

A tutorial on hidden markov models and selected applications in speech recognition



Author

Rabiner L.

Journal

Proceeding of the IEEE, 1989

Speaker

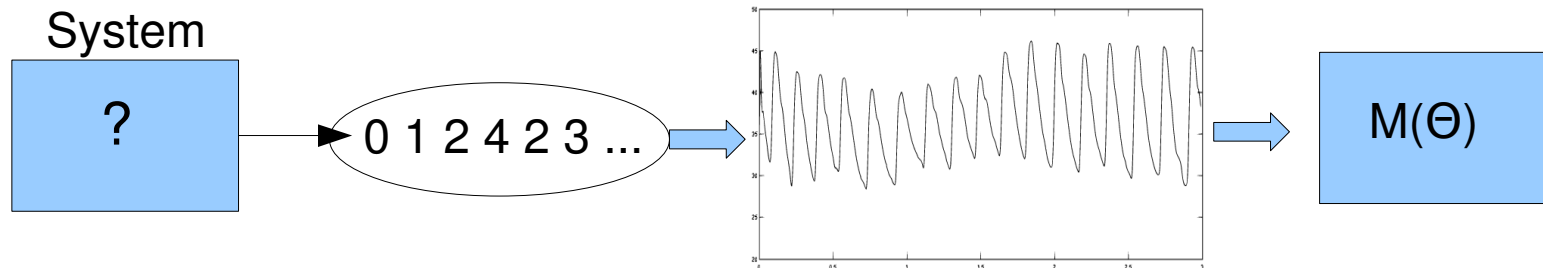
Simone Tognetti



- Introduction
 - Graph models for signal
- Learning a model
 - ML
 - EM
- Markov Chain
- HMM
 - Solve HMMs
- Extension of HMMs
- Applications
- Conclusion

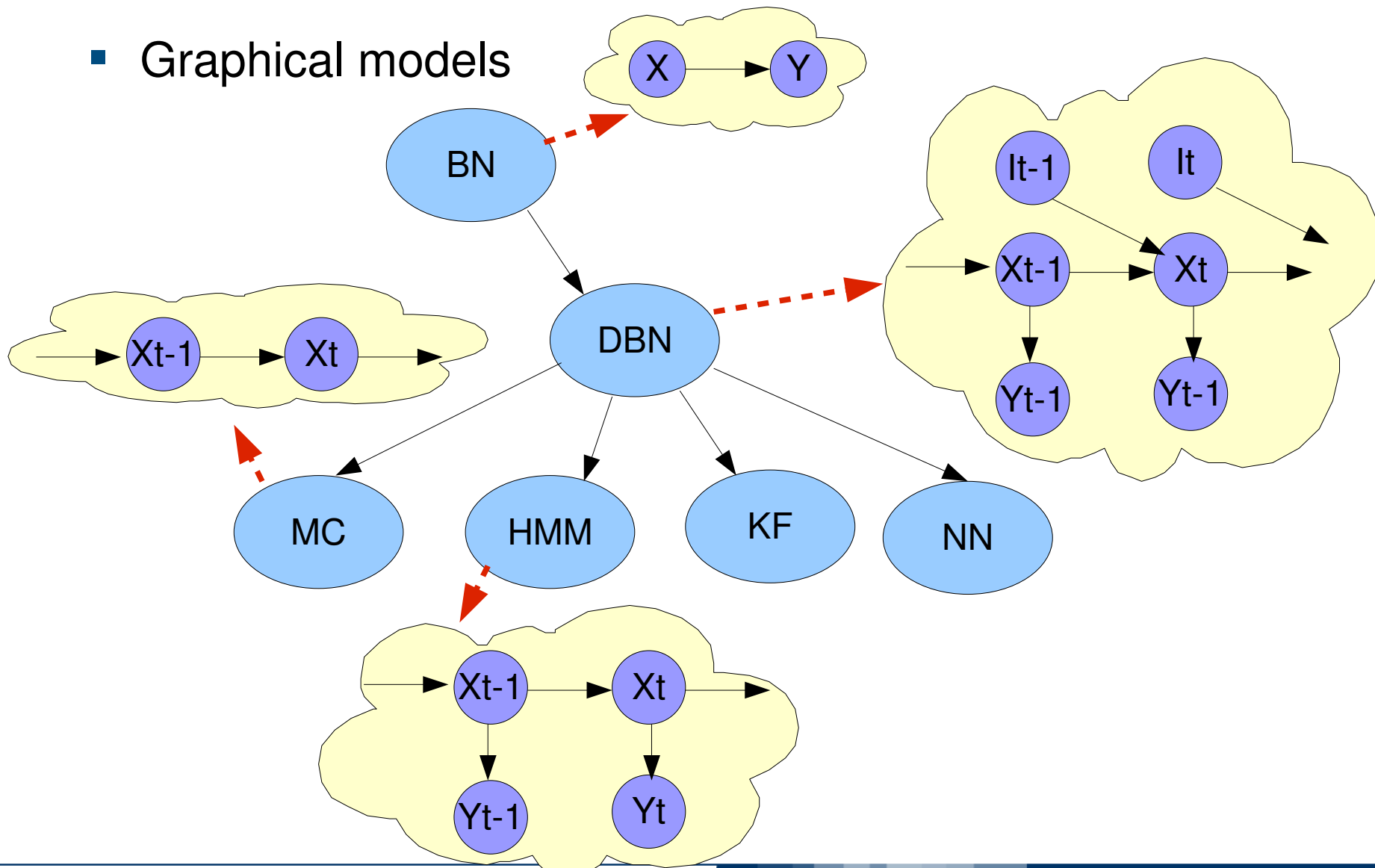


- Real world system generate observable outputs
- Outputs can be viewed as stochastic signals
- We can build a model of that signal because :
 - the model provide information concerning how to do the processing of the signal (ex. how to remove noise)
 - the model give informations about the system that have generated the signal
 - from the model we can make predictions.





- Graphical models





- Introduction
 - Graph models for signal
- Learning a model
 - ML
 - EM
- Markov Chain
- HMM
 - Solve HMMs
- Extension of HMMs
- Applications
- Conclusion



- Maximum likelihood estimation
 - Estimation of the parameters of a stochastic model when all the random variable are observable.
 - Find the parameters that maximize the likelihood

Sequence of random variable iid $\sim N(\mu, \sigma^2)$

Y_{t-1}

Y_t

Y_{t+1}

$$\mathcal{L}(\theta) = \log P(\{Y_t\}|\theta) = \log \prod_{t=1}^T \left(\frac{1}{2\pi\sigma^2} \right) \exp \left\{ -\frac{(Y_t - \mu)^2}{2\sigma^2} \right\}$$

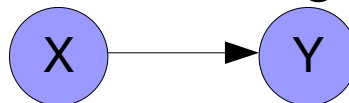
$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}(\theta)}{\partial \mu} = \frac{1}{T} \sum_{t=1}^T Y_t \\ \frac{\partial \mathcal{L}(\theta)}{\partial \sigma^2} = \frac{1}{T} \sum_{t=1}^T (Y_t - \mu)^2 \end{array} \right.$$



- Expectation Maximization (EM)

Learn parameters of a stochastic model when some variables are hidden. (ie. HMM)

- Learn means to find something that we don't see
- A simple example



$$\mathcal{L}(\theta) = \log P(Y|\theta) = \log \sum_X P(Y, X|\theta) = \log \sum_X Q(X) \frac{P(Y, X|\theta)}{Q(X)}$$

$$\mathcal{L}(\theta) \geq E(Q, \theta) - H(Q) = F(Q, \theta)$$

- EM maximize the likelihood in a iterative with two steps

- Expectation $Q_{k+1} \leftarrow \arg \max_Q F(Q, \theta_k)$

- Maximization $\theta_{k+1} \leftarrow \arg \max_{\theta} F(Q_{k+1}, \theta)$

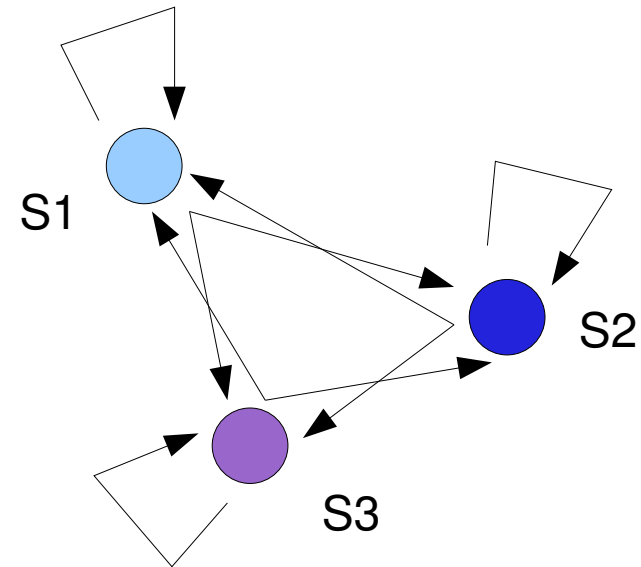


- Introduction
 - Graph models for signal
- Learning a model
 - ML
 - EM
- Markov Chain
- HMM
 - Solve HMMs
- Extension of HMMs
- Applications
- Conclusion

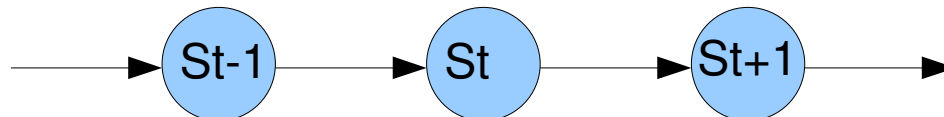


- Set of state S
- Transition Matrix M

$$M = \begin{pmatrix} P(S1|S1) & P(S2|S1) & P(S3|S1) \\ P(S1|S2) & P(S2|S2) & P(S3|S2) \\ P(S1|S3) & P(S2|S3) & P(S3|S3) \end{pmatrix}$$

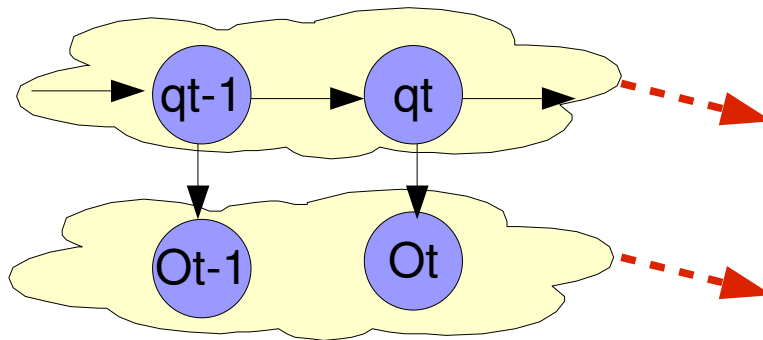


- Prior Probability
 $\pi(0) = [P(S1) \ P(S2) \ P(S3)]$
- Marginal (Posterior probability)
 $\pi(t) = \pi(0) M^t$
- DBN





- DBN



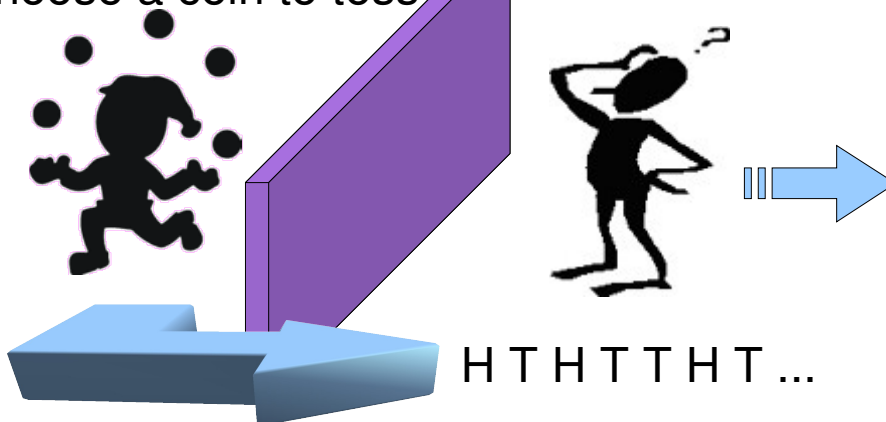
Underlying Markov Process that is not Observable

Output Markov Process that is Observable

- How to think about them?

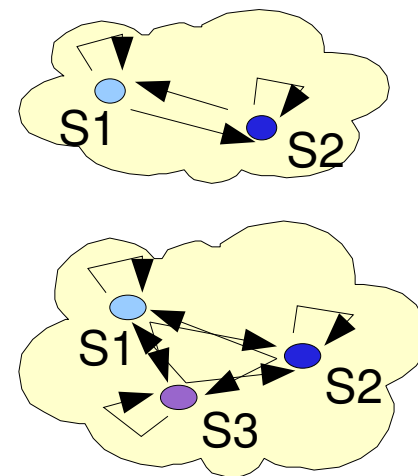
.. i.e. coin tossing

Underlying process:
choose a coin to toss



Output process: depends on the chosen coin

Hypothesis on the underlying model:





- Set of state: q_t in $\{ S_1, S_2, S_3, \dots, S_N \}$
- Underlying process transition matrix A
 $a_{ij} = P(q_t = S_i | q_{t-1} = S_j)$
- Output process matrix B (Produce a set of observation O)
 - Discrete: O_t in $\{ v_1, v_2, \dots, v_D \}$
 $b_i(k) = P(O_t = v_k | q_t = S_i)$, B is $N \times D$
 - Continuous: O_t in \mathbb{R}
 $b_i = P(O_t | q_t = S_i)$, B is a vector of distribution
- Prior probability on states
 $\pi(0) = [P(q_0 = S_1), P(q_0 = S_2), \dots, P(q_0 = S_N)]$



- Basic problems

We have an HMM $\lambda = (A, B, \pi)$ and a sequence of observation $O = O_1, O_2, \dots, O_n$

1. $P(O|\lambda)$: Probability that an observation sequence came from a given HMM
2. $P(q_1, q_2, \dots, q_t|O)$: Probability to have a specific state sequence given the observation sequence
3. How to adjust λ to maximize $P(O|\lambda)$: Learning the HMM



- Given a state sequence $Q = q_1 q_2 \cdots q_T$

- Probability of the observation sequence

$$P(O|Q, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdots b_{q_T}(O_T).$$

- Probability of a state sequence

$$P(Q|\lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \cdots a_{q_{T-1} q_T}.$$

- Join probability of Q and O $P(O, Q|\lambda) = P(O|Q, \lambda) P(Q, \lambda).$

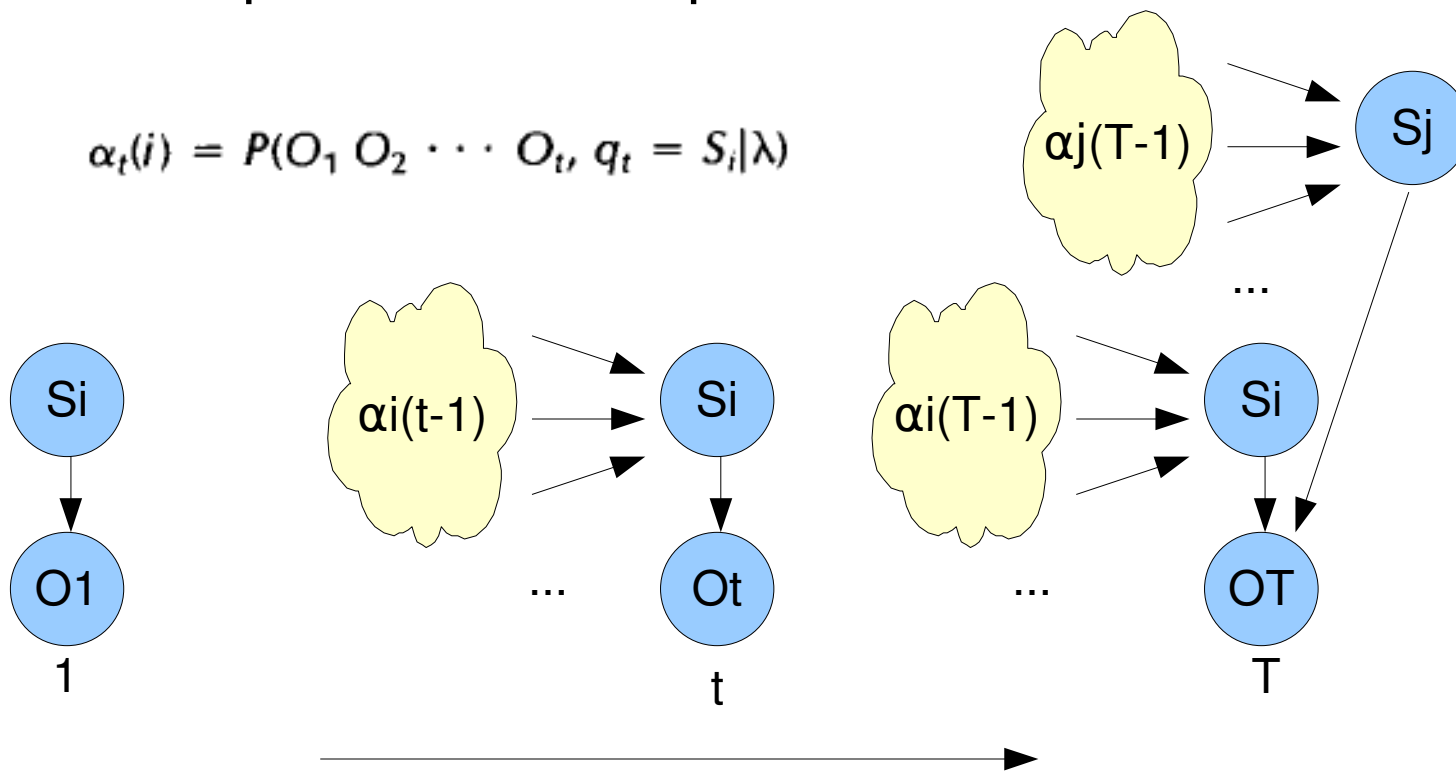
- Probability of O by summation over all state sequence

$$P(O|\lambda) = \sum_{\text{all } Q} P(O|Q, \lambda) P(Q|\lambda)$$



- Forward procedure: compute a forward variable

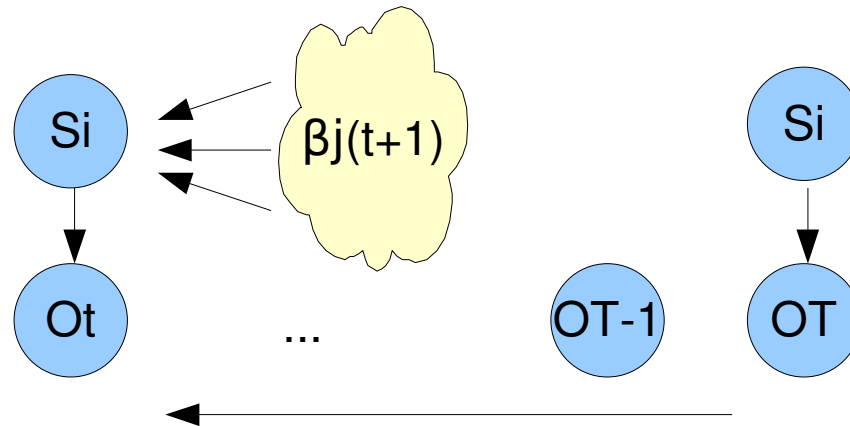
$$\alpha_t(i) = P(O_1 O_2 \cdots O_t, q_t = S_i | \lambda)$$



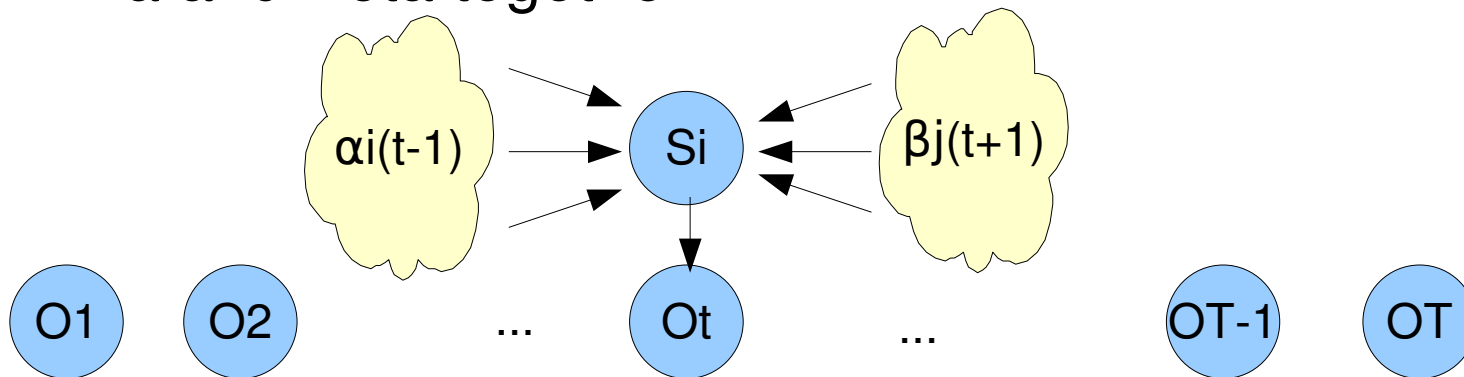
$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i).$$



- Forward Variable $\alpha_t(i) = P(O_1 O_2 \cdots O_t, q_t = S_i | \lambda)$
 - Backward variable $\beta_t(i) = P(O_{t+1} O_{t+2} \cdots O_T | q_t = S_i, \lambda)$
- } E-step



- Alfa and Beta together





- Probability of being in state S_i at time t given the observation sequence

$$\gamma_t(i) = P(q_t = S_i | O, \lambda)$$

- Expression in terms of forward and backward variable

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{P(O | \lambda)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

- Possible solution: find local best state

$$q_t = \operatorname{argmax}_{1 \leq i \leq N} [\gamma_t(i)], \quad 1 \leq t \leq T.$$

- but possible solution are not ammissible state sequence

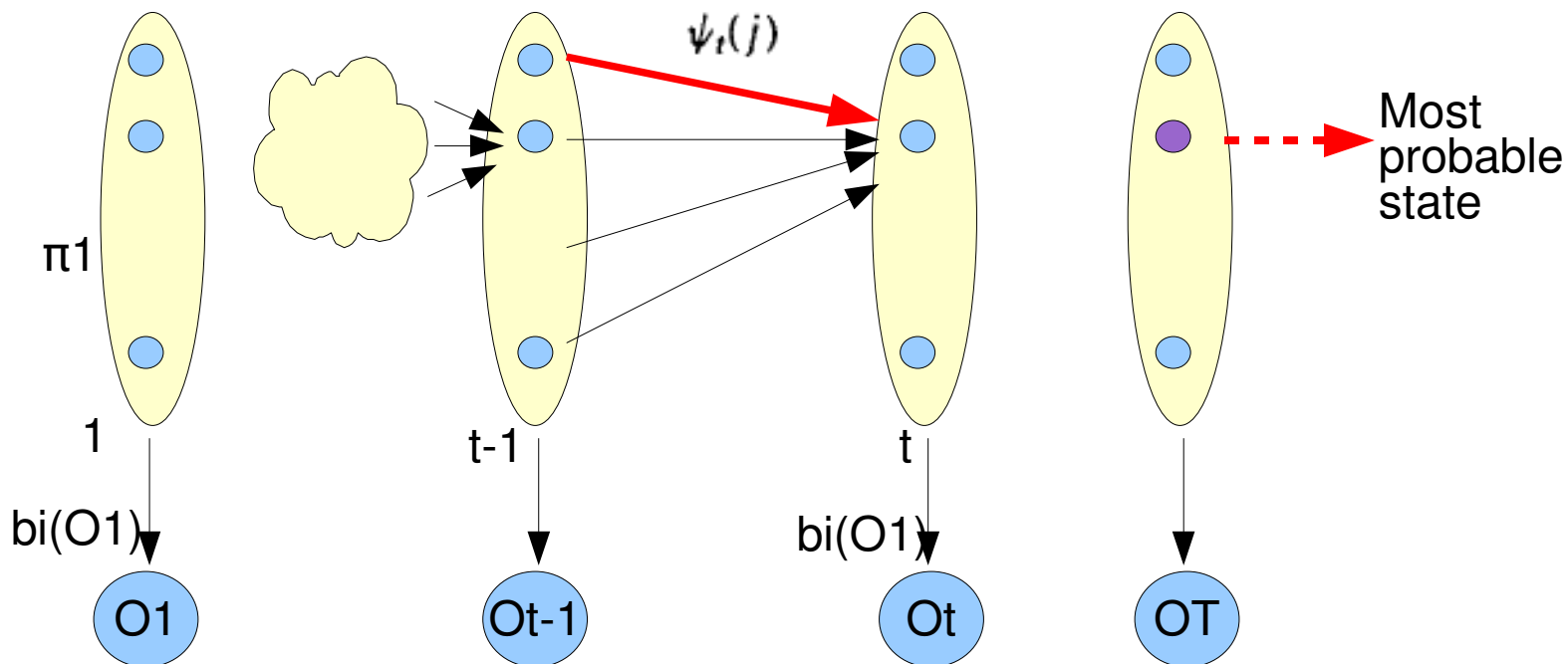


- Viterbi algorithm

- Best score along a single path

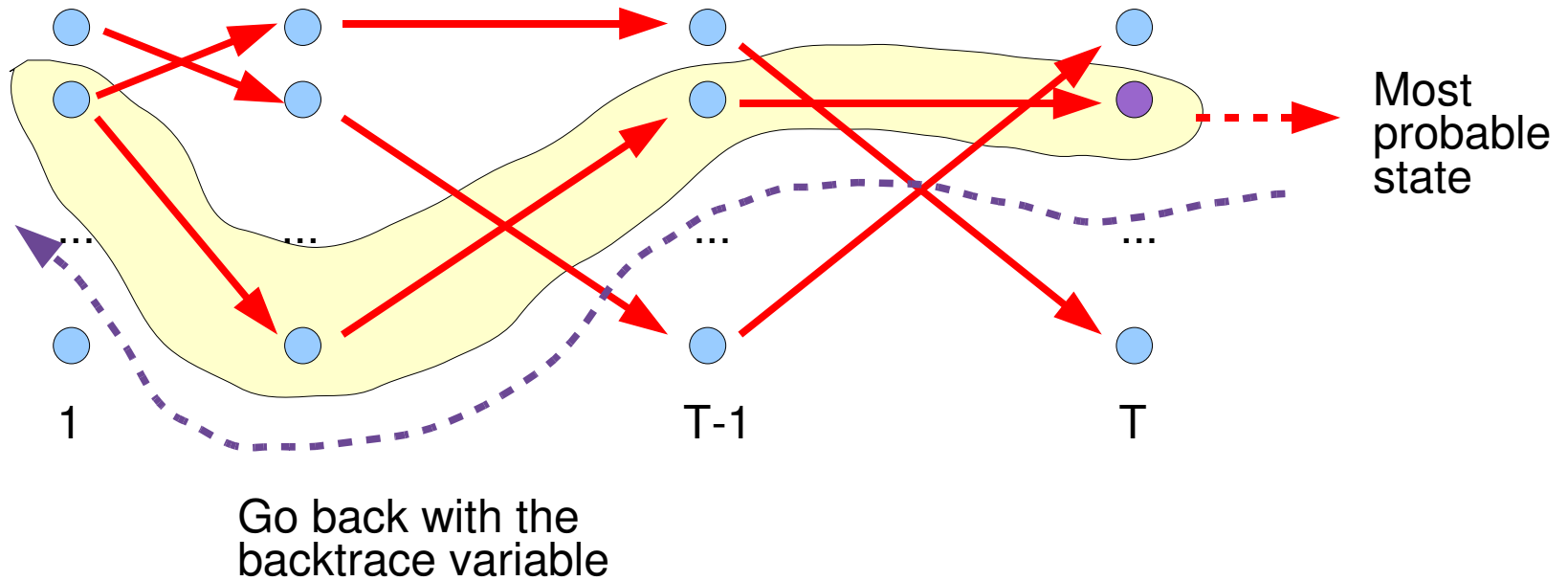
$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1 q_2 \dots q_t = i, O_1 O_2 \dots O_t | \lambda]$$

- Keep a track of the best path that can reach the state j at time t





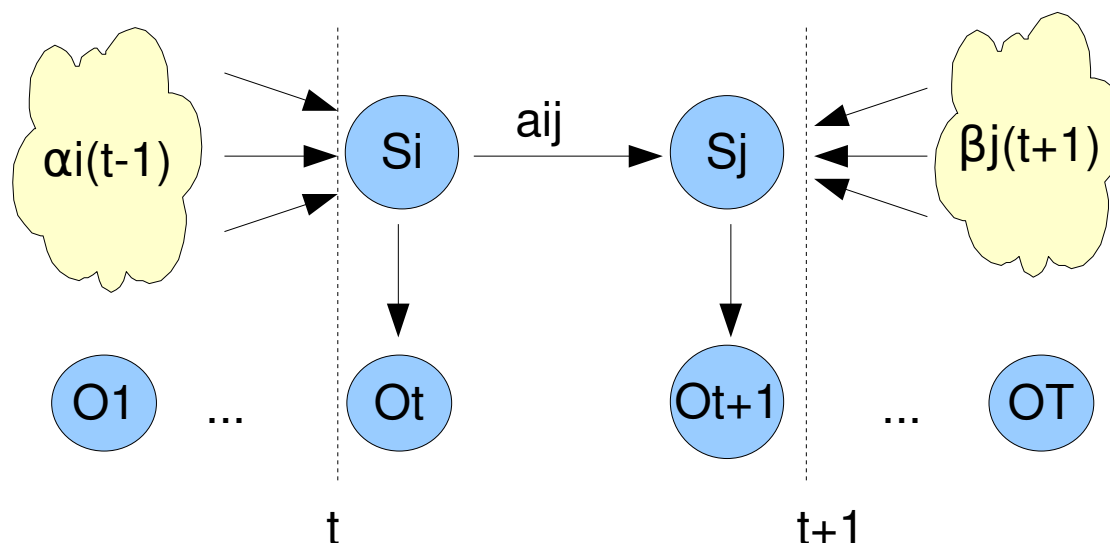
- Global behaviour





- Baum-Welch (EM algorithm for HMM)
 - Probability of being in state i at time t and in state j at time $t+1$

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda).$$



Extract the effect of parameter a_{ij}



- Baum-Welch (EM algorithm for HMM)

- Expression with forward and backward variable

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O|\lambda)} = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}$$

- Relation between $\gamma_t(i)$ and $\xi_t(i, j)$ $\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$.

- Expected number of transition from S_i $\sum_{t=1}^{T-1} \gamma_t(i)$

- Expected number of transition from S_i to S_j $\sum_{t=1}^{T-1} \xi_t(i, j)$



- Baum-Welch: estimation of parameters

$\bar{\pi}_i$ = expected frequency (number of times) in state S_i at time ($t = 1$) = $\gamma_1(i)$

$$\bar{a}_{ij} = \frac{\text{expected number of transitions from state } S_i \text{ to state } S_j}{\text{expected number of transitions from state } S_i} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\bar{b}_j(k) = \frac{\text{expected number of times in state } j \text{ and observing symbol } v_k}{\text{expected number of times in state } j} = \frac{\sum_{t=1}^T \gamma_t(j) \text{ s.t. } O_t = v_k}{\sum_{t=1}^T \gamma_t(j)}$$



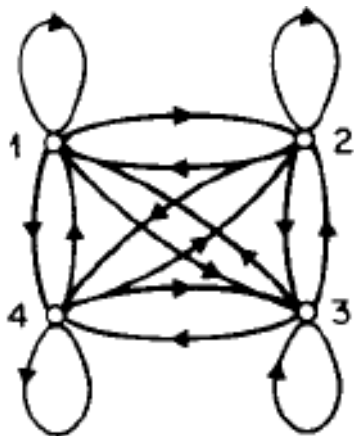
- Baum-Welch: Algorithm
 1. Initial estimation of λ
 2. E-step: Compute α , β
 3. M-Step: New estimation of parameter λ'
 4. Repeat from 2 until convergence
- EM : if we compute the likelihood we obtain same equation for the M-step
- Problem: Learning is sensible to initial parameter value when we have continuous observations



- Introduction
 - Graph models for signal
- Learning a model
 - ML
 - EM
- Markov Chain
- HMM
 - Solve HMMs
- **Extension of HMMs**
- Applications
- Conclusion



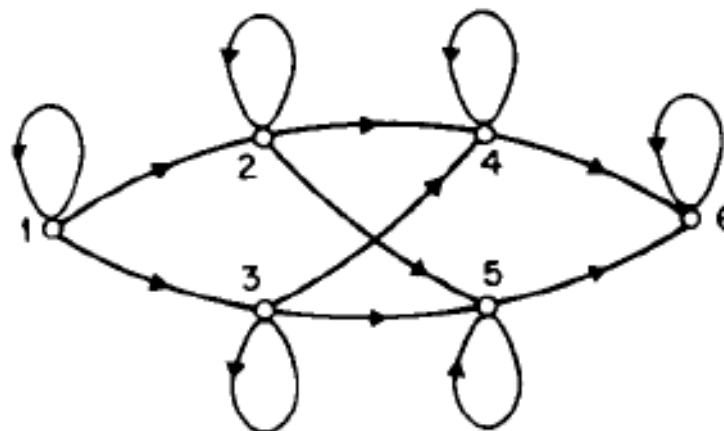
- Type change with the structure of the transition matrix of the underlying process



Ergodic HMM with 4 state



Left to right HMM with 4 state



Parallel path left to right HMM with 6 state

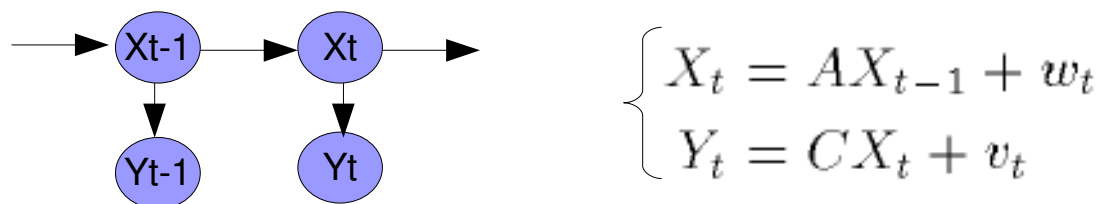


- Continuous observations: we can model the distribution probability of the output with a mixture

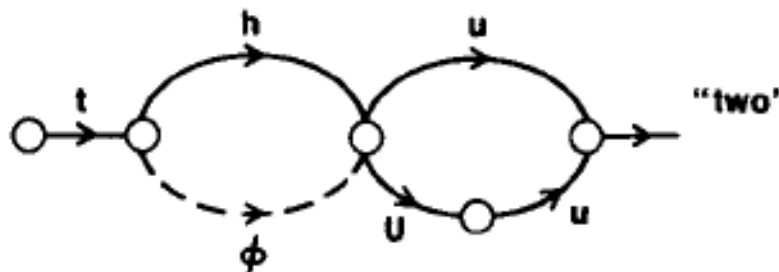
$$b_j(\mathbf{O}) = \sum_{m=1}^M c_{jm} \mathcal{N}[\mathbf{O}, \boldsymbol{\mu}_{jm}, \mathbf{U}_{jm}], \quad 1 \leq j \leq N$$

Prior vector Mean vector Covariance matrix } New parameter to estimate

- Kalman filter

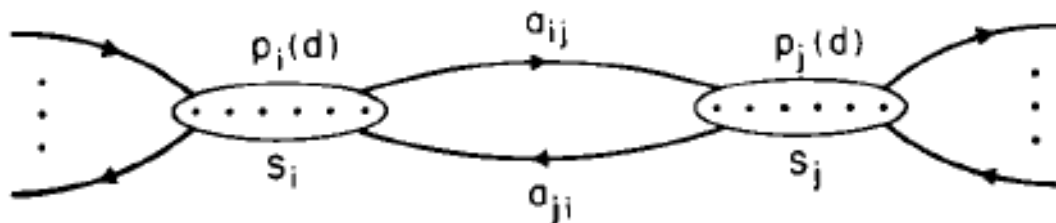


- Null transition





- Explicit state duration



- Distance within two HMM $\lambda_1 = (A_1, B_1, \pi_1)$ $\lambda_2 = (A_2, B_2, \pi_2)$

$$D(\lambda_1, \lambda_2) = \frac{1}{T} [\log P(O^{(2)}|\lambda_1) - \log P(O^{(2)}|\lambda_2)]$$

$$D_s(\lambda_1, \lambda_2) = \frac{D(\lambda_1\lambda_2) + D(\lambda_2, \lambda_1)}{2}. \quad \text{Symmetric version}$$

- Scaling

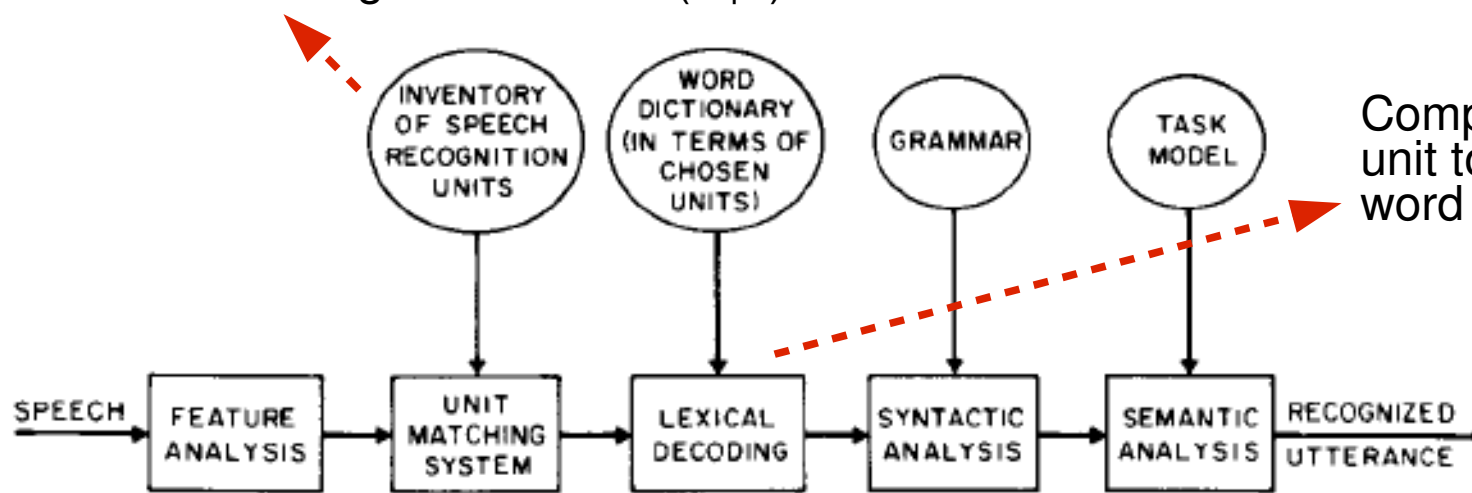


- Introduction
 - Graph models for signal
- Learning a model
 - ML
 - EM
- Markov Chain
- HMM
 - Solve HMMs
- Extension of HMMs
- **Applications**
- Conclusion



- A generic framework

Mapping from HMM to a voice unit.
Given the observation choose the HMM that give the best $P(O|\lambda)$



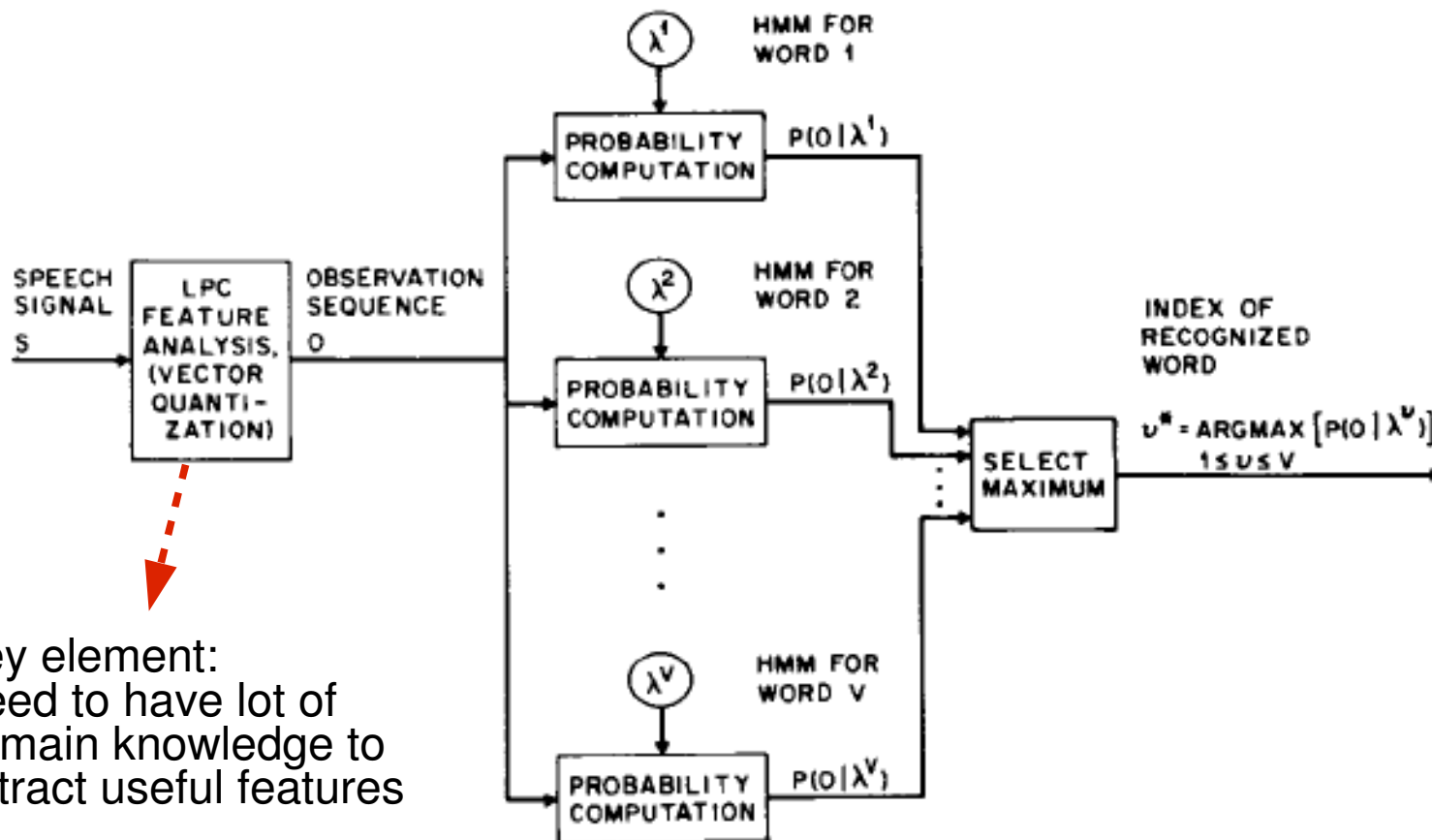
Composition of unit to have a word

Temporal and spectral analysis to obtain a observation sequence

Recognition of unit of the language. Word are divided into small unit to have a small set of models



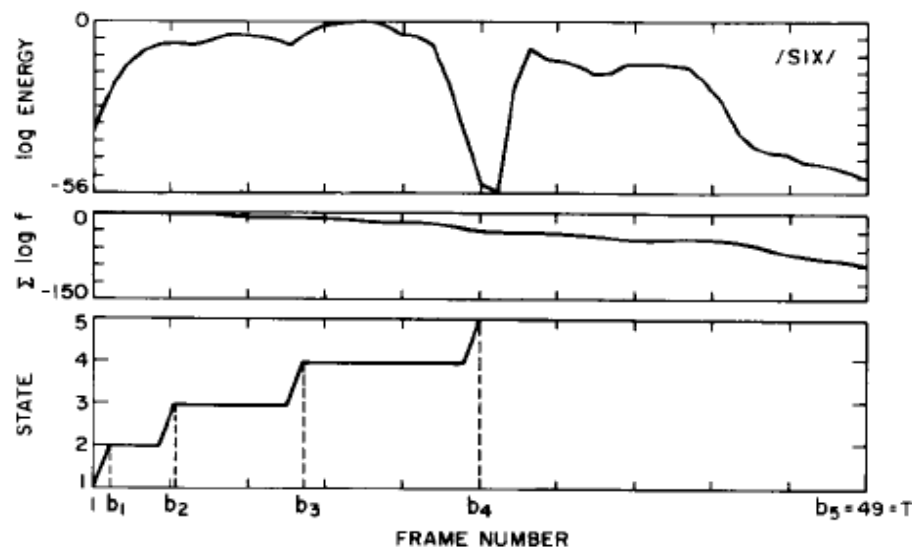
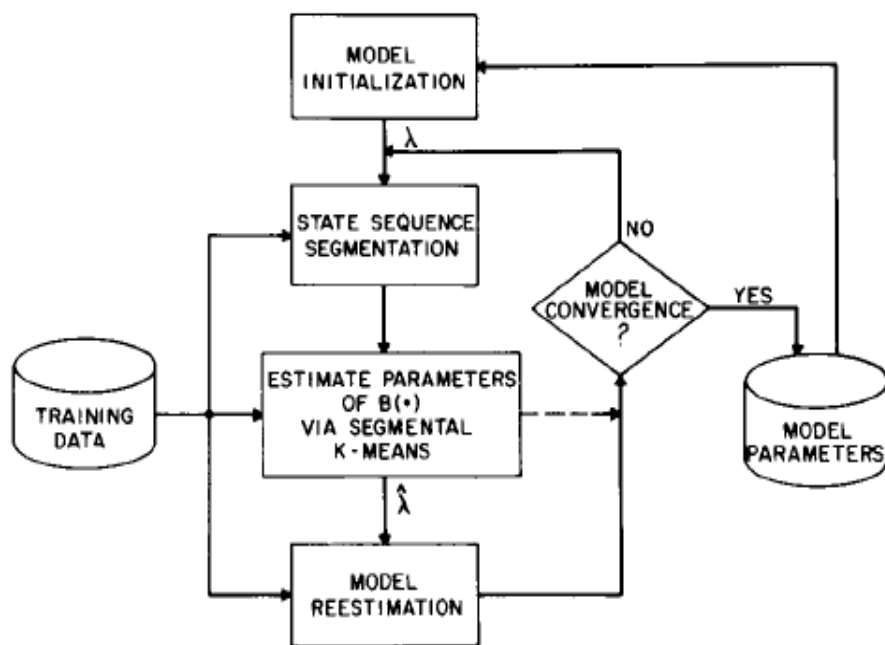
- Pronunciation of a single word and recognition of it



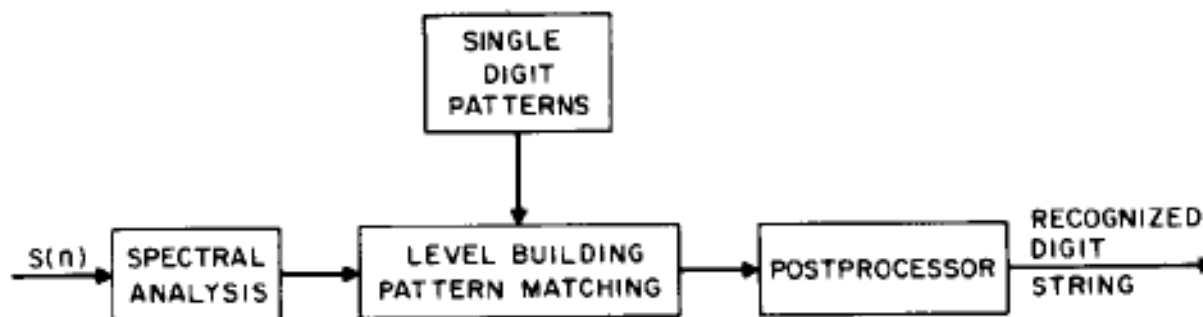
Key element:
Need to have lot of
domain knowledge to
extract useful features



- When continuous observations is needed (i.e. speech recognition), the initial estimation of the observation distribution is important for the convergence



- The observation sequence matched at each step with a single word recognition system



- The observation sequences are non pre-classified and we have no information about the ending of each word
- The building level match the observation sequence to a digit sequence with some probability
- An alternative way is to use the state segmentation with an higher level model



- HMMs are general stochastic models
- EM is a good algorithm to learn such models
- We need prior knowledge to define the structure of the model
- Lot of parameters needs lot of data
- They perform very well for many applications if they are applied in the correct way
 - Signal segmentation and classification
 - Clustering of signals
 - Prediction



Questions ?

