

CS 5100: Foundations of Artificial Intelligence

Bayesian Networks

Prof. Amy Sliva

November 10, 2011

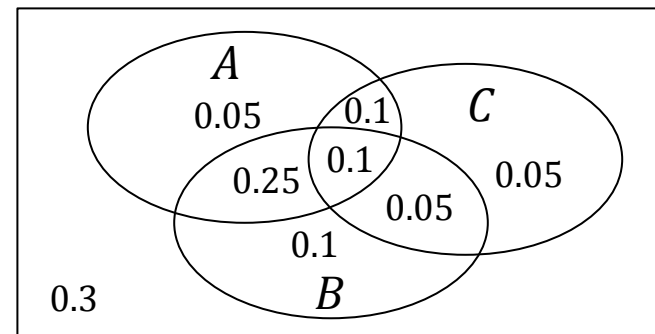
Outline

- Review joint and conditional probability
 - Bayes rule
 - Conditional independence
 - Bayesian networks
 - Naïve Bayes classification
-
- Thanks to Andrew Moore from CMU for some great slides (used with permission).

Joint probability distribution (JPD) review

- Recipe for making a joint distribution of M variables:
 1. Make a truth table listing **all combinations of values** of the variables (if there are M boolean variables, then the table will have 2^M rows)
 2. For each combination of values, say how probable it is
 3. Following the **axioms of probability**, those numbers must sum to 1

<i>A</i>	<i>B</i>	<i>C</i>	Probability
F	F	F	0.3
F	F	T	0.05
F	T	F	0.1
F	T	T	0.05
T	F	F	0.05
T	F	T	0.1
T	T	F	0.25
T	T	T	0.1



Using the joint distribution

- Once you have the JPD you can calculate the probability of **any logical expression** involving your variables

Gender	Hours Worked	Wealth	
Female	< 40.5	poor	0.253122
		rich	0.0245895
	> 40.5	poor	0.0421768
		rich	0.0116293
Male	< 40.5	poor	0.331313
		rich	0.0971295
	> 40.5	poor	0.134106
		rich	0.105933

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Using the joint distribution

$$\begin{aligned} P(\text{Poor Male}) \\ &= 0.3313 + 0.1341 \\ &= 0.4654 \end{aligned}$$

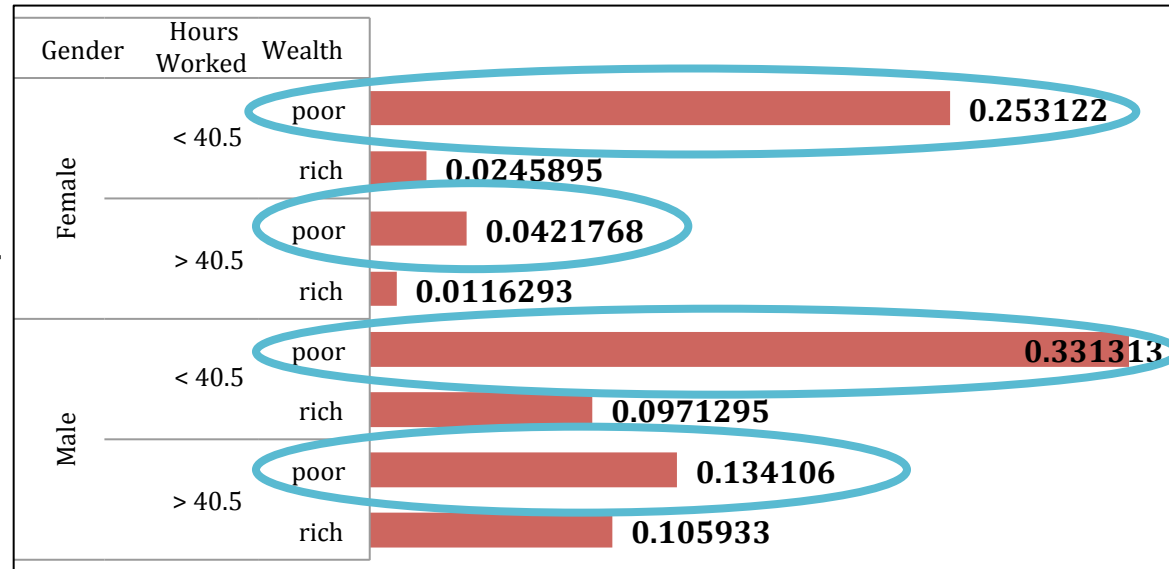
Gender	Hours Worked	Wealth	
Female	< 40.5	poor	0.253122
		rich	0.0245895
	> 40.5	poor	0.0421768
		rich	0.0116293
Male	< 40.5	poor	0.331313
		rich	0.0971295
	> 40.5	poor	0.134106
		rich	0.105933

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Using the joint distribution

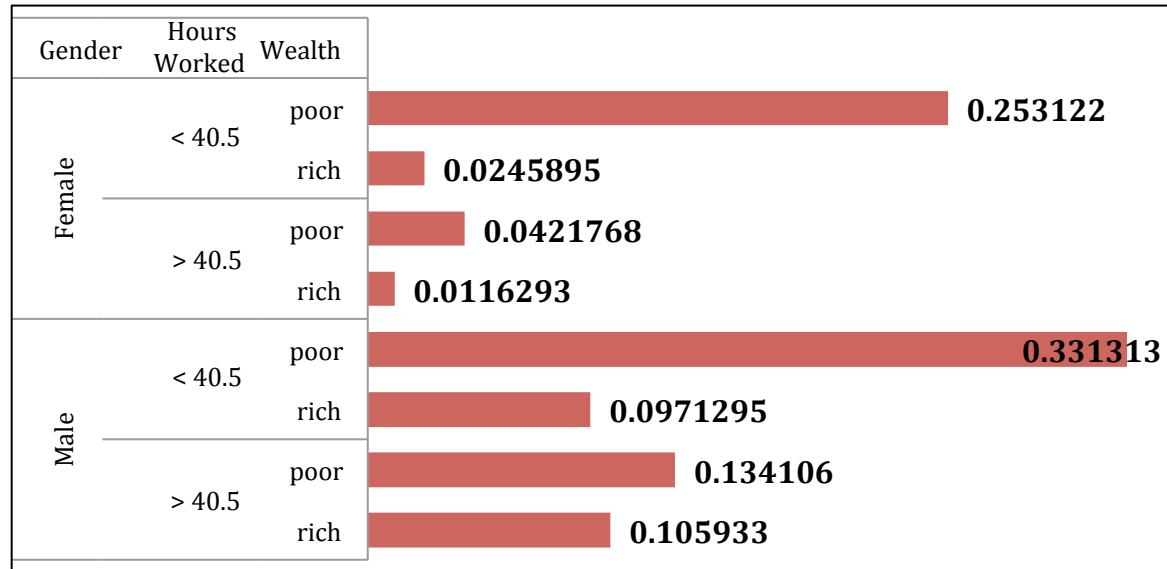
$P(\text{Poor})$

$$\begin{aligned} &= 0.2531 + 0.0421 + \\ &0.3313 + 0.1341 \\ &= 0.7604 \end{aligned}$$



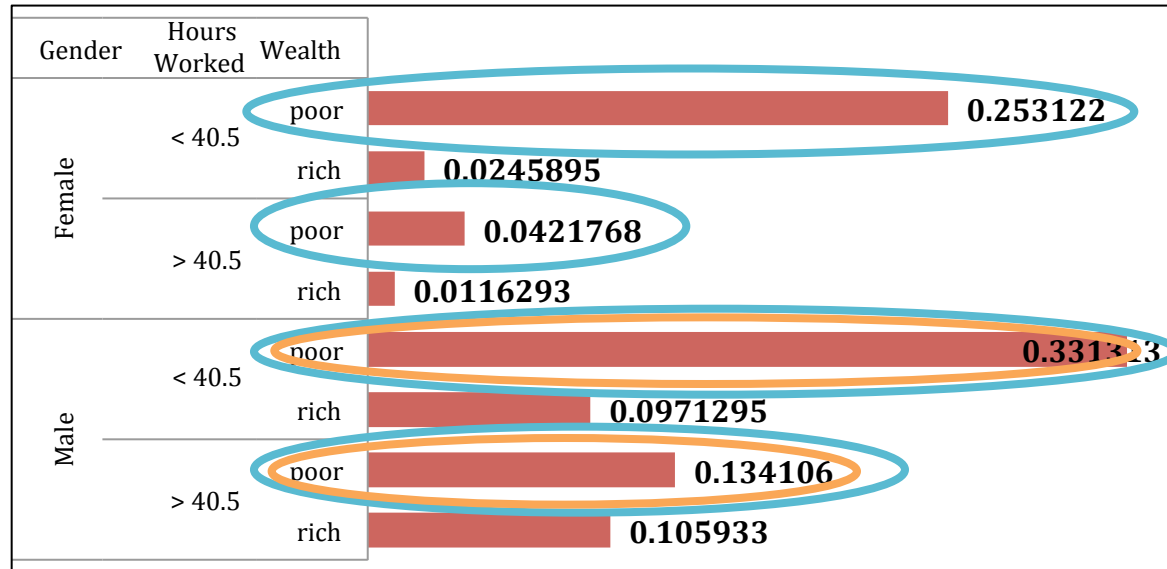
$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Inference with the joint distribution



$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

Inference with the joint distribution



$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

$$P(\text{Male} | \text{Poor}) = 0.4654 / 0.7604 = 0.612$$

The good and the bad with JPDs

- **Good news** 😊

Once you have a joint distribution, you can ask important questions about stuff that involves a lot of uncertainty!

- **Bad news** ☹️

Impossible to create for more than about 10 variables because there are so many numbers needed when you build the thing!

Solution—using fewer variables

- What **assumption** can we make to reduce the size of the problem?
- Suppose two events:
A = Amy teaches class (otherwise it is cancelled!)
S = it is sunny
- The joint pdf will have 4 entries
- To build the joint distribution we have to invent those 4 numbers....**OR WILL WE?**
 - We don't have to specify bottom level conjunctive events—i.e., $P(\neg A \wedge S)$ —IF...
 - ...instead may be more convenient to specify things like $P(A)$ or $P(S)$
- But just these marginal probs do not derive the joint distribution so cannot answer all questions

Independence

- The sunshine levels do not depend on and do not influence who is teaching
- This can be specified very simply:
$$P(S | A) = P(S)$$
- **This is a very powerful statement!**
- Requires **domain knowledge** and an understanding of causation, not just numerical probabilities

General definitions of independence

- From $P(S | A) = P(S)$, the axioms of probability imply:
 1. $P(\neg S | A) = P(\neg S)$
 2. $P(A | S) = P(A)$
 3. $P(A \wedge S) = P(A) P(S)$
 4. $P(\neg A \wedge S) = P(\neg A) P(S)$
 5. $P(A \wedge \neg S) = P(A) P(\neg S)$
 6. $P(\neg A \wedge \neg S) = P(\neg A) P(\neg S)$
- In **general** $P(A = u \wedge S = v) = P(A = u) P(S = v)$ for each combination of $u = \text{True/False}$ and $v = \text{True/False}$

Joint distribution assuming independence

- We know:

$$P(A) = 0.6$$

$$P(S) = 0.3$$

$$P(S | A) = P(S)$$

- From these statements we can derive the full joint distribution (**assuming independence**)

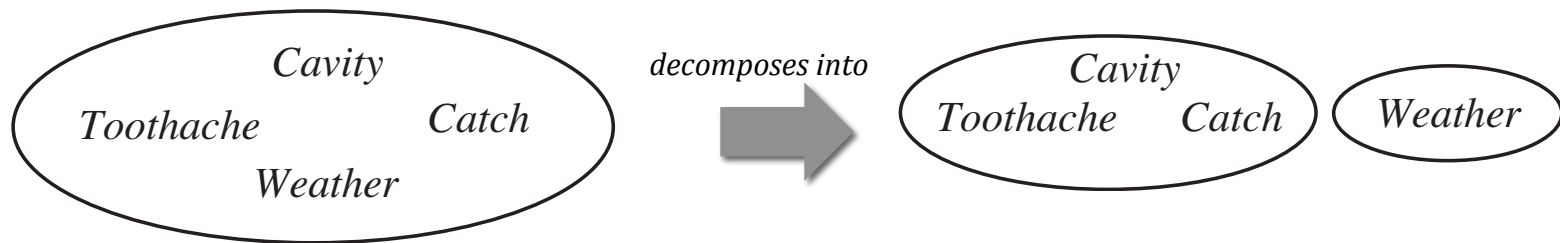
<i>A</i>	<i>S</i>	Probability
T	T	0.18
T	F	0.42
F	T	0.12
F	F	0.28

- Since we have the joint pdf, we can make any query!

Reduction in the joint pdf representation

- A and B are independent iff
 $\mathbf{P}(A | B) = \mathbf{P}(A)$ or $\mathbf{P}(B|A) = \mathbf{P}(B)$ or $\mathbf{P}(A,B) = \mathbf{P}(A)\mathbf{P}(B)$
- E.g., $\mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather})$ has $2 \times 2 \times 2 \times 4$ entries in JPD
 $= \mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity})\mathbf{P}(\textit{Weather})$ has $2 \times 2 \times 2 + 4$ entries

32 entries reduced to 12!



- **Absolute independence** powerful, but rare
 - In general, reduces exponential to linear complexity

Conditional probability review

- **Conditional** or **posterior** probabilities—based on known information
- Definition of **conditional probability**:
$$P(a | b) = \frac{P(a \wedge b)}{P(b)} \quad \text{if } P(b) > 0$$
- Corollary: the **chain rule**
$$P(a \wedge b) = P(a | b) P(b)$$
- Extend to $P(a \wedge b \wedge c \wedge \dots) = ?$

Chain rule follows from this definition

- **Product rule** gives an alternative formulation:

$$P(a \wedge b) = P(a | b)P(b) = P(b | a)P(a)$$

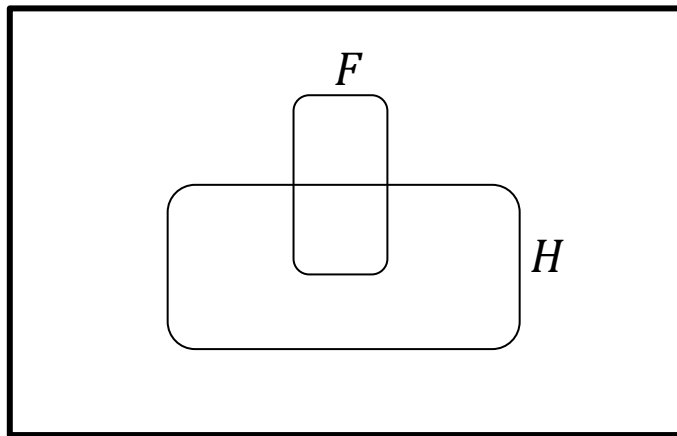
- **Chain rule** is derived by successive application of product rule

$$\begin{aligned} \mathbf{P}(X_1, \dots, X_n) &= \mathbf{P}(X_1, \dots, X_{n-1}) \mathbf{P}(X_n | X_1, \dots, X_{n-1}) \\ &= \mathbf{P}(X_1, \dots, X_{n-2}) \mathbf{P}(X_{n-1} | X_1, \dots, X_{n-2}) \mathbf{P}(X_n | X_1, \dots, X_{n-1}) \\ &= \dots \\ &= \mathbf{P}(X_1) \mathbf{P}(X_2 | X_1) \mathbf{P}(X_3 | X_1, X_2) \dots \mathbf{P}(X_n | X_1, \dots, X_{n-1}) \end{aligned}$$

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

Conditional probability example

- $P(A | B)$ = Fraction of worlds in which B is true that also have A true



H = Have a headache

F = Coming down with flu

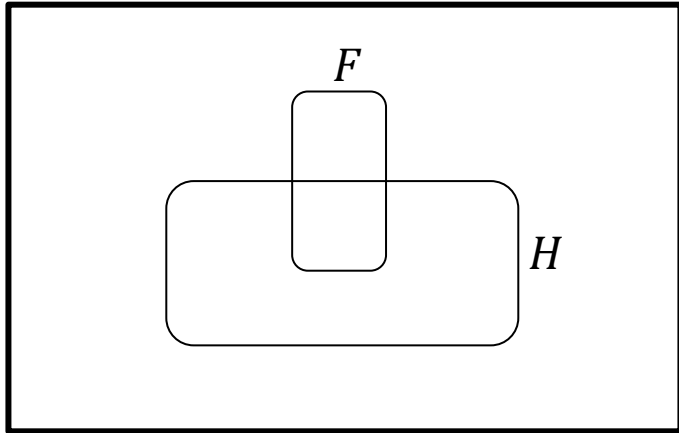
$$P(H) = 0.1$$

$$P(F) = 0.025$$

$$P(H | F) = 0.5$$

Headaches are rare and flu is rare, but if you are coming down with flu there is a 50-50 chance you'll have a headache.

Conditional probability example



H = Have a headache

F = Coming down with flu

$$P(H) = 0.1$$

$$P(F) = 0.025$$

$$P(H | F) = 0.5$$

$P(H | F) =$ **Fraction of flu-infected worlds in which you have a headache**

$$= \frac{\# \text{ worlds with flu and headache}}{\# \text{ worlds with flu}}$$

$$= \frac{\text{Area of "H and F" region}}{\text{Area of "F" region}}$$

$$= \frac{P(H \wedge F)}{P(F)}$$

Conditional probability exercise

	Likes Football	Dislikes Football	Neutral
<i>Male</i>	0.25	0.1	0.15
<i>Female</i>	0.1	0.3	0.1

- Calculate:

$$P(\text{Likes Football} \mid \text{Male}) = ?$$

$$P(\neg \text{Likes} \mid \text{Female}) = ?$$

Conditional probability exercise

	Likes Football	Dislikes Football	Neutral
<i>Male</i>	0.25	0.1	0.15
<i>Female</i>	0.1	0.3	0.1

- Calculate:

$$P(\text{Likes Football} \mid \text{Male}) = \mathbf{0.5}$$

$$P(\neg \text{Likes} \mid \text{Female}) = \mathbf{0.8}$$

Conditional probability-based AI models

- Given a set of random variables X we are interested in the posterior joint distribution of the **query variables** Y given specific values e for the **evidence variables** E

$$\mathbf{P}(Y \mid E = e) = \alpha \mathbf{P}(Y \wedge E = e)$$

- Might have some hidden variables $H = X - Y - E$ giving

$$\mathbf{P}(Y \mid E = e) = \alpha \sum_h \mathbf{P}(Y \wedge E = e \wedge H = h)$$

- Note: what is α ?**

- From the definition $\mathbf{P}(A \mid B) = \frac{\mathbf{P}(A \wedge B)}{\mathbf{P}(B)}$

α is the denominator $1 / \mathbf{P}(E = e)$. $\mathbf{P}(E = e)$ can be computed from joint distribution as $\sum_h \mathbf{P}(E = e \wedge H = h)$

Example (medical diagnosis)

- **Causal model:**

$$D \rightarrow I \rightarrow S \quad (Y \rightarrow H \rightarrow E)$$

Cancer \rightarrow anemia \rightarrow fatigue

Kidney \rightarrow disease anemia \rightarrow fatigue

$$P(Y = \textit{cancer} \mid E = \textit{fatigue}) =$$

$$\alpha [P(Y = \textit{cancer} \wedge E = \textit{fatigue} \wedge \textit{anemia}) + \\ P(Y = \textit{cancer} \wedge E = \textit{fatigue} \wedge \neg \textit{anemia})]$$

$$\alpha = 1/P(E = \textit{fatigue}) \text{ or } 1/[P(E = \textit{fatigue} \wedge \textit{anemia}) + \\ P(E = \textit{fatigue} \wedge \neg \textit{anemia})]$$

Conditional probability analysis

- $\mathbf{P}(Y \mid E = e) = \alpha \mathbf{P}(Y \wedge E = e) = \alpha \sum_h \mathbf{P}(Y \wedge E = e \wedge H = h)$
- Terms in the summation are **joint distribution entries** because Y , E , and H together exhaust the full set of random variables
- **Obvious problems**
 1. Time and space complexity $O(d^n)$ where d is the largest arity
 2. Where do we get the numbers to solve real world problems?
- Independence assumption provides a solution to 1...

...but where do we find the numbers?

- Assuming independence, doctors may be able to **estimate** $P(\textit{symptom} \mid \textit{disease})$ for each S/D pair (causal reasoning)
- Hard to estimate what we really **need** to know: $P(\textit{disease} \mid \textit{symptom})$
- This is why **Bayes rule** is so important in probabilistic AI!

Bayes Rule

- Product rule $P(a \wedge b) = P(a | b)P(b) = P(b | a)P(a)$

$$\Rightarrow \text{Bayes rule: } P(a | b) = \frac{P(b | a)P(a)}{P(b)}$$

- Or in distribution form:

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)} = \alpha P(X | Y) P(Y)$$

- Useful for assessing **diagnostic** probability from **causal** probability:

- $P(\text{Cause} | \text{Effect}) = \frac{P(\text{Effect} | \text{Cause}) P(\text{Cause})}{P(\text{Effect})}$

- E.g., Let M be meningitis, S be stiff neck

$$P(M | S) = P(S | M) P(M) / P(S) = 0.8 \times 0.0001 / 0.1 = 0.0008$$

- Note: posterior probability of meningitis is still very small!

More general forms of Bayes rule

- $$P(A | B) = \frac{P(B | A)P(A)}{P(B | A)P(A) + P(B | \neg A)P(\neg A)}$$

- $$P(A | B) = \frac{P(B | A, e)P(A | e)}{P(B | e)}$$

- $$P(A = v_i | B) = \frac{P(B | A = v_i)P(A = v_i)}{\sum_{k=1}^n P(B | A = v_k)P(A = v_k)}$$

Conditional independence

- $\mathbf{P}(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$ has $2^3 - 1 = 7$ independent entries
 - If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
 - $\mathbf{P}(\textit{catch} \mid \textit{toothache}, \textit{cavity}) = \mathbf{P}(\textit{catch} \mid \textit{cavity})$
 - The same independence holds if I haven't got a cavity
 - $\mathbf{P}(\textit{catch} \mid \textit{toothache}, \neg\textit{cavity}) = \mathbf{P}(\textit{catch} \mid \neg\textit{cavity})$
- Catch is **conditionally independent** of Toothache given Cavity: $\mathbf{P}(\textit{Catch} \mid \textit{Toothache}, \textit{Cavity}) = \mathbf{P}(\textit{Catch} \mid \textit{Cavity})$
 - Equivalent statements (from original definitions of independence)
 - $\mathbf{P}(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) = \mathbf{P}(\textit{Toothache} \mid \textit{Cavity})$
 - $\mathbf{P}(\textit{Toothache}, \textit{Catch} \mid \textit{Cavity}) = \mathbf{P}(\textit{Toothache} \mid \textit{Cavity})\mathbf{P}(\textit{Catch} \mid \textit{Cavity})$

Conditional independence (cont.)

- Write out full joint distribution using chain rule:

$$\begin{aligned}\mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) &= \mathbf{P}(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity})\mathbf{P}(\textit{Catch}, \textit{Cavity}) \\ &= \mathbf{P}(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity})\mathbf{P}(\textit{Catch} \mid \textit{Cavity})\mathbf{P}(\textit{Cavity}) \\ &= \mathbf{P}(\textit{Toothache} \mid \textit{Cavity})\mathbf{P}(\textit{Catch} \mid \textit{Cavity})\mathbf{P}(\textit{Cavity})\end{aligned}$$

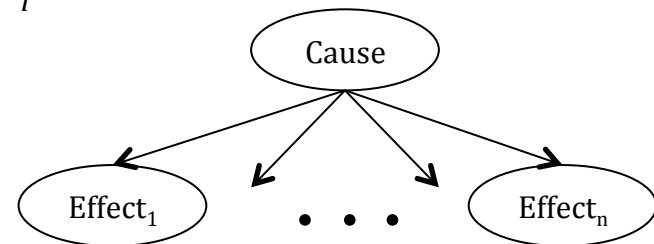
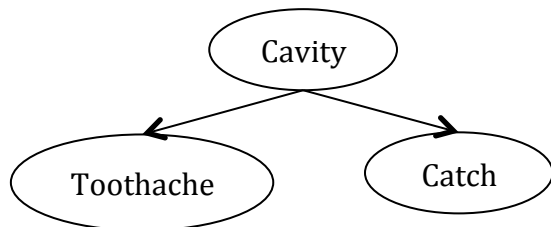
i.e., $2 + 2 + 1 = 5$ independent numbers

- In most cases, conditional independence reduces size of representation from **exponential** in n to **linear** in n
- Conditional independence is most basic and robust form of knowledge about uncertain environments

Bayes rule and conditional independence

- $\mathbf{P}(\text{Cavity} \mid \text{toothache catch})$
= $\mathbf{P}(\text{toothache catch} \mid \text{Cavity})\mathbf{P}(\text{Cavity})$
= $\mathbf{P}(\text{toothache} \mid \text{Cavity})\mathbf{P}(\text{catch} \mid \text{Cavity})\mathbf{P}(\text{Cavity})$
- We say: “toothache and catch are independent, given cavity”
 - Cavity **separates** Toothache and Catch because it is a direct cause of both
 - Example of a **naïve Bayes** model

- $\mathbf{P}(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = \mathbf{P}(\text{Cause}) \prod_i \mathbf{P}(\text{Effect}_i \mid \text{Cause})$



- Total number of parameters is **linear** in n (number of systems)
- This is our first Bayesian inference net!

Revisiting our joint distribution example

- A : Amy teaches the class
- S : It is sunny
- L : The lecturer arrives *slightly* late
- Assume lecturers are sometimes delayed by bad weather
- Start by writing down knowledge we're happy about:
 - $P(S | A) = P(S)$, $P(S) = 0.3$, $P(A) = 0.6$
 - **Lateness is not independent of the weather and is not independent of the lecturer**
 - **Weather and lecturer are independent**

Reduce complexity with independence

- A : Amy teaches the class
- S : It is sunny
- L : The lecturer arrives *slightly* late
- Assume lecturers are sometimes delayed by bad weather
- Start by writing down knowledge we're happy about:
 - $P(S | A) = P(S)$, $P(S) = 0.3$, $P(A) = 0.6$
 - **Lateness is not independent of the weather and is not independent of the lecturer**
 - **Weather and lecturer are independent**
- Know the joint probability of S and A , so now need:
 - $P(L | S, A)$ for the 4 cases where S and A are true/false

Reduce complexity with independence

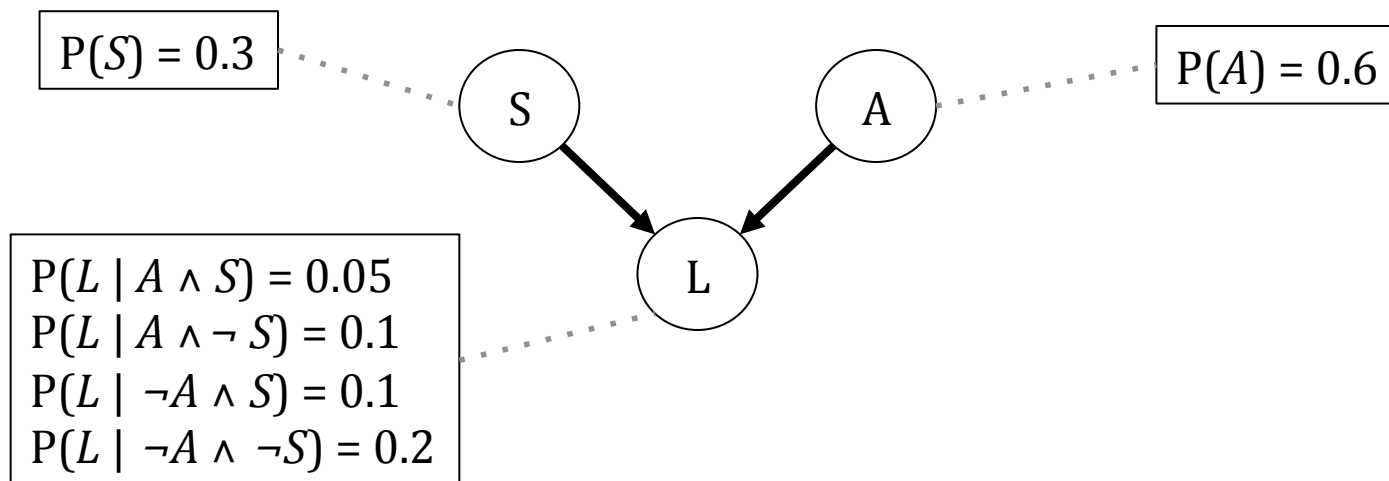
- A : Amy teaches the class
- S : It is sunny
- L : The lecturer arrives *slightly* late
- Assume lecturers are sometimes delayed by bad weather

$P(S A) = P(S)$	$P(L A \wedge S) = 0.05$
$P(S) = 0.3$	$P(L A \wedge \neg S) = 0.1$
$P(A) = 0.6$	$P(L \neg A \wedge S) = 0.1$
	$P(L \neg A \wedge \neg S) = 0.2$

- Now we can **derive** a full joint distribution with a “mere” six numbers instead of seven
 - NOTE: Savings are larger for larger numbers of variables

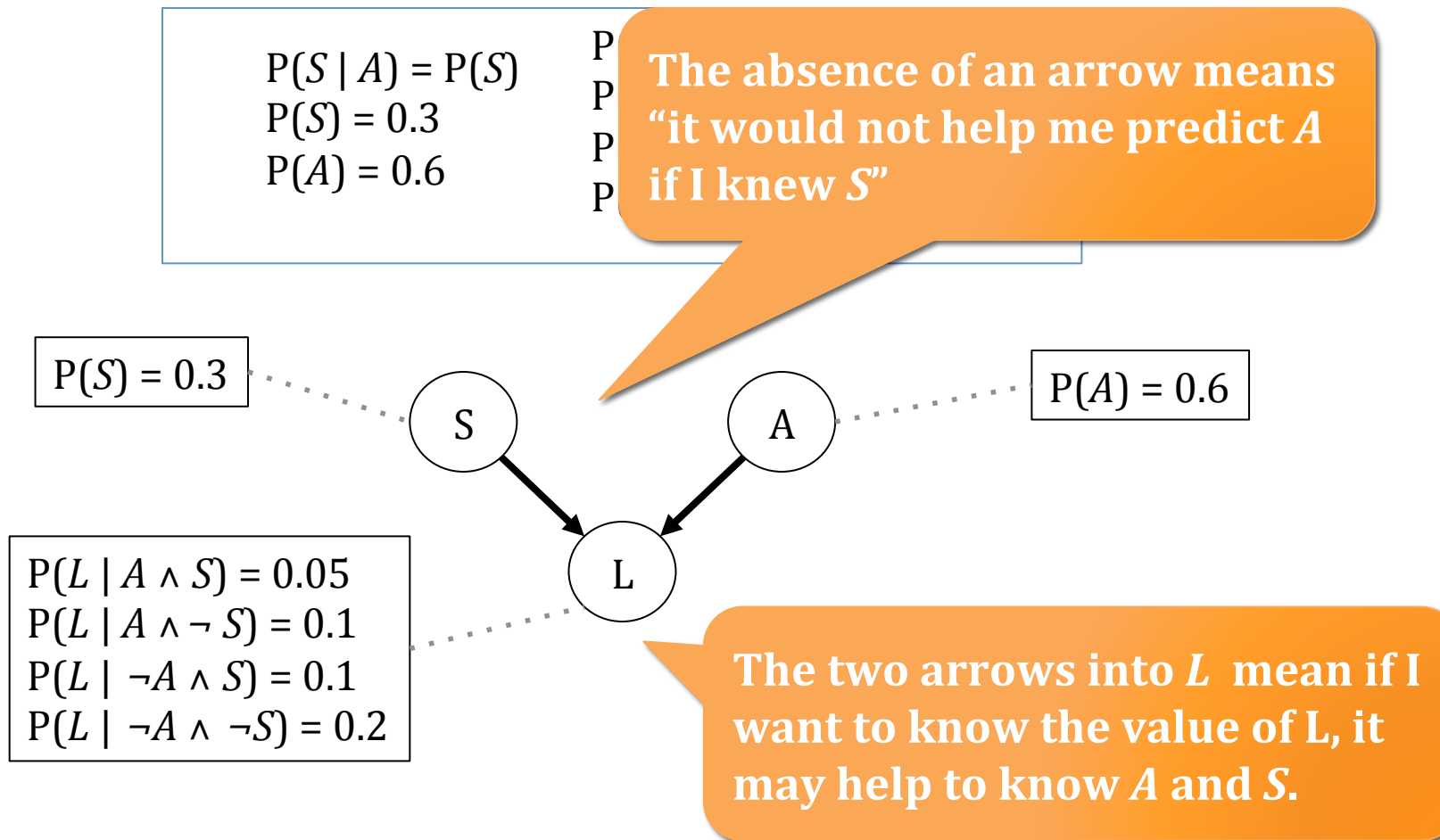
Graphical representation of independence

$P(S A) = P(S)$	$P(L A \wedge S) = 0.05$
$P(S) = 0.3$	$P(L A \wedge \neg S) = 0.1$
$P(A) = 0.6$	$P(L \neg A \wedge S) = 0.1$
	$P(L \neg A \wedge \neg S) = 0.2$



- This is our second Bayesian inference net!

Graphical representation of independence



- This is our second Bayesian inference net!

Example of conditional independence

- A : Amy teaches the class
- L : The lecturer arrives *slightly* late
- P : The lecture concerns probability
- Assume:
 - Oliver has a higher chance of being late than Amy
 - Amy has a higher chance of giving probability lectures.
- **Now what kind of independence can we find?**
- How about...
 - $P(L | A) = P(L)$?
 - $P(P | A) = P(P)$?
 - $P(L | P) = P(L)$?

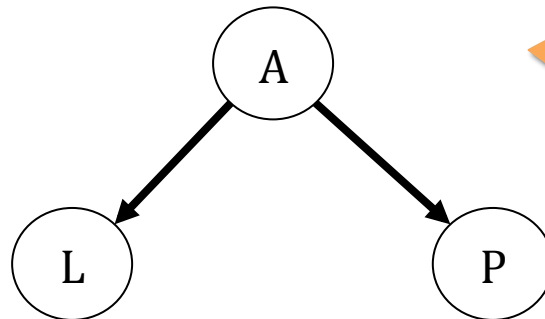
Example of conditional independence

- Once you know who the lecturer is, then whether they arrive late does not affect whether the lecture concerns probability:

$$P(P | A, L) = P(P | A)$$

$$P(P | \neg A, L) = P(P | \neg A)$$

- P and L are **conditionally independent** given A



Given knowledge of A , knowing anything else in the diagram won't help us with L

Conditional independence reduces complexity

- We can write down $P(A)$ —already know this
- Since we know L is only directly influenced by A , we can write down values of $P(L | A)$ and $P(L | \neg A)$ to fully specify L 's behavior.
- Same goes for P .

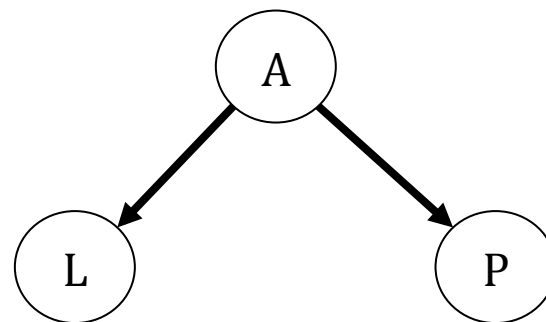
$$P(A) = 0.6$$

$$P(L | A) = 0.085$$

$$P(L | \neg A) = 0.17$$

$$P(P | A) = 0.6$$

$$P(P | \neg A) = 0.3$$



- Again, we can obtain any member of the **full joint distribution!**

How to construct a Bayes net

- Expand the example to include 5 variables:
 - T : The lecture started by 6:05
 - A : Amy teaches the class
 - L : The lecturer arrives *slightly* late
 - P : The lecture concerns probability
 - S : It is sunny
- **T only directly influenced by L (conditionally independent of A , P , and S given L)**
- **L only directly influenced by A and S (conditionally independent of P given A and S)**
- **P only directly influenced by A (conditionally independent of T , L , and S given A)**
- **A and S are independent**

Making a Bayes net

S

A

L

T

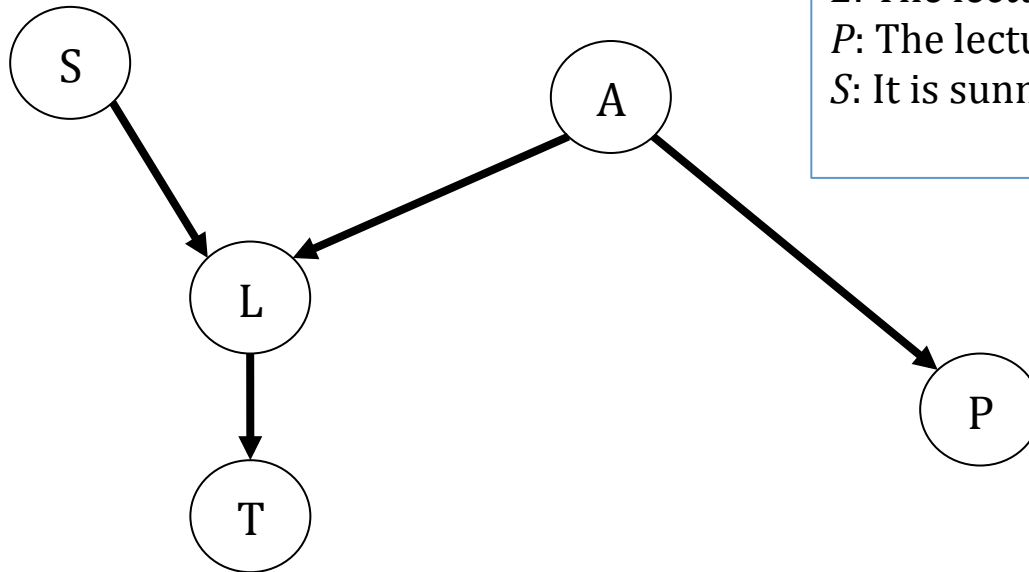
P

T: The lecture started by 6:05
A: Amy teaches the class
L: The lecturer arrives *slightly* late
P: The lecture concerns probability
S: It is sunny

- **Step 1: add variables**

- Choose the variables you want to include in the net

Making a Bayes net

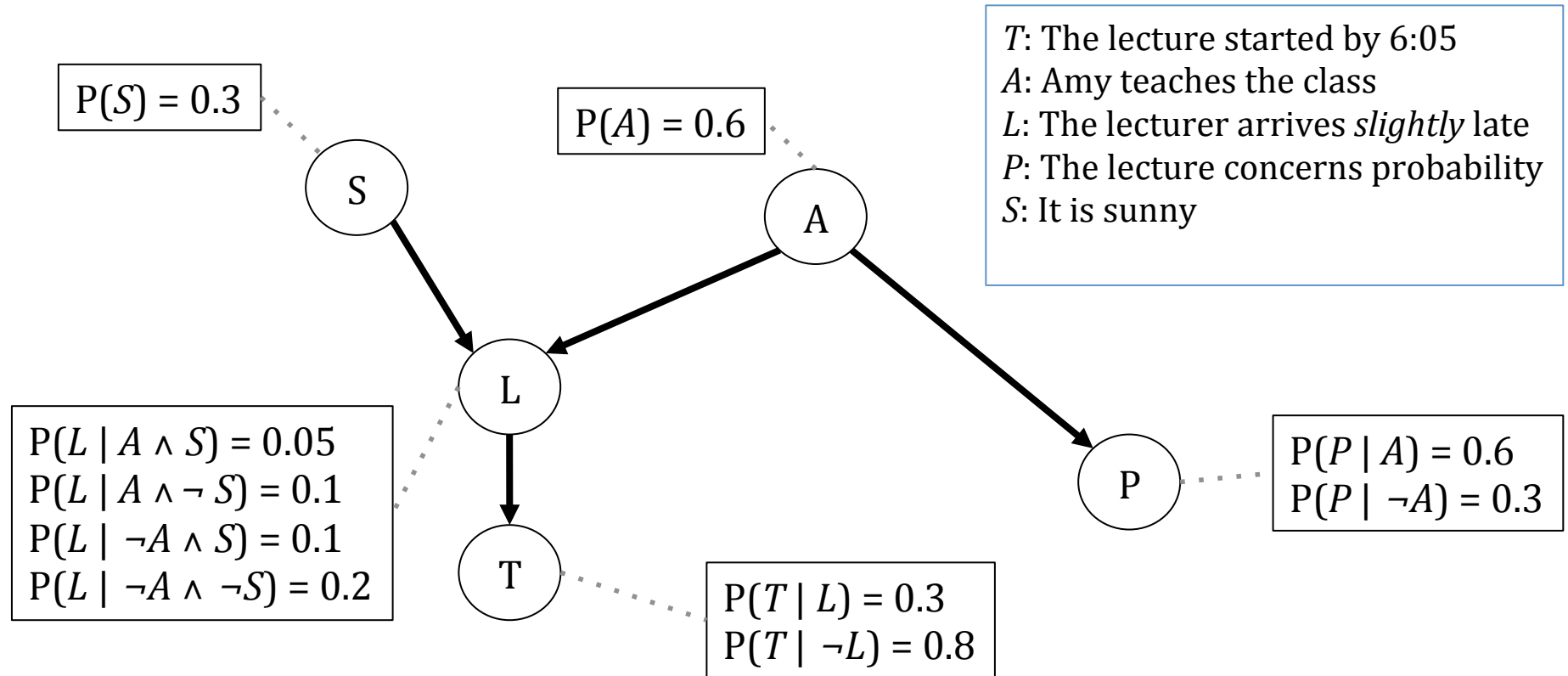


T: The lecture started by 6:05
A: Amy teaches the class
L: The lecturer arrives *slightly* late
P: The lecture concerns probability
S: It is sunny

- **Step 2: add links**

- Link structure must be acyclic
- If node X is given parents $\{Q_1, Q_2, \dots, Q_n\}$, then any non-descendant of X is conditionally independent of X given $\{Q_1, Q_2, \dots, Q_n\}$

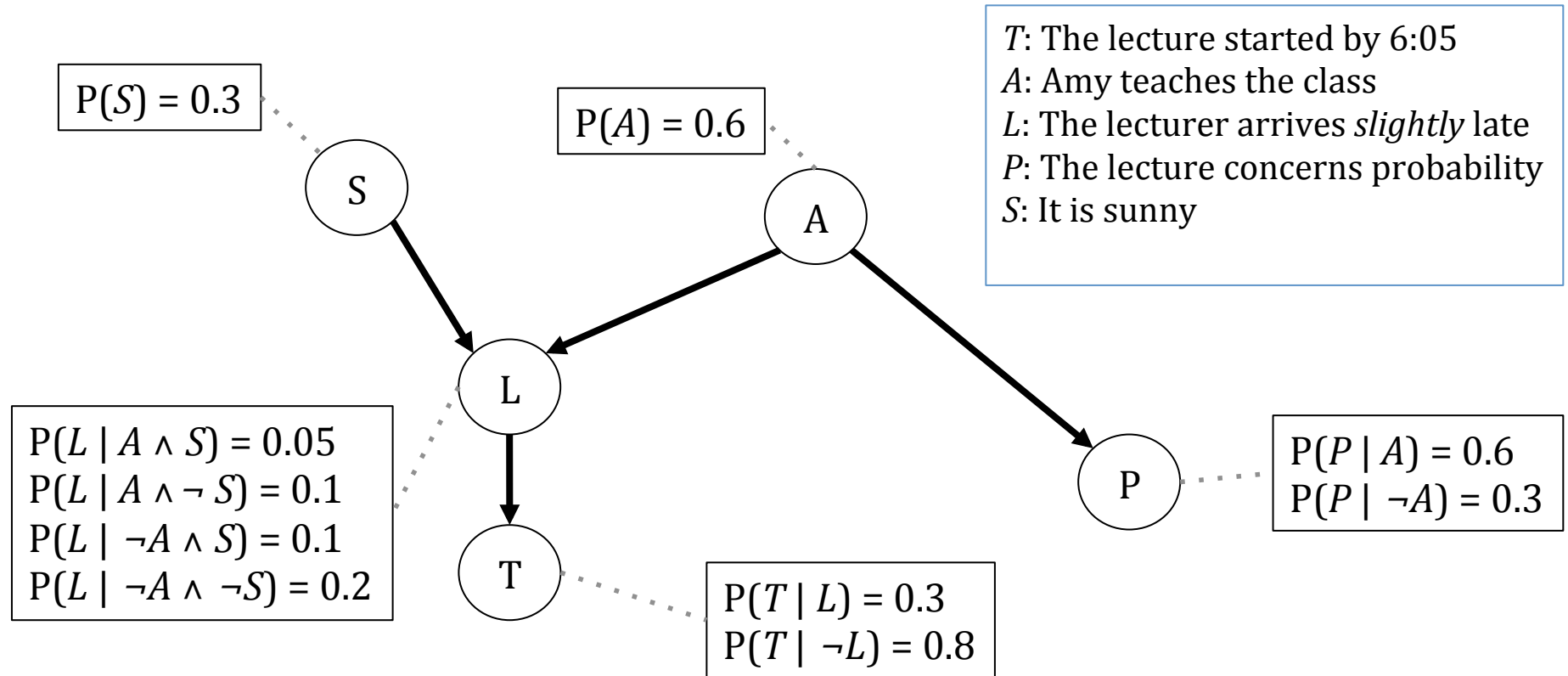
Making a Bayes net



- **Step 3: add a probability table for each node**

- The table for node X must list $P(X | \text{Parent values})$ for each possible combination of parent values

Making a Bayes net



- **Two unconnected variables may still be correlated**
- Each node is conditionally independent of all non-descendants in the tree **given its parents**
- We can deduce many other conditional independence relations from a Bayes net

Bayes nets formalized

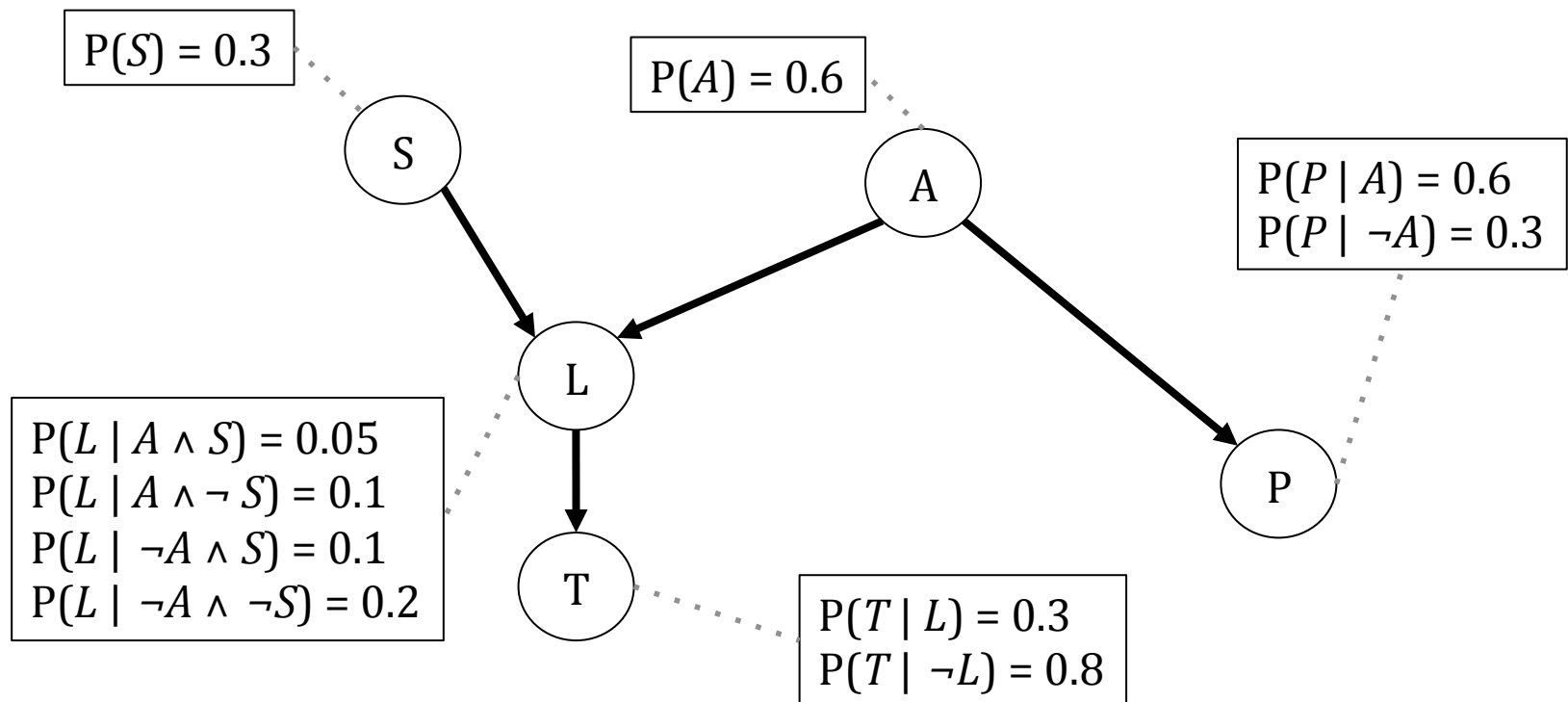
- Bayes net (or **belief network**)—augmented directed, acyclic graph represented by the pair V, E where:
 - V is a set of vertices
 - E is a set of directed edges (no loops of any length are allowed!)
- Each vertex contains the following information
 - Name of a random variable
 - **Probability distribution table** indicating how the probability of this variable's values depends on all possible combinations of parental values

Formal procedure for building a Bayes net

1. Choose a set of relevant variables
2. Choose an ordering for them (often requires domain knowledge)
 - Assume the variables are X_1, \dots, X_m where X_1 is the first in the ordering, X_2 is the second, etc.
3. For $i = 1$ to m :
 - a) Add the node X_i to the network
 - b) Set $Parents(X_i)$ to be the **minimal subset** $\{X_1, \dots, X_{i-1}\}$ s.t. we have conditional independence of X_i and all other members of $\{X_1, \dots, X_{i-1}\}$ given $Parents(X_i)$
 - c) Define the conditional probability table of $P(X_i \mid \text{Assignments of } Parents(X_i))$

Calculating using Bayes nets

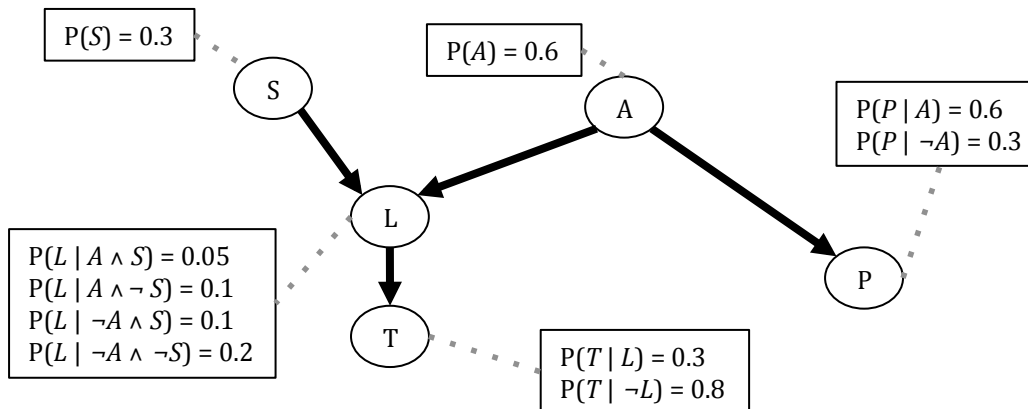
- How to compute an entry in the joint distribution?
E.g., what is $P(S \wedge \neg A \wedge L \wedge \neg P \wedge T)$?



Calculating using Bayes nets

- Use the chain rule and conditional independence!

$$\begin{aligned} & P(S \wedge \neg A \wedge L \wedge \neg P \wedge T) \\ &= P(T \mid \neg P \wedge L \wedge \neg A \wedge S) P(\neg P \wedge L \wedge \neg A \wedge S) \\ &= P(T \mid L) P(\neg P \wedge L \wedge \neg A \wedge S) \\ &= P(T \mid L) P(\neg P \mid L \wedge \neg A \wedge S) P(L \wedge \neg A \wedge S) \\ &= P(T \mid L) P(\neg P \mid \neg A) P(L \wedge \neg A \wedge S) \\ &= P(T \mid L) P(\neg P \mid \neg A) P(L \mid \neg A \wedge S) P(\neg A \wedge S) \\ &= P(T \mid L) P(\neg P \mid \neg A) P(L \mid \neg A \wedge S) P(\neg A \mid S) P(S) \\ &= P(T \mid L) P(\neg P \mid \neg A) P(L \mid \neg A \wedge S) P(\neg A) P(S) \end{aligned}$$



The general case

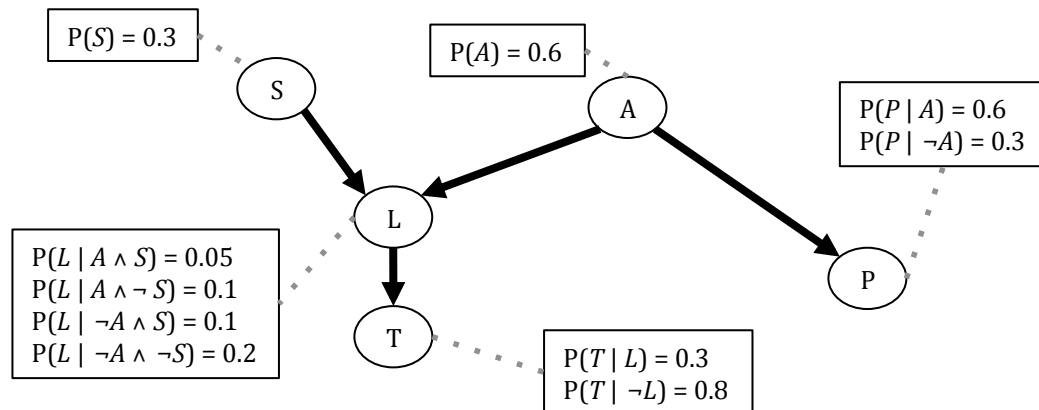
- $\mathbf{P}(X_1, \dots, X_n) = \mathbf{P}(X_1, \dots, X_{n-1})\mathbf{P}(X_n | X_1, \dots, X_{n-1})$
= $\mathbf{P}(X_1, \dots, X_{n-2})\mathbf{P}(X_{n-1} | X_1, \dots, X_{n-2}) \mathbf{P}(X_n | X_1, \dots, X_{n-1})$
= ...
= $\mathbf{P}(X_1)\mathbf{P}(X_2 | X_1)\mathbf{P}(X_3 | X_1, X_2) \dots \mathbf{P}(X_n | X_1, \dots, X_{n-1})$

= $\prod_{i=1}^n \mathbf{P}(X_i | X_1, \dots, X_{i-1})$

= $\prod_{i=1}^n \mathbf{P}(X_i | \textit{Assignment of Parents}(X_i))$
- Any entry in the joint distribution can be computed and **any conditional probability** can be computed

Where are we now?

- Have a methodology for building Bayes nets
- Reduced storage necessary for probability tables
 - No longer exponential (like for joint distribution), but only exponential in the maximum number of parents of any node
- Can compute probabilities of any assignment of values to variables
 - Linear time w.r.t. number of nodes
- Can use Bayes net to determine answer to any question!



- **What could we do to compute $P(P | T, \neg S)$**

Example Bayes net query

- $P(P \mid T, \neg S)$

- Step 1: compute $P(P \wedge T \wedge \neg S)$

- Sum of all rows in the joint distribution that match $P \wedge T \wedge \neg S$ (4 joint computes)
- Use Bayes net computation method

$$P(P \wedge T \wedge \neg S) = P(P \wedge A \wedge T \wedge L \wedge \neg S)$$

$$\begin{aligned} &= (P(P \mid A) P(A) P(T \mid L) P(L \mid A \wedge \neg S) P(\neg S)) + \\ &\quad (P(P \mid \neg A) P(\neg A) P(T \mid L) P(L \mid \neg A \wedge \neg S) P(\neg S)) + \\ &\quad (P(P \mid A) P(A) P(T \mid \neg L) P(\neg L \mid A \wedge \neg S) P(\neg S)) + \\ &\quad (P(P \mid \neg A) P(\neg A) P(T \mid \neg L) P(\neg L \mid \neg A \wedge \neg S) P(\neg S)) \end{aligned}$$

$$= 0.00756 + 0.00504 + 0.18144 + 0.05376$$

$$= 0.2478$$

Example Bayes net query

- $P(P | T, \neg S)$

- Step 2: compute $P(\neg P \wedge T \wedge \neg S)$

- Sum of all rows in the joint distribution that match $\neg P \wedge T \wedge \neg S$ (4 joint computes)
- Use Bayes net computation method

$$P(\neg P \wedge T \wedge \neg S) = (\neg P \wedge A \wedge T \wedge L \wedge \neg S)$$

$$\begin{aligned} &= (P(\neg P | A) P(A) P(T | L) P(L | A \wedge \neg S) P(\neg S)) + \\ &\quad (P(\neg P | \neg A) P(\neg A) P(T | L) P(L | \neg A \wedge \neg S) P(\neg S)) + \\ &\quad (P(\neg P | A) P(A) P(T | \neg L) P(\neg L | A \wedge \neg S) P(\neg S)) + \\ &\quad (P(\neg P | \neg A) P(\neg A) P(T | \neg L) P(\neg L | \neg A \wedge \neg S) P(\neg S)) \end{aligned}$$

$$= 0.0054 + 0.01176 + 0.12096 + 0.12544$$

$$= 0.26356$$

Example Bayes net query

- $P(P \mid T, \neg S)$

- Step 3: Return
$$P(P \mid T, \neg S) = \frac{P(P \wedge T \wedge \neg S)}{P(P \wedge T \wedge \neg S) + P(\neg P \wedge T \wedge \neg S)}$$
$$= \frac{0.2478}{0.2478 + 0.26356}$$
$$= 0.4846$$

- **Inference by enumeration** algorithm

The good news

- We can do inference!
 - We can compute any conditional probability
P(Some variable | Some other variable values)

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{joint entries matching } E_1 \text{ and } E_2} P(\text{joint entry})}{\sum_{\text{joint entries matching } E_2} P(\text{joint entry})}$$

The good news

- We can do inference!
 - We can compute any conditional probability
P(Some variable | Some other variable values)

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{joint entries matching } E_1 \text{ and } E_2} P(\text{joint entry})}{\sum_{\text{joint entries matching } E_2} P(\text{joint entry})}$$

- Suppose you have m binary-valued variables in your Bayes net and expression E_2 mentions k variables.
 - **How much work is the above computation?**

The sad, bad news

- Conditional probabilities by enumerating all matching entries in the joint are expensive
 - **Exponential in the number of variables**
- Perhaps there are faster ways of querying Bayes nets?
 - In fact, if you ever manually do a Bayes net inference, you'll find there are often many tricks to save you time
 - All we have to do is program our computer to do those tricks too, right?
- Sadder and worse news:
 - **General querying of Bayes nets is NP-complete**

Naïve Bayes classifiers

- Our first example of machine learning!
- A supervised learning method
- Make (naïve) **independence assumption**
 - Can explore a simple subset of Bayesian nets s.t. it is easy to estimate the CPTs from sample data
- **Maximum likelihood estimation**
 - Given a set of correctly classified representative examples
 - **Q:** What estimates of conditional probabilities maximize the likelihood of the data that was observed?
 - **A:** The estimates that reflect the sample proportions

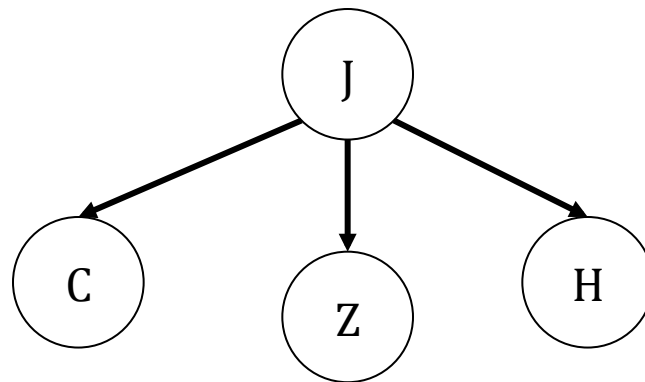
A simple Bayes net

J: Person is a junior

C: Brought coat to classroom

Z: Lives in zip code 12345

H: Saw *Harry Potter* more than once



A simple Bayes net

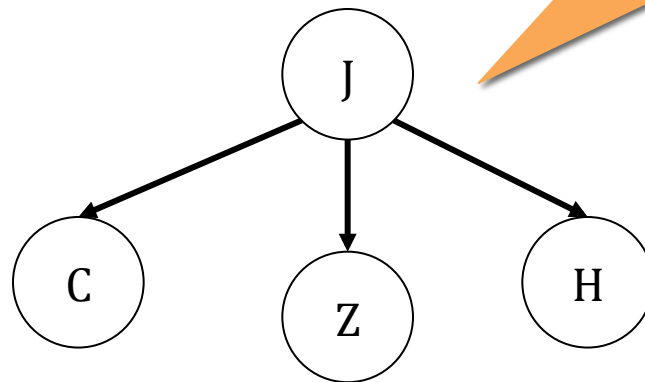
J: Person is a junior

C: Brought coat to classroom

Z: Lives in zip code

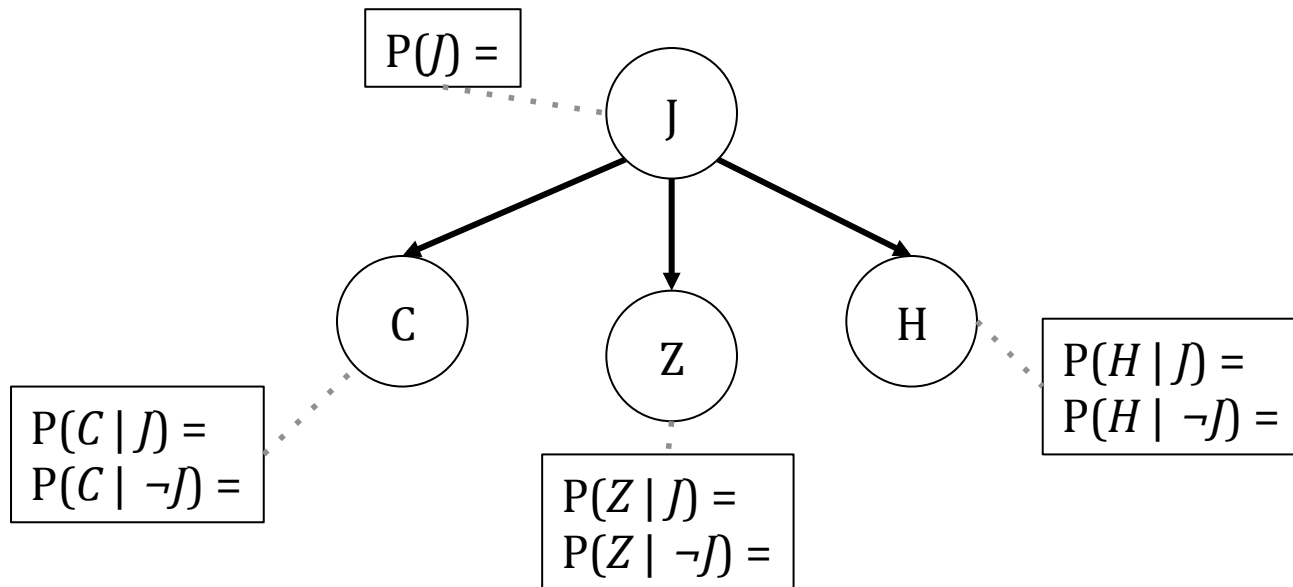
H: Saw *Harry Potter*

What parameters are stored in the CPTs of this Bayes net?



A simple Bayes net

J: Person is a junior
C: Brought coat to classroom
Z: Lives in zip code 12345
H: Saw *Harry Potter* more than once



A simple Bayes net

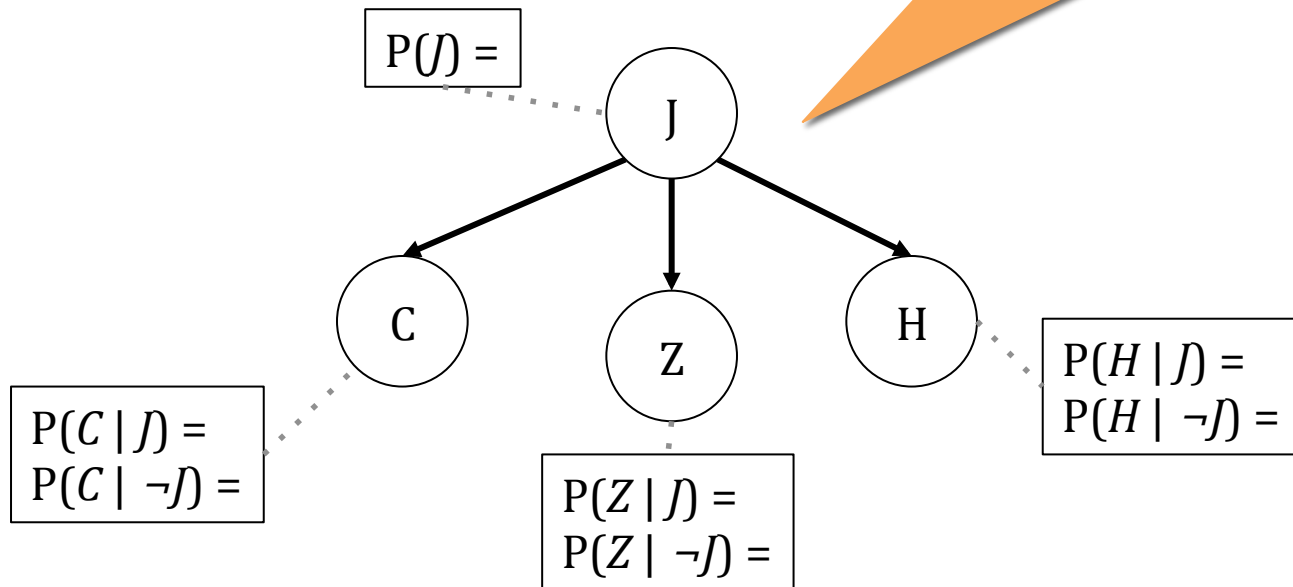
J: Person is a junior

C: Brought coat to classroom

Z: Lives in zip code

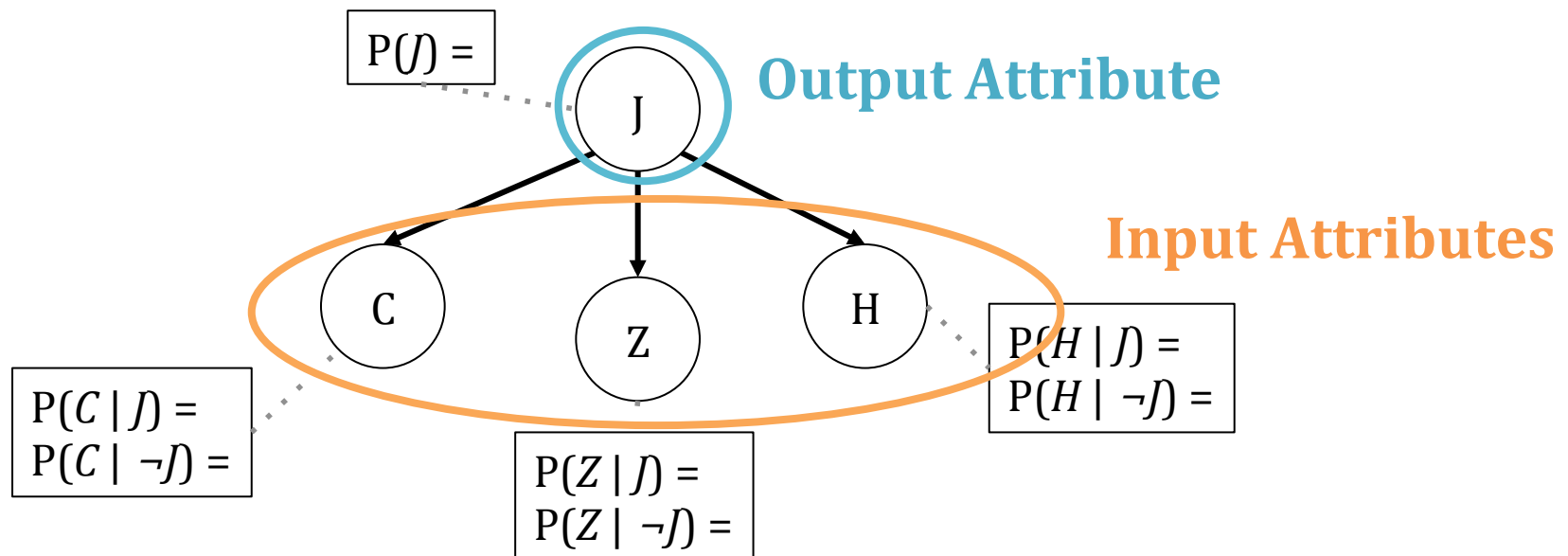
H: Saw *Harry Potter*

Suppose we had a database from 20 people. How could we use that to estimate the values in this CPT?



Naïve Bayes classifier

J : Person is a junior
 C : Brought coat to classroom
 Z : Lives in zip code 12345
 H : Saw *Harry Potter* more than once

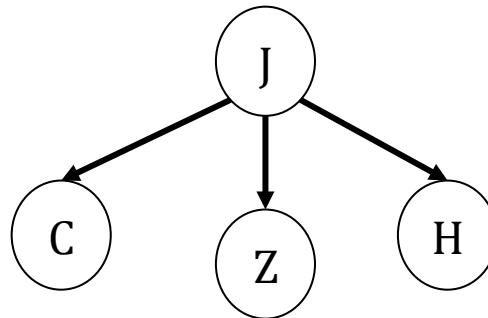


Inference with a Naïve Bayes classifier

- $P(J | C \wedge \neg Z \wedge H) = \frac{P(J \wedge C \wedge \neg Z \wedge H)}{P(C \wedge \neg Z \wedge H)}$

$$= \frac{P(J \wedge C \wedge \neg Z \wedge H)}{P(J \wedge C \wedge \neg Z \wedge H) + P(\neg J \wedge C \wedge \neg Z \wedge H)}$$

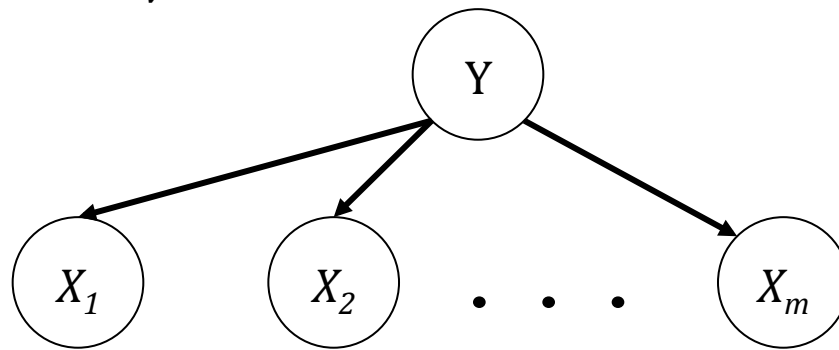
$$= \frac{P(C | J) P(\neg Z | J) P(H | J) P(J)}{\left(P(C | J) P(\neg Z | J) P(H | J) P(J) + P(C | \neg J) P(\neg Z | \neg J) P(H | \neg J) P(\neg J) \right)}$$



Naïve Bayes—the general case

- Estimate $P(Y = v)$ as a fraction of records with $Y = v$
- Estimate $P(X_i = u_i \mid Y = v)$ as fraction of $Y = v$ records that also have $X_i = u_i$
- To **predict** the Y value given observations of all the X_i values, compute:

$$Y^{predict} = \underset{v}{\operatorname{argmax}} P(Y = v \mid X_1 = u_1, \dots, X_m = u_m)$$



- Assume features X_1, \dots, X_m are independent of each other given the class Y

Computing Naïve Bayes classification

$$Y^{predict} = \underset{v}{\operatorname{argmax}} P(Y = v \mid X_1 = u_1, \dots, X_m = u_m)$$

- With **Bayes rule** we have

$$Y^{predict} = \underset{v}{\operatorname{argmax}} \frac{P(X_1 = u_1, \dots, X_m = u_m \mid Y = v) P(Y = v)}{P(X_1 = u_1, \dots, X_m = u_m)}$$

- We really only need the numerator to make a classification

$$= \underset{v}{\operatorname{argmax}} P(X_1 = u_1, \dots, X_m = u_m \mid Y = v) P(Y = v)$$

$$= \underset{v}{\operatorname{argmax}} P(X_1 = u_1 \mid Y = v) \dots P(X_m = u_m \mid Y = v) P(Y = v)$$

$$= \underset{v}{\operatorname{argmax}} P(Y = v) \prod_i P(X_i = u_i \mid Y = v)$$