
HMM tutorial 5

by Dr Philip Jackson

- Implementing B-W formulae
 - Review of α , β , γ and ξ
 - Forward & backward procedures
 - Accumulate & update
- Emission or output pdfs
 - Multivariate Gaussian
 - Alternative functions
 - Gaussian mixture pdf
 - Re-estimation with mixtures
- Practical issues



<http://www.ee.surrey.ac.uk/Personal/P.Jackson/ce.adsp/>

Implementing B-W re-estimation

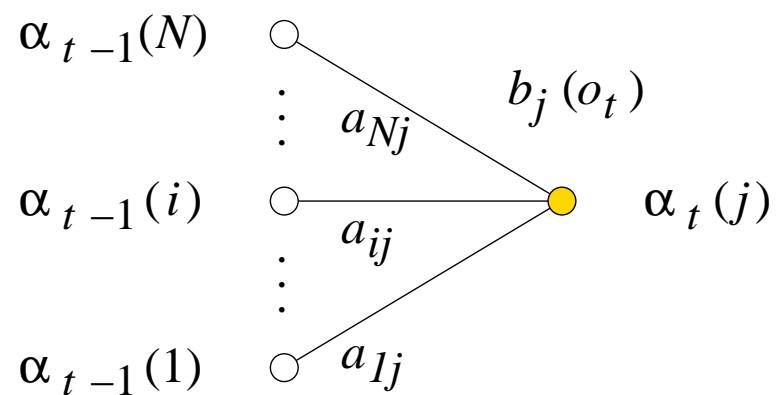
Overview

- **Forward:** compute likelihoods $\alpha_t(i)$
- **Backward:** compute likelihoods $\beta_t(i)$
- **Parallel:**
 - compute occupation $\gamma_t(i)$
 - compute transition $\xi_t(i, j)$
 - accumulate $\underline{\pi}_i$, \underline{a}_{ij} and \bar{a}_i
 - accumulate $\underline{\mu}_i$, $\underline{\Sigma}_i$ and \bar{b}_i
- *Repeat* for all files in the training set, $r \in \{1, 2, \dots, R\}$.

Forward likelihood

The forward likelihood is a joint probability,

$$\alpha_t(i) = P(x_t = i, \mathbf{o}_1^t | \lambda). \quad (1)$$

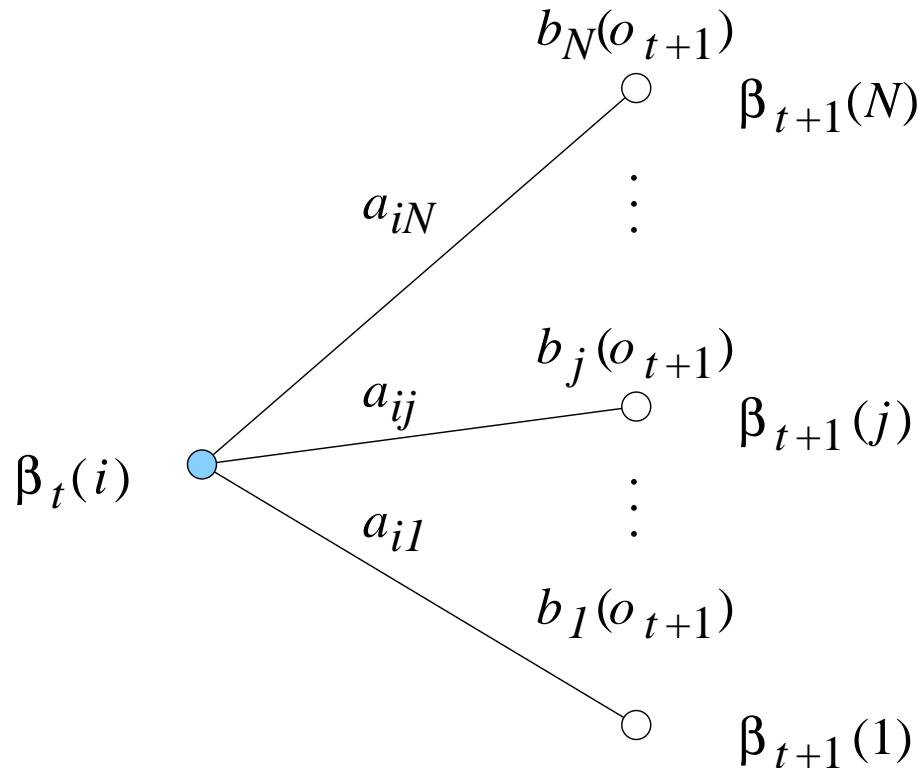


Trellis fragment depicting calculation of a forward likelihood.

Backward likelihood

The backward likelihood is a conditional probability,

$$\beta_t(j) = P(\mathbf{o}_{t+1}^T | x_t = j, \lambda). \quad (2)$$

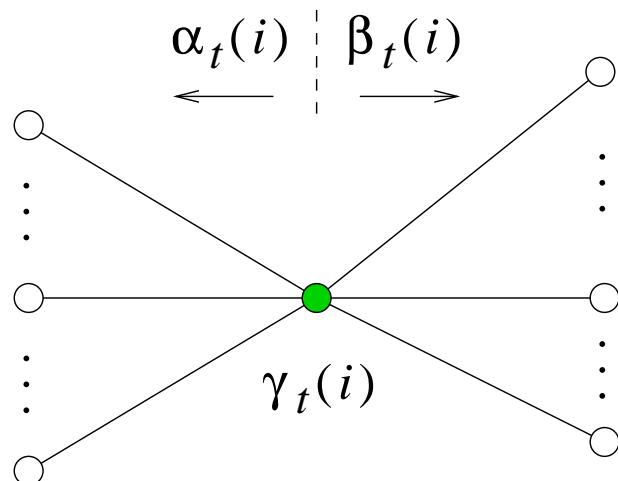


Trellis fragment depicting calculation of a backward likelihood.

Occupation likelihood

The occupation likelihood is a conditional probability,

$$\gamma_t(i) = P(x_t = i | \mathcal{O}, \lambda). \quad (3)$$

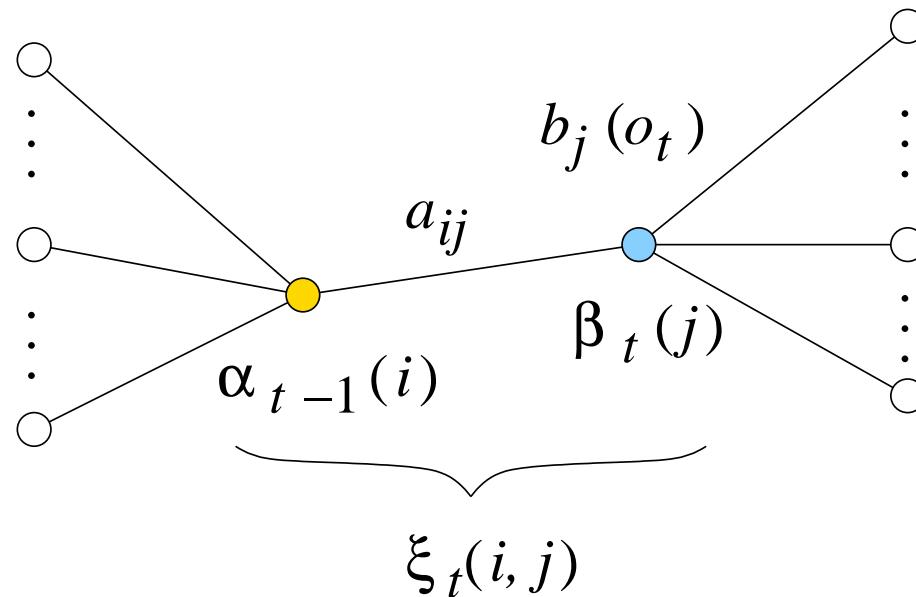


Trellis fragment depicting calculation of occupation likelihood.

Transition likelihood

The transition likelihood is the joint probability of states i and j , conditional on the observations,

$$\xi_t(i, j) = P(x_{t-1} = i, x_t = j | \mathcal{O}, \lambda). \quad (4)$$



Trellis fragment depicting calculation of a transition likelihood.

Forward procedure (Problem 1)

1. Initially,

$$\alpha_1(i) = \pi_i b_i(\mathbf{o}_1), \quad \text{for } 1 \leq i \leq N;$$

2. For $t = 2, 3, \dots, T$,

$$\alpha_t(j) = \left[\sum_{i=1}^N \alpha_{t-1}(i) a_{ij} \right] b_j(\mathbf{o}_t), \quad \text{for } 1 \leq j \leq N;$$

3. Finally, for the r th file,

$$P^r = P(\mathcal{O}|\lambda) = \sum_{i=1}^N \alpha_T(i).$$

Backward procedure (Problem 1)

1. Initially,

$$\beta_T(i) = 1, \quad \text{for } 1 \leq i \leq N;$$

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

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(\mathbf{o}_{t+1}) \beta_{t+1}(j), \quad \text{for } 1 \leq i \leq N.$$

3. Alternatively, we have $P^r = \sum_{i=1}^N \pi_i b_i(\mathbf{o}_1) \beta_1(i)$.

Parallel procedure

Occupation likelihoods:

$$\gamma_t(j) = \frac{\alpha_t(j) \beta_t(j)}{Pr}$$

Transition likelihoods:

$$\xi_t(i, j) = \frac{\alpha_{t-1}(i) a_{ij} b_j(\mathbf{o}_t) \beta_t(j)}{Pr}$$

Transition accumulators:

$$\underline{\pi}_i \mapsto \underline{\pi}_i + \gamma_1(i)$$

$$\underline{a}_{ij} \mapsto \underline{a}_{ij} + \sum_{t=2}^T \xi_t(i, j)$$

$$\bar{a}_i \mapsto \bar{a}_i + \sum_{t=2}^T \gamma_{t-1}(i)$$

Parallel (continued)

Emission accumulators:

$$\underline{\mu}_i \leftarrow \underline{\mu}_i + \sum_{t=1}^T \gamma_t(i) \mathbf{o}_t$$

$$\underline{\Sigma}_i \leftarrow \underline{\Sigma}_i + \sum_{t=1}^T \gamma_t(i) (\mathbf{o}_t - \underline{\mu}_i)(\mathbf{o}_t - \underline{\mu}_i)'$$

$$\bar{b}_i \leftarrow \bar{b}_i + \sum_{t=1}^T \gamma_t(i)$$

Repeat

For all files in the training set:

1. recompute the forward and backward likelihoods;
2. recompute the occupation and transition likelihoods;
3. increment the transition and emission accumulators.

Update (Problem 3)

Finally, we update the models:

$$\hat{\pi}_i = \frac{1}{R} \sum_{r=1}^R \underline{\pi}_i$$

$$\hat{a}_{ij} = \frac{\sum_{r=1}^R \underline{a}_{ij}}{\sum_{r=1}^R \bar{a}_i}$$

$$\hat{b}_i \left\{ \begin{array}{l} \hat{\mu}_i = \frac{\sum_{r=1}^R \underline{\mu}_i}{\sum_{r=1}^R \bar{b}_i} \\ \hat{\Sigma}_i = \frac{\sum_{r=1}^R \underline{\Sigma}_i}{\sum_{r=1}^R \bar{b}_i} \end{array} \right.$$

where $\hat{\cdot}$ denotes new values of the model parameters,
 $\lambda = \{\pi, A, B\}$.

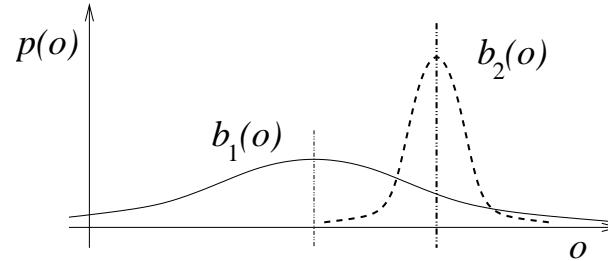
Summary of the training procedure

1. initialise accumulators for all HMMs' parameters
2. read the next training utterance
3. join HMMs in sequence to make composite HMM
4. calculate forward & backward probabilities
5. increment accumulators using forwd/backwd probs
6. repeat from step 2 until all utterances processed
7. use accumulators to update parameters for all HMMs

Parametric emission pdfs

Univariate Gaussian

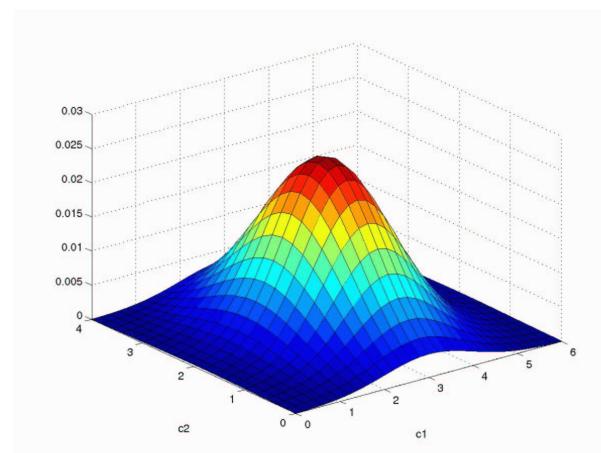
$$b_i(o_t) = \frac{1}{\sqrt{2\pi\Sigma_i}} \exp \left[-\frac{(o_t - \mu_i)^2}{2\Sigma_i} \right].$$



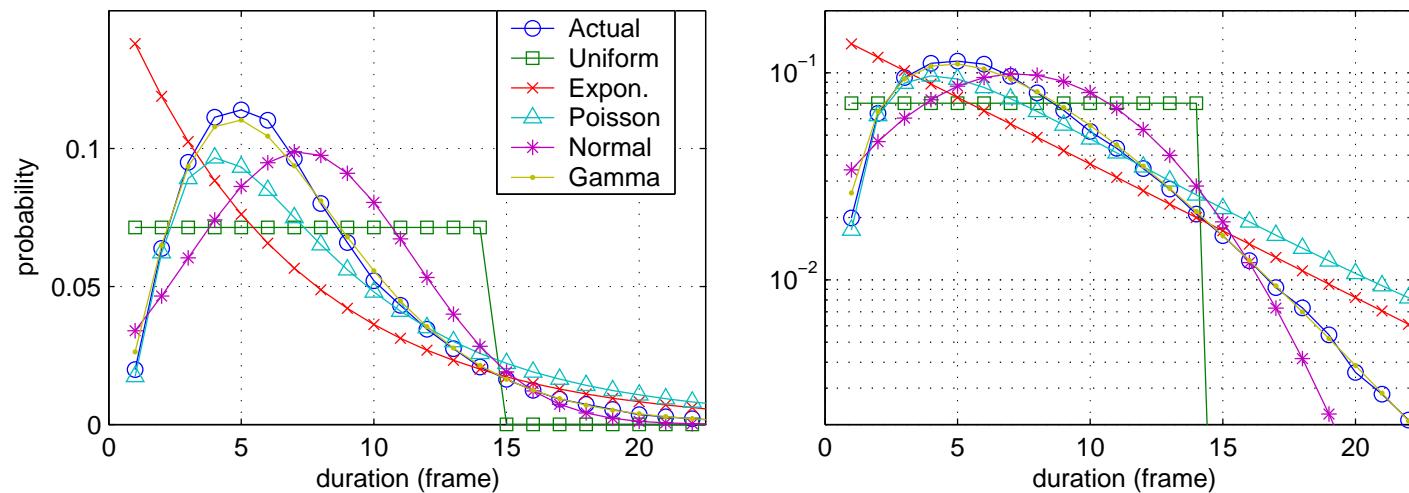
Multivariate Gaussian

$$b_i(o_t) = \frac{1}{\sqrt{(2\pi)^K |\Sigma_i|}} \exp \left[-\frac{1}{2}(o_t - \mu_i)\Sigma_i^{-1}(o_t - \mu_i)' \right],$$

where the dimensionality of the observation space is K .



Other distributions



Various functions on (left) linear and (right) logarithmic scales.

- Exponential
- Poisson
- Gamma
- Log-normal
- Gaussian mixture
- Ricean
- Rayleigh
- Beta
- student-*t*
- Cauchy

Gaussian mixture pdfs

Univariate Gaussian mixture

$$\begin{aligned} b_i(o_t) &= \sum_{m=1}^M c_{im} \mathcal{N}(o_t; \mu_{im}, \Sigma_{im}) \\ &= \sum_{m=1}^M \frac{c_{im}}{\sqrt{2\pi\Sigma_{im}}} \exp\left[-\frac{(o_t - \mu_{im})^2}{2\Sigma_{im}}\right], \end{aligned} \quad (5)$$

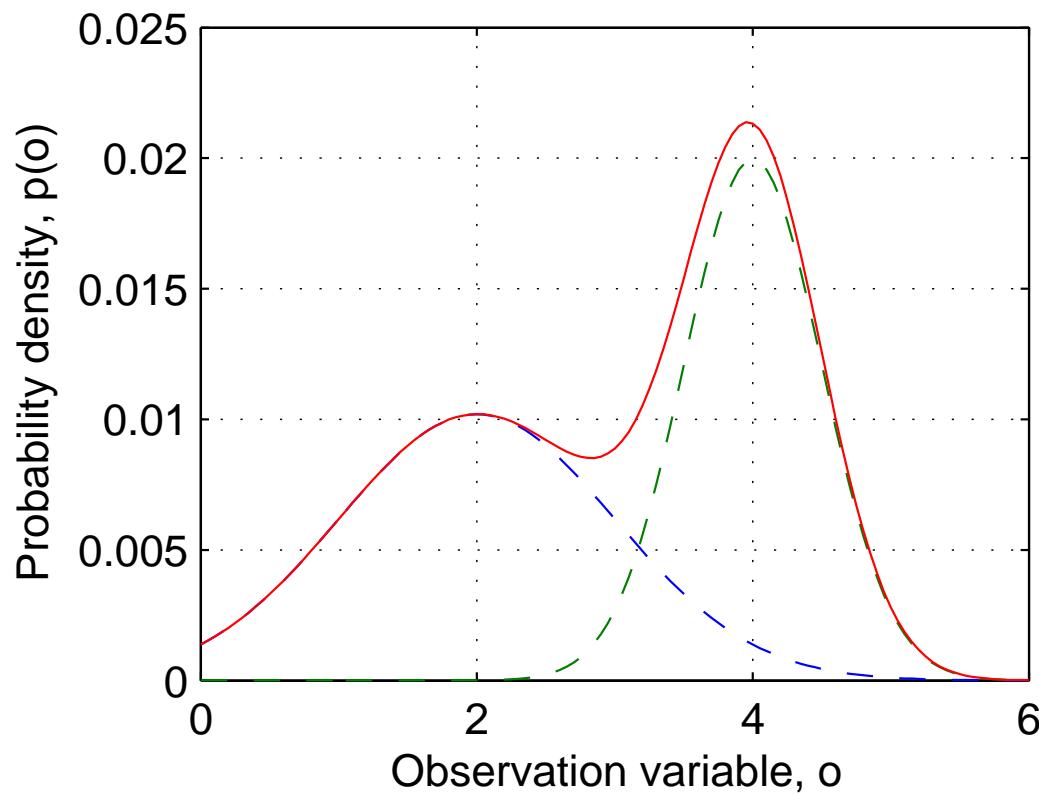
where M is the number of mixture components in the Gaussian mixture (i.e. M -mix), and $\sum_{m=1}^M c_{im} = 1$.

Multivariate Gaussian mixture

$$b_i(\mathbf{o}_t) = \sum_{m=1}^M c_{im} \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{im}, \boldsymbol{\Sigma}_{im}), \quad (6)$$

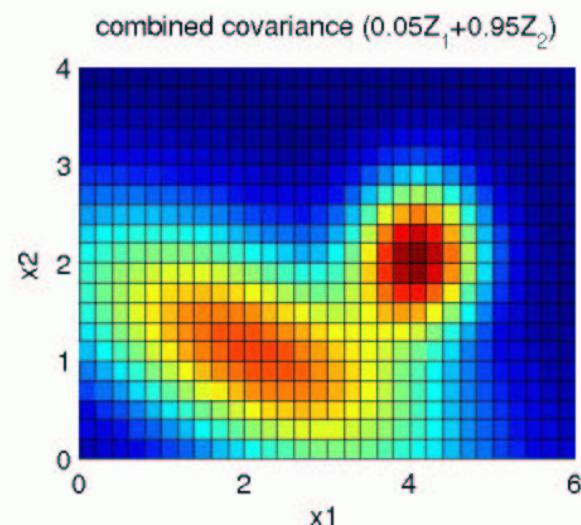
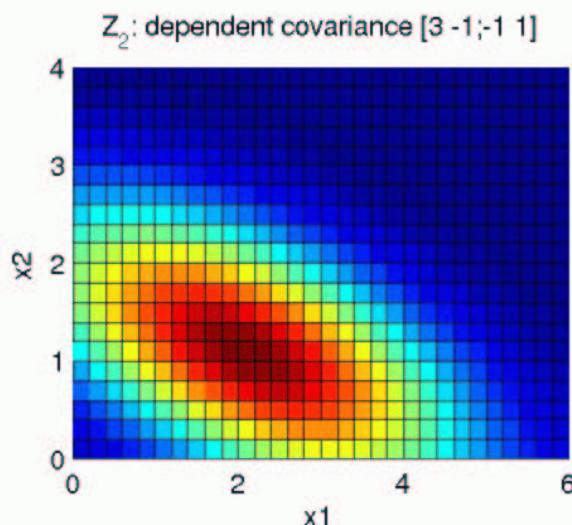
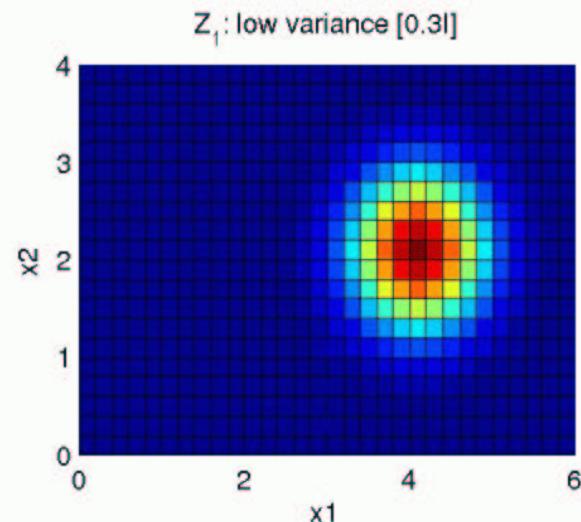
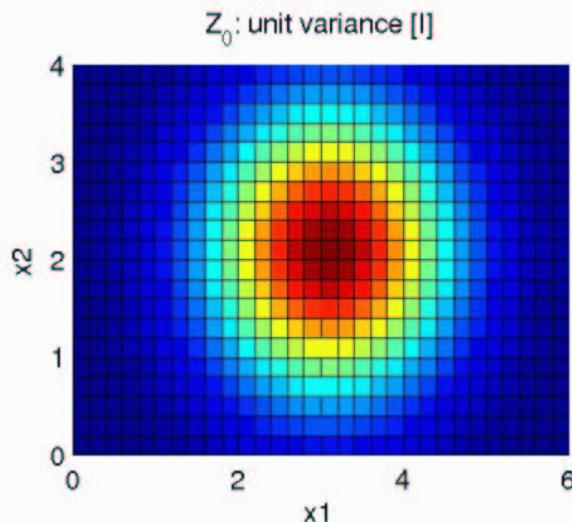
where $\mathcal{N}(\mathbf{o}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ evaluated at \mathbf{o} .

Univariate Gaussian mixture



Example of a univariate Gaussian mixture pdf.

Multivariate Gaussian mixture



Examples of multivariate Gaussian pdfs
and a multivariate Gaussian mixture.

Viterbi training with mixtures

Mean: $\hat{\mu}_{jm} = \frac{\sum_{t=1}^T q_t(j,m) \mathbf{o}_t}{\sum_{t=1}^T q_t(j,m)},$

Covariance: $\hat{\Sigma}_{jm} = \frac{\sum_{t=1}^T q_t(j,m) (\mathbf{o}_t - \hat{\mu}_{jm})(\mathbf{o}_t - \hat{\mu}_{jm})'}{\sum_{t=1}^T q_t(j,m)},$

Weights: $\hat{c}_{jm} = \frac{\sum_{t=1}^T q_t(j,m)}{\sum_{t=1}^T q_t(j)},$

where $q_t(j, m)$ and $q_t(j)$ are occupation indicators:

$$q_t(j) = \begin{cases} 1 & \text{for } i = x_t; \\ 0 & \text{otherwise,} \end{cases}$$
$$q_t(j, m) = \begin{cases} 1 & \text{for } i = x_t, m = \hat{m}; \\ 0 & \text{otherwise,} \end{cases}$$

and the occupied mixture is

$$\hat{m} = \arg \max_m c_{jm} \mathcal{N}(\mathbf{o}_t; \hat{\mu}_{jm}, \hat{\Sigma}_{jm}).$$

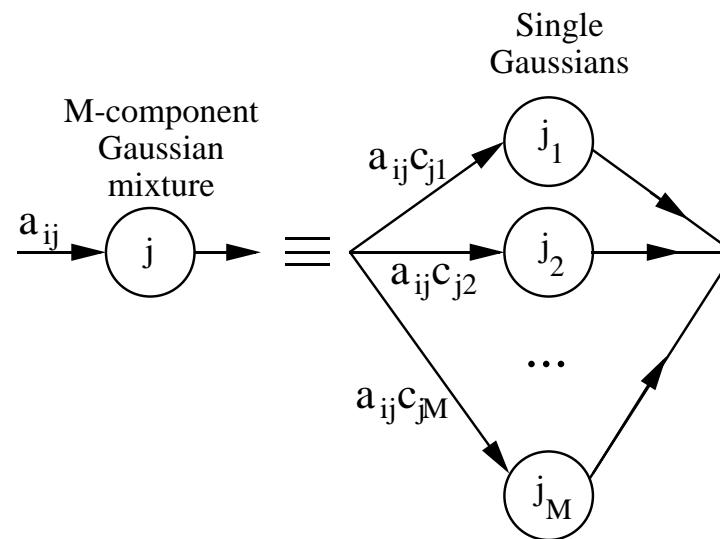
Baum-Welch training with mixtures

Mixture-occupation likelihood

$$\gamma_t(j, m) = \frac{\alpha_t^\emptyset(j) c_{jm} \mathcal{N}(o_t; \mu_{jm}, \Sigma_{jm}) \beta_t(j)}{P(\mathcal{O}|\lambda)}, \quad (7)$$

where

$$\gamma_t(j) = \sum_{m=1}^M \gamma_t(j, m), \quad \text{and} \quad \alpha_t^\emptyset(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right].$$



Representing a mixture of Gaussians.*

B-W re-estimation of mixture parameters

Mean:

$$\hat{\mu}_{jm} = \frac{\sum_{t=1}^T \gamma_t(j, m) \mathbf{o}_t}{\sum_{t=1}^T \gamma_t(j, m)}, \quad (8)$$

Covariance:

$$\hat{\Sigma}_{jm} = \frac{\sum_{t=1}^T \gamma_t(j, m) (\mathbf{o}_t - \hat{\mu}_{jm})(\mathbf{o}_t - \hat{\mu}_{jm})'}{\sum_{t=1}^T \gamma_t(j, m)}, \quad (9)$$

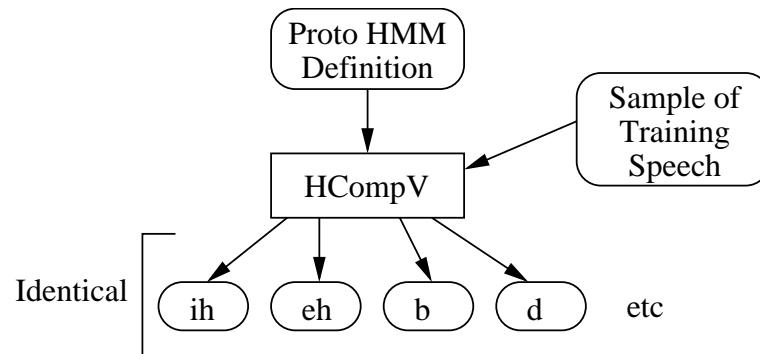
Mixture weights:

$$\hat{c}_{jm} = \frac{\sum_{t=1}^T \gamma_t(j, m)}{\sum_{t=1}^T \gamma_t(j)}. \quad (10)$$

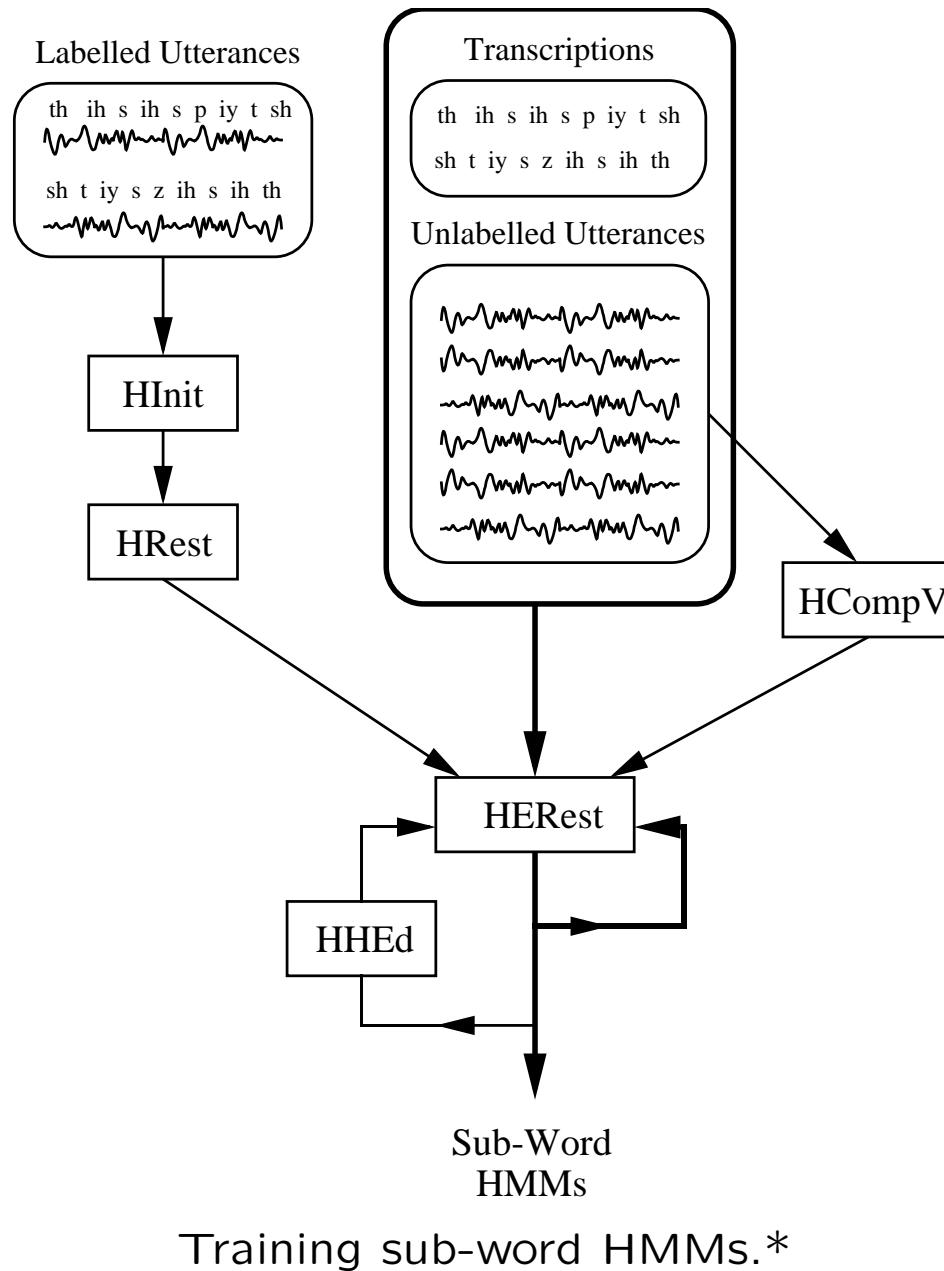
Practical issues

Model initialisation

1. Random
2. Flat start
3. Least squares
4. Viterbi
5. Baum-Welch
(supervised)
6. Baum-Welch
(unsupervised)



Re-estimation and embedded re-estimation

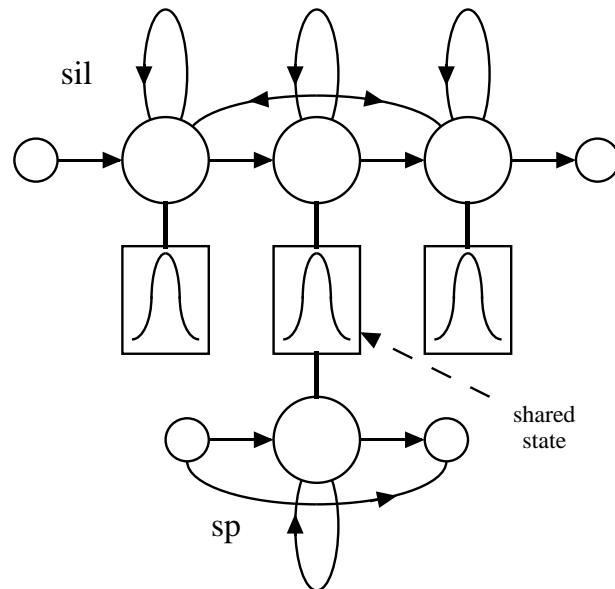


Number of parameters

- Parsimony
 - Occam's razor
- Amount of training data

Regularisation

- variance floor
- parameter tying
 - context sensitivity
 - agglomerative clustering
 - phonetic decision tree



Linking the silence and short pause models.*

Context-sensitive models

- **Problem:** coarticulation causes phones to vary
- **Solution:** train separate models for each context
 - monophone:
sil - t - e - n - sil
 - triphone:
 $\text{sil}_t - \text{sil} \text{t}_e - t \text{e}_n - e \text{n}_{\text{sil}} - n \text{sil}$

Tutorial 5 summary

- Review of likelihoods α_t , β_t , γ_t and ξ_t
- Implementation of B-W formulae:
 - forward, backward, accumulate & update
- Output probability distribution functions
 - Multivariate Gaussian mixture, etc.
- Practical issues
 - Annotations, initialisation & tying
 - Context-sensitive models