

# Machine Learning 10-701

Tom M. Mitchell  
Machine Learning Department  
Carnegie Mellon University

October 11, 2011

## Today:

- Graphical models
- Bayes Nets:
  - Representing distributions
  - Conditional independencies
  - Simple inference
  - Simple learning

## Readings:

- Required:
- Bishop chapter 8, through 8.2

## Graphical Models

- Key Idea:
  - Conditional independence assumptions useful
  - but Naïve Bayes is extreme!
  - Graphical models express sets of conditional independence assumptions via graph structure
  - Graph structure plus associated parameters define joint probability distribution over set of variables
- Two types of graphical models:
  - Directed graphs (aka Bayesian Networks)
  - Undirected graphs (aka Markov Random Fields)

today



## Graphical Models – Why Care?

- Among most important ML developments of the decade
- Graphical models allow combining:
  - Prior knowledge in form of dependencies/independencies
  - Prior knowledge in form of priors over parameters
  - Observed training data
- Principled and ~general methods for
  - Probabilistic inference
  - Learning
- Useful in practice
  - Diagnosis, help systems, text analysis, time series models, ...

## Conditional Independence

*Definition:* X is conditionally independent of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Which we often write  $P(X|Y, Z) = P(X|Z)$

E.g.,  $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

## Marginal Independence

*Definition:*  $X$  is marginally independent of  $Y$  if

$$(\forall i, j) P(X = x_i, Y = y_j) = P(X = x_i)P(Y = y_j)$$

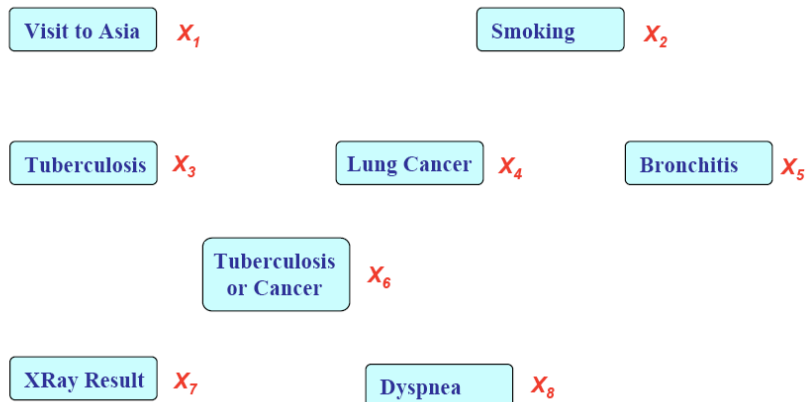
Equivalently, if

$$(\forall i, j) P(X = x_i | Y = y_j) = P(X = x_i)$$

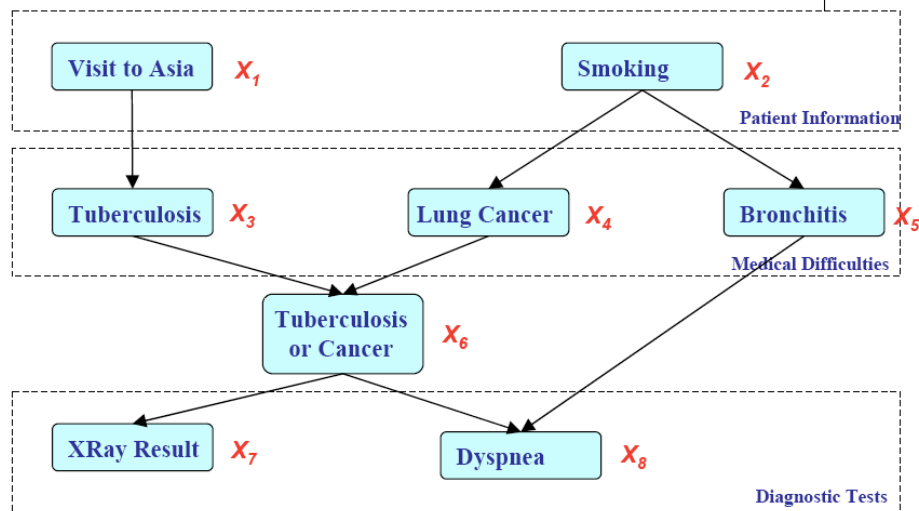
Equivalently, if

$$(\forall i, j) P(Y = y_j | X = x_i) = P(Y = y_j)$$

## Represent Joint Probability Distribution over Variables



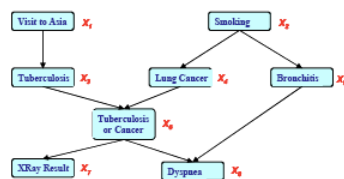
## Describe network of dependencies



Eric Xing

4

## Bayes Nets define Joint Probability Distribution in terms of this graph, plus parameters

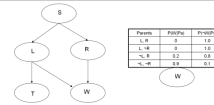


$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= P(X_1) P(X_2) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_2) \\
 &\quad P(X_6 | X_3, X_4, X_5) P(X_7 | X_6) P(X_8 | X_6)
 \end{aligned}$$

### Benefits of Bayes Nets:

- Represent the full joint distribution in fewer parameters, using prior knowledge about dependencies
- Algorithms for inference and learning

## Bayesian Networks Definition



A Bayes network represents the joint probability distribution over a collection of random variables

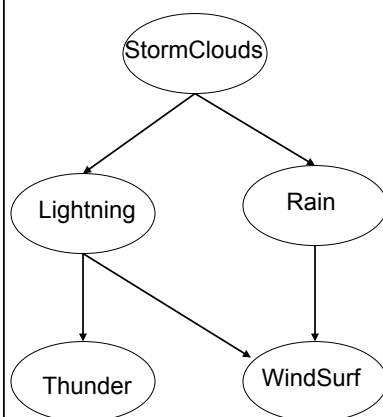
A Bayes network is a directed acyclic graph and a set of conditional probability distributions (CPD's)

- Each node denotes a random variable
- Edges denote dependencies
- For each node  $X_i$  its CPD defines  $P(X_i | Pa(X_i))$
- The joint distribution over all variables is defined to be

$$P(X_1 \dots X_n) = \prod_i P(X_i | Pa(X_i))$$

$Pa(X)$  = immediate parents of X in the graph

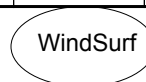
## Bayesian Network



Nodes = random variables

A conditional probability distribution (CPD) is associated with each node N, defining  $P(N | Parents(N))$

Parents	P(W Pa)	P(¬W Pa)
L, R	0	1.0
L, ¬R	0	1.0
¬L, R	0.2	0.8
¬L, ¬R	0.9	0.1



The joint distribution over all variables:

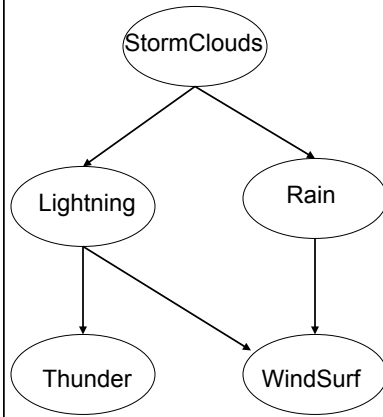
$$P(X_1 \dots X_n) = \prod_i P(X_i | Pa(X_i))$$

## Bayesian Network

What can we say about conditional independencies in a Bayes Net?

One thing is this:

Each node is conditionally independent of its non-descendants, given only its immediate parents.



Parents	$P(W Pa)$	$P(\neg W Pa)$
L, R	0	1.0
L, $\neg R$	0	1.0
$\neg L$ , R	0.2	0.8
$\neg L$ , $\neg R$	0.9	0.1

WindSurf

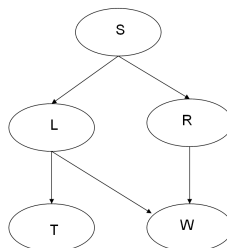
## Some helpful terminology

Parents =  $Pa(X)$  = immediate parents

Antecedents = parents, parents of parents, ...

Children = immediate children

Descendants = children, children of children, ...

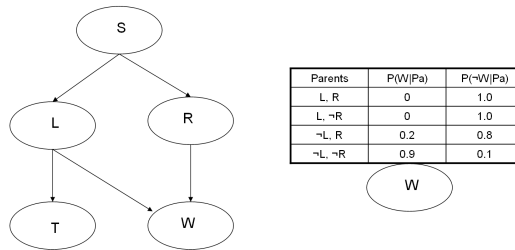


Parents	$P(W Pa)$	$P(\neg W Pa)$
L, R	0	1.0
L, $\neg R$	0	1.0
$\neg L$ , R	0.2	0.8
$\neg L$ , $\neg R$	0.9	0.1

W

## Bayesian Networks

- CPD for each node  $X_i$  describes  $P(X_i | Pa(X_i))$

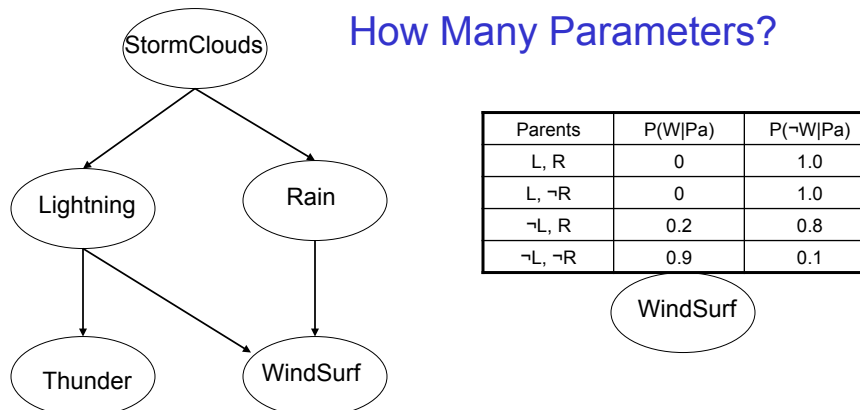


Chain rule of probability says that in general:

$$P(S, L, R, T, W) = P(S)P(L|S)P(R|S, L)P(T|S, L, R)P(W|S, L, R, T)$$

But in a Bayes net:  $P(X_1 \dots X_n) = \prod_i P(X_i | Pa(X_i))$

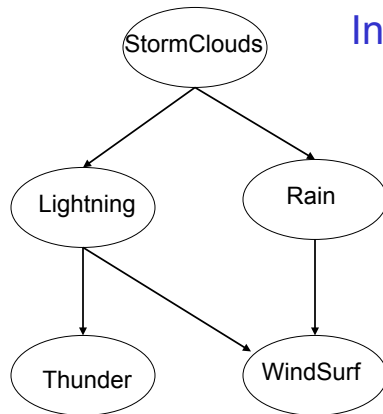
## How Many Parameters?



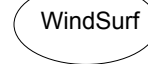
To define joint distribution in general?

To define joint distribution for this Bayes Net?

## Inference in Bayes Nets

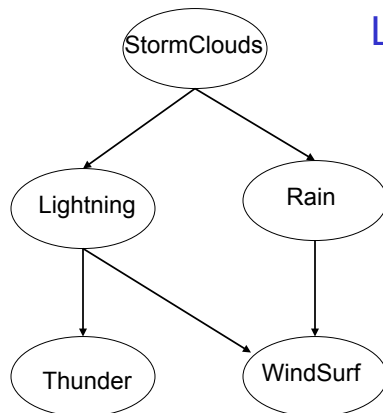


Parents	$P(W Pa)$	$P(\neg W Pa)$
L, R	0	1.0
L, $\neg R$	0	1.0
$\neg L$ , R	0.2	0.8
$\neg L$ , $\neg R$	0.9	0.1

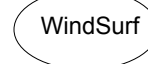


$$P(S=1, L=0, R=1, T=0, W=1) =$$

## Learning a Bayes Net



Parents	$P(W Pa)$	$P(\neg W Pa)$
L, R	0	1.0
L, $\neg R$	0	1.0
$\neg L$ , R	0.2	0.8
$\neg L$ , $\neg R$	0.9	0.1



Consider learning when graph structure is given, and data = { <s,l,r,t,w> }

What is the MLE solution? MAP?



## Algorithm for Constructing Bayes Network

- Choose an ordering over variables, e.g.,  $X_1, X_2, \dots, X_n$
- For  $i=1$  to  $n$ 
  - Add  $X_i$  to the network
  - Select parents  $Pa(X_i)$  as minimal subset of  $X_1 \dots X_{i-1}$  such that
$$P(X_i|Pa(X_i)) = P(X_i|X_1, \dots, X_{i-1})$$

Notice this choice of parents assures

$$\begin{aligned} P(X_1 \dots X_n) &= \prod_i P(X_i|X_1 \dots X_{i-1}) \quad (\text{by chain rule}) \\ &= \prod_i P(X_i|Pa(X_i)) \quad (\text{by construction}) \end{aligned}$$

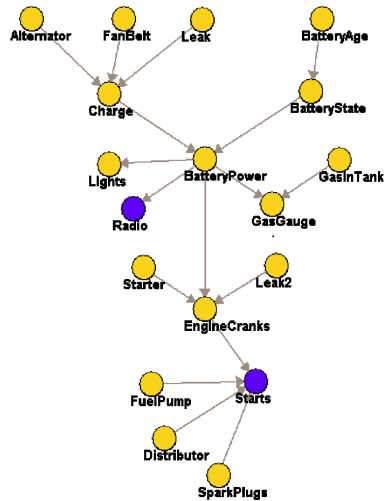
## Example

- Bird flu and Allergies both cause Nasal problems
- Nasal problems cause Sneezes and Headaches

What is the Bayes Network for  $X_1, \dots, X_4$  with NO assumed conditional independencies?

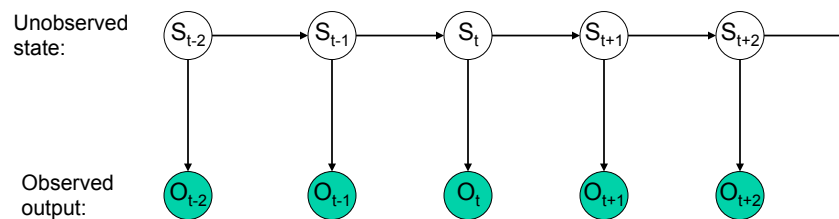
What is the Bayes Network for Naïve Bayes?

What do we do if variables are mix of discrete and real valued?



## Bayes Network for a Hidden Markov Model

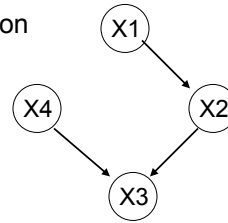
Implies the future is conditionally independent of the past, given the present



$$P(S_{t-2}, O_{t-2}, S_{t-1}, \dots, O_{t+2}) =$$

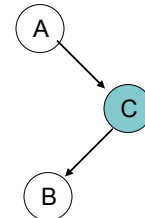
## Conditional Independence, Revisited

- We said:
  - Each node is conditionally independent of its non-descendants, given its immediate parents.
- Does this rule give us all of the conditional independence relations implied by the Bayes network?
  - No!
  - E.g.,  $X1$  and  $X4$  are conditionally indep given  $\{X2, X3\}$
  - But  $X1$  and  $X4$  not conditionally indep given  $X3$
  - For this, we need to understand D-separation



## Easy Network 1: Head to Tail

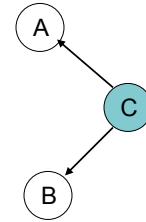
prove A cond indep of B given C?  
ie.,  $p(a,b|c) = p(a|c) p(b|c)$



let's use  $p(a,b)$  as shorthand for  $p(A=a, B=b)$

### Easy Network 2: Tail to Tail

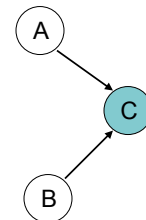
prove A cond indep of B given C? ie.,  $p(a,b|c) = p(a|c) p(b|c)$



let's use  $p(a,b)$  as shorthand for  $p(A=a, B=b)$

### Easy Network 3: Head to Head

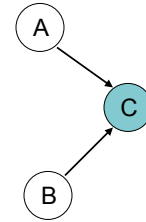
prove A cond indep of B given C? ie.,  $p(a,b|c) = p(a|c) p(b|c)$



let's use  $p(a,b)$  as shorthand for  $p(A=a, B=b)$

## Easy Network 3: Head to Head

prove A cond indep of B given C? NO!



Summary:

- $p(a,b)=p(a)p(b)$
- $p(a,b|c) \neq p(a|c)p(b|c)$

Explaining away.

e.g.,

- A=earthquake
- B=breakIn
- C=motionAlarm

**X and Y are conditionally independent given Z,  
if and only if X and Y are D-separated by Z.**

[Bishop, 8.2.2]

Suppose we have three sets of random variables: X, Y and Z

X and Y are **D-separated** by Z (and therefore conditionally indep, given Z)  
iff every path from every variable in X to every variable in Y is **blocked**

A path from variable A to variable B is **blocked** if it includes a node such that either

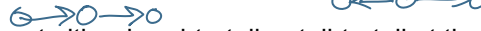


1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from every variable in X to every variable in Y is **blocked**

A path from variable A to variable B is **blocked** if it includes a node such that either



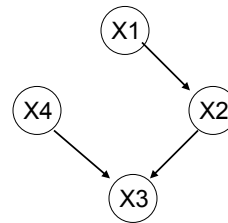
1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X1 indep of X3 given X2?

X3 indep of X1 given X2?

X4 indep of X1 given X2?



X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked** by Z

A path from variable A to variable B is **blocked** by Z if it includes a node such that either

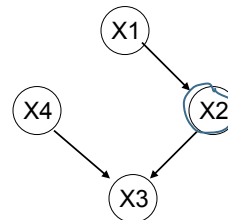
1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X4 indep of X1 given X3?

X4 indep of X1 given {X3, X2}?

X4 indep of X1 given {}?



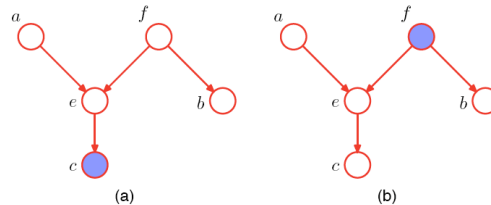
X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked**

A path from variable A to variable B is **blocked** if it includes a node such that either

1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z
2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

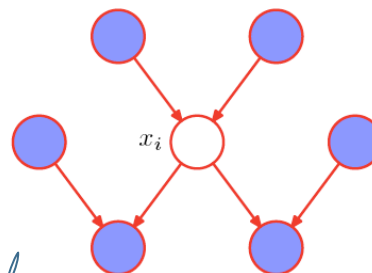
a indep of b given c?

a indep of b given f ?



## Markov Blanket

The Markov blanket of a node  $x_i$  comprises the set of parents, children and co-parents of the node. It has the property that the conditional distribution of  $x_i$ , conditioned on all the remaining variables in the graph, is dependent only on the variables in the Markov blanket.



co-parent = other side  
of  $x_i$ 's colliders

from [Bishop, 8.2]



## What You Should Know

- Bayes nets are convenient representation for encoding dependencies / conditional independence
- BN = Graph plus parameters of CPD's
  - Defines joint distribution over variables
  - Can calculate everything else from that
  - Though inference may be intractable
- Reading conditional independence relations from the graph
  - Each node is cond indep of non-descendants, given only its parents
  - D-separation
  - 'Explaining away'

See Bayes Net applet: <http://www.cs.cmu.edu/~javabayes/Home/applet.html>