

Review

- Learning CPTs from incomplete data

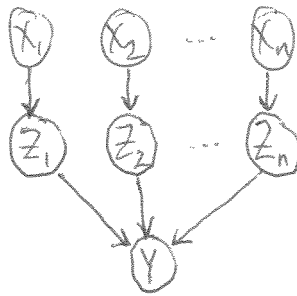
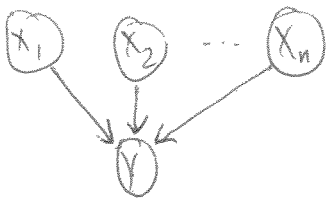
EM update $P(X_i=x | p_{ai}=\pi) \leftarrow \frac{\sum_t P(X_i=x, p_{ai}=\pi | v^{(t)})}{\sum_t P(p_{ai}=\pi | v^{(t)})}$ for nodes w/ parents

$P(X_i=x) \leftarrow \frac{1}{T} \sum_t P(X_i=x | v^{(t)})$ for root nodes

where $v^{(t)}$ denotes visible nodes

- Ex: Noisy-OR

Hidden variable model



$$P(Y=1 | \vec{X}) = 1 - \prod_{i=1}^n (1-p_i)^{X_i}$$

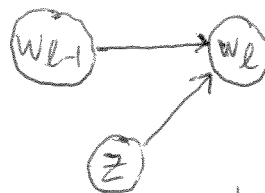
Same as noisy-OR!

How to estimate p_i from data $\{(\vec{x}_t, y_t)\}_{t=1}^T$?

EM update: $p_i = P(Z_i=1 | X_i=1) \leftarrow \frac{1}{T_i} \sum_{t=1}^T \frac{y_t x_{it} p_i}{1 - \prod_{j=1}^n (1-p_j)^{x_{jt}}}$ from posterior $P(Z_i=1 | \vec{X}, Y)$, computed in terms of CPTs

- Ex: Mixture of n-gram models

Hidden variable model



$$P(w_l | w_{l-1}) = \lambda P_1(w_l) + (1-\lambda) P_2(w_l | w_{l-1}) = P_m(w_l | w_{l-1})$$

How to estimate λ from data $\{(w_{l-1}, w_l)\}_{l=1}^L$?

EM update: $\lambda = P(Z=1) \leftarrow \frac{1}{L} \sum_{l=1}^L P(Z=1 | w_{l-1}, w_l)$ posterior computed in terms of CPTs

Hidden Markov Models (HMMs)

• Random variables

$S_t \in \{1, 2, \dots, n\}$ state at time t

$O_t \in \{1, 2, \dots, m\}$ observation at time t

"noisy" reflection of hidden state S_t

• Ex: puppy training

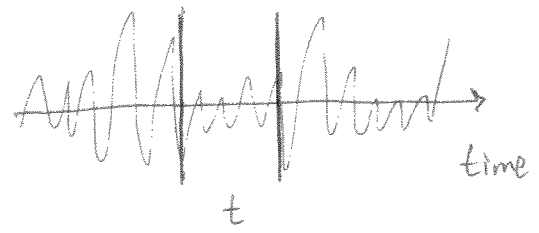
$S = \{ \text{"have-to-go"}, \text{"don't-have-to-go"}, \text{"went"} \}$

$O = \{ \text{"wagging tail"}, \text{"whimpering"}, \text{"running in circles"}, \text{"hiding in corners"} \}$

• Ex: speech recognition

$S =$ units of language: words, syllables, phonemes

$O =$ acoustic measurements

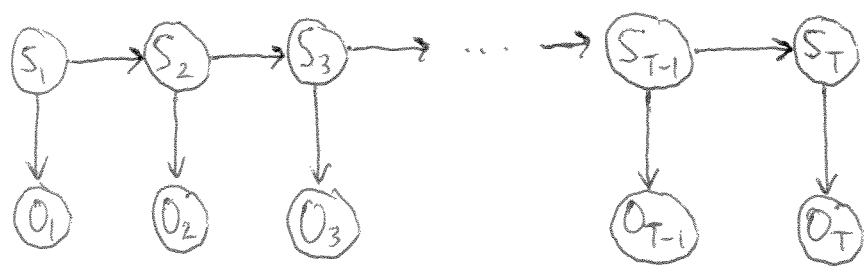


• Ex: robotics

$S_t =$ location

$O_t =$ sensor readings

• Belief Network of HMM



Polytree?
YES!

• Markov assumptions

- finite context

$$P(S_t | S_1, S_2, \dots, S_{t-1}) = P(S_t | S_{t-1})$$

$$P(O_t | S_1, S_2, \dots, S_{t-1}, S_t, S_{t+1}, \dots; S_T) = P(O_t | S_t)$$

- shared CPTs

$$P(S_{t+1} = s' | S_t = s) = P(S_t = s' | S_{t-1} = s)$$

$$P(O_t = o | S_t = s) = P(O_{t+1} = o | S_{t+1} = s)$$

• Joint distribution

$$P(\vec{S}, \vec{O}) = P(S_1) \cdot \left\{ \prod_{t=2}^T P(S_t | S_{t-1}) \right\} \cdot \left\{ \prod_{t=1}^T P(O_t | S_t) \right\}$$

$\downarrow \quad \downarrow \quad \downarrow$
 $(S_1, S_2, \dots, S_T) \quad (O_1, O_2, \dots, O_T) \quad \text{initial state}$

• Parameters (CPTs)

$\pi_i = P(S_1 = i)$ initial state distribution

$a_{ij} = P(S_{t+1} = j | S_t = i)$ transition matrix ($n \times n$)

$b_{ik} = P(O_t = k | S_t = i)$ emission matrix ($n \times m$)

For clarity: $b_{ik} = b_i(k)$ alternate notation

Ex: isolated word speech recognizer

recognize word "CAT"

build HMM that assigns high probability to "CAT" utterances
low probability to other utterances

use HMM with 5 states

state #	sound
1	initial silence
2	"C"
3	"A"
4	"T"
5	final silence

$$\pi = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

must start in state 1

$$a_{ij} = \begin{bmatrix} 0.9 & 0.1 & \emptyset & \emptyset \\ & 0.9 & 0.1 & \emptyset \\ & & 0.99 & 0.01 \\ \emptyset & & & 0.4 & 0.6 \\ & & & & 1 \end{bmatrix}$$

upper diagonal transition matrix
special case: left to right HMM

Key Questions for HMMs

Inference: given $\{\pi_i, a_{ij}, b_{ik}\}$ parameters

1) How to compute likelihood $P(o_1, o_2, \dots, o_T)$?

2) How to compute most likely state sequence?

$$(s_1^*, s_2^*, \dots, s_T^*) = \underset{s_1, s_2, \dots, s_T}{\operatorname{argmax}} P(\underbrace{s_1, s_2, \dots, s_T}_{\text{sequence of states}} | o_1, o_2, \dots, o_T)$$

sequence of states

that maximize posterior probability

3) How to compute $P(S_t = i | o_1, o_2, \dots, o_t)$?

updating belief in real-time

Learning: given $\{o_1, o_2, \dots, o_T\}$ observations

4) How to estimate parameters $\{\pi_i, a_{ij}, b_{ik}\}$ that maximize likelihood $P(o_1, o_2, \dots, o_T)$?

EM algorithm!

1) Computing likelihood

$$P(o_1, o_2, \dots, o_T) = \sum_{\vec{s}} \overbrace{P(s_1, s_2, \dots, s_T)}^{\text{hidden}} \overbrace{P(o_1, o_2, \dots, o_T)}^{\text{observed}}$$

Sum over n^T hidden state sequences

$$= \sum_{\vec{s}} P(s_1) \prod_{t=2}^T P(s_t | s_{t-1}) \prod_{t=1}^T P(o_t | s_t)$$

• Efficient recursion

$$P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1}=j) = \sum_{i=1}^n P(o_1, \dots, o_t, o_{t+1}, s_{t+1}=j, s_t=i) \quad \text{marginalization}$$

$$= \sum_{i=1}^n P(o_1, \dots, o_t, s_t=i) P(s_{t+1}=j, o_{t+1} | s_t=i, o_1, \dots, o_t) \quad \text{product rule}$$

$$= \sum_{i=1}^n P(o_1, \dots, o_t, s_t=i) P(s_{t+1}=j, o_{t+1} | s_t=i) \quad \text{conditional independence}$$

$$= \sum_{i=1}^n P(o_1, \dots, o_t, s_t=i) P(s_{t+1}=j | s_t=i) P(o_{t+1} | s_{t+1}=j, s_t=i) \quad \text{product rule}$$

$$= \sum_{i=1}^n \underbrace{P(o_1, \dots, o_t, s_t=i)}_{\text{recursive instance}} \underbrace{P(s_{t+1}=j | s_t=i) P(o_{t+1} | s_{t+1}=j)}_{\text{CPTs}} \quad \text{conditional independence}$$

- Shorthand notation

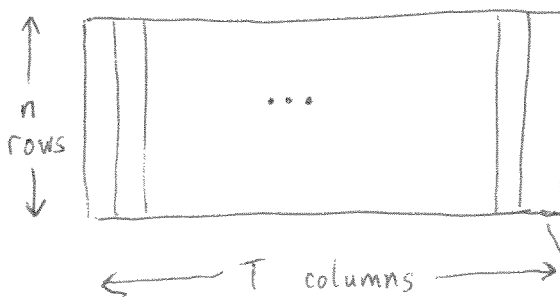
$$\alpha_{it} = P(o_1, o_2, \dots, o_t, S_t = i)$$

↑

$i = 1, \dots, n$ # hidden states

$t = 1, \dots, T$ sequence length

α matrix



sum of last column
is $P(\vec{O})$

- Forward algorithm

recursive step:
$$\alpha_{jt+1} = \sum_{i=1}^n \alpha_{it} a_{ij} b_j(o_{t+1})$$

- Initial condition (1st column of α)

$$\alpha_{i1} = P(o_1, S_1 = i) = P(S_1 = i) P(o_1 | S_1 = i) = \pi_i b_i(o_1) \quad \text{for } i = 1, \dots, n$$

- Back to likelihood computation

$$\begin{aligned} P(o_1, o_2, \dots, o_T) &= \sum_{i=1}^n P(o_1, o_2, \dots, o_T, S_T = i) && \text{marginalization} \\ &= \sum_{i=1}^n \alpha_{iT} \end{aligned}$$

- Scales as $O(n^2 T)$

linear, not exponential, in sequence length
quadratic in # of states

- Warning: naive calculations will underflow for long sequences

because $P(o_1, o_2, \dots, o_T) \ll 1$