

# Bayes Nets

**Russell and Norvig: Chapter 12, 13**

**CSE 240: Winter 2023**

**Lecture 13**

# Announcements

- Quiz 2 grading in progress
- Assignment 3 is due tonight at 5pm
- Working on regrades

# Agenda and Topics

- Probability
  - Independence and conditional independence
- Naives Bayes
- Bayesian Networks

# Conditional Independence

- Unfortunately, random variables of interest are rarely independent of each other
- A more suitable notion is that of **conditional independence**
- Two variables  $X$  and  $Y$  are **conditionally independent** given  $Z$  if
  - $P(x|y,z) = P(x|z)$  for all values  $x,y,z$
  - That is, learning the values of  $Y$  does not change prediction of  $X$  once we know the value of  $Z$
  - Equivalently,  $P(x,y|z) = P(x|z)P(y|z)$  for all values  $x,y,z$
- notation:  $I(X; Y | Z)$ ,  $X \perp\!\!\!\perp Y | Z$

# The Chain Rule

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots$$

- Trivial decomposition:

$$P(\text{Traffic}, \text{Rain}, \text{Umbrella}) = \\ P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain}, \text{Traffic})$$

- With assumption of conditional independence:

$$P(\text{Traffic}, \text{Rain}, \text{Umbrella}) = \\ P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain})$$

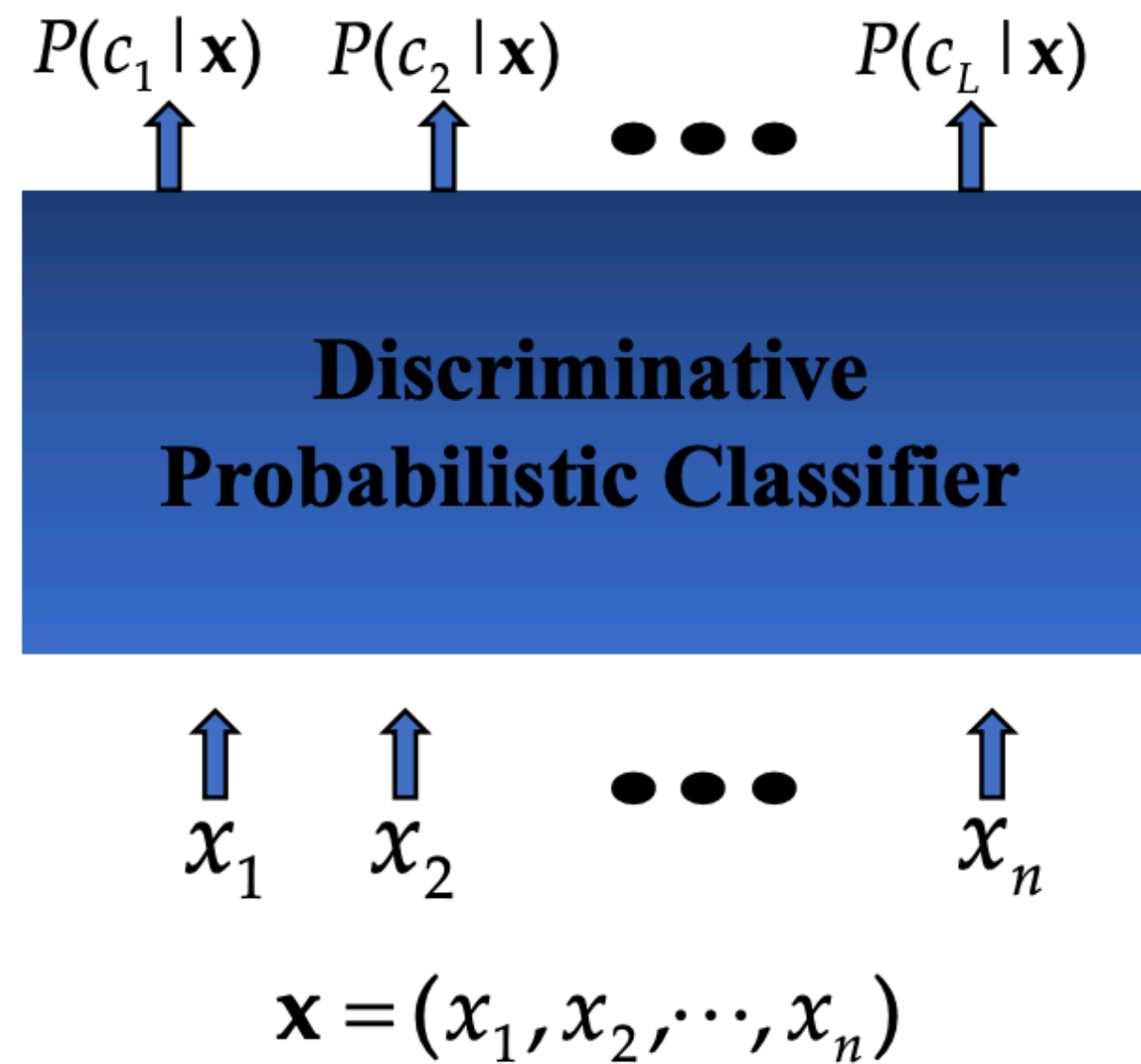
- Bayes' nets / graphical models help us express conditional independence assumptions

# Naïve Bayes

# Probabilistic Classification Principle

Establishing a probabilistic model for classification (**Discriminative model**)

$$P(c | \mathbf{x}) \quad c = c_1, \dots, c_L, \mathbf{x} = (x_1, \dots, x_n)$$

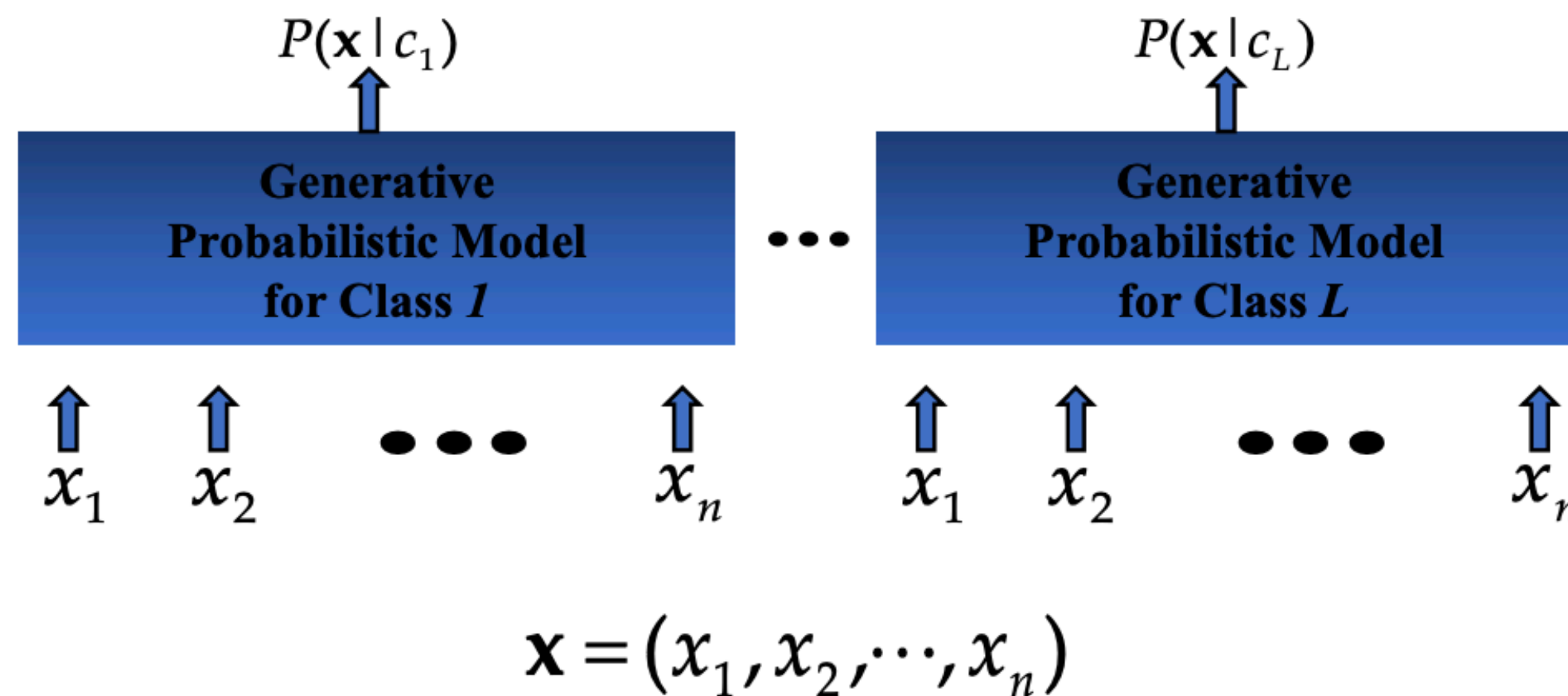


Output  $L$  probabilities for  $L$  class labels in a probabilistic classifier

# Probabilistic Classification Principle

Establishing a probabilistic model for classification (**Generative model**)

$$P(\mathbf{x} | c) \quad c = c_1, \dots, c_L, \mathbf{x} = (x_1, \dots, x_n)$$



- $L$  probabilistic models have to be trained **independently**
- Each is trained on only the examples of the same label
- Output  $L$  probabilities for a given input with  $L$  models



# Probabilistic Classification Principle

- Maximum A Posterior (MAP) classification rule
  - For an input  $x$ , find the largest one from  $L$  probabilities output by a discriminative probabilistic classifier  $P(c_1 | x), \dots, P(c_L | x)$ .
  - Assign  $x$  to label  $c^*$  if  $P(c^* | x)$  is the largest
- Generative classification with MAP rule
  - Apply Bayesian rule to convert them into posterior probabilities
  - $$P(c_i | x) = \frac{P(x | c_i)P(c_i)}{P(x)} \propto P(x | c_i)p(c_i)$$
    - (for  $i = 1, 2, \dots, L$ )
  - Then apply the MAP rule to assign a label

# Naïve Bayes

- Bayes classification
  - $P(c | x) \propto P(x | c)P(c) = P(x_1, \dots, x_n | c)P(c)$ 
    - (for  $c = c_1, \dots, c_L$ )
  - Difficulty: learning the joint probability  $P(x_1, \dots, x_n | c)$  is often infeasible!
- Naïve Bayes classification

Assume all input features are conditionally independent!

$$\begin{aligned} P(x_1, \dots, x_n | c) &= P(x_1 | x_2, \dots, x_n, c)P(x_2, \dots, x_n | c) \\ &= P(x_1 | c)P(x_2, \dots, x_n | c) \\ &= P(x_1 | c)P(x_2 | c) \dots P(x_n | c) \end{aligned}$$

$$\left[ P(a_1 | c^*) \dots P(a_n | c^*) \right] P(c^*) > \left[ P(a_1 | c) \dots P(a_n | c) \right] P(c), \quad c \neq c^*, c = c_1, \dots, c_L$$

# The Naïve Bayes Model

- The *Naïve Bayes Assumption*:
  - Assume that all features are independent **given the class label Y**.
- Equationally speaking:

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

# Why is This Useful?

- # of parameters for modeling  $P(X_1, \dots, X_n | Y)$ :
  - $2(2^n - 1)$
- # of parameters for modeling  $P(X_1 | Y), \dots, P(X_n | Y)$ :
  - $2n$

# Naïve Bayes Training

- Training in Naïve Bayes is **easy**:
  - Estimate  $P(Y = v)$  as the fraction of records with  $Y = v$

$$P(Y = v) = \frac{\text{Count}(Y = v)}{\#records}$$

- Estimate  $P(X_i = u \mid Y = v)$  as the fraction of records with  $Y = v$  for which  $X_i = u$ .

$$P(X_i \mid Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v)}{\text{Count}(Y = v)}$$

# Naïve Bayes Training

- In practice, some of these counts can be zero
- Fix this by adding “virtual” counts

$$P(X_i | Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v) + 1}{\text{Count}(Y = v) + 2}$$

- This is called Smoothing

# Example

The weather data, with counts and probabilities													
outlook			temperature			humidity			windy			play	
		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								
sunny	2/9	3/5	hot	2/9	2/5	high	3/9	4/5	false	6/9	2/5	9/14	5/14
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								
A new day													
outlook			temperature			humidity			windy			play	
sunny			cool			high			true			?	

# Example

- Learning phase

Outlook	Play=Yes	Play=No
<i>Sunny</i>	2/9	3/5
<i>Overcast</i>	4/9	0/5
<i>Rain</i>	3/9	2/5

Temperature	Play=Yes	Play=No
<i>Hot</i>	2/9	2/5
<i>Mild</i>	4/9	2/5
<i>Cool</i>	3/9	1/5

Humidity	Play=Yes	Play=No
<i>High</i>	3/9	4/5
<i>Normal</i>	6/9	1/5

Wind	Play=Yes	Play=No
<i>Strong</i>	3/9	3/5
<i>Weak</i>	6/9	2/5

$$P(\text{Play=Yes}) = 9/14$$

$$P(\text{Play=No}) = 5/14$$



# Solution

- Test Phase

- Given a new instance, predict its label

$\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$

- Look up tables achieved in the learning phrase

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{Yes}) = 2/9$$

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{No}) = 1/5$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{No}) = 4/5$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Play}=\textit{Yes}) = 9/14$$

$$P(\text{Play}=\textit{No}) = 5/14$$

$$P(\text{Yes} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{Yes})P(\textit{Cool} \mid \textit{Yes})P(\textit{High} \mid \textit{Yes})P(\textit{Strong} \mid \textit{Yes})]P(\text{Play}=\textit{Yes}) = 0.0053$$

$$P(\text{No} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{No})P(\textit{Cool} \mid \textit{No})P(\textit{High} \mid \textit{No})P(\textit{Strong} \mid \textit{No})]P(\text{Play}=\textit{No}) = 0.0206$$

Given the fact  $P(\text{Yes} \mid \mathbf{x}') < P(\text{No} \mid \mathbf{x}')$ , we label  $\mathbf{x}'$  to be “No”.

Decision making with the MAP rule

# CE 13: Other Applications

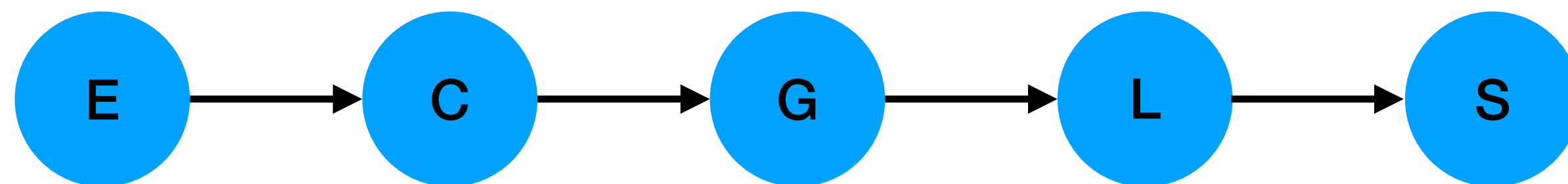
- What are some other applications where we would use Naïve Bayes? Why?

# Bayesian Network

# Exploiting Conditional Independence

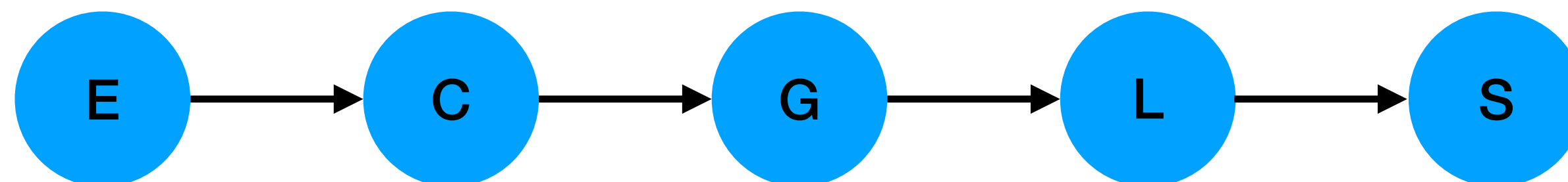
Consider a story:

If Leilani woke up too early  $E$ , Leilani probably needs a coffee  $C$ ; if Leilani needs coffee, she's likely grumpy  $G$ . If she is grumpy then it's possible that the lecture won't go smoothly  $L$ . If the lecture does not go smoothly then the students will likely be sad  $S$ .



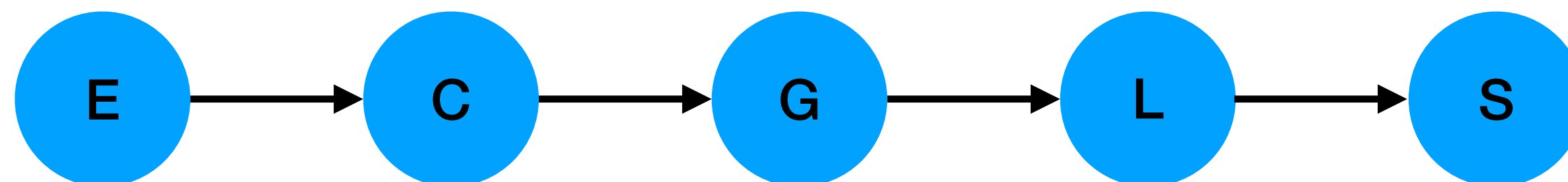
# Conditional Independence

- If you learned any of  $E$ ,  $C$ ,  $G$ , or  $L$ , would your assessment of  $p(S)$  change?
- If any of these are not seen to true, you would increase  $p(S)$  and decrease  $p(\neg S)$ . So  $S$  is not independent of  $E$ , or  $C$ , or  $G$ , or  $L$ .
- If you knew the value of  $L$  (true or false), would learning the value of  $E$ ,  $C$ , or  $G$  influence  $p(S)$ ?
- Influence that these factors have on  $S$  is mediated by their influence on  $L$ .
- Student's aren't sad because Leilani woke up early, they are sad because of the lecture. So  $S$  is independent of  $E$ ,  $C$ , and  $G$ , given  $L$ .



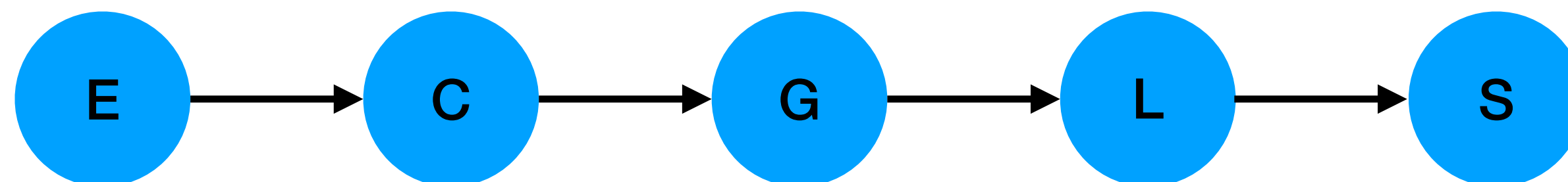
# Conditional Independence

- So  $S$  is independent of  $E$ ,  $C$ , and  $G$ , given  $L$ .
- Similarly
  - $L$  is independent of  $E$ , and  $C$ , given  $G$ .
  - $G$  is independent of  $E$ , given  $C$
- This means that:
  - $p(S | L, \{G, C, E\}) = p(S | L)$
  - $p(L | G, \{C, E\}) = p(L | G)$
  - $p(G | C, \{E\}) = p(G | C)$
  - $p(C | E)$  and  $p(E)$  does not simplify further.



# Conditional Independence

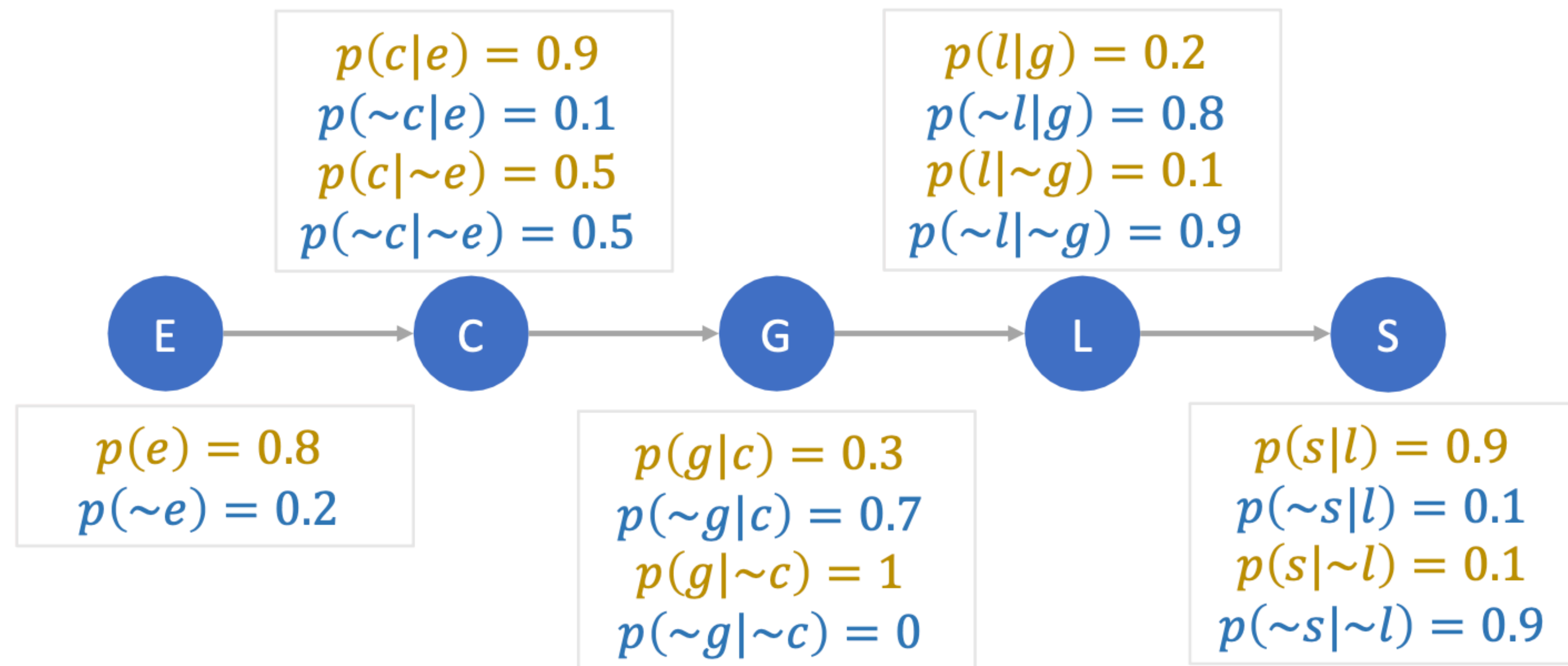
- By the chain rule (for any instantiation of S...E):
  - $P(S, L, G, C, E) = p(S | L, G, C, E)p(L | G, C, E)p(G | C, E)p(C | E)p(E)$
- By our independence assumptions:
  - $P(S, L, G, C, E) = p(S | L)p(L | G)p(G | C)p(C | E)p(E)$
- We can specify the full joint probability by specifying five local conditional probabilities
  - $p(S | L)$
  - $p(L | G)$
  - $p(G | C)$
  - $p(C | E)$  and  $p(E)$



# Example Quantification

Specifying the joint requires only 9 parameters

What is  $p(g)$





# Bayesian Networks

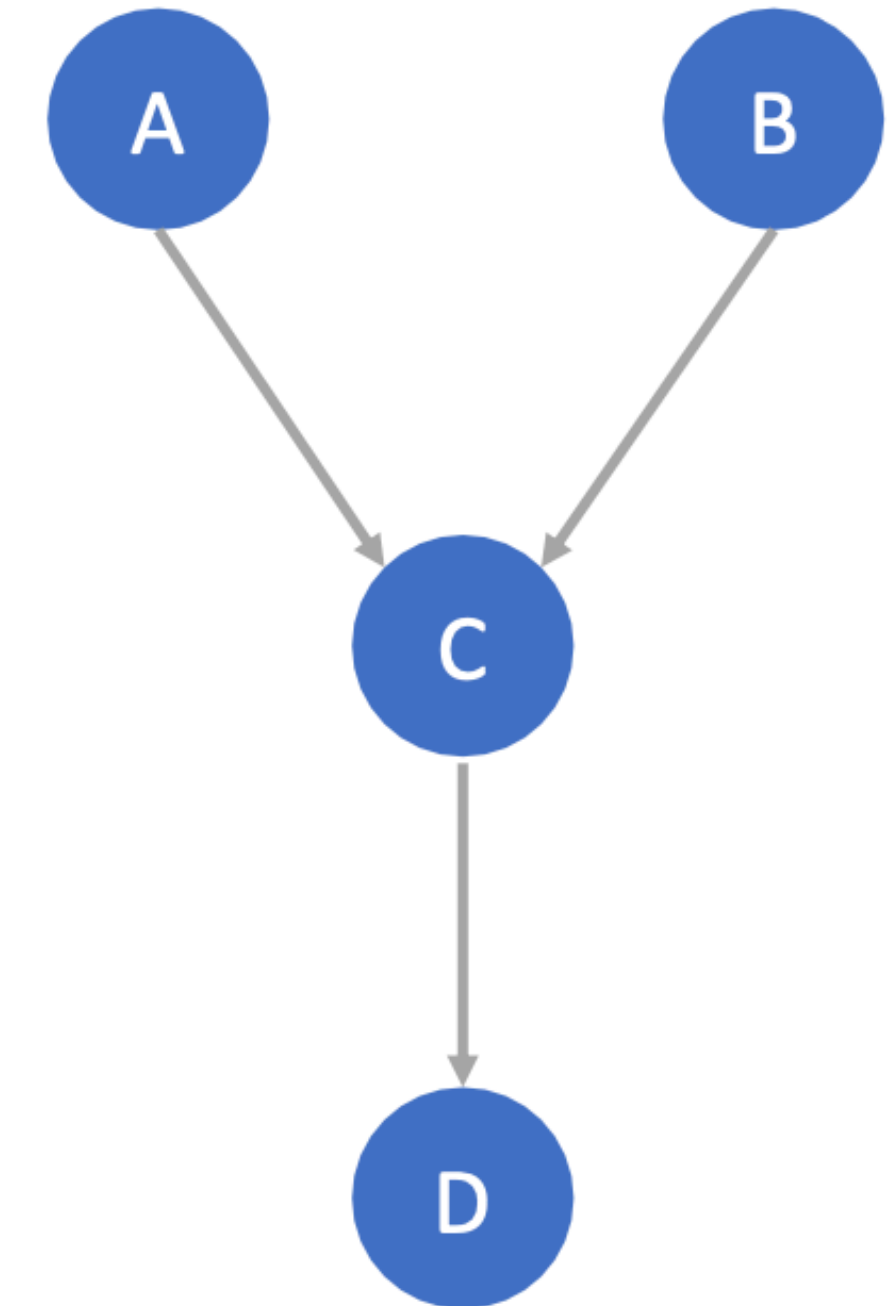
- The structure we just mentioned is a Bayesian network.
- Graphical representation of the direct dependencies over a set of variables + a set of conditional probability tables (CPTs) quantifying the strength of those influences.
- Bayesian Networks generalize the above ideas in very interesting ways, leading to effective means of representation and inference under uncertainty.

# Bayesian Networks

- A simple, graphical notation for conditional independence assertions resulting in a compact representation for the full joint distribution
- Syntax:
  - a set of nodes, one per random variable
  - a directed, acyclic graph (link = ‘direct influences’)
  - a conditional distribution (CPT) for each node given its parents:  
 $P(X_i | \text{Parents}(X_i))$

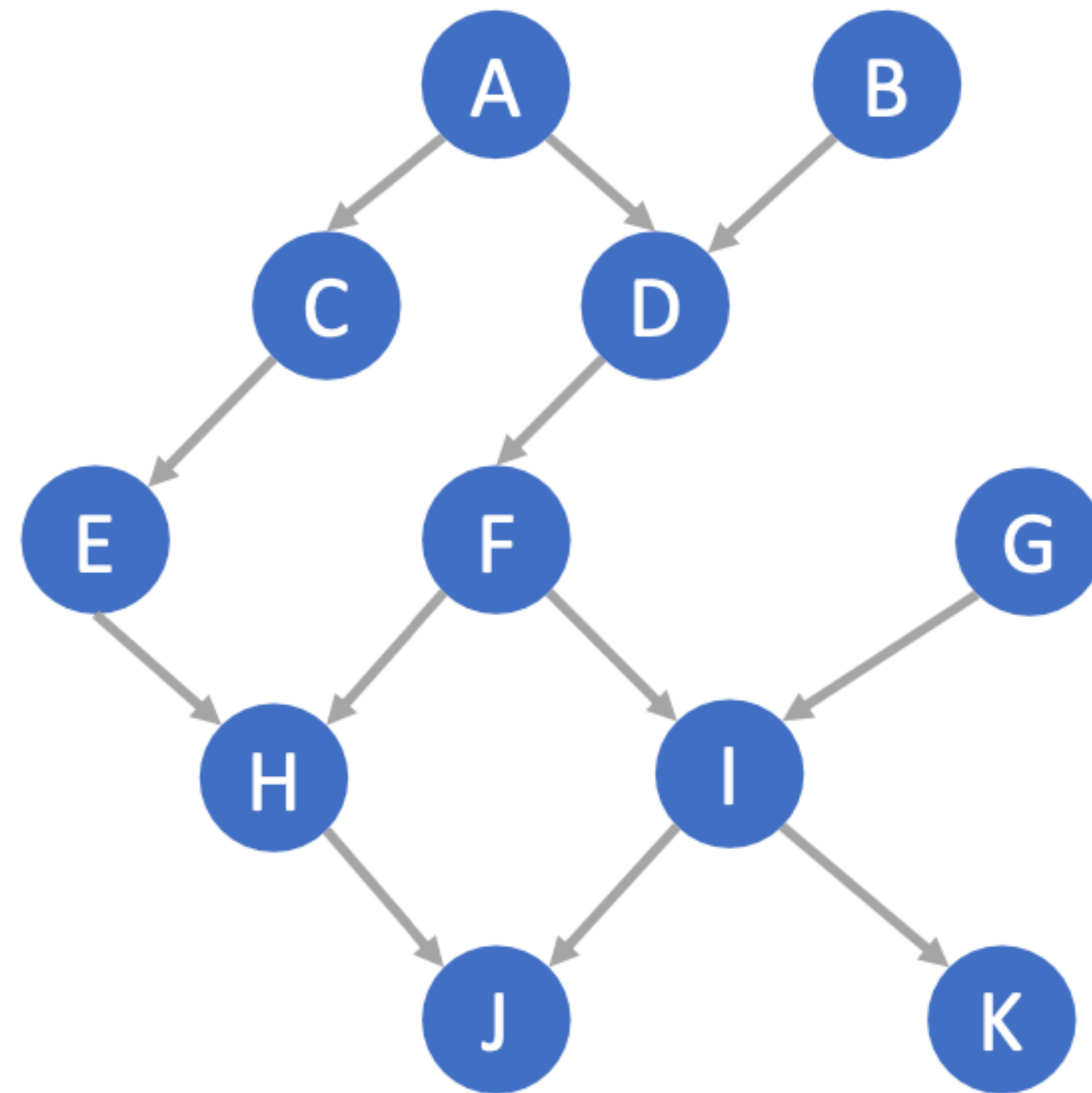
# Key Notions

- Some definitions:
  - **parents** of a node  $\rightarrow par(C) = \{A, B\}$
  - **children** of node  $\rightarrow children(A) = \{C\}$
  - **descendants** of a node  $\rightarrow descendants(B) = \{C, D\}$
  - **ancestors** of a node  $\rightarrow ancestors(D) = \{A, B, C\}$
  - **family**: set of nodes consisting of  $x_i$  and its parents  $\rightarrow family(C) = \{C, A, B\}$
- CPTs are defined over families in the BN



# An Example of a Bayes Net

- How many parameters do we need for the following BN?



# Semantics of a Bayesian Network

- The structure of the BN means: every  $x_i$  is conditionally independent of all of its non-descendants given its parents:
  - $p(x_i | X \cup \text{par}(x_i)) = p(\text{par}(x_i))$
  - For any subset  $S \subseteq \text{non} - \text{descendants}(x_i)$

# How to build a Bayesian Network

1. Define a total order over the random variables:  $(x_1, \dots, x_n)$

2. Apply the chain rule:

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

3. For each  $x_i$ , select the smallest set of predecessors  $par(x_i)$  such that:

$$p(x_i | x_1, \dots, x_{i-1}) = p(x_i | par(x_i))$$

4. Then we can rewrite

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | par(x_i))$$

# How to build a Bayesian Network (2)

5. This is a compact representation of the initial JPD. Factorization of the JPD based on existing conditional independencies among the variables
6. Construct the Bayesian Net (BN)
  - ✓ Nodes are the random variables
  - ✓ Draw a directed edge from each variable in  $par(x_i)$  to  $x_i$
  - ✓ Define a conditional probability table (CPT) for each variable  $x_i$ :  
$$p(x_i | par(x_i))$$

# Example for BN Construction

- You want to diagnose whether there is a fire in a building
- You can receive reports (possibly noisy) about whether everyone is **leaving** the building
- If everyone is leaving, this may have been caused by a **firealarm**
- If there is a fire alarm, it may have been caused by a **fire** or by **tampering**
- If there is a fire, there may be **smoke**



# Fire Diagnosis: Step 1

- Start by choosing the random variables for this domain:
  - Tampering (T) is true when the alarm has been tampered with
  - Fire (F) is true when there is a fire
  - Alarm (A) is true when there is an alarm
  - Smoke (S) is true when there is smoke
  - Leaving (L) is true if there are lots of people leaving the building
  - Report (R) is true if the sensor reports that lots of people are leaving the building

# Fire Diagnosis: Step 2

- Define total ordering of variables.
- Let's choose an order that follows the causal sequence of events:
  - Fire(F), Tampering (T), Alarm (A), Smoke (S), Leaving (L), Report (R)

# Fire Diagnosis: Step 3

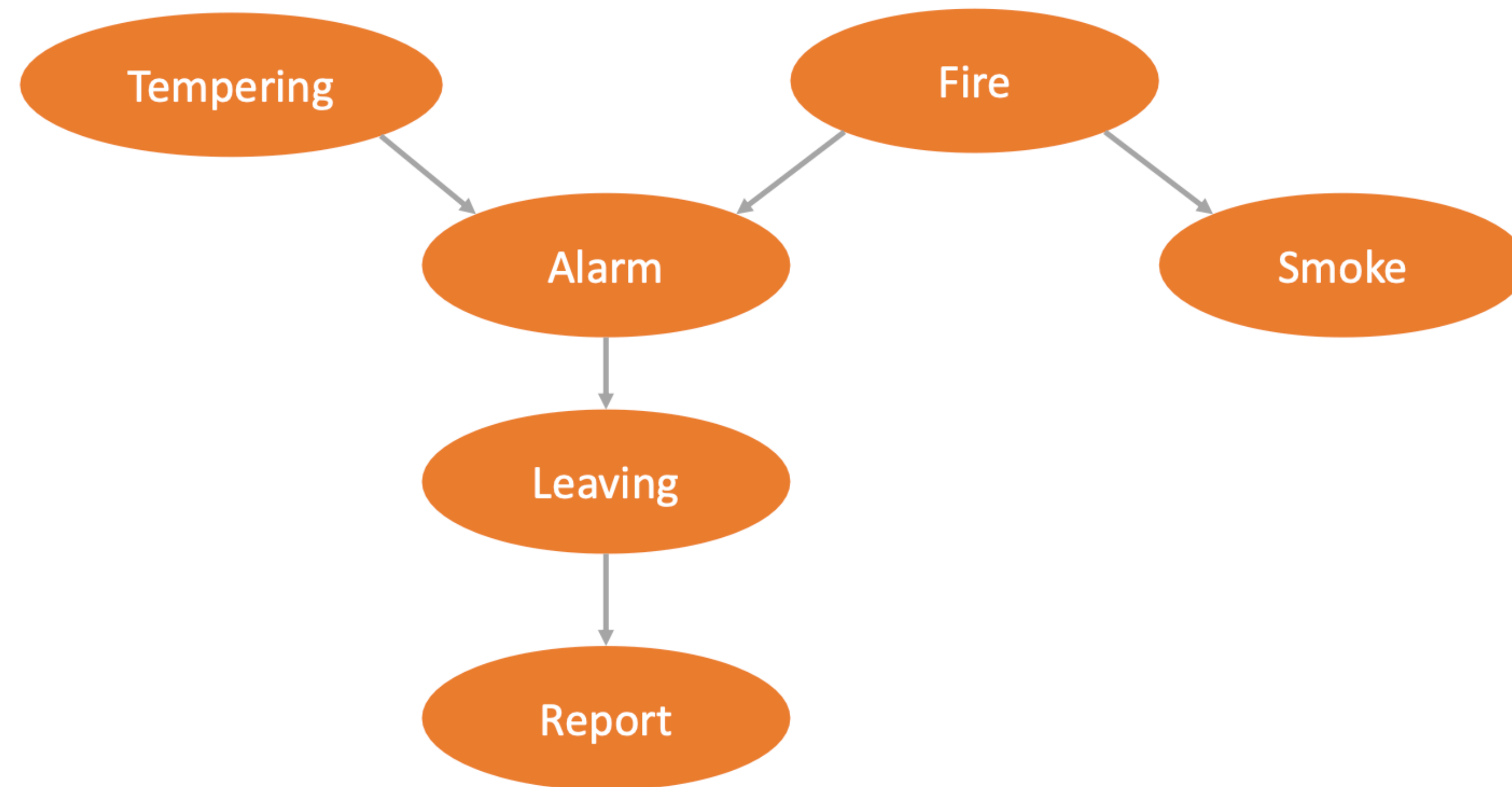
- Ordering:
  - Fire(F), Tampering (T), Alarm (A), Smoke (S), Leaving (L), Report (R)
- Apply the chain rule:
- $p(F, T, A, S, L, R) = p(F)p(T|F)p(A|F, T)p(S|F, T, A)p(L|F, T, A, S)p(R|F, T, A, S, L)$

# Fire Diagnosis: Step 4 & 5

- $p(F, T, A, S, L, R) = p(F)p(T|F)p(A|F, T)p(S|F, T, A)p(L|F, T, A, S)p(R|F, T, A, S, L)$
- For each variable,  $x_i$  choose parents  $par(x_i)$  and re-write the joint probability distribution:
  - $p(F, T, A, S, L, R) = p(F)p(T)p(A|F, T)p(S|F)p(L|A)p(R|L)$
- Now we need to build the BN based on the above JPD.

# Fire Diagnosis: Drawing BN

- $p(F, T, A, S, L, R) = p(F)p(T)p(A | F, T)p(S | F)p(L | A)p(R | L)$



# Summary and Next Time

- Today:
  - Naïve Bayes
  - Bayes Networks
- Thursday
  - Continue Bayes Networks
  - Markov Decision Processes