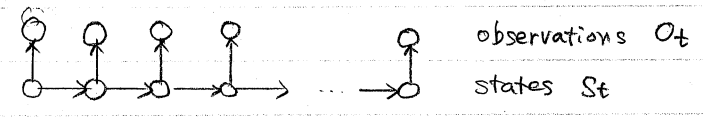


1/4 Review

* Hidden Markov models (HMMs)



* Parameters

$$\pi_i = P(S_1 = i)$$

$$a_{ij} = P(S_{t+1} = j \mid S_t = i)$$

$$b_{ik} = P(O_t = k \mid S_t = i)$$

* Key computations

- 1) how to compute likelihood $P(O_1, O_2, \dots, O_T)$?
- 2) how to compute $\arg \max_{s_1, \dots, s_T} P(s_1, \dots, s_T \mid O_1, \dots, O_T)$?
- 3) how to estimate (learn) $\{\pi_i, a_{ij}, b_{ik}\}$?

* HW problem: belief updating

- recursion for $P(S_t \mid O_1, O_2, \dots, O_t)$
- important for real-time monitoring.

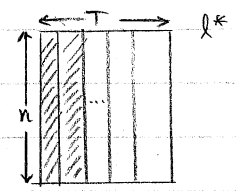
(2) Computing most likely state sequence

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

Recursion: $l_{j,t+1}^* = \max_i [l_{it}^* + \log a_{ij}] + \log b_j(O_{t+1})$

How to derive s^* from l^* ?

discrete state sequence \leftarrow real-valued $n \times l$ matrix



* Record most likely state transitions:

$$\Phi_{t+1}(j) = \arg \max_i [l_{it}^* + \log a_{ij}]$$

What is most likely state at time t given state j at time $t+1$ with observations O_1, O_2, \dots, O_t .

* Compute s^* by backtracking:

$$s_T^* = \arg \max_i [l_{iT}^*]$$

$$s_t^* = \Phi_{t+1}(s_{t+1}^*) \quad t = T-1, T-2, \dots, 1$$

$s^* = \{s_1^*, s_2^*, \dots, s_T^*\}$ is known as Viterbi path.

Viterbi algorithm is instance of dynamic programming.

(3) Learning in HMMs.

Given: sequence of observations $\{O_1, O_2, \dots, O_T\}$ (just one for simplicity)

Goal: estimate $\{\pi_i, a_{ij}, b_{ik}\}$ to maximize $P(O_1, O_2, \dots, O_T)$

Fixed: number of hidden states n .

* Shared CPTs to estimate: $\pi_i = P(S_1 = i)$

$$a_{ij} = P(S_{t+1} = j | S_t = i)$$

$$b_{ik} = P(O_t = k | S_t = i)$$

* E-step

Compute $P(S_t = i | O_1, \dots, O_T)$

$$P(S_t = i, S_{t+1} = j | O_1, \dots, O_T)$$

$$P(S_t = i, O_t = k | O_1, \dots, O_T) = P(S_t = i | O_1, \dots, O_T) I(O_t, k)$$

(analogy: p_a, x, v)

* Compute posterior probabilities

Analogous to $\alpha_{it} = P(O_1, O_2, \dots, O_t, S_t = i)$ before and up to time t

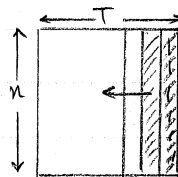
define $\beta_{it} = P(O_{t+1}, O_{t+2}, \dots, O_T | S_t = i)$ after time t

Starting at $t+1$ Conditioning on $S_t = i$

- Recursion for β_{it} :

(1) base case $\beta_{iT} = P(\dots | S_T = i) ?$

$\beta_{iT} = 1$ for all i .



(2) backwards step

$$\beta_{it} = P(O_{t+1}, \dots, O_T | S_t = i)$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T, S_{t+1} = j | S_t = i) \text{ marginalization.}$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T | S_{t+1} = j, S_t = i) P(S_{t+1} = j | S_t = i) \text{ product rule}$$

conditional independence

$$= \sum_{j=1}^N P(O_{t+1} | S_{t+1} = j) P(O_{t+2}, \dots, O_T | S_{t+1} = j, O_{t+1}) P(S_{t+1} = j | S_t = i) \text{ product rule}$$

$$= \sum_{j=1}^N b_j(O_{t+1}) \beta_{j,t+1} a_{ij}$$

Computing α & β matrix \leftrightarrow "forward-backward" algorithm (a.k.a Baum-Welch)

- Posterior probabilities in E-step:

$$P(S_t = i, S_{t+1} = j | O_1, \dots, O_T) = \frac{P(S_t = i, S_{t+1} = j, O_1, \dots, O_T)}{P(O_1, \dots, O_T)}$$

$$= \frac{P(O_1, \dots, O_t, S_t = i) P(S_{t+1} = j | S_t = i, O_1, \dots, O_t) P(O_{t+1} | S_{t+1} = j, S_t = i, O_1, \dots, O_t) \times P(O_{t+2}, \dots, O_T | S_{t+1} = j, S_t = i, O_1, \dots, O_t)}{P(O_1, O_2, \dots, O_T)}$$

$$= \frac{\alpha_{it} a_{ij} b_j(O_{t+1}) \beta_{j,t+1}}{\sum_i \alpha_{iT}}$$

Other posteriors in E-step

$$P(S_t = i | O_1, \dots, O_T) = \sum_{j=1}^N P(S_t = i, S_{t+1} = j | O_1, \dots, O_T)$$

$P(S_1 = i | O_1, \dots, O_T)$ special case with $t=1$.

* M-step updates

$$\pi_i \leftarrow P(S_1 = i \mid O_1, \dots, O_T)$$

$$a_{ij} \leftarrow \frac{\sum_{t=1}^{T-1} P(S_t = i, S_{t+1} = j \mid O_1, \dots, O_T)}{\sum_{t=1}^{T-1} P(S_t = i \mid O_1, \dots, O_T)}$$

$$b_{ik} \leftarrow \frac{\sum_{t=1}^T P(S_t = i \mid O_1, \dots, O_T) I(k, O_t)}{\sum_{t=1}^T P(S_t = i \mid O_1, \dots, O_T)}$$

* Complexity of HMM computations

(i) to compute $P(O_1, O_2, \dots, O_T)$

(ii) to decode $S^* = \arg \max_{\mathcal{S}} P(\mathcal{S} \mid \mathcal{O})$

(iii) parameter update of EM

$O(n^2 T)$ $T = \text{sequence length}$
 $n = \# \text{ states.}$

Multivariate Gaussian distributions

* Random variable $\vec{x} \in \mathbb{R}^n$ (real-valued vector)

* Probability density function (PDF)

$$P(\vec{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\}$$



* Parameters

$$\vec{\mu} = E[\vec{x}] = \int P(\vec{x}) \vec{x} d^n \vec{x}$$

$$\Sigma_{ij} = E[(\vec{x} - \vec{\mu})_i (\vec{x} - \vec{\mu})_j] = \int P(\vec{x}) (x_i - \mu_i) (x_j - \mu_j) d^n \vec{x}$$

Shorthand: $P(\vec{x}) = N(\vec{x}; \vec{\mu}, \Sigma)$

* Mathematical properties:

(i) if $P(\vec{x})$ and $P(\vec{y})$ are gaussian PDFs over $\vec{x}, \vec{y} \in \mathbb{R}^n$,

then $P(\alpha \vec{x} + \beta \vec{y})$ is also gaussian.

α, β are linear scalar coefficients.

(ii) if $P(\vec{x})$ is gaussian over $\vec{x} \in (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$,

then all marginals $\{P(x_i), P(x_i, x_j), \dots\}$ and all conditionals $\{P(x_i | x_j), P(x_i, x_j | x_k), \dots\}$ are also gaussian.

* Maximum likelihood estimation.

Given i.i.d data $\{\vec{x}_t\}_{t=1}^T$ where $\vec{x}_t \in \mathbb{R}^n$, how to choose $\vec{\mu}, \Sigma$

to maximize $P(\text{data}) = \prod_{t=1}^T N(\vec{x}_t; \vec{\mu}, \Sigma)$

- log-likelihood

$$\mathcal{L} = \log P(\text{data}) = \sum_{t=1}^T \log N(\vec{x}_t; \vec{\mu}, \Sigma)$$

To maximize: $\frac{\partial R}{\partial \bar{\mu}} = 0$, $\frac{\partial R}{\partial \Sigma_{ij}} = 0$.

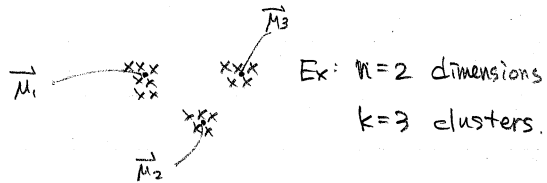
$\bar{\mu} = \frac{1}{T} \sum_{t=1}^T \bar{x}_t$ sample mean

$\Sigma_{ij} = \frac{1}{T} \sum_{t=1}^T (\bar{x}_t - \bar{\mu})_i (\bar{x}_t - \bar{\mu})_j$ sample covariance.

Clustering:

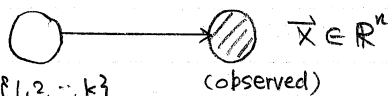
* Inputs $\{ \bar{x}_1, \bar{x}_2, \dots, \bar{x}_T \}$ with $\bar{x}_t \in \mathbb{R}^n$

* Goal: partition inputs into k clusters.



Gaussian mixture model

* DAG



$Z \in \{1, 2, \dots, k\}$
(hidden)
cluster label

* CPTs

$P(Z=i)$ fraction of data in cluster i

$P(\bar{x} | Z=i) = N(\bar{x}; \bar{\mu}_i, \Sigma_i)$ cluster-dependent means and covariance matrices.

* Aside: ML estimation for complete data $\prod_{t=1}^T P(\bar{x}_t, z_t)$

Let $T_i = \sum_{t=1}^T I(z_t, i)$ count of label i .

$P(z_i) = \frac{T_i}{T}$

$\bar{\mu}_i = \frac{1}{T_i} \sum_{t=1}^T \bar{x}_t I(z_t, i)$

$\Sigma_{\alpha\beta}^i = \frac{1}{T_i} \sum_{t=1}^T (\bar{x}_t - \bar{\mu}_i)_\alpha (\bar{x}_t - \bar{\mu}_i)_\beta I(z_t, i)$