

→ Learning CPTs from incomplete data

EM update:

$$P(x_i = \alpha | pa_i = \pi) \leftarrow \frac{\sum_t P(x_i = \alpha, pa_i = \pi | v^t)}{\sum_t P(pa_i = \pi | v^t)} \quad (\text{for roots with parents})$$

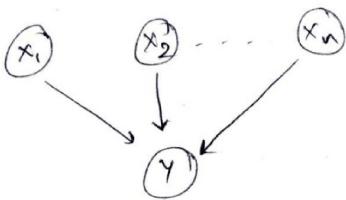
$$P(x_i = \alpha) \leftarrow \frac{1}{T} \sum_t P(x_i = \alpha | v^t) \quad (\text{root nodes})$$

where v^t denotes visible nodes.

→ HN4 | NOISY-OR

$$P(y=1 | \vec{x}) = 1 - \prod_{i=1}^n (1 - p_i)^{x_i}$$

we learn from data $\{(\vec{x}_t, y_t)\}_{t=1}^T$



EM update of this model is given by:

$$p_i \leftarrow \frac{1}{T_i} \sum_{t=1}^T \frac{y_t x_{it} p_i}{1 - \prod_{j=1}^n (1 - p_j)^{x_{jt}}}$$

where T_i is the number of examples in which $x_i = 1$

→ Linear interpolation (mixing) of n-gram models?

$$p_m(\omega_t | \omega_{t-1}) = \lambda p_1(\omega_t) + (1-\lambda) p_2(\omega_t | \omega_{t-1})$$

To estimate λ , we rewrite this as a hidden variable model

$$p(\omega_t | \omega_{t-1}, z) = \begin{cases} p_1(\omega_t) & \text{if } z=1 \\ p_2(\omega_t | \omega_{t-1}) & \text{if } z=2 \end{cases}$$

$z \in \{1, 2\}^Y$

$$p(z=1) = \lambda$$

$$p(z=2) = 1 - \lambda$$

In this model,

$$\begin{aligned} p(\omega_1 | \omega_{1-1}) &= \sum_{z'} p(\omega_1, z' | \omega_{1-1}) \\ &= \sum_{z'} p(z' | \omega_{1-1}) p(\omega_1 | z', \omega_{1-1}) \quad [\text{Product Rule}] \\ &= \sum_{z'} p(z') p(\omega_1 | z', \omega_{1-1}) \quad [\text{conditional independence}] \\ &= \lambda p_1(\omega_1) + (1-\lambda) p_2(\omega_1 | \omega_{1-1}) \end{aligned}$$

E-STEP

Compute posterior probability

$$* p(z | \omega_1, \omega_{1-1}) = \frac{p(\omega_1 | z, \omega_{1-1}) p(z | \omega_{1-1})}{p(\omega_1 | \omega_{1-1})} \quad [\text{Conditional Bayes' Rule}]$$

$$* p(z=1 | \omega_1, \omega_{1-1}) = \frac{\lambda p_1(\omega_1 | \omega_{1-1})}{\lambda p_1(\omega_1) + (1-\lambda) p_2(\omega_1 | \omega_{1-1})}$$

$$* p(z=2 | \omega_1, \omega_{1-1}) = 1 - p(z=1 | \omega_1, \omega_{1-1})$$

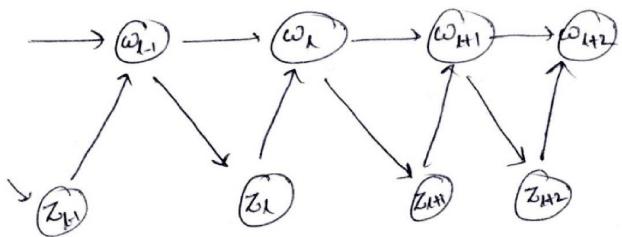
M-STEP

In this model

$$p(z=1) \leftarrow \frac{1}{L} \sum_{l=1}^L p(z=1 | \omega_{l-1}, \omega_l)$$



The above model over simplifies the belief network. In real world applications, smoothing parameter, λ would depend on the previous word.



$$P(z_t=1/\omega_{t-1}) = \lambda(\omega_{t-1}) \quad P(z_t=2/\omega_{t-1}) = 1 - \lambda(\omega_{t-1})$$

$$P(\omega_t/\omega_{t-1}) = \lambda(\omega_{t-1}) P_1(\omega_t) + 1 - \lambda(\omega_{t-1}) P_2(\omega_t/\omega_{t-1})$$

Using EM, we can estimate λ , which has the same size as the words in the vocabulary.

HIDDEN MARKOV MODELS

→ Random variables:

$$S_t = \{1, 2, \dots, n\} \quad \text{State at time } t$$

$$O_t = \{1, 2, \dots, m\} \quad \text{Observation at time } t$$

* Note that n, m have no relation ($n > m, n = m, n < m \leftarrow$ all possible)

* Observations O_t can be thought of as a noisy reflection of the hidden state S_t

→ Ex: puppy training.

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

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

→ Ex: Speech recognition!

s_t = units of language: words, syllables, phonemes

o_t = acoustic measurement: energy in a window, no. of peaks in a window...



→ Ex: Robotics

s_t = location

o_t = sensor reading

MARKOV ASSUMPTIONS IN HMM!

finite Context:

$$P(s_t | s_1, s_2, \dots, s_{t-1}) = P(s_t | s_{t-1})$$

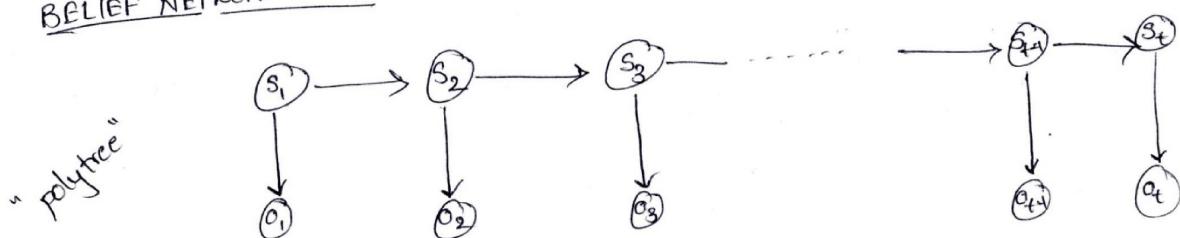
$$P(o_t | s_1, s_2, \dots, s_{t-1}, s_t, \dots, s_T) = P(o_t | s_t)$$

shared CPTs

$$P(s_{t+1} = s' | s_t = s) = P(s_t = s' | s_{t-1} = s)$$

$$P(o_{t+1} = o' | s_t = s) = P(o_{t+1} = o' | s_{t+1} = s')$$

BELIEF NETWORK OF HMM



PARAMETERS OF HMM:

$$\pi_i = \prod_{s_1=i} \quad [Initial State Distribution]$$

$$nxn \text{ matrix } a_{ij} = p(s_{t+1}=j | s_t=i) \quad [Transition matrix] \quad \begin{matrix} i=1,2,\dots,n \\ j=1,2,\dots,n \end{matrix}$$

$$nxm \text{ matrix } b_{ik} = p(o_t=k | s_t=i) \quad [Emission matrix] \quad \begin{matrix} i=1,2,\dots,n \\ k=1,2,\dots,m \end{matrix}$$

JOINT DISTRIBUTION OF HMM:

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

$\downarrow \quad \downarrow$

$s_1, s_2, \dots, s_T \quad o_1, o_2, \dots, o_T$

→ Ex: * Isolated word speech recognition

* Consider the problem of recognising the word "CAT"

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

* "

→ UND HMM with 5 states

State #	Sound
1	initial silence
2	"C"
3	"A"
4	"T"
5	final silence

$$\pi_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \leftarrow \text{always start with state 1}$$

$$a_{ij} = \begin{pmatrix} 0.9 & 0.1 & 0 & 0 & 0 \\ 0 & 0.9 & 0.1 & 0 & 0 \\ 0 & 0 & 0.99 & 0.01 & 0 \\ 0 & 0 & 0 & 0.4 & 0.6 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

→ this is an upper diagonal transition matrix
→ special case: left to right HMM

KEY QUESTIONS FOR HMMs

perform 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

$$(s_1^*, s_2^*, \dots, s_T^*) = \underset{s_1, s_2, \dots, s_T}{\operatorname{argmax}} P(s_1, s_2, \dots, s_T | o_1, o_2, \dots, o_T)$$

n^T possible settings.

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

↳ this can be thought of as 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 maximises likelihood $P(o_1, o_2, \dots, o_T)$ } EM algorithm.

$$\begin{aligned} 1) P(o_1, o_2, \dots, o_T) &= \sum_{s_1, s_2, \dots, s_T} p(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T) \\ &\quad \xrightarrow{\text{sum over } n^T \text{ hidden state sequences}} \\ &= \sum_s p(s_1) \prod_{t=2}^T p(s_t | s_{t-1}) \prod_{t=1}^T p(o_t | s_t) \end{aligned}$$