# Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

February 18, 2015

Today:

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

Readings:

- 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 $G = <V, E>$
  - Graph structure plus associated parameters define *joint probability distribution over set of variables*

    Vertex: (node)
    Edge

- Two types of graphical models: 10-601
  - Directed graphs (aka Bayesian Networks)
  - Undirected graphs (aka Markov Random Fields)

# 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 <u>conditionally independent</u> 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., $P(Thunder | Rain, Lightning) = P(Thunder | Lightning)$

# Marginal Independence

*Definition*: X is <u>marginally independent</u> of Y if

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

$$= P(Y = y_j | X = x_i) P(X = x_i) =$$
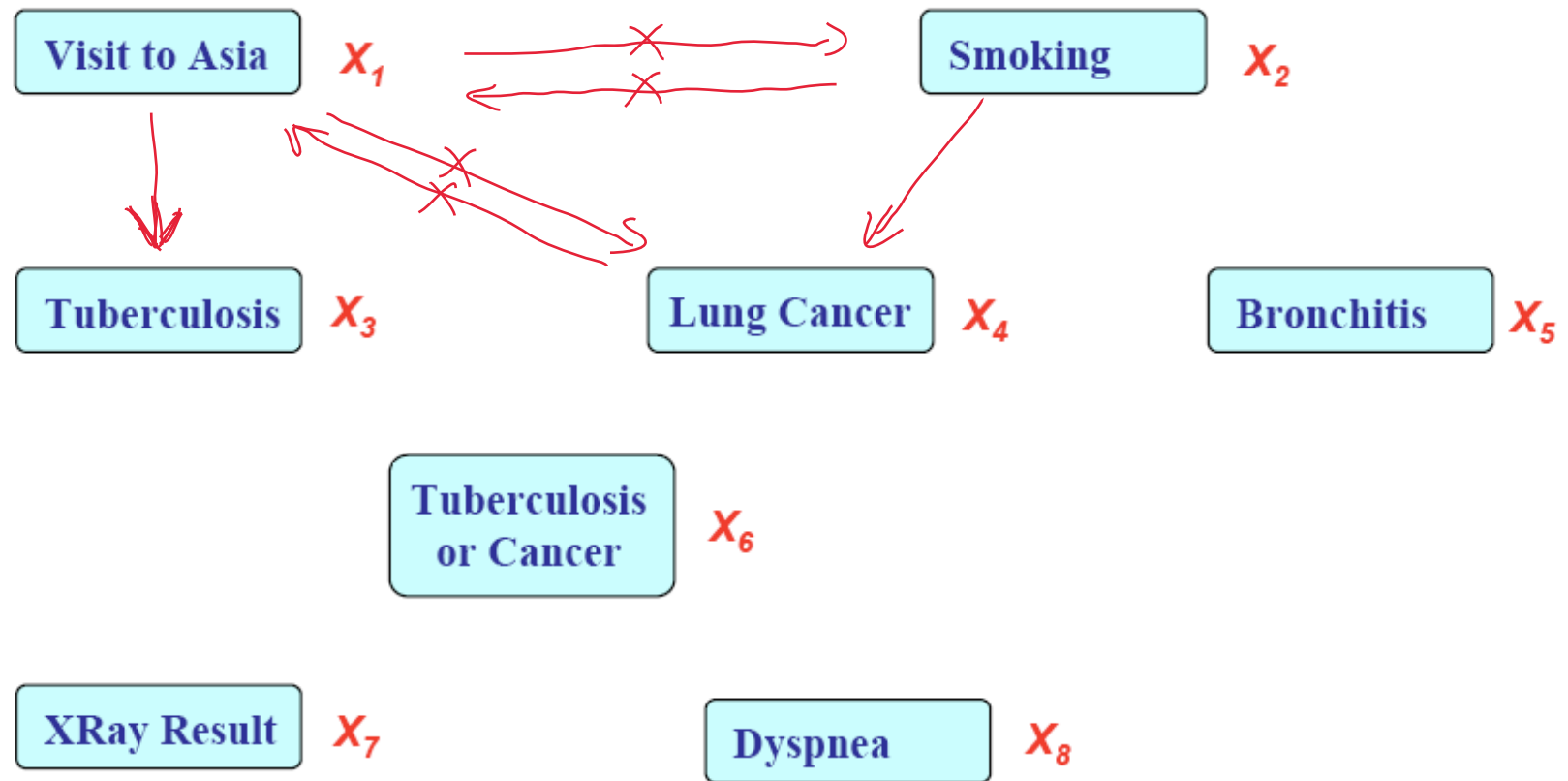
$$= P(X = x_i | Y = y_j) P(Y = y_j)$$

Equivalently, if

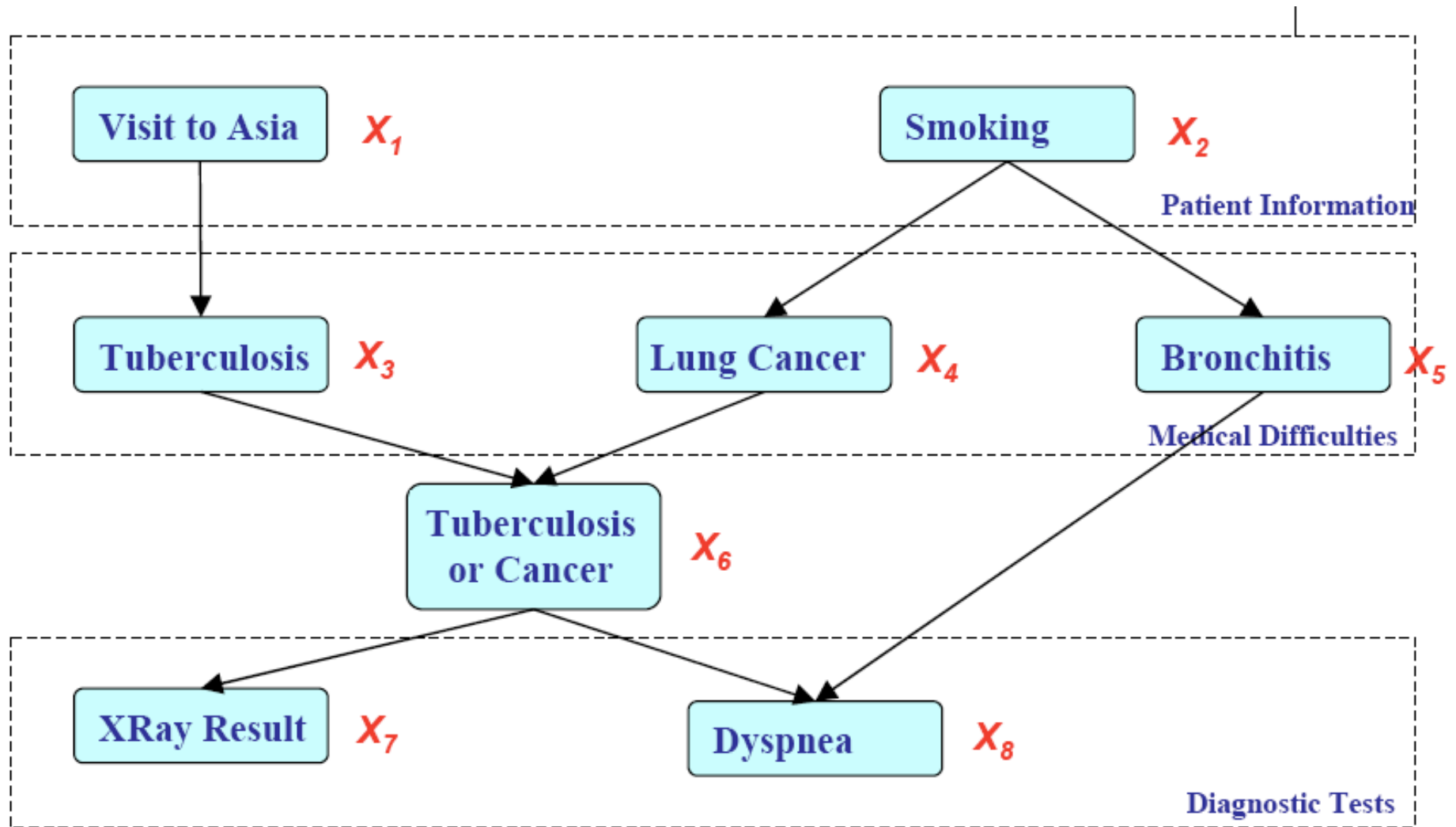$$(\forall i, j) P(X = x_i | Y = y_j) = P(X = x_i)$$

Equivalently, if

$$(\forall i, j) P(Y = y_i | X = x_j) = P(Y = y_i)$$

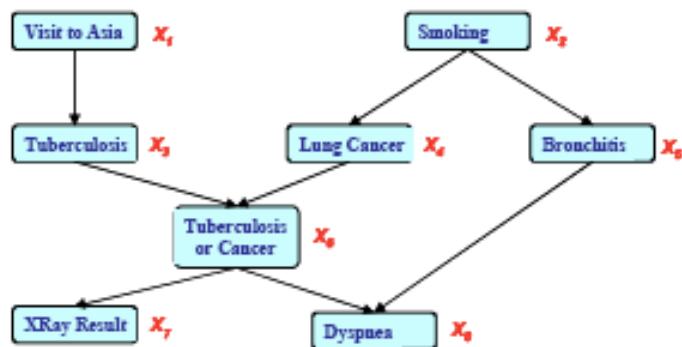# Represent Joint Probability Distribution over Variables



| Visit to Asia $X_1$ | | Smoking $X_2$ |
| Tuberculosis $X_3$ | Lung Cancer $X_4$ | Bronchitis $X_5$ |
| | Tuberculosis or Cancer $X_6$ | |
| XRay Result $X_7$ | Dyspnea $X_8$ | |

$$P(X_1, X_2, \ldots, X_8) = P(X_1)\,P(X_2)\,P(X_3|X_1)\,P(X_4|X_2)\,P(X_5|X_2)\,P(X_6|X_3,X_4)$$
$$P(X_7|X_6)\,P(X_8|X_5,X_6)$$

# Describe network of dependencies



Visit to Asia — $X_1$

Smoking — $X_2$

Patient Information

Tuberculosis — $X_3$

Lung Cancer — $X_4$

Bronchitis — $X_5$

Medical Difficulties

Tuberculosis or Cancer — $X_6$

XRay Result — $X_7$

Dyspnea — $X_8$

Diagnostic Tests

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



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= P(X_1) \, P(X_2) \, P(X_3|X_1) \, P(X_4|X_2) \, P(X_5|X_2)$$
$$P(X_6|X_3, X_4) \, P(X_7|X_6) \, P(X_8|X_5, X_6)$$

(annotations: $1$, $1$, $2$, $2$, $2$, $2^2$, $2$, $2^2$)

\#params: 18

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 (DAG) 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 $P(X_i \mid Pa(X_i))$
- The joint distribution over all variables is defined to be

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

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

*(handwritten annotations:)* (DAG)

BN —
(DAG).
Graph ← prior.
CPD ← MLE
MAP

# Bayesian Network

Nodes = random variables

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

StormClouds (S)

Lightning (L)

Rain (R)

Thunder (T)

WindSurf (W)

$Pa(W) = \{L, R\}$

| Parents | P(W\|Pa) | P(¬W\|Pa) |
|---------|----------|-----------|
| L, R | 0 $\theta_1$ | 1.0 $L\ \theta_1$ |
| L, ¬R | 0 $\theta_2$ | 1.0 $1-\theta_2$ |
| ¬L, R | 0.2 $\theta_3$ | 0.8 $1-\theta_3$ |
| ¬L, ¬R | 0.9 $\theta_4$ | 0.1 $1-\theta_4$ |

WindSurf

The joint distribution over all variables:

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

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

# Bayesian Network

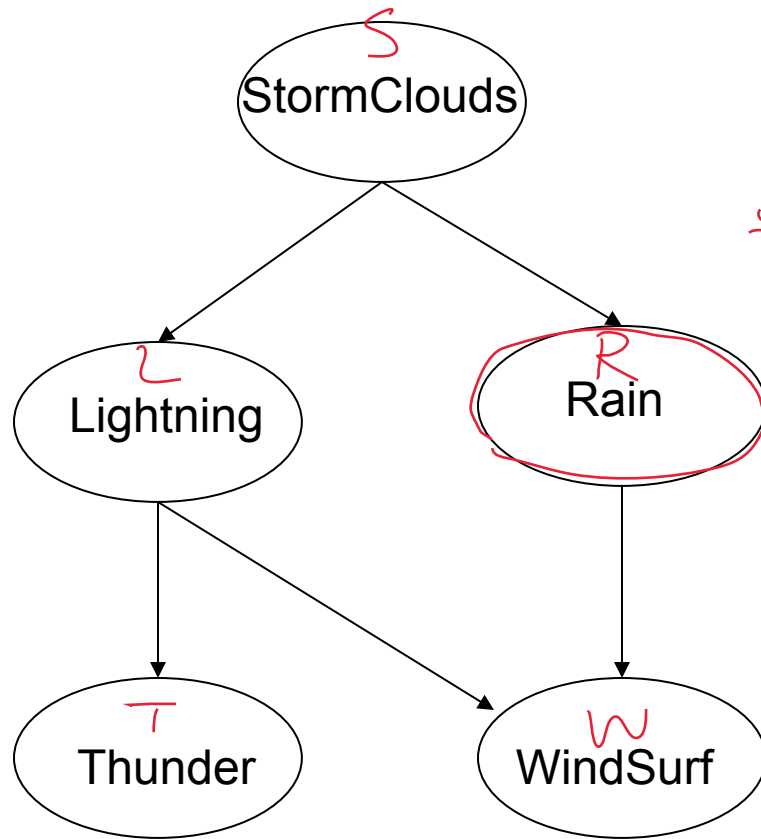$W \perp\!\!\!\perp T \mid \{L, R\}$

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

One thing is this:

( Each node is conditionally independent of its non-descendents, given only its immediate parents. )

$R \perp\!\!\!\perp \{L, T\} \mid S$

S

StormClouds

L

Lightning

R

Rain

T

Thunder

W

WindSurf

| 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 |

WindSurf

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

$\Rightarrow P(S, L, R) = P(S) \, P(L|S) \, P(R|S)$
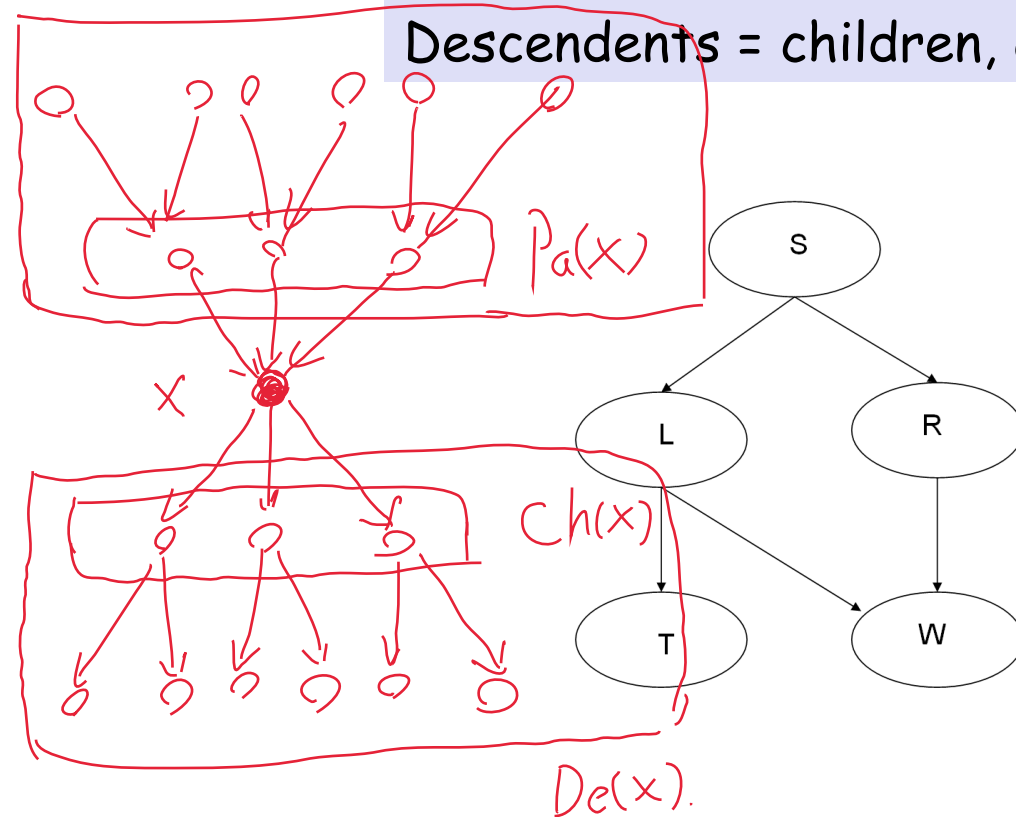
$= P(L, R|S) P(S)$

$L \perp\!\!\!\perp R \mid S$

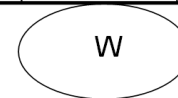# Some helpful terminology

Parents = Pa(X) = immediate parents

$An(x)$ Antecedents = parents, parents of parents, ...

Children = immediate children   $Ch(x)$

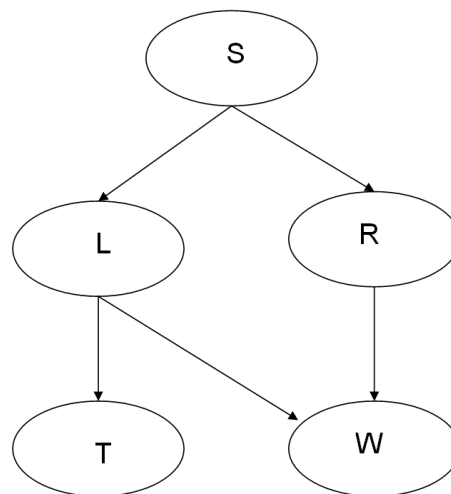Descendents = children, children of children, ...   $De(x)$



| 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      |

# Bayesian Networks

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

$$P(X \mid Y, Z) = P(X \mid Z) : \quad X \perp\!\!\!\perp Y \mid Z$$



| 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 |

$$W \perp\!\!\!\perp \{S, T\} \mid \{L, R\}$$

Chain rule of probability says that in general:

No Conditional Independence $(S \rightarrow L \rightarrow R \rightarrow T \rightarrow W)$

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

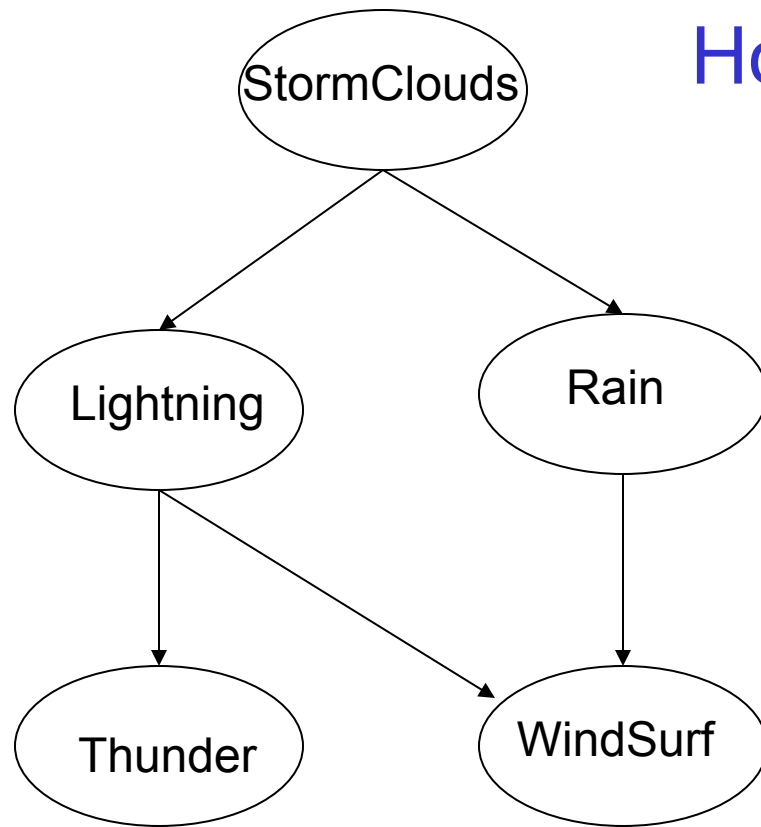$$P(R|S,L) = P(R|S) \Rightarrow R \perp\!\!\!\perp L \mid S$$

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

Conditional Independ.

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

$$P(T|S,L,R) = P(T|L) \Rightarrow T \perp\!\!\!\perp \{S, R\} \mid L$$

# How Many Parameters?

StormClouds

Lightning          Rain

Thunder          WindSurf

| 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       |

WindSurf

#nodes=n ,(boolean)

To define joint distribution in general?  #paras: $2^n - 1$  (expenontial)

1-order:  $2^1 n$

To define joint distribution for this Bayes Net?  2-order  $2^2 n$

3-order  $2^3 n$

(Linear)

# (Exact) Inference in Bayes Nets

**NP-Hard**

$P(T=n)$



| 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 |

0.2

$(P(W=1|L=0, R=1)$

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

$P(S=1, L=0, R=1) = \sum_{T, W \in \{0,1\}} P(S=1, L=0, R=1, W, T)$

$P(S=1 \mid L=0) = \dfrac{P(S=1, L=0)}{P(L=0)} = \dfrac{\sum_{R, T, W} P(S=1, L=0, R, T, W)}{\sum_{S, R, T, W} P(L=0, R, S, T, W)}$

$BN \Big\langle$ Graph (DAG) ← $prior$
$CPD \leftarrow MLE/MAP$

# Learning a Bayes Net

**StormClouds**

**Lightning**

**Rain**

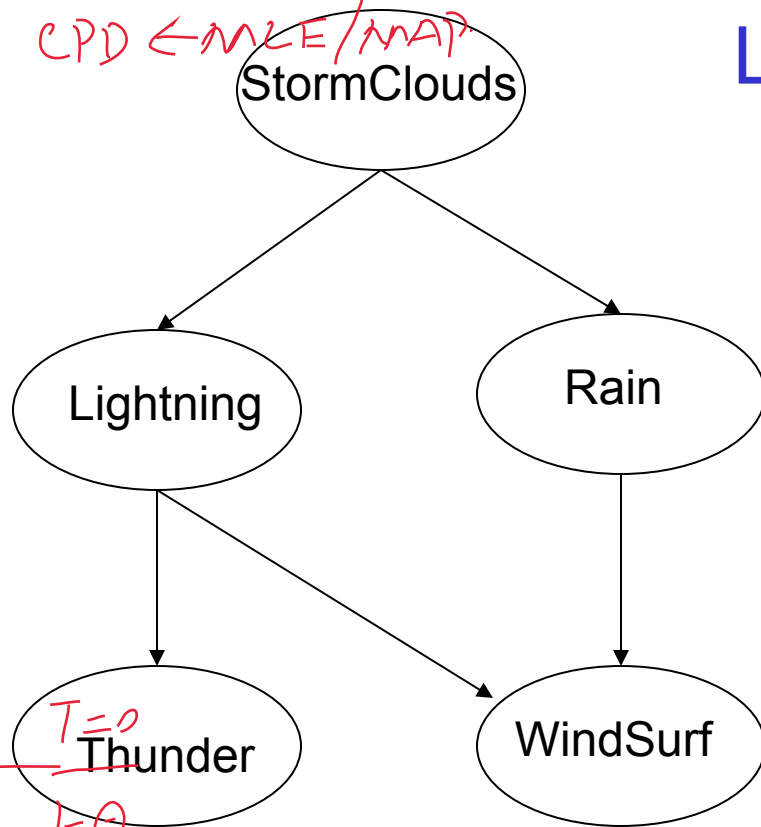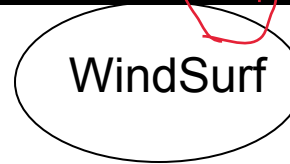| Parents | P(W\|Pa) | | P(¬W\|Pa) | |
|---|---|---|---|---|
| L, R | 0 | $\theta_1$ | 1.0 | $1-\theta_1$ |
| L, ¬R | 0 | $\theta_2$ | 1.0 | $1-\theta_2$ |
| ¬L, R | 0.2 | $\theta_3$ | 0.8 | $1-\theta_3$ |
| ¬L, ¬R | 0.9 | $\theta_4$ | 0.1 | $1-\theta_4$ |

**Thunder**

**WindSurf**

**WindSurf**

| P | T=1 | T=0 |
|---|---|---|
| L=1 | $\theta_1$ | $1-\theta_1$ |
| L=0 | $\theta_0$ | $1-\theta_0$ |

Training data:

$$D = \{ (s_j, \ell_j, r_j, t_j, w_j) \}_{j=1}^{M} \quad (i.i.d.)$$

$$\ell_1 \qquad x_j$$

$\theta = \{\theta_1, \theta_0\}$

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

What is the MLE solution?  MAP?

$$P(t | \theta_1, \theta_0) = \theta_1^{LT}(1-\theta_1)^{L(1-T)} \times$$

$$\theta_0^{(1-L)T}(1-\theta_0)^{(1-L)(1-T)}$$

$$\max_\theta L(\theta) = P(D|\theta) \Rightarrow \prod_{j=1}^{M} P(t_j | \theta_1, \theta_0) \overset{t_j}{}$$

$$lng-\text{Likelihood} \cdot \ell(\theta) \Rightarrow \sum_{j=1}^{M} \ln P(t_j | \theta_1, \theta_0) \begin{cases} \dfrac{\partial \ell(\theta)}{\partial \theta_1} = 0, & \Rightarrow \hat{\theta}_1 = \dfrac{|D_{T=1, L=1}|}{|D_{T=1}|} \\[4mm] \dfrac{\partial \ell(\theta)}{\partial \theta_0} = 0 & \Rightarrow \hat{\theta}_0 = \dfrac{|D_{T=1, L=0}|}{|D_{T=1}|} \end{cases}$$

# Algorithm for Constructing Bayes Network

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

  $$P(X_i|Pa(X_i)) = P(X_i|X_1, \ldots, X_{i-1})$$

  $$\left( X_i \perp\!\!\!\perp \overline{Pa(X_i)} \mid Pa(X_i) \right)$$

  $$P\left(X_i \mid Pa(X_i), \overline{Pa(X_i)}\right)$$
  $$= P\left(X_i \mid Pa(X_i)\right)$$
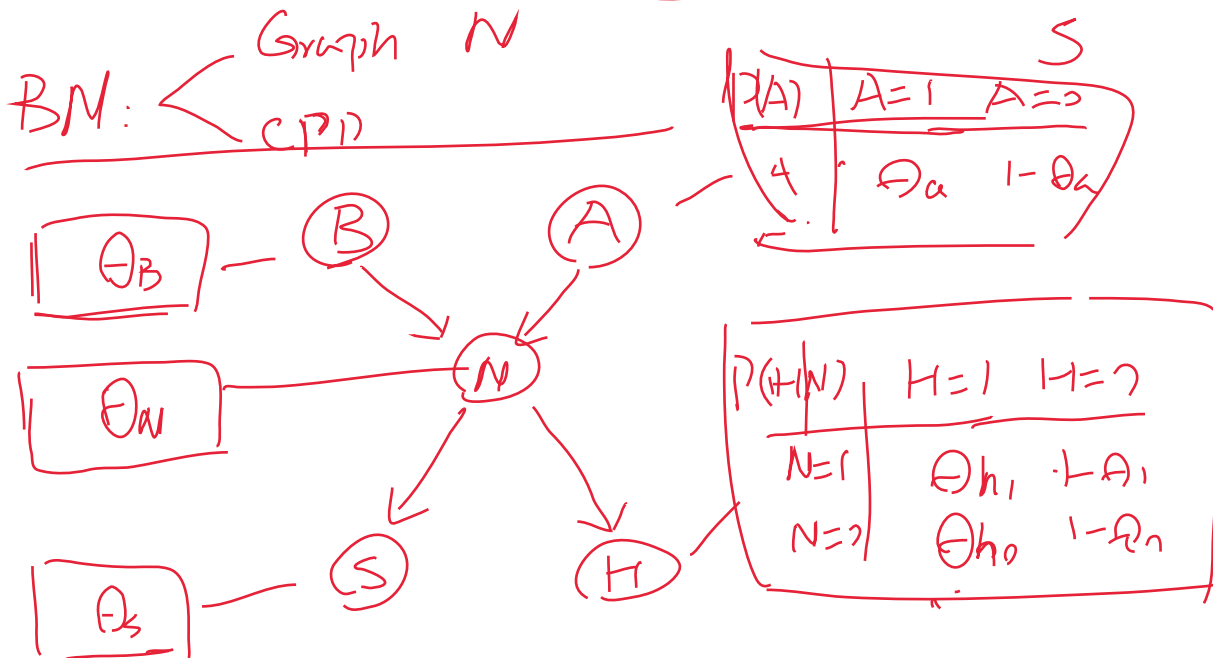
  $$Pa(X_i), \overline{Pa(X_i)}$$

  Notice this choice of parents assures

  $$P(X_1 \ldots X_n) = \prod_i P(X_i|X_1 \ldots X_{i-1}) \quad \text{(by chain rule)}$$

  $$= \prod_i P(X_i|Pa(X_i)) \quad \text{(by construction)}$$

# Example



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

B   A   N

Graph N

BN: { Graph N
       CPD

S   H

BM:

$\theta_B$

$\theta_{av}$

$\theta_s$

B   A

N

S   H

| P(A) | A=1 | A=0 |
|------|-----|-----|
| 4 | $\theta_a$ | $1-\theta_a$ |

| P(H|N) | H=1 | H=0 |
|--------|-----|-----|
| N=1 | $\theta_{h_1}$ | $1-\theta_{h_1}$ |
| N=0 | $\theta_{h_0}$ | $1-\theta_{h_0}$ |

Training data

$$D = \{ (b_j, a_j, n_j, s_j, h_j) \}_{j=1}^{m}$$
$$\underbrace{\qquad}_{x_j}$$

$$\ell(\theta) = \ln P(D | \theta)$$

$$= \sum_{j=1}^{m} \ln P(x_j | \theta_a, \theta_b, \theta_n, \theta_s, \theta_H)$$

$$= \sum_{j=1}^{m} \ln P(x_j | \theta_a) + \sum_{j=1}^{m} \ln P(x_j | \theta_b) +$$
$$\sum_{j=1}^{m} \ln P(x_j | \theta_n) + \sum_{j=1}^{m} \ln P(x_j | \theta_s) +$$
$$\sum_{j=1}^{m} \ln P(x_j | \theta_h)$$

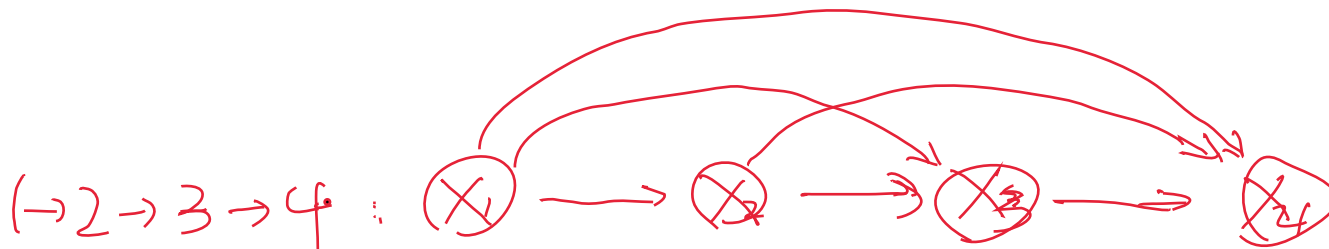$$\frac{\partial \ell(\theta)}{\partial \theta} = 0 \quad \longleftarrow$$

# What is the Bayes Network for X1,…X4 with NO assumed conditional independencies?

$$P(X_1, X_2, X_3, X_4) = P(X_1) \, P(X_2 | X_1) \, P(X_3 | X_1, X_2) \, P(X_4 | X_1, X_2, X_3)$$

$1 \rightarrow 2 \rightarrow 3 \rightarrow 4$

$1 \rightarrow 2 \rightarrow 4 \rightarrow 3 \implies P(X_1) \, P(X_2 | X_1) \, P(X_4 | X_1, X_2) \, P(X_3 | X_1, X_2, X_4)$

$4!$ 

(equivalent but not unique)

$\vdots$

$4 \rightarrow 3 \rightarrow 2 \rightarrow 1 \implies P(X_4) \, P(X_3 | X_4) \, P(X_2 | X_3, X_4) \, P(X_1 | X_2, X_3, X_4)$

$1 \rightarrow 2 \rightarrow 3 \rightarrow 4$ :

# What is the Bayes Network for Naïve Bayes?

$$(x_i \perp\!\!\!\perp x_j \mid Y, \forall i \neq j)$$

| $P(Y)$ | $Y=1$ | $Y=0$ |
|---|---|---|
| | $\pi$ | $1-\pi$ |

| $P(x_1 \mid Y)$ | $X_1=1$ | $X_1=0$ |
|---|---|---|
| $Y=0$ | $\theta_0$ | $1-\theta_0$ |
| $Y=1$ | $\theta_1$ | $1-\theta_1$ |

$$P(X_1, X_2, \ldots, X_n, Y) \cdot$$
$$= P(Y) \, P(X_1 \mid Y) \, P(X_2 \mid Y) \cdots P(X_n \mid Y)$$
$$= P(Y) \prod_{i=1}^{n} P(X_n \mid Y)$$

$$D = \{(x_j, y_j)\}_{j=1}^{M}$$

$$\ell(\theta; \pi) = \ell(\theta) + \ell(\pi)$$

$$\text{MLE} \begin{cases} \dfrac{\partial \ell(\theta)}{\partial \theta} = 0 \Rightarrow \hat{\theta}_0 = \dfrac{|D_{Y=0, X_1=1}|}{|D_{Y=0}|} \,, \quad \hat{\theta}_1 = \dfrac{|D_{Y=1, X_1=1}|}{|D_{Y=1}|} \\[4ex] \dfrac{\partial \ell(\pi)}{\partial \pi} = 0 \Rightarrow \hat{\pi} = \dfrac{|D_{Y=1}|}{|D|} \end{cases}$$

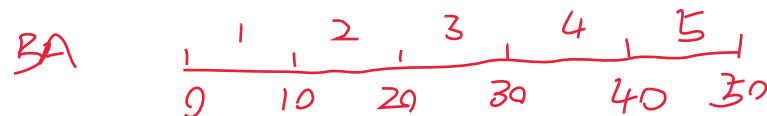# What do we do if variables are mix of discrete and real valued?



Alternator   FanBelt   Leak   BatteryAge ← continuos

Charge   BatteryState ← discrete

Lights   BatteryPower   GasInTank

Radio   GasGauge

Starter   Leak2

EngineCranks

FuelPump   Starts

Distributor

SparkPlugs

| $P(BS\mid BA)$ | $BS=0$ | $BS=1$ |
|---|---|---|
| $BA=0.1$ | $\theta_1$ | |
| $BA=0.15$ | $\theta_2$ | |
| $BA=0.2$ | $\theta_2$ | |
| $\vdots$ | $\vdots$ | |

(infinite rows ⟶ intractable)

Solutions

① Discretization

$BA \in \{1, 2, 3, 4, 5\}$

BA

|  1  |  2  |  3  |  4  |  5  |
|---|---|---|---|---|
| 0 | 10 | 20 | 30 | 40 | 50 |

② Parameterized distribution.

$$P(BS\mid BA) = G(\beta \cdot BA + \beta_0) = \frac{1}{1+e^{-(\beta BA + \beta_0)}}$$

• CPD of BS can be calculated. by $\beta$ and $\beta_0$