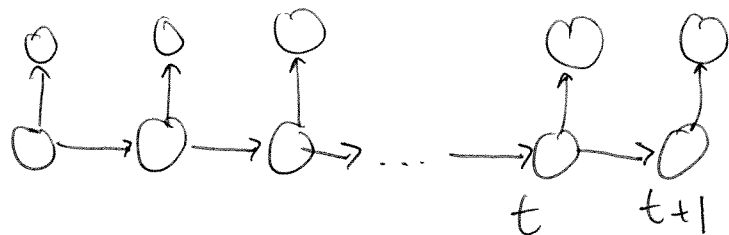


Review

* Hidden Markov models (HMMs)



observations $O_t \in \{1, 2, \dots, m\}$

states $S_t \in \{1, 2, \dots, n\}$

* Joint distribution

$$P(\vec{S}, \vec{O}) = P(s_1) \left[\prod_{t=2}^T P(s_t | s_{t-1}) \right] \left[\prod_{t=1}^T P(o_t | s_t) \right]$$

* Parameters

$$\pi_i = P(s_1 = i)$$

initial state distribution

$$a_{ij} = P(s_{t+1} = j | s_t = i)$$

transition matrix

$$b_{ik} = P(o_t = k | s_t = i)$$

emission matrix

* Key questions

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

2) How to decode hidden states $\vec{S}^* = \underset{\vec{S}}{\operatorname{argmax}} P(\vec{S} | \vec{O})$

3) How to update beliefs $P(s_t = i | o_1, o_2, \dots, o_t)$

4) How to estimate $\{\pi_i, a_{ij}, b_{ik}\}$ from data

} inference
[learning

1) Computing likelihood

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

\swarrow sum over n^T sequences of hidden states

$$= \sum_{\vec{s}} P(s_1) \prod_t P(s_t | s_{t-1}) \prod_t P(o_t | s_t)$$

* Efficient recursion — if know state at time t , can calc $t+1$
(don't recurse over all $t \dots$)

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

magnification

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

$$= \sum_{i=1}^n \underbrace{P(o_1, o_2, \dots, o_t, s_t=i)}_{\text{recursive instance}} \cdot \underbrace{P(s_{t+1}=j | s_t=i)}_{a_{ij}} \cdot \underbrace{P(o_{t+1} | s_{t+1}=j)}_{b_j(o_{t+1})}$$

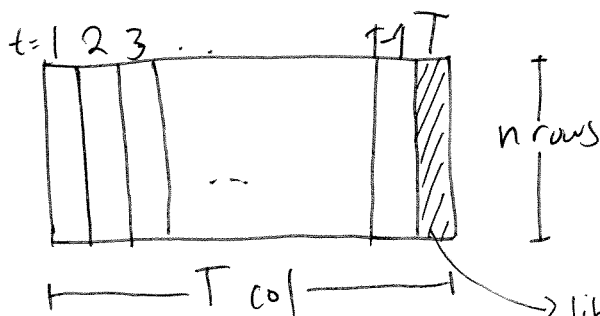
cond. indep

CPTs of HMM

* shorthand notation

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

$i=1 \dots n$ # hidden states
 $t=1 \dots T$ seq. length



likelihood is sum of last column
 $P(o_1, o_2, \dots, o_T)$

* Forward algorithm

- base case ($t=1$) first column of matrix

$$\begin{aligned}\alpha_{i1} &= P(o_1, s_1=i) && \text{by def} \\ &= P(s_1=i) \cdot P(o_1 | s_1=i) && \text{product rule} \\ &= \underbrace{\pi_i}_{\text{initial state}} \cdot \underbrace{b_i(o_1)}_{\text{emission matrix}}\end{aligned}$$

- recursive step from time $t-1$ to time t

$$\alpha_{j,t+1} = \sum_{i=1}^n \alpha_{i,t} \cdot a_{ij} \cdot b_j(o_{t+1})$$

* Back to likelihood computation

$$\begin{aligned}P(o_1, o_2, \dots, o_T) &= \sum_{i=1}^n P(o_1, o_2, \dots, o_T, s_T=i) && \text{marginalization} \\ &= \sum_{i=1}^n \alpha_{iT} && \begin{array}{l} \text{sum of elements} \\ \text{(last column of } \alpha \text{ matrix)} \end{array}\end{aligned}$$

* Scales as $O(Tn^2)$

- linear, not exponential, in sequence length T
- quadratic in # of states n

* Warning: naive calculation will underflow for long sequences $T \gg 1$

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

(fix by rescaling each iteration, eventually taking log)

2) How to compute most likely state sequence?

$$S^* = \{s_1^*, s_2^*, \dots, s_{T-1}^*, s_T^*\}$$

$$= \operatorname{argmax}_{\vec{s}} P(s_1, s_2, \dots, s_T | o_1, o_2, \dots, o_T)$$

$$= \operatorname{argmax}_{\vec{s}} \left[\frac{P(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T)}{P(o_1, o_2, \dots, o_T)} \right] \quad \text{product rule}$$

constant w.r.t \vec{s} , can drop for argmax

$$S^* = \operatorname{argmax}_{\vec{s}} P(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T)$$

How to compute s^* ?

Define: $l_{it}^* = \max_{\substack{s_1:t-1 \\ \{s_1, s_2, \dots, s_{t-1}\}}} \log P(s_1, s_2, \dots, s_{t-1}, \underline{s_t = i}, o_1, o_2, \dots, o_t)$

- log-prob. of most likely sequence of t hidden states that ends in state i at time t that explains observations o_1, \dots, o_t

* Recursive algorithm

- base case ($t=1$) leftmost column

$$l_{i1}^* = \log P(s_1 = i, o_1)$$

$$= \log [P(s_1 = i) \cdot P(o_1 | s_1 = i)] \quad \text{product rule}$$

$$= \log [\pi_i \cdot b_i(o_1)]$$

$$= \log \pi_i + \log b_i(o_1)$$

• recursive step from $t-1$ to t

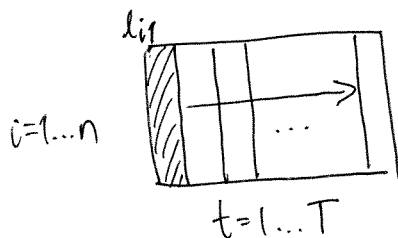
$$\begin{aligned}
 l_{j,t+1}^* &= \max_{s_1, s_2, \dots, s_t} \log P(s_1, s_2, \dots, s_t, s_{t+1}=j, o_1, o_2, \dots, o_t, o_{t+1}) \\
 &= \max_{s_1, s_2, \dots, s_{t-1}} \max_i \log P(s_1, s_2, \dots, s_{t-1}, s_t=i, s_{t+1}=j, o_1, o_2, \dots, o_t, o_{t+1}) \\
 &= \max_{s_1, s_2, \dots, s_{t-1}} \max_i \log \left[P(s_1, s_2, \dots, s_{t-1}, s_t=i, o_1, o_2, \dots, o_t) \cdot \right. \\
 &\quad \left. P(s_{t+1}=j | s_1, \dots, s_t=i, o_1, \dots, o_t) \cdot P(o_{t+1} | s_1, \dots, s_{t-1}, s_t=i, s_{t+1}=j, o_1, \dots, o_t) \right] \\
 &= \max_{s_1, s_2, \dots, s_{t-1}} \max_i \log P(s_1, \dots, s_{t-1}, s_t=i, o_1, \dots, o_t) \cdot P(s_{t+1}=j | s_t=i) \cdot P(o_{t+1} | s_{t+1}=j) \\
 &= \max_{s_1, \dots, s_{t-1}} \left[\max_i \left\{ \log P(s_1, \dots, s_{t-1}, s_t=i, o_1, \dots, o_t) + \log P(s_{t+1}=j | s_t=i) \right\} \right. \\
 &\quad \left. + \log P(o_{t+1} | s_{t+1}=j) \right]
 \end{aligned}$$

Reassembling:

$$\begin{aligned}
 l_{j,t+1}^* &= \max_i \left[\max_{s_1, \dots, s_{t-1}} \log P(s_1, \dots, s_{t-1}, s_t=i, o_1, \dots, o_t) + \log P(s_{t+1}=j | s_t=i) \right] \\
 &\quad + \log P(o_{t+1} | s_{t+1}=j)
 \end{aligned}$$

$$\boxed{l_{j,t+1}^* = \max_i \left[l_{it}^* + \log a_{ij} \right] + \log b_j(o_{t+1})} \quad \text{forward pass}$$

algorithm: fill in l_{it}^* column by column



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

Record most likely transitions:

$$\Phi_{t+1}(j) = \operatorname{argmax}_i \left[l_{it}^* + \log a_{ij} \right]$$

(if at state j at time $t+1$, what state i at time t got you there~~s~~, for observations o_1, o_2, \dots, o_{t+1})

* Compute s^* from backtracking:

• At time T : $s_T^* = \operatorname{argmax}_i [l_{iT}^*]$ (what is largest element of last column in l^*)

• For time $t=T-1$ to $t=1$:

$$s_t^* = \Phi_{t+1}(s_{t+1}^*)$$

backward pass

* Jargon:

- s^* is known as "Viterbi" state sequence or "path"
- Viterbi algorithm is example of dynamic programming