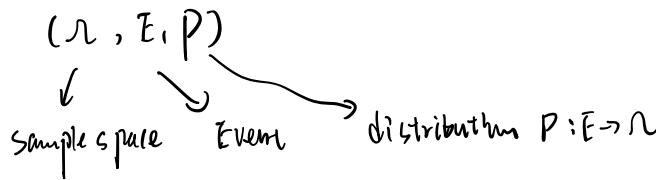


Probability Space:



Random Variables: X (upper case)

- $X: N \rightarrow S$.
- Set of values X can take: $\text{Val}(X)$
- $P(x) \Leftrightarrow P(X=x)$
- x^i : a special value of X

Prob Review:

1. Joint Prob: 2 or more events: X, Y, Z, \dots

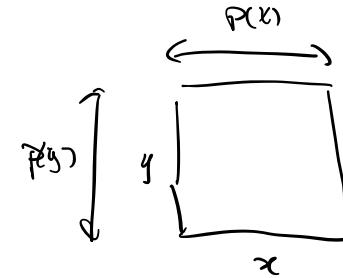
$$P(X=x, Y=y)$$

2. Marginal:

$$P(x) = \int p(x, y) dy, \quad p(y) = \int p(x, y) dx$$

3. Conditional:

• $p(x|Y=y) := \text{prob of } X \text{ given } Y=y.$



• Calculation:

$$P(x|Y=y) = \frac{p(x, Y=y)}{p(Y=y)}$$

$$\Rightarrow p(x|y) = \frac{p(x,y)}{p(y)}$$

$$\Rightarrow p(w, v, y, z) = p(w, x, y | z) \cdot p(z)$$

$$= p(w, x | y, z) \cdot p(y | z) \cdot p(z)$$

$$= p(w | x, y, z) \cdot p(x | y, z) \cdot p(y | z) \cdot p(z)$$

chain rule

likelyhood

$$p(y|x) = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

posterior

Evidence (constant, normalizer)

(y is latent variable,
x is observation)

4. independence: $p(x|y) = p(x)$, $p(y|x) = p(y)$

$$\Rightarrow p(x,y) = p(x)p(y)$$

5. Expectation: $E[f(x)] = \sum_x f(x)p(x)$

$$E[f(x)] = \int_x f(x)p(x)dx$$

Rules: ① $E(c) = c$

③ $E(f(x) + g(x)) = E(f(x)) + E(g(x))$

② $E(c \cdot f(x)) = c \cdot E(f(x))$

④ $E(f(x) \cdot g(y)) = E(f(x)) \cdot E(g(y)) \Leftrightarrow x, y \text{ independent}$



$$\int f(x) dx \cdot \int g(y) dy$$

6. common distributions:

① Bernoulli: $p(x|\lambda) = \begin{cases} 1-\lambda, & x=0 \\ \lambda, & x=1 \end{cases}$

denote as: $p(x) = \text{Bern}_x[\lambda]$

② Gaussian:

I. single R.V., 2 para:

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

denote as: $\text{Norm}_x[\mu, \sigma^2]$

II. D-dimensional R.V. X , para: $\mu \in \mathbb{R}^D$, $\Sigma \in \mathbb{R}^{D \times D}$

symmetrical, positive
covariance matrix

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$$

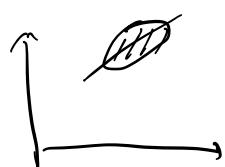
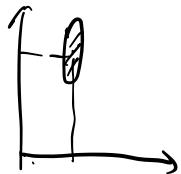
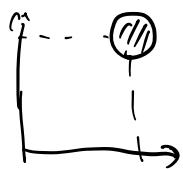
denote as: $\text{Norm}_x[\mu, \Sigma]$

Type of covariance: spherical, diagonal, full

$$\Sigma_{\text{sphere}} = \begin{bmatrix} b^2 & 0 \\ 0 & b^2 \end{bmatrix}$$

$$\Sigma_{\text{diag}} = \begin{bmatrix} \zeta_1^2 & 0 \\ 0 & \zeta_2^2 \end{bmatrix}$$

$$\Sigma_{\text{full}} = \begin{bmatrix} \zeta_{11}^2 & \zeta_{12}^2 \\ \zeta_{21}^2 & \zeta_{22}^2 \end{bmatrix}$$



Summary: (FDU AI - probabilistic Reasoning and decision making).

① Probability:

- { Marginalization, independence (marginal, conditional)
- product
- Bayesian

② Bayes Net : (represent relationship between variables)

• inference :

- calculate prob with full joint

- variable elimination. \leftarrow multiply sum polytree
conditioning

\hookrightarrow order of leaves
{ non-query
non-evidence . }

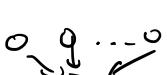
- independence:

- { d-separation.
- markov blanket

1. $x \rightarrow E \rightarrow Y$
2. $x \leftarrow E \rightarrow Y$
3. $x \rightarrow O \leftarrow Y$

(all Os are non-evidence)

- Noisy-OR: (n para vs 2^n para)



$$P(Y|X_1, \dots, X_n) = 1 - \prod_{i=1}^n (1 - P_i)^{x_i}$$

• learning :

- MLE for complete data

- EM for incomplete data (EM for Noisy-OR)

• special BN:

- naive Bayes

- HMM (forward Algorithm)

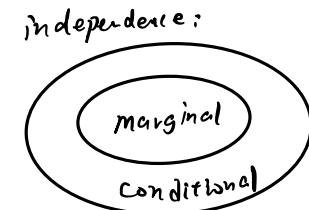
↳:

• Conditional Independence:

$$A, B \text{ conditional independent} \Leftrightarrow P(A|BC) = P(A|C) . P(B|AC) = P(B|C)$$

then we can get : $P(ABC) = P(A|C) \cdot P(B|C)$

$$P(ABC) = P(C) \cdot P(A|C) \cdot P(B|C)$$



• marginally Independent:

$$A, B \text{ marginally independent} \Leftrightarrow P(A|B) = P(A) . P(B|A) = P(B)$$

• CPT: conditional prob table

★ Bayes Net:

BN(def): A DAG with:

- Node — Variables
- edge — conditional dependence
- CPT at each node — how a node depends on its parents

E: Earthquake happens

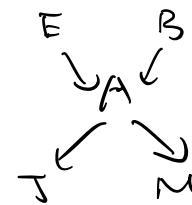
B: Burglar breaks in

A: alarm goes off

J: Jamal calls

M: Maya calls

BN: (dependency)



⇒ no CPT:

$$P(E), P(B)$$

$$P(A|EB)$$

$$P(J|A), P(M|A)$$

↳ Binary Events

1. compare: $P(B=1)$, $P(B=1|A=1)$, $P(B=1|A=1, E=1)$

$$P(B=1|A=1) > P(B=1|A=1, E=1) > P(B=1)$$

Intuitive Reason: $P(B=1) < P(B=1|A=1, E=1) < P(B=1|A=1)$

$$\cancel{P(B=1)} < P(B=1|A=1)$$

i.e. "I generally believe there is a very low chance that a burglar breaks in. But if I get notified that my alarm is going off, then I think the alarm is probably caused by the burglar and my belief that a burglar broke in goes up. However, if I then hear about an earthquake, I no longer am as fearful that a burglar broke in, but I'm still more worried than I would have been if the alarm hadn't gone off at all!"

$$P(B=1|A=1) > P(B=1|A=1, E=1) \quad P(B=1|A=1, E=1) > P(B=1)$$

2. calculate: $P(B=1|A=1)$, $P(B=1|A=1, E=1)$

Ans: E, B: marginal independent . i.e. $P(E|B) = P(E)$. $P(B|E) = P(B)$

$$\textcircled{1} \quad P(B=1 | A=1) = \frac{P(A=1 | B=1) \cdot P(B=1)}{P(A=1)} = \frac{\sum_{i,j} P(A=1, E=e_i, B=b_j)}{\sum_{i,j} P(A=1, E=e_i, B=b_j)} \cdot P(B=1)$$

$$= \frac{\left(P(A=1 | E=1, B=1) \cdot P(E=1 | B=1) + P(A=1 | E=0, B=1) \cdot P(E=0 | B=1) \right) \cdot P(B=1)}{\sum_{i,j} P(A=1, E=e_i, B=b_j)}$$

✓

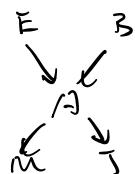
$$= \frac{\left(P(A=1 | E=1, B=1) \cdot P(E=1) + P(A=1 | E=0, B=1) \cdot P(E=0) \right) \cdot P(B=1)}{\sum_{i,j} P(A=1, E=e_i, B=b_j)}$$

✓

$$\textcircled{2} \quad P(B=1 | A=1, E=1) = \frac{P(B=1, A=1, E=1)}{P(A=1, E=1)} = \frac{P(A=1 | E=1, B=1) \cdot P(E=1) \cdot P(B=1)}{\sum_i P(A=1, B=b_i, E=1)}$$

3. compare: $P(A=1), P(A=1 | S=1), P(A=1 | S=1, M=1)$

$$P(A=1 | S=1, M=1) > P(A=1 | S=1) > P(A=1)$$



4. calculate: $P(A=1 | S=1), P(A=1 | S=1, M=1)$

条件: S, M are conditional independent from each other and from E, B given A

$$\text{i.e.: } P(S | A, B, E, M) = P(S | A)$$

$$\textcircled{1} \quad P(A=1 | S=1) = \frac{P(S=1 | A=1) \cdot P(A=1)}{P(S=1)} = \frac{\sum_i P(S=1 | A=1) \cdot P(A=1)}{\sum_i P(S=1 | A=1) \cdot P(A=1) + P(S=1 | A=0) \cdot P(A=0)}$$

$$\textcircled{2} \quad P(A=1 | S=1, M=1) = \frac{P(A=1, S=1, M=1)}{P(S=1, M=1)} = \frac{\sum_{i,j} P(A=1, S=1, M=1, B=b_i, E=e_j)}{\sum_i P(M=1 | S=1) \cdot P(S=1)}$$

$$= \frac{\sqrt{P(M=1, A=1 | S=1) + P(M=1, A=0 | S=1)}}{\sqrt{(P(M=1, A=1 | S=1) + P(M=1, A=0 | S=1)) \cdot (P(M=1, A=1 | S=1) + P(M=1, A=0 | S=1))}}$$

$$= \frac{\sqrt{(P(M=1 | A=1) \cdot P(A=1 | S=1) + P(M=1 | A=0) \cdot P(A=0 | S=1)) \cdot (P(M=1 | A=1) \cdot P(A=1 | S=1) + P(M=1 | A=0) \cdot P(A=0 | S=1))}}{\sqrt{(P(M=1 | A=1) \cdot P(A=1 | S=1) + P(M=1 | A=0) \cdot P(A=0 | S=1)) \cdot (P(M=1 | A=1) \cdot P(A=1 | S=1) + P(M=1 | A=0) \cdot P(A=0 | S=1))}}$$

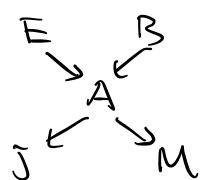
$$\textcircled{5} \quad P(A=1 | S=1, B=1) = \frac{P(A=1, S=1, B=1)}{P(S=1, B=1)} = \frac{P(S=1 | A=1) \cdot P(A=1 | B=1) \cdot P(B=1)}{P(S=1 | A=1) \cdot P(A=1 | B=1) \cdot P(B=1) + P(S=1 | A=0) \cdot P(A=0 | B=1) \cdot P(B=1)}$$

$$= \frac{P(S=1 | A=1) \cdot P(A=1 | B=1) \cdot P(B=1)}{P(S=1 | A=1) \cdot P(A=1 | B=1) \cdot P(B=1) + P(S=1 | A=0) \cdot P(A=0 | B=1) \cdot P(B=1)}$$

6. Full Joint Distribution: $P(A, M, J, B, E)$

↳ which is the most easiest prob to compute with BN.

Eg:



$$\begin{aligned} P(A, M, J, B, E) \\ = P(M|A) \cdot P(J|A) \cdot P(A|B, E) \cdot P(B) \cdot P(E) \end{aligned}$$

Complexity: (disadvantage)

for $\{x_1, \dots, x_n\}$ n variables.

- ① # of values need to specify Full Joint: 2^n
- ② # of values needed to specify the BN: $O(n \cdot 2^k)$
(k: max # of parents any node has)

7. BN, BN inference.

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n | x_{n-1}, x_{n-2}, \dots, x_1) \cdot P(x_{n-1} | x_{n-2}, \dots, x_1) \cdots P(x_2 | x_1) \cdot P(x_1) \\ &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) \\ &= \prod_{i=1}^n P(x_i | \text{parent}(x_i)), \quad \text{parent}(x_i) \subseteq \{x_1, \dots, x_{i-1}\} \end{aligned}$$



8. How to build a BN?

1. choose R.V. $\{x_1, \dots, x_n\}$

2. choose an ordering

3. while there are variables left:

I. add node x_i to BN

II. set $\text{parents}(x_i)$: parents

non-parents: conditional independent to x_i

III. Define $P(x_i | \text{parent}(x_i))$ CPT.



9. Better Ordering: Root events first (reason: reduce complexity of BN.)

10. Model: x_1, x_2, \dots, x_k



i. Representing CPTs:

advantage: easy to read

disadvantage: size goes

x_1, x_2, \dots, x_k	$P(Y=1)$
0 0 .. 0	$P(Y=1 x_1=0, \dots, x_k=0)$
0 0 .. 1	$P(Y=1 x_1=1, \dots, x_k=0)$
⋮ ⋮ ⋮	⋮

$$\text{exponentially } 1 \cdot 1 \cdots 1 \quad | \quad p(Y=1 | X_1=1, \dots, X_k=1)$$

II. Representing OR logic: make $p(Y=1)$ deterministic of its parents.

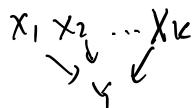
$$P(Y=1 | X_1 \dots X_k) = \text{OR}(X_1 \dots X_k)$$

$$\hookrightarrow \text{OR}(X_1 \dots X_k) = \begin{cases} 1, & \text{if any } X_i = 1 \\ 0, & \text{otherwise.} \end{cases}$$

advantage: Save space (just calculation)

disadvantage: deterministic

Noisy-OR:



$$P(Y=1 | X_1 \dots X_k) = 1 - \prod_{i=1}^k (1-p_i)^{X_i} \quad \text{—— a k-param model.}$$

$$\hookrightarrow p_i \stackrel{\text{def}}{=} p(X_i=1, \text{others}=0)$$

Intuition:

- all parents off \Rightarrow child is off certainly
- more parents on \Rightarrow more likely child is on
- all parents on doesn't ensure that child is on (but very likely it's)

III. 3 equations for X, Y conditional independence given E :

$$\left\{ \begin{array}{l} P(X|Y, E) = P(X|E) \\ P(Y|X, E) = P(Y|E) \\ P(X, Y|E) = P(X|E) \cdot P(Y|E) \end{array} \right.$$

IV. BN and conditional independence:

X, Y, E — sets and disjoint. When X, Y independent given E ?

Thm: $P(X, Y|E) = P(X|E) \cdot P(Y|E) \Leftrightarrow$ every undirected path from a node in X to a node in Y is d-separated (blocked) by E

D-Separated (def): A path π is d-separated if $\exists z \in \pi$ holds 1 of these 3 conditions:

① $\circ \rightarrow \dots \circ \rightarrow \circ \rightarrow \dots \circ$, $z \in E$

② $\circ \leftarrow \dots \circ \leftarrow \circ \rightarrow \dots \circ$, $z \in E$

③ $\circ \rightarrow \dots \circ \rightarrow \circ \leftarrow \circ \rightarrow \dots \circ$, $z \notin E$, descendants(z) $\notin E$

Checking Steps:

I. list X, Y, E

II. list all path $X - Y$, regardless directions

III. put back the directions, check if all paths are d-sep

$$\text{ex. } P(F \wedge H \mid A, \bar{E}) = P(F \mid A) P(H \mid A, \bar{E}) ?$$

$$X = \{F\}, Y = \{H, \bar{E}\}, E = \{\bar{A}\}$$

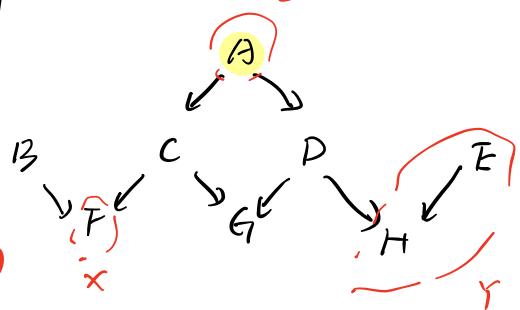
$$\text{With d-sep, we can know } P(X \mid E) \cdot P(Y \mid \bar{E}) = P(X, Y \mid \bar{E})$$

$$\therefore P(F, H \mid \bar{E} \mid A) = P(F \mid A) \cdot P(H \mid \bar{E} \mid A)$$

$$\stackrel{H}{P(F, H \mid E \mid A)} \cdot P(E \mid A)$$

$$\stackrel{H}{P(H \mid A \mid \bar{E})} \cdot P(\bar{E} \mid A)$$

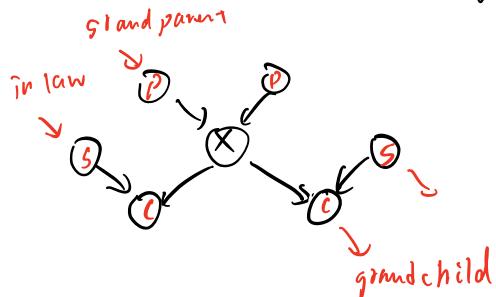
$$\Rightarrow P(F, H \mid E, A) = P(F \mid A) \cdot P(H \mid A, \bar{E}).$$



13. Markov Blanket: (Bx)

$$B_x := \{ \text{nodes consist of } \text{parent}(x), \text{children}(x), \text{spouses}(x) \}$$

$$\text{Thm: } P(x \mid B_x, \bar{x}) = P(x \mid B_x), \text{ if } \bar{x} \notin \text{parent}(B_x)$$



14. Enumeration:

add & # of multiplications of $P(B=1 \mid I, M=1)$

$$P(B, I, M) = \sum_{i,j} P(B=1, I=1, M=1, E=e_i, A=a_j) = \sum_{i,j} P(B=1) \cdot P(E=e) \cdot P(A=1 \mid B=1, E=e) \cdot P(I=1 \mid A=a)$$

$$P(B \mid I, M) = \frac{P(B, I, M)}{P(B, I, M) + P(\bar{B}, I, M)}$$

→ redundant

\downarrow \downarrow
16* 3+ 16* 3+

6 multiplications, 3 additions
· $P(M=1 \mid A=a)$

$\therefore \text{total: } 32* . 7+$.

query \uparrow Evidence

Day 6 Variable Elimination:

$$P(Q | E)$$

- Factor: a table that describes a function maps some values to vars

$v_1 v_2 \dots$	Values	(not necessarily just for prob)

- Operations on Factors:

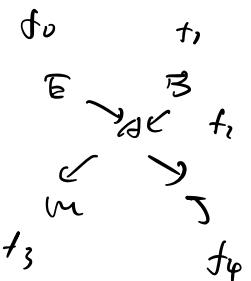
} Multiplication
 } Addition
 } Conditioning

$$\text{eg: } P(B | J=1, M=1) = \frac{P(B, J=1, M=1)}{P(J=1, M=1)}$$

$$= \frac{\sum_e \sum_a P(B, J=1, M=1, E=e, A=a)}{\dots + \dots}$$

$$P(B, J=1, M=1, E=e, A=a) = P(E=e) \cdot P(B) \cdot P(A=a | B, E=e) \\ P(J=1 | A=a) \cdot P(M=1 | A=a)$$

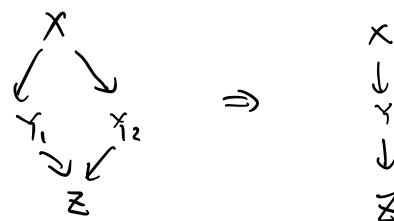
$$= f_0(e) \cdot f_1(B) \cdot f_2(a | B, e) f_3(J=1 | A) \\ f_4(M=1 | A)$$



- Order of Elimination: (non-query, non-evidence) leaves \rightarrow root, and head \rightarrow bottom.

- polytree: A graph with no undirected loops.

collapsing nodes:



- How to eliminate variables?

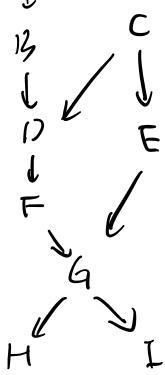
① create factors for CPTs

② eliminating order (leaves \rightarrow parents)

③ operate
 { multiple
 summing out $\rightarrow \sum_{e,a} P(B, J=1, M=1, A=a, E=e)$ (marginal)
 conditioning }

$$\rightarrow P(B, J=1)$$

Eg:



Calculate $P(G|H=1)$

I. # of ops we need to do if we just use joint prob

$$P(G|H=1) \leftarrow \frac{P(G=0|H=1)}{P(G=1|H=1)} \leftarrow P(G=0, H=1) \cdot P(G=1, H=1)$$

$$P(G=0, H=1) = \sum_{a,b,c,d,e,f,i} P(A=a, B=b, C=c, D=d, E=e, F=f, G=0, H=1, I=i)$$

$$\stackrel{\text{CPTs}}{=} \sum_{a,b,c,d,e,f,i} P(A=a) \cdot P(B=b|A=a) \cdot P(C=c) \cdot P(D=d|B=b, C=c) \cdot P(E=e|C=c) \cdot P(F=f|D=d) \\ \cdot P(G=0|F=f, E=e) \cdot P(H=1|G=0) \cdot P(I=i|G=0)$$

each variable is

binary

we have all CPTs

$\therefore 2^7$ terms, each term has $8\times$

$$\# \text{ of multiplication} = 2 \times (2^7 \cdot 2^3) = 2^9$$

$$\# \text{ of addition} = 2 \times (2^7 - 1) = 2^8 - 2$$

II. Calculate $P(G|H=1)$

$$P(G, H=1) = \sum_{a,b,c,d,e,f,i} P(G, H=1, A=a, B=b, C=c, D=d, E=e, F=f, I=i)$$

$$\stackrel{\text{CPTs}}{=} \sum_{a,b,c,d,e,f,i} P(A=a) \cdot P(B=b|A=a) \cdot P(C=c) \cdot P(D=d|B=b, C=c) \cdot P(E=e|C=c) \cdot P(F=f|D=d) \\ \cdot P(G|F=f, E=e) \cdot P(H=1|G) \cdot P(I=i|G)$$

(Elimination Order: I, H, A, B, C, D, E, F ...)

④ I, H, A, B, ...

$$\sum_{B=b}^{4x} \left(P(D=d|B=b, C=c) \sum_{A=a}^{2x} P(A=a) \cdot P(B=b|A=a) \right) \cdot P(H=1|G) \sum_{I=i}^{2x} P(I=i|G)$$

$$\underbrace{F(B, A)}_{F(B)}$$

condition

$$\sum_{I=i}^{2x} P(I=i|G)$$

$$\hookrightarrow F(G)$$

$$F(C, D)$$

C, D	F(C, D)
0 0	...
1 1	...

② C, D, E, F :

$$\begin{aligned}
 & \sum_{e,f} P(G|F=f, E=e) \cdot \sum_{D=d} P(F=f|I=d) \sum_{C=c} P(C=c) \cdot P(E=e|C=c) \cdot F(C,D) \cdot F(G) \\
 & = F(G) \cdot \sum_{e,f} F(G, F, E) \cdot \sum_D F(F, D) \cdot \sum_C F(C) \cdot F(E, C) \cdot F(C, D) \\
 & \quad \downarrow \qquad \qquad \qquad \uparrow \qquad \qquad \qquad \downarrow \\
 & \quad F(C, D) \qquad \qquad \qquad F(C, E) \\
 & \quad \qquad \qquad \qquad \downarrow \\
 & F(G, F, E) \leftarrow F(F, E) \\
 & \qquad \qquad \qquad \downarrow \\
 & F'(G)
 \end{aligned}$$

$\therefore 40 \times 20^4$ (decrease from $2^{11} \times 2^{8+}$)

• Does Elimination order matter? Yes.

I. Time. Space of VE is dominated by the largest factor

II. Eliminate variables that lead to the smallest next factor.

III. In polytree: linear time inference

Day 7 . Learning Bayes Net. } learning parameters : CPTs
 } learning structure : NP-hard.

• Problem: Bayes Net includes DAG / CPTs.

but sometimes we don't know them fully. We may need to learn the DAG or CPTs

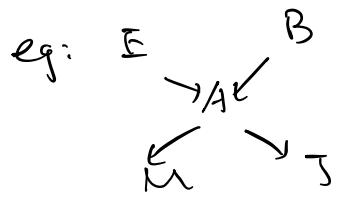
- In this lecture, we assume a given DAG and learn the CPTs from data :

- with complete data : MLE
- with incomplete data : EM

• ML (Maximum likelihood) learning:

• Idea: choose the model (CPTs) that maximizes the prob of data.

i.e. $P_{\text{model}}(\text{data})$



Sample	E	B	A	J	M
1	0	0	1	1	1
2	1	0	1	0	0
3	1	1	0	0	0
⋮	⋮	⋮	⋮	⋮	⋮

that is: $\underset{\text{CPTs}}{\operatorname{argmax}} P(\text{data} | \text{CPTs})$

DAG:



$X: \{\text{cherry}, \text{jalapeno}\}$

CPT: $P(X_i = \text{cherry}) = p$

data: $\{X^{(1)}, X^{(2)}, \dots, X^{(T)}\}$

— T samples.

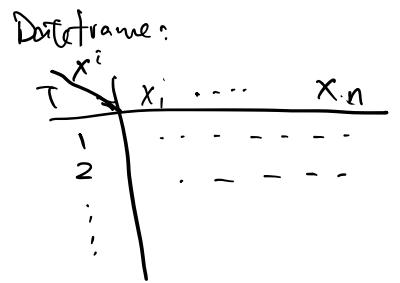
$$\begin{aligned}
 & \underset{P}{\operatorname{argmax}} P(\text{data} | p) \\
 &= \underset{P}{\operatorname{argmax}} \underbrace{P(X^1 \sim X^T | p)}_{\text{likelyhood}} \\
 &= \underset{P}{\operatorname{argmax}} \prod_{i=1}^T P(X^{(i)} | p) \\
 &= \underset{P}{\operatorname{argmax}} \log \prod_{i=1}^T P(X^{(i)} | p) \\
 &= \underset{P}{\operatorname{argmax}} \log p^{N_c} (1-p)^{N_j} \\
 &= \underset{P}{\operatorname{argmax}} (N_c \cdot \log p + N_j \cdot \log(1-p))
 \end{aligned}$$

∴ log-likelihood: $L(p) = N_c \log p + N_j \log(1-p)$

$\therefore \frac{\partial L(p)}{\partial p} = 0 \Leftrightarrow p = \frac{N_c}{N_c + N_j} = \frac{N_c}{T}$

Estimating para for B.N.: $p(X_i=x | P_a(X_i)=\pi)$

Data: $\{X_1^{(t)}, \dots, X_n^{(t)}\}_{t=1}^T$



$$L = \log p(\text{data}) \stackrel{i.i.d.}{=} \log \prod_{t=1}^T p(X_i=x_i^{(t)}, \dots, X_n=x_n^{(t)})$$

$$p(X_1^{(t)}, \dots, X_n^{(t)}) = \prod_{i=1}^n p(X_i=x_i^{(t)} | P_a(X_i)=P_a(x_i^{(t)}))$$

$$\Rightarrow L = \log \prod_{t=1}^T \prod_{i=1}^n p(X_i=x_i^{(t)} | P_a(X_i)=P_a(x_i^{(t)}))$$

$$= \sum_{t=1}^T \sum_{i=1}^n \log p(X_i=x_i^{(t)} | P_a(X_i)=P_a(x_i^{(t)}))$$

Sum all samples Sum all variables

$$(\text{for single variable} ; L = \sum_{t=1}^T \log(p(X=x^{(t)})))$$

• Def: Count ($X_i=x_i^{(t)}$, $P_a(X_i)=P_a(x_i^{(t)})$) $\stackrel{\text{def}}{=}$ # of samples where $X_i=x_i^{(t)}$, $P_a(X_i)=P_a(x_i^{(t)})$

$$\Rightarrow L = \sum_{i=1}^n \sum_{t=1}^T \log p(X_i=x^{(t)} | P_a(x^{(t)}))$$

$$= \sum_{i=1}^n \underbrace{\sum_x \sum_{\pi} \text{count}(X_i=x, P_{ai}=\pi) \log p(X_i=x | P_{ai}=\pi)}_{\text{constant}}$$

$$\text{we need: } \frac{\partial L}{\partial P(X_i=x | P_{ai}=\pi)} = 0$$

$$\Rightarrow P_{ML}(X_i=x | P_{ai}=\pi) = \frac{\text{Count}(X_i=x, P_{ai}=\pi)}{\sum_{x'} \text{Count}(X_i=x', P_{ai}=\pi)}$$

$$\text{or } P_{ML}(X_i=\pi) = \frac{\text{Count}(X_i=\pi)}{T}$$

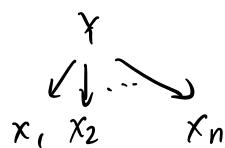
Def: $I(x, y) \stackrel{\text{def}}{=} \begin{cases} 1 & , x=y \\ 0 & , \text{otherwise.} \end{cases}$

$$\text{then we know: } P_{ML}(X_i=x) = \frac{1}{T} \text{Count}(X_i=x) = \frac{1}{T} \sum_{t=1}^T I(x_i^{(t)}, x)$$

Eg: $P(F=f | E=e, D=d)$

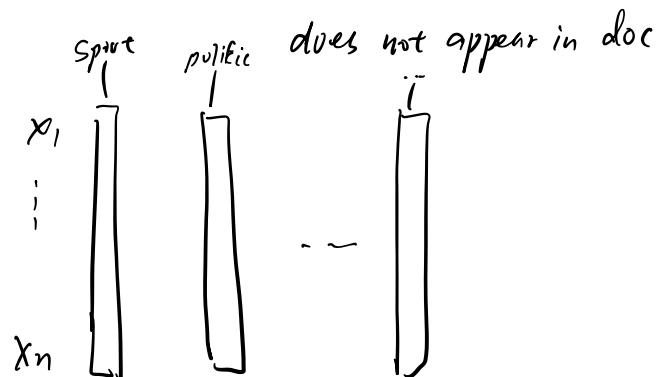
$$= \frac{\sum_t I(f_t=f) I(e_t=e) I(d_t=d)}{\sum_t I(c_t=c) I(d_t=d)}$$

* Naive Bayes: a Bayes Net with a single parent and many children.



e.g.: $Y \in \{1, 2, \dots, k\}$, 1=sports, 2=fashion, 3=...

$x_i \in \{0, 1\}$, for $i = 1, \dots, n$, where 0 means i th word in dictionary



learning:

- goal: observe a set of data $\{x_1, \dots, x_n\}$, try to inference Y .

from data: $P(x_i=x_i | Y=y) = \frac{\sum_t I(y_t, y) \cdot I(x_i^t, x_i)}{\sum_t I(y_t, y)}$

Classification:

$$P(Y=y | \{x_1=x_1, \dots, x_n=x_n\}) = \frac{P(Y=y) \prod_{i=1}^n P(x_i=x_i | Y=y)}{\sum_{y'} P(Y=y') \prod_{i=1}^n P(x_i=x_i | Y=y')}$$

* likelihood and KL distance:

log-likelihood: $L(p) = \log P_{\text{model}}(\text{data})$

$$KL(q, p) = \sum_{k=1}^n q_k \log \frac{q_k}{p_k}$$

Situation: An n -sided die is tossed T times, data: $\{x^{(1)}, \dots, x^{(T)}\}$.

In the course of T tosses: $C_k \stackrel{\text{def}}{=} \# \text{ of the } k^{\text{th}} \text{ side observed.}$

$$\text{so: } L(p) = \log P_{\text{model}}(\text{data}) = \log P(x^{(1)}, \dots, x^{(T)}) = \log \prod_{t=1}^T P(X=x^{(t)}) = \sum_{t=1}^T \log P(X=x^{(t)})$$

define $q_k = \frac{C_k}{T}$ as a distribution:

$$= \sum_{k=1}^n C_k \log P_k$$

$$KL(q, p) = \sum_{k=1}^n q_k \log q_k - \sum_{k=1}^n q_k \log p_k$$

$$= \sum_{k=1}^n q_k \log q_k - \frac{\sum_{k=1}^n C_k \log P_k}{T}$$

$$= \sum_{k=1}^n q_k \log q_k - \frac{L(p)}{T}$$

∴ conclusion:

$$\text{maximizing } \mathcal{L}(p) \iff \text{minimizing } KL(q, p)$$

Day 9: EM Algorithm (Expectation Algorithm):

· goal: estimate all the CPTs.

· from MLE (all data observed):

$$\begin{aligned} \text{result: } & \left\{ \begin{array}{l} P(X=x_i) = \frac{\text{count}(X=x_i)}{T} \\ P(X=x_i | P_{ai}=\pi) = \frac{\text{count}(X=x_i, P_{ai}=\pi)}{\text{count}(P_{ai}=\pi)} = \frac{\sum_i I(X_i^{(t)}, X_i) I(P_{ai}^{(t)}, \pi)}{\sum_i I(P_{ai}^{(t)}, \pi)} \end{array} \right. \end{aligned}$$

· learning with partial data, generally:

Nodes: $\{X_1, \dots, X_n\}$ in BN.

unobserved nodes: H

observed nodes: V

$$H \cup V = \{X_1, \dots, X_n\} . H \cap V = \emptyset$$

→ goal: estimate CPTs to maximize Prob of given data.

i.e. $\underset{CPTs}{\operatorname{argmax}} \text{p(data} | \text{CPTs)}$

$$\begin{aligned} \rightarrow \text{method: } & \underset{CPTs}{\operatorname{argmax}} \log L , \quad L = \log \prod_{t=1}^T P(V=v^t) = \sum_{t=1}^T \log P(V=v^t) \\ & = \sum_{t=1}^T \log \sum_h P(V=v^t, H=h) \\ & = \sum_{t=1}^T \log \sum_h \left(\prod_{i=1}^n P(X_i=x_i | P_{ai}=\pi_i) \right) \end{aligned}$$

$$\frac{\partial L}{\partial P(X_i=x_i | P_{ai}=\pi)} = 0 \Rightarrow \text{no closed form solution.}$$

→ Idea: we use a iterative algorithm.

* EM (Expectation Maximization) alg:

$$\left. \begin{array}{l} \text{Node with parents: } P(X_i=x | P_{a_i}=\pi_i) = \frac{\text{count}(X_i=x, P_{a_i}=\pi_i)}{\text{count}(P_{a_i}=\pi_i)} \\ \text{root nodes: } P(X_i=x) = \frac{\text{count}(X_i=x)}{T} \end{array} \right.$$

* EM alg:

Init para of the model

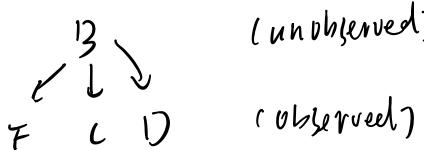
while para not converged:

For each para p :

- E-step: Calculate expected counts by using para from last iteration
- M-step: Use expected counts to update P

* Initialization:

- random guess.
- in $[0, 1]$
- EM may not produce global max

e.g.  we want $\{P(B)$
 $P(F|B), P(C|B), P(D|B)\}$

$$P(B=1) = \frac{\text{count}(B=1)}{T} \rightarrow \text{we don't have.}$$

$$\rightarrow \text{use expected count} = \sum_{t=1}^T P(B=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

$$P(B=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

$$= \frac{P(B=1, F=f^{(t)}, C=c^{(t)}, D=d^{(t)})}{P(B=1, F=f^{(t)}, C=c^{(t)}, D=d^{(t)}) + P(B=0, F=f^{(t)}, C=c^{(t)}, D=d^{(t)})}$$

$$\hookrightarrow P(B=1, F=f^{(t)}, C=c^{(t)}, D=d^{(t)}) = P(B) \cdot P(F=f^{(t)}|B) \cdot P(C=c^{(t)}|B) \cdot P(D=d^{(t)}|B)$$



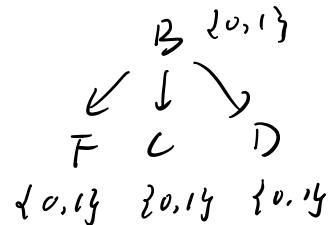
previous guess

$$\begin{aligned} \text{count}(B=1) &= \sum_{t=1}^T P(B=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)}) \\ &= \sum_{t=1}^T \frac{P(B=1) \cdot P(F=f^{(t)} | B=b) \cdot P(C=c^{(t)} | B=b) \cdot P(D=d^{(t)} | B=b)}{\sum_b P(B=b) \cdot P(F=f^{(t)} | B=b) \cdot P(C=c^{(t)} | B=b) \cdot P(D=d^{(t)} | B=b)} \end{aligned}$$

E-step : use above method to calculate count($B=1$)

M-step : update $P(B=1) \leftarrow \frac{\text{count}(B=1)}{T}$

$P(F=1 | B=1) = \frac{\text{count}(B=1, F=1)}{\text{count}(B=1)}$



$$\text{Count}(B=1, F=1) = \sum_{t=1}^T P(B=1, F=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

$$= \sum_{t=1}^T P(B=1 | F=1, C=c^{(t)}, D=d^{(t)}) \cdot I(f^{(t)}, 1)$$

$$= \sum_{t=1}^T \frac{P(B=1, F=1, C=c^{(t)}, D=d^{(t)}) \cdot I(f^{(t)}, 1)}{\sum_b P(B=b, F=1, C=c^{(t)}, D=d^{(t)})}$$

$$= \sum_{t=1}^T \frac{P(B=1) \cdot P(F=1 | B=1) \cdot P(D=d^{(t)} | B=1) \cdot P(C=c^{(t)} | B=1) \cdot I(f^{(t)}, 1)}{\sum_b P(B=b) \cdot P(F=1 | B=b) \cdot P(D=d^{(t)} | B=b) \cdot P(C=c^{(t)} | B=b)}$$

~~~~~

↳ from last iteration.

Generally, EM :  $H$  — hidden nodes

$V$  — visible nodes

$$\begin{aligned} \text{Nodes with parents : } P(X_i=x | \text{Pa}_i=\pi) &= \frac{\text{count}(X_i=x, \text{Pa}_i=\pi)}{\text{count}(\text{Pa}_i=\pi)} \\ &= \frac{\sum_t P(X_i=x, \text{Pa}_i=\pi | \vec{V}^{(t)})}{\sum_t P(\text{Pa}_i=\pi | \vec{V}^{(t)})} \end{aligned}$$

$$\text{Roots : } P(X_v=x) = \frac{\text{count}(X_v=x)}{T} = \frac{\sum_t P(X_v=x | \vec{V}^{(t)})}{T}$$

• EM properties:

1. Monotonic convergence: each iteration increases log-likelihood of observed data
2. No tuning para, lr. backtracking.
3. converges to a local or global maximum.

e.g.:  $\begin{matrix} \textcircled{1} & \rightarrow & \textcircled{0} & \rightarrow & \textcircled{1} \\ x & & x & & z \end{matrix}$  use EM alg to update  $P(Z=z | Y=y)$

$$V = \{X, Z\}, H = \{Y\}$$

$$P(Z=z | Y=y) = \frac{\text{Count}(Z=z, Y=y)}{\text{Count}(Y=y)}$$

$$= \frac{\sum_t P(Z=z, Y=y | X=x^{(t)}, Z=z^{(t)})}{\sum_t P(Y=y | X=x^{(t)}, Z=z^{(t)})}$$

$$P(Y=y | X=x^{(t)}, Z=z^{(t)})$$

$$= \frac{P(Y=y, X=x^{(t)}, Z=z^{(t)})}{\sum_y P(Y=y, X=x^{(t)}, Z=z^{(t)})}$$

$$= \frac{\sum_t P(Y=y | X=x^{(t)}, Z=z^{(t)}) \cdot I(Z, z^{(t)})}{\sum_t P(Y=y | X=x^{(t)}, Z=z^{(t)})}$$

$$P(Y=y, X=x^{(t)}, Z=z^{(t)})$$

$$= P(X=x^{(t)}) \cdot P(Y=y | X=x^{(t)}) \cdot P(Z=z^{(t)} | Y=y)$$

$$\therefore P(Z=z | Y=y) = \frac{\sum_t \frac{P(X=x^{(t)}) \cdot P(Y=y | X=x^{(t)}) \cdot P(Z=z^{(t)} | Y=y)}{\sum_y P(X=x^{(t)}) \cdot P(Y=y | X=x^{(t)}) \cdot P(Z=z^{(t)} | Y=y)} \cdot I(Z, z^{(t)})}{\sum_t \frac{P(X=x^{(t)}) \cdot P(Y=y | X=x^{(t)}) \cdot P(Z=z^{(t)} | Y=y)}{\sum_y P(X=x^{(t)}) \cdot P(Y=y | X=x^{(t)}) \cdot P(Z=z^{(t)} | Y=y)}}$$

(binary variables)

↳ CPTs from last iteration.

e.g2:  $\begin{matrix} \textcircled{1} & \leftarrow & \textcircled{0} & \rightarrow & \textcircled{1} \\ A & & B & & C \end{matrix}$

$$\text{CPTs: } P(B=1), P(A=1 | B=1), P(C=1 | B=1), P(A=1 | B=0), P(C=1 | B=0)$$

A, C observed

B unobserved

$$P(B=1) = \frac{\text{Count}(B=1)}{T} = \frac{\sum_t P(B=1 | A=a^{(t)}, C=c^{(t)})}{T}$$

$$\text{Data: } \{(A_t, C_t)\}_{t=1}^T$$

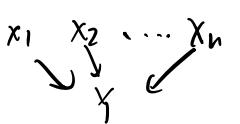
$$\begin{array}{c|cc} & a_1 & c_2 \\ \hline 1 & 0 & 0 \\ 2 & 1 & 0 \\ \vdots & & \end{array}$$

$$P(A=1 | B=1) = \frac{\text{Count}(A=1, B=1)}{\text{Count}(B=1)} = \frac{\sum_t P(A=1, B=1 | A=a^{(t)}, C=c^{(t)})}{\sum_t P(B=1 | A=a^{(t)}, C=c^{(t)})}$$

$$\text{CPTs } \leftarrow \text{Full Joint} = \frac{\sum_t P(B=1 | A=a^{(t)}, C=c^{(t)}) \cdot I(A, a^{(t)})}{\sum_t P(B=1 | A=a^{(t)}, C=c^{(t)})}$$

Other CPTs can be updated in the same way.

\* EM for learning Noisy-OR model.



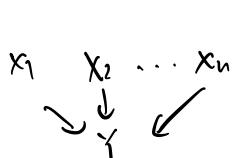
- Data:  $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$
- $P(Y=1 | x_1, x_2, \dots, x_n) = \prod_{i=1}^n (1 - p_i)^{x_i}$

- problem: estimate  $p_i \in [0, 1]$

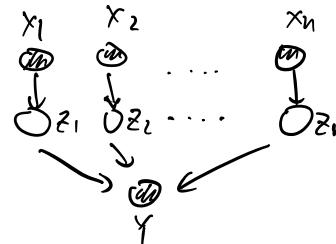
• Why don't we just use MLE since we have fully observed data?

→ we can't get a closed form.

• Method:



Transformed



CPTs:

$$P(z_i=0 | x_i=0) = 1$$

$$P(z_i=1 | x_i=1) = p_i$$

$$P(Y|z) = \text{OR}(z)$$

$$\cdot \text{First show that: } P(Y=1 | \vec{x}) = \sum_{\vec{z}} P(Y=1 | \vec{z}) \cdot P(\vec{z} | \vec{x})$$

$$\begin{aligned} P(Y=1 | \vec{x}) &= \sum_{\vec{z}} P(Y=1, \vec{z} | \vec{x}) = \sum_{\vec{z}} P(Y=1 | \vec{x}, \vec{z}) \cdot P(\vec{z} | \vec{x}) \\ &= \sum_{\vec{z}} P(Y=1 | \vec{z}) \cdot P(\vec{z} | \vec{x}) \end{aligned}$$

$$\cdot Q: \text{when does } P(Y=1 | \vec{z}) \cdot P(\vec{z} | \vec{x}) = 0 ?$$

→ all  $z_i \in \vec{z}$  equals to 0.

$$\cdot P(Y=1 | \vec{x}) = \sum_{\vec{z} \neq 0} P(Y=1 | \vec{z}) \cdot P(\vec{z} | \vec{x}) = 1 - \sum_{\vec{z}=0} P(\vec{z} | \vec{x})$$

• update  $P(z_i=1 | x_i=1)$  by EM alg:

$$P(z_i=1 | x_i=1) = \frac{\text{Count}(z_i=1, x_i=1)}{\text{Count}(x_i=1)} \rightarrow \text{from data}$$

$$\text{Count}(z_i=1, x_i=1) = \sum_{t=1}^T P(z_i=1, x_i=1 | \vec{x}^{(t)}, y^{(t)})$$

$$= \sum_{t=1}^T x_i^{(t)} \cdot P(z_i=1 | \vec{x}^{(t)}, y^{(t)})$$

( $x_i^{(t)} \Leftrightarrow I(x_i^{(t)}, 1)$ , numerically)

$$P(z_i=1 | \vec{x}^{(t)}, y^{(t)}) = \frac{P(x_i=y^{(t)} | z_i=1) \cdot P(z_i=1 | \vec{x}^{(t)})}{P(Y=y^{(t)} | \vec{x}^{(t)})}$$

| Data: |          | Y     |   |
|-------|----------|-------|---|
| $x_1$ | $\dots$  | $x_n$ |   |
| 0     | $\dots$  | 1     | 1 |
| 0     | $\dots$  | 0     | 0 |
| :     | $\vdots$ | :     | : |
| 0     | $\vdots$ | 1     | 0 |
|       |          |       | T |

$$= \frac{y^{(t)} \cdot p_i}{1 - \prod_{j=1}^n (1-p_j)^{x_j}}$$

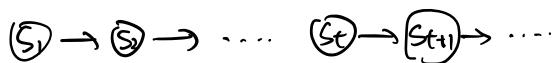
$$\Rightarrow \text{count}(z_i=1, x_i=1) = \sum_{t=1}^T \left[ \frac{p_i \cdot y^{(t)} \cdot x_i^{(t)}}{1 - \prod_{j=1}^n (1-p_j)^{x_j}} \right]$$

$y^{(t)}$ : OK relationship  
 $p_i = P(z_i=1 | \vec{x})$ ;  $x_i=1$  and  
d-sep

$$\Rightarrow p_i \leftarrow \frac{\text{Count}(z_i=1, x_i=1)}{\text{Count}(x_i=1)} = \frac{p_i}{T_i} \sum_{t=1}^T \left[ \frac{y^{(t)} \cdot x_i^{(t)}}{1 - \prod_{j=1}^n (1-p_j)^{x_j}} \right]$$

## Day 10 HMM (hidden markov model):

- Markov chain:



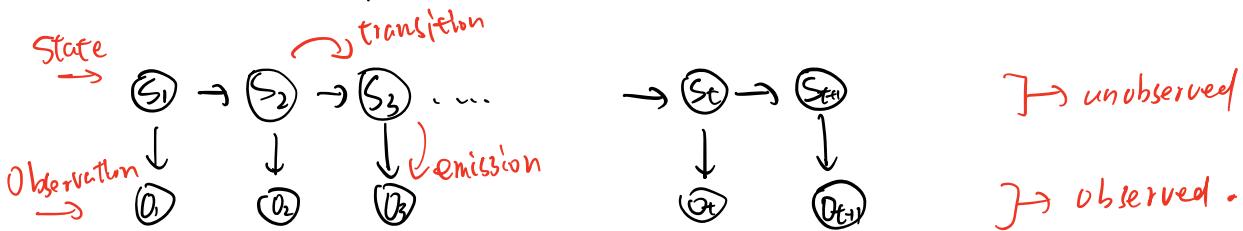
2 assumptions:

I.  $P(S_t | S_1, \dots, S_{t-1}) = P(S_t | S_{t-1})$  ————— finite context

II.  $P(S_{t+k} = j | S_t = i) = P(S_{t+k+1} = j | S_{t+k} = i)$  ————— position invariant.

(we can save lots of CPTs with these 2 assumptions)

- HMM: extension of Markov chain.



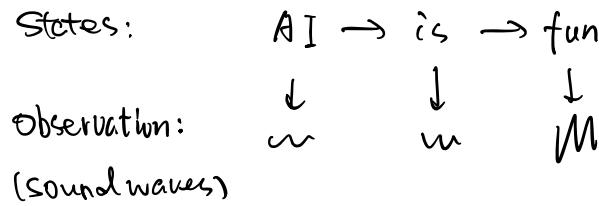
Assumptions:

I.  $\begin{cases} P(S_t | S_1, \dots, S_{t-1}) = P(S_t | S_{t-1}) \\ P(O_t | S_1, \dots, S_t) = P(O_t | S_t) \end{cases}$

II.  $\begin{cases} P(S_{t+1} | S_t) = P(S_{t+k+1} | S_{t+k}) \\ P(O_t | S_t) = P(O_{t+k} | S_{t+k}) \end{cases}$

(CPTs are matrices)

eg: HMM for language:



\* Conditional Independence in HMM:

- $p(s_t | s_{t-1}) \neq p(s_t | s_{t-1}, o_t)$  is obvious, which means observation affects the prob.
- $P(s_2, s_3, \dots, s_T | s_1) = \prod_{t=2}^T p(s_t | s_{t-1})$

Show it:

$$\begin{aligned}
 P(s_2 \dots s_T | s_1) &= p(s_T | s_{T-1} \dots s_1) \cdot p(s_{T-1} | s_{T-2} \dots s_1) \\
 &= p(s_T | s_{T-1} \dots s_1) \cdot p(s_{T-1} | s_{T-2} \dots s_1) \cdot p(s_{T-2} \dots s_2 | s_1) \\
 &\quad \vdots \\
 &= p(s_T | s_{T-1} \dots s_1) \cdot p(s_{T-1} | s_{T-2} \dots s_1) \cdot \dots \cdot p(s_2 | s_1) \\
 &= p(s_T | s_{T-1}) \dots p(s_2 | s_1) \\
 &= \prod_{t=2}^T p(s_t | s_{t-1})
 \end{aligned}$$

Compute:  $P(O_1, O_2, \dots, O_T)$

$$P(O_1, O_2, \dots, O_T) = \sum_{i=1}^n P(O_1, \dots, O_T, s_i=i)$$

$$\text{def: } \alpha_{it} = P(O_1, \dots, O_t, s_t=i)$$

$$\text{transition prob : } \alpha_{ij} = p(s_t=i | s_{t-1}=j)$$

$$\text{emission prob : } b_{jk} = p(O_t=o_k | s_t=j)$$

$$\text{init prob : } \pi_i = p(s_1=i)$$

We want to get a recurrent relationship of  $\alpha_{it}$

$$\alpha_T(s_T) = P(O_1, \dots, O_T, s_T)$$

$$= \sum_{k=1}^n P(O_1, \dots, O_T, s_T, s_{T-1}=k)$$

$$= \sum_{k=1}^n P(O_T | s_T, s_{T-1}=k, O_1, \dots, O_{T-1}) \cdot P(s_T | s_{T-1}=k, O_1, \dots, O_{T-1}) \cdot P(O_1, \dots, O_{T-1}, s_{T-1}=k)$$

$$\begin{aligned}
 \alpha_{t+1}^{\text{sep}} &= \sum_{k=1}^n p(o_1 | s_1) \cdot p(s_T | s_{T-1}=k) \cdot p(o_1 \dots o_{T-1}, s_{T-1}=k) \\
 &= \sum_{k=1}^n p(o_T | s_T) \cdot p(s_T | s_{T-1}=k) \cdot \alpha_{T-1}(s_{T-1}=k)
 \end{aligned}$$

↓                    ↓                    ↓  
 emission      transition      recursion.

$$\therefore \alpha_{t+1} = \sum_{i=1}^n \alpha_i \alpha_{i,t} a_{ij} b_j o_{t+1} , \quad \alpha_{i,1} = \pi_i b_{i,0} \quad \longrightarrow p(o_1, s_1) = p(o_1 | s_1) \cdot p(b_1)$$

$$\therefore p(o_1, \dots, o_T) = \sum_{i=1}^n \alpha_i \alpha_{T-1}(i) = \sum_{i=1}^n \sum_{k=1}^n \alpha_{T-1}(k) \alpha_{i,T-1} a_{ij} b_j o_T$$

for  $N$  states,  $N$  observation, complexity:  $\Theta(n^2N)$