4.0 More about Hidden Markov Models

Reference: 1. 6.1-6.6, Rabiner and Juang

2. 4.4.1 of Huang

Markov Model

- Markov Model (Markov Chain)
 - First-order Markov chain of N states is a triplet (S,A,π)
 - S is a set of N states
 - A is the $N \times N$ matrix of state transition probabilities

 $P(q_t=j|q_{t-1}=i, q_{t-2}=k, ...)=P(q_t=j|q_{t-1}=i) \equiv \mathbf{a}_{ij}$

• π is the vector of initial state probabilities

 $\pi_j = P(q_0 = j)$

- The output for any given state is an observable event (deterministic)
- The output of the process is a sequence of observable events



A Markov chain with 5 states (labeled $S_1 \mbox{ to } S_5)$ with state transitions.

Markov Model

• An example : a 3-state Markov Chain λ

 State 1 generates symbol A *only*, State 2 generates symbol B **only**, and State 3 generates symbol C **only**

	0.6	0.3	0.1
$\mathbf{A} =$	0.1	0.7	0.2
	0.3	0.2	0.5
$\pi = [0.4]$		0.5	0.1]



- Given a sequence of observed symbols $\mathbf{O}=\{CABBCABC\}$, the **only** one corresponding state sequence is $\{S_3S_1S_2S_2S_3S_1S_2S_3\}$, and the corresponding probability is $P(\mathbf{O}|\lambda)=P(q_0=S_3)$ $P(S_1/S_3)P(S_2/S_1)P(S_2/S_2)P(S_3/S_2)P(S_1/S_3)P(S_2/S_1)P(S_3/S_2)$ $=0.1 \times 0.3 \times 0.3 \times 0.7 \times 0.2 \times 0.3 \times 0.3 \times 0.2 = 0.00002268$

Hidden Markov Model

- HMM, an extended version of Markov Model
 - The observation is a probabilistic function (discrete or continuous) of a state instead of an one-to-one correspondence of a state
 - The model is a doubly embedded stochastic process with an underlying stochastic process that is not directly observable (hidden)
 - What is hidden? *The State Sequence* According to the observation sequence, we never know which state sequence generates it
- Elements of an HMM {S,A,B,π}
 - S is a set of N states
 - A is the $N \times N$ matrix of state transition probabilities
 - B is a set of N probability functions, each describing the observation probability with respect to a state
 - $-\pi$ is the vector of initial state probabilities

Simplified HMM



RGBGGBBGRRR.....

Hidden Markov Model

- **Two types of HMM's according to the observation functions** <u>Discrete and finite observations :</u>
 - The observations that all distinct states generate are finite in number

$$\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_M\}, \mathbf{v}_k \in \mathbf{R}^D$$

- the set of observation probability distributions $B = \{b_j(\mathbf{v}_k)\}$ is defined as $b_j(\mathbf{v}_k) = P(\mathbf{o}_t = \mathbf{v}_k | \mathbf{q}_t = j), 1 \le k \le M, 1 \le j \le N$

 \mathbf{o}_t : observation at time t, \boldsymbol{q}_t : state at time t

 \Rightarrow for state *j*, $b_j(\mathbf{v}_k)$ consists of only *M* probability values

Continuous and infinite observations :

- The observations that all distinct states generate are infinite and continuous, $\mathbf{V} = \{ \mathbf{v} | \mathbf{v} \in \mathbf{R}^D \}$
- the set of observation probability distributions $B = \{b_j(\mathbf{v})\}$ is defined as $b_j(\mathbf{v}) = P(\mathbf{o}_t = \mathbf{v} | \mathbf{q}_t = j), 1 \le j \le N$

 $\Rightarrow b_j(\mathbf{v})$ is a continuous probability density function and is often assumed to be a mixture of Gaussian distributions

$$b_{j}(v) = \sum_{k=1}^{M} C_{jk} \left[\frac{1}{\sqrt{2\pi}^{D} |\Sigma_{jk}|^{\frac{1}{2}}} exp\left(-\frac{1}{2} \left(\left(v - \mu_{jk} \right)^{t} \Sigma_{jk}^{-1} \left(v - \mu_{jk} \right) \right) \right) \right] = \sum_{k=1}^{M} C_{jk} b_{jk}(v)$$

• An example : a 3-state discrete HMM λ

$$\mathbf{A} = \begin{bmatrix} 0.6 & 0.3 & 0.1 \\ 0.1 & 0.7 & 0.2 \\ 0.3 & 0.2 & 0.5 \end{bmatrix}$$

$$b_1(\mathbf{A}) = 0.3, b_1(\mathbf{B}) = 0.2, b_1(\mathbf{C}) = 0.5$$

$$b_2(\mathbf{A}) = 0.7, b_2(\mathbf{B}) = 0.1, b_2(\mathbf{C}) = 0.2$$

$$b_3(\mathbf{A}) = 0.3, b_3(\mathbf{B}) = 0.6, b_3(\mathbf{C}) = 0.1$$

$$\pi = \begin{bmatrix} 0.4 & 0.5 & 0.1 \end{bmatrix}$$



- Given a sequence of observations \overline{O} ={ABC}, there are **27 possible** corresponding state sequences, and therefore the corresponding probability is

$$P(\overline{\mathbf{O}}|\lambda) = \sum_{i=1}^{27} P(\overline{\mathbf{O}}, \mathbf{q}_i | \lambda) = \sum_{i=1}^{27} P(\overline{\mathbf{O}}|\mathbf{q}_i, \lambda) P(\mathbf{q}_i | \lambda), \quad \mathbf{q}_i : \text{state sequence}$$

e.g. when $\mathbf{q}_i = \{S_2 S_2 S_3\}, P(\overline{\mathbf{O}}|\mathbf{q}_i, \lambda) = P(\mathbf{A}|S_2) P(\mathbf{B}|S_2) P(\mathbf{C}|S_3) = 0.7 * 0.1 * 0.1 = 0.007$
$$P(\mathbf{q}_i | \lambda) = P(q_0 = S_2) P(S_2|S_2) P(S_3|S_2) = 0.5 * 0.7 * 0.2 = 0.07$$

Hidden Markov Model

Three Basic Problems for HMMs

Given an observation sequence $\overline{O} = (o_1, o_2, \dots, o_T)$, and an HMM

- $\lambda = (A, B, \pi)$
 - Problem 1:

How to *efficiently* compute $P(\overline{\mathbf{O}} | \lambda)$?

- \Rightarrow Evaluation problem
- Problem 2:

How to choose an optimal state sequence $\mathbf{q} = (q_1, q_2, \dots, q_T)$?

- \Rightarrow Decoding Problem
- Problem 3:

Given some observations \overline{O} for the HMM λ , how to adjust the model parameter $\lambda = (A, B, \pi)$ to maximize $P(\overline{O} | \lambda)$?

⇒ Learning /Training Problem

$$1 \rightarrow 2 \rightarrow \cdots \rightarrow N \quad \lambda = (A, B, \pi)$$

 $\overline{O} = O_1 O_2 O_3 \dots O_t \dots O_T$ observation sequence

 $\overline{q} = q_1 q_2 q_3 \dots q_t \dots q_T$ state sequence

- Problem 1: Given λ and \overline{O} , find $P(\overline{O}|\lambda)=Prob[observing \overline{O} \text{ given } \lambda]$
- Direct Evaluation: considering all possible state sequence \overline{q}

$$P(\overline{O}|\lambda) = \sum_{all \,\overline{q}} P(\overline{O}, \overline{q}|\lambda) = \sum_{all \,\overline{q}} P(\overline{O}|\overline{q}, \lambda) P(\overline{q}|\lambda)$$

$$P(\overline{O}|\overline{q}, \lambda)$$

$$\widehat{T}$$

$$P(\overline{O}|\lambda) = \sum_{all \,\overline{q}} ([b_{q_1}(o_1) \bullet b_{q_2}(o_2) \bullet \dots ... b_{q_T}(o_T)] \bullet [\pi_{q_1} \bullet a_{q_1q_2} \bullet a_{q_2q_3} \bullet \dots ... a_{q_{T-1}q_T}])$$

$$\square$$

$$P(\overline{q}|\lambda)$$
total number of different $\overline{q} : N^T$
huge computation requirements

• Forward Algorithm: defining a forward variable $\alpha_t(i)$

$$\alpha_t(i) = P(o_1 o_2 \dots o_t, q_t = i | \lambda)$$

=Prob[observing $o_1 o_2 \dots o_t$, state i at time t $|\lambda$]

- Initialization

$$\alpha_1(i) = \pi_i b_i(o_1) , \quad 1 \le i \le N$$

- Induction $\alpha_{t+1}(j) = \left[\sum_{i=1}^{N} \alpha_t(i)a_{ij}\right] b_j(o_{t+1})$ $1 \le j \le N$ $1 \le t \le T-1$
- Termination

$$P(\overline{O}|\lambda) = \sum_{i=1}^{N} \alpha_{T}(i)$$

See Fig. 6.5 of Rabiner and Juang

- All state sequences, regardless of how long previously, merge to the N state at each time instant t







Forward Algorithm

- **Problem 2:** Given λ and $\overline{O} = o_1 o_2 \dots o_T$, find a best state sequence $\overline{q} = q_1 q_2 \dots q_T$
- **Backward Algorithm :** defining a backward variable $\beta_t(i)$

$$\beta_{t}(i) = P(o_{t+1}, o_{t+2}, ..., o_{T} | q_{t} = i, \lambda)$$

- = Prob[observing o_{t+1} , o_{t+2} ,..., o_T |state i at time t, λ]
- Initialization

$$\beta_{\mathbf{T}}(i) = 1, \ 1 \le i \le N$$
 ($\beta_{\mathbf{T-1}}(i) = \sum_{j=1}^{N} a_{ij} b_j(o_{\mathbf{T}})$)

- Induction $\beta_{t}(i) = \sum_{j=1}^{N} a_{ij} b_{j}(o_{t+1})\beta_{t+1}(j)$ $t = T-1, T-2, \dots, 2, 1, \qquad 1 \le i \le N$ See Fig. 6.6 of Rabiner and Juang
- Combining Forward/Backward Variables

 $P(\overline{O}, q_t = i | \lambda)$

= Prob [observing $o_1, o_2, ..., o_t, ..., o_T, q_t = i |\lambda]$

$$= \alpha_{t}(i)\beta_{t}(i)$$

$$P(\overline{O}|\lambda) = \sum_{i=1}^{N} P(\overline{O}, q_{t} = i |\lambda) = \sum_{i=1}^{N} [\alpha_{t}(i)\beta_{t}(i)]$$



$$B_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = i, \lambda)$$



 $\beta_{t}(i) = \sum_{j=1}^{N} a_{ij} b_{j}(o_{t+1})\beta_{t+1}(j)$ $t = T-1, T-2, ..., 2, 1, 1 \le i \le N$

Backward Algorithm



(*i*) $P(\overline{O}, q_t = i | \lambda)$ $= Prob [observing o_1, o_2, ..., o_t, ..., o_T, q_t = i | \lambda]$ $= \alpha_t(i)\beta_t(i)$

$$A: (o_1 \ o_2 \ \cdots \ o_t | \lambda)$$

$$B: (o_{t+1}, o_{t+2}, \cdots \ o_T | \lambda)$$

$$C: (q_t = i | \lambda)$$

$$P(A, B, C) = P(A, C) P(B | A, C)$$

$$// \qquad // \qquad (B \perp A)$$

$$P(\overline{O}, q_t = i | \lambda) \ \alpha_t(i) \ P(B | C)$$

$$// \qquad \beta_t(i)$$



- Approach 1 Choosing state q_t * individually as the most likely state at time t
 - Define a new variable $\gamma_t(i) = P(q_t = i | \overline{O}, \lambda)$

$$\gamma_{t}(i) = \frac{\alpha_{t}(i)\beta_{t}(i)}{\sum_{i=1}^{N} \alpha_{t}(i)\beta_{t}(i)} = \frac{P(O, q_{t}=i|\lambda)}{P(O|\lambda)}$$

- Solution

$$q_t^* = \arg \max_{1 \le i \le N} [\forall t(i)], \ 1 \le t \le T$$

in fact

$$q_t^* = \arg \max_{1 \le i \le N} [P(\overline{O}, q_t = i | \lambda)]$$

 $= \arg \max_{1 \le i \le N} [\alpha_t(i)\beta_t(i)]$

- Problem

maximizing the probability at each time t individually

$$\overline{q}^* = q_1^* q_2^* \dots q_T^*$$
 may not be a valid sequence

(e.g.
$$a_{q_t * q_{t+1} *} = 0$$
)

- Approach 2 —Viterbi Algorithm finding the single best sequence $\overline{q}^* = q_1^* q_2^* \dots q_T^*$
 - Define a new variable $\delta_t(i)$

$$\delta_{t}(i) = \max_{q_{1},q_{2},\dots,q_{t-1}} P[q_{1},q_{2},\dots,q_{t-1}, q_{t} = i, o_{1},o_{2},\dots,o_{t} | \lambda]$$

= the highest probability along a certain single path ending at state i at time t for the first t observations, given λ

- Induction

 $\delta_{t+1}(j) = \max_{i} \left[\delta_t(i) a_{ij} \right] \bullet b_j(o_{t+1})$

- Backtracking $\psi_{t+1}(j) = \arg \max_{1 \le i \le N} [\delta_t(i)a_{ij}]$

the best previous state at t–1 given at state j at time t keeping track of the best previous state for each j and t

Viterbi Algorithm



 $\delta_{t}(i) = \max_{q_{1},q_{2},...,q_{t-1}} P[q_{1},q_{2},...,q_{t-1}, q_{t} = i, o_{1},o_{2},...,o_{t} |\lambda]$

Viterbi Algorithm



- Complete Procedure for Viterbi Algorithm
 - Initialization

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

- Recursion

$$\begin{split} \delta_{t+1}(j) &= \max_{1 \leq i \leq N} \left[\delta_t(i) a_{ij} \right] \bullet b_j(o_{t+1}) \\ &1 \leq t \leq T\text{-}1, \quad 1 \leq j \leq N \\ \psi_{t+1}(j) &= \arg\max_{1 \leq i \leq N} \left[\delta_t(i) a_{ij} \right] \\ &1 \leq t \leq T\text{-}1, \quad 1 \leq j \leq N \end{split}$$

- Termination

 $P^* = \max_{1 \le i \le N} [\delta_T(i)]$ $q_T^* = \arg \max_{1 \le i \le N} [\delta_T(i)]$

- Path backtracking

 $q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, t-2, \dots, 2, 1$

• Application Example of Viterbi Algorithm

- Isolated word recognition

$$\lambda_0 = (\mathbf{A}_0, \mathbf{B}_0, \boldsymbol{\pi}_0)$$
$$\lambda_1 = (\mathbf{A}_1, \mathbf{B}_1, \boldsymbol{\pi}_1)$$
$$\vdots$$
$$\lambda_n = (\mathbf{A}_n, \mathbf{B}_n, \boldsymbol{\pi}_n)$$

observation

 $\overline{O} = (o_1, o_2, \dots o_T)$ $k^* = \arg \max_{1 \le i \le n} P[\overline{O} \mid \lambda_i] \approx \arg \max_{1 \le i \le n} [P^* \mid \lambda_i]$ $\widehat{\Box}$ Basic Problem 1
Basic Problem 1
Forward Algorithm
(for all paths)
(for a single best path)

-The model with the highest probability for the most probable path usually also has the highest probability for all possible paths.

- **Problem 3:** Give \overline{O} and an initial model $\lambda = (A, B, \pi)$, adjust λ to maximize $P(\overline{O}|\lambda)$
 - Baum-Welch Algorithm (Forward-backward Algorithm)
 - Define a new variable

$$\begin{split} \boldsymbol{\epsilon}_{t}(i,j) &= P(q_{t} = i, q_{t+1} = j \mid \overline{O}, \lambda) \\ &= \frac{\alpha_{t}(i) a_{ij} b_{j}(o_{t+1})\beta_{t+1}(j)}{\sum\limits_{i=1}^{N} \sum\limits_{j=1}^{N} [\alpha_{t}(i)a_{ij} b_{j}(o_{t+1})\beta_{t+1}(j)]} \\ &= \frac{Prob[\overline{O}, q_{t} = i, q_{t+1} = j \mid \lambda]}{P(\overline{O} \mid \lambda)} \end{split}$$

See Fig. 6.7 of Rabiner and Juang

- Recall $\gamma_t(i) = P(q_t = i \mid \overline{O}, \lambda)$

 $\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of times that state i is visited in } \overline{O} \text{ from } t = 1 \text{ to } t = T-1$

= expected number of transitions from state i in O $\sum_{t=1}^{T-1} \varepsilon_t(i, j) = \text{expected number of transitions from state i to state j in }\overline{O}$

$\alpha_t(i) a_{ij} b_j(o_{t+1})\beta_{t+1}(j)$





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

$$\varepsilon_{t}(i,j) = \frac{\alpha_{t}(i) \ a_{ij} \ b_{j}(o_{t+1}) \ \beta_{t+1}(j)}{\sum_{j=1}^{N} \sum_{i=1}^{N} \alpha_{t}(i) \ a_{ij} \ b_{j}(o_{t+1}) \ \beta_{t+1}(j)}$$
$$= \frac{P(\bar{O}, q_{t} = i, q_{t+1} = j | \lambda)}{P(\bar{O} | \lambda)} = P(q_{t} = i, q_{t+1} = j | \bar{O}, \lambda)$$



- Results

Continuous Density HMM

$$b_{j}(o) = \sum_{k=1}^{M} c_{jk} N(o; \mu_{jk}, U_{jk})$$

N(): Multi-variate Gaussian

 μ_{jk} : mean vector for the k-th mixture component

 U_{jk} : covariance matrix for the k-th mixture component $\sum_{k=1}^{M} c_{jk} = 1$ for normalization

• Continuous Density HMM

- Define a new variable
 - $\gamma_t(j, k) = \gamma_t(j)$ but including the probability of o_t evaluated in the k-th mixture component out of all the mixture components

$$= \left[\frac{\alpha_{t}(j)\beta_{t}(j)}{\sum\limits_{j=1}^{N} \alpha_{t}(j)\beta_{t}(j)}\right] \left[\frac{c_{jk} N(o_{t}; \mu_{jk}, U_{jk})}{\sum\limits_{m=1}^{M} c_{jm} N(o_{t}; \mu_{jm}, U_{jm})]}\right]$$

- Results

$$\overline{c}_{jk} = \frac{\sum_{t=1}^{T} \gamma_t(j, k)}{\sum_{t=1}^{T} \sum_{k=1}^{M} \gamma_t(j, k)}$$

See Fig. 6.9 of Rabiner and Juang

• Continuous Density HMM

$$\overline{\mu}_{jk} = \frac{\sum_{t=1}^{T} [\gamma_t(j, k) \bullet o_t]}{\sum_{t=1}^{T} \gamma_t(j, k)}$$
$$\overline{U}_{jk} = \frac{\sum_{t=1}^{T} [\gamma_t(j, k)(o_t - \mu_{jk}) (o_t - \mu_{jk})']}{\sum_{t=1}^{T} \gamma_t(j, k)}$$

• Iterative Procedure

$$\lambda = (A, B, \pi) \xrightarrow{} \overline{\lambda} = (\overline{A}, \overline{B}, \overline{\pi})$$

$$\overrightarrow{O} = o_1 o_2 \dots o_T$$

- It can be shown (by EM Theory (or EM Algorithm)) $P(\overline{O}|\overline{\lambda}) \ge P(\overline{O}|\lambda)$ after each iteration

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

$$\int_{-\infty}^{\infty} x \, f_X(x) \, dx = \bar{x}$$

$$\overline{U}_{jk} = \frac{\sum_{t=1}^{T} [\gamma_t(j,k) (o_t - \mu_{jk})(o_t - \mu_{jk})']}{\sum_{t=1}^{T} \gamma_t(j,k)} \int_{-\infty}^{\infty} (x - \overline{x})^2 f_X(x) dx = \sigma_X^2$$

 $f_X(x)$: prob. density function

$$\int \sigma_{x}^{2} k \left[\begin{array}{c} o_{t_{1}} - \mu_{jk_{1}} \\ o_{t_{2}} - \mu_{jk_{2}} \\ \vdots \\ o_{t_{D}} - \mu_{jk_{D}} \end{array} \right] \left[o_{t_{1}} - \mu_{jk_{1}} o_{t_{2}} - \mu_{jk_{2}} \cdots o_{t_{D}} - \mu_{jk_{D}} \right]$$

$$\overline{U} = \begin{bmatrix} m \\ \vdots \end{bmatrix} = E\left(\begin{bmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} [x_1 - \bar{x}_1, x_2 - \bar{x}_2, \cdots]\right)$$

 $\bar{u}_{lm} = E[(x_l - \bar{x}_l)(x_m - \bar{x}_m)]$

- No closed-form solution, but approximated iteratively
- An initial model is needed-model initialization
- May converge to local optimal points rather than global optimal point
 - heavily depending on the initialization
- Model training





Vector Quantization (VQ)

- An Efficient Approach for Data Compression
 - replacing a set of real numbers by a finite number of bits
- An Efficient Approach for Clustering Large Number of Sample Vectors
 - grouping sample vectors into clusters, each represented by a single vector (codeword)

Scalar Quantization

- replacing a single real number by an R-bit pattern
- a mapping relation



$$S = \bigcup_{k=1}^{L} J_k, V = \{ v_1, v_2, ..., v_L \}$$

O:S \rightarrow V

$$Q(x[n]) = v_k \text{ if } x[n] \in J_k$$
$$L = 2^R$$

Each v_k represented by an R-bit pattern

- -Quantization characteristics (codebook) { J₁, J₂, ..., J_L } and { v₁, v₂, ..., v_L } designed considering at least 1. error sensitivity
 - 2. probability distribution of x[n]

Vector Quantization

Scalar Quantization : Pulse Coded Modulation (PCM)



Vector Quantization



Vector Quantization (VQ)

2-dim Vector Quantization (VQ)

Example: $\overline{x}_n = (x[n], x[n+1])$ $S = \{\overline{x}_n = (x[n], x[n+1]); |x[n]| < A, |x[n+1]| < A\}$ •VQ S divided into L 2 dim ragions

- S divided into L 2-dim regions { J_1 , J_2 , ..., J_k , ..., J_L } S = $\bigcup_{k=1}^{L} J_k$

each with a representative vector $\overline{v}_k \in J_k, V = \{\overline{v}_1, \overline{v}_2, ..., \overline{v}_L\}$ $-Q : S \rightarrow V$ $Q(\overline{x}_n) = \overline{v}_k \text{ if } \overline{x}_n \in J_k$ $L = 2^R$

each \overline{v}_k represented by an R-bit pattern

– Considerations

- 1.error sensitivity may depend on x[n], x[n+1] jointly
- 2.distribution of x[n] , x[n+1] may be correlated statistically
- 3.more flexible choice of J_k
- Quantization Characteristics (codebook)

{ J₁, J₂, ..., J_L } and { $\overline{v_1}$, $\overline{v_2}$, ..., $\overline{v_L}$ }

Vector Quantization A A -A-A



$$(256)^2 = (2^8)^2 = 2^{16}$$

$$1024=2^{10}$$

Vector Quantization



Vector Quantization (VQ)

N-dim Vector Quantization $\overline{\mathbf{x}} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ $\mathbf{S} = \{ \overline{\mathbf{x}} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) ,$ $|x_{k}| < A, k = 1, 2, ... N$ $\mathbf{S} = \mathbf{U}_{\mathbf{M}} \mathbf{J}_{\mathbf{k}}$ $\mathbf{V} = \{\overline{\mathbf{v}}_1, \overline{\mathbf{v}}_2, \dots, \overline{\mathbf{v}}_L\}$ $Q: S \rightarrow V$ $Q(\overline{x}) = \overline{v}_k$ if $\overline{x} \in J_k$ $L = 2^{\mathbf{R}}$, each $\overline{v}_{\mathbf{k}}$ represented by an R-bit pattern

Codebook Trained by a Large **Training** Set [•]Define distance measure between two vectors x, y $d(\overline{x}, \overline{y}) : S \times S \rightarrow R^+$ (non-negative real numbers) -desired properties $d(\overline{x},\overline{y}) \ge 0$ $d(\overline{x}, \overline{x}) = 0$ $d(\overline{x}, \overline{y}) = d(\overline{y}, \overline{x})$ $d(\overline{x}, \overline{y}) + d(\overline{y}, \overline{z}) \ge d(\overline{x}, \overline{z})$ examples : $d(\overline{x}, \overline{y}) = \sum_{i} (x_i - y_i)^2$ $d(\overline{x}, \overline{y}) = \sum_{i} |x_i - y_i|$ d($\overline{x}, \overline{y}$) = $(\overline{x} - \overline{y})^t \sum^{-1} (\overline{x} - \overline{y})$ Mahalanobis Distance Σ : Co-variance Matrix

Distance Measures



 $d(\bar{x}, \bar{y}) = (\bar{x} - \bar{y})^t \Sigma^{-1} (\bar{x} - \bar{y}) \quad \text{Mahalanobis distance}$ $\sum = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}, d(\bar{x}, \bar{y}) = \sum_i (x_i - y_i)^2$ $\sum = \begin{bmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \sigma_2^2 & \vdots \\ 0 & \cdots & \ddots & \sigma_n^2 \end{bmatrix}, d(\bar{x}, \bar{y}) = \sum_i \frac{(x_i - y_i)^2}{\sigma_i^2}$

Vector Quantization (VQ)

K-Means Algorithm/Lloyd-Max Algorithm



(1)
$$J_{k} = \{ \overline{x} \mid d(\overline{x}, \overline{v}_{k}) < d(\overline{x}, \overline{v}_{j}), j \neq k \}$$

 $\rightarrow D = \sum_{all \overline{x}} d(\overline{x}, Q(\overline{x})) = min$
nearest neighbor condition
(2) For each k
 $\overline{v}_{k} = \frac{1}{M} \sum_{x \in V} \overline{x}$

$$\rightarrow D_{\mathbf{k}} = \sum_{\overline{\mathbf{x}} \in \mathbf{J}_{\mathbf{k}}}^{\mathbf{M}} \mathbf{d}(\overline{\mathbf{x}}, \overline{\mathbf{v}}_{\mathbf{k}}) = \min$$

centroid condition

(3) Convergence condition $D = \sum_{k=1}^{L} D_{k}$

> after each iteration D is reduced, but $D \ge 0$ $| D^{(m+1)} - D^{(m)} | < \in, m : iteration$

• Iterative Procedure to Obtain Codebook from a Large Training Set



Vector Quantization (VQ)

• K-means Algorithm may Converge to Local Optimal Solutions

- depending on initial conditions, not unique in general

Training VQ Codebook in Stages— LBG Algorithm

- step 1: Initialization. L = 1, train a 1-vector VQ codebook

$$\overline{\mathbf{v}} = \frac{1}{N} \sum_{j} \overline{\mathbf{x}}_{j}$$

- step 2: Splitting.

Splitting the L codewords into 2L codewords, L = 2L

- example 1 $\overline{v}_{k}^{(1)} = \overline{v}_{k}(1+\varepsilon)$ $\overline{v}_{k}^{(2)} = \overline{v}_{k}(1-\varepsilon)$ • example 2 $\overline{v}_{k}^{(1)} = \overline{v}_{k}$ $\overline{v}_{k}^{(2)} : \text{the vector most}$ far apart
- step 3: K-means Algorithm: to obtain L-vector codebook
- step 4: Termination. Otherwise go to step 2
- Usually Converges to Better Codebook

LBG Algorithm



Initialization in HMM Training

• An Often Used Approach— Segmental K-Means

 Assume an initial estimate of all model parameters (e.g. estimated by segmentation of training utterances into states with equal length)

•For discrete density HMM

 $b_j(k) = \frac{\text{number of vectors in state } j \text{ associated with codeword } k}{\text{total number of vectors in state } j}$

•For continuous density HMM (M Gaussian mixtures per state)

 \Rightarrow cluster the observation vectors within each state j into a set of M clusters

(e.g. with vector quantiziation)

 c_{jm} = number of vectors classified in cluster m of state j

divided by number of vectors in state j

 μ_{jm} = sample mean of the vectors classified in cluster m of state j

 \sum_{jm} = sample covariance matrix of the vectors classified in cluster m of state j

- Step 1 : re-segment the training observation sequences into states based on the initial model by Viterbi Algorithm
- Step 2 : Reestimate the model parameters (same as initial estimation)
- Step 3: Evaluate the model score $P(\overline{O}|\lambda)$:

If the difference between the previous and current model scores exceeds a threshold, go back to Step 1, otherwise stop and the initial model is obtained

Segmental K-Means



Initialization in HMM Training

• An example for Continuous HMM

- 3 states and 4 Gaussian mixtures per state



Initialization in HMM Training

- An example for discrete HMM
 - 3 states and 2 codewords

