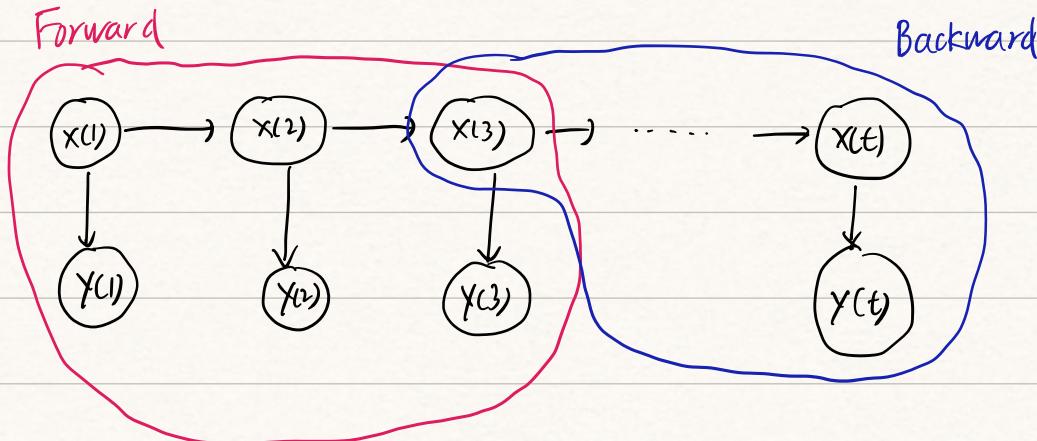


Recall  $\rightarrow$  Forward & Backward Algorithm  $\Rightarrow$  calculate  $P(y_1, \dots, y_n)$

$$\begin{cases} \alpha_t(j) = P(y_1, \dots, y_t, X(t)=j), j \in S \sim \text{forward message} \\ \beta_t(j) = P(y_{t+1}, \dots, y_n | X(t)=j), j \in S \sim \text{Backward message} \end{cases}$$



Task 3: (B-W Alg.)

Given observation from HMM. We want to estimate the paras

$(y_1, \dots, y_n)$

of the HMM.



an example of EM

paras  $\left\{ \begin{array}{l} \{q(z)\}_{i=1}^k \sim \text{initial} \\ \{P_{ij}\}_{i,j=1}^k \sim \text{transition} \\ \{P(y|z)\}_{z=1}^k, y \in \mathcal{Y} \sim \text{emission} \end{array} \right.$

What is GOOD to have (to estimate paras)?

$\rightarrow$  It would be good to also have  $x_1, \dots, x_n$  in addition to our  
Hidden Variables

Real Observation.

$(y_1, \dots, y_n)$

如果 Incomplete  $\rightarrow \{y_1, \dots, y_n\}$  likelihood

$$L(x, \theta) = \prod_{i=1}^k P(x_i)$$

$$\text{Complete Dataset} = \{ \underbrace{x_1, \dots, x_n}_{\text{Hidden Variable}}, y_1, \dots, y_n \} = \prod_{i=1}^K \sum_j P(s) P(x_i | s)$$

$$\left\{ \begin{array}{l} \delta(i|t) = \begin{cases} 1 & \text{if } x_t = i \in S \\ 0 & \text{else} \end{cases} \\ \delta(i,j|t) = \begin{cases} 1 & \text{if } x_t = i, x_{t+1} = j \\ 0 & \text{else} \end{cases} \end{array} \right. \rightsquigarrow \left\{ \begin{array}{l} \delta(i|t) : i \in S, t=1, \dots, n \\ \delta(i,j|t), i \in S, t=1, \dots, n \end{array} \right.$$

$$\log P = \log q(x) \cdot \prod_{t=1}^m P_{x_t|x_{t-1}} \prod_{t=1}^n P(y_t|x_t)$$

Given the complete dataset = { $x_1, \dots, x_n, y_1, \dots, y_n$ }

We can write down the complete log-likelihood:

$$\ell(\{x_t\}, \{y_t\}; \theta) = \sum_{i=1}^K \delta(i|1) \log(q(i)) + \sum_{i=1}^K \left( \sum_{t=1}^n \delta(i|t) \log P(y_t | i) \right) + \sum_{i,j=1}^K \left( \sum_{t=1}^{n-1} \delta(i,j|t) \right) \log P_{ij}$$

initial                          emission  
transition

X is hidden

$$M\text{-step: } \mathbb{E}_{p(x|y, \theta^{(e)})} [\ell(\{x_t\}, \{y_t\}; \theta)] \rightarrow \underset{\theta^{(t+1)}}{\text{maximize}}$$

$\downarrow$

the quantity we care

$$E\text{-step: } p^{(e)}(i|t) = P(x(t)=i | y_1, \dots, y_n) = \mathbb{E}[\delta(i|t) | \{y_t\}, \underline{\theta}^{(e)}]$$

$$p^{(e)}(i,j|t) = P(x(t)=i, x(t+1)=j | y_1, \dots, y_n) = \mathbb{E}[\delta(i,j|t) | \{y_t\}, \underline{\theta}^{(e)}]$$

Q: How to compute  $p(i|t, \{y_k\}, \underline{\theta}^{(e)})$  }  $\Rightarrow \text{Posterior}$   
 $p(j|i, t, \{y_k\}, \underline{\theta}^{(e)})$  }

Ans: ①  $P^{(1)}(X(t)=i | y_1, \dots, y_n, \underline{\theta}^{(e)})$

$$= \frac{P(X(t)=i, y_1, \dots, y_n | \underline{\theta}^{(e)})}{P(y_1, \dots, y_n | \underline{\theta}^{(e)})}$$

$$= \frac{P(X(t)=i, y_1, \dots, y_n | \underline{\theta}^{(e)})}{\sum_{i'} P(X(t)=i', y_1, \dots, y_n | \underline{\theta}^{(e)})} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i'} \alpha_t(i') \beta_t(i')}$$

Numerator  $P(X(t)=i, y_1, \dots, y_n; \underline{\theta}^{(e)})$

$$= \alpha_t(i) \beta_t(i) \leadsto \text{according to [pars] } \underline{\theta}^{(e)}$$

②

$$P^{(2)}(i, j | t) = P(X(t)=i, X(t+1)=j | y_1, \dots, y_n; \underline{\theta}^{(e)})$$

$$= \frac{P(y_1, \dots, y_n, X(t)=i, X(t+1)=j; \underline{\theta}^{(e)})}{\sum_{i'j'} P(y_1, \dots, y_n, X(t)=i', X(t+1)=j'; \underline{\theta}^{(e)})}$$

Numerator:

$$P(y_1, \dots, y_n, X(t)=i, X(t+1)=j) \quad (\text{omit } \underline{\theta}^{(e)})$$

$$= P(y_{t+2}, \dots, y_n | X(t+1)=j) \cdot P(X(t+1)=j, X(t)=i, y_1, \dots, y_{t+1})$$

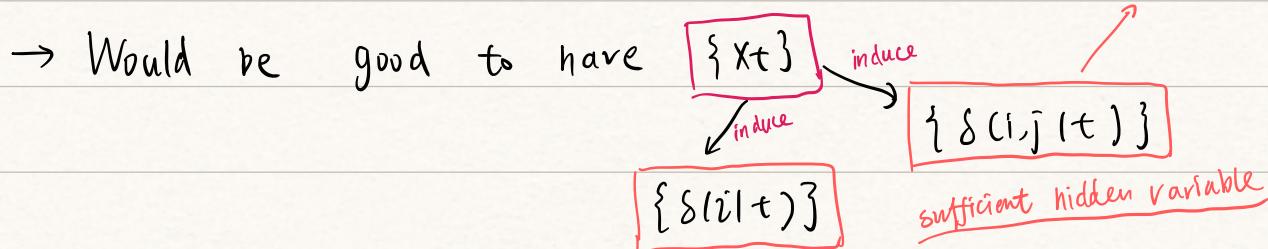
$$= P(y_{t+2}, \dots, y_n | X(t+1) = j) \cdot P(y_{t+1} | j) \underbrace{P(i, j)}_{P(y_{t+1}, X(t+1) = j | X(t) = i)} \cdot P(y_1, \dots, y_t, X(t) = i)$$

$$= \alpha_t(i) P_{ij} P(y_{t+1} | j) \beta_{t+1}(j)$$

Thus,  $P(i, j | t) = \frac{\alpha_t(i) P_{ij} P(y_{t+1} | j) \beta_{t+1}(j)}{\sum_{i', j'} \alpha_t(i') P_{i'j'} P(y_{t+1} | j') \beta_{t+1}(j')}$

Baum-Welch Alg.

Given  $\{y_t\}_{t=1}^n$ , estimate  $\{q(i)\}$   $\{p_{ij}\}$   $\{P(y|z)\}$   
 origin hidden variable  $\mathbb{E}[\delta(i,j|t) | \theta^{(e)}, y]$



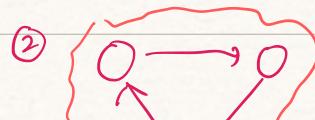
→ But  $\underline{p(i|t)}$  &  $\underline{p(i,j|t)}$  can be estimated recursively using  
 $\alpha_t$  &  $\beta_t$

Bayesian Networks

i) A directed acyclic graph (DAG)

ii) An associated prob. distribution that respects the graph

Example





DAG provides a qualitative description of the (conditional) independence relations among

Prob. dist<sup>?</sup> provides a quantitative description.

DAG can be used to ① explain properties of dist<sup>?</sup>

② perform efficient computations

$$\text{HMM : } \stackrel{\textcircled{1}}{P(y_1, \dots, y_n)} \rightarrow \text{F-B alg.} \rightarrow O(k^2 n)$$

$$\stackrel{\textcircled{2}}{(x_1^*, \dots, x_n^*)} = \underset{x_1, \dots, x_n}{\operatorname{arg\,max}} P(x_1, \dots, x_n, y_1, \dots, y_n) \rightarrow \text{Viterbi} \rightarrow O(n)$$

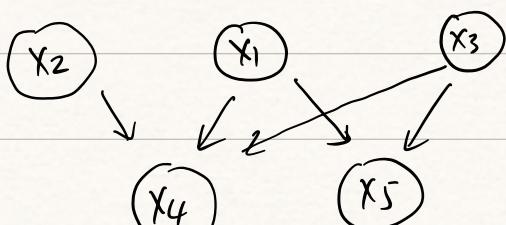


this can be generalized to any DAG



Exploit the DAG structure to do the above ① & ② efficiently!

[E.g.] Form a joint dist<sup>?</sup>.



x<sub>1</sub> & x<sub>2</sub> are parents of x<sub>4</sub>

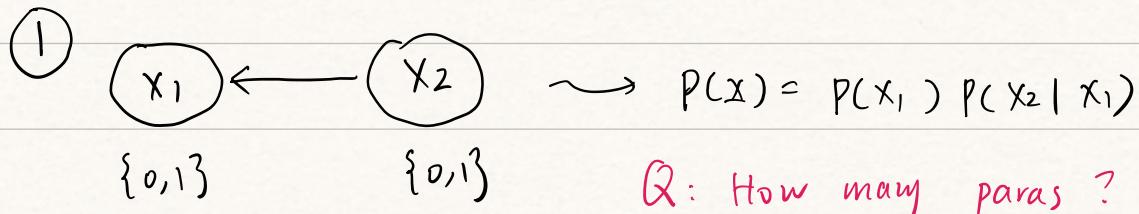
$$P(X_1, \dots, X_5) = P(X_1) P(X_2) P(X_3) P(X_4 | X_1, X_2, X_3) P(X_5 | X_1, X_3)$$

More generally,  $\underline{X} = (X_1, \dots, X_d)$

$$P(\underline{X}) = \prod_{k=1}^d P(X_k | \text{pa}(k)) \quad \textcircled{1} \text{ pa}(k) \text{ is the set of parents of } k$$

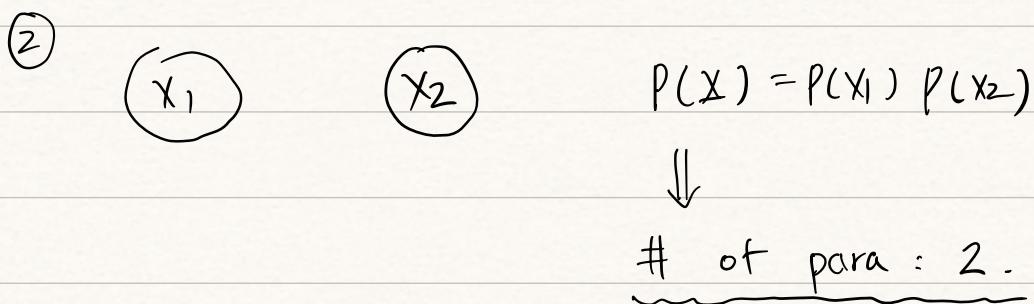
$$\textcircled{2} \quad X_u = (X_i : i \in u)$$

Consider the case where all variables are binary, i.e.,  $\{0, 1\}$

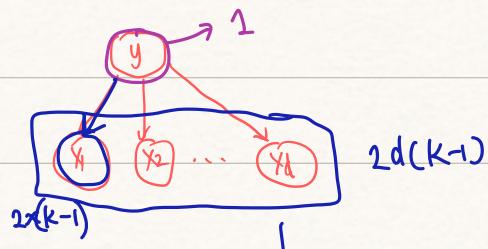


Ans: 3 params.

$$\begin{cases} P(X_1 = 0) \\ P(X_2 = 0 | X_1 = 0) \\ P(X_2 = 0 | X_1 = 1) \end{cases}$$



③  $X_i \in \{1, \dots, K\} \quad y \in \{0, 1\}$   
 $i = 1, 2, \dots, d$



# Naive Bayes Model for Binary classification

$y \in \{0, 1\}$

Q: How many params?  $\Rightarrow$  # of params:  $2d(k-1) + 1$

$$\text{Ans: } \underbrace{P(X_i | y=0)}_{k-1} \quad \underbrace{P(X_i | y=1)}_{k-1}$$

1 feature  $\leadsto 2(k-1)$  params

d features  $\leadsto 2d(k-1)$  params

$$P(X, y) = \underbrace{P(X|y)}_{2d(k-1)} \underbrace{P(y)}_1$$

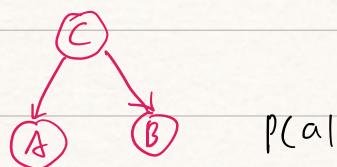
Conditional Independence & Explain away.

Def: A & B are conditional independent given C if

$$P(A=a, B=b | C=c) = P(A=a | C=c) P(B=b | C=c)$$

for all a, b, c

E.g:



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

Qn 1: A & B independent ? → NO

$$\text{Pf: } P(a, b) = \sum_c P(c) P(a|c) P(b|c) \neq P(a) P(b)$$

Qn 2: A & B cond. independent given C ? → YES!

$$\text{Pf: } P(a, b, c) = P(c) P(a|c) P(b|c)$$

$$\Leftrightarrow P(a, b|c) = P(a|c) P(b|c)$$

$$\Rightarrow \underbrace{A \perp\!\!\!\perp B \mid C}$$

E.g.



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

Qn 1: A & B independent ? → NO

$$\text{Pf: } P(a, b) = \sum_c P(a) P(c|a) P(b|c) = P(a) \sum_c P(c|a) P(b|c) \neq P(a) P(b)$$

Qn 2: A & B cond. independent given C ? → YES!

$$\text{Pf: } P(a, b|c) = \frac{P(a, b, c)}{P(c)}$$

$$= \frac{P(a) P(c|a) P(b|c)}{P(c)}$$

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

$$\Rightarrow A \perp\!\!\!\perp B \mid C$$

E.g. :



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

Qn 1: A & B independent?  $\Rightarrow \underline{\text{YES!}}$

Pf:  $P(a, b) = \sum_c P(a, b, c)$

$$= P(a) \cdot P(b) \cdot \sum_c P(c|a, b)$$

$$= P(a) P(b)$$

Qn 2: A & B cond. indep. given C?  $\Rightarrow \text{NO!}$

Pf:  $P(a, b | c) = \frac{P(a, b, c)}{P(c)}$

$$= \frac{P(a) P(b) P(c|a, b)}{P(c)}$$

$$\neq P(a|c) P(b|c) \quad \underline{\text{in general}}$$

Intuition:  $(A) \rightarrow (C) \leftarrow (B)$



C is a head-to-head node

Conditioned on C, path between A & B become unblocked  $\rightarrow$  A & B are DEPENDENT

Explain Away

[E.g.]

Rain  
Last Night

Sprinkle  
On

$R=1 \Rightarrow$  Rain Last Night  
 $S=1 \Rightarrow$  Sprinkle On  
 $G=1 \Rightarrow$  Grass Wet



$$P(R=1) = P(S=1) = 0.9$$

$S$	$R$	0	1
0		0.1	0.2
1		0.2	0.8

①  $P(S=0) = 0.1$

② Suppose we observe  $G=0$

$$\underline{P(G=1 | R=r, S=s)}$$

$$P(S=0 | G=0) = 0.25 \dots$$

↓ 说明  $S$  和  $G$  不独立  $S \rightarrow G$

③ Suppose we observe  $G=0$  &  $R=0$

$$P(S=0 | G=0, R=0) = 0.111 \dots \stackrel{\textcircled{1}}{>} 0.1$$

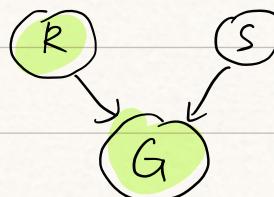
说明 Given  $G$

$$\stackrel{\textcircled{2}}{<} 0.25$$

$S$  和  $R$  不独立  
 $S \rightarrow G$

Explain Away the fact that  $S=0$

make the prob closer to prior 0.1



If  $R \perp\!\!\!\perp S | G$  (which is not true), then

$$\underline{P(S=0 | G=0)} = P(S=0 | G=0, R=0)$$



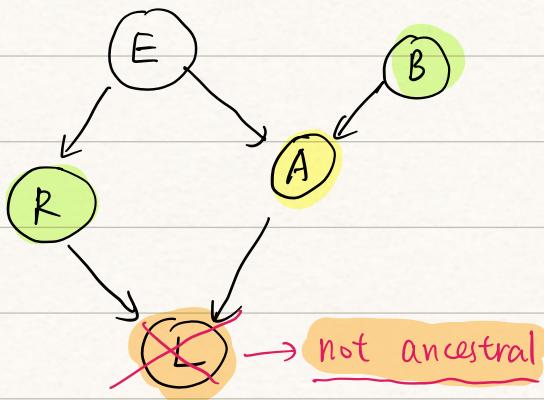
But this equality does not hold.

Instead of writing down the joint dist & marginalizing to answer

Cond. indep. questions, we want a systematic way to answer

C. I. question !

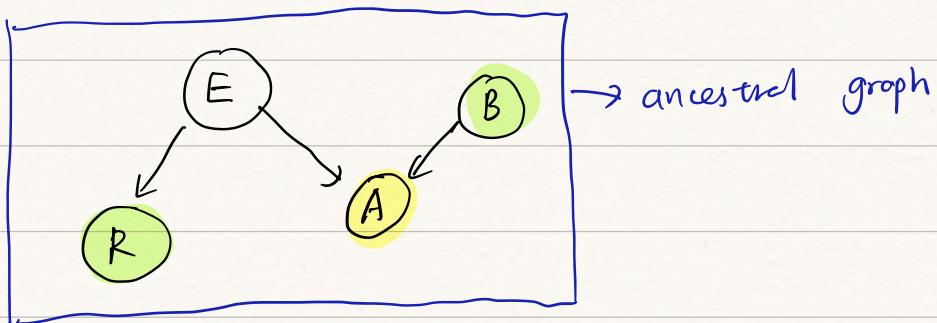
## Moralization



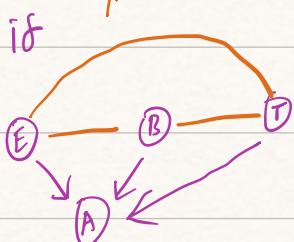
Qn : Are R & B C.I given A ?

Ans: Four steps :

Step 1 : Construct ancestral graph of nodes / variable of interest (R, A, B)

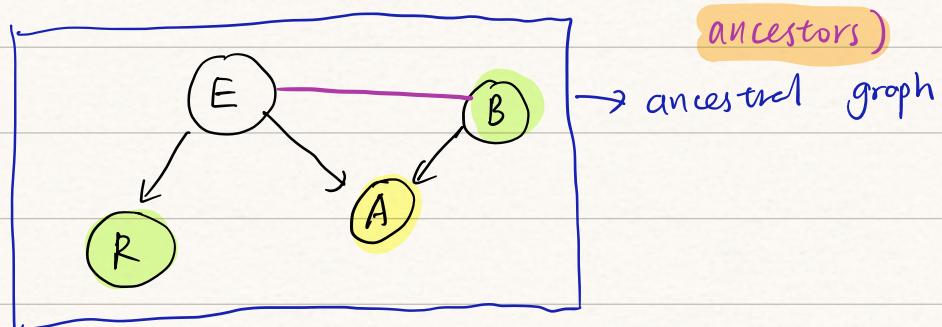


Pair-wise many !



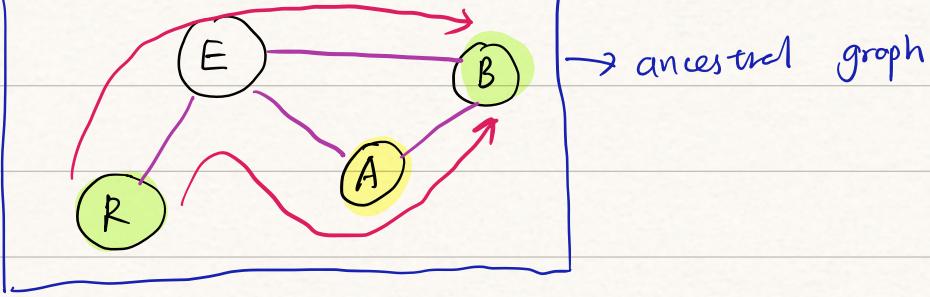
Step 2 : Moralize the ancestral graph

by marrying the parents  $\Rightarrow$  only the parents (not all ancestors)



(Takes care of head - to - head explaining away phenomenon)

Step 3 : Change arcs to edges



Step 4: Answer !



Is A a **seperator** of R & B in the  
undirected graph

Defn: S is a seperator of A and B

NO!

↔ every path from A to B passes



through some nodes in S

Conclusion:  $R \not\perp\!\!\!\perp B \mid A$



$R \perp\!\!\!\perp B \mid (A, E)$  is TRUE !