

Reasoning Patterns and D-Separation

Sargur Srihari

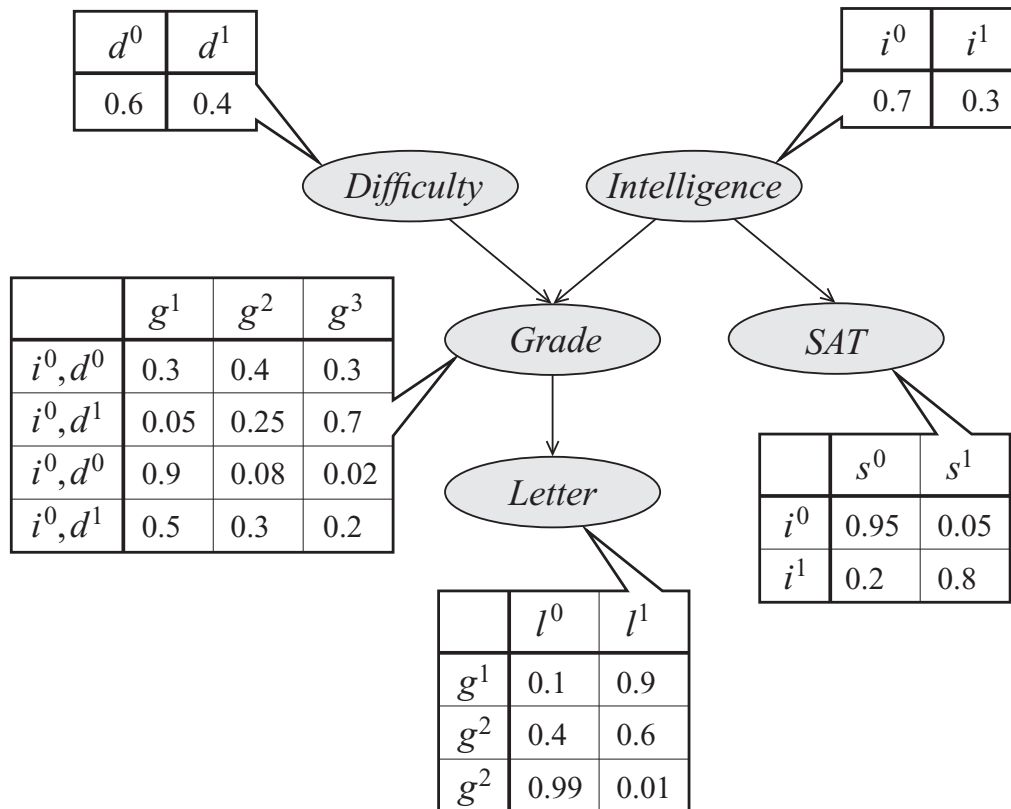
srihari@cedar.buffalo.edu

Topics

- Reasoning Patterns
 - Causal and Evidential Reasoning
- D-separation
 - Direct Connection
 - Indirect Connection
- Soundness and Completeness
- Algorithm for D-separation
- I-Equivalence
- Markov Blanket

Bayesian Network: Student Model

Graph and CPDs



$Val(I) = \{i^0 = \text{low intelligence}, i^1 = \text{high intelligence}\}$

$Val(D) = \{d^0 = \text{easy}, d^1 = \text{hard}\}$

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

$Val(S) = \{s^0 = \text{low}, s^1 = \text{high}\}$

$Val(L) = \{l^0 = \text{weak}, l^1 = \text{strong}\}$

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

$$P(i^1, d^0, g^2, s^1, l^0) = P(i^1)P(d^0)P(g^2 | i^1, d^0)P(s^1 | i^1)P(l^0 | g^2)$$

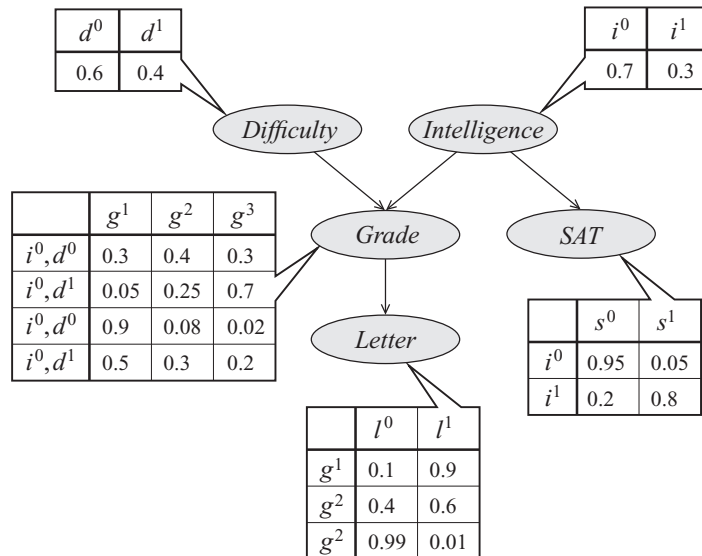
$$= 0.3 \cdot 0.6 \cdot 0.08 \cdot 0.8 \cdot 0.4 = 0.004608$$

**Chain rule for
Bayesian network₃**

Local Independencies in a BN

- A BN \mathcal{G} is a directed acyclic graph whose nodes represent random variables X_1, \dots, X_n .
- Let $Pa(X_i)$ denote parents of X_i in \mathcal{G}
- Let Non-descendants (X_i) denote variables in \mathcal{G} that are not descendants of X_i
- Then \mathcal{G} encodes the following set of conditional independence assumptions denoted $I_c(\mathcal{G})$
- For each X_i : $(X_i \perp\!\!\!\perp \text{NonDescendants } X_i \mid Pa(X_i))$
- Also known as *Local Markov Independencies*

Example of Local Independencies



- ◆ \mathcal{G} is equivalent to a set of *local* independence assertions

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

$$I_\ell(\mathcal{G}) = \{(L \perp I, D, S \mid G),$$

$$(S \perp D, G, L \mid I),$$

$$(G \perp S \mid D, I),$$

$$(I \perp D \mid \phi),$$

$$(D \perp I, S \mid \phi)\}$$

L is conditionally independent of all other nodes given parent G

S is conditionally independent of all other nodes given parent I

Even given parents, G is NOT independent of descendant L

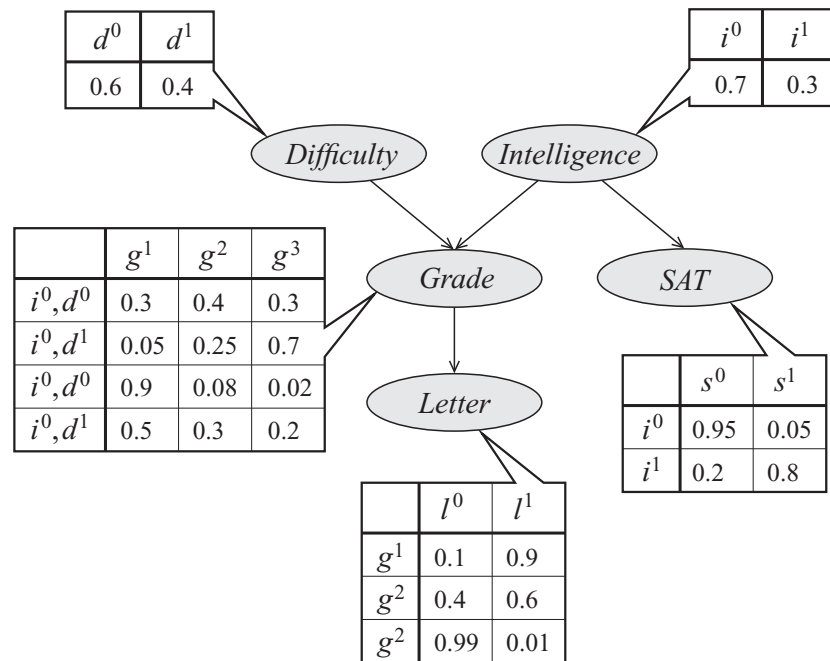
Nodes with no parents are marginally independent

D is independent of non-descendants I 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

Reasoning Patterns

Reasoning about a student George using the model



- Causal Reasoning

- George is interested in knowing as to how likely he is to get a strong *Letter* (based on *Intelligence*, *Difficulty*)?

- Evidential Reasoning

- Recruiter is interested in knowing whether George is *Intelligent* (based on *Letter*, *SAT*)

Causal Reasoning

1. How likely is George to get a strong *Letter* (knowing nothing else)?

- $P(l^1) = 0.502$
- Obtained by summing-out other variables in joint distribution

2. But George is not so *Intelligent* (i^0)

- $P(l^1 | i^0) = 0.389$

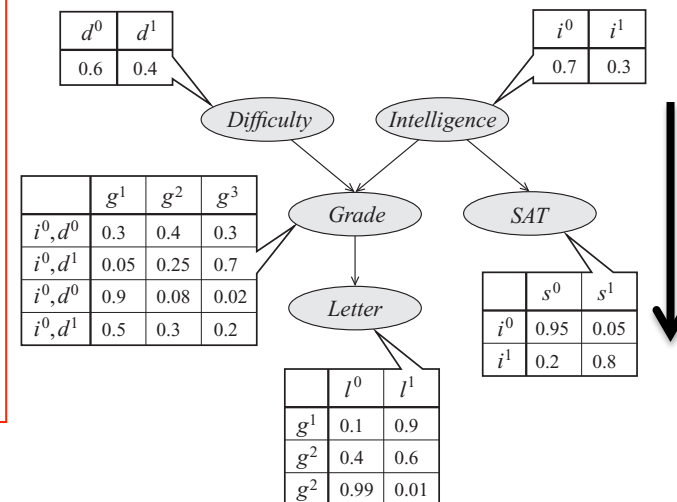
3. Next we find ECON101 is not *Difficult* (d^0)

- $P(l^1 | i^0, d^0) = 0.513$

Observe how probabilities change as evidence is obtained

$$P(D, I, G, S, l^1) =$$

$$\sum_{D, I, G, S} P(D)P(I)P(G|D, I)P(S|I)P(l^1|G)$$

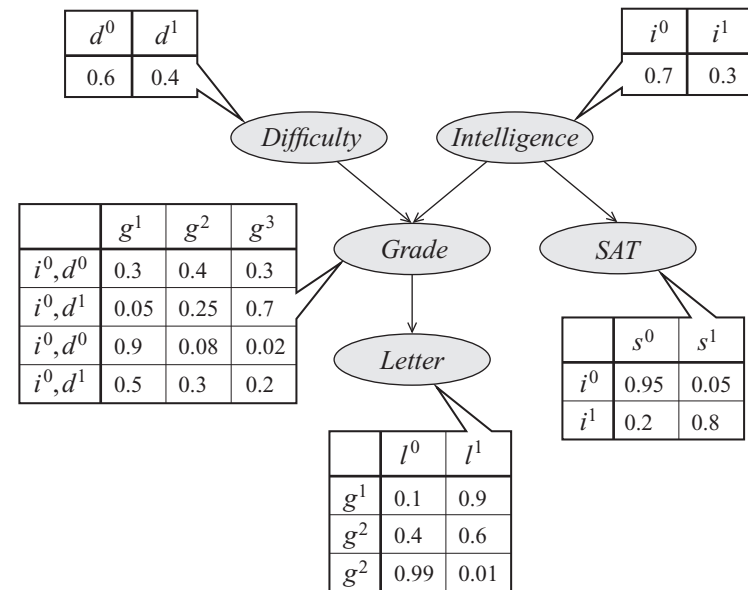


Query is Example of Causal Reasoning:

Predicting downstream effects of factors such as intelligence

Evidential Reasoning

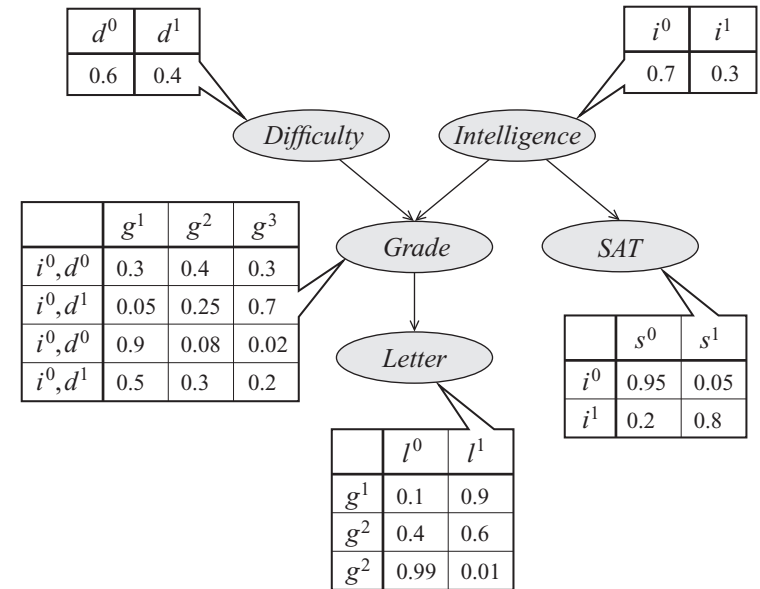
- Recruiter wants to hire *Intelligent* student
- A priori George is 30% likely to be *Intelligent*
- $$P(i^1) = 0.3$$
- Finds that George received *Grade C* (g^3) in ECON101
- $$P(i^1 | g^3) = 0.079$$
- Similarly probability class is *Difficult* goes up from 0.4 to
- $$P(d^1 | g^3) = 0.629$$
- If recruiter has lost *Grade* but has *Letter*
- $$P(i^1 | l^0) = 0.14$$



- Recruiter has both *Grade* and *Letter*
- $$P(i^1 | l^0, g^3) = 0.079$$
 - Same as if he had only *Grade*
 - Letter* is immaterial
- Reasoning from effects to causes is called evidential reasoning

Intercausal reasoning

- Recruiter has grade (letter does not matter)
- $P(i^1|g^3)=P(i^1|l^0,g^3)=0.079$
- Recruiter receives high SAT score (leads to dramatic increase)
- $P(i^1|g^3,s^1)=0.578$
- Intuition:
 - High SAT score outweighs poor grade since low intelligence rarely gets good SAT scores
 - Smart students more likely to get Cs in hard classes
- Probability of class is difficult also goes up from
- $P(d^1|g^3)=0.629$ to
- $P(d^1|g^3,s^1)=0.76$



Information about SAT score gave us information about Intelligence which with Grade told us about difficulty of course

One causal factor for *Grade* (*Intelligence*) gives us information about another (*Difficulty*)

Explaining Away

An example:

Given $grade=C$ $letter=weak$

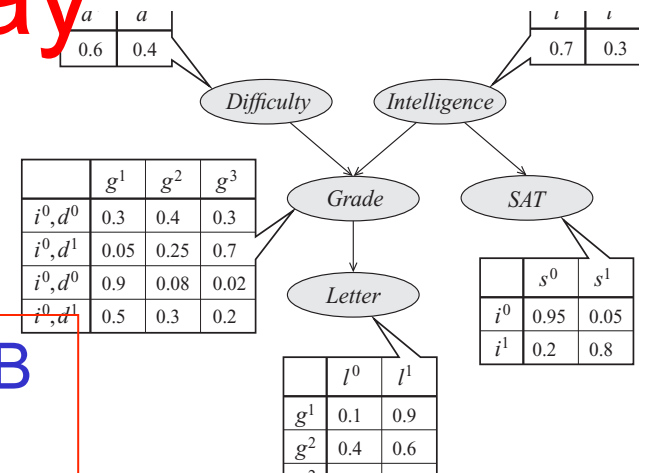
- $P(i^1 | l^0, g^3) = 0.079$

If we observe ECON101 is hard $difficulty=high$

- $P(i^1 | g^3, d^1) = 0.11$
- We have provided partial explanation for George's performance in ECON101

Another example:

- If George gets a B in ECON101
- $P(i^1 | g^2) = 0.175$
- If we observe ECON101 is hard
- $P(i^1 | g^2, d^1) = 0.34$
- We have explained away the poor grade via the difficulty of the class

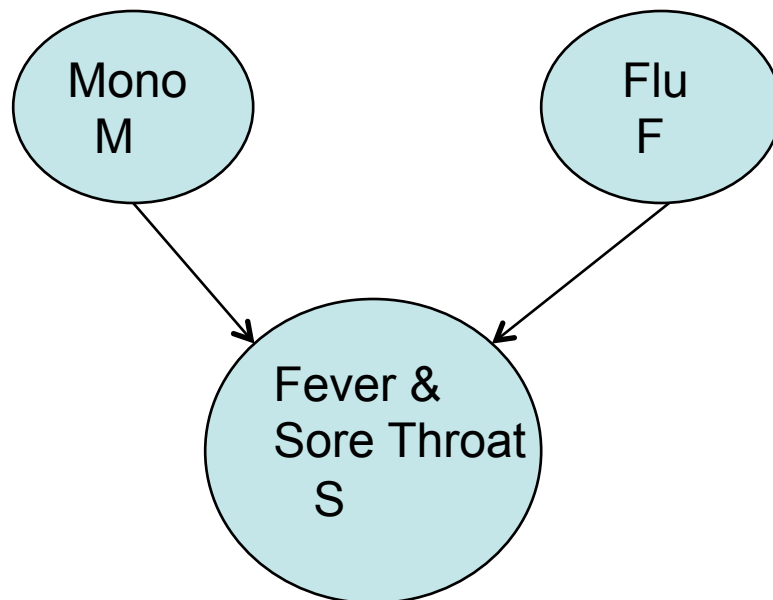


Explaining away is one type of intercausal reasoning

- Different causes of the same effect can interact
- All determined by probability calculation rather than heuristics

Intercausal Reasoning is Common in Human Reasoning

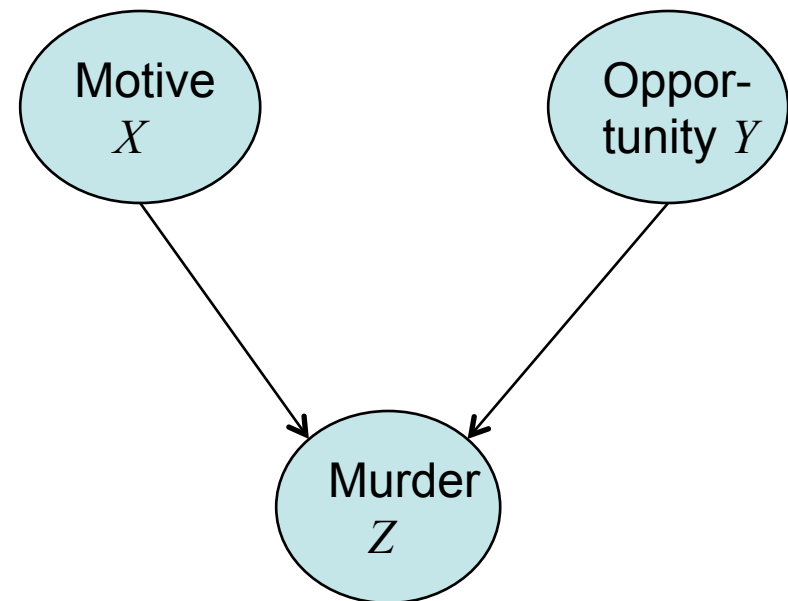
Another example of explaining away



- Binary Variables
- Fever & Sore Throat can be caused by mono and flu
- When flu is diagnosed probability of mono is reduced (although mono could still be present)
- It provides an alternative explanation of symptoms
- $P(m^1|s^1) > P(m^1|s^1, f^1)$

Another Type of Intercausal Reasoning

- Binary Variables
 - Murder (leaf node)
 - Motive and Opportunity are causal nodes
- Binary Variables X, Y, Z
- X and Y both increase the probability of Murder
 - $P(z^1|x^1) > P(z^1)$
 - $P(z^1|y^1) > P(z^1)$
- Each of X and Y increase probability of other
 - $P(x^1|z^1) < P(x^1|y^1, z^1)$
 - $P(y^1|z^1) < P(y^1|x^1, z^1)$



**Can go in any direction
Different from Explaining
Away**

Independencies in Graphs

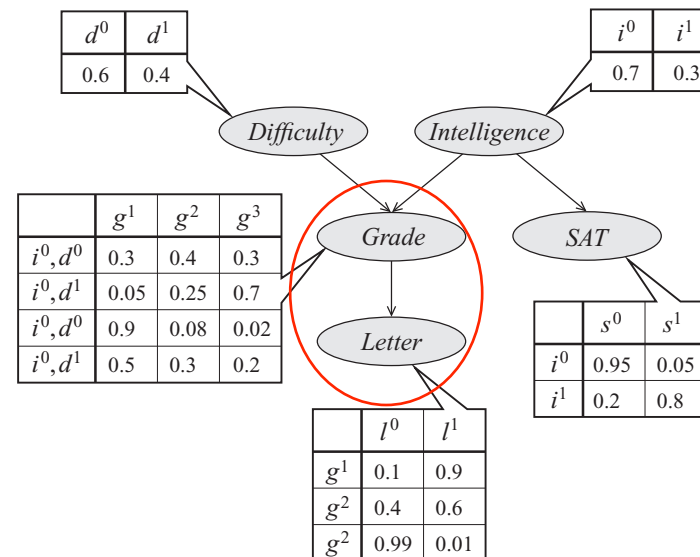
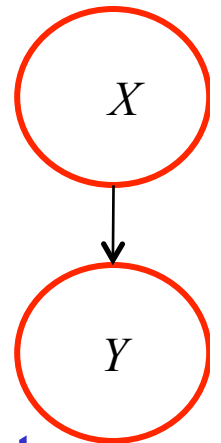
- A graph structure \mathcal{G} encodes a set of conditional independence assumptions $I(\mathcal{G})$
- Are there other independencies that we can read-off?
 - i.e., are there other independencies that hold for every distribution that factorizes over \mathcal{G} ?
- D-separation holds the key

Dependencies and Independencies

- Crucial for understanding network behavior
- Independence properties are important for answering queries
 - Exploited to reduce computation of inference
 - A distribution P that factorizes over G satisfies $I(G)$
 - Are there other independencies that can be *read off* directly from G ?
 - That hold for every P that factorizes over G

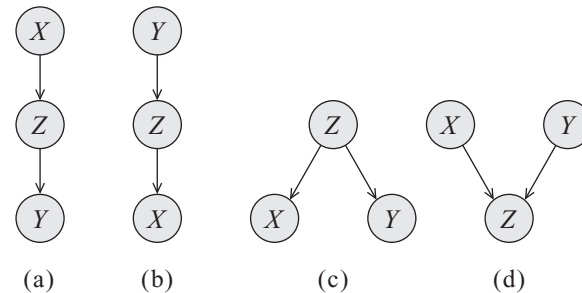
Direct Connection between X and Y

- X and Y are correlated regardless of any evidence about any other variables
 - E.g., Feature Y and character X are correlated
 - Grade G and Letter L are correlated
- If X and Y are directly connected we can get examples where they influence each other regardless of Z



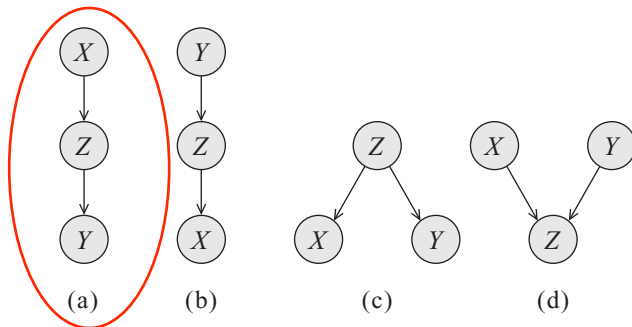
Indirect Connection betwn X and Y

- Four cases where X and Y are connected via Z
 - Indirect causal effect
 - Indirect evidential effect
 - Common cause
 - Common effect
- We will see that first three cases are similar while fourth case (V -structure) is different



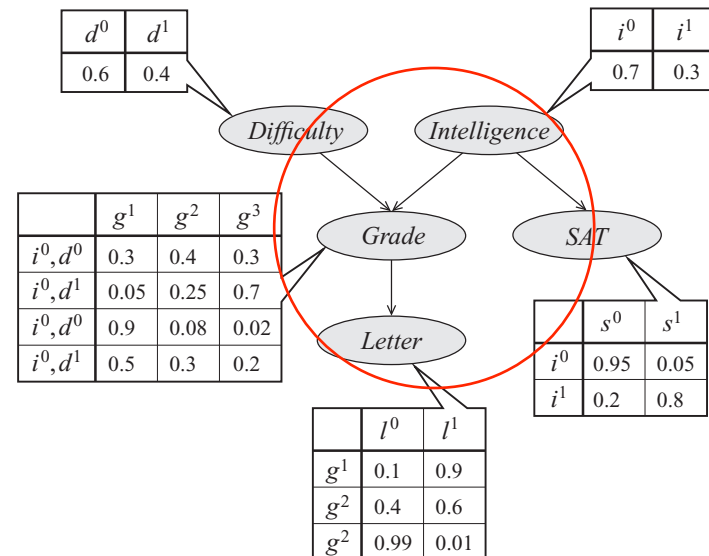
1. Indirect Causal Effect: $X \rightarrow Z \rightarrow Y$

- Cause X cannot influence effect Y if Z observed
 - Observed Z blocks influence
- If $Grade$ observed then I does not influence L
 - $Intelligence$ influences $Letter$ if $Grade$ is unobserved



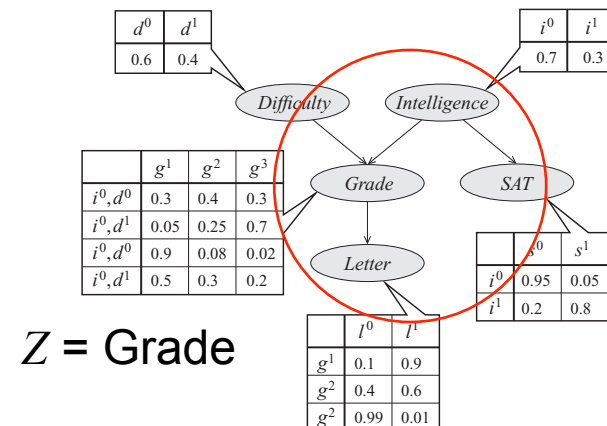
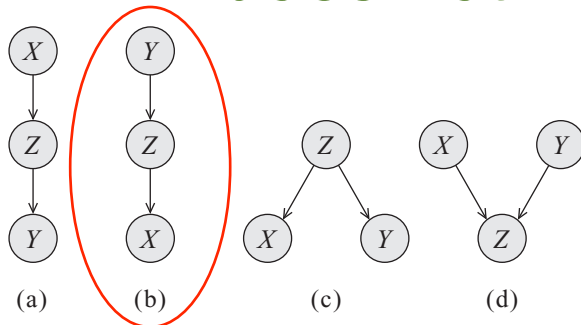
$Z = Grade$

$I \perp L | G$



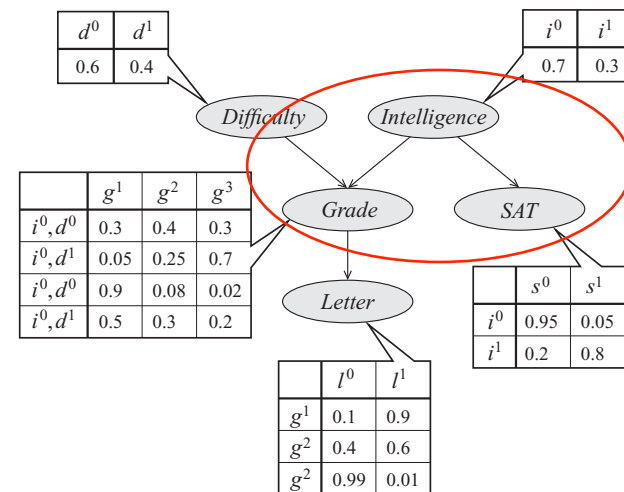
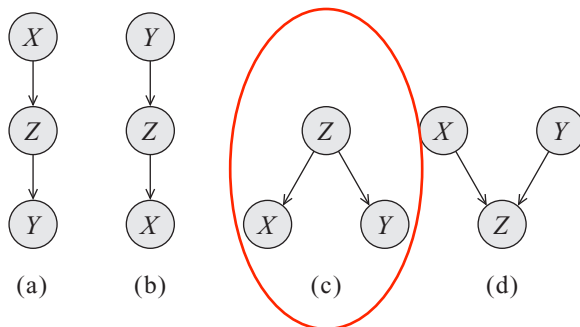
2. Indirect Evidential Effect: $Y \rightarrow Z \rightarrow X$

- Evidence X can influence Y via Z only if Z is unobserved
 - Observed Z blocks influence
- If Grade unobserved, Letter influences assessment of Intelligence
- Dependency is a symmetric notion
 - $X \perp Y$ does not hold then $Y \perp X$ does not hold either



3. Common Cause: $X \leftarrow Z \rightarrow Y$

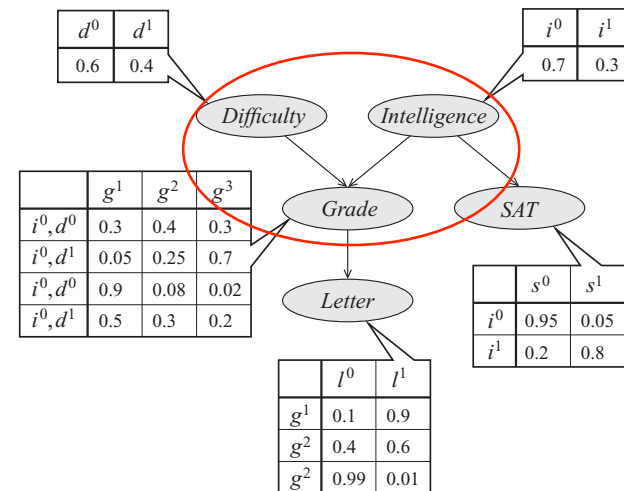
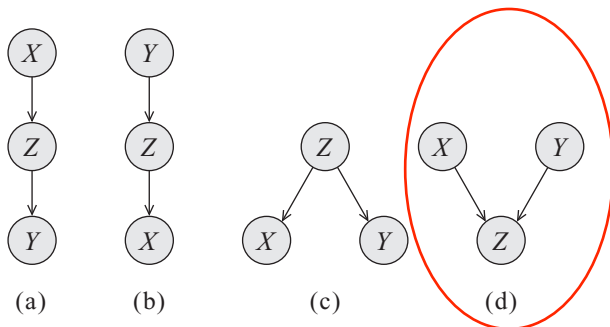
- X can influence Y if and only if Z is not observed
 - Observed Z blocks
- Grade is correlated with SAT score
- But if Intelligence is observed then SAT provides no additional information



S ind G||

4. Common Effect (V-structure) $X \rightarrow Z \leftarrow Y$

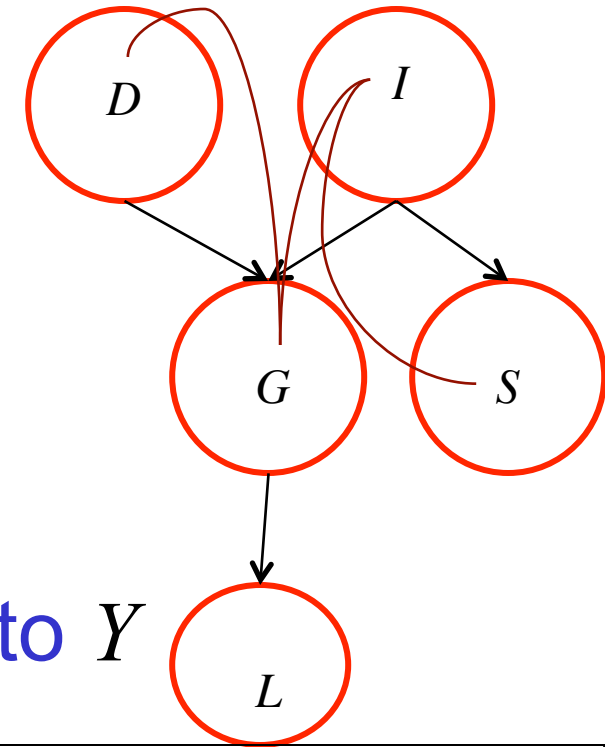
- Influence cannot flow on trail $X \rightarrow Z \leftarrow Y$ if Z is not observed
 - Observed Z enables
 - Opposite to previous 3 cases (Observed Z blocks)
- When G not observed I and D are independent
- When G is observed, I and D are correlated



$$I \perp D | \sim G$$

Active Trail

- *Grade* is not observed
- Observe weak letter L
 - Which indicates low *Grade*
 - Suffices to correlate D and I
- When influence can flow from X to Y via Z then trail $X—Z—Y$ is active
- Summary



Consider Trail $D \rightarrow G \leftarrow I \rightarrow S$
 When observed $Z = \{\phi\}$: trail is *inactive* because v-structure $D \rightarrow G \leftarrow I$ is inactive

$Z = \{L\}$: *active* ($D \rightarrow G \leftarrow I$ active) since L is descendant of G

$Z = \{L, I\}$: *inactive* because observing I blocks $G \leftarrow I \rightarrow S$.

Causal trail: $X \rightarrow Z \rightarrow Y$: active iff Z not observed

Evidential Trail: $X \leftarrow Z \leftarrow Y$: active iff Z is not observed

Common Cause: $X \leftarrow Z \rightarrow Y$: active iff Z is not observed

Common Effect: $X \rightarrow Z \leftarrow Y$: active iff either Z or one of its descendants is observed

D-separation definition

- Let X, Y and Z be three sets of nodes in \mathcal{G} .
- X and Y are d-separated given Z denoted $d\text{-sep}_{\mathcal{G}}(X:Y|Z)$ if there is no active trail between any node $X \in X$ and $Y \in Y$ given Z
- That is, nodes in X cannot influence nodes in Y
- Provides notion of separation between nodes in a directed graph (“directed” separation)

Independencies from D-separation

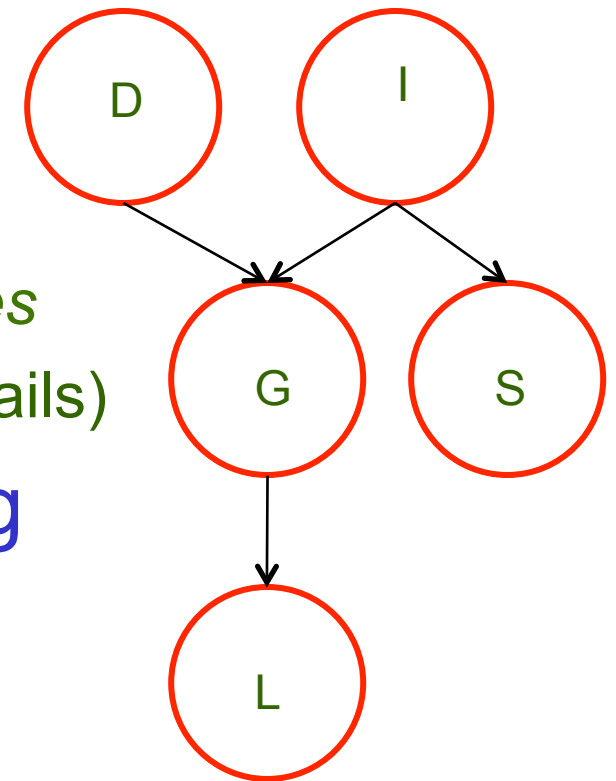
- **Definition:**

$$I(\mathcal{G}) = \{(X \perp Y | Z) : d\text{-sep}_{\mathcal{G}}(X:Y|Z)\}$$

- Also called *Global Markov independencies*
- Note: Derived purely from graph (using trails)

- **Ex: *Global* independence using D-separation**

$$I(\mathcal{G}) = \{(L \perp I, D, S | G)\}$$



- **Compare with *local* independencies**

- For each X_i : $(X_i \perp \text{NonDescendants } X_i | \text{Pa}(X_i))$

$$I_{\ell}(\mathcal{G}) = \{(L \perp I, D, S | G), (S \perp D, G, L | I), (G \perp S | D, I), (D \perp I, S | \phi)\}$$

Soundness and Completeness

- Formalizing notion of D-separation
- Soundness Theorem
 - If a distribution P factorizes according to \mathcal{G} then $I(\mathcal{G}) \subseteq I(P)$
- A distribution P is faithful to graph \mathcal{G} if any independence in P is reflected in \mathcal{G}
 - \mathcal{G} is then called a Perfect Map
- Completeness Theorem
 - Definition of $I(\mathcal{G})$ is the maximal one
- Thus D-separation test precisely characterizes independencies that hold for P

Algorithm for D-Separation

- Enumerate all trails between X & Y is inefficient
 - No. of trails is *exponential* with graph size
- Linear time algorithm has two phases
 - Algorithm $Reachable(G, X, Z)$ returns nodes reachable from X given Z
 - Phase 1 (simple)
 - Traverse bottom-up from leaves marking all nodes in Z or descendants in Z ; to enable v-structures
 - Phase 2 (subtle)
 - Traverse top-down from X to Y stopping when blocked by a node

Algorithm for finding nodes reachable from X given Z via active trails

```

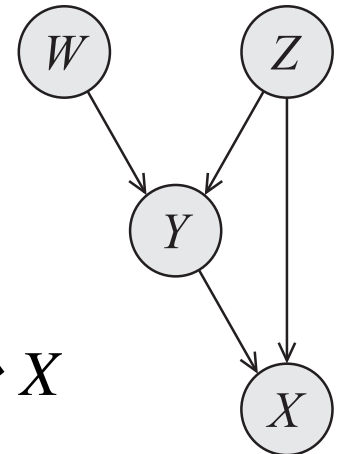
Procedure Reachable (
   $\mathcal{G}$ , // Bayesian network graph
   $X$ , // Source variable
   $Z$  // Observations
)
1 // Phase I: Insert all ancestors of  $Z$  into  $V$ 
2  $L \leftarrow Z$  // Nodes to be visited
3  $A \leftarrow \emptyset$  // Ancestors of  $Z$ 
4 while  $L \neq \emptyset$ 
5   Select some  $Y$  from  $L$ 
6    $L \leftarrow L - \{Y\}$ 
7   if  $Y \notin A$  then
8      $L \leftarrow L \cup \text{Pa}_Y$  //  $Y$ 's parents need to be visited
9      $A \leftarrow A \cup \{Y\}$  //  $Y$  is ancestor of evidence
10
11 // Phase II: traverse active trails starting from  $X$ 
12  $L \leftarrow \{(X, \uparrow)\}$  // (Node,direction) to be visited
13  $V \leftarrow \emptyset$  // (Node,direction) marked as visited
14  $R \leftarrow \emptyset$  // Nodes reachable via active trail
15 while  $L \neq \emptyset$ 
16   Select some  $(Y, d)$  from  $L$ 
17    $L \leftarrow L - \{(Y, d)\}$ 
18   if  $(Y, d) \notin V$  then
19     if  $Y \notin Z$  then
20        $R \leftarrow R \cup \{Y\}$  //  $Y$  is reachable
21        $V \leftarrow V \cup \{(Y, d)\}$  // Mark  $(Y, d)$  as visited
22       if  $d = \uparrow$  and  $Y \notin Z$  then // Trail up through  $Y$  active if  $Y$  not in  $Z$ 
23         for each  $Z \in \text{Pa}_Y$ 
24            $L \leftarrow L \cup \{(Z, \uparrow)\}$  //  $Y$ 's parents to be visited from bottom
25         for each  $Z \in \text{Ch}_Y$ 
26            $L \leftarrow L \cup \{(Z, \downarrow)\}$  //  $Y$ 's children to be visited from top
27       else if  $d = \downarrow$  then // Trails down through  $Y$ 
28         if  $Y \notin Z$  then
29           // Downward trails to  $Y$ 's children are active
30           for each  $Z \in \text{Ch}_Y$ 
31              $L \leftarrow L \cup \{(Z, \downarrow)\}$  //  $Y$ 's children to be visited from top
32         if  $Y \in A$  then // v-structure trails are active
33           for each  $Z \in \text{Pa}_Y$ 
34              $L \leftarrow L \cup \{(Z, \uparrow)\}$  //  $Y$ 's parents to be visited from bottom
35 return  $R$ 

```

Example for D-separation algorithm

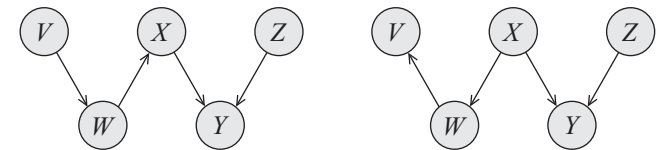
- Task: Find all nodes reachable from X
- Assume that Y is observed, i.e., $Y \in \mathbf{Z}$
 - Assume algorithm first encounters Y via edge $Y \rightarrow X$
 - Any extension of this trail is blocked by Y
 - Trail $X \leftarrow Z \rightarrow Y \leftarrow W$ is not blocked by Y

Because of v-structure
 - Thus when we encounter Y for the second time via the edge $Z \rightarrow Y$ we should not ignore it
 - Therefore after the first visit to Y we can mark it as visited
 - For trails coming from children of Y
 - Not for purpose of trails coming from parents of Y



I-Equivalence

- Conditional independence assertion statements can be the same with different structures
- Two graphs \mathcal{K}_1 and \mathcal{K}_2 are *I-equivalent* if $I(\mathcal{K}_1) = I(\mathcal{K}_2)$
- Skeleton of a BN graph \mathcal{G} is an undirected graph with an edge for every edge in \mathcal{G}
- If two BN graphs have the same set of skeletons and v-structures then they are *I-equivalent*



Same skeleton
Same v-structure $X \rightarrow Y \leftarrow Z$

Markov Blanket

- Joint distribution $p(\mathbf{x}_1, \dots, \mathbf{x}_D)$ represented by directed graph with D nodes
- Distribution of \mathbf{x}_i conditioned on rest is

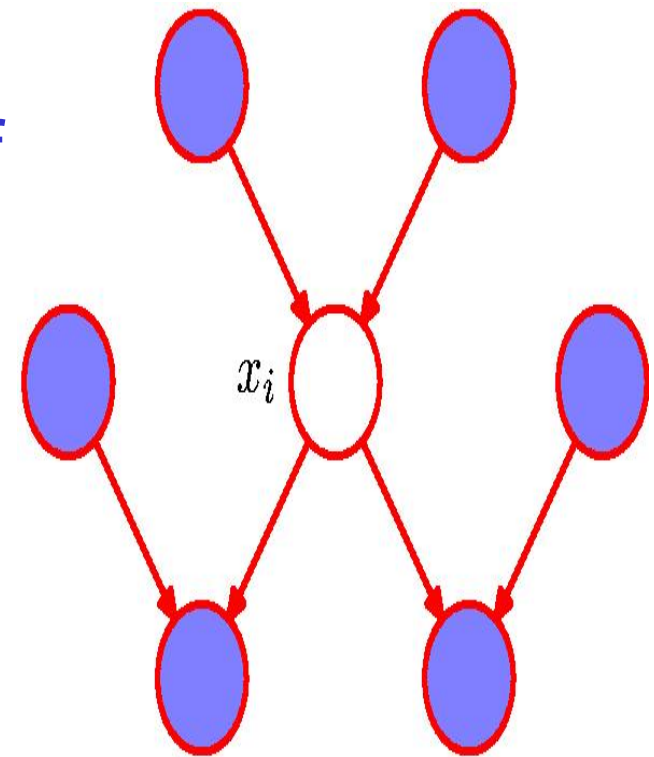
$$p(\mathbf{x}_i | \mathbf{x}_{\{j \neq i\}}) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_D)}{\int p(\mathbf{x}_1, \dots, \mathbf{x}_D) d\mathbf{x}_i} \quad \text{from product rule}$$

$$= \frac{\prod_k p(\mathbf{x}_k | \text{pa}_k)}{\int \prod_k p(\mathbf{x}_k | \text{pa}_k) d\mathbf{x}_i} \quad \text{from factorization property}$$

- Factors not depending on \mathbf{x}_i can be taken outside integral over \mathbf{x}_i and cancel between numerator and denominator
- Remaining will depend only on parents, children and *co-parents*

Markov Blanket of a Node

- Set of parents, children and co-parents of a node
- The conditional distribution of x_i conditioned on all the remaining variables in the graph is dependent on only the variables in the Markov blanket
- Also called *Markov boundary*
- Can be used to efficiently implement Gibbs sampling



Gibbs Sampling Algorithm

- **Initialize** $\{z_i, i=1, \dots, M\}$, $M = \text{no of variables}$
- **For** $t=1, \dots, T$, $T = \text{no of samples}$
 - **Sample** $z_1^{(t+1)} \sim p(z_1 | z_2^{(t)}, z_3^{(t)}, \dots, z_M^{(t)})$
 - **Sample** $z_2^{(t+1)} \sim p(z_2 | z_1^{(t+1)}, z_3^{(t)}, \dots, z_M^{(t)})$
 -
 - **Sample** $z_j^{(t+1)} \sim p(z_j | z_1^{(t+1)}, \dots, z_{j-1}^{(t+1)}, z_{j+1}^{(t)}, \dots, z_M^{(t)})$
 -
 - **Sample** $z_M^{(t+1)} \sim p(z_M | z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{M-1}^{(t+1)})$