# 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
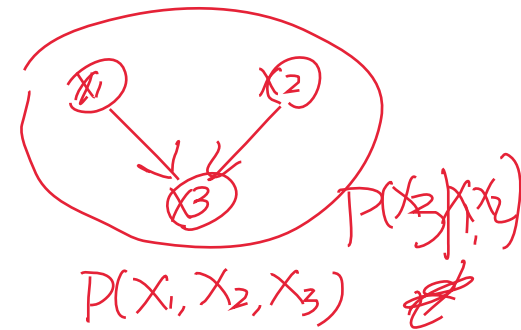  - Simple inference
  - Simple learning

## 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
  - 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)

10-601

# 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   $P(X|Y, Z) = P(X|Z)$

$P(Y|Z)$

$$P(X, Y | Z) = P(X|Z) P(Y|Z)$$

$$P(X|Z) = \sum_{i=1}^{n} P(X_i | Z) \quad \longleftarrow \quad \text{Naive Bayes}$$

E.g.,  $P(Thunder | Rain, Lightning) = P(Thunder | Lightning)$

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

# Marginal Independence

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

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

$P(X, Y) = P(X) \cdot P(Y)$

$P(X) = \prod_{i=1}^{n} P(X_i)$

$P(X|Y) P(Y)$

$\implies$ $P(X|Y) = P(X)$
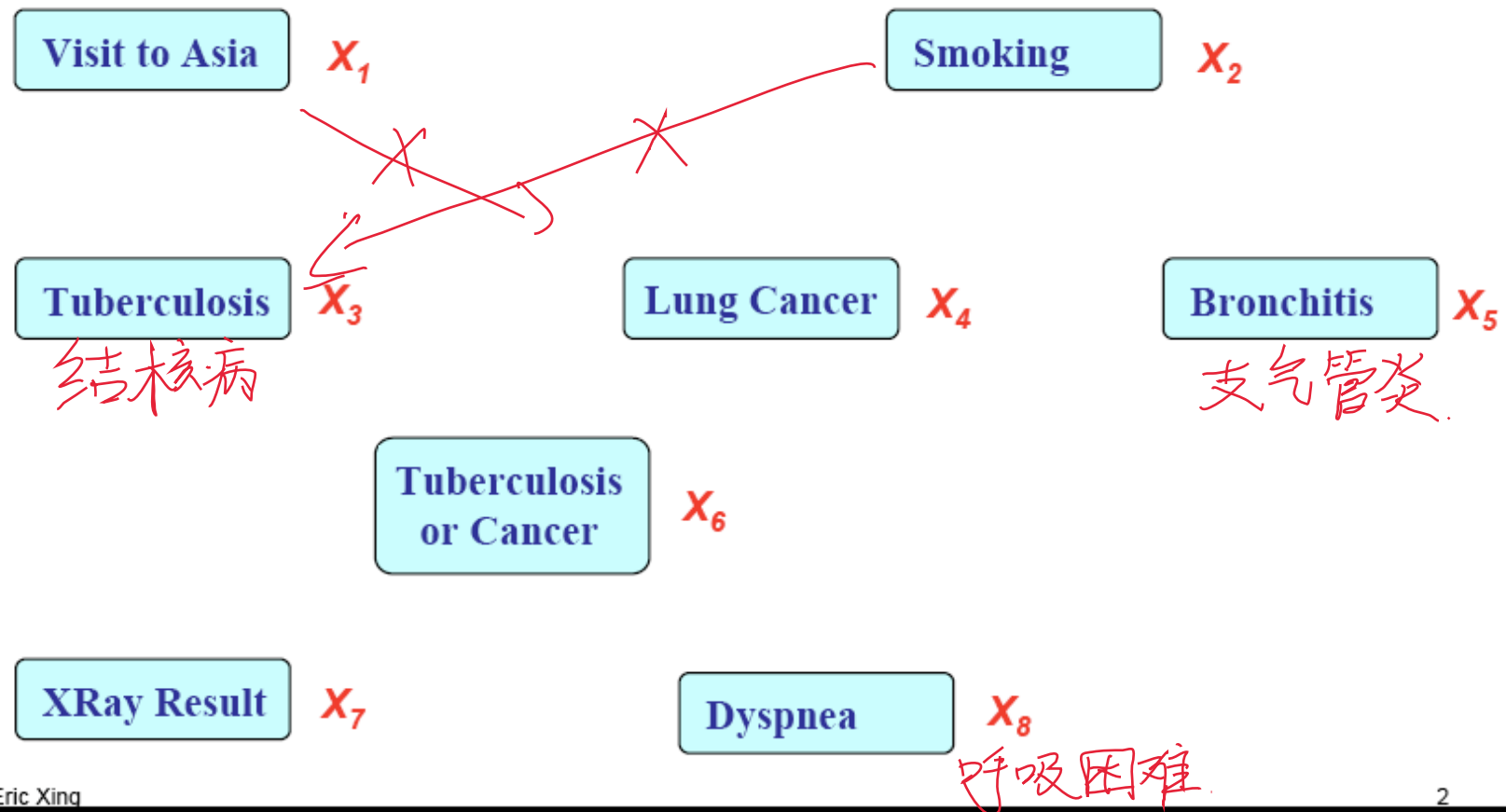
Equivalently, if

$P(Y|X) P(X)$

$\implies$ $P(Y|X) = P(Y)$

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

Equivalently, if

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

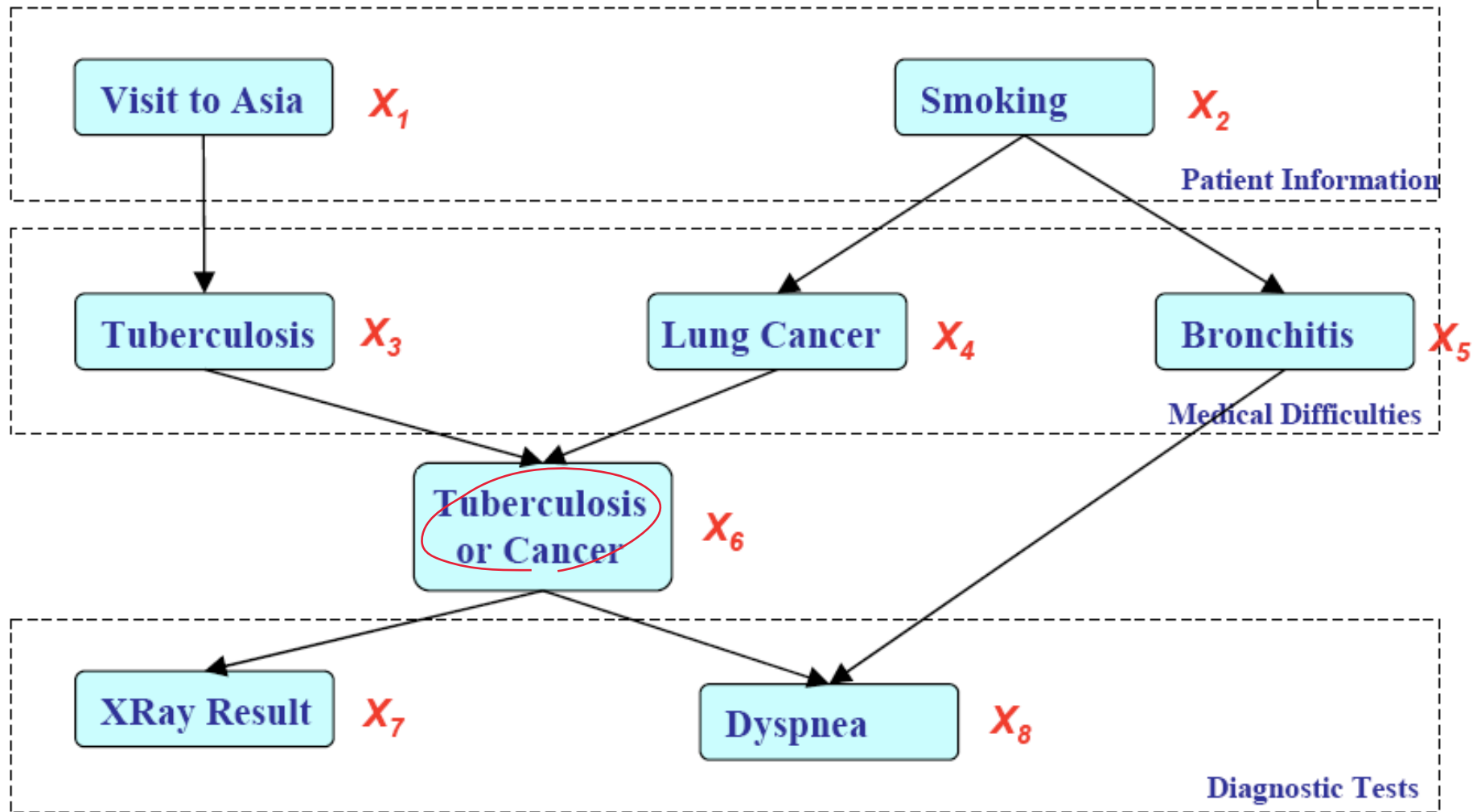# Represent Joint Probability Distribution over Variables

| | | |
|---|---|---|
| **Visit to Asia** $X_1$ | | **Smoking** $X_2$ |

| | | |
|---|---|---|
| **Tuberculosis** $X_3$ | **Lung Cancer** $X_4$ | **Bronchitis** $X_5$ |

结核病                                            支气管炎

**Tuberculosis or Cancer** $X_6$

| | |
|---|---|
| **XRay Result** $X_7$ | **Dyspnea** $X_8$ |

呼吸困难

$$\boxed{P(X_1, X_2, \ldots, X_8)}$$

$$P(X_8 \mid X_1, X_2, \ldots, X_7)$$

$$= \frac{P(X_8, X_1, X_2, \ldots, X_7)}{P(X_1, X_2, \ldots, X_7)}$$

$$P(X_1, X_2, \ldots, X_8) = P(X_1)\,P(X_2 \mid X_1) \cdots \quad P(X_8 \mid X_1, X_2, \ldots, X_7)$$

$(1 \rightarrow 2 \rightarrow \cdots \rightarrow 8)$

# 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

$P(X_2|X_1) = P(X_2)$

$X_1 \perp\!\!\!\perp X_2$

| | | |
|---|---|---|
| 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 $X_8$ |

$P(X_5|X_2)$

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

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

$P(X_1, X_2, \ldots, X_8) = P(X_1)\,P(X_2|X_1)\,P(X_3|X_1, X_2) \cdots P(X_8|X_1, \ldots, X_7)$

$1 \rightarrow 2 \rightarrow 3 \rightarrow \cdots \rightarrow 8$

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

**Benefits of Bayes Nets:**

- Represent the full joint distribution in fewer parameters, using prior knowledge about dependencies

$X_2 \perp\!\!\!\perp X_3 \mid X_1$

- Algorithms for inference and learning

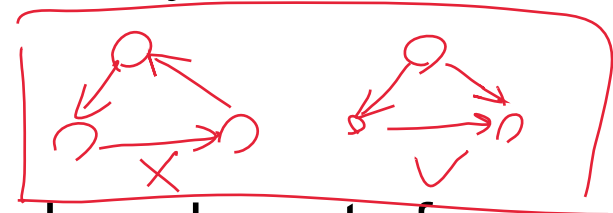$P(X_8|X_1, X_2, \ldots, X_7)$

$P(X_8|X_5, X_6)$

# Bayesian Networks Definition

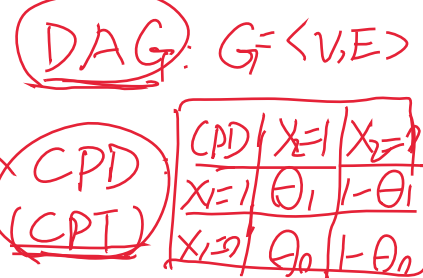A Bayes network represents the joint probability distribution over a collection of random variables

*BN*     *(DAG)*

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

*DAG: $G = \langle V, E \rangle$*

*BN*

*CPD (CPT)*

| CPD | $X_2 = 1$ | $X_2 = ?$ |
|-----|-----------|-----------|
| $X_1 = 1$ | $\theta_1$ | $1 - \theta_1$ |
| $X_1 = ?$ | $\theta_n$ | $1 - \theta_n$ |

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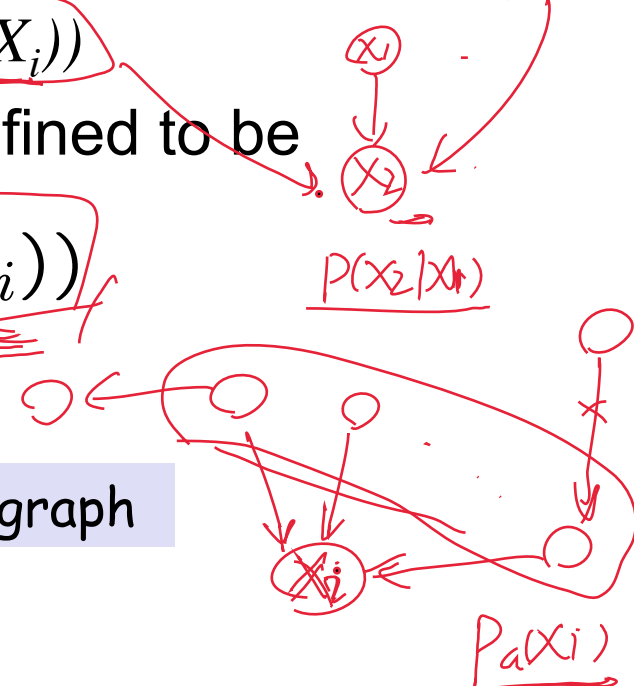- The joint distribution over all variables is defined to be

$P(X_2 \mid X_1)$

$$BN \rightarrow P(X_1 \ldots X_n) = \prod_i P(X_i \mid Pa(X_i))$$

$= P(X_1) \cdot P(X_2 \mid X_1) \cdots P(X_n \mid X_1, \ldots, X_{n-1})$

$= \prod_{i=1}^{n} P(X_i \mid Pa(X_i))$

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

$Pa(X_i)$

# Bayesian Network

*DAG*
*CPD*

*boolean*

*S*

StormClouds

*CPD | S=1  S=0*
*P(S) | θ  1-θ*

*L*

Lightning

*R*

Rain

*T*

Thunder

*W*

WindSurf

*P(W|L,R)*  *2²*

WindSurf
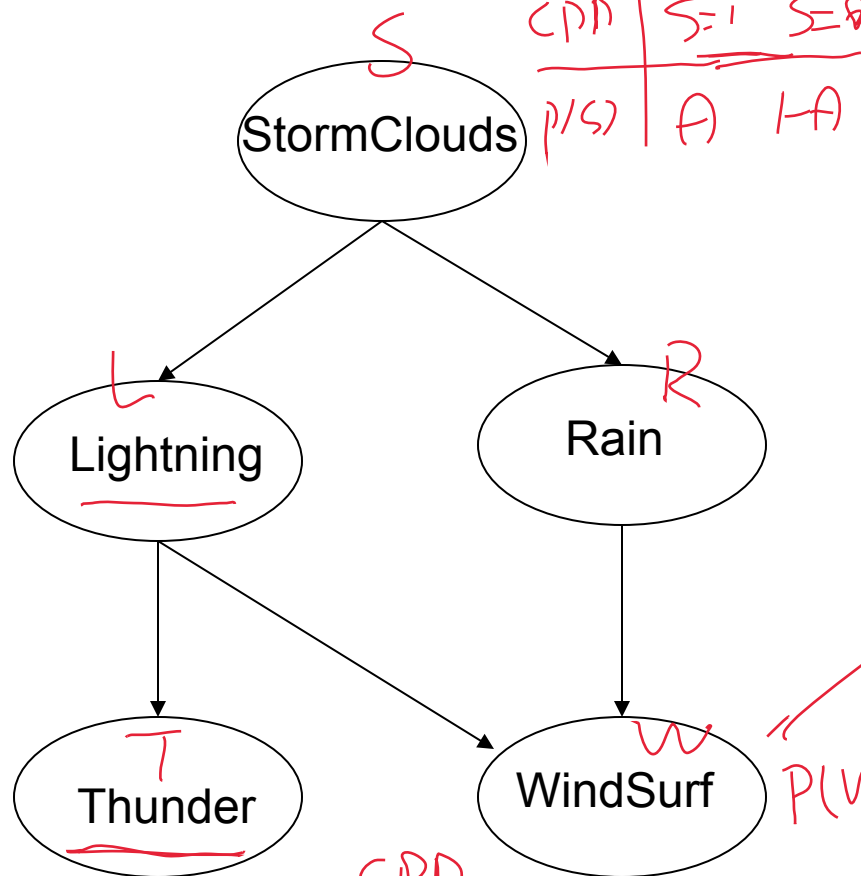
Nodes = random variables

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

*P*   *1-P*

| 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      |

*#i-parus = 4*

*CPD:*
*P(T|L)*

|      | T=1 | T=0 |
|------|-----|-----|
| L=1  | θ₁  | 1-θ₁ |
| L=0  | θ₂  | 1-θ₂ |

*#i-parus = 2*

*CPD*
*↓ P(W|L,R)*
*2² = 4*

The joint distribution over all variables:

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

*CPD*

*P(W=1, T=0, R=0, L=1, S=1) =*
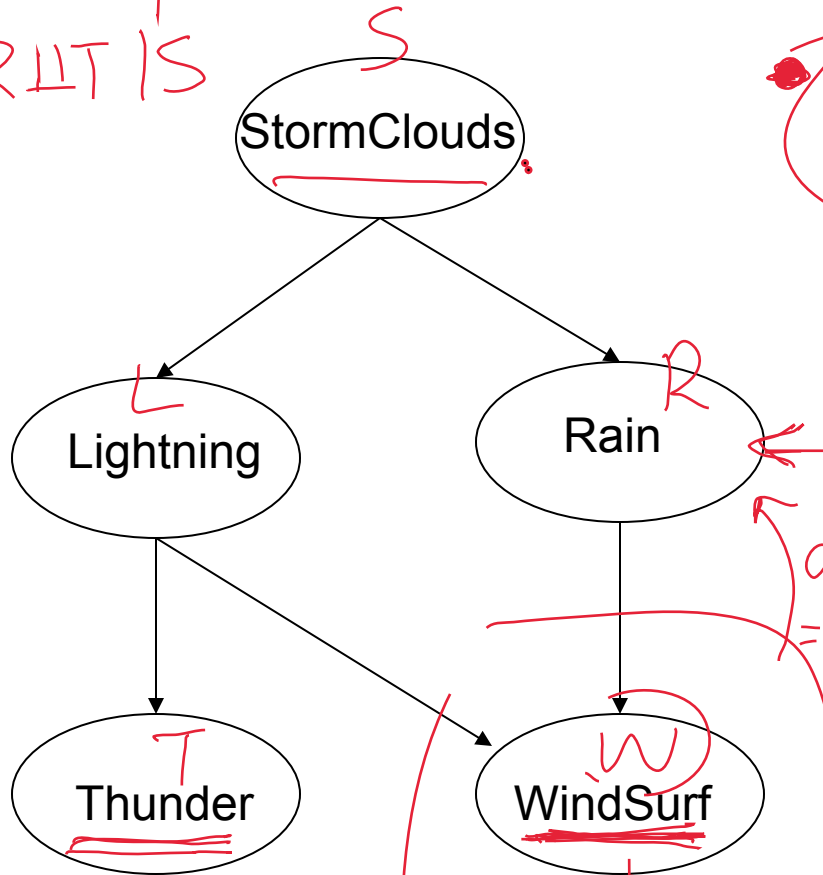*P(S) · P(R|S)  P(L|S) P(W|L,R) P(T|L) = ...*

# 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-descendents, given only its immediate parents.

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

$R \perp\!\!\!\perp T \mid S$

S

**StormClouds**

L

**Lightning**

R

**Rain**

T

**Thunder**

W

**WindSurf**

$\} de(R) = \{W, X\}$

| 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       |

**WindSurf**

X

$R \not\perp\!\!\!\perp X \mid S$

Cond. independ.

$P(X, Y \mid Z) = P(X \mid Z) P(Y \mid Z)$

$P(X \mid Y, Z) = P(X \mid Z)$

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

$(= P(T \mid L))$

$Pa(x) \subseteq An(x)$

# Some helpful terminology

$Pa(x)$.

$An(x)$.

$Ch(x)$.

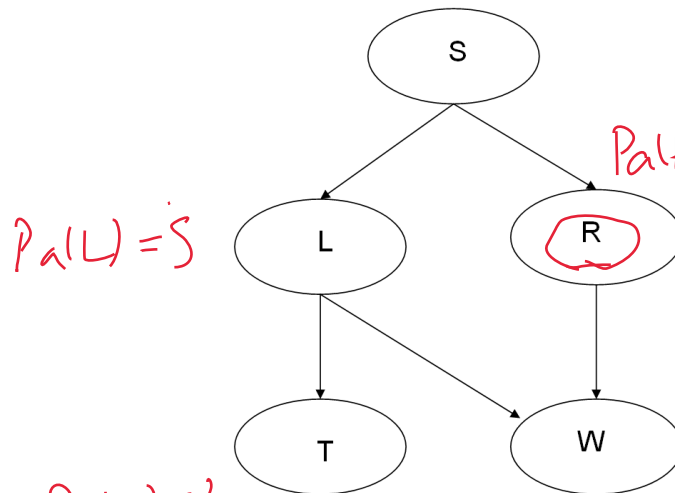$De(x)$.

Parents = Pa(X) = immediate parents

Antecedents = parents, parents of parents, ...

Children = immediate children

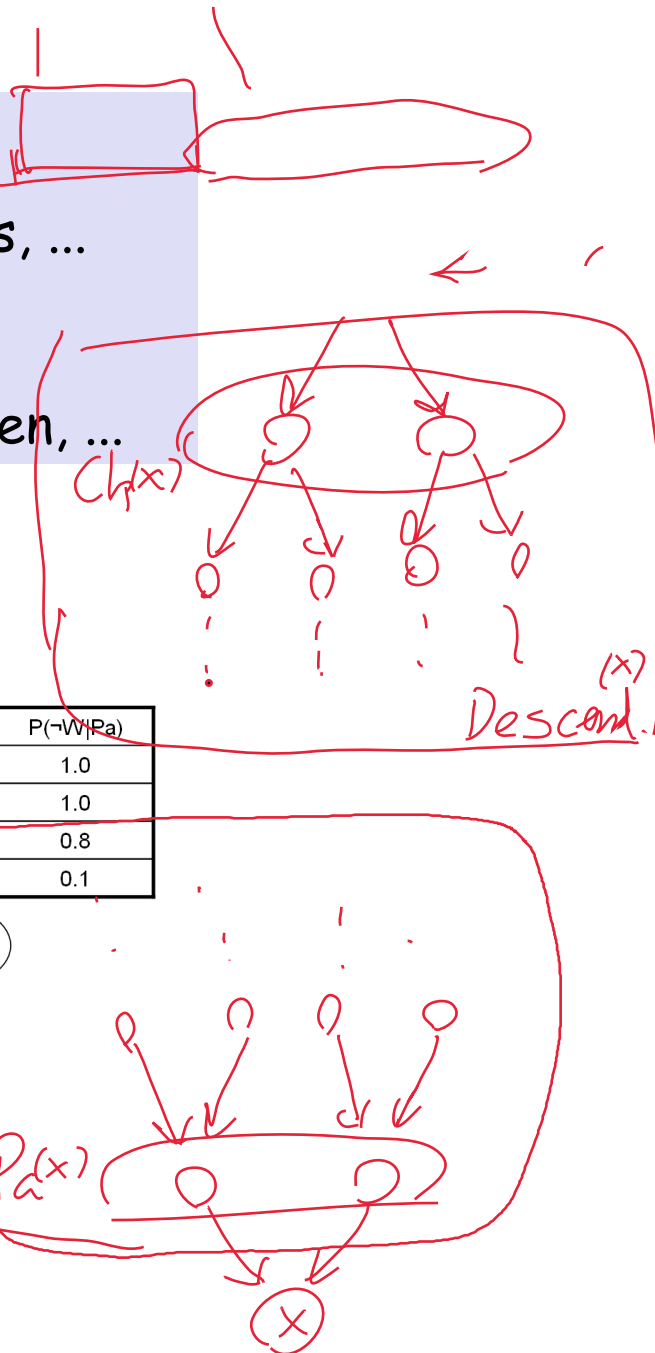Descendents = children, children of children, ...

$Ch(x) \subseteq De(X)$

Ante(x)

Ch(x)

Descend.(x)

$Pa(L) = S$

$Pa(R) = S$

$Pa(T) = L$

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

$Pa(x)$

| 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 |

(Diagram: S → L, S → R, L → T, L → W, R → W)

# Bayesian Networks

Boolean



| 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 |

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

$$\#iparms = 2^5 - 1 = 31$$

Chain rule of probability says that in general:

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

no condⁱ independ.

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

$$1 \quad 2 \quad 2 \quad 2 \quad 2^2 = 4$$

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

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

cond. independ.

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

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

$\#iparas = 11$

## 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

$n$ boolean

To define joint distribution in general?     $2^n$     exponential

To define joint distribution for this Bayes Net?

1 order :   $2n$

2 order :   $2^2 n$     $\longrightarrow$ linear

3 order :   $2^3 n$

NP complete $P(W = 1)$

[2^n]

# Inference in Bayes Nets

$BN \big\langle \begin{array}{l} DAG \\ CPD \end{array}$

S

**StormClouds**

L

**Lightning**

R

**Rain**

| 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**

T

**Thunder**

W

**WindSurf**

CPD

$$P(s, L, r, t) = \sum_{w \Rightarrow \{0,1\}} P(W = w, s, L, r, t)$$

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

$$P(T=0 \mid L=0)\ P(W=1 \mid L=0, R=1)$$
$$0.2$$

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

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

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

BN  $P(W=1 \mid L=1, R=0) = 0$

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

Training Set: $D = \{x_1, x_2, \ldots, x_m\}$  $x_m \in \mathbb{B}^5$

# Learning a Bayes Net

$x^T = (S, L, R, T, W)$

#params = 4

StormClouds

Lightning          Rain

Thunder            WindSurf

| Parents | P(W\|Pa) | P(¬W\|Pa) |
|---------|---------|-----------|
| L, R | 0 $\theta_{11}$ | 1.0 $1-\theta_{11}$ |
| L, ¬R | 0 $\theta_{10}$ | 1.0 $1-\theta_{10}$ |
| ¬L, R | 0.2 $\theta_{01}$ | 0.8 $1-\theta_{01}$ |
| ¬L, ¬R | 0.9 $\theta_{00}$ | 0.1 $1-\theta_{00}$ |

WindSurf

$$\hat{\theta}_1 = \frac{\alpha_{11} + \beta_{11}}{(\alpha_{11} + \alpha_{10}) + (\beta_{11} + \beta_{10})} \qquad \hat{\theta}_0 = \frac{\alpha_{01} + \beta_{01}}{(\alpha_{01} + \alpha_{00}) + \beta_0}$$

$$\Theta = \langle \theta_1, \theta_0 \rangle$$

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

What is the MLE solution?  MAP?

$L T$   $L(1-T)$   $(1-L)T$   $(1-L)(1-T)$

CPD

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

$P(T|L) = \theta_1 \qquad (1-\theta_1) \qquad \theta_0 \qquad (1-\theta_0)$

i.i.d.

$$L(\theta) = P(D|\theta)$$

$$\frac{\partial L(\theta)}{\partial \theta_1} = \theta_1^{\sum_{i=1}^m \{L=1\}T} (1-\theta_1)^{\sum \{L=1\}(\bar{T})} \cdots$$

$$\hat\theta_0 = \cdots (1-\theta_1)^{\alpha_{00}} = 0$$

# 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

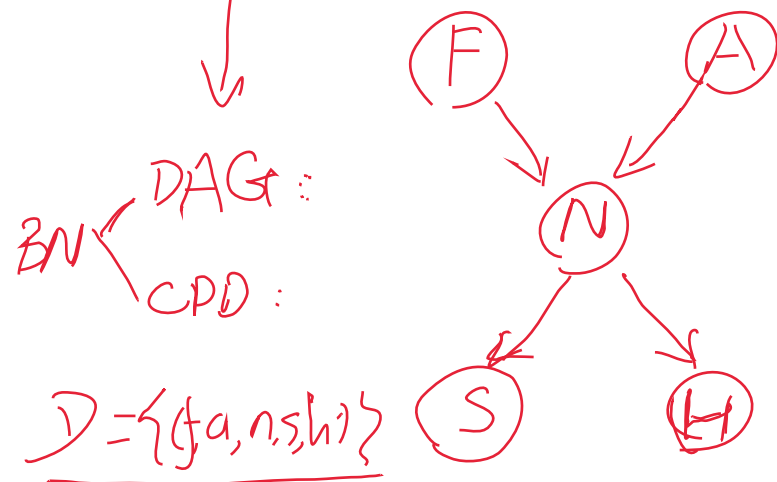  $$P(X_i|Pa(X_i)) = P(X_i|X_1, \ldots, X_{i-1})$$

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

Prior:

Boolean

F

A

N

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

S

H



BN { DAG :
     CPD :

$D = \{(f, a, n, s, h)\}$

$A \perp\!\!\!\perp F$

$A \not\!\perp\!\!\!\perp F \mid N$

MLE :
MAP

① PDF
② Likelihood
③ Derivative

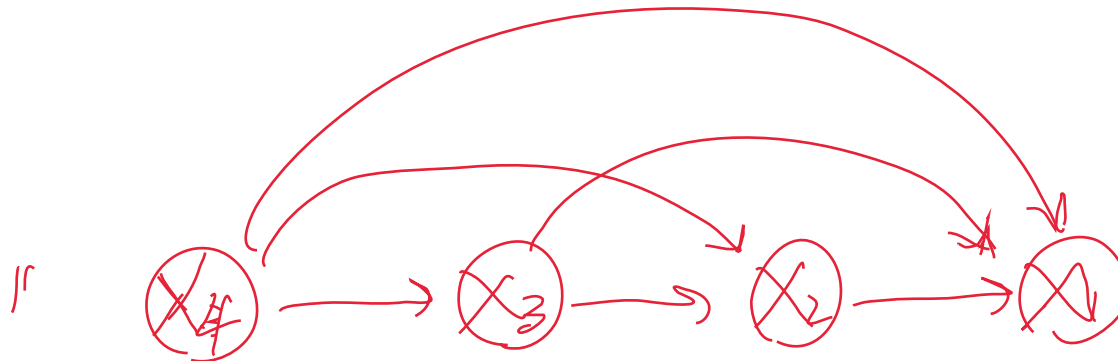| $P(H \mid N)$ | $H=1$ | $H=0$ |
|---|---|---|
| $N=1$ | $\theta_1$ | $1-\theta_1$ |
| $N=0$ | $\theta_0$ | $1-\theta_0$ |

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

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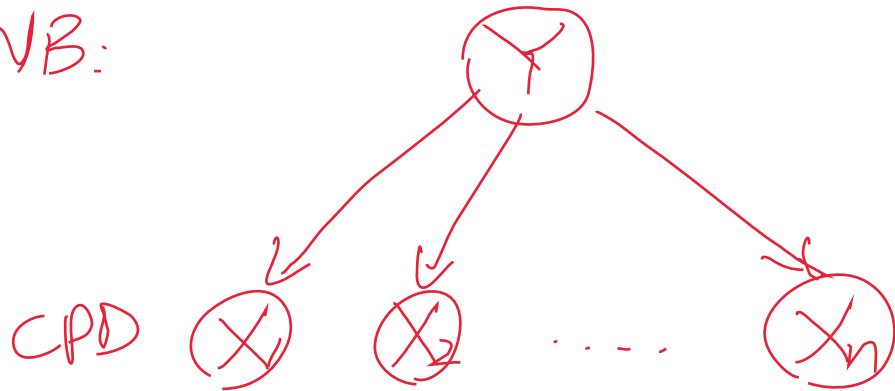


fully-connected BN

$$4! = 4 \times 3 \times 2 \times 1$$



$$P(X_1, X_2, X_3, X_4) = \quad - \quad \underline{\quad\quad} \quad -$$

# What is the Bayes Network for Naïve Bayes?

NB:



CPD

$X_i \perp\!\!\!\perp X_j \mid Y$ . $\forall\ i \neq j$

$P\left(Y=1 \mid X_1=1, X_2=0, \sim, X_n=1\right)$

$= \dfrac{\left(P(Y=1, x_1, \ldots, x_n\right)}{P(x_1, x_2 \ldots x_n)}$

$= \sum\limits_{y \in \{0, 1\}} P(Y=y, x_1, \ldots, x_n)$.

$P(X \mid Y) = \overset{n}{\underset{i=1}{\prod}} P(X_i \mid Y)$.

$P(X, Y) = P(X \mid Y) P(Y)$

$= \overset{n}{\underset{i=1}{\prod}} P(X_i \mid Y) P(Y)$

### GNB

$P(Y \mid X) \propto P(X \mid Y) P(Y)$

$\dfrac{}{\prod}$

$\overset{n}{\underset{i=1}{\prod}} P(X_i \mid Y)$

Q.

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

$[0, 5.0]$  $\{0, 10\}$  $\{10, 20\}$  $- - \{40, 50\}$

$P(S|A)$

$P(S|A)$ | $S=1$ | $S=0$

$A = 10$
$A = 2$
$A \doteq 3$
$\parallel$
$A = 5$

infinite

BatteryAge — A. ← continuous

BatteryState $\{0, 1\}$ ← discrete
S

Alternator  FanBelt  Leak

Charge

Lights  BatteryPower  GasInTank

Radio  GasGauge

Starter  Leak2

EngineCranks

FuelPump  Starts

Distributor

SparkPlugs

② parameterized model.

$\sigma(\beta A + \beta_0)$

Logistic / sigmoid.

$$P(x) = \frac{1}{1 + e^{-\theta(x)}}$$

$P(Y|X) \leftarrow \sigma(\beta^T X + \beta_0)$

$\beta^T X + \beta_0$