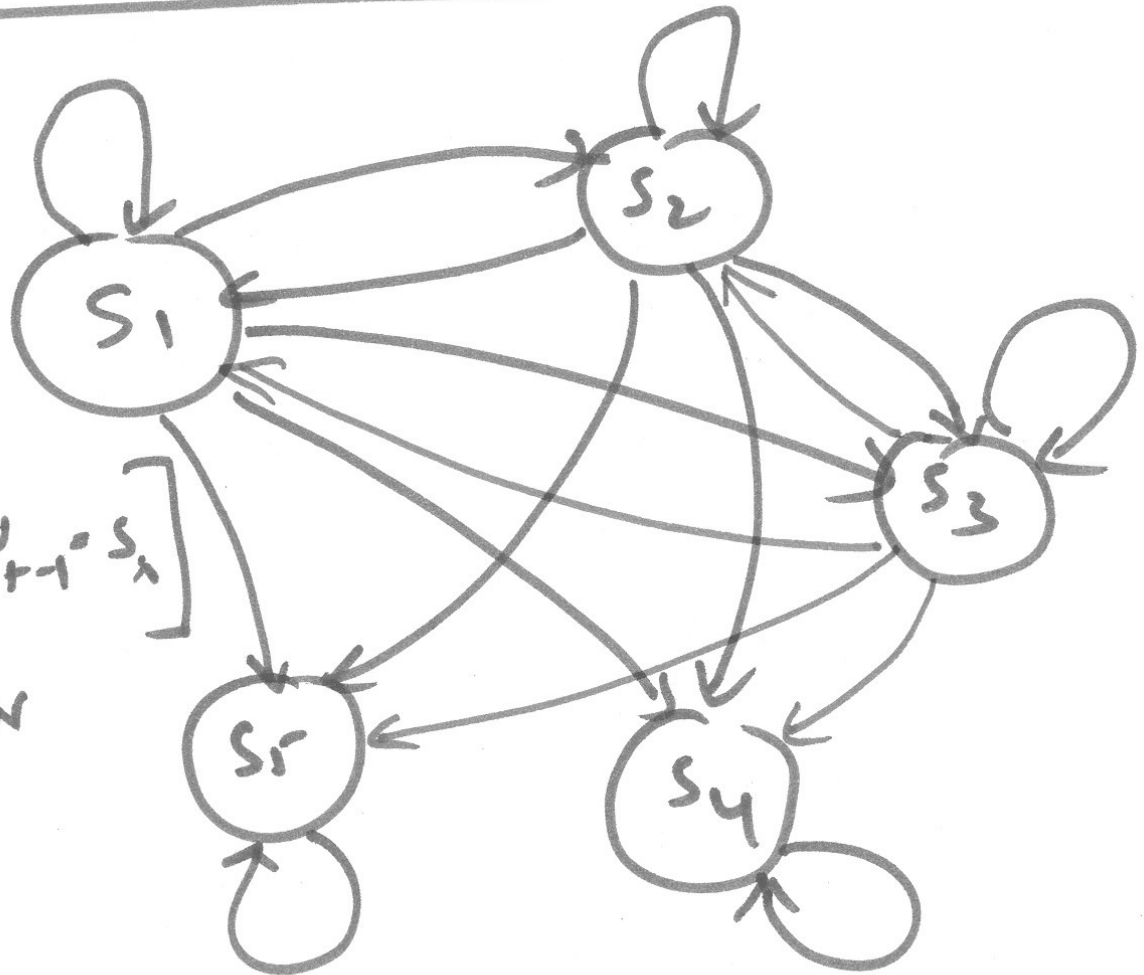


Hidden Markov Models (HMM)

Discrete Markov Process

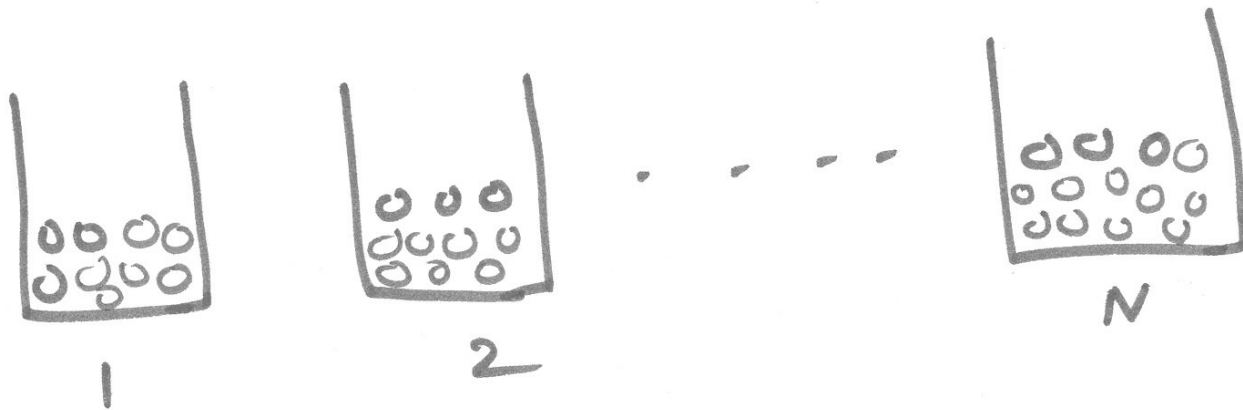


$$a_{ij} = P[u_t = s_j | u_{t-1} = s_i]$$

$$i \geq 1; j \leq N$$

$$\sum_{j=1}^N a_{ij} = 1$$

Hidden Markov Model



$$O = (b, G, R, R, b, G, b, b, R, \dots)$$

\Rightarrow # boxes, Sequence of boxes

Symbol prob

Spectral vector \Leftrightarrow Sequence of VT shapes

Basic Elements of HMM

1. # states — N
2. # distinct symbols / state — M
$$V = \{v_1, v_2, \dots, v_M\}$$
3. State transition probability $A = \{a_{ij}\}_{N \times N}$
$$a_{ij} = P[q_{t+1} = j \mid q_t = i]; \quad \begin{matrix} i \geq 1 \\ j \leq N \end{matrix}$$
4. Observation symbol prob $B = \{b_j(k)\}$
$$b_j(k) = P[v_t = v_k \mid q_t = j]$$
5. Initial prob $\pi_i = \pi$; $\pi_{\hat{i}} = P[q_1 = \hat{i}]$
$$\text{MODEL} = \lambda = [A, B, \pi]$$

Basic Problems in HMM

① How do we efficiently compute $P(O/\lambda)$

$$P(o_1 o_2 \dots o_T / \lambda) \Rightarrow \begin{array}{l} \text{Testing, evaluation} \\ \text{Recognition, calculation} \end{array}$$

$$\left. \begin{array}{l} O = o_1 o_2 \dots o_T \\ \lambda = (A, B, \pi) \end{array} \right\} \Rightarrow Q = v_1 v_2 \dots v_T \quad (\text{Seq of states})$$

best seq that explains the observations in better way

optimal seq for the given observation

states & # symbols can be modified using ~~step~~ step-2

③ How do we adjust $\lambda = (A, B, \pi)$ to max $P(O/\lambda)$

Training step, building models

Solution to Problem-1 (computation of $P(u/\lambda)$)

$$O = O_1, O_2, \dots, O_T, \quad \lambda = A, B, \pi$$

$$\text{Possible state } u = Q = (u_1, u_2, \dots, u_T) \Rightarrow N^T$$

$$P(O/\lambda) = \sum_{u_1, u_2, \dots, u_T} \cancel{P(O/u, \lambda)} P(O, u/\lambda)$$

$$P(O, u/\lambda) = P(O/u, \lambda) P(u/\lambda)$$

$$P(O/u, \lambda) = \prod_{t=1}^T P(O_t/u_t, \lambda)$$

$$= b_{u_1}(O_1) b_{u_2}(O_2) \dots b_{u_T}(O_T)$$

$$P(u/\lambda) = \prod_{v_1} a_{v_1, v_2} a_{v_2, v_3} \dots a_{v_{T-1}, v_T}$$

$$P(o, v/\lambda) = P(o/v, \lambda) P(v/\lambda)$$

$$= \prod_{v_1} b_{v_1}(o_1) a_{v_1, v_2} b_{v_2}(o_2) \dots a_{v_{T-1}, v_T} b_{v_T}(o_T)$$

$\rightarrow (2T-1)$ multiplications

$$P(o/\lambda) = \sum P(o, v/\lambda)$$

$v_1, v_2, \dots, v_T \Rightarrow N_T$ state sequences

Total no. of computations $\begin{cases} 2T N^T - \text{multiplications} \\ N^T - 1 - \text{additions} \end{cases}$

Ex

$$N = 5, \quad T = 100$$

$$\# \text{ computations} = 2 \times 100 \times 5^{100} \Rightarrow \underline{\underline{10^{72}}}$$

Forward Procedure $\{\alpha_T(\bar{i})\}$

$\alpha_T(\bar{i}) \rightarrow$ Prob of partial observation sequence

$$\alpha_T(\bar{i}) = P(o_1 o_2 \dots o_T, q_T = \bar{i} \mid \lambda)$$

1. Initialization $\alpha_1(\bar{i}) = P(o_1, q_1 = \bar{i} \mid \lambda)$
 $= \pi_{\bar{i}} b_{\bar{i}}(o_1) ; 1 \leq \bar{i} \leq N$

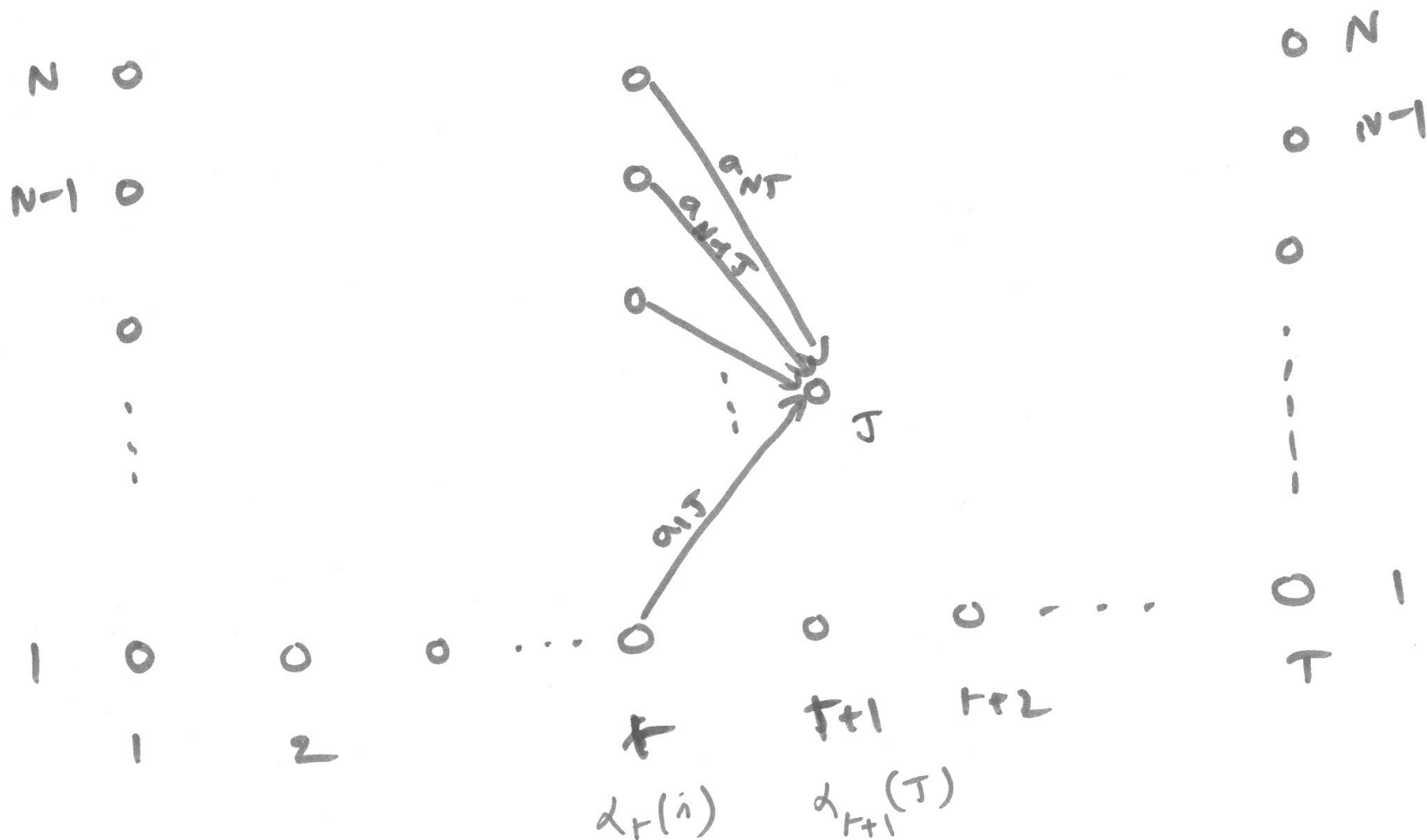
2. Induction

$$\alpha_{t+1}(\bar{j}) = \left[\sum_{\bar{i}=1}^N \alpha_t(\bar{i}) a_{\bar{i}\bar{j}} \right] b_{\bar{j}}(o_{t+1}) ; \begin{matrix} 1 \leq t \leq T-1 \\ 1 \leq \bar{j} \leq N \end{matrix}$$

3. Termination

$$P(o/\lambda) = \sum_{\bar{i}=1}^N \alpha_T(\bar{i})$$

Illustration of computation of forward variable $d_i(t)$



computations $\Rightarrow N^2 T$

Backward Procedure

$$\beta_T(\bar{i})$$

$$\beta_T(\bar{i}) = P(o_{t+1} o_{t+2} \dots o_T | a_t = \bar{i}, \lambda)$$

Initialization

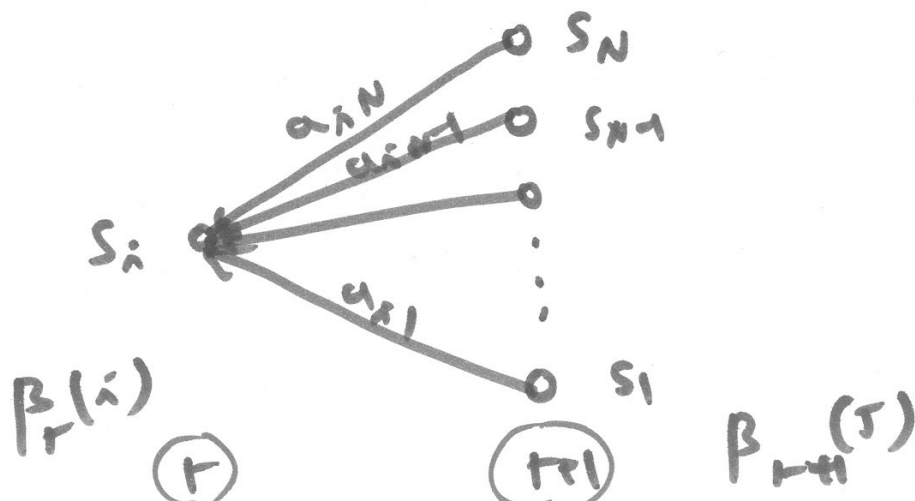
$$\beta_T(\bar{i}) = 1$$

$$1 \leq \bar{i} \leq N$$

Induction

$$\beta_t(\bar{i}) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

$$t = T-1, T-2, \dots, 1; \quad 1 \leq \bar{i} \leq N$$



$N^2 T$ computations

Solution to Problem-2 (Optimal state sequence)

- Individually most likely states at each time
- Map expectation of pairs of states (u_t, u_{t+1})
- Map expectation of triplets of states (u_{t-1}, u_t, u_{t+1})

$$\gamma_t(i) = P(u_t = i \mid o, \lambda) = \frac{P(o, u_t = i \mid \lambda)}{\sum_{\tilde{i}=1}^N P(o, u_t = \tilde{i} \mid \lambda)}$$

$$= \frac{P(o, u_t = i \mid \lambda)}{P(o \mid \lambda)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{\tilde{i}=1}^N \alpha_t(\tilde{i}) \beta_t(\tilde{i})}$$

$$u_t^* = \arg \max_i (\gamma_t(i))$$

For speech some $a_{1T} = 0$; \therefore optimal seq with u_t^* may not be valid sequence

Viterbi Algorithm (dynamic programming)

$$\delta_T(i) = \max_{v_1, v_2, \dots, v_{T-1}} P[o_1, v_1, o_2, v_2, \dots, v_{T-1}, v_T = i | \lambda]$$

By induction v_1, v_2, \dots, v_{T-1}

$$\rightarrow \delta_{T+1}(j) = \max_i [\delta_T(i) a_{ij}] b_j(o_{T+1})$$

$\psi_T(j)$ = tracking the argument that max for each $t \leq T$
 $t = 2 \text{ to } T; \bullet 1 \leq j \leq N$

Initialization

$$\delta_1(i) = \pi_i b_i(o_1) ; 1 \leq i \leq N$$

$$\psi_1(i) = 0$$

Recursion

$$\delta_T(j) = \max_i [\delta_{T-1}(i) a_{ij}] b_j(o_T) \quad 2 \leq T \leq T$$

$$\psi_T(j) = \arg \max_i [\delta_{T-1}(i) a_{ij}] \quad 1 \leq j \leq N$$

Viterbi Algorithm

Soln to problem-2

Termination

$$f_T(i) = \max_j \left[f_{T-1}(j) a_{ji} \right] b_i(o_T)$$

$$p^* = \max_i f_T(i) ; a_T^* = \operatorname{argmax}_i [f_T(i)]$$

Path backtracking (state sequence)

$$q_t^* = \psi_{t+1}(q_{t+1}^*)$$

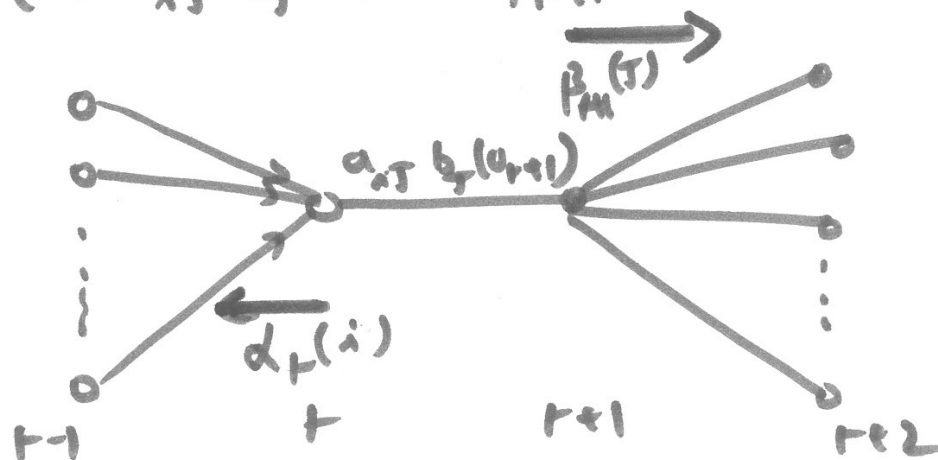
$$t = T-1, T-2, \dots, 1$$

Soln to Problem-3 (Parameter estimation)

Baum-welch method (structure approach)

$$\begin{aligned}\xi_T(i, T) &= P(u_T = i, u_{T+1} = T \mid o, \lambda) \\ &= P(o, u_T = i, u_{T+1} = T \mid \lambda) / P(o \mid \lambda)\end{aligned}$$

$$= \frac{\alpha_T(\lambda) a_{iT} b_T(o_{T+1}) \beta_{T+1}(T)}{\sum_i \sum_T \alpha_T(i) a_{iT} b_T(o_{T+1}) \beta_{T+1}(T)}$$



$$\delta_T(i) = P[a_T = i \mid o, \lambda] = \sum_{T=1}^N \xi_T(i, T)$$

$$\text{Avg \# transitions from state 'i' in 'o'} = \sum_{t=1}^{T-1} \delta_T(i)$$

$$\text{Avg \# transition from state } i \rightarrow j \text{ in 'o'} = \sum_{t=1}^{T-1} \xi_T(i, j)$$

$$\overline{\pi}_i = \text{expected \# times in state 'i' at } t=1 = \delta_1(i)$$

$$\overline{a_{ij}} = \frac{\text{Expected \# transitions from } i \rightarrow j}{\text{Expected \# transition from } i} = \frac{\sum_{t=1}^{T-1} \xi_T(i, j)}{\sum_{t=1}^{T-1} \xi_T(i)}$$

$$\overline{b_j}(k) = \frac{\text{Expected \# times in state 'j' } o_T = v_k}{\text{Expected \# times in state j}} = \frac{\sum_{t=1}^T \xi_T(j) / o_T = v_k}{\sum_{t=1}^T \xi_T(j)}$$

$$\overline{\lambda} = (\overline{A}, \overline{B}, \overline{\pi}); \quad P(o/\overline{\lambda}) > P(o/\lambda)$$

Types of HMMs

Based on structure of transition matrix

- Ergodic
- Left-right

(1) Ergodic (fully connected) HMM

$$\forall a_{iT} > 0$$

(2) Left-right model (Bakis model)

$$a_{iT} = 0 \quad T < \bar{i}$$

$$\pi_{\bar{i}} = 0 \quad \bar{i} \neq 1$$

$$= 1 \quad \bar{i} = 1$$

$$a_{iT} = 0 \quad T > \bar{i} + \Delta \bar{i}$$

$$\Delta \bar{i} = 2$$



(3) Discrete vs Continuous HMMs

Continuous observation densities in HMMs

Prob density of observation ~~seq~~ symbols at state $J = b_J(o)$

$$b_J(o) = \sum_{k=1}^M c_{Jk} N(o, \mu_{Jk}, \Sigma_{Jk})$$

Mixture coeff
for k^{th} component

Mean vector for k^{th} component

Covariance matrix for k^{th} component

Gaussian density

$$\sum_{k=1}^M c_{Jk} = 1$$

$$c_{Jk} \geq 0$$

$$1 \leq J \leq N ; 1 \leq k \leq M$$

$$\int_{-\infty}^{\infty} b_J(o) do = 1$$

$$1 \leq J \leq N$$

Reestimation formulas for mixture parameters in Continuous density HMM

$$\bar{C}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^K \gamma_t(j, k)}$$

$$\bar{\mu}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot o_t}{\sum_{t=1}^T \gamma_t(j, k)}$$

$$\bar{\Sigma}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) (o_t - \mu_{jk})(o_t - \mu_{jk})^T}{\sum_{t=1}^T \gamma_t(j, k)}$$

$$\gamma_t(j, k) = \left[\frac{\alpha_t(j) \beta_t(j)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)} \right] \left[\frac{C_{jk} N(o_t, \mu_{jk}, \Sigma_{jk})}{\sum_{k=1}^K C_{jk} N(o_t, \mu_{jk}, \Sigma_{jk})} \right]$$

$$\gamma_t(j, k) = \gamma_t(j)$$

Mixture with only one component