

REVIEW

- EM updates

• root nodes : $P(X_i = x) \leftarrow \frac{1}{T} \sum_t P(X_i = x | V^{(t)})$
(nodes without parents)

• nodes with parents : $P(X_i = x | P_{A_i} = \pi) \leftarrow \frac{\sum_t P(X_i = x, P_{A_i} = \pi | V^{(t)})}{\sum_t P(P_{A_i} = \pi | V^{(t)})}$

Updates converge monotonically in log-likelihood $\sum_t \log P(V^{(t)})$

Example #1



Incomplete data set : $\{(a_t, c_t)\}_{t=1}^T$

EM updates : $P(B=b | A=a) \leftarrow \frac{\sum_t I(a, a_t) P(b | a_t, c_t)}{\sum_t I(a, a_t)}$

$P(C=c | B=b) \leftarrow \frac{\sum_t I(c, c_t) P(b | a_t, c_t)}{\sum_t P(b | a_t, c_t)}$

Application #1

word clustering



$W, W' \in \{1, 2, \dots, V\}$
 $Z \in \{1, 2, \dots, k\}$ ($k \ll V$)

Example #2



Hidden : H

Observed : A, B, C

• Posterior $P(H | A, B, C) = \frac{P(C | H, A, B) P(H | A, B)}{P(C | A, B)}$

Bayes Rule

$= \frac{P(C | H, A, B) P(H)}{\sum_h P(C | H=h, A, B) P(H=h)}$

marginal independence
normalization

* Incomplete data set $\{(a_t, b_t, c_t)\}_{t=1}^T$

log (conditional) likelihood: $\mathcal{L} = \sum_t \log P(c_t | a_t, b_t)$

$$= \sum_t \log \left(\sum_h P(c_t, h | a_t, b_t) \right) \quad \text{marginalization}$$

$$= \sum_t \log \sum_h \left[P(h | a_t, b_t) \cdot P(c_t | a_t, b_t, h) \right] \quad \begin{array}{l} \text{Product rule} \\ \text{marginal independence} \end{array}$$

* EM update for CPT at node H : $P(H=h) \leftarrow \frac{1}{T} \sum_t P(H=h | a_t, b_t, c_t)$

* Aside: $P(V_1=v, H=h | V_1^{(t)}, V_2^{(t)}) = I(v, V_1^{(t)}) \cdot P(h | V_1^{(t)}, V_2^{(t)})$

Application #2: Linear Interpolation of Markov Models

$$P_M(W_L | W_{L-1}, W_{L-2}) = \lambda_1 P_1(W_L) + \lambda_2 P_2(W_L | W_{L-1}) + \lambda_3 P(W_L | W_{L-1}, W_{L-2})$$

mixture model

* Suppose n -gram models are trained on large corpus A

* How to estimate λ_i ; where $\lambda_i \geq 0$ and $\sum_i \lambda_i = 1$?

* Methodology: Train P_1, P_2, P_3 on corpus A ; fix these models

Train $\lambda_1, \lambda_2, \lambda_3$ on corpus C

Estimate λ_i to maximize log-likelihood of mixture model on corpus C .

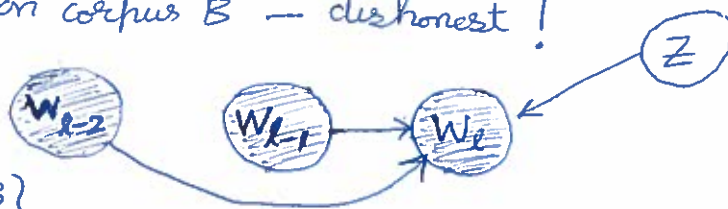
* Don't use corpus A to estimate λ_i .

- otherwise you find: $\lambda_3 = 1$; $\lambda_1 = \lambda_2 = 0$.

* Test $P_M = \sum_{i=1}^3 \lambda_i P_i$ on corpus B .

- Don't estimate λ_i on corpus B — dishonest!

* Hidden Variable Model:



$$P(Z=i) = \lambda_i ; Z \in \{1, 2, 3\}$$

Define CPT at node W_L :

$$P(W_L | W_{L-1}, W_{L-2}, Z) = \begin{cases} P_1(W_L) & ; \text{if } Z=1 \\ P_2(W_L | W_{L-1}) & ; \text{if } Z=2 \\ P_3(W_L | W_{L-1}, W_{L-2}) & ; \text{if } Z=3 \end{cases}$$

In this model,

$$P(W_L | W_{L-1}, W_{L-2}, Z) = \sum_{z=1}^3 P(W_L, Z=z | W_{L-1}, W_{L-2}) \quad \text{marginalization}$$

$$= \lambda_1 P_1(W_L) + \lambda_2 P_2(W_L | W_{L-1}) + \lambda_3 P_3(W_L | W_{L-1}, W_{L-2})$$

EM update : $\frac{P(Z=i)}{\lambda_i} \leftarrow \frac{1}{L_c} \sum_{\ell} P(Z=i | W_{\ell}, W_{\ell-1}, W_{\ell-2})$
 \downarrow
 Length of corpus C

• If we change the BNs :



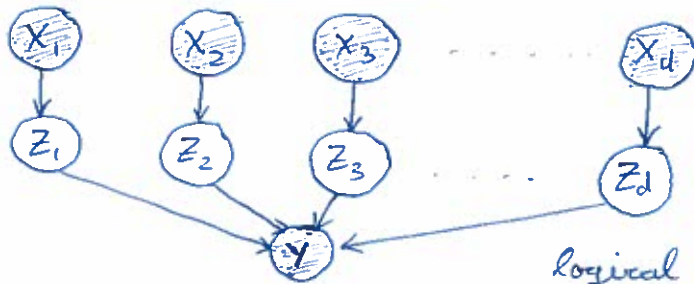
Example / Application #3

: NOISY-OR (HW 6)



• One approach to ML : gradient ascent / Newton's method for complete data and parameterized CPTs

Consider Hidden Variable Model :



logical-OR : $Y = \text{OR}(Z_1, Z_2, \dots, Z_d)$

Show $P(Y=1 | X_1, \dots, X_d)$ in above model matches noisy-OR.

Hidden Markov Models

* Random Variables : $S_t \in \{1, 2, 3, \dots, n\}$ (hidden) state at time t

$O_t \in \{1, 2, 3, \dots, m\}$ observed at time t

Observations O_t are noisy, partial reflections of true underlying state S_t at time t .

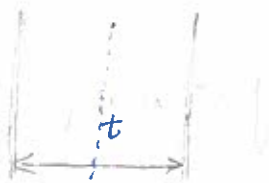
Ex : house-training puppy

$S_t \in \{ \text{has-to-go}, \text{doesn't need to go}, \text{went} \}$

$O_t \in \{ \text{barking}, \text{waiting by door}, \text{sleeping}, \text{hiding}, \dots \}$

$P(S_t=i | O_1, O_2, \dots, O_t)$ do I take the puppy outside?

Ex: Speech Recognition



O_t = acoustic measurements of windowed signal around time t .

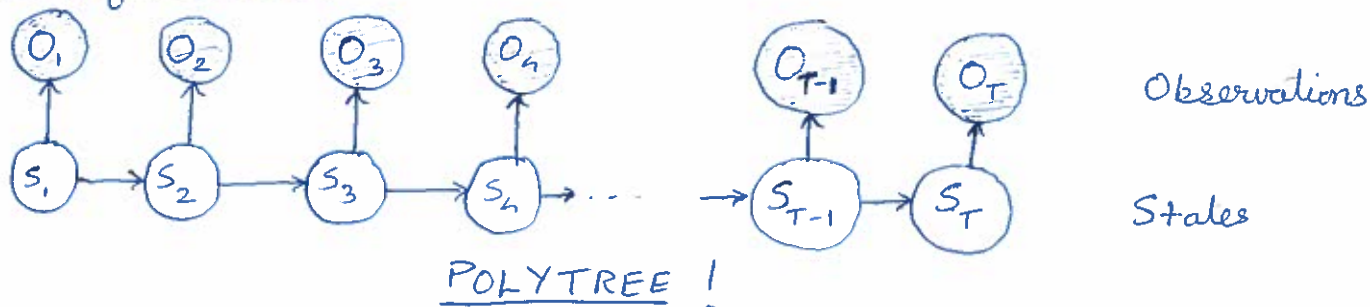
S_t = unit of language (eg. phoneme, letter, syllable) being uttered at time t .

Ex: Robotics

O_t : sensor readings

S_t : location / orientation

* Belief Network



* Conditional Independence assumptions

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

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

* CPTs are shared across time

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

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

* Joint Distribution

$$P(\underbrace{S_1, S_2, \dots, S_{T-1}, S_T}_{\vec{S}}, \underbrace{O_1, O_2, \dots, O_{T-1}, O_T}_{\vec{O}})$$

$$= P(\underbrace{S_1}_{\substack{\downarrow \\ \text{initial} \\ \text{state}}}) \left\{ \prod_{t=2}^T P(S_t | S_{t-1}) \right\} \left\{ \prod_{t=1}^T P(O_t | S_t) \right\}$$

* Parameters

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

$$b_i(k) = b_{ik} = P(O_t = k \mid S_t = i) \quad n \times m \text{ emission matrix}$$

$$\pi_i = P(S_1 = i) \quad \text{initial state distribution. (n \times 1 vector)}$$

* Key Computations / Questions in HMMs

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

2) how to compute most likely hidden state sequence ?

$$\arg \max_{s_1, s_2, \dots, s_T} P(s_1, s_2, \dots, s_T \mid o_1, o_2, \dots, o_T)$$

} Inference

3) How to estimate parameters $\{\pi_i, a_{ij}, b_{ik}\}$ that maximize

$$P(o_1, o_2, \dots, o_T) \text{ or maybe } \prod_{\text{Sequence } i} P(o_1^{(i)}, o_2^{(i)}, \dots, o_T^{(i)}) ?$$

} Learning