

1/2 Review - Markov models

* Random variables $S_t = \{1, 2, \dots, n\}$ state at time t .

* Belief network $(S_1) \rightarrow (S_2) \rightarrow \dots \rightarrow (S_{t-1}) \rightarrow (S_t)$

* Assumptions

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

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

* Weaknesses

- modeling k^{th} order correlations requires CPTs with $O(n^k)$ elements.

- assumes that the true state of the world can be observed.

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 .

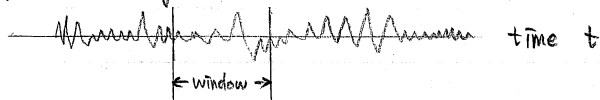
Observations O_t are noisy, partial reflections of states S_t .

Ex: toilet training

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

$O = \{ \text{neutral, funny walk, intense concentration, squat} \}$

Ex: speech recognition



O_t : acoustic measurements on windowed waveform at time t

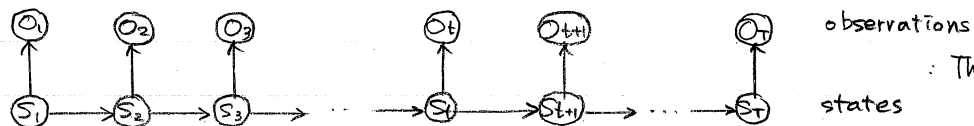
S_t : unit of language (word, syllable, phoneme)

Ex: robotics

O_t : sensor readings

S_t : location, orientation.

* Belief network



: This is a polytree!

* Assumptions

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

$P(O_t | S_1, S_2, \dots, S_T) = P(O_t | S_t)$

- 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

$$P(\underbrace{S_1, S_2, \dots, S_T}_S, \underbrace{O_1, O_2, \dots, O_T}_O) = P(S_1) \left[\prod_{t=2}^T P(S_t | S_{t-1}) \right] \left[\prod_{t=1}^T P(O_t | S_t) \right]$$

initial state

* Parameters

$$\pi_i = P(S_1 = i) \quad \text{initial state distribution}$$

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

$$b_{ik} = P(O_t = k | S_t = i) \quad \text{emission matrix}$$

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

* Key computations / questions in HMMs.

1) How to compute likelihood $P(O_1, O_2, \dots, O_T)$?

Ex: isolated word recognition.

2) How to compute most likely (hidden) state sequence

$$\arg \max_S P(S_1, S_2, \dots, S_T | O_1, O_2, \dots, O_T)$$

Ex: continuous speech recognition.

3) How to estimate parameters $\{\pi_i, a_{ij}, b_{ik}\}$ that maximize $P(O_1, O_2, \dots, O_T)$?
[or maybe multiple observation sequences]

inference
(assume parameters given)

learning

(1) Computing likelihood

$$P(O_1, O_2, \dots, O_T) = \sum_S P(S_1, S_2, \dots, S_T, O_1, \dots, O_T) \quad \text{marginalization}$$

$$= \sum_{S_1} P(S_1) \prod_{t=2}^T P(S_t | S_{t-1}) \prod_{t=1}^T P(O_t | S_t)$$

* Efficient recursion

$$P(O_1, O_2, \dots, O_t, O_{t+1}, S_{t+1} = j)$$

$$= \sum_{i=1}^n P(O_1, O_2, \dots, O_{t+1}, S_{t+1} = j, S_t = i) \quad \text{marginalization}$$

$$= \sum_{i=1}^n P(O_1, O_2, \dots, O_t, S_t = i) P(S_{t+1} = j | S_t = i, O_1, \dots, O_t) P(O_{t+1} | S_{t+1} = j, S_t = i, O_1, \dots, O_t)$$

conditional independence

product rule

conditional independence

$$= \sum_{i=1}^n \underbrace{P(O_1, O_2, \dots, O_t, S_t = i)}_{\text{recursion}} \underbrace{P(S_{t+1} = j | S_t = i) P(O_{t+1} | S_{t+1} = j)}_{\text{CPTs of HMM}}$$

* Shorthand notation

$$\alpha_{it} \triangleq P(O_1, O_2, \dots, O_t, S_t = i) \quad [n \times T \text{ matrix}]$$

$$\alpha_{j,t+1} = \sum_{i=1}^n \alpha_{it} a_{ij} b_j(O_{t+1}) \quad \text{"forward algorithm"}$$

\uparrow $b_{j,O_{t+1}}$



sum last column to get likelihood

base case: $\alpha_{i1} = P(O_1, S_1 = i) = P(S_1 = i) P(O_1 | S_1 = i) = \pi_i b_T(O_1)$

* likelihood

$$P(O_1, \dots, O_T) = \sum_{i=1}^N P(O_1, O_2, \dots, O_T, S_T = i) \quad \text{marginalization.}$$

$$= \sum_{i=1}^N \alpha_{iT}$$

* Warning: for long sequence, watch out for underflow.

(2) Computing most likely state sequence

$$S^* = \{s_1^*, s_2^*, \dots, s_T^*\}$$

$$= \operatorname{argmax}_S P(s_1, s_2, \dots, s_T | O_1, O_2, \dots, O_T) \quad \begin{array}{l} \text{constant with respect} \\ \text{to hidden states.} \end{array}$$

$$= \operatorname{argmax}_S \left[\frac{P(s_1, s_2, \dots, s_T, O_1, O_2, \dots, O_T)}{P(O_1, O_2, \dots, O_T)} \right]$$

$$= \operatorname{argmax}_S P(s_1, s_2, \dots, s_T, O_1, O_2, \dots, O_T)$$

Define $l_{it}^* = \max_{\{s_1, s_2, \dots, s_{t-1}\}} \log P(O_1, O_2, \dots, O_t, s_1, s_2, \dots, s_{t-1}, s_t = i)$



$= \log$ probability of most likely t -step "path"

that ends in state i at time t for observations $\{O_1, O_2, \dots, O_t\}$

* Form recursion

(i) base case ($t=1$)

$$l_{i1}^* = \log P(S_1 = i, O_1) = \log [P(S_1 = i) P(O_1 | S_1 = i)]$$

$$= \log \pi_i + \log b_i(O_1).$$

product rule +
conditional independence

(ii) from time t to time $t+1$:

$$l_{j,t+1}^* = \max_{s_1, \dots, s_t} \log P(s_1, s_2, \dots, s_t, s_{t+1} = j, O_1, O_2, \dots, O_{t+1})$$

$$= \max_{s_1, \dots, s_{t-1}} \max_i \left[\log P(s_1, \dots, s_{t-1}, s_t = i, O_1, \dots, O_t) P(s_{t+1} = j | s_t = i) \right]$$

\leftarrow representing state s_t

$$= \max_i \left[\max_{s_1, \dots, s_{t-1}} \log P(s_1, \dots, s_{t-1}, s_t = i, O_1, \dots, O_t) \right] + \log P(s_{t+1} = j | s_t = i)$$

$$+ \log P(O_{t+1} | s_{t+1} = j)$$

$$= \max_i \left[l_{it}^* + \log a_{ij} \right] + \log b_j(O_{t+1})$$

* How to derive S^* from l^* ?