

# Learning Analytics

Dr. Bowen Hui

Computer Science

University of British Columbia Okanagan

# Last Class

- Review of probability
  - Basic terminology: random variables, joint distribution
  - Conditional probability, sum-out rule, product rule
  - A few calculation examples
- All in the context of multiagent interaction
  - Inference to model our world
  - Estimate values of hidden variables using observations

# Introduced Bayes Rule

- Real world problems typically requires us to compute  $\Pr(H|e)$ 
  - Recall Asian flu example: given  $\Pr(A)$ ,  $\Pr(F)$ ,  $\Pr(F|A)$
- Bayes rule rewrites  $\Pr(H|e) \propto \Pr(e|H)\Pr(H)$

Posterior probability  $\propto$  Likelihood  $\times$  Prior probability

# Belief Perseverance

I saw a red  
fire truck!

Actually it  
was blue.

No, really,  
it was red.



Image taken from [www.fotosearch.com](http://www.fotosearch.com)



Image taken from [www.shutterstock.com](http://www.shutterstock.com)



# Changes in Representation

- Propositional logic
  - E.g.  $P = \text{John sees Mary.}$
  - E.g. If  $P$  is true then  $Q$  is also true

# Changes in Representation

- Propositional logic
  - E.g.  $P = \text{John sees Mary.}$
  - E.g. If  $P$  is true then  $Q$  is also true
- Predicate logic
  - E.g.  $\text{sees}(\text{John}, \text{Mary})$
  - E.g.  $\forall x \exists y \text{ s.t. } \text{loves}(x, y)$

# Changes in Representation

- **Propositional logic**
  - E.g.  $P = \text{John sees Mary.}$
  - E.g. If  $P$  is true then  $Q$  is also true
- **Predicate logic**
  - E.g.  $\text{sees}(\text{John}, \text{Mary})$
  - E.g.  $\forall x \exists y \text{ s.t. } \text{loves}(x, y)$
- **Bayesian inference**: Reasoning under uncertainty
  - E.g.  $\text{Pr}(\text{JohnSeesMary}) = p_1$
  - E.g.  $\text{Pr}(\text{JohnSeesMary} \wedge \text{MarySeesJohn}) = p_2$

# Probabilistic Inference

- Formally:
  - Given a prior distribution  $Pr$  over some variables (represents degrees of belief over variables)
  - Given new evidence  $E = e$  for some variable  $E$
  - Revise your degrees of belief to get the posterior distribution,  $Pr_e$

# Probabilistic Inference

- Formally:
  - Given a prior distribution  $Pr$  over some variables (represents degrees of belief over variables)
  - Given new evidence  $E = e$  for some variable  $E$
  - Revise your degrees of belief to get the posterior distribution,  $Pr_e$
- Intuition:
  - How do your degrees of belief change as a result of learning  $E = e$ ?  
(or more generally,  $E = e$ , for set  $\mathbf{E}$ )

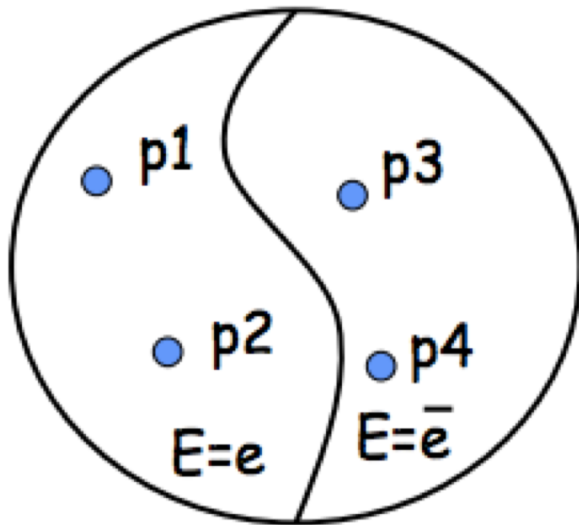
# Conditioning

- We define:

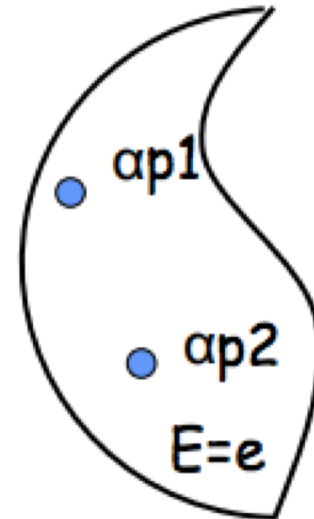
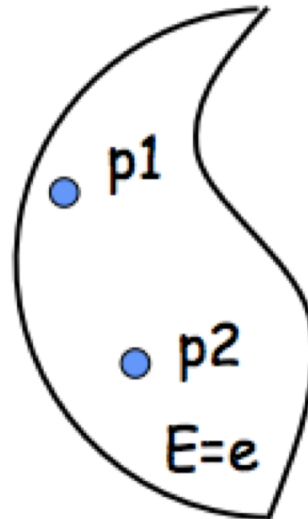
$$Pr_e(a) = Pr(a | e)$$

- That is, we produce  $Pr_e$  by conditioning the prior distribution on the observed evidence  $e$

# Semantics of Conditioning



$Pr$



$Pr_e$

$\alpha = 1/(p1+p2)$   
normalizing constant

# Computational Bottleneck

- How do we specify the full joint distribution over a set of RVs  $X_1, \dots, X_n$ ?
- Inference in this representation is frightfully slow



# Computational Bottleneck

- How do we specify the full joint distribution over a set of RVs  $X_1, \dots, X_n$ ?
  - Exponential number of possible worlds
  - These numbers are not robust/stable
  - These numbers are not natural to assess
- Inference in this representation is frightfully slow

# Computational Bottleneck

- How do we specify the full joint distribution over a set of RVs  $X_1, \dots, X_n$ ?
- Inference in this representation is frightfully slow
  - Must sum over **exponential** number of worlds to answer query  $Pr(a)$  or to condition on evidence  $e$  to determine  $Pr_e(a)$

# Recall Headache Example

	sunny		~sunny	
	cold	~cold	cold	~cold
headache	0.108	0.012	0.072	0.008
~headache	0.016	0.064	0.144	0.576

$$\Pr(\text{headache}) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2$$

$$\begin{aligned}\Pr(\text{headache} \wedge \text{cold} | \text{sunny}) &= \Pr(\text{headache} \wedge \text{cold} \wedge \text{sunny}) / \Pr(\text{sunny}) \\ &= 0.108 / (0.108 + 0.012 + 0.016 + 0.064) = 0.54\end{aligned}$$

$$\begin{aligned}\Pr(\text{headache} \wedge \text{cold} | \sim \text{sunny}) &= \Pr(\text{headache} \wedge \text{cold} \wedge \sim \text{sunny}) / \Pr(\sim \text{sunny}) \\ &= 0.072 / (0.072 + 0.008 + 0.144 + 0.576) = 0.09\end{aligned}$$

# Practical Solution

- How to avoid these two bottlenecks?
  - No solution in general
  - In practice, we will **exploit structure**
- Use independence assumptions

# Independence

- Two variables  $A$  and  $B$  are independent if knowledge of  $A$  does not change the uncertainty of  $B$  (and vice versa)
  - $\Pr(A | B) = \Pr(A)$
  - $\Pr(B | A) = \Pr(B)$
  - $\Pr(AB) = \Pr(A)\Pr(B)$

# Independence Example


- Consider: Bennett smiles and squint eyes
- If  $\Pr(\text{Smile} | \text{Squint}) = \Pr(\text{Smile})$ 
  - Chance of him smiling when he squints
  - Chance of him smiling in anyway
- And  $\Pr(\text{Squint} | \text{Smile}) = \Pr(\text{Squint})$ 
  - Chance of him squinting when he smiles
  - Chance of him squinting no matter what else he's doing
- Then Smile and Squint are independent




Image taken from iemoji.com

# What does Independence Buy Us?

- Product rule changes:


$$\Pr(ab) = \Pr(a | b)\Pr(b)$$
$$\Pr(ab) = \Pr(a)\Pr(b)$$

- Chain rule changes:


$$\Pr(abcd) = \Pr(a | bcd)\Pr(b | cd)\Pr(c | d)\Pr(d)$$
$$\Pr(abcd) = \Pr(a)\Pr(b)\Pr(c)\Pr(d)$$

# Conditional Independence

- To loosen the independence assumption, we can use conditional independence
- Two variables  $A$  and  $B$  are conditionally independent given  $C$  if:
  - $\Pr(a | b, c) = \Pr(a | c) \quad \forall a, b, c$
- Knowing the value of  $B$  does not change the prediction of  $A$  given the presence of  $C$



# Conditional Independence Example

- Consider: Want tea, pink cup, and rainy
- If  $\Pr(\text{Tea} \mid \text{Pink}, \text{Rainy}) = \Pr(\text{Tea} \mid \text{Rainy})$ 
  - Chance of wanting tea on rainy days in pink cup is the same as chance of wanting tea on rainy days in any cup
- And  $\Pr(\text{Tea} \mid \text{Pink}, \sim \text{Rainy}) = \Pr(\text{Tea} \mid \sim \text{Rainy})$   
And  $\Pr(\text{Tea} \mid \sim \text{Pink}, \text{Rainy}) = \Pr(\text{Tea} \mid \text{Rainy})$   
And  $\Pr(\text{Tea} \mid \sim \text{Pink}, \sim \text{Rainy}) = \Pr(\text{Tea} \mid \sim \text{Rainy})$   
And  $\Pr(\sim \text{Tea} \mid \text{Pink}, \text{Rainy}) = \Pr(\sim \text{Tea} \mid \text{Rainy})$   
And ...
  - Check equivalence for all other combinations
- Then Tea is independent of Pink given Rainy



Image taken from pinterest.com

# Formal Definitions

- $x$  and  $y$  are **independent** iff:

$$\Pr(x) = \Pr(x|y) \Leftrightarrow \Pr(y) = \Pr(y|x) \Leftrightarrow \Pr(xy) = \Pr(x)\Pr(y)$$

– Intuitively, learning  $y$  doesn't influence beliefs about  $x$

- $x$  and  $y$  are **conditionally independent** given  $z$  iff:

$$\Pr(x|z) = \Pr(x|yz) \Leftrightarrow \Pr(y|z) = \Pr(y|xz) \Leftrightarrow$$

$$\Pr(xy|z) = \Pr(x|z)\Pr(y|z) \Leftrightarrow \dots$$

– Intuitively, learning  $y$  doesn't influence beliefs about  $x$  **if you already know  $z$**

# What Good is Independence?

- Given (say, Boolean) variables  $X_1, \dots, X_n$  are mutually independent
- How to specify the full joint distribution  $Pr(X_1, \dots, X_n)$ ?

# What Good is Independence?

- Given (say, Boolean) variables  $X_1, \dots, X_n$  are mutually independent
- How to specify the full joint distribution  $Pr(X_1, \dots, X_n)$ ?
  - Joint is simplified as:  $\prod_{i=1}^n Pr(X_i)$
  - Can specify the full joint using only  $n$  parameters (**linear**) instead of  $2^n - 1$  (**exponential**)

# Example

- Given 4 mut. Indep. Boolean RVs:  $X_1, X_2, X_3, X_4$   
 $\Pr(x_1) = 0.4, \Pr(x_2) = 0.2, \Pr(x_3) = 0.5, \Pr(x_4) = 0.8$
- $\Pr(x_1, \sim x_2, x_3, x_4) = ?$
- $\Pr(x_1, x_2, x_3 | x_4) = ?$

# The Value of Independence

- Complete independence reduces both representation of joint distribution and inference from  $O(2^n)$  to  $O(n)$
- **Unfortunately**, complete independence is very rare
  - Most realistic domains don't exhibit this property
- **Fortunately**, most domains exhibit a fair amount of conditional independence
  - Can exploit conditional independence for representation and inference too
  - **Bayesian networks** do just this

# An Aside on Notation

- $Pr(X)$  for variable  $X$  (or set of variables) refers to the (marginal) distribution over  $X$ 
  - Distinguish from  $Pr(x)$  or  $Pr(\sim x)$  (or  $Pr(x_i)$  for non-Boolean vars) which are numbers
  - Think of  $Pr(X)$  as a function that accepts any  $x_i \in Dom(X)$  as an argument and returns  $Pr(x_i)$

# An Aside on Notation

- $Pr(X)$  for variable  $X$  (or set of variables) refers to the (marginal) distribution over  $X$ 
  - Distinguish from  $Pr(x)$  or  $Pr(\sim x)$  (or  $Pr(x_i)$  for non-Boolean vars) which are numbers
  - Think of  $Pr(X)$  as a function that accepts any  $x_i \in Dom(X)$  as an argument and returns  $Pr(x_i)$
- $Pr(X/Y)$  refers to family of conditional distributions over  $X$ , one for each  $y \in Dom(Y)$ 
  - Think of  $Pr(X/Y)$  as a function that accepts any  $x_i$  and  $y_k$  and returns  $Pr(x_i/y_k)$

Think truth tables



# Exploiting Conditional Independence

- Consider the following story:
  - If Bowen woke up too early (E), she needs caffeine (C)
  - If Bowen needs caffeine, she's likely to be grumpy (G)
  - If she is grumpy, then her lecture won't be as good (L)
  - If lecture doesn't go smoothly, then students will be disappointed (S)



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- If you learned any of E,C,G,L, would your assessment of  $\Pr(S)$  change?



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- If you learned any of E,C,G,L, would your assessment of  $\Pr(S)$  change?
  - If any of E,C,G,L are true, you would increase  $\Pr(s)$  and decrease  $\Pr(\sim s)$
  - Therefore, S is not independent of E,C,G,L



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- If you knew the value of L (true or false), would learning the value of E, C, or G influence your assessment of  $\Pr(S)$ ?



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- If you knew the value of L (true or false), would learning the value of E, C, or G influence your assessment of  $\Pr(S)$ ?
  - Influence that E, C, G has on S is mediated by L
  - E.g. Students aren't disappointed because Bowen is grumpy, it's because the lecture wasn't smooth
  - So S is independent of E, C, G, given L



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- We have: S is independent of E,C,G, given L
- Similarly:
  - L is independent of E,C given G
  - G is independent of E given C
- This translates to:
  - $\Pr(S|L,G,C,E) = \Pr(S|L)$
  - $\Pr(L|G,C,E) = \Pr(L|G)$
  - $\Pr(G|C,E) = \Pr(G|C)$
  - $\Pr(C|E)$  % doesn't simplify further
  - $\Pr(E)$  % doesn't simplify further



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- Specifying the full joint distribution  $\Pr(S,L,G,C,E)$ ?



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Conditional Independence

- Specifying the full joint distribution  $\Pr(S,L,G,C,E)$ ?
- By the chain rule:  
 $\Pr(S,L,G,C,E) = \Pr(S|L,G,C,E)\Pr(L|G,C,E)\Pr(G|C,E)\Pr(C|E)\Pr(E)$



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed



# Conditional Independence

- Specifying the full joint distribution  $\Pr(S,L,G,C,E)$ ?
- By the chain rule:  
$$\Pr(S,L,G,C,E) = \Pr(S|L,G,C,E)\Pr(L|G,C,E)\Pr(G|C,E)\Pr(C|E)\Pr(E)$$
- By our independence assumptions:  
$$\Pr(S,L,G,C,E) = \Pr(S|L)\Pr(L|G)\Pr(G|C)\Pr(C|E)\Pr(E)$$
- The full joint is specified by 5 **local conditional distributions!**



E = Woke up too early

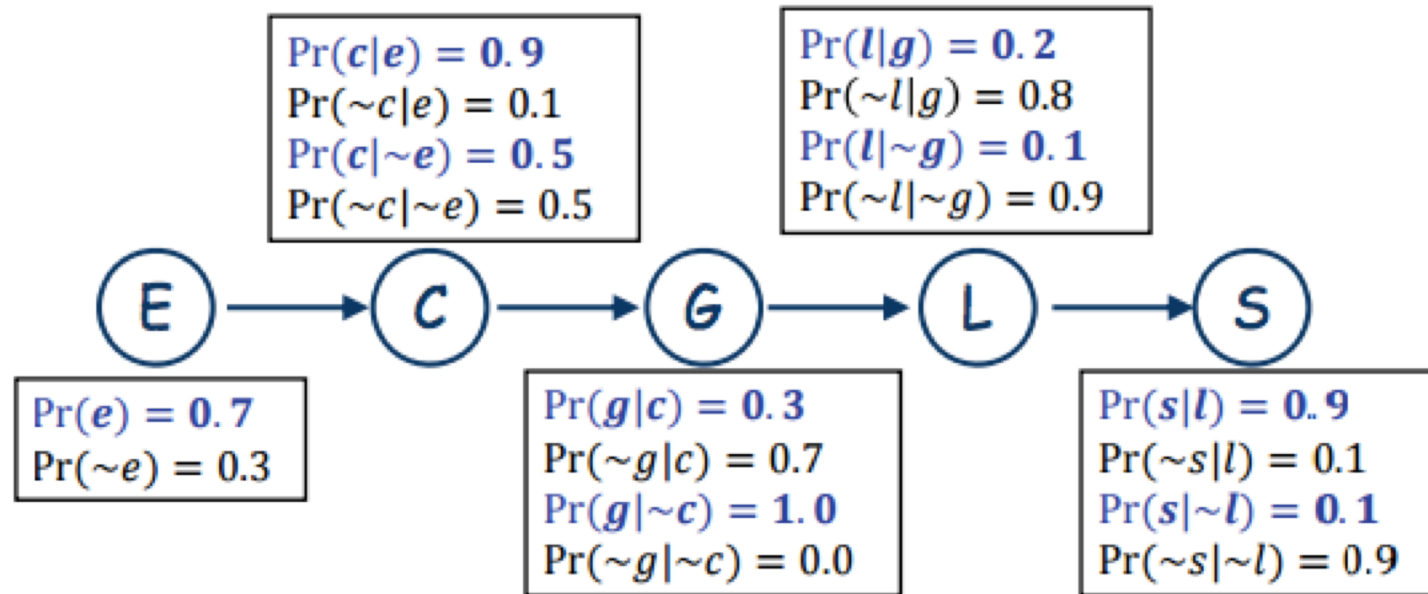
C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Example Quantification



- Specifying the joint requires only 9 parameters!
  - Instead of 31 ( $= 2^5 - 1$ ) for explicit representation

# Inference is Easy

- How to compute  $\Pr(g)$ ?



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Inference is Easy

- How to compute  $\Pr(g)$ ?
  - Apply the sum-out rule

$$P(g) = \sum_{c_i \in \text{Dom}(C)} \Pr(g | c_i) \Pr(c_i)$$
$$= \sum_{c_i \in \text{Dom}(C)} \Pr(g | c_i) \sum_{e_i \in \text{Dom}(E)} \Pr(c_i | e_i) \Pr(e_i)$$

Terms available in our local distributions!



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Inference is Easy

- Concrete example to compute  $\Pr(g)$ :



E = Woke up too early

C = need caffeine

G = gets grumpy

L = lecture not smooth

S = students disappointed

# Modeling Example

- Suppose you have a simple world with 3 variables: weather, sprinkler, and grass condition
  - If it's rainy, the grass is wet.
  - If the sprinkler is on, the grass is wet.
  - If it's cloudy, the sprinkler should be off.
- How to model these interactions?

# Modeling with a Bayes Net

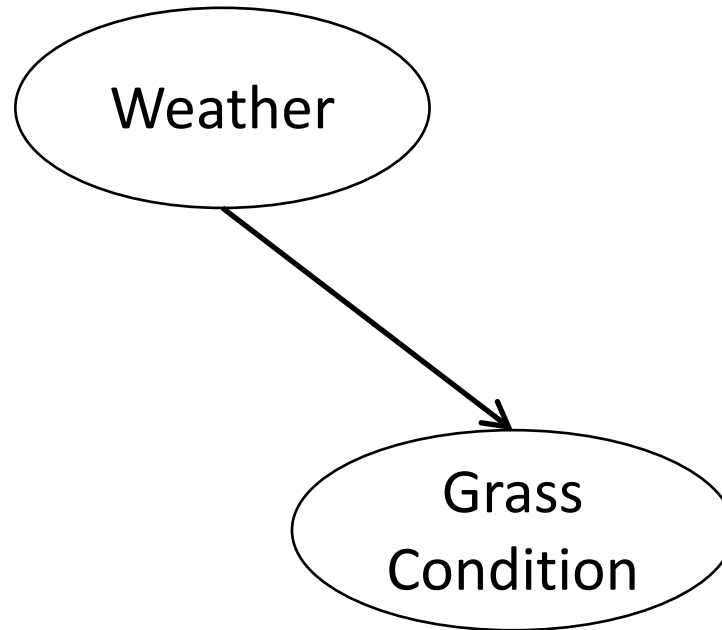
Weather = {sunny, cloudy, rainy}

Sprinkler = {on, off}

GrassCondition = {wet, dry}

# Modeling with a Bayes Net

If it's rainy, the grass is wet.



Weather = {sunny, cloudy, rainy}

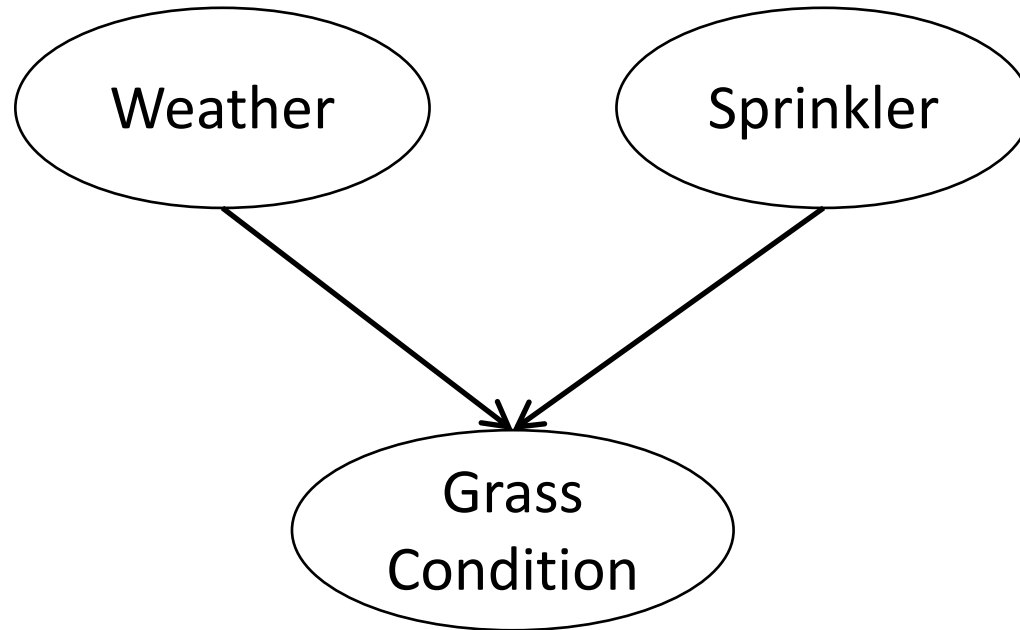
Sprinkler = {on, off}

GrassCondition = {wet, dry}



# Modeling with a Bayes Net

If the sprinkler is on, the grass is wet.



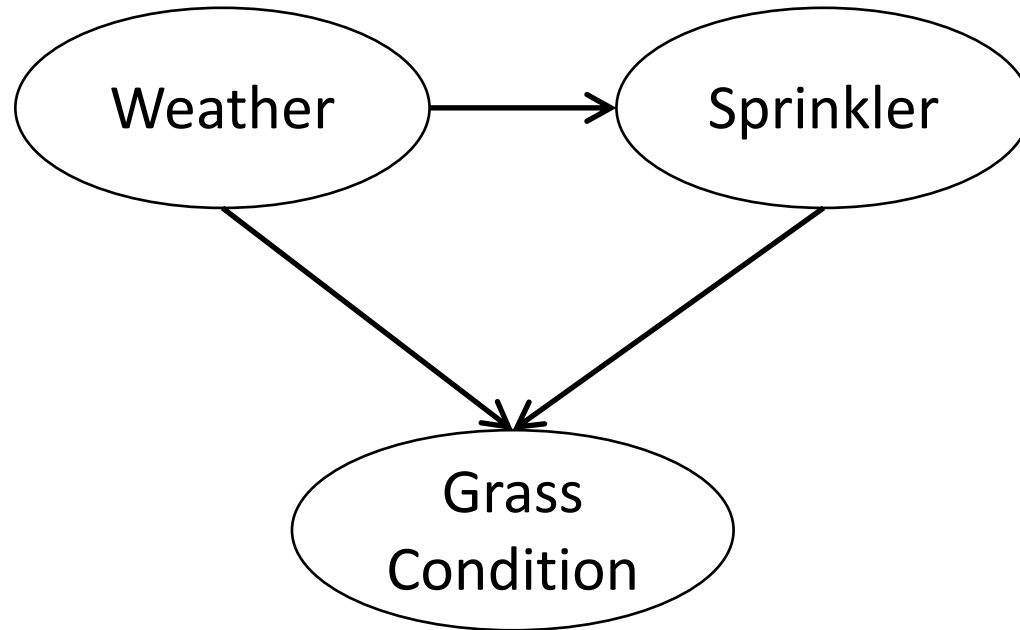
Weather = {sunny, cloudy, rainy}

Sprinkler = {on, off}

GrassCondition = {wet, dry}

# Modeling with a Bayes Net

If it's cloudy, the sprinkler should be off.

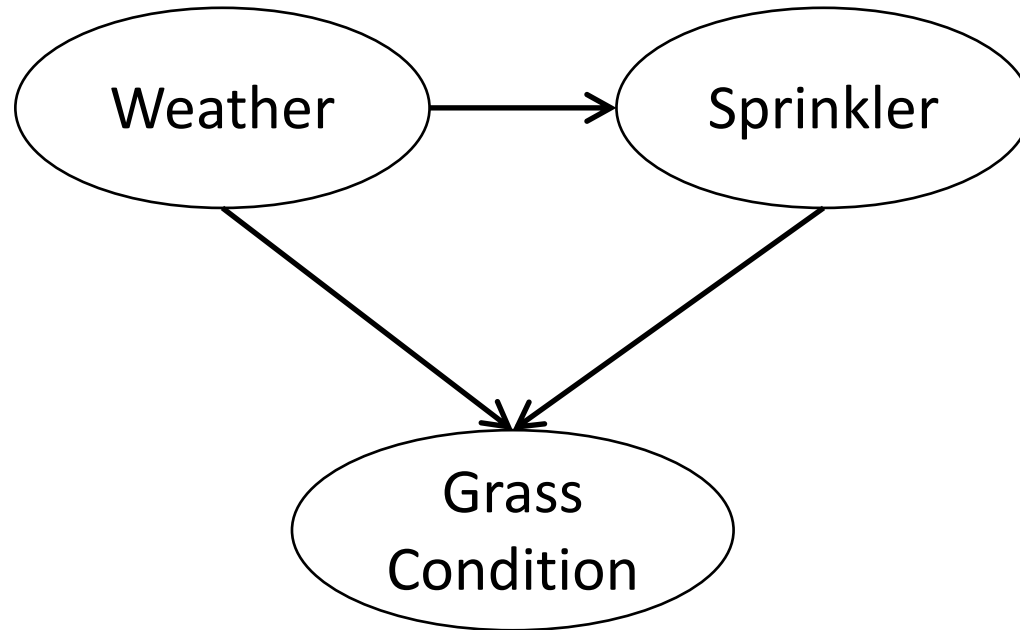


Weather = {sunny, cloudy, rainy}

Sprinkler = {on, off}

GrassCondition = {wet, dry}

# Most Popular Bayes Net Example



Weather = {sunny, cloudy, rainy}

Sprinkler = {on, off}

GrassCondition = {wet, dry}

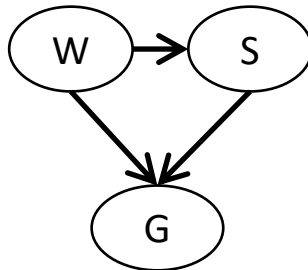
# What is a Bayes Net (BN)

- Also called **Bayesian network**, belief network
- A graphical representation of the direct dependencies over a set of variables
- Directed dependencies express the **causality** between the variables
- Each variable has an associated **conditional probability tables (CPTs)** quantifying the strength of those influences

# BN Definition

- A BN over variables  $\{X_1, X_2, \dots, X_n\}$  consists of:
  - A directed acyclic graph whose nodes are variables
  - A set of CPTs  $Pr(X_i | Parents(X_i))$  for each  $X_i$

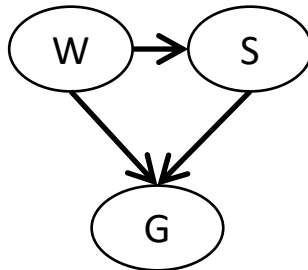
$Pr(W=\text{sunny})$	$Pr(W=\text{cloudy})$	$Pr(W=\text{rainy})$
0.6	0.3	0.1



# BN Definition

- A BN over variables  $\{X_1, X_2, \dots, X_n\}$  consists of:
  - A directed acyclic graph whose nodes are variables
  - A set of CPTs  $Pr(X_i | Parents(X_i))$  for each  $X_i$

$Pr(W=\text{sunny})$	$Pr(W=\text{cloudy})$	$Pr(W=\text{rainy})$
0.6	0.3	0.1



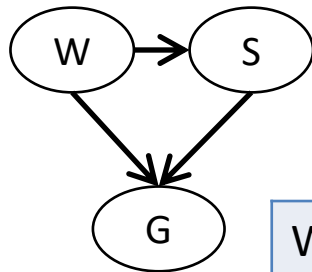
	$Pr(S=\text{on}   W)$	$Pr(S=\text{off}   W)$
W=sunny	0.1	0.9
W=cloudy	0.8	0.2
W=rainy	0.001	0.999

# BN Definition

- A BN over variables  $\{X_1, X_2, \dots, X_n\}$  consists of:
  - A directed acyclic graph whose nodes are variables
  - A set of CPTs  $Pr(X_i | Parents(X_i))$  for each  $X_i$

$Pr(W=\text{sunny})$	$Pr(W=\text{cloudy})$	$Pr(W=\text{rainy})$
0.6	0.3	0.1

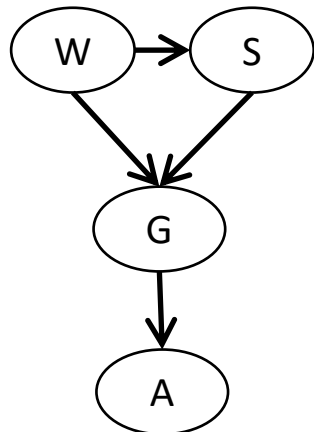
	$Pr(S=\text{on}   W)$	$Pr(S=\text{off}   W)$
W=sunny	0.1	0.9
W=cloudy	0.8	0.2
W=rainy	0.001	0.999



		$Pr(G=\text{wet}   W,S)$	$Pr(G=\text{dry}   W,S)$
W=sunny	S=on	0.9	0.1
W=sunny	S=off	0.001	0.999
W=cloudy	S=on	0.99	0.01
W=cloudy	S=off	0.2	0.8
W=rainy	S=on	1	0
W=rainy	S=off	0.9	0.1

# Key Terminology

- **Parents** of a node:  $Parents(X_i)$
- **Children** of a node
- **Descendants** of a node
- **Ancestors** of a node
- **Family**: set of nodes consisting of  $X_i$  and its parents
  - CPTs are defined over families in the BN



Parents(W) = ?

Children(S) = ?

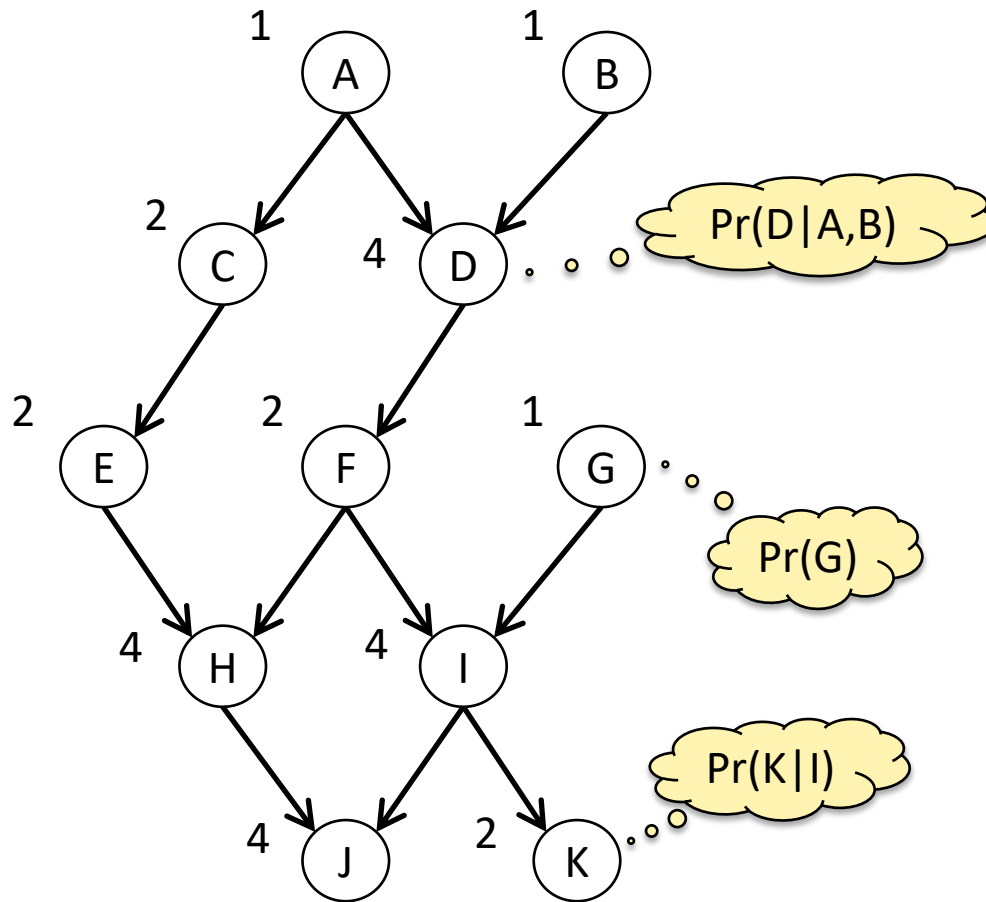
Descendants(W) = ?

Ancestors(A) = ?

Family(G) = ?



# An Example Bayes Net



- A few CPTs “shown”
- Explicit joint requires  $2^{11} - 1 = 2047$  params
- BN requires only 27 params (the number of entries for each CPT is written)

# Semantics of Bayes Nets

- The structure of the BN means:
  - Every  $X_i$  is conditionally independent of all its non-descendants given its parents
  - Intuition: your parents is the only ones who has influence on you
- Formally:

$$Pr(X_i / S \cup Par(X_i)) = Pr(X_i / Par(X_i))$$

for any subset  $S \subseteq NonDescendants(X_i)$

# Semantics of Bayes Nets

- If we ask for  $Pr(x_1, \dots, x_n)$ 
  - Assuming an ordering consistent with the network
- By the chain rule, we have:
$$\begin{aligned} & Pr(x_1, \dots, x_n) \\ &= Pr(x_n | x_{n-1}, \dots, x_1) Pr(x_{n-1} | x_{n-2}, \dots, x_1) \dots Pr(x_1) \\ &= Pr(x_n | \text{Par}(x_n)) Pr(x_{n-1} | \text{Par}(x_{n-1})) \dots Pr(x_1) \end{aligned}$$
- Thus, the joint is recoverable using the parameters (CPTs) specified in an arbitrary BN

# Key Ideas

- Main concept
  - Computational bottlenecks in computing joint probability distributions appear in representation and inference
  - Exploit independence and conditional independence
  - Computation is linear rather than exponential
- Representation:
  - Bayes net is a directed acyclic graph whose nodes are random variables with associated CPTs
  - Expresses the joint probability distribution using the product of local distributions, i.e.  $Pr(X_i | Par(X_i))$