a},b_2) \\]. [\sum P(\mathbf{d}|x_1)]= \alpha 0.114$$

$$P(x_2|c_1,b_2)=P(x_2,c_1,b_2)/P(c_1,b_2) \\ =\alpha P(b_2).P(c_1|x_2) [\sum P(\mathbf{a}). \\ P(x_2|\mathbf{a},b_2)][\sum P(\mathbf{d}|x_2)]= \alpha 0.042$$

Normalize:

$$P(x_1|c_1,b_2)=0.73$$

$$P(x_2|c_1,b_2)=0.27$$

Decision: x_1 =salmon

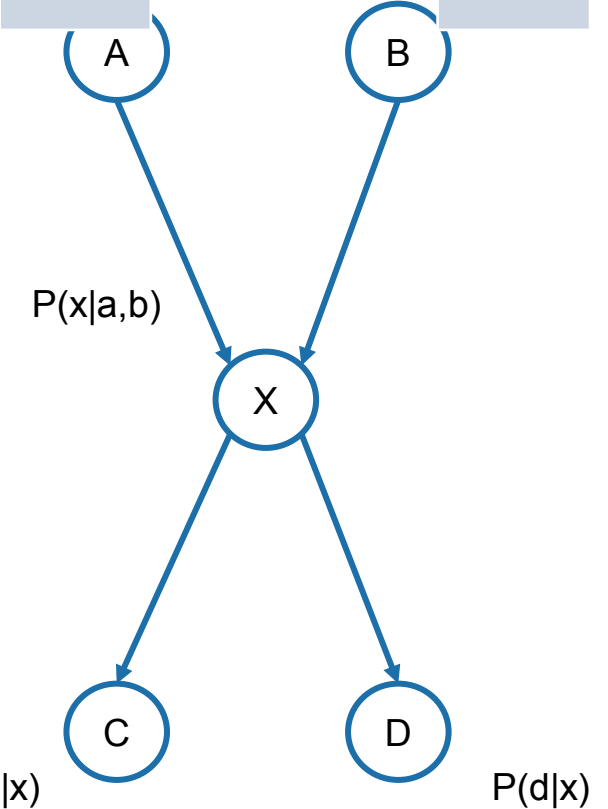
P(a1)	P(a2)	P(a3)	P(a4)
0.25	0.25	0.25	0.25

	P(x1 a,b)
a1,b1	0.3
a1,b2	0.7
a2,b1	0.6
a2,b2	0.8
a3,b1	0.4
a3,b2	0.1
a4,b1	0.2
a4,b2	0.3

P(a)

P(b)

P(b1)
0.6



	P(c1 x)	P(c2 x)	P(c3 x)
x1	0.6	0.2	0.2
x2	0.2	0.3	0.5

	P(d1 x)
x1	0.3
x2	0.6

Knowledge Representation and Reasoning

Learning from Data

Masayu Leylia Khodra
(masayu@informatika.org)

KK IF – Teknik Informatika – STEI ITB

Inteligensi Buatan
(*Artificial Intelligence*)

Problems in BN constructed by human expert

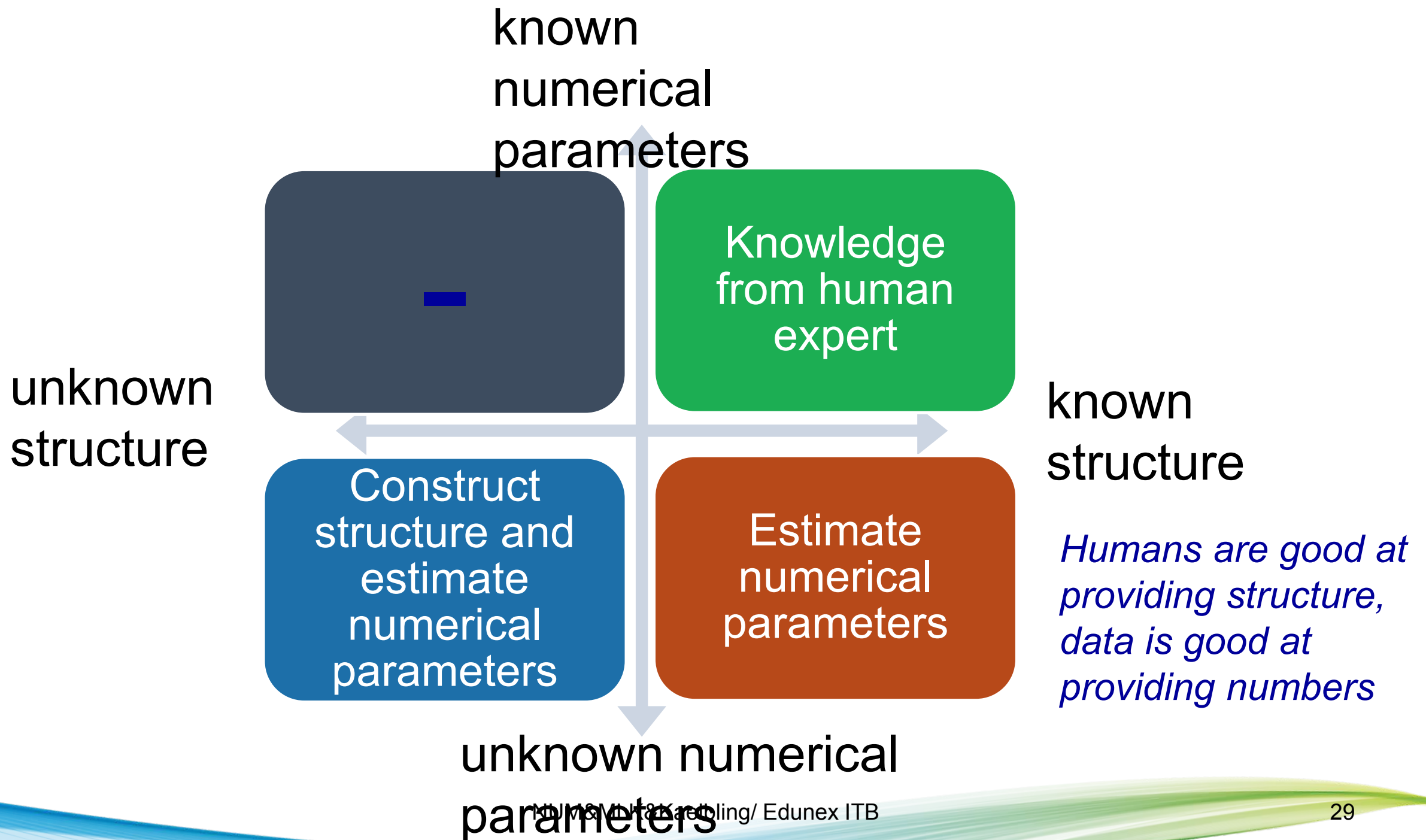
Knowledge acquisition bottleneck problem

Knowledge elicitation problem: slow speed, and inability of expert to express the knowledge they possess.

Very difficult in getting reliable probability estimates.

Ideally, determine probability by data distribution, experience, and assumption (subjective).

What is difference between $P(a)=0.5$ and $P(a)=0.6$?



Parameter Estimation

Given structure
with m nodes

Given a data set

$D =$
 $\{ \langle v_1^1, \dots, v_m^1 \rangle, \dots, \langle v_1^k, \dots, v_m^k \rangle \}$

Count $\#(V_i=T)$, $\#(V_i=F)$,
 $\#(V_i=T, V_j=T)$, $\#(V_i=T, V_j=F)$

Variable V_i with no parent
$$P(V_i) \approx \frac{\#(V_i = T)}{k}$$

Variable V_i with parent V_j
$$P(V_i|V_j) \approx \frac{\#(V_i = T, V_j = T)}{\#(V_j = T)}$$

$$P(V_i|\neg V_j) \approx \frac{\#(V_i = T, V_j = F)}{\#(V_j = F)}$$

Variable V_i with no parent

$$P(V_i) \approx \frac{\#(V_i = T) + 1}{k + 2}$$

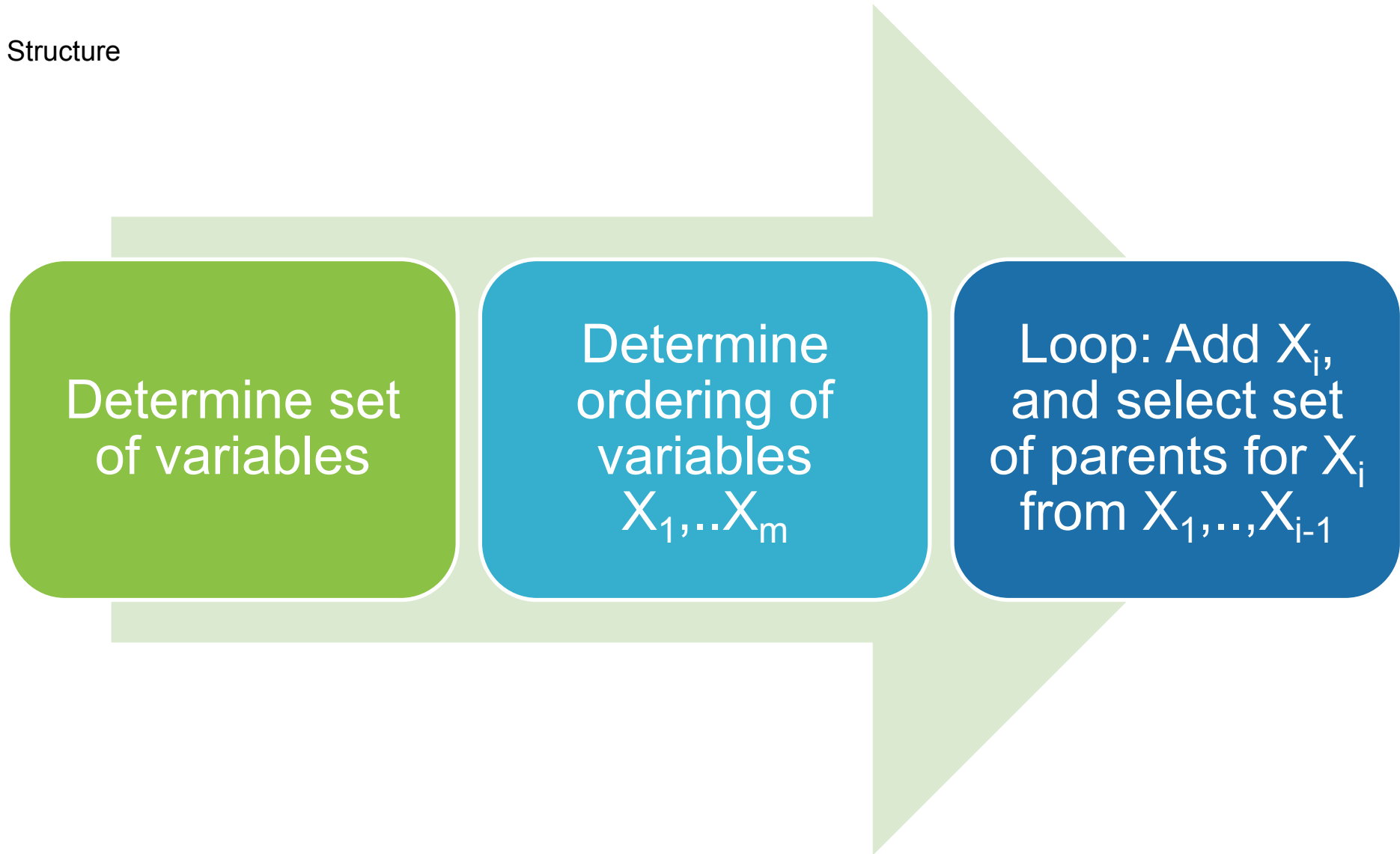
*kemungkinan
nilainya True/False
⇒ ada 2*

Variable V_i with parent V_j

$$P(V_i|V_j) \approx \frac{\#(V_i = T, V_j = T) + 1}{\#(V_j = T) + 2}$$

$$P(V_i|\neg V_j) \approx \frac{\#(V_i = T, V_j = F) + 1}{\#(V_j = F) + 2}$$

Construct BN Structure



Construct BN Structure: Example

Based on causal knowledge:

- Set of variables: A,B,C,D,X
- Ordering of variables: A,B,C,D,X
- Add A: no parent.
- Add B: is A parent of B ?
- Add C: parent of C ?
- Add D: parent of D ?
- Add X: parent of X ?

Based on data:

- Set of variables: A,B,C,D,X
- Ordering of variables: A,B,C,D,X
- Add A: no parent.
- Add B: $P(B|A)=P(B)$?
- Add C:
 - $P(C|A)=P(C)$?
 - $P(C|B)=P(C)$?
 - $P(C|A,B)=P(C|A)$?
 - $P(C|A,B)=P(C|B)$?
 - $P(C|A,B)=P(C)$?
- Etc.

is A the parent of B??

cara ceknya computer akan hitung $P(B|A) = P(B)$
kalau iya berarti mereka saling independen dan A bukan parent dari B

$P(C|A,B) = P(C,B)$
A tidak mempengaruhi C, A bukan parent C.
 $P(C|B) \neq P(C)$.
berarti parent C adalah B.

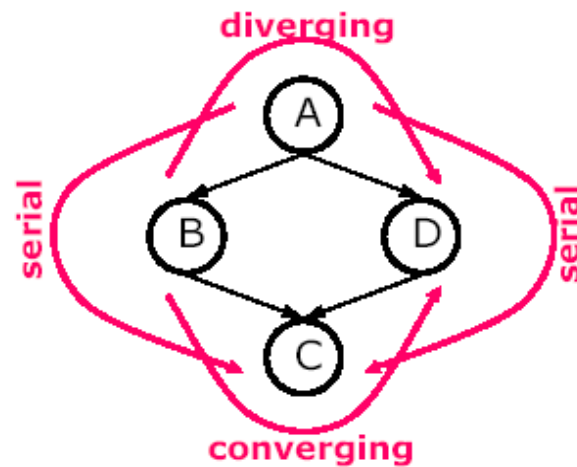
D Separation

Two variables A and B are d-separated iff for every path between them, there is an intermediate variable V such that either

- The connection is serial or diverging and V is known
- The connection is converging and neither V nor any descendant is instantiated

Two variables are d-connected iff they are not d-separated

dalam BN,
boleh ubah² panah asal
ga bikin converging
baru + ga menghilangkan
converging lama



- A-B-C: serial, blocked when B is known, connected otherwise
- A-D-C: serial, blocked when D is known, connected otherwise
- B-A-D: diverging, blocked when A is known, connected otherwise
- B-C-D: converging, blocked when C has no evidence, connected otherwise

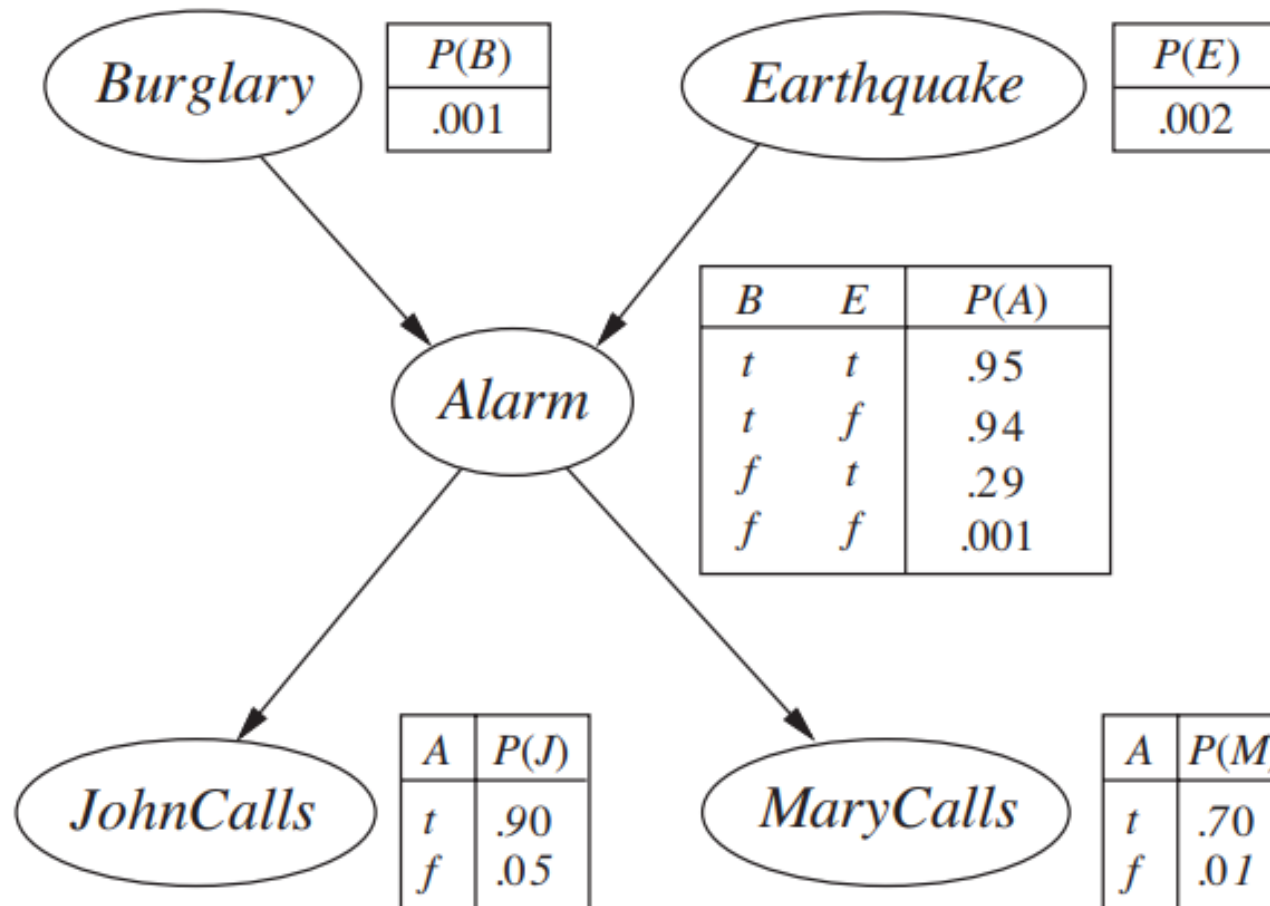
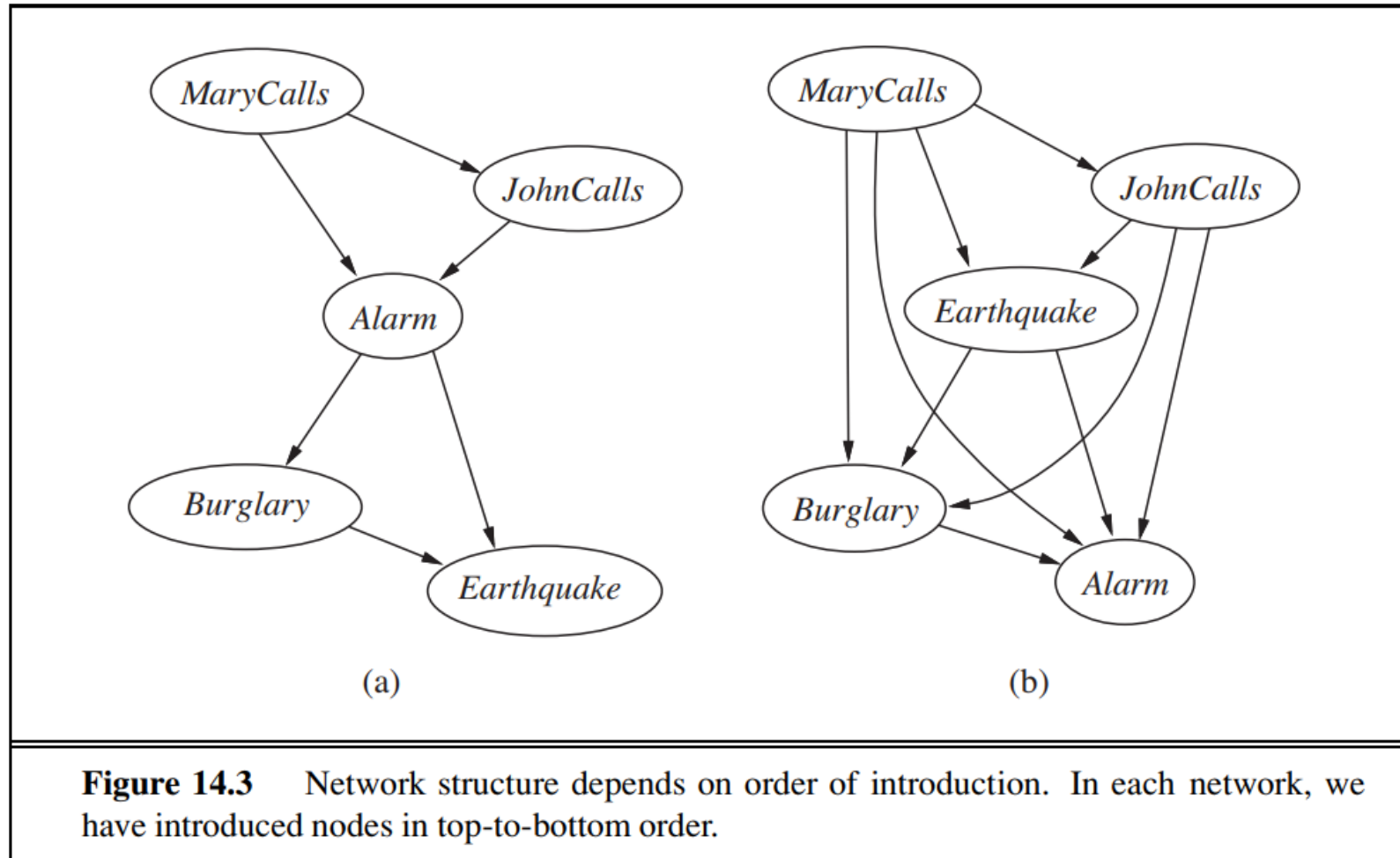


Figure 14.2 A typical Bayesian network, showing both the topology and the conditional probability tables (CPTs). In the CPTs, the letters B , E , A , J , and M stand for *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, and *MaryCalls*, respectively.

Variable Ordering Influences Structure

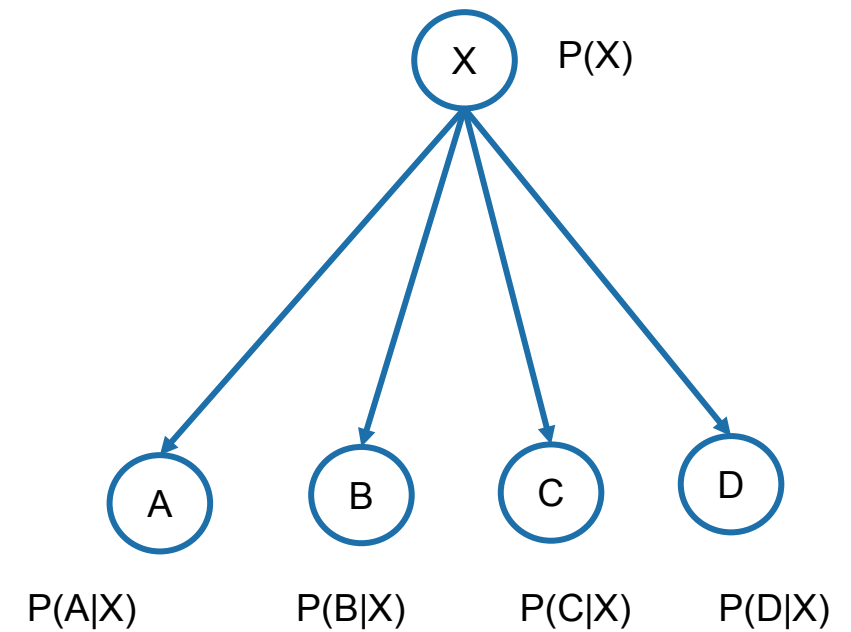


Russel & Norvig (2010). AI.

Construct BN Structure: Example 2

Based on causal knowledge:

- Set of variables: X, A, B, C, D
- Ordering of variables: X, A, B, C, D
- Add X : no parent.
- Add A : $\text{parent}(A)=X$.
- Add B : $\text{parent}(B)=X$. $P(B|A,X)=P(B|X)$
- Add C : $\text{parent}(C)=X$. $P(C|A,B,X)=P(C|X)$
- Add D : $\text{parent}(D)=X$. $P(D|A,B,C,X)=P(D|X)$

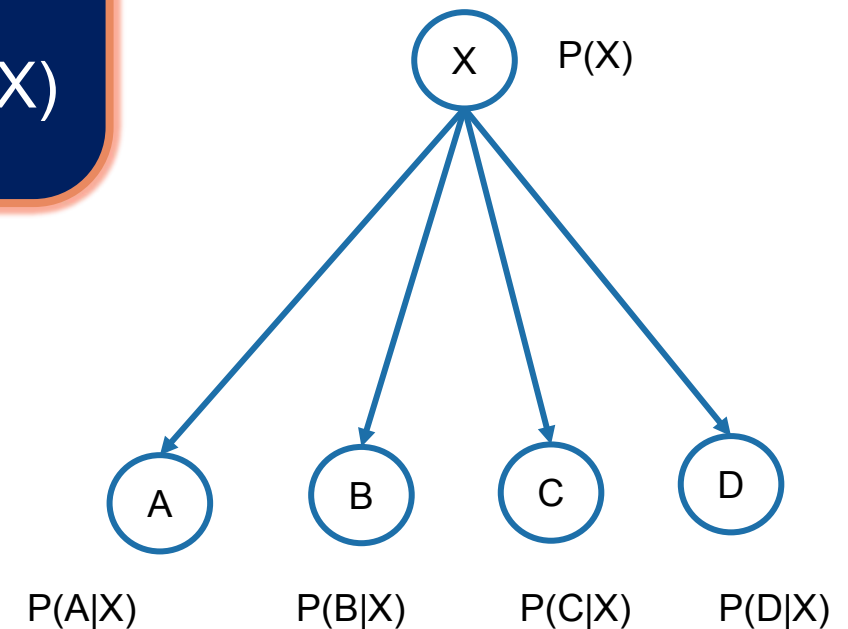


Naive Bayes as “Special” Bayesian Network

What is the Bayes Network for X_1, \dots, X_n with NO assumed conditional independencies?

Probability Model of $P(X)$ and $P(a_i|X)$
 $P(X, A, B, C, D) = P(X) \cdot P(A|X) \cdot P(B|X) \cdot P(C|X) \cdot P(D|X)$

Classification:
Find the maximum $P(X | A, B, C, D)$



Example: Play Tennis Dataset

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes

outlook	temp.	humidity	windy	play
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

Frequency of
(sunny|yes) \rightarrow 2

Frequency of
(overcast|yes) \rightarrow 4

Frequency of
(hot|yes) \rightarrow 2

...

Frequency of
class 'yes' \rightarrow 9

Frequency of
(sunny|no) \rightarrow 3

Frequency of
(overcast|no) \rightarrow 0

Frequency of
(hot|no) \rightarrow 2

Frequency of
class 'no' \rightarrow 5

Example: Play Tennis Dataset

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes

outlook	temp.	humidity	windy	play
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

Frequency of
(sunny|yes) →

Frequency of
(sunny|no) →

outlook			temperature			humidity			windy			play	
	yes	no		yes	no		yes	no		yes	no	yes	no
sunny	2	3	hot	2	2	high	3	4	false	6	2	9	5
overcast	4	0	mild	4	2	normal	6	1	true	3	3		
rainy	3	2	cool	5	1								

Frequency of
class 'yes' → 9

Frequency of
class 'no' → 5

Example: Play Tennis Dataset

Frequency

outlook			temperature			humidity			windy			play	
	yes	no		yes	no		yes	no		yes	no	yes	no
sunny	2	3	hot	2	2	high	3	4	false	6	2	9	5
overcast	4	0	mild	4	2	normal	6	1	true	3	3		
rainy	3	2	cool	3	1								

$P(a_i|v_j)$

outlook			temperature			humidity			windy		
yes no			yes no			yes no			yes no		
sunny	2/9	3/5	hot	2/9	2/5	high	3/9	4/5	false	6/9	2/5
overcast	4/9	0/5	mild	4/9	2/5	normal	6/9	1/5	true	3/9	3/5
rainy	3/9	2/5	cool	3/9	1/5						

$P(v_j)$

play	
yes	no
9/14	5/14

Probability Model

Classify New Instance: <Sunny, Cool, High, True>

$P(a_i|v_j)$

outlook		temperature		humidity		windy	
	yes no		yes no		yes no		yes no
sunny	2/9 3/5	hot	2/9 2/5	high	3/9 4/5	false	6/9 2/5
overcast	4/9 0/5	mild	4/9 2/5	normal	6/9 1/5	true	3/9 3/5
rainy	3/9 2/5	cool	3/9 1/5				

$P(v_j)$

play
yes no
9/14 5/14

$$P(v_j | a_1, a_2, \dots, a_n) = P(v_j) \cdot \prod_i P(a_i | v_j)$$

$P(\text{yes}|\text{sunny, cool, high, true})$

$= P(\text{yes}) \cdot P(\text{sunny}|\text{yes}) \cdot P(\text{cool}|\text{yes}) \cdot P(\text{high}|\text{yes}) \cdot P(\text{true}|\text{yes})$

$= 9/14 \cdot 2/9 \cdot 3/9 \cdot 3/9 \cdot 3/9 = 0.0053$

$P(\text{no}|\text{sunny, cool, high, true})$

$= P(\text{no}) \cdot P(\text{sunny}|\text{no}) \cdot P(\text{cool}|\text{no}) \cdot P(\text{high}|\text{no}) \cdot P(\text{true}|\text{no})$

$= 5/14 \cdot 3/5 \cdot 1/5 \cdot 4/5 \cdot 3/5 = 0.0206$

THANK YOU