initial parameters are then computed as

$$\pi_m = N_m / N_{\text{tot}}$$

$$\mu_m = N_m^{-1} \sum_{n=1}^{N} \sum_{t=1}^{T^{(n)}} o_t^{(n)} h_m(o_t^{(n)}) \qquad (15)$$

$$\sigma_m^2 = N_m^{-1} \sum_{n=1}^{N} \sum_{t=1}^{T^{(n)}} (o_t^{(n)} - \mu_m)^2 h_m(o_t^{(n)}) \qquad (16)$$

*(handwritten: $t=0$?)*

This approximation is then improved ~~upon~~ by utilizing the expectation-maximization procedure described by Bilmes [22].

$$\pi_m' = N_{\text{tot}}^{-1} \sum_{n=1}^{N} \sum_{t=0}^{T^{(n)}} \chi_m(o_t^{(n)}, \mathbf{E}, \boldsymbol{\pi})$$

$$\mu_m' = (\pi_m N_{\text{tot}})^{-1} \sum_{n=1}^{N} \sum_{t=0}^{T^{(n)}} o_t^{(n)} \chi_m(o_t^{(n)}, \mathbf{E}, \boldsymbol{\pi})$$

$$\sigma_m'^2 = (\pi_m N_{\text{tot}})^{-1} \sum_{n=1}^{N} \sum_{t=0}^{T^{(n)}} (o_t^{(n)} - \mu_m')^2 \chi_m(o_t^{(n)}, \mathbf{E}, \boldsymbol{\pi}) \quad (17)$$

*(handwritten left margin: Is this right? Should these be $\chi_m$ instead of $\pi_m$?)*

*(handwritten: $P_0'$?)*

where the function $p(m|o, \mathbf{E}, \boldsymbol{\pi})$ is given by the fuzzy membership function:

$$\chi_m(o, \mathbf{E}, \boldsymbol{\pi}) = \frac{\pi_m \varphi(o|e_m)}{\sum_{l=1}^{M} \pi_l \varphi(o|e_l)} \qquad (18)$$

This iterative procedure is terminated when the change in the parameters $\{\boldsymbol{\pi}, \mu, \sigma^2\}$ falls below a certain relative threshold, such as $\|\boldsymbol{\pi}^{(n)} - \boldsymbol{\pi}^{(n-1)}\| / \|\boldsymbol{\pi}^{(n)}\| < 10^{-4}$.

*(handwritten: index (i) so not to confuse w/ data trace?)*

### 2. *Transition matrix estimation*

Once initial state observable emission parameters $\mathbf{E}$ are determined, an initial transition matrix is estimated using an iterative likelihood maximization approach that enforces detailed balance [23]. First, a matrix of fractional transition counts $\mathbf{C} \equiv (c_{ij})$ is estimated using the membership function:

*(handwritten: bold)*

$$c_{ij} = \sum_{n=1}^{N} \sum_{t=1}^{T^{(n)}} \chi_i(o_{t-1}^{(n)}, \mathbf{E}, \boldsymbol{\pi}) \chi_j(o_t^{(n)}, \mathbf{E}, \boldsymbol{\pi}) \qquad (19)$$

A symmetric $M \times M$ matrix $\mathbf{X} \equiv (x_{ij})$ is initialized by

$$x_{ij} = x_{ji} = c_{ij} + c_{ji} \qquad (20)$$

and a vector of row sums

$$x_i = \sum_{j=1}^{M} x_{ij} \qquad (21)$$

Then, the iterative procedure described in Algorithm 1 of [23] is applied. For each update iteration, we first update the diagonal elements of $\mathbf{X}$:

$$x_{ii}' = \frac{c_{ii}(x_i - x_{ii})}{c_i - c_{ii}}$$

$$\cancel{x_i' = \sum_{j=1}^{M} x_{ij}'} \qquad (22)$$

followed by the off-diagonal elements:

$$x_{ij}' = x_{ji}' = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

$$\cancel{x_i = \sum_{j=1}^{M} x_{ij}'} \qquad (23)$$

where the quantities $a$, $b$, and $c$ are computed from $\mathbf{X}$ and $\mathbf{C}$ as

$$a \equiv c_i - c_{ij} + c_j - c_{ji}$$
$$b \equiv c_i(x_j - x_{ji}) + c_j(x_i - x_{ij}) - (c_{ij} + c_{ji})(x_i - x_{ij} + x_j - x_{ji})$$
$$c \equiv -(c_{ij} + c_{ji})(x_i - x_{ij})(x_j - x_{ji}) \qquad (24)$$

[JDC: Can we merge the updates for $x_{ii}'$ and $x_{ij}'$ to simplify this description?] Once a sufficient number of iterations have been completed to compute a stable estimate of $\mathbf{X}$ (such as the relative convergence criteria $\|\mathbf{X}^{(n)} - \mathbf{X}^{(n-1)}\| / \|\mathbf{X}^{(n)}\| < 10^{-4}$, the maximum likelihood transition matrix estimate $\mathbf{T}$ is computed as

$$T_{ij} = \frac{x_{ij}}{x_i}. \qquad (25)$$

Note that the equilibrium probability vector $\boldsymbol{\pi}$ computed during the Gaussian mixture model fitting is not respected during this step.

### B. Fitting a maximum likelihood HMM

The HMM model parameters $\boldsymbol{\Theta} \equiv \{\mathbf{T}, \mathbf{E}\}$ are fit to the observed data $\mathbf{O}$ through use of the expectation-maximization (EM) algorithm [24]. This is an iterative procedure, where the model parameters are subsequently refined through successive iterations.

During each iteration, the Baum-Welch algorithm [12] is used to compute, for each trace $n$, $\boldsymbol{\Xi}^{(n)} \equiv (\xi_{tij}^{(n)})$, which represents the probability that the system transitions from hidden state $i$ at time $t-1$ to hidden state $j$ at time $t$, and $\gamma_{ti}^{(n)}$, the probability that the system ~~was~~ *is* in state $i$ at time $t$. This is accomplished by first executing the *forward algorithm*, which proceeds (suppressing the superscript $(n)$) as

*(handwritten: $\rho_j$?)*

$$\alpha_{tj} = \begin{cases} \rho_i \varphi(o_0|e_j) & t = 0 \\ \varphi(o_t|e_j) \sum_{i=1}^{M} \alpha_{(t-1)i} T_{ij} & t = 1, \ldots, T_n \end{cases} \qquad (26)$$

*(handwritten: (n)? or suppress?)*

followed by the *backward algorithm*,

$$\beta_{ti} = \begin{cases} 1 & t = T_n \\ \sum_{j=1}^{M} T_{ij} \varphi(o_{t+1}|e_j) \beta_{(t+1)j} & t = T_n - 1, \ldots, 0 \end{cases} \qquad (27)$$

The $M \times M \times T_n$ matrix $\boldsymbol{\Xi}$ is then computed for $t = 0, \ldots, (T_n - 1)$ as

*(handwritten: $T_n \times M \times M$; or $T_n - 1$?)*

$$\xi_{tij} = \alpha_{ti} \varphi(o_{t+1}|e_i) T_{ij} \beta_{(t+1)j} / \sum_{i=1}^{M} \alpha_{Ti} \qquad (28)$$

*(handwritten: $e_j$?)*

$$\gamma_{ti} = \sum_{j=1}^{M} \xi_{tij} \qquad (29)$$

In practice, the logarithms of these quantities are computed instead to avoid numerical underflow. [JDC: Not sure what happened to $\gamma_{ti}$ for $t = T_n$. Check this?]

*(handwritten bottom: like notation better! could we use the notation $x_i.$ and $c_i.$ to indicate $\sum_j x_{ij}$? You shouldn't need to update them then.)*