

# Bayesian Network Representation

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Joint and Conditional Distributions
- I-Maps
  - I-Map to Factorization
  - Factorization to I-Map
  - Perfect Map
- Knowledge Engineering
  - Picking Variables
  - Picking Structure
  - Picking Probabilities
  - Sensitivity Analysis

# Joint Distribution

- Company is trying to hire recent graduates
- Goal is to hire intelligent employees
  - No way to test intelligence directly
  - But have access to Student's SAT score
    - Which is informative but not fully indicative
- Two random variables
  - Intelligence:  $Val(I)=\{i^1, i^0\}$ , high and low
  - Score:  $Val(S)=\{s^1, s^0\}$ , high and low
- Joint distribution has 4 entries
  - Need three parameters

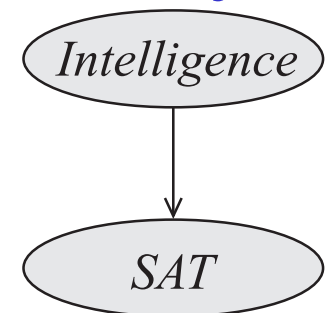
$I$	$S$	$P(I, S)$
$i^0$	$s^0$	0.665
$i^0$	$s^1$	0.035
$i^1$	$s^0$	0.06
$i^1$	$s^1$	0.24

# Alternative Representation: Conditional Parameterization

- $P(I, S) = P(I)P(S | I)$ 
  - Representation more compatible with causality
    - Intelligence influenced by Genetics, upbringing
    - Score influenced by Intelligence
- Note: BNs are not required to follow causality but they often do

- Need to specify  $P(I)$  and  $P(S|I)$

$i^0$	$i^1$	$I$	$s^0$	$s^1$
0.7	0.3	$i^0$	0.95	0.05
		$i^1$	0.2	0.8



- Three binomial distributions (3 parameters) needed

- One marginal, two conditionals  $P(S|I=i^0)$ ,  $P(S|i=i^1)$

# Naïve Bayes Model

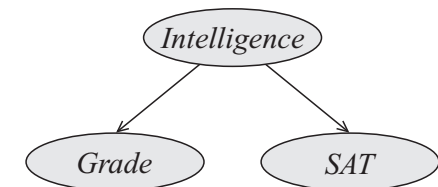
- Conditional Parameterization combined with Conditional Independence assumptions

–  $Val(G) = \{g^1, g^2, g^3\}$  represents grades  $A, B, C$

$I$	$g^1$	$g^2$	$g^3$
$i^0$	0.2	0.34	0.46
$i^1$	0.74	0.17	0.09

- SAT and Grade are independent given Intelligence (assumption)

- Knowing intelligence, SAT gives no information about class grade  $P \models (S \perp G \mid I)$ .



- Assertions

- From probabilistic reasoning  $P(I, S, G) = P(S, G \mid I)P(I)$

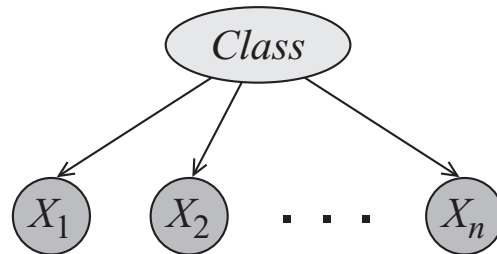
- From assumption  $P(S, G \mid I) = P(S \mid I)P(G \mid I)$ .

– **Combining**  $P(I, S, G) = P(S \mid I)P(G \mid I)P(I)$ .

$$\begin{aligned}
 P(i^1, s^1, g^2) &= P(i^1)P(s^1 \mid i^1)P(g^2 \mid i^1) \\
 &= 0.3 \cdot 0.8 \cdot 0.17 = 0.0408.
 \end{aligned}$$

Three binomials, two  
3-value multinomials:  
7 params  
More compact than  
joint distribution

# BN for General Naïve Bayes Model



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

Encoded using a very small number of parameters

Linear in the number of variables

# Application of Naiive Bayes Model

- Medical Diagnosis
  - Pathfinder expert system for lymph node disease (Heckerman et.al., 1992)
- Full BN agreed with human expert 50/53 cases
- Naiive Bayes agreed 47/53 cases

# Graphs and Distributions

- Relating two concepts:
  - Independencies in distributions
  - Independencies in graphs
- I-Map is a relationship between the two



# Independencies in a Distribution

- Let  $P$  be a distribution over  $X$
- $I(P)$  is set of conditional independence assertions of the form  $(X \perp Y | Z)$  that hold in  $P$

$X$	$Y$	$P(X,Y)$
$x^0$	$y^0$	0.08
$x^0$	$y^1$	0.32
$x^1$	$y^0$	0.12
$x^1$	$y^1$	0.48

$X$  and  $Y$  are independent in  $P$ , e.g.,

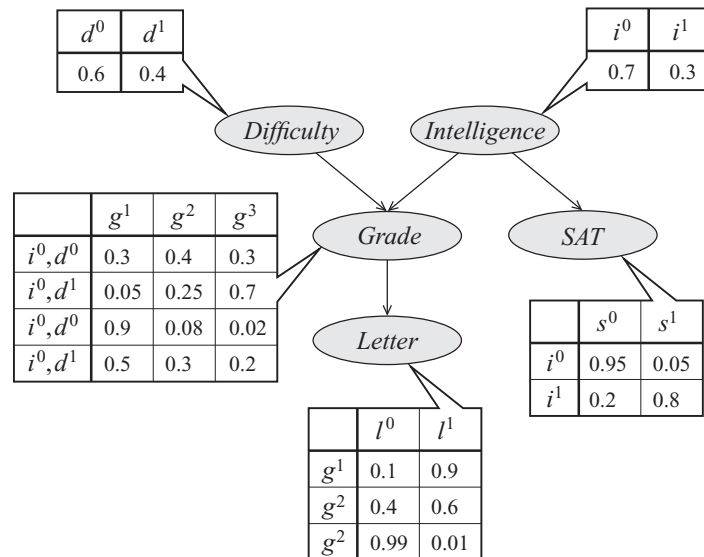
$$P(x^1) = 0.48 + 0.12 = 0.6$$

$$P(y^1) = 0.32 + 0.48 = 0.8$$

$$P(x^1, y^1) = 0.48 = 0.6 \times 0.8$$

Thus  $(X \perp Y | \phi) \in I(P)$

# Independencies in a Graph



- Graph with CPDs is equivalent to a set of independence assertions

$$P(D, I, G, S, L) = P(D)P(I)P(G | D, I)P(S | I)P(L | G)$$

- Local Conditional Independence Assertions** (starting from leaf nodes):

$$I(G) = \{ (L \perp I, D, S | G), \quad L \text{ is conditionally independent of all other nodes given parent } G$$

$$(S \perp D, G, L | I), \quad S \text{ is conditionally independent of all other nodes given parent } I$$

$$(G \perp S | D, I), \quad \text{Even given parents, } G \text{ is NOT independent of descendant } L$$

$$(I \perp D | \phi), \quad \text{Nodes with no parents are marginally independent}$$

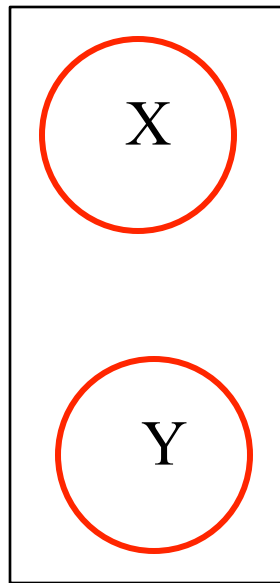
$$(D \perp I, S | \phi) \} \quad D \text{ is independent of non-descendants } I \text{ and } S$$

- Parents of a variable shield it from probabilistic influence
  - Once value of parents known, no influence of ancestors
- Information about descendants can change beliefs about a node

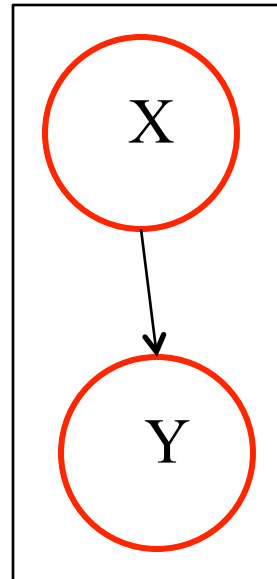
# I-MAP

- Let  $G$  be a graph associated with a set of independencies  $I(G)$
- Let  $P$  be a probability distribution with a set of independencies  $I(P)$
- Then  $G$  is an  $I$ -map of  $I$  if  $I(G) \subseteq I(P)$
- From direction of inclusion
  - distribution can have more independencies than the graph
  - Graph does not mislead in independencies existing in  $P$

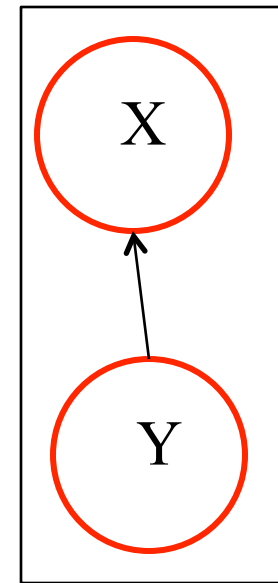
# Example of I-MAP



$G_0$  encodes  
 $X \perp Y$  or  
 $I(G_0) = \{X \perp Y\}$



$G_1$   
encodes no  
Independence  
or  
 $I(G_1) = \{\Phi\}$



$G_2$   
encodes no  
Independence  
 $I(G_2) = \{\Phi\}$

$X$	$Y$	$P(X,Y)$
$x^0$	$y^0$	0.08
$x^0$	$y^1$	0.32
$x^1$	$y^0$	0.12
$x^1$	$y^1$	0.48

$X$  and  $Y$  are independent  
in  $P$ , e.g.,

$G_0$  is an I-map of  $P$   
 $G_1$  is an I-map of  $P$   
 $G_2$  is an I-map of  $P$

$X$	$Y$	$P(X,Y)$
$x^0$	$y^0$	0.4
$x^0$	$y^1$	0.3
$x^1$	$y^0$	0.2
$x^1$	$y^1$	0.1

$X$  and  $Y$  are not  
independent in  $P$   
Thus  $(X \perp Y) \notin I(P)$

$G_0$  is not an I-map of  $P$   
 $G_1$  is an I-map of  $P$   
 $G_2$  is an I-map of  $P$

If  $G$  is an I-map of  $P$  then it captures some of the independences, not all

# I-map to Factorization

- A Bayesian network  $G$  encodes a set of conditional independence assumptions  $I(G)$
- Every distribution  $P$  for which  $G$  is an I-map should satisfy these assumptions
  - Every element of  $I(G)$  should be in  $I(P)$
- This is the key property to allowing a compact representation

# I-map to Factorization

- From chain rule of probability

$$P(I, D, G, L, S) = P(I)P(D|I)P(G|I, D)P(L|I, D, G)P(S|I, D, G, L)$$

- Relies on no assumptions
- Also not very helpful

- Last factor requires evaluation of 24 conditional probabilities

- Apply conditional independence assumptions induced from the graph

$D \perp I \in I(P)$  therefore  $P(D|I) = P(D)$

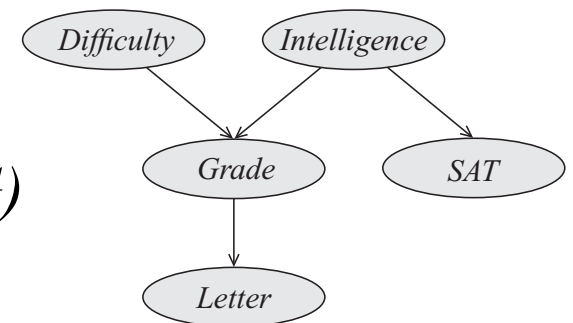
$(L \perp I, D) \in I(P)$  therefore  $P(L|I, D, G) = P(L|G)$

- Thus we get

$$P(D, I, G, S, L) = P(D)P(I)P(G|D, I)P(S|I)P(L|G)$$

- Which is a factorization into local probability models

- Thus we can go from graphs to factorization of  $P$



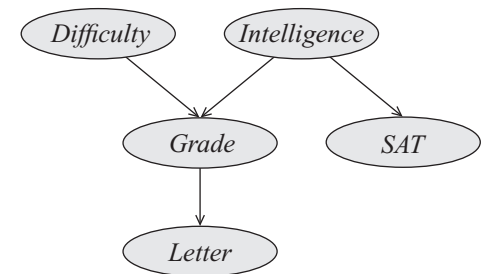
# Factorization to I-map

- We have seen that we can go from the independences encoded in  $G$ , i.e.,  $I(G)$ , to Factorization of  $P$
- Conversely, Factorization according to  $G$  implies associated conditional independences
  - If  $P$  factorizes according to  $G$  then  $G$  is an I-map for  $P$
  - Need to show that if  $P$  factorizes according to  $G$  then  $I(G)$  holds in  $P$
  - Proof by example

# Example that independences in $G$ hold in $P$

- $P$  is defined by set of CPDs
- Consider independences for  $S$  in  $G$ , i.e.,

$$P(S \perp\!\!\!\perp D, G, L | I)$$



- Starting from factorization induced by graph

$$P(D, I, G, S, L) = P(D)P(I)P(G | D, I)P(S | I)P(L | G)$$

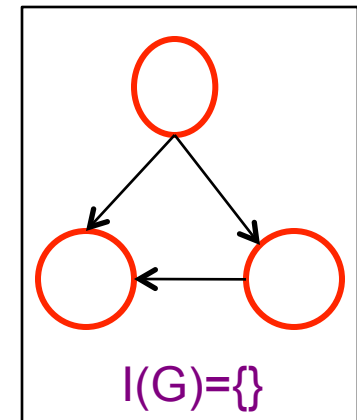
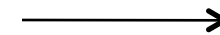
- Can show that  $P(S | I, D, G, L) = P(S | I)$
- Which is what we had assumed for  $P$



# Perfect Map

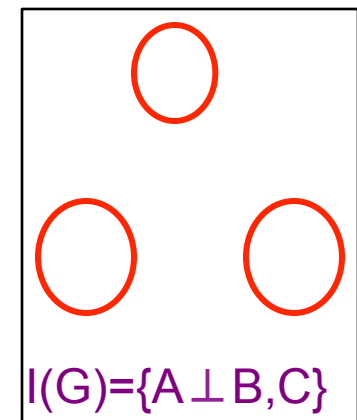
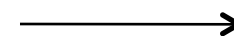
- I-map

- All independencies in  $I(G)$  present in  $I(P)$
- Trivial case: all nodes interconnected



- D-Map

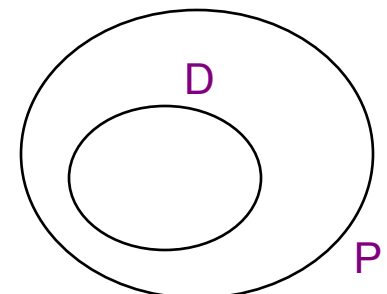
- All independencies in  $I(P)$  present in  $I(G)$
- Trivial case: all nodes disconnected



- Perfect map

- Both an I-map and a D-map
- Interestingly not all distributions  $P$  over a given set of variables can be represented as a perfect map

- Venn Diagram where  $D$  is set of distributions that can be represented as a perfect map



# Knowledge Engineering

- Going from given distribution to Bayesian network is more complex
- We have a vague model of the world
  - Need to crystallize it into network structure and parameters
- Task has several components
  - Each is subtle
  - Mistakes have consequences in quality of answers

# Three tasks in model building

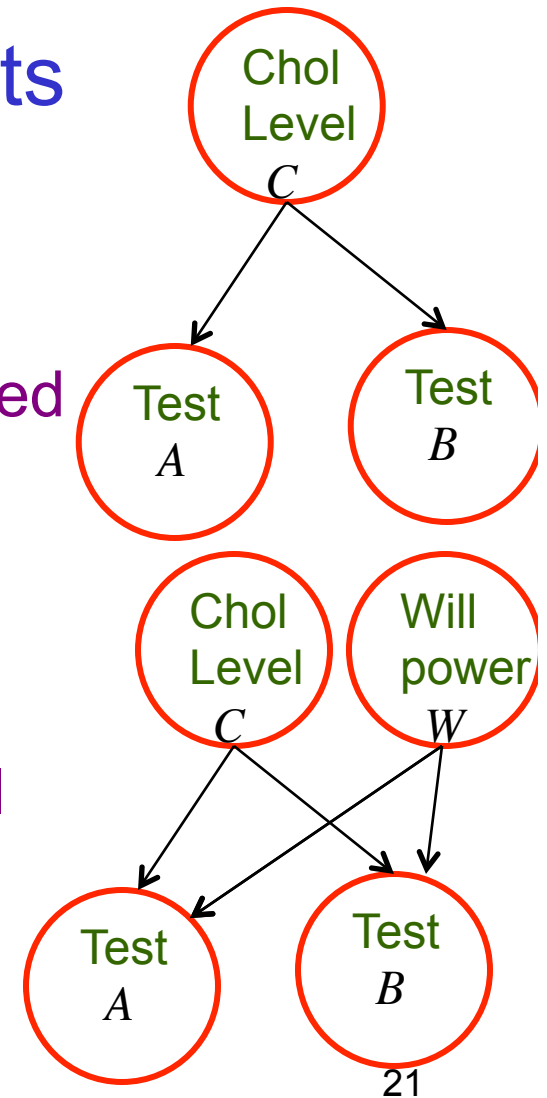
- All three tasks are hard:
  1. Picking variables
    - Many ways to pick entities and attributes
  2. Determining structure
    - Many structures hold
  3. Determining probabilities
    - Eliciting probabilities from people is hard

# 1. Picking Variables

- Model should contain variables
  - we can observe or that we will query
- Choosing variables is one of the hardest tasks
  - There are implications throughout the model
- Common problem: ill-defined variables
  - In medical domain: variable “Fever”
    - Temperature at time of admission?
    - Over prolonged period?
    - Thermometer or internal temperature?
  - Interaction of fever with other variables depend on specific interpretation

# Need for Hidden Variables

- There are several Cholestorol Tests
- For accurate answers:
  - Nothing to eat after 10:00pm
  - If person eats, all tests become correlated
- Hidden variable: willpower
  - Including it will render:
    - cholestorol tests conditionally independent given true cholestorol level and willpower
- Hidden variables: to avoid all variables being correlated



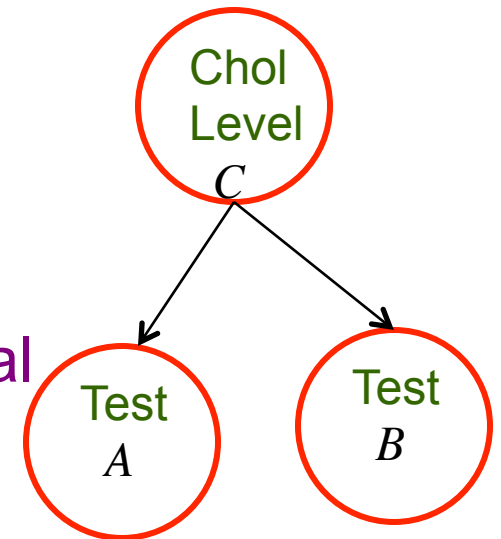
$$A \perp B | C, W$$

# Some variables not needed

- Not necessary to include every variable
  - SAT score may depend on partying previous night
  - Probability already accounts for poor score despite intelligence

# Picking Domain for Variables

- Reasonable domain of values to be chosen
- If partitions not fine enough conditional independence assumptions may be false
- Task of determining cholesterol level ( $C$ )
  - Two tests  $A$  and  $B$
  - $(A \perp B | C)$ 
    - $C$ : Normal if  $\leq 200$ , High if  $> 200$
    - Both tests fail if chol level has a marginal value( say 210)
      - Conditional independence assump. is false!
    - Introduce marginal value



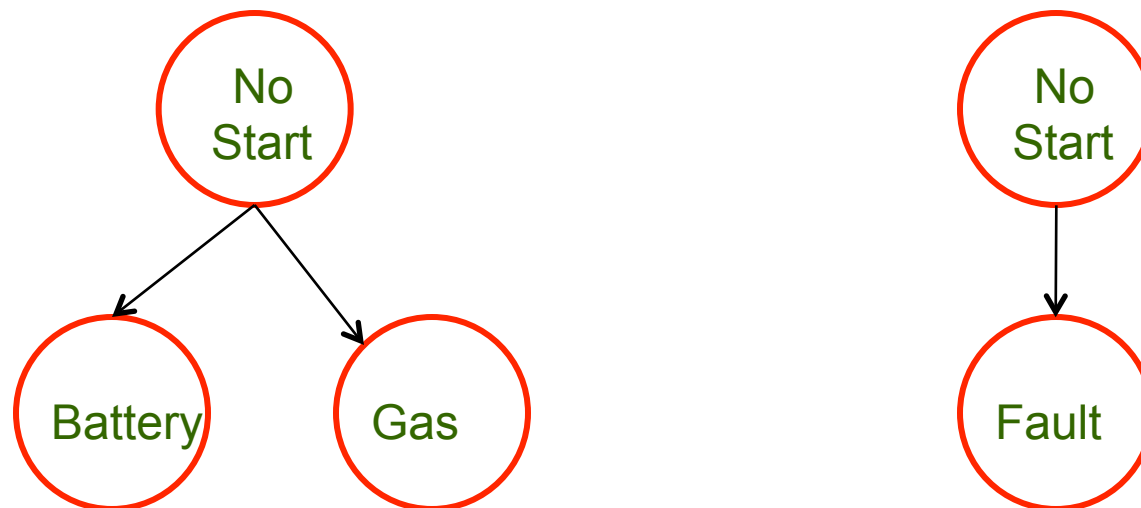
## 2. Picking Structure

- Many structures are consistent if we pick same set of independences
- Choose structure that reflects causal order and dependencies
  - Causes are parents of the effect
  - Causal graphs tend to be sparser
- Backward Construction Process
  - Lung cancer should have smoking as a parent
  - Smoking should have gender as a parent



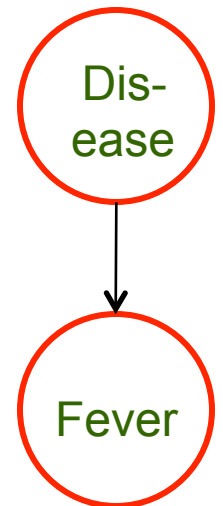
# Modeling weak influences

- Reasoning in a Bayesian network strongly depends on connectivity
- Adding edges can make it expensive to use
- Make approximations to decrease complexity



# 3. Picking Probabilities

- Zero Probabilities
  - Common mistake
    - Event extremely unlikely but not impossible
    - Can never condition away: irrecoverable errors
- Orders of Magnitude
  - Small diffs in low probs can make large differences in conclusions
    - $10^{-4}$  is very different from  $10^{-5}$
- Relative Values
  - Probability of fever higher with pneumonia than with flu



Disease \ Fever	High	Lo
Pneum	0.9	0.1
Flu	0.6	0.4

# Sensitivity Analysis

- Useful tool for estimating network parameters
- Determine extent to which a given probability parameter affects outcome
- Allows us to determine whether it is important to get a particular CPD entry right
- Helps figure out which CPD entries are responsible for an answer that does not match our intuition