

# Ch 04. Parametric Model

# Part 1 Hidden Markov Model

# Markov Chain

---

- **State**  $\omega_i, i = 1, 2, \dots$
- The state at time  $t$   $\omega(t)$
- State sequence in discrete time of length  $T$

$$\omega^T = \{\omega(1), \omega(2), \dots, \omega(T)\}$$

For example:  $\omega^6 = \{\omega_1, \omega_4, \omega_2, \omega_2, \omega_1, \omega_4\}$

- **Transition Probability (Matrix)**

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$a_{ij} = P(\omega(t+1) = \omega_j \mid \omega(t) = \omega_i)$$

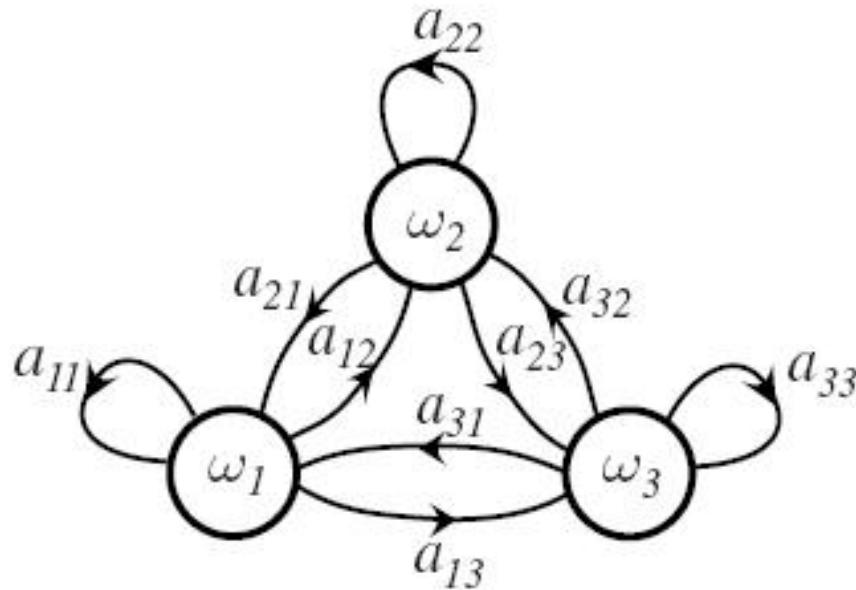
$$\sum_j a_{ij} = 1$$

$a_{ij}$  is the transition probability from state  $\omega_i$  to state  $\omega_j$

# Markov Chain

---

- State Transition Diagram



# Markov Chain

---

- **j-order Markov process**

- The probability of being a certain state at next moment is only related to the nearest j states

$$\begin{aligned} P(\omega(t+1) \mid \omega(1), \omega(2), \dots, \omega(t)) \\ = P(\omega(t+1) \mid \underbrace{\omega(t-j+1), \omega(t-j+2), \dots, \omega(t)}_{\text{only related to the nearest j states}}) \end{aligned}$$

- **First-order Markov process**

- The probability of being a state at any moment is only related to the state at the previous moment

$$P(\omega(t+1) \mid \omega(1), \omega(2), \dots, \omega(t)) = P(\omega(t+1) \mid \underbrace{\omega(t)}_{\text{only related to the state at the previous moment}})$$

# Hidden Markov Model

---

- **Hidden Markov Model** ( abbreviated to **HMM**, 隐马尔可夫模型 )
- State is **invisible**
- At time  $t$ , the hidden state excites the visible symbol  $x(t)$  with a certain probability, whose value is expressed as  $v_1, v_2, v_3, \dots$
- A sequence of visible symbols in discrete time of length  $T$   $\mathbf{X}^T = \{x(1), x(2), \dots, x(T)\}$

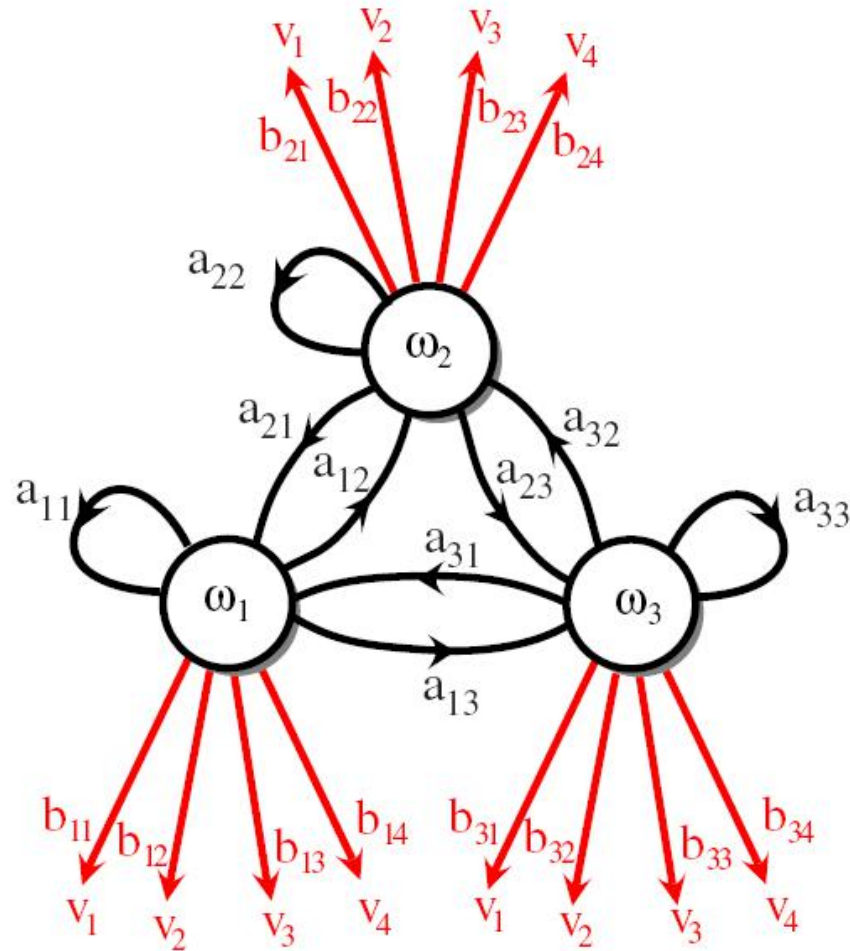
For example:  $\mathbf{X}^6 = \{v_5, v_1, v_1, v_5, v_2, v_3\}$

- The probability of observing visible symbols

$$b_{jk} = P(x(t) = v_k \mid \omega(t) = \omega_j) \quad \sum_k b_{jk} = 1$$

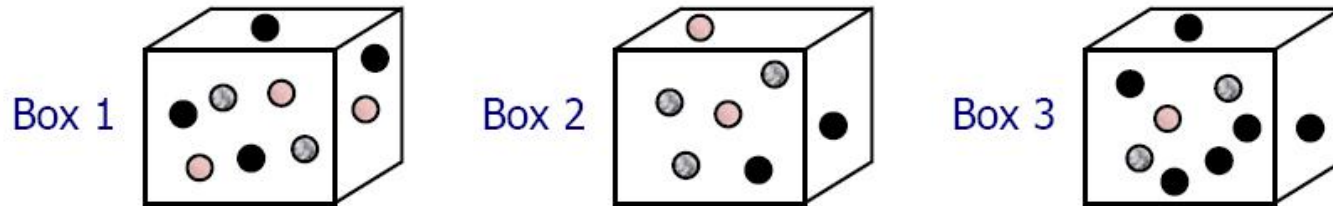
# Hidden Markov Model

- State Transition Diagram



# One example

---



- The box number is not visible
- Take one small ball out of any box at a time
- **Hidden State**: box number
- **Visible symbol**: small ball
- The probability of getting various small balls out of box  $i$   
 $P(\bullet | i)$      $P(\bullet | i)$      $P(\bullet | i)$
- What's the probability of getting a particular sequence of small balls?  
● ● ● ● ● ● ● ● ● ●



# Symbolic Representation of Discrete HMM

---

- Hidden state set

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$$

- Visible state set

$$V = \{v_1, v_2, \dots, v_m\}$$

- State sequence

$$\omega = \omega(1) \omega(2) \dots \omega(T)$$

- Observed sequence

$$\mathbf{X} = x(1) x(2) \dots x(T)$$

- State transition probability

$$\mathbf{A} = \{a_{ij} \mid a_{ij} = P(\omega(t+1) = \omega_j \mid \omega(t) = \omega_i)\}$$

- The probability of observing a visible sign

$$\mathbf{B} = \{b_{jk} \mid b_{jk} = P(x(t) = v_k \mid \omega(t) = \omega_j)\}$$

- Initial state probability

$$\Pi = \{\pi_i \mid \pi_i = P(\omega(1) = \omega_i)\}$$

Complete HMM  
parameter vector

$$\theta = (\mathbf{A}, \mathbf{B}, \Pi)$$

# Three Core Problems of HMM

---

- **Valuation problem**

- **Given**

- A specific symbol sequence  $\mathbf{X}$  is observed
    - Parameter vector of HMM  $\boldsymbol{\theta}$

- **To solve**

- Likelihood function  $P(\mathbf{X} | \boldsymbol{\theta})$

- **Decoding problem**

- **Given**

- A specific symbol sequence  $\mathbf{X}$  is observed
    - Parameter vector of HMM  $\boldsymbol{\theta}$

- **To solve**

- The hidden state sequence most likely to produce  $\mathbf{X}$

# Three Core Problems of HMM

---

- **Learning (or parameter estimation) problem**

- **Given**

- A specific symbol sequence  $\mathbf{X}$  is observed

- **To Solve**

- The estimated value of the model parameter vector  $\theta$

For example: ML estimation

$$\theta_{\text{ML}} = \arg \max_{\theta} P(\mathbf{X} | \theta)$$

# Summary

---

- First-order Markov chain
- Hidden Markov Model (HMM)
- Three core problems of HMM
  - Valuation Problem
    - **HMM forward algorithm**
    - **HMM backward algorithm**

# Valuation Problem

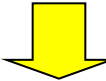
---

- Directly calculate the probability that the HMM model produces a symbol sequence  $\mathbf{X}$  of visible length  $T$

$$P(\mathbf{X} | \boldsymbol{\theta}) = \sum_{\boldsymbol{\omega}} P(\mathbf{X} | \boldsymbol{\omega}, \boldsymbol{\theta}) P(\boldsymbol{\omega} | \boldsymbol{\theta})$$

$$P(\mathbf{X} | \boldsymbol{\omega}, \boldsymbol{\theta}) = b_{\omega(1)x(1)} b_{\omega(2)x(2)} \cdots b_{\omega(T)x(T)}$$

$$P(\boldsymbol{\omega} | \boldsymbol{\theta}) = \pi_{\omega(1)} a_{\omega(1)\omega(2)} a_{\omega(2)\omega(3)} \cdots a_{\omega(T-1)\omega(T)}$$


$$P(\mathbf{X} | \boldsymbol{\theta}) = \sum_{\boldsymbol{\omega}} \prod_{t=1}^T a_{\omega(t-1)\omega(t)} b_{\omega(t)x(t)}$$

where  $a_{\omega(0)\omega(1)}$  represents the initial probability  $\pi_{\omega(1)}$  of the state  $\omega(1)$

Suppose that there are  $c$  hidden states in HMM, the computational complexity is  $O(c^T T)$  !

For example:  $c=10$ ,  $T=20$ , basic operations  $10^{21}$  times!

# Valuation Problem

---

- Solution

- Recursive Computation

The calculation at time  $T$  involves only the results of the previous step,  $\omega(t)$ ,  $\omega(t-1)$  and  $x(t)$

- **HMM forward algorithm**
    - **HMM backward algorithm**

# Valuation Problem

- HMM forward algorithm

Define  $\alpha_i(t)$  : The probability that state is  $i$  at time  $t$ , and  $x(1), x(2), \dots$  has been observed

- Initialization

for each hidden state  $i$ , calculate  $\alpha_i(1) = \pi_i b_{ix(1)}$

- Recursion

for  $t=2$  to  $T$

For each hidden state  $j$ , calculate  $\alpha_j(t) = \left[ \sum_{i=1}^c \alpha_i(t-1) a_{ij} \right] b_{jx(t)}$

end

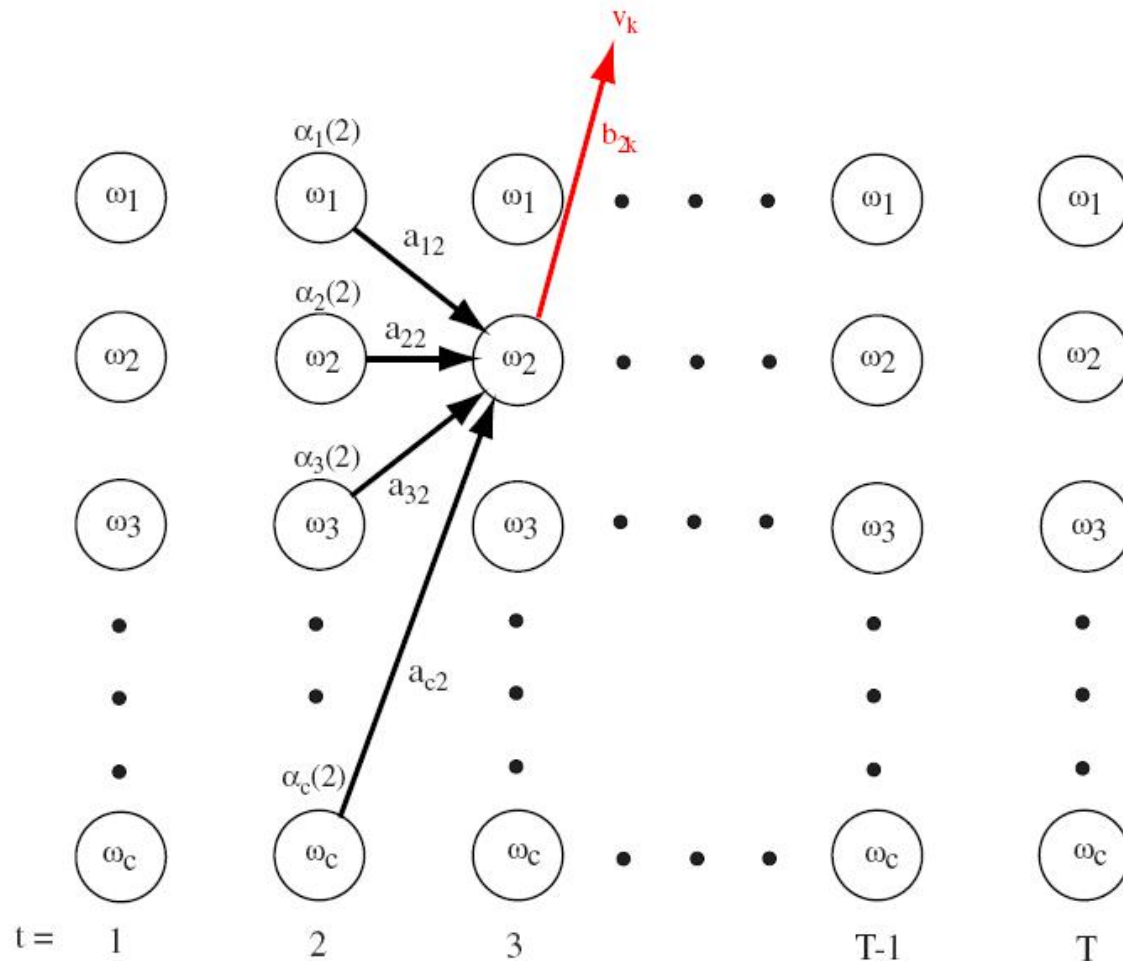
- End

$$P(\mathbf{X} | \boldsymbol{\theta}) = \sum_{i=1}^c \alpha_i(T)$$

Compute Complexity  $O(c^2T)$   $O(c^T T)$

# Valuation Problem

- **HMM forward algorithm**





# Valuation Problem

- **HMM backward algorithm** (time inversion version of forward algorithm)

Define  $\beta_i(t)$ : The probability that state is  $i$  at time  $t$ , and  $x(T), x(T-1), \dots, x(t)$  has been observed in reverse

- **Initialization**

for each hidden state, calculate  $\beta_i(T) = \frac{b_{ix(T)}}{c}$

(suppose that the probability of each state at time  $T$  are the same)

- **Recursion**

for  $t=T-1$  to  $1$

For each hidden state  $i$ , compute  $\beta_i(t) = \left[ \sum_{j=1}^c a_{ij} \beta_j(t+1) \right] b_{ix(t)}$

end

- **End**

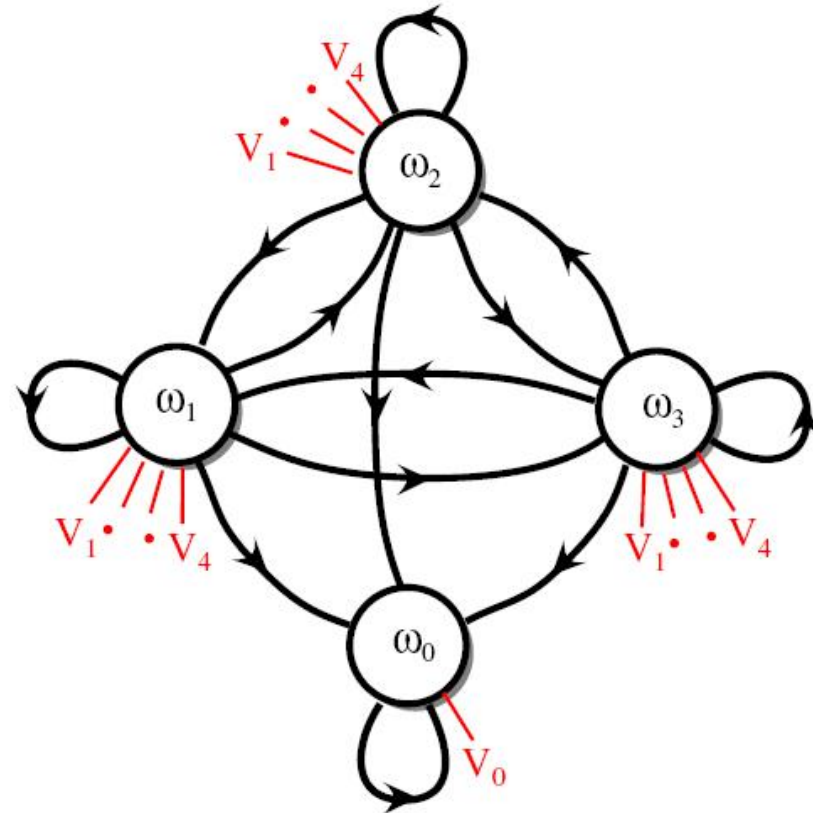
$$P(\mathbf{X} | \boldsymbol{\theta}) = \sum_{i=1}^c \pi_i \beta_i(1)$$

Compute Complexity  $O(c^2T)$   $O(c^T T)$

# Example

- HMM is

$$a_{ij} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.2 & 0.3 & 0.1 & 0.4 \\ 0.2 & 0.5 & 0.2 & 0.1 \\ 0.8 & 0.1 & 0.0 & 0.1 \end{pmatrix}$$
$$b_{jk} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0.3 & 0.4 & 0.1 & 0.2 \\ 0 & 0.1 & 0.1 & 0.7 & 0.1 \\ 0 & 0.5 & 0.2 & 0.1 & 0.2 \end{pmatrix}$$



- $\omega_0$ : **The absorbed state** that the inevitable state at the end of the sequence. This state produces a unique particular visible symbol  $V_0$ , indicating the end of the HMM process

# Example

---

- Given state is  $\omega_1$  at time  $t$ , that is

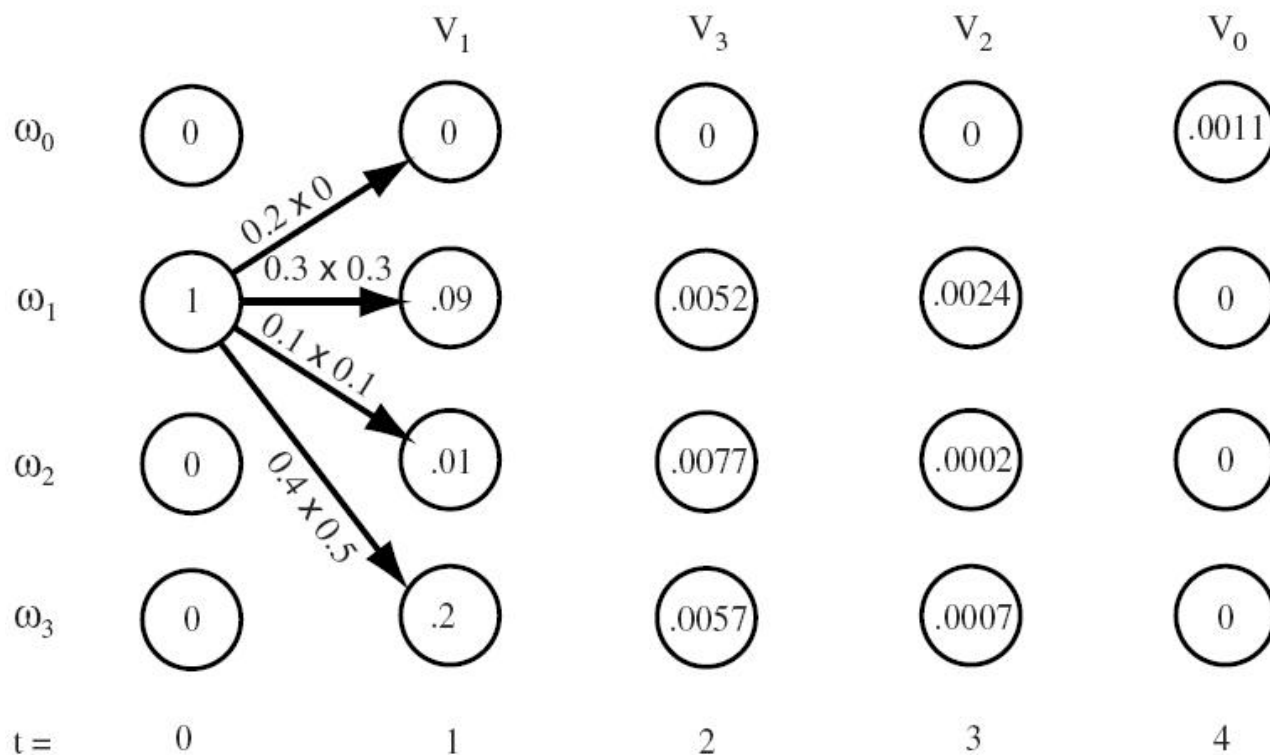
$$\pi_0 = a_{10} = 0.2, \quad \pi_1 = a_{11} = 0.3,$$

$$\pi_2 = a_{12} = 0.1, \quad \pi_3 = a_{13} = 0.4$$

- The observed sequence is  $V^4 = \{v_1, v_3, v_2, v_0\}$
- Calculate the probability that the HMM produce this particular observation sequence

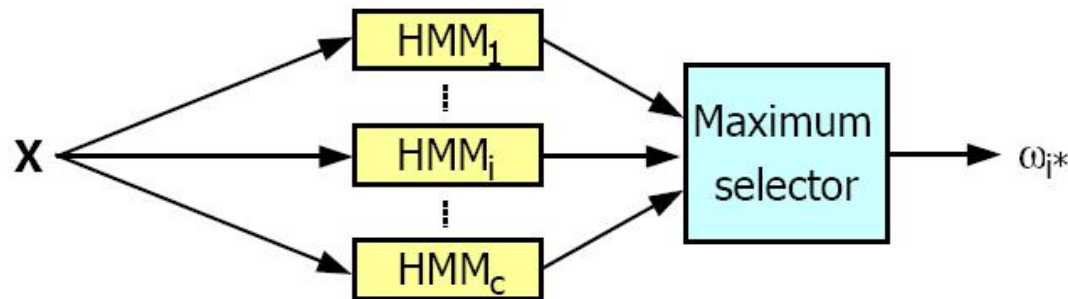
# Example

- Solution



# HMM for classification

- Build a HMM for each category



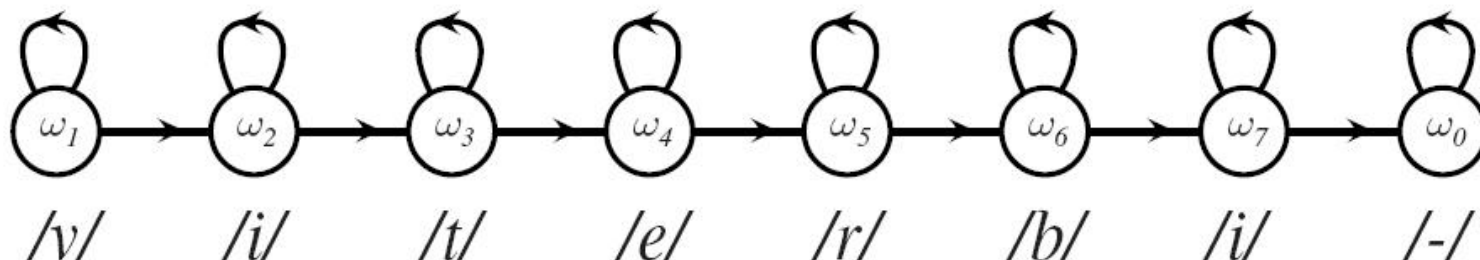
- Each HMM has its own parameter vector  $\boldsymbol{\theta}_i$ , which can be learned (estimated) from the samples belonging to category  $i$

- Bayes Decision  $P(\boldsymbol{\theta}_i | \mathbf{X}) = \frac{P(\mathbf{X} | \boldsymbol{\theta}_i)P(\boldsymbol{\theta}_i)}{\sum_{i=1}^c P(\mathbf{X} | \boldsymbol{\theta}_i)P(\boldsymbol{\theta}_i)}$ 
  - Decision Result

$$i^* = \arg \max_i (P(\mathbf{X} | \boldsymbol{\theta}_i)P(\boldsymbol{\theta}_i))$$

# HMM for Automatic Speech Recognition (ASR)

- left-to-right (从左到右) HMM



left-to-right HMM for pronunciation of **viterbi**

- Build a HMM for each word pronunciation, whose parameter is  $\theta_i$
- Use forward algorithm to calculate the class-conditional probability  $P(\mathbf{X}|\theta_i)$  of pronunciation sequence  $\mathbf{X}$
- $P(\theta_i)$  depends on the language itself and the context semantics
- Use Bayes Formula to calculate the posterior probability  $P(\theta_i | \mathbf{X})$  of  $\mathbf{X}$
- The maximum posterior probability indicates the speech content

# Decoding Problem

---

- Given a sequence of observations  $X^T$ , look for the most likely sequence of hidden states
- Method of exhaustion
  - Calculate the probabilities of all the possible sequences of hidden states
  - Compute complexity  $O(c^T T)$

# Decoding Problem

- **Viterbi algorithm**

- **Initialization**

for each hidden state  $i$ , compute  $\delta_i(1) = \pi_i b_{ix(1)}$

- **Recursion**

for  $t=2$  to  $T$ :

For each hidden state  $j$ , compute

$$\delta_j(t) = \left[ \max_{1 \leq i \leq c} \delta_i(t-1) a_{ij} \right] b_{jx(t)} \quad \psi_j(t) = \arg \max_{1 \leq i \leq c} \delta_i(t-1) a_{ij}$$

end

- **End**

$$\omega^*(T) = \arg \max_{1 \leq i \leq c} \delta_i(T)$$

for  $t=T-1$  to  $1$  (path backtracking) :

$$\omega^*(t) = \psi_{\omega^*(t+1)}(t+1)$$

end

**Compute**  $O(c^2T)$      $O(c^T T)$   
**Complexity**

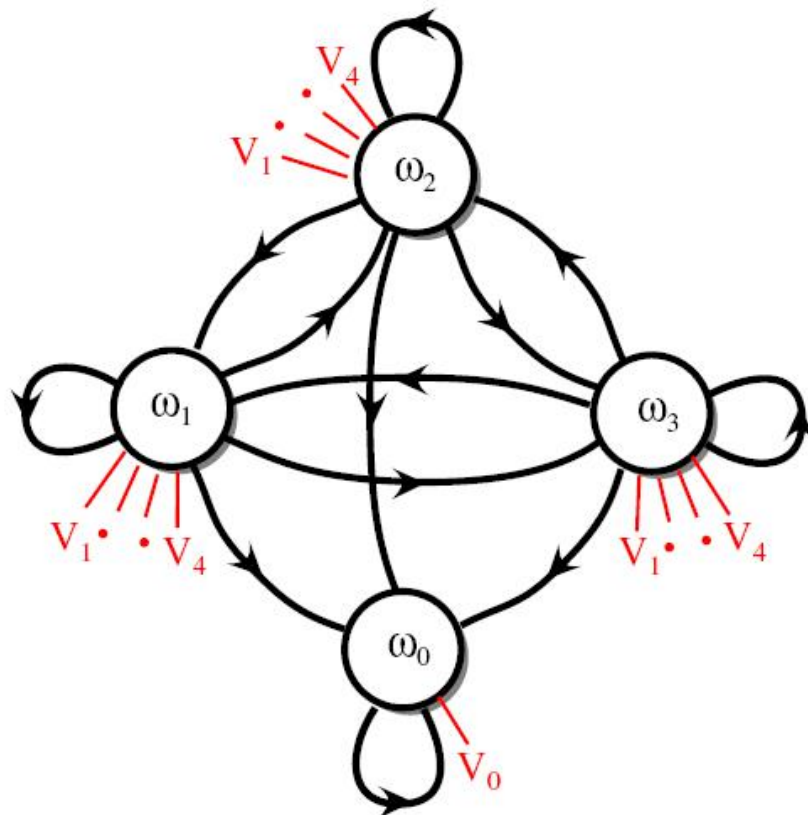


# Example

- HMM is

$$a_{ij} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.2 & 0.3 & 0.1 & 0.4 \\ 0.2 & 0.5 & 0.2 & 0.1 \\ 0.8 & 0.1 & 0.0 & 0.1 \end{pmatrix}$$

$$b_{jk} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0.3 & 0.4 & 0.1 & 0.2 \\ 0 & 0.1 & 0.1 & 0.7 & 0.1 \\ 0 & 0.5 & 0.2 & 0.1 & 0.2 \end{pmatrix}$$



# Example

---

- Given that state is  $\omega_1$  at  $t=0$ , that is

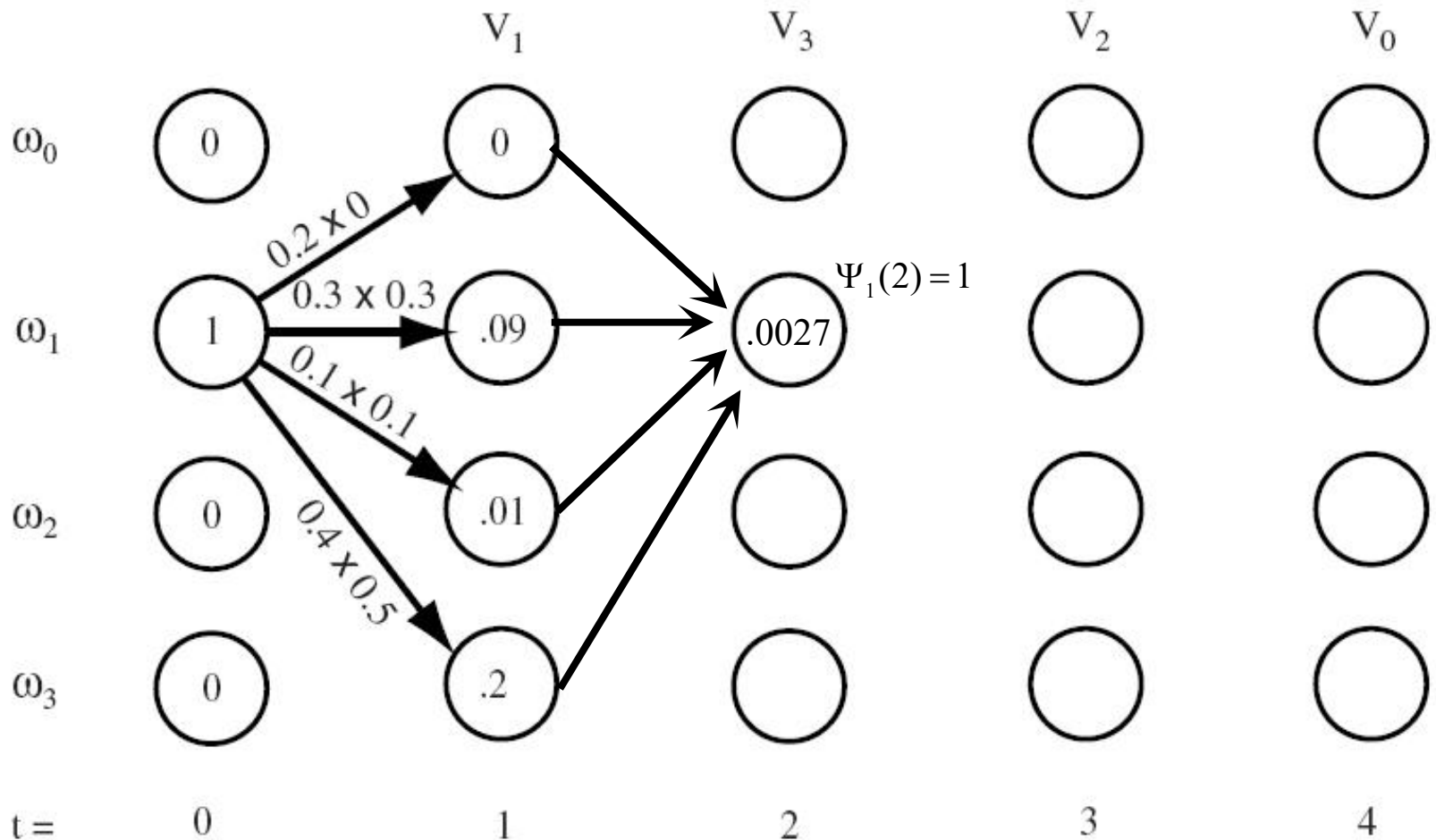
$$\pi_0 = a_{10} = 0.2, \quad \pi_1 = a_{11} = 0.3,$$

$$\pi_2 = a_{12} = 0.1, \quad \pi_3 = a_{13} = 0.4$$

- The observed sequence is  $V^4 = \{v_1, v_3, v_2, v_0\}$
- Compute the most likely sequence of hidden states

# Example

- Solution



**Exercise:** Fill out the diagram and trace back the optimal path

# Decoding Problem

---

- For long sequences, the Viterbi algorithm may cause the computer **underflow**
- **Improvement:** Viterbi algorithm based on logarithm
  - Advantage
    - Change multiplication to addition
    - Avoid underflow
    - The results are same as Viterbi algorithm

$$\tilde{a}_{ij} = \log a_{ij}$$

$$\tilde{b}_{iX(t)} = \log b_{iX(t)}$$

$$\tilde{\pi}_i = \log \pi_i$$

$$\tilde{\xi}_i(t) = \log \delta_i(t)$$

# Decoding Problem

---

- **Log-Viterbi algorithm**

- **Initialization**

- for each hidden state  $i$ , compute  $\tilde{\delta}_i(1) = \tilde{\pi}_i + \tilde{b}_{ix(1)}$

- **Recursion**

- for  $t=2$  to  $T$ :

- for each hidden state  $i$ . compute

- $$\tilde{\delta}_j(t) = \max_{1 \leq i \leq c} [\tilde{\delta}_i(t-1) + \tilde{a}_{ij}] + \tilde{b}_{jx(t)} \quad \psi_j(t) = \arg \max_{1 \leq i \leq c} [\tilde{\delta}_i(t-1) + \tilde{a}_{ij}]$$

- end

- **End**

- $$\omega^*(T) = \arg \max_{1 \leq i \leq c} \tilde{\delta}_i(T)$$

- for  $t=T-1$  to  $1$  (path backtracking) :

- $$\omega^*(t) = \psi_{\omega^*(t+1)}(t+1)$$

- end

# Learning Problem

---

- Learning parameter vector  $\theta$  of HMM from a set of training data  $D=\{X_1, X_2, \dots, X_n\}$
- There is no algorithm to determine the optimal HMM parameter based on the training set
- Common algorithms

**forward-backward algorithm** (向前向后算法)

also known as

**Baum-Welch re-estimation algorithm**

( **Baum-Welch** 重估算法 )

- **Core idea**

- Recursively update the HMM parameters to get the HMM parameters that best explain the training sample

# Learning Problem

- **Baum-Welch re-estimation algorithm**
- Given  $\mathbf{X}$  and  $\boldsymbol{\theta}$ , the posterior probability that is state  $i$  at time  $t$  and is state  $j$  at time  $t+1$

$$\gamma_{ij}(t) = \frac{\alpha_i(t-1) a_{ij} b_{jk} \beta_j(t)}{P(\mathbf{X}^T | \boldsymbol{\theta})}$$

**forward** **backward**

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T \gamma_{ij}(t)}{\sum_{t=1}^T \sum_k \gamma_{ik}(t)}$$

$$\hat{b}_{jk} = \frac{\sum_{t=1}^T \sum_l \gamma_{jl}(t)}{\sum_{t=1}^T \sum_l \gamma_{jl}(t) \text{ where } v(t)=v_k}$$

# Learning Problem

---

- **forward-backward algorithm**

- **Initiate  $\theta$**

- **repeat**

- use **Baum-Welch re-estimation algorithm** to compute  $\hat{\theta}$   
based on  $\theta$  and  $X$

- $\theta \leftarrow \hat{\theta}$

- until**  $\theta$  convergence

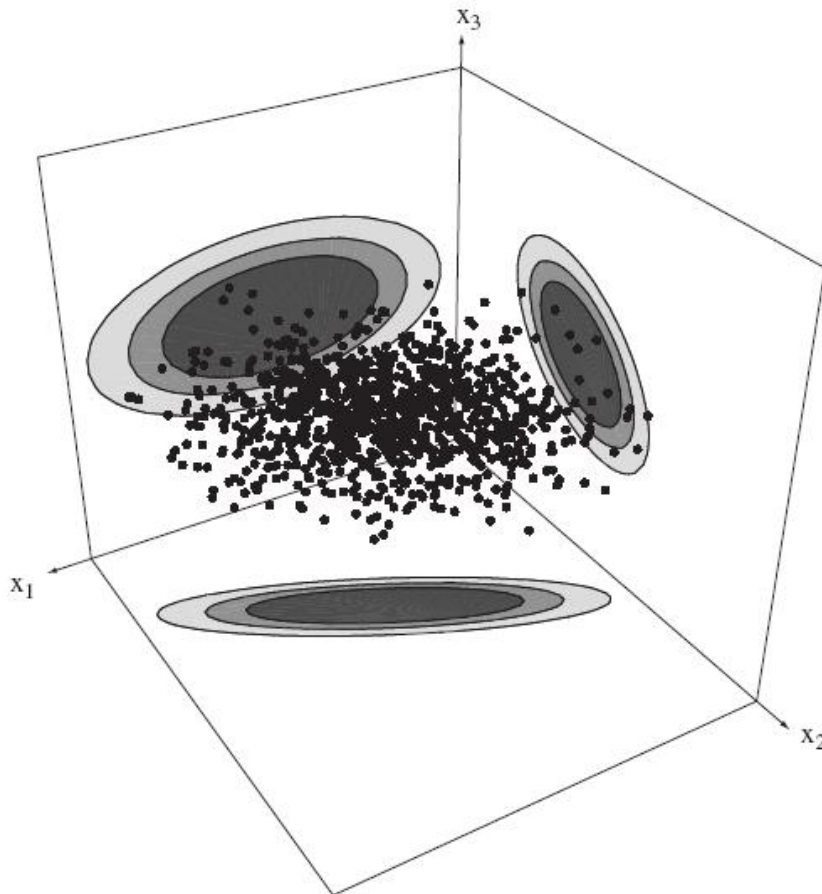
- returns the result of parameter estimate  $\theta$



# Part 2 Bayesian Belief Net

# Feature correlation

- In some cases, the prior knowledge about the distribution is not directly in the form of probability distribution, but related to the statistical correlation (or independence) relationship between each characteristic component



$x_1$  and  $x_3$  are statistically independent, but the other features are not

# Example of Correlation

---

- state of the car
  - engine temperature
  - oil temperature
  - oil pressure
  - tire pressure
- correlation
  - oil pressure and tire pressure are **independent** of each other
  - oil temperature is **related** to engine temperature

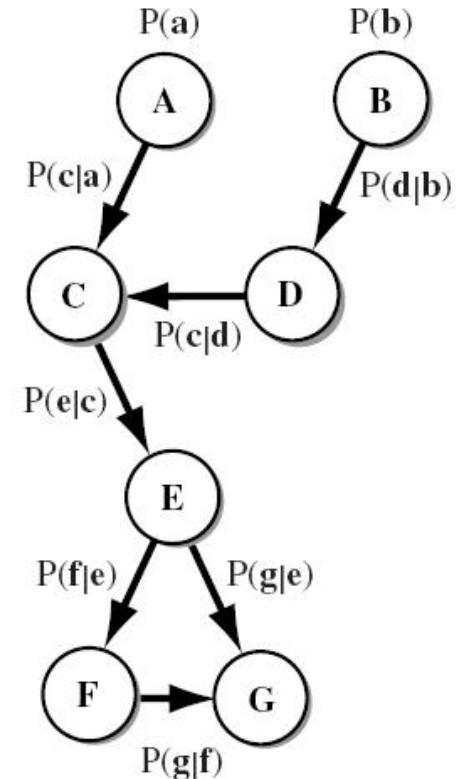
# Bayesian Belief Net

---

- Represent the causal dependence between features by graph
  - **Bayesian belief net** （贝叶斯置信网）
  - **causal network** （因果网）
  - **belief net** （置信网）
- Direct Acyclic Graph (DAG)
  - The connection between nodes is **directional**
  - There is **no cyclic** path in the graph
- Only discrete cases are discussed

# Bayesian Belief Net

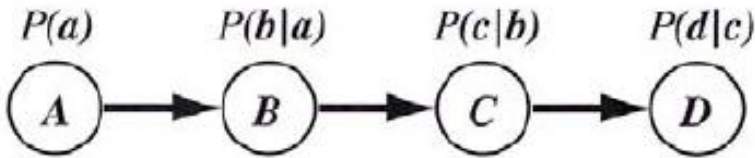
- Each node A, B, C... represents a system variable (feature)
  - The possible discrete values for each node
    - Values of A:  $a_1, a_2, a_3, \dots$
    - For example
      - A represents the state of the lamp
      - $a_1 = \text{on}, a_2 = \text{off}, P(a_1) = 0.7, P(a_2) = 0.3$
- Directed connections between nodes represent dependencies between variables
  - The connection from A to C represent  $P(c_i | a_j)$  or  $P(\mathbf{c} | \mathbf{a})$
- The state of any node can be inferred from the state of its adjacent nodes



# Joint Probability

---

- Linear Chain



$$P(a, b, c, d) = P(a)P(b | a)P(c | b)P(d | c)$$

$$P(b, c, d) = P(c | b)P(d | c) \sum_a P(a)P(b | a)$$

$$P(c, d) = P(d | c) \sum_a \sum_b P(a)P(b | a)P(c | b)$$

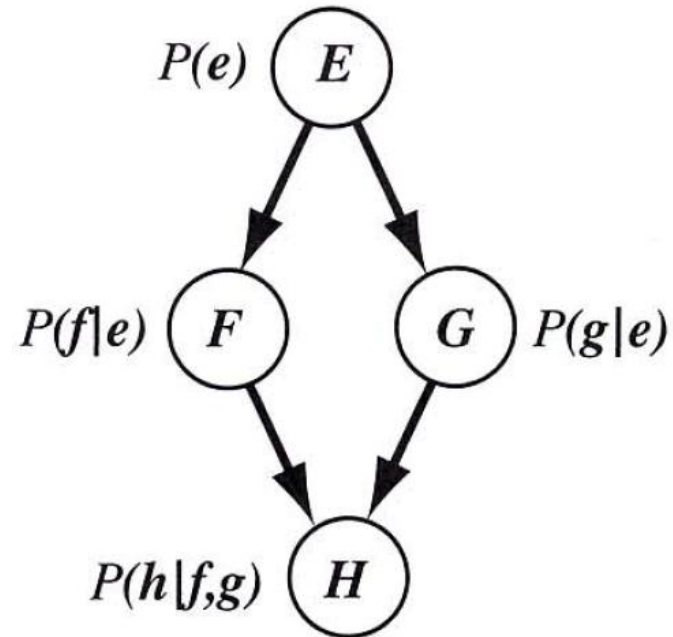
# Joint Probability

- Simple circuit

$$P(e, f, g, h) = P(e)P(f | e)P(g | e)P(h | f, g)$$

$$P(f, g, h) = P(h | f, g) \sum_e P(e)P(f | e)P(g | e)$$

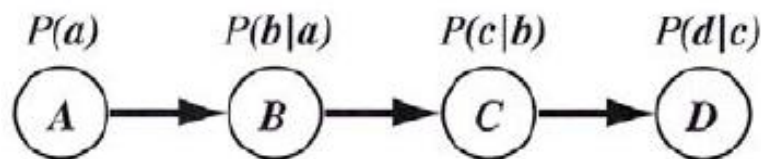
$$P(g, h) = \sum_e \sum_f P(e)P(f | e)P(g | e)P(h | f, g)$$



# The probability of any node taking a Specific value

---

- Linear chain



$$\begin{aligned} P(\mathbf{d}) &= \sum_{\mathbf{a}, \mathbf{b}, \mathbf{c}} P(\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}) \\ &= \sum_{\mathbf{a}, \mathbf{b}, \mathbf{c}} P(\mathbf{a}) P(\mathbf{b}|\mathbf{a}) P(\mathbf{c}|\mathbf{b}) P(\mathbf{d}|\mathbf{c}) \\ &= \sum_{\mathbf{c}} P(\mathbf{d}|\mathbf{c}) \underbrace{\sum_{\mathbf{b}} P(\mathbf{c}|\mathbf{b}) \underbrace{\sum_{\mathbf{a}} P(\mathbf{b}|\mathbf{a}) P(\mathbf{a})}_{P(\mathbf{b})}}_{P(\mathbf{c})} \\ &\quad \underbrace{\hspace{10em}}_{P(\mathbf{d})} \end{aligned}$$

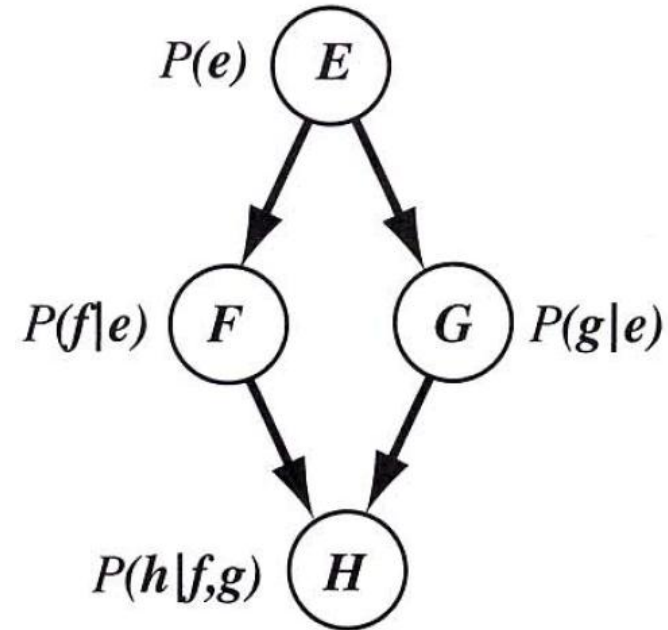


# The probability of any node taking a Specific value

---

- Simple circuit

$$\begin{aligned} P(h) &= \sum_{e,f,g} P(e, f, g, h) \\ &= \sum_{e,f,g} P(e)P(f | e)P(g | e)P(h | f, g) \end{aligned}$$



# Example One

- fish classification belief net

$P(a)$

$P(a_1)$	$P(a_2)$	$P(a_3)$	$P(a_4)$
0.25	0.25	0.25	0.25

$a_1 = \text{winter}$   
 $a_2 = \text{spring}$   
 $a_3 = \text{summer}$   
 $a_4 = \text{autumn}$

$P(b)$

$P(b_1)$	$P(b_2)$
0.6	0.4

$b_1 = \text{north Atlantic}$   
 $b_2 = \text{south Atlantic}$

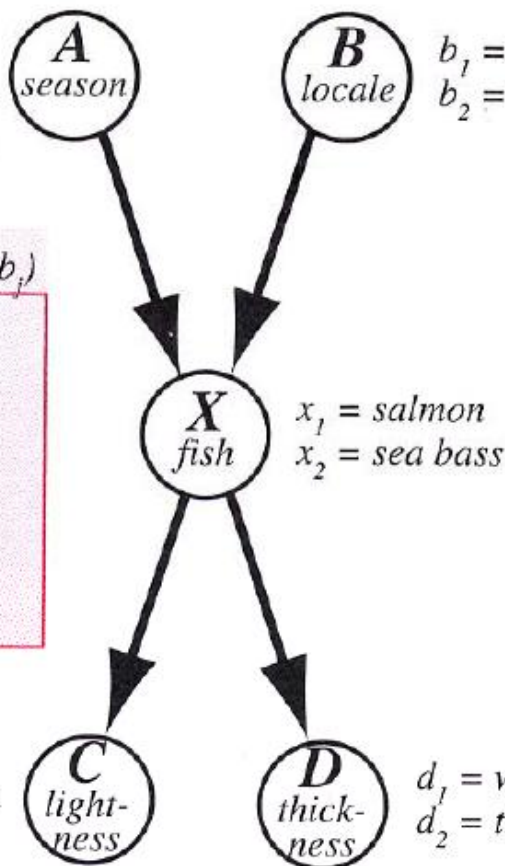
$P(x|a,b)$

$a_i, b_j$	$P(x_1 a_i, b_j)$	$P(x_2 a_i, b_j)$
$a_1, b_1$	0.5	0.5
$a_1, b_2$	0.7	0.3
$a_2, b_1$	0.6	0.4
$a_2, b_2$	0.8	0.2
$a_3, b_1$	0.4	0.6
$a_3, b_2$	0.1	0.9
$a_4, b_1$	0.2	0.8
$a_4, b_2$	0.3	0.7

$P(c|x)$

	$P(c_1 x_k)$	$P(c_2 x_k)$	$P(c_3 x_k)$
$x_1$	0.6	0.2	0.2
$x_2$	0.2	0.3	0.5

$c_1 = \text{light}$   
 $c_2 = \text{medium}$   
 $c_3 = \text{dark}$



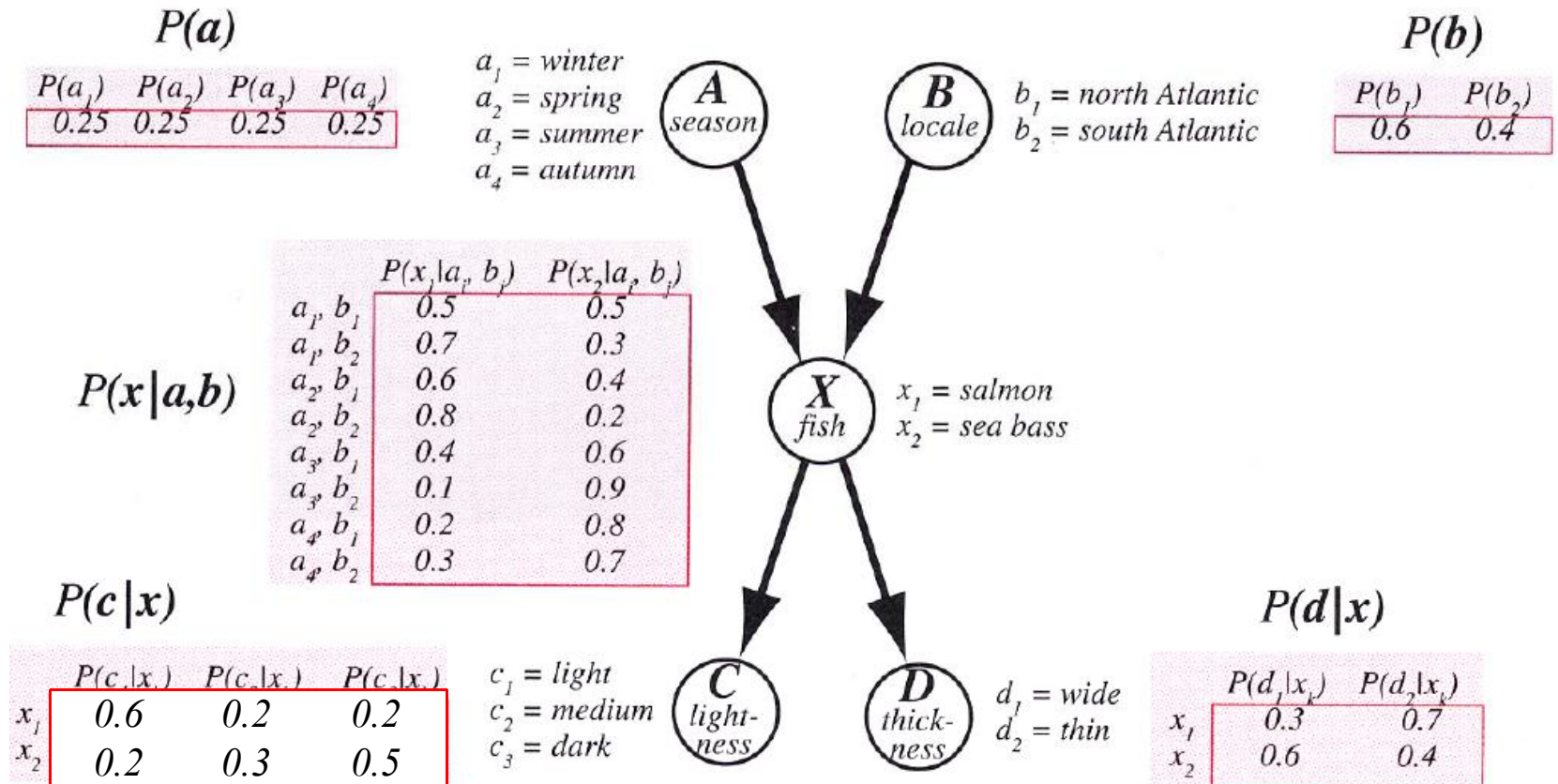
$P(d|x)$

	$P(d_1 x_k)$	$P(d_2 x_k)$
$x_1$	0.3	0.7
$x_2$	0.6	0.4

$d_1 = \text{wide}$   
 $d_2 = \text{thin}$

# Example One

- Calculate the probability that "a fish caught in the North Atlantic in summer is a bass with a dim gloss and a narrow width"



# Example One

---

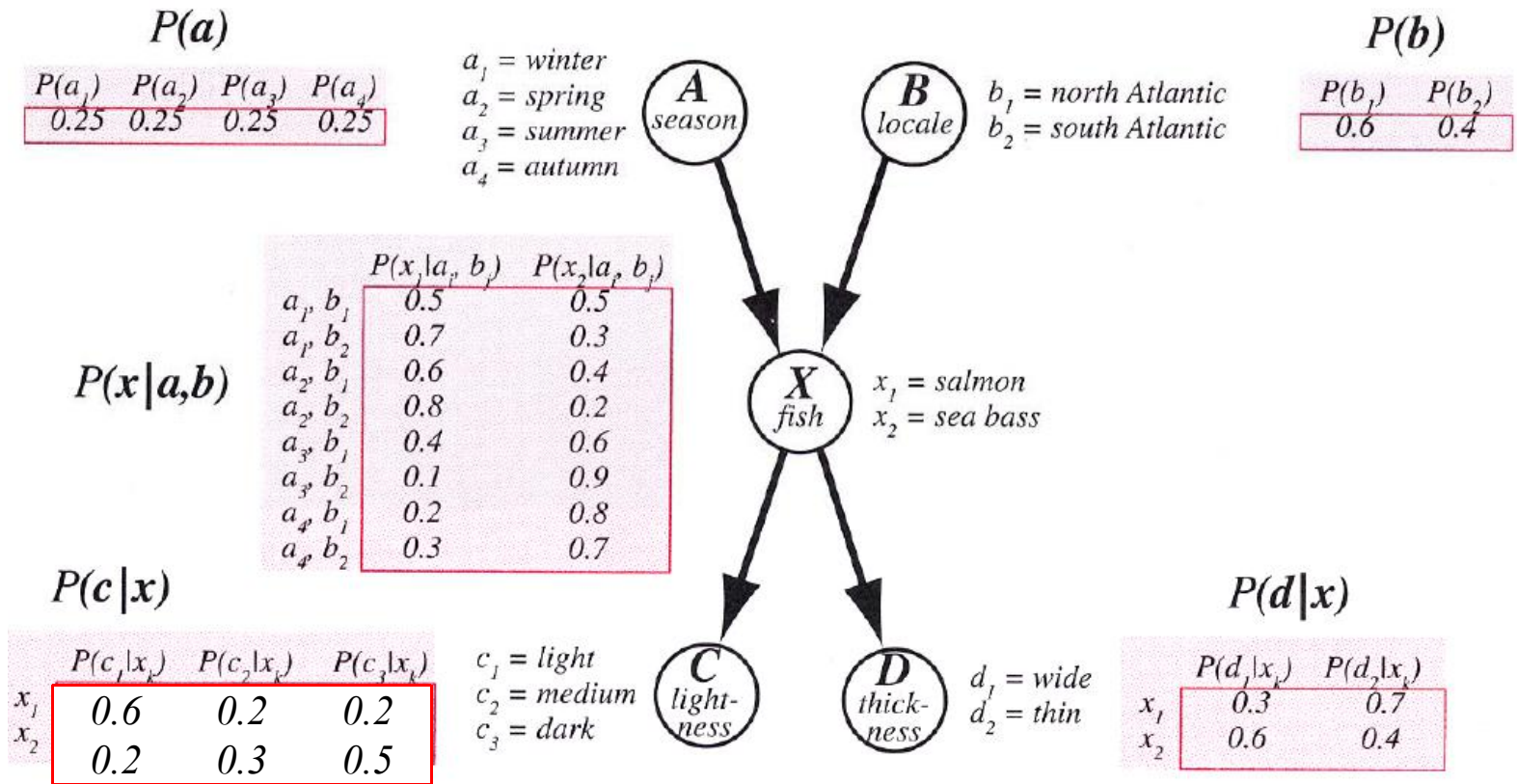
- Calculate the probability that "a fish caught in the North Atlantic in summer is a bass with a dim gloss and a narrow width"

- Summer:  $a_3$
- North Atlantic :  $b_1$
- dim gloss :  $c_3$
- narrow width :  $d_2$
- Bass:  $x_2$

$$\begin{aligned}P(a_3, b_1, x_2, c_3, d_2) &= P(a_3)P(b_1)P(x_2 | a_3, b_1)P(c_3 | x_2)P(d_2 | x_2) \\&= 0.25 \times 0.6 \times 0.6 \times 0.5 \times 0.4 \\&= 0.018\end{aligned}$$

# Example One

1. The probability of catching salmon in the South Atlantic in winter
2. The probability of catching bright bass in the South Atlantic
3. The probability of catching a wide, bright fish in the North Atlantic in summer



# Evidence

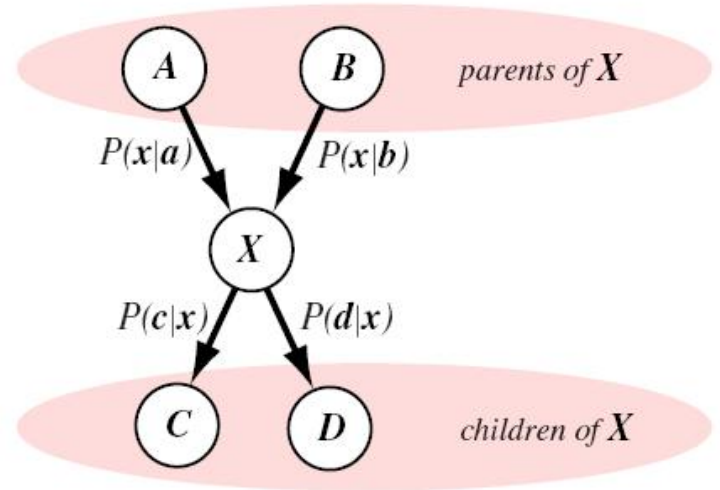
- Given the values of variables other than the target variable  $X$ , determine the probability of the other variables
- Evidence  $\{e_A, e_B, e_C, e_D\}$ , where  $e_i$  indicate the value of the variable
- For example, fish classification belief net
  - Available evidence  $\mathbf{e} = \{e_A, e_B, e_C, \dots\}$ 
    - $e_A$ : winter now  
 $P(a_1|e_A) = 1 \quad P(a_i|e_A) = 0 \text{ for } i = 2, 3, 4$
    - $e_B$ : fishermen prefer the South Atlantic  
 $P(b_1|e_B) = 0.2 \quad P(b_2|e_B) = 0.8$
    - $e_C$ : the fish has a lighter sheen  
 $P(e_C|c_1) = 1 \quad P(e_C|c_2) = 0.5 \quad P(e_C|c_3) = 0$
    - $e_D$ : the width cannot be measured because of occlusion  
 $P(e_D|d_1) = P(e_D|d_2)$

Pay attention to  
the position of  $e_i$ !



# Confidence

- Consider a certain node  $X$
- The set of nodes before  $X$  is called the **parent node**  $P$  of  $X$  and the set of nodes after  $X$  is called the **child node**  $C$  of  $X$
- For example:
  - parent node of  $X$ :  $\{A, B\}$
  - child node of  $X$  :  $\{C, D\}$
- When estimating the probability of  $X$ , the parent node and the child node of  $X$  should be treated differently
  - Evidence  $\mathbf{e}$ : Values of variables at nodes other than  $X$
  - Given  $\mathbf{e}$ , the Confidence Belief of  $x = (x_1, x_2, \dots)$ 
$$P(\mathbf{x}|\mathbf{e}) \propto P(\mathbf{e}^C|\mathbf{x})P(\mathbf{x}|\mathbf{e}^P)$$
  - Must be normalized so that the sum of the probabilities of all values of  $\mathbf{x}$  is 1



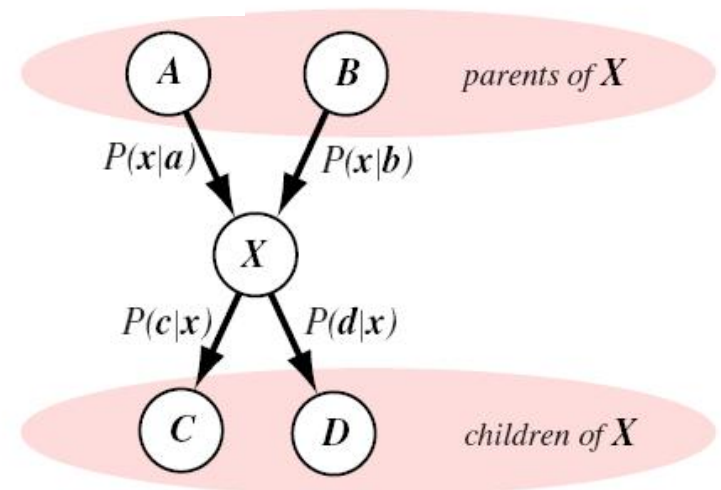
# Confidence

- For child node of  $X$ 
  - Suppose there are no connections between child nodes

$$\begin{aligned}P(\mathbf{e}^C | \mathbf{x}) &= P(\mathbf{e}_{C_1}, \mathbf{e}_{C_2}, \dots, \mathbf{e}_{C_{|C|}} | \mathbf{x}) \\&= P(\mathbf{e}_{C_1} | \mathbf{x}) P(\mathbf{e}_{C_2} | \mathbf{x}) \cdots P(\mathbf{e}_{C_{|C|}} | \mathbf{x}) \\&= \prod_{j=1}^{|C|} P(\mathbf{e}_{C_j} | \mathbf{x}),\end{aligned}$$

- For example:

$$P(\mathbf{e}_C, \mathbf{e}_D | \mathbf{x}) = P(\mathbf{e}_C | \mathbf{x}) P(\mathbf{e}_D | \mathbf{x})$$





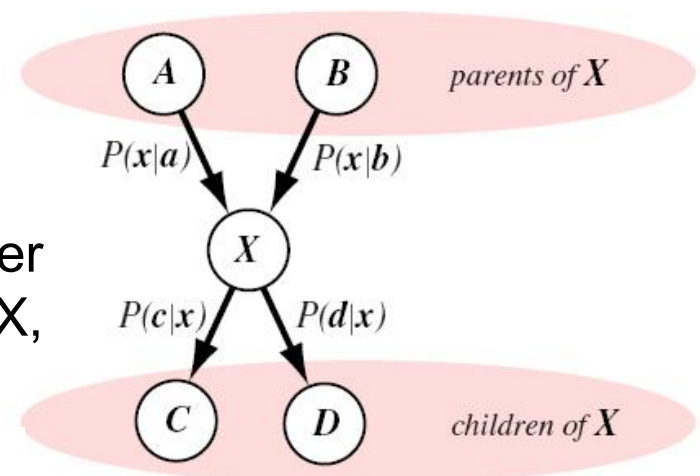
# Confidence

- For parent node of  $X$ 
  - Suppose there are no connections between parent nodes

$$\begin{aligned}
 P(x|e^{\mathcal{P}}) &= P(x|e_{\mathcal{P}_1}, e_{\mathcal{P}_2}, \dots, e_{\mathcal{P}_{|\mathcal{P}|}}) \\
 &= \sum_{\text{all } i,j,\dots,k} P(x|\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k}) P(\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k} | e_{\mathcal{P}_1}, \dots, e_{\mathcal{P}_{|\mathcal{P}|}}) \\
 &= \sum_{\text{all } i,j,\dots,k} P(x|\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k}) P(\mathcal{P}_{1i} | e_{\mathcal{P}_1}) \cdots P(\mathcal{P}_{|\mathcal{P}|k} | e_{\mathcal{P}_{|\mathcal{P}|}}),
 \end{aligned}$$

- $\mathcal{P}_{mn}$  represents the value of the parent node  $\mathcal{P}_m$  in state  $n$
- Ignore node interdependencies other than the parent and child nodes of  $X$ , and simplify above equation

$$P(x|e^{\mathcal{P}}) = \sum_{\text{all } \mathcal{P}_{mn}} P(x|\mathcal{P}_{mn}) \prod_{i=1}^{|\mathcal{P}|} P(\mathcal{P}_i | e_{\mathcal{P}_i})$$

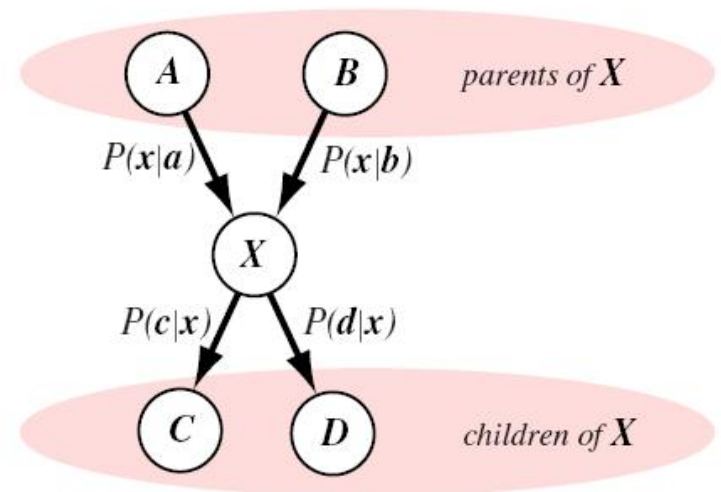


# Confidence

- The confidence of proposition  $X$

$$P(\mathbf{x}|\mathbf{e}) \propto \underbrace{\prod_{j=1}^{|\mathcal{C}|} P(\mathbf{e}_{\mathcal{C}_j}|\mathbf{x})}_{P(\mathbf{e}^{\mathcal{C}}|\mathbf{x})} \underbrace{\left[ \sum_{\text{all } \mathcal{P}_{mn}} P(\mathbf{x}|\mathcal{P}_{mn}) \prod_{i=1}^{|\mathcal{P}|} P(\mathcal{P}_i|\mathbf{e}_{\mathcal{P}_i}) \right]}_{P(\mathbf{x}|\mathbf{e}^{\mathcal{P}})}$$

- The probability that node  $X$  takes a particular value is equal to the product of two factors
  - The first depends on the child nodes
  - The first depends on the parent nodes



# Evidence

---

- Simple cases
  - $e_i$  directly represent the value of variable
  - Confidence

$$P(x | \mathbf{e}) = \frac{P(x, \mathbf{e})}{P(\mathbf{e})} = \alpha P(x, \mathbf{e})$$

For a fixed  $\mathbf{e}$ ,  $\alpha$  is a constant

# Example One

- fish classification belief net

$P(a)$

$P(a_1)$	$P(a_2)$	$P(a_3)$	$P(a_4)$
0.25	0.25	0.25	0.25

$a_1 = \text{winter}$   
 $a_2 = \text{spring}$   
 $a_3 = \text{summer}$   
 $a_4 = \text{autumn}$



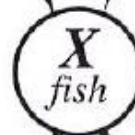
$b_1 = \text{north Atlantic}$   
 $b_2 = \text{south Atlantic}$

$P(b)$

$P(b_1)$	$P(b_2)$
0.6	0.4

$P(x|a,b)$

	$P(x_1 a_i, b_j)$	$P(x_2 a_i, b_j)$
$a_1, b_1$	0.5	0.5
$a_1, b_2$	0.7	0.3
$a_2, b_1$	0.6	0.4
$a_2, b_2$	0.8	0.2
$a_3, b_1$	0.4	0.6
$a_3, b_2$	0.1	0.9
$a_4, b_1$	0.2	0.8
$a_4, b_2$	0.3	0.7



$x_1 = \text{salmon}$   
 $x_2 = \text{sea bass}$

$P(c|x)$

	$P(c_1 x_k)$	$P(c_2 x_k)$	$P(c_3 x_k)$
$x_1$	0.6	0.2	0.2
$x_2$	0.2	0.3	0.5

$c_1 = \text{light}$   
 $c_2 = \text{medium}$   
 $c_3 = \text{dark}$



$d_1 = \text{wide}$   
 $d_2 = \text{thin}$

$P(d|x)$

	$P(d_1 x_k)$	$P(d_2 x_k)$
$x_1$	0.3	0.7
$x_2$	0.6	0.4

# Example One

- Catch a shiny fish in the South Atlantic, salmon or bass?

$P(a)$

$P(a_1)$	$P(a_2)$	$P(a_3)$	$P(a_4)$
0.25	0.25	0.25	0.25

$a_1 = \text{winter}$   
 $a_2 = \text{spring}$   
 $a_3 = \text{summer}$   
 $a_4 = \text{autumn}$

$P(b)$

$P(b_1)$	$P(b_2)$
0.6	0.4

$b_1 = \text{north Atlantic}$   
 $b_2 = \text{south Atlantic}$

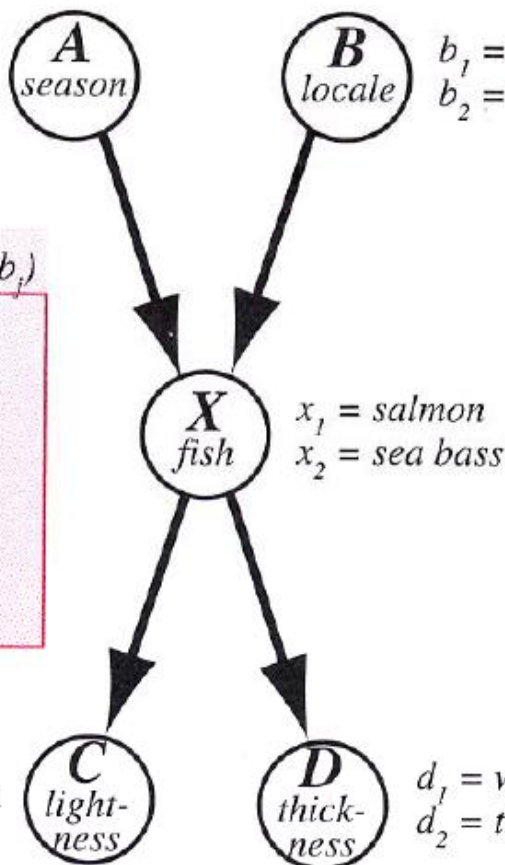
$P(x|a,b)$

	$P(x_1 a_i, b_j)$	$P(x_2 a_i, b_j)$
$a_1, b_1$	0.5	0.5
$a_1, b_2$	0.7	0.3
$a_2, b_1$	0.6	0.4
$a_2, b_2$	0.8	0.2
$a_3, b_1$	0.4	0.6
$a_3, b_2$	0.1	0.9
$a_4, b_1$	0.2	0.8
$a_4, b_2$	0.3	0.7

$P(c|x)$

	$P(c_1 x_k)$	$P(c_2 x_k)$	$P(c_3 x_k)$
$x_1$	0.6	0.2	0.2
$x_2$	0.2	0.3	0.5

$c_1 = \text{light}$   
 $c_2 = \text{medium}$   
 $c_3 = \text{dark}$



$P(d|x)$

	$P(d_1 x_k)$	$P(d_2 x_k)$
$x_1$	0.3	0.7
$x_2$	0.6	0.4

$d_1 = \text{wide}$   
 $d_2 = \text{thin}$

# Example One

- Catch a shiny fish in the South Atlantic, salmon or bass?

**e:** a is unknown                       $b_2$ =South Atlantic  
 $c_1$ =shinny                              d is unknown

$$\begin{aligned} P(x_1 | \mathbf{e}) &= \alpha P(x_1, b_2, c_1) \\ &= \alpha \sum_{a,d} P(x_1, a, b_2, c_1, d) \\ &= \alpha \sum_{a,d} P(a)P(b_2)P(x_1 | a, b_2)P(c_1 | x_1)P(d | x_1) \\ &= \alpha P(b_2)P(c_1 | x_1) \sum_a P(a)P(x_1 | a, b_2) \sum_d P(d | x_1) \\ &= \alpha P(b_2)P(c_1 | x_1) \times [P(a_1)P(x_1 | a_1, b_2) + P(a_2)P(x_1 | a_2, b_2) \\ &\quad + P(a_3)P(x_1 | a_3, b_2) + P(a_4)P(x_1 | a_4, b_2)] \times [P(d_1 | x_1) + P(d_2 | x_1)] \\ &= \alpha \times 0.4 \times 0.6 \times [0.25 \times 0.7 + 0.25 \times 0.8 + 0.25 \times 0.1 + 0.25 \times 0.3] \times 1.0 \\ &= 0.114\alpha \end{aligned}$$

**First calculate the  
probability of  
 $x_1$ =salmon**

# Example One

---

- Catch a shiny fish in the South Atlantic, salmon or bass?

a is unknown

$b_2$ =South Atlantic

$c_1$ =shinny

d is unknown

$$P(x_2 | \mathbf{e}) = 0.042\alpha$$

**Then calculate the  
probability of  $x_2$ =bass**

- normalization (make  $P(x_1 | \mathbf{e}) + P(x_2 | \mathbf{e}) = 1$  )

$$P(x_1 | \mathbf{e}) = 0.63$$

$$P(x_2 | \mathbf{e}) = 0.27$$

Because  $P(x_1 | \mathbf{e}) > P(x_2 | \mathbf{e})$  , so it's salmon



# Example Two

---

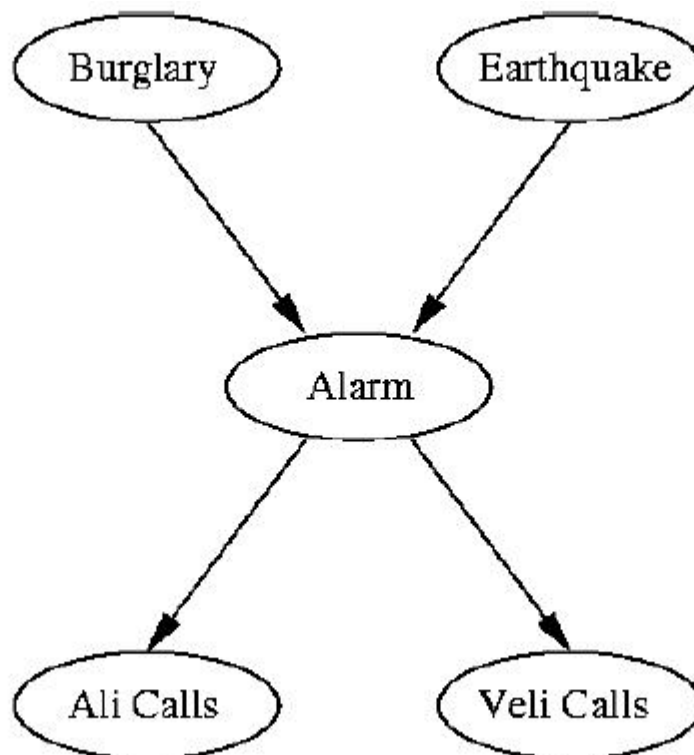
- You have installed an anti - theft system in your house
- The system is sensitive to burglary detection, but sometimes earthquakes can trigger alarms
- You have two neighbors: Ali and Veli. When you are not at home, they would call you if they heard the alarm
- Ali would call you when he heard the alarm, but sometimes he would call you because he thought phone ringing was an alarm
- Veli often listens to music at home, so sometimes she doesn't hear the alarm
- Can you estimate the true probability of a real burglary based on which neighbor called you?



# Example Two

- Modeling

$P(B=T)$	$P(B=F)$
0.001	0.999



$P(E=T)$	$P(E=F)$
0.002	0.998

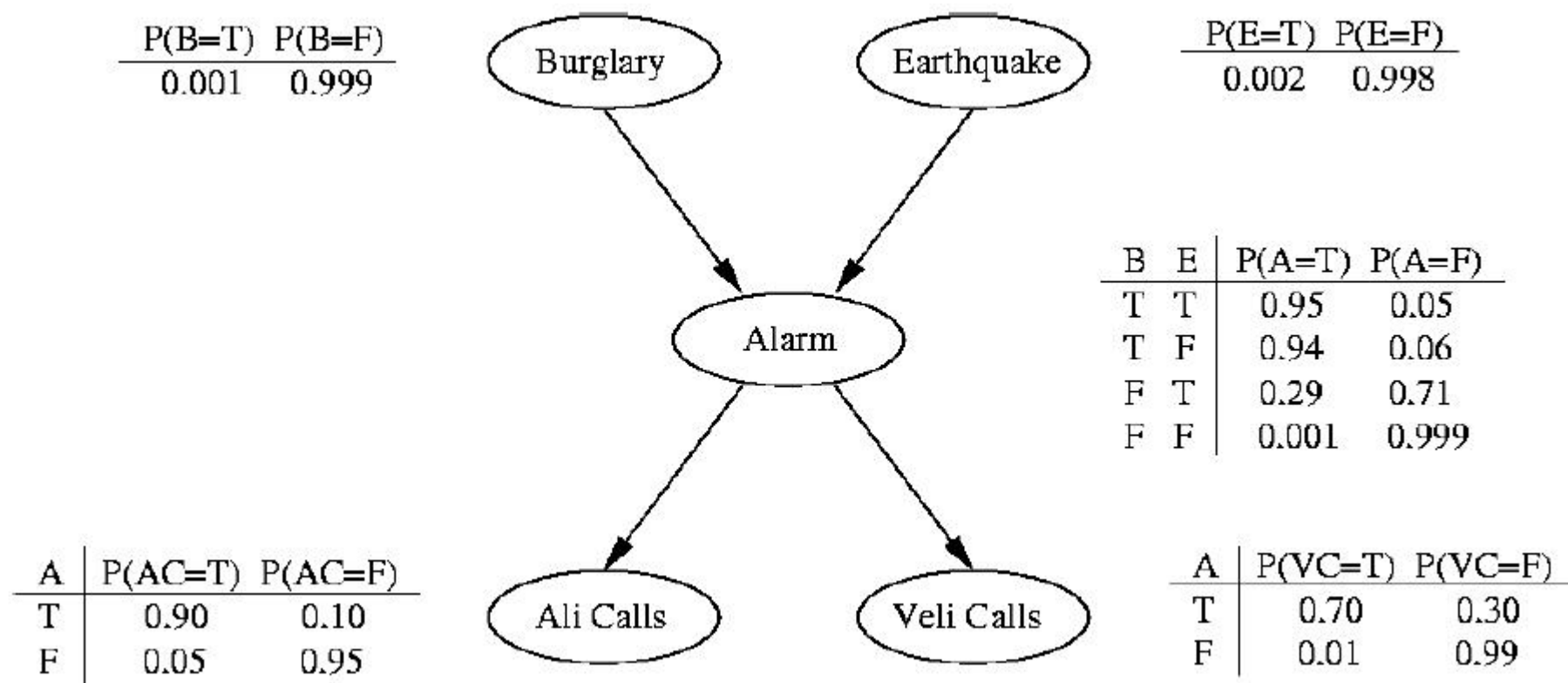
B	E	$P(A=T)$	$P(A=F)$
T	T	0.95	0.05
T	F	0.94	0.06
F	T	0.29	0.71
F	F	0.001	0.999

A	$P(AC=T)$	$P(AC=F)$
T	0.90	0.10
F	0.05	0.95

A	$P(VC=T)$	$P(VC=F)$
T	0.70	0.30
F	0.01	0.99

# Example Two

- The system alarms, but neither burglary nor earthquake occurred, besides both Ali and Veli call you



# Example Two

---

- Calculate the probability of the following event
  - The system alarms, but neither burglary nor earthquake occurred, besides both Ali and Veli call you

$$P(AC, VC, A, \neg B, \neg E)$$

$$= P(AC|A)P(VC|A)P(A|\neg B, \neg E)P(\neg B)P(\neg E)$$

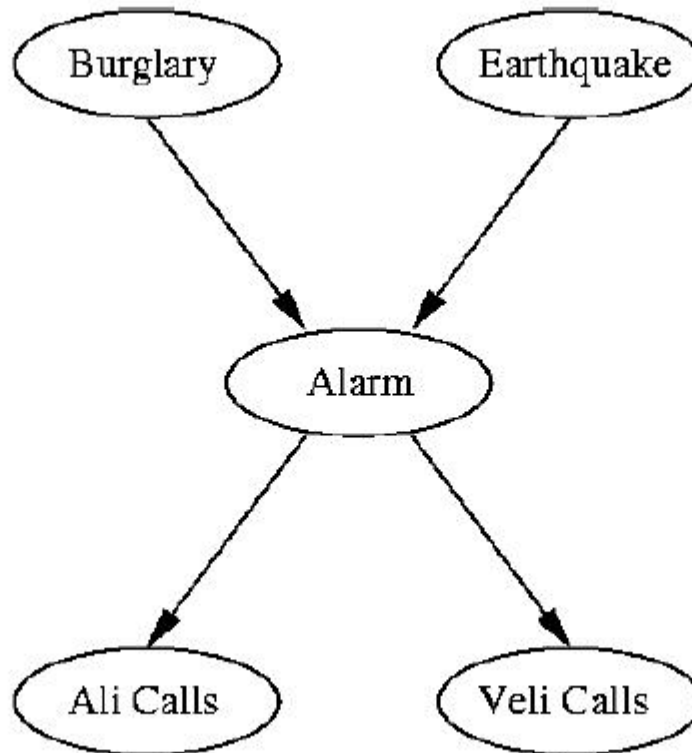
$$= 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998$$

$$= 0.00062$$

# Example Two

- If Ali calls you, calculate the confidence that a burglary has occurred

$P(B=T)$	$P(B=F)$
0.001	0.999



$P(E=T)$	$P(E=F)$
0.002	0.998

B	E	$P(A=T)$	$P(A=F)$
T	T	0.95	0.05
T	F	0.94	0.06
F	T	0.29	0.71
F	F	0.001	0.999

A	$P(AC=T)$	$P(AC=F)$
T	0.90	0.10
F	0.05	0.95

A	$P(VC=T)$	$P(VC=F)$
T	0.70	0.30
F	0.01	0.99

# Example Two

---

- If Ali calls you, calculate the confidence that a burglary has occurred

- **Method One**

$$P(B \mid AC) = \alpha P(B, AC)$$

$$= \alpha \sum_{vc} \sum_a \sum_e P(AC \mid a) P(vc \mid a) P(a \mid B, e) P(B) P(e)$$

$$= 0.00084632 \alpha$$

$$P(\neg B \mid AC) = \alpha P(\neg B, AC)$$

$$= \alpha \sum_{vc} \sum_a \sum_e P(AC \mid a) P(vc \mid a) P(a \mid \neg B, e) P(\neg B) P(e)$$

$$= 0.0513 \alpha$$

$$\text{Normalization } P(B \mid AC) = \frac{0.00084632 \alpha}{0.00084632 \alpha + 0.0513 \alpha} = 0.0162$$

# Example Two

---

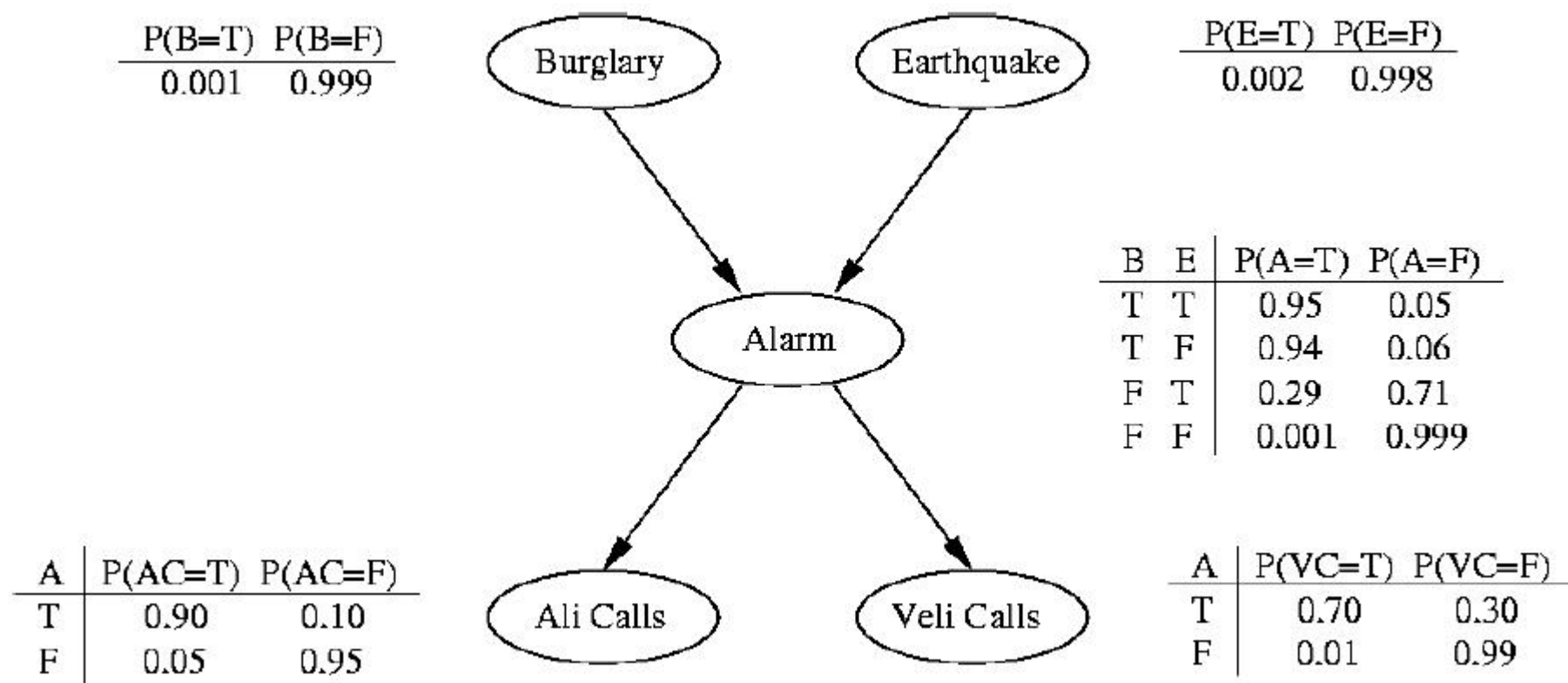
- If Ali calls you, calculate the confidence that a burglary has occurred

- **Method Two**

$$\begin{aligned}P(B|AC) &= \frac{P(B, AC)}{P(AC)} \\&= \frac{\sum_{vc} \sum_a \sum_e P(AC|a)P(vc|a)P(a|B, e)P(B)P(e)}{P(B, AC) + P(\neg B, AC)} \\&= \frac{0.00084632}{0.00084632 + 0.0513} \\&= 0.0162\end{aligned}$$

# Example Two

- If Both Ali and Veli call you, calculate the confidence that a burglary has occurred



# Example Two

---

- If Both Ali and Veli call you, calculate the confidence that a burglary has occurred

$$\begin{aligned}P(B \mid AC, VC) &= \frac{P(B, AC, VC)}{P(AC, VC)} \\&= \frac{\sum_a \sum_e P(AC \mid a)P(VC \mid a)P(a \mid B, e)P(B)P(e)}{P(B, AC, VC) + P(\neg B, AC, VC)} \\&= 0.29\end{aligned}$$



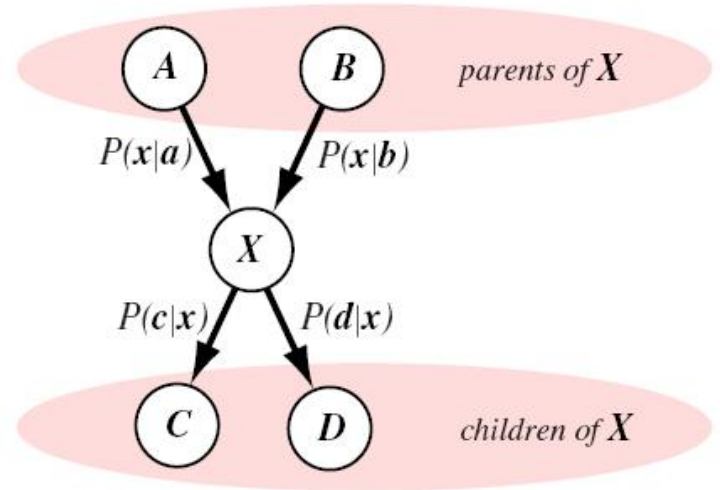
# Evidence

- Given the values of variables other than the target variable  $X$ , determine the probability of the other variables
- Evidence  $\{e_A, e_B, e_C, e_D\}$ , where  $e_i$  indicate the value of the variable
- For example, fish classification belief net
  - Available evidence  $\mathbf{e} = \{e_A, e_B, e_C, \dots\}$ 
    - $e_A$ : winter now  
 $P(a_1|e_A) = 1 \quad P(a_i|e_A) = 0 \text{ for } i = 2, 3, 4$
    - $e_B$ : fishermen prefer the South Atlantic  
 $P(b_1|e_B) = 0.2 \quad P(b_2|e_B) = 0.8$
    - $e_C$ : the fish has a lighter sheen  
 $P(e_C|c_1) = 1 \quad P(e_C|c_2) = 0.5 \quad P(e_C|c_3) = 0$
    - $e_D$ : the width cannot be measured because of occlusion  
 $P(e_D|d_1) = P(e_D|d_2)$

Pay attention to  
the position of  $e_i$ !

# Confidence

- Consider a certain node  $X$
- The set of nodes before  $X$  is called the **parent node**  $P$  of  $X$  and the set of nodes after  $X$  is called the **child node**  $C$  of  $X$
- For example:
  - parent node of  $X$ :  $\{A, B\}$
  - child node of  $X$  :  $\{C, D\}$
- When estimating the probability of  $X$ , the parent node and the child node of  $X$  should be treated differently
  - Evidence  $\mathbf{e}$ : Values of variables at nodes other than  $X$
  - Given  $\mathbf{e}$ , the Confidence Belief of  $x = (x_1, x_2, \dots)$ 
$$P(\mathbf{x}|\mathbf{e}) \propto P(\mathbf{e}^C|\mathbf{x})P(\mathbf{x}|\mathbf{e}^P)$$
  - Must be normalized so that the sum of the probabilities of all values of  $\mathbf{x}$  is 1



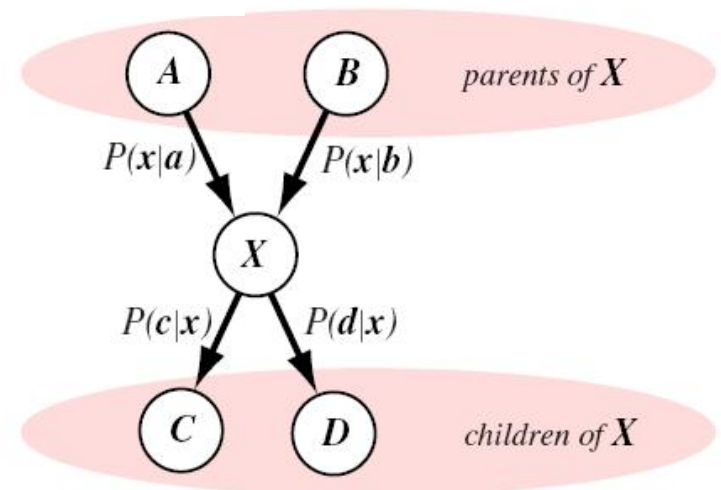
# Confidence

- For child node of  $X$ 
  - Suppose there are no connections between child nodes

$$\begin{aligned}P(\mathbf{e}^C | \mathbf{x}) &= P(\mathbf{e}_{C_1}, \mathbf{e}_{C_2}, \dots, \mathbf{e}_{C_{|C|}} | \mathbf{x}) \\&= P(\mathbf{e}_{C_1} | \mathbf{x}) P(\mathbf{e}_{C_2} | \mathbf{x}) \cdots P(\mathbf{e}_{C_{|C|}} | \mathbf{x}) \\&= \prod_{j=1}^{|C|} P(\mathbf{e}_{C_j} | \mathbf{x}),\end{aligned}$$

- For example:

$$P(\mathbf{e}_C, \mathbf{e}_D | \mathbf{x}) = P(\mathbf{e}_C | \mathbf{x}) P(\mathbf{e}_D | \mathbf{x})$$



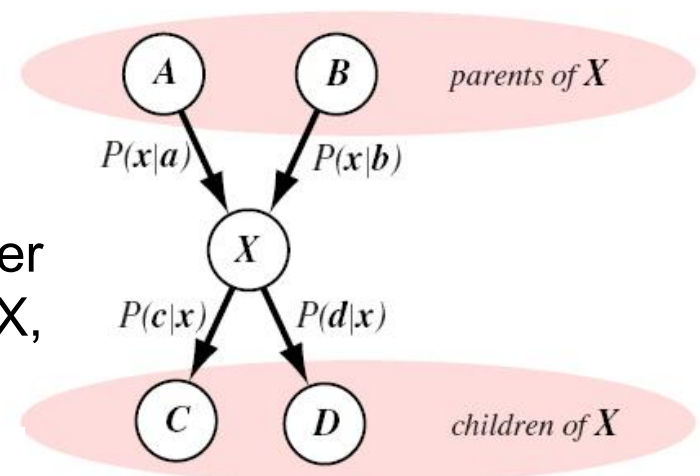
# Confidence

- For parent node of  $X$ 
  - Suppose there are no connections between parent nodes

$$\begin{aligned}
 P(x|e^{\mathcal{P}}) &= P(x|e_{\mathcal{P}_1}, e_{\mathcal{P}_2}, \dots, e_{\mathcal{P}_{|\mathcal{P}|}}) \\
 &= \sum_{\text{all } i,j,\dots,k} P(x|\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k}) P(\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k} | e_{\mathcal{P}_1}, \dots, e_{\mathcal{P}_{|\mathcal{P}|}}) \\
 &= \sum_{\text{all } i,j,\dots,k} P(x|\mathcal{P}_{1i}, \mathcal{P}_{2j}, \dots, \mathcal{P}_{|\mathcal{P}|k}) P(\mathcal{P}_{1i} | e_{\mathcal{P}_1}) \cdots P(\mathcal{P}_{|\mathcal{P}|k} | e_{\mathcal{P}_{|\mathcal{P}|}}),
 \end{aligned}$$

- $\mathcal{P}_{mn}$  represents the value of the parent node  $\mathcal{P}_m$  in state  $n$
- Ignore node interdependencies other than the parent and child nodes of  $X$ , and simplify above equation

$$P(x|e^{\mathcal{P}}) = \sum_{\text{all } \mathcal{P}_{mn}} P(x|\mathcal{P}_{mn}) \prod_{i=1}^{|\mathcal{P}|} P(\mathcal{P}_i | e_{\mathcal{P}_i})$$

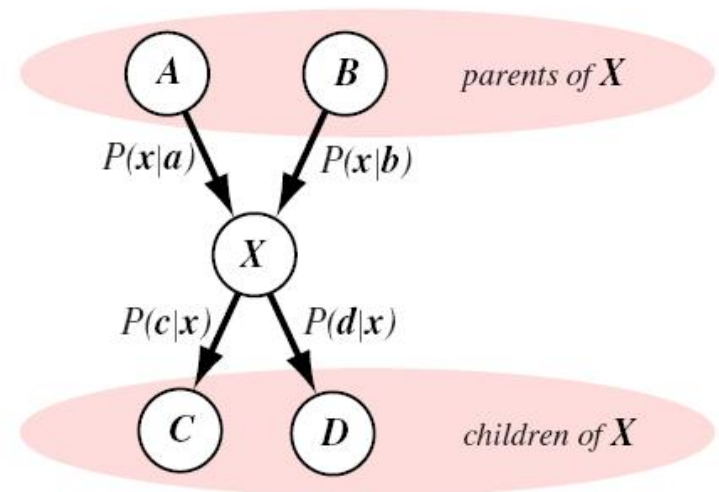


# Confidence

- The confidence of proposition  $X$

$$P(\mathbf{x}|\mathbf{e}) \propto \underbrace{\prod_{j=1}^{|\mathcal{C}|} P(\mathbf{e}_{\mathcal{C}_j}|\mathbf{x})}_{P(\mathbf{e}^{\mathcal{C}}|\mathbf{x})} \underbrace{\left[ \sum_{\text{all } \mathcal{P}_{mn}} P(\mathbf{x}|\mathcal{P}_{mn}) \prod_{i=1}^{|\mathcal{P}|} P(\mathcal{P}_i|\mathbf{e}_{\mathcal{P}_i}) \right]}_{P(\mathbf{x}|\mathbf{e}^{\mathcal{P}})}$$

- The probability that node  $X$  takes a particular value is equal to the product of two factors
  - The first depends on the child nodes
  - The first depends on the parent nodes



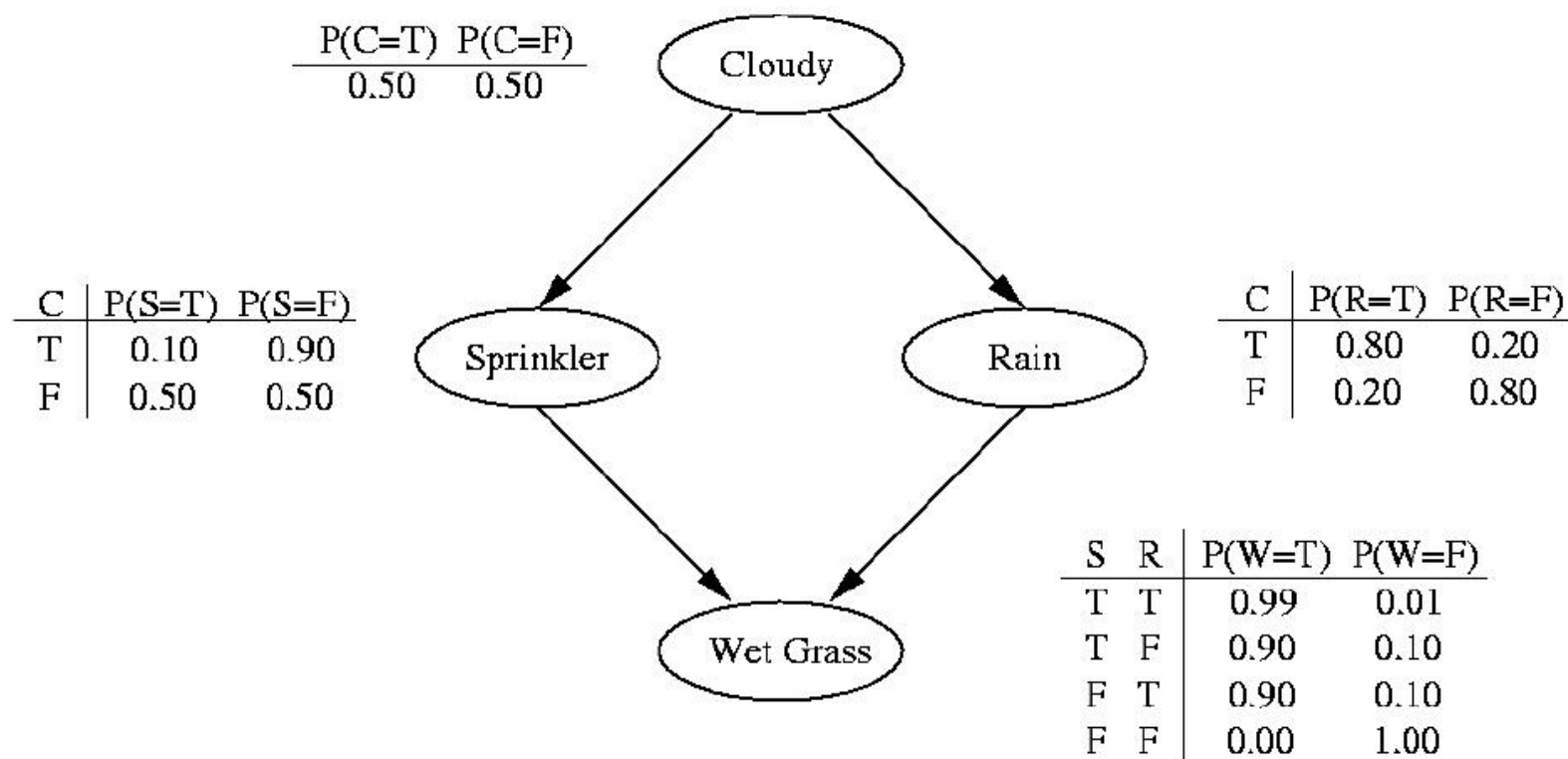
# Example Three

---

- The grass can get wet for two reasons: the sprinklers have been turned on, or it has rained
- If it's cloudy, it's more likely to rain than it is sunny
- If it is cloudy, it's less likely to turn on the sprinklers
- Suppose it's equally likely to be cloudy or sunny

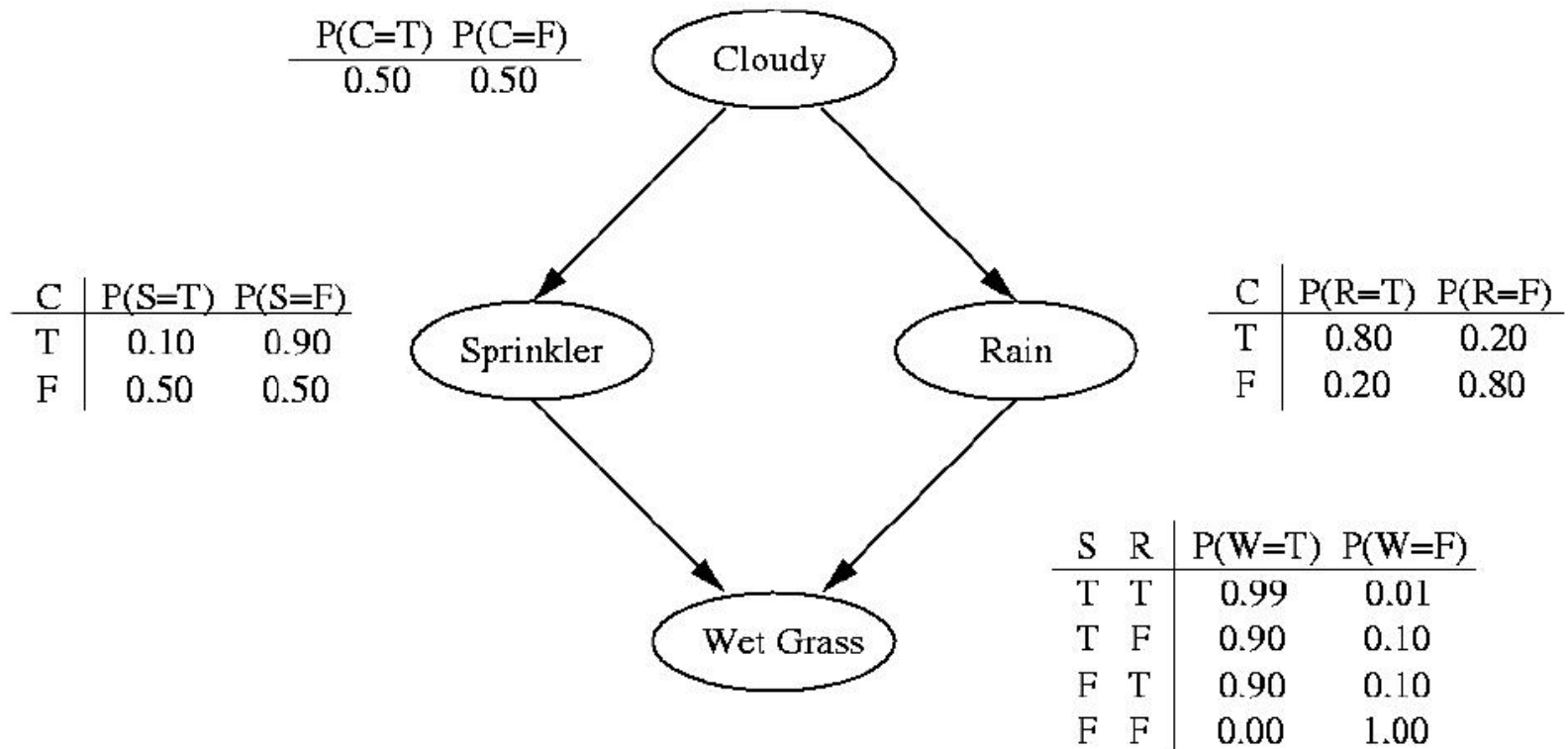
# Example Three

- Modeling



# Example Three

- If you see that the grass is wet, which reason is more likely to be the sprinklers or the rain?





# Example Three

---

- If you see that the grass is wet, which reason is more likely to be the sprinklers or the rain?

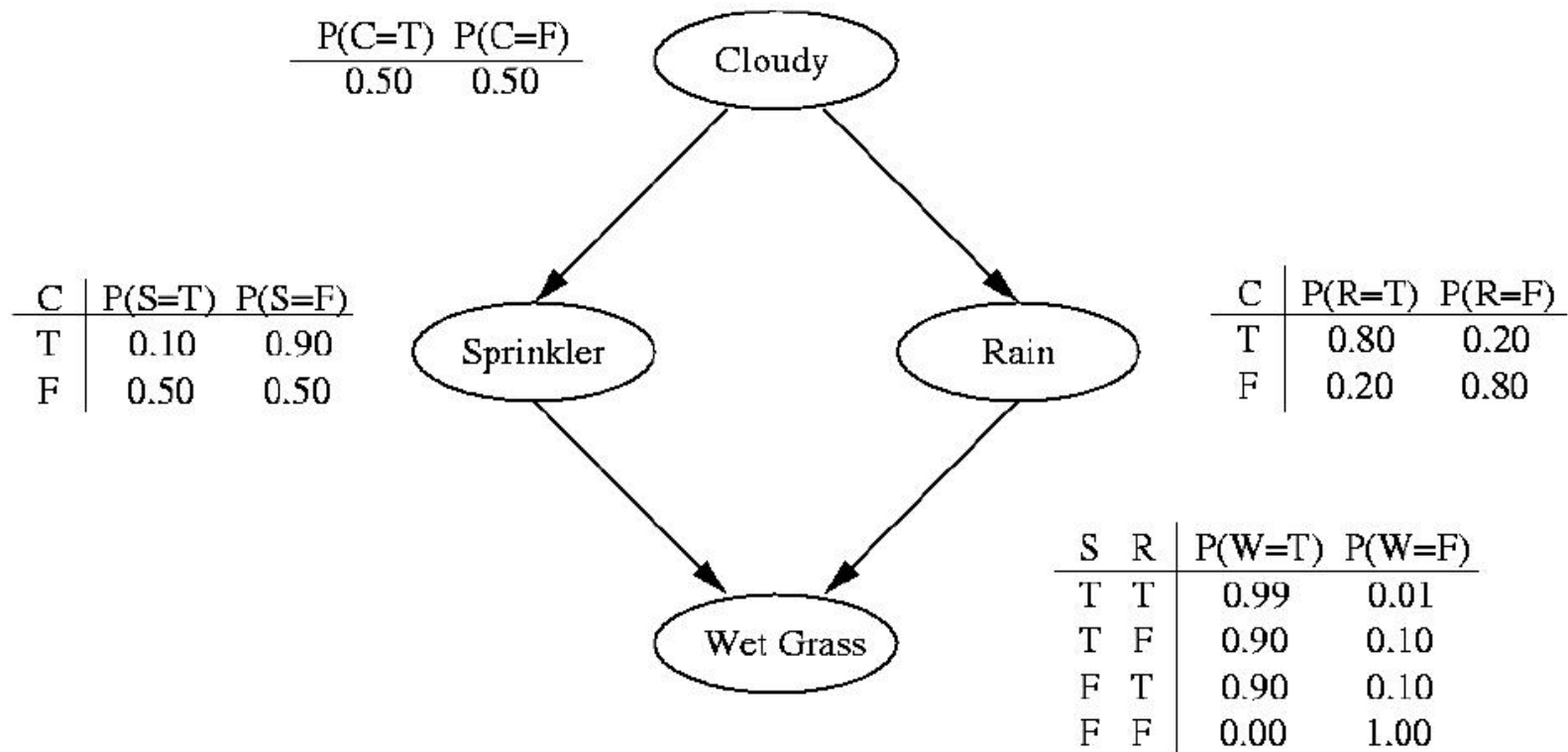
$$\begin{aligned} P(S | W) &= \frac{P(S, W)}{P(W)} = \frac{\sum_c \sum_r P(W | S, r) P(S | c) P(r | c) P(c)}{P(S, W) + P(\neg S, W)} \\ &= \frac{0.2781}{0.2781 + 0.369} = 0.430 \end{aligned}$$

$$\begin{aligned} P(R | W) &= \frac{P(R, W)}{P(W)} = \frac{\sum_c \sum_s P(W | s, R) P(s | c) P(R | c) P(c)}{P(R, W) + P(\neg R, W)} \\ &= \frac{0.4581}{0.6471} = 0.708 \end{aligned}$$

**Because  $P(S | W) < P(R | W)$ , So it is more likely that the grass would be wet by rain**

# Example Three

- What if the grass is wet and the weather is sunny?



# Example Three

---

- What if the grass is wet and the weather is sunny?

$$P(S | W, \neg C) = \frac{P(S, W, \neg C)}{P(W, \neg C)} = \frac{P(S | \neg C)P(\neg C) \sum_r P(W | S, r)P(r | \neg C)}{P(S, W, \neg C) + P(\neg S, W, \neg C)}$$
$$= \frac{0.2295}{0.2295 + 0.045} = 0.836$$

$$P(R | W, \neg C) = \frac{P(R, W, \neg C)}{P(W, \neg C)} = \frac{P(R | \neg C)P(\neg C) \sum_s P(W | s, R)P(s | \neg C)}{P(R, W, \neg C) + P(\neg R, W, \neg C)}$$
$$= \frac{0.0945}{0.2745} = 0.344$$

**Because**  $P(S | W, \neg C) > P(R | W, \neg C)$ ,  
**So it is more likely that the grass**  
**would be wet by sprinkling**

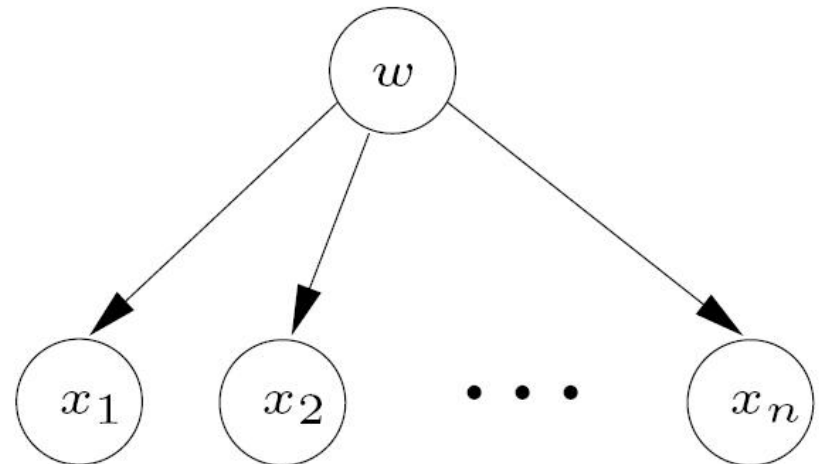
# Naïve Bayes Rule

---

- When the dependencies between features are unknown, it is often assumed that each feature is conditionally independent under a given category condition

$$p(\mathbf{x} | \omega_k) = \prod_{i=1}^d p(x_i | \omega_k)$$

- assumption is called naïve Bayes rule (朴素贝叶斯规则) or idiot Bayes rule (傻瓜贝叶斯规则)
- Naïve Bayes belief net



# Part 3 Expected Maximization (EM) algorithm

# Missing Feature

---

- Suppose there is a Bayesian classifier based on the eigenvector  $\mathbf{x}$ , in which a part of the feature  $\mathbf{x}_g$  is visible in each  $\mathbf{x}$  and the rest of the feature  $\mathbf{x}_b$  is missing. The missing features may be different in different samples
- How to make decisions?
  - Method One
    - simply throw away the sample containing the missing value
  - Method Two
    - replace  $\mathbf{x}_b$  with some other sample mean  $\bar{\mathbf{x}}_b$  known to have this feature, that is  $\mathbf{x} = (\mathbf{x}_g \bar{\mathbf{x}}_b)$

# Expectation Maximization

---

- **Mehod Three**

- By extending the maximum likelihood method, the model parameters can be learned in the case of missing values in the training set
- **expectation-maximization, EM** (期望最大化) Algorithm can estimate likelihood function based on existing data recursively

- Two main applications of EM algorithm

- Learn when the data is incomplete or has missing values
- When the likelihood equation is difficult to solve directly, initialize some unknown parameters to simplify the problem

# Expectation Maximization

---

- Suppose sample  $\mathbf{x}$  obeys a certain distribution  $p(\mathbf{x} | \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is the unknown parameter vector
- Sample set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where  $\mathbf{x}_k = \{\mathbf{x}_{kg}, \mathbf{x}_{kb}\}$ ,  $\mathbf{x}_{kg}$  is the complete (or good) part and  $\mathbf{x}_{kb}$  is the missing (or damaged) part
- represent the set of  $\mathbf{x}_{kg}$  and  $\mathbf{x}_{kb}$  separately  $\mathcal{D} = \mathcal{D}_g \cup \mathcal{D}_b$
- Given some assumption  $\boldsymbol{\theta}^i$  (not necessarily accurate) for  $\boldsymbol{\theta}$ , calculate the expectation of the log-likelihood function under a distribution determined by  $\boldsymbol{\theta}^i$  about missing features, and get a function of  $\boldsymbol{\theta}$

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i) &= \mathcal{E}_{\mathcal{D}_b} [\ln p(\mathcal{D}_g, \mathcal{D}_b; \boldsymbol{\theta}) | \mathcal{D}_g; \boldsymbol{\theta}^i] \\ &= \int_{-\infty}^{+\infty} \ln p(D_g, D_b | \boldsymbol{\theta}) p(D_b | \boldsymbol{\theta}^i; D_g) dD_b \end{aligned}$$



# Expectation Maximization

---

$$\begin{aligned} Q(\theta; \theta^i) &= \mathcal{E}_{\mathcal{D}_b} [\ln p(\mathcal{D}_g, \mathcal{D}_b; \theta) | \mathcal{D}_g; \theta^i] \\ &= \int_{-\infty}^{+\infty} \ln p(D_g, D_b | \theta) p(D_b | \theta^i; D_g) dD_b \end{aligned}$$

- **Meaning:** the parameter vector  $\theta^i$  is the current estimate of the parameter and  $\theta$  is a candidate parameter vector to improve the current estimate. For each  $\theta$ , the log-likelihood function of the training set can be computed. Due to the presence of missing values  $\mathcal{D}_b$ , the likelihood function needs to marginalize  $\mathcal{D}_b$ , and the distribution of  $\mathcal{D}_b$  in the edge integral is determined by the current estimate  $\theta^i$

# Expectation Maximization

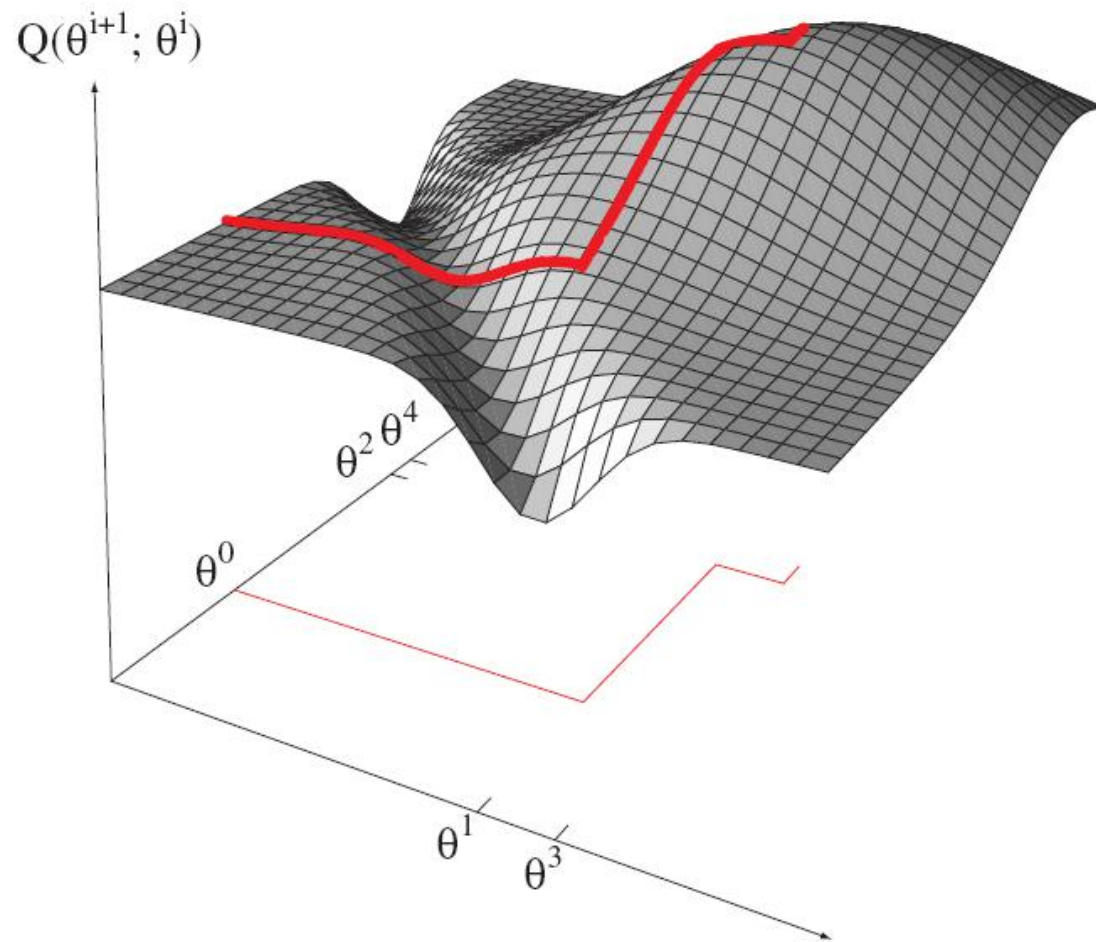
---

- **expectation-maximization(EM) algorithm**

```
1 begin initialize  $\theta^0, T, i = 0$   
2       do  $i \leftarrow i + 1$   
3           E step: compute  $Q(\theta; \theta^i)$   
4           M step:  $\theta^{i+1} \leftarrow \arg \max_{\theta} Q(\theta; \theta^i)$   
5           until  $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \leq T$   
6       return  $\hat{\theta} \leftarrow \theta^{i+1}$   
7 end
```

# Expectation Maximization

---



# Example

---

- EM algorithm for two-dimensional normal distribution
  - Dataset  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\}$
  - Probabilistic Model
    - two-dimensional normal distribution,  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$
- Target: To estimate the parameter vectors of a normal distribution

$$\theta = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \sigma_1^2 \\ \sigma_2^2 \end{pmatrix}$$

# Example

---

- Solution

- initialization
- The mean is at the origin, the covariance matrix is the identity matrix, that is
- Calculate  $\theta^1$

$$\theta^0 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$

- **Step E**

# Example

---

$$\begin{aligned}
 Q(\boldsymbol{\theta}; \boldsymbol{\theta}^0) &= \mathcal{E}_{x_{41}} [\ln p(\mathbf{x}_g, \mathbf{x}_b; \boldsymbol{\theta} | \boldsymbol{\theta}^0; \mathcal{D}_g)] \\
 &= \int_{-\infty}^{\infty} \left[ \sum_{k=1}^3 \ln p(\mathbf{x}_k | \boldsymbol{\theta}) + \ln p(\mathbf{x}_4 | \boldsymbol{\theta}) \right] p(x_{41} | \boldsymbol{\theta}^0; x_{42} = 4) dx_{41} \\
 &= \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] + \int_{-\infty}^{\infty} \ln p \left( \binom{x_{41}}{4} \middle| \boldsymbol{\theta} \right) \underbrace{\frac{p \left( \binom{x_{41}}{4} \middle| \boldsymbol{\theta}^0 \right)}{\left( \int_{-\infty}^{\infty} p \left( \binom{x'_{41}}{4} \middle| \boldsymbol{\theta}^0 \right) dx'_{41} \right)}}_{\equiv K} dx_{41} \\
 &= \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] + \frac{1}{K} \int_{-\infty}^{\infty} \ln p \left( \binom{x_{41}}{4} \middle| \boldsymbol{\theta} \right) \frac{1}{2\pi \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix}} \exp \left[ -\frac{1}{2}(x_{41}^2 + 4^2) \right] dx_{41} \\
 &= \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] - \frac{1 + \mu_1^2}{2\sigma_1^2} - \frac{(4 - \mu_2)^2}{2\sigma_2^2} - \ln (2\pi\sigma_1\sigma_2).
 \end{aligned}$$

# Example

---

- **Step M**

$$\nabla_{\theta} Q(\theta | \theta^0) = 0$$

$$\theta^1 = \begin{pmatrix} 0.75 \\ 2.0 \\ 0.938 \\ 2.0 \end{pmatrix}$$

- Iterate Step E and Step M, and after 3 iterations  $\theta$  converges to

$$\mu = \begin{bmatrix} 1.0 \\ 2.0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 0.667 & 0 \\ 0 & 2.0 \end{bmatrix}$$

# Example

---

