

M step for HMM

We assume the E step computed distributions for

The degree each state explains each data point (analogous to GMM responsibilities). $\gamma(z_n) = p(z_n | X, \theta^{(s)})$

The degree that the combination of a state, and a previous one explain two data points.

$$\xi(z_{n-1}, z_n) = p(z_{n-1}, z_n | X, \theta^{(s)})$$

← “xi”

EM for HMM (sketch)

$$\log(p(X, Z | \theta)) =$$

$$\sum_{k=1}^K \gamma_{1k} \log(\pi) + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \gamma_{n-1,j} \gamma_{n,k} \log(p(z_n | z_{n-1}, A)) + \sum_{n=1}^N \sum_k \gamma_{nk} \log(p(x_n | z_n, \phi))$$

We define

$$\gamma(z_n) = p(z_n | X, \theta^{(s)})$$

$$\xi(z_{n-1}, z_n) = p(z_{n-1}, z_n | X, \theta^{(s)})$$

← “xi”

EM for HMM (sketch)

$$\log(p(X, Z | \theta)) =$$

$$\sum_{k=1}^K \gamma_{1k} \log(\pi) + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \gamma_{n-1,j} \gamma_{n,k} \log(p(z_n | z_{n-1}, A)) + \sum_{n=1}^N \sum_k \gamma_{nk} \log(p(x_n | z_n, \phi))$$

By analogy with the GMM

$$\begin{aligned} Q(\theta^{(s+1)}, \theta^{(s)}) &= \sum_z p(Z | \theta^{(s)}) \log(p(X, Z | \theta^{(s+1)})) \\ &= \sum_{k=1}^K \gamma_{1k} \log(\pi) + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{n,k}) \log(p(z_n | z_{n-1}, A)) \\ &\quad + \sum_{n=1}^N \sum_k \gamma_{nk} \log(p(x_n | z_n, \phi)) \end{aligned}$$

EM for HMM (sketch)

Doing the maximization using Lagrange multipliers gives us

$$\pi_k = \frac{\gamma_{1k}}{\sum_{k'} \gamma_{1k'}}$$

Much like the GMM. Taking the partial derivative for π_k kills second and third terms.

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{n,k})}{\sum_{k'} \sum_{n=2}^N \xi(z_{n-1,j}, z_{n,k'})}$$

EM for HMM (sketch)

The maximization of $p(x_n | \phi)$ is exactly the same as the mixture model.

For example, if we have Gaussian emissions, then

$$\mu_k = \frac{\sum_n x_n \gamma(z_{nk})}{\sum_n \gamma(z_{nk})}$$

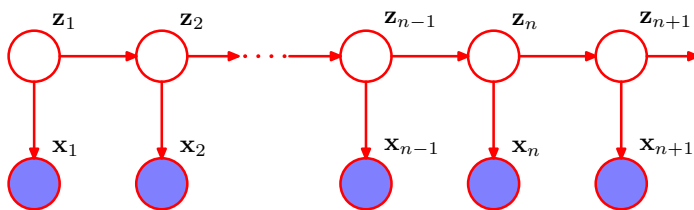
E step for EM for HMM

Computing the E step is a bit more involved.

Recall that in the mixture case it was easy because we only needed to consider the relative likelihood that each cluster independently explain the observations.

However, here the sequence also must play a role.

Graphical model for the E step



Note that our task is to compute marginal probabilities

Computing marginals in an HMM

Various names, flavors, notations, ...

Forward-Backward algorithm

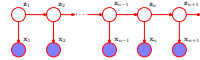
Alpha-beta algorithm

Sum-product for HMM

(Bishop also says “Baum Welch” but that is a synonym for the EM algorithm as whole).

Alpha-beta algorithm

$$\begin{aligned}
 \gamma(z_n) &= p(z_n | X) \\
 &= \frac{p(X | z_n) p(z_n)}{p(X)} \\
 &= \frac{p(x_1, \dots, x_n | z_n) p(x_{n+1}, \dots, x_N | z_n) p(z_n)}{p(X)} \\
 &= \frac{p(x_1, \dots, x_n, z_n) p(x_{n+1}, \dots, x_N | z_n)}{p(X)} \\
 &= \frac{\alpha(z_n) \beta(z_n)}{p(X)}
 \end{aligned}$$



Where we define

$$\alpha(z_n) = p(x_1, \dots, x_n, z_n)$$

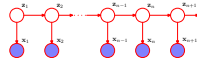
$$\beta(z_n) = p(x_{n+1}, \dots, x_N | z_n)$$

Expressing alpha recursively



$$\begin{aligned}
 \alpha(z_n) &= p(x_1, \dots, x_n, z_n) \\
 &= p(x_1, \dots, x_n | z_n) p(z_n) && \text{(definition of "I")} \\
 &= p(x_n | z_n) p(x_1, \dots, x_{n-1} | z_n) p(z_n) && \text{(conditional independence)} \\
 &= p(x_n | z_n) p(x_1, \dots, x_{n-1}, z_n) && \text{(definition of "I")} \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_{n-1}, z_n) && \text{(marginal)} \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1} | z_{n-1}) p(z_{n-1}) && \text{(definition of "I")} \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1} | z_{n-1}) p(z_n | z_{n-1}) p(z_{n-1}) && \text{(conditional independence)} \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_{n-1}) p(z_n | z_{n-1}) && \text{(definition of "I")} \\
 &= p(x_n | z_n) \sum_{z_{n-1}} \alpha(z_{n-1}) p(z_n | z_{n-1}) && \text{(definition of } \alpha(z_n) \text{)}
 \end{aligned}$$

Expressing alpha recursively



$$\alpha(z_n) = p(x_n | z_n) \sum_{z_{n-1}} \alpha(z_{n-1}) p(z_n | z_{n-1})$$

This is a recursive evaluation of alpha. So we can compute all of them easily if we know the first one, $\alpha(z_1)$.

$$\begin{aligned}
 \alpha(z_1) &= p(x_1, z_1) && \text{(we defined } \alpha(z_n) = p(x_1, \dots, x_n, z_n) \text{)} \\
 &= p(z_1) p(x_1 | z_1) && \text{(this is a K dimensional vector for fixed } x_1 \text{)}
 \end{aligned}$$

$$\alpha(z_1)_k = \pi_k p(x_1 | \phi_k)$$

Alpha-beta algorithm

Similarly, we can derive a recurrence relation for beta

$$\beta(z_n) = \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n)$$

Alpha-beta algorithm

The details for $\beta(z_n) = \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n)$

$$\begin{aligned}
 \beta(z_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_n) \\
 &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, z_{n+1} | z_n) \\
 &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_n, z_{n+1}) p(z_{n+1} | z_n) \\
 &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_{n+1}) p(z_{n+1} | z_n) \\
 &= \sum_{z_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | z_{n+1}) p(\mathbf{x}_{n+1} | z_{n+1}) p(z_{n+1} | z_n).
 \end{aligned}$$

Alpha-beta algorithm

Our recurrence relation for beta

$$\beta(z_n) = \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n)$$

We can compute the betas if we know the last one.

$$\begin{aligned}
 p(z_N | X) &= \frac{\alpha(z_N) \beta(z_N)}{p(X)} \\
 &= \frac{p(X, z_N) \beta(z_N)}{p(X)} \quad (\text{we defined } \alpha(z_n) = p(x_1, \dots, x_n, z_n)) \\
 &= p(z_N | X) \beta(z_N)
 \end{aligned}$$

So $\beta(z_N) = 1$

Alpha-beta algorithm

Given the alphas and betas, we can compute all the quantities we need for the E step.

$$\gamma(z_n) = \frac{\alpha(z_n) \beta(z_n)}{p(X)} \quad (\text{our definition})$$

We know that $\sum_{z_n} \gamma(z_n) = 1$

$$\text{so } \sum_{z_n} \frac{\alpha(z_n) \beta(z_n)}{p(X)} = 1$$

$$\text{and } p(X) = \sum_{z_n} \alpha(z_n) \beta(z_n)$$

We do not need $p(X)$ for EM, but it is the likelihood which we want to monitor ($p(X) = p(X | \theta^{(s)})$).

Alpha-beta algorithm

Given the alphas and betas, we can compute all the quantities we need for the E step.

$$\begin{aligned}
 \xi(z_{n-1}, z_n) &= p(z_{n-1}, z_n | \mathbf{X}) \\
 &= \frac{p(\mathbf{X} | z_{n-1}, z_n) p(z_{n-1}, z_n)}{p(\mathbf{X})} \\
 &= \frac{p(x_1, \dots, x_{n-1} | z_{n-1}) p(x_n | z_n) p(x_{n+1}, \dots, x_N | z_n) p(z_n | z_{n-1}) p(z_{n-1})}{p(\mathbf{X})} \\
 &= \frac{\alpha(z_{n-1}) p(x_n | z_n) p(z_n | z_{n-1}) \beta(z_n)}{p(\mathbf{X})} \quad (13.43) \\
 &\quad (\text{in Bishop})
 \end{aligned}$$