

ECE 8873
Data Compression and Modeling

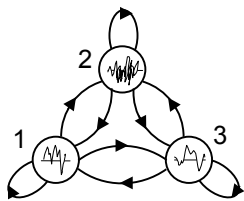
Lecture 12:
Hidden Markov Model

School of Electrical and Computer Engineering
Georgia Institute of Technology
Spring, 2004

Non-stationarity & Non-memoryless

- Many real world signals have characteristics that change with time – it is what makes them interesting and “informative.”
- Recall the structural component of information – it often changes at irregular times.
- Finite state Markov chain is one of the simplest process that allows us to model varying characteristics with “memory.”
- A simple model that can be used to approximate non-stationary, non-memoryless sources is the hidden Markov model.

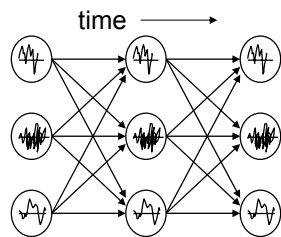
Hidden Markov Models



$$P(X | \Lambda) = \sum_{\mathbf{q}} P(X, \mathbf{q} | \Lambda)$$

$$P(X, \mathbf{q} | \Lambda) = \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$$

$$X = (x_1, x_2, \dots, x_T) \quad \mathbf{q} = (q_0, q_1, q_2, \dots, q_T)$$



- Each state represents a process of measurable observations.
- Inter-process transition is governed by a finite state Markov chain.
- Processes are stochastic and individual observations do not immediately identify the state.

Hidden Markov Models - Specifications

$X = (x_1, x_2, \dots, x_T)$ is the sequence of observations
 $\mathbf{q} = (q_0, q_1, \dots, q_T)$ is the sequence of states the system is in

- Number of states of the Markov chain, N
- State transition probability matrix, $A = [a_{ij}]_{N \times N}$
 $a_{ij} = \Pr[q_t = j | q_{t-1} = i] \quad \sum_{j=1}^N a_{ij} = 1$ for all i
- In-state observation probability distribution functions
 $B = \{b_i(x)\}_{i=1}^N \quad b_i(x) \Rightarrow b_i(x, \lambda_i)$ i.e., parameterized by λ_i
- Initial state probability distribution,
 $\boldsymbol{\pi}^t = [\pi_1, \pi_2, \dots, \pi_N]$ where $\pi_i = \Pr[q_0 = i] \quad \sum_{i=1}^N \pi_i = 1$

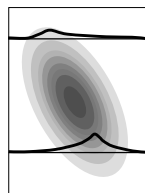
The triple $\Lambda = (\boldsymbol{\pi}, A, B)$ defines a hidden Markov model.

In-State (Local) Observation Distributions

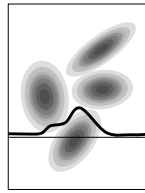
- Discrete distributions
 $x \in \{s_k\}_{k=1}^M$ and $b_i(x = s_k) = \Pr\{x = s_k, q = i\} = b_{ik}$
- Log-concave probability density functions
 x is continuous - valued and $\log b_i(x) = \log f_i(x)$ is a concave function
- Elliptically symmetric probability density functions
 $b_i(x) = \int f(x, g) d\mu(g)$
- General mixture probability density functions

$$b_i(x) = \sum_{k=1}^M c_{ik} f_{ik}(x)$$

where $\sum_{k=1}^M c_{ik} = 1$



Elliptically Symmetric Distribution



Mixture of Elliptically Symmetric Distribution

marginal

Three Basic Problems of HMM

- Given the observation sequence $X = (x_1, x_2, \dots, x_T)$ and a model $\Lambda = (\pi, A, B)$, how do we efficiently compute $P(X | \Lambda)$?
- Given the observation sequence $X = (x_1, x_2, \dots, x_T)$ and the model $\Lambda = (\pi, A, B)$, how do we find a corresponding state sequence $\mathbf{q} = (q_0, q_1, q_2, \dots, q_T)$ that is optimal in some sense?
- Given an observation sequence X , or a number of sequences $\{X^{(i)}\}_i$, how to estimate parameters in the model set $\Lambda = (\pi, A, B)$?

Evaluation of HMM Probability

$$P(X | \Lambda) = \sum_{\mathbf{q}} P(X, \mathbf{q} | \Lambda) \quad P(X, \mathbf{q} | \Lambda) = \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$$

$$P(X | \Lambda) = \sum_{\mathbf{q}} P(X, \mathbf{q} | \Lambda) = \sum_{\mathbf{q}} \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$$

Direct evaluation will involve $2T \cdot N^T$ calculations.

The Forward Procedure

Define $\alpha_t(i) = \Pr(x_1, x_2, \dots, x_t, q_t = i | \Lambda)$

$\alpha_t(i)$ is the probability of the partial observation sequence x_1, x_2, \dots, x_t , up to time t , and the system is at state i at time t .

Forward Procedure

Initialization: $\alpha_0(i) = \pi_i, \quad i = 1, 2, \dots, N$

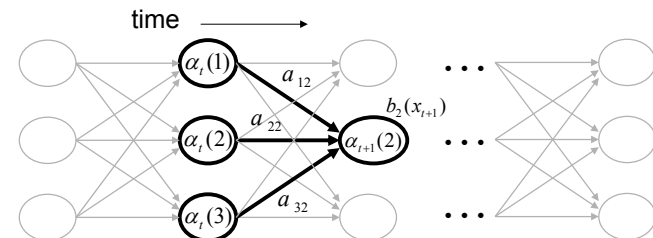
Induction:

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(x_{t+1}), \quad 0 \leq t \leq T-1, 1 \leq j \leq N$$

Termination:

$$P(X | \Lambda) = \sum_{i=1}^N \alpha_T(i)$$

$\sim N^2 T$ calculations



Backward Procedure

Define $\beta_t(i) = \Pr\{(x_{t+1}, x_{t+2}, \dots, x_T | q_t = i, \Lambda)$

$\beta_t(i)$ is the probability of the partial observation sequence $x_{t+1}, x_{t+2}, \dots, x_T$, given that the system is at state i at time t .

Initialization: $\beta_T(i) = 1, \quad i = 1, 2, \dots, N$

Induction:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(x_{t+1}) \beta_{t+1}(j), \quad t = T-1, T-2, \dots, 1, \quad 1 \leq i \leq N$$

Optimal State Sequence

- Several possibilities

- The state sequence that maximizes the joint state-observation probability

$$\mathbf{q}_{opt} = \arg \max_{\mathbf{q}} P(X, \mathbf{q} | \Lambda) = \arg \max_{\mathbf{q}} \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$$

- The state sequence that consists of individual states maximizing the a posteriori probability given the observation

$$\gamma_t(i) = P(q_t = i | X, \Lambda) \quad \text{i.e. probability of being in state } i \text{ at time } t, \text{ given } X$$

$$\begin{aligned} \gamma_t(i) &= P(q_t = i | X, \Lambda) = P(X, q_t = i | \Lambda) [P(X | \Lambda)]^{-1} \\ &= P(X, q_t = i | \Lambda) \left[\sum_{i=1}^N P(X, q_t = i | \Lambda) \right]^{-1} \\ &= \alpha_t(i) \beta_t(i) \left[\sum_{j=1}^N \alpha_t(j) \beta_t(j) \right]^{-1} \end{aligned}$$

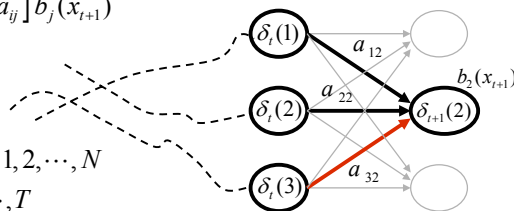
The Viterbi Algorithm

Find $\mathbf{q}^* = \arg \max_{\mathbf{q}} P(X, \mathbf{q} | \Lambda) = \arg \max_{\mathbf{q}} \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$

$$\delta_t(i) \equiv \max_{q_1, q_2, \dots, q_{t-1}} P\{q_1, q_2, \dots, q_{t-1}, [q_t = i], x_1, x_2, \dots, x_t | \Lambda\}$$

$$\delta_{t+1}(j) \equiv \max_{q_1, q_2, \dots, q_t} P\{q_1, q_2, \dots, q_t, [q_{t+1} = j], x_1, x_2, \dots, x_{t+1} | \Lambda\}$$

$$\begin{aligned} &= \max_i \left\{ \max_{q_1, q_2, \dots, q_{t-1}} P\{q_1, q_2, \dots, q_{t-1}, [q_t = i], x_1, x_2, \dots, x_t | \Lambda\} a_{ij} \right\} b_j(x_{t+1}) \\ &= \max_i \left[\delta_t(i) a_{ij} \right] b_j(x_{t+1}) \end{aligned}$$



Do for all $j = 1, 2, \dots, N$
and $t = 1, 2, \dots, T$

The Viterbi Algorithm

- Initialization

$$\delta_0(i) = \pi_i, \quad i = 1, 2, \dots, N$$

- Recursion

$$\begin{aligned} \delta_t(j) &= \max_i \left[\delta_{t-1}(i) a_{ij} \right] b_j(x_t) \quad j = 1, 2, \dots, N \\ \psi_t(j) &= \arg \max_i \left[\delta_{t-1}(i) a_{ij} \right] \quad t = 1, 2, \dots, T \end{aligned}$$

- Termination

$$\begin{aligned} P^* &= \max_i \left[\delta_T(i) \right] \\ \mathbf{q}^* &= \arg \max_i \left[\delta_T(i) \right] \end{aligned}$$

- Backtracking

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

Parameter Estimation

- Maximum likelihood (ML)
- Maximum mutual information (MMI)
- Minimum discrimination information (MDI)
- Minimum classification error (MCE)

Maximum Likelihood:

Given X , find $\bar{\Lambda} = \arg \max_{\Lambda} P(X | \Lambda)$

Or, given $\{X^{(i)}\}_{i=1}^N$ find

$$\bar{\Lambda} = \arg \max_{\Lambda} P(\{X^{(i)}\}_{i=1}^N | \Lambda)$$

X

Baum-Welch Re-estimation

Define $Q(\Lambda, \Lambda') = \sum_{\mathbf{q}} P(X, \mathbf{q} | \Lambda) \log P(X, \mathbf{q} | \Lambda')$ Q = Auxiliary function

Theorem: If $Q(\Lambda, \Lambda') \geq Q(\Lambda, \Lambda)$ then $P(X | \Lambda') \geq P(X | \Lambda)$. The inequality is strict unless $P(X, \mathbf{q} | \Lambda') \geq P(X, \mathbf{q} | \Lambda)$ almost everywhere.

Reestimation:

1. Given Λ , define the auxiliary function as a function of Λ' ;
2. Maximize the auxiliary function over Λ' and obtain $\bar{\Lambda}$

$$\bar{\Lambda} = T(\Lambda) \in \Psi = \{\hat{\Lambda} | Q(\Lambda, \hat{\Lambda}) = \max_{\Lambda'} Q(\Lambda, \Lambda')\}$$
3. Replace Λ with $\bar{\Lambda}$ and repeat the above until a stationary point is reached.

A general hill-climbing algorithm; similar to EM (expectation-maximization) algorithm

Reestimation Transformation

$$\bar{\Lambda} = T(\Lambda) \in \Psi = \{\hat{\Lambda} | Q(\Lambda, \hat{\Lambda}) = \max_{\Lambda'} Q(\Lambda, \Lambda')\}$$

For Gaussian mixture density HMM: $b_i(x) = \sum_{k=1}^M c_{ik} f(x; \mu_{ik}, \Sigma_{ik})$

Initial state probability: $\bar{\pi}_i = P(X, q_0 = i | \Lambda) [P(X | \Lambda)]^{-1}$

State transition probability:

$$\bar{a}_{ij} = \sum_{t=1}^T P(X, q_{t-1} = i, q_t = j | \Lambda) \left[\sum_{t=1}^T P(X, q_{t-1} = i | \Lambda) \right]^{-1}$$

Mixture weights:

$$\bar{c}_{ik} = \sum_{t=1}^T P(X, q_t = i, u_t = k | \Lambda) \left[\sum_{t=1}^T P(X, q_t = i | \Lambda) \right]^{-1}$$

Gaussian parameters:

$$\bar{\mu}_{ik} = \sum_{t=1}^T x_t P(X, q_t = i, u_t = k | \Lambda) \left[\sum_{t=1}^T P(X, q_t = i, u_t = k | \Lambda) \right]^{-1}$$

$$\bar{\Sigma}_{ik} = \sum_{t=1}^T (x_t - \mu_{ik})(x_t - \mu_{ik})' P(X, q_t = i, u_t = k | \Lambda) \left[\sum_{t=1}^T P(X, q_t = i, u_t = k | \Lambda) \right]^{-1}$$

Forward & Backward Probability

$$\alpha_t(i) = \Pr\{x_1, x_2, \dots, x_t, q_t = i | \Lambda\} \quad \beta_t(i) = \Pr\{x_{t+1}, x_{t+2}, \dots, x_T | q_t = i, \Lambda\}$$

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(x_{t+1}), \quad 0 \leq t \leq T-1, 1 \leq j \leq N$$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(x_{t+1}) \beta_{t+1}(j), \quad t = T-1, T-2, \dots, 1, \quad 1 \leq i \leq N$$

$$\gamma_t(i, k) = P(x_1, x_1, \dots, x_t, q_t = i, u_t = k | \Lambda)$$

$$= \sum_{j=1}^N \alpha_{t-1}(j) a_{ji} c_{jk} f_{jk}(x_t), \quad t = 1, 2, \dots, T, \quad 1 \leq i \leq N, 1 \leq k \leq M$$

$$P(X, q_t = i | \Lambda) = \alpha_t(i) \beta_t(i) \quad P(X | \Lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i) = \sum_{i=1}^N \alpha_T(i)$$

$$P(X, q_{t-1} = i, q_t = j | \Lambda) = \alpha_{t-1}(i) a_{ij} \left[\sum_{k=1}^M c_{jk} f_{jk}(x_t) \right] \beta_t(j)$$

$$P(X, q_t = i, u_t = k | \Lambda) = \sum_{j=1}^N \alpha_{t-1}(j) a_{ji} c_{jk} f_{jk}(x_t) \beta_t(i)$$

Reestimation formulas can be expressed all in terms of the forward & backward probabilities.

Hidden Markov Model

- A reasonable modeling tool
- Useful in extracting the structural component of the information (signal)
- Can be used in coding schemes (e.g., finite state coding schemes)
- Certainly useful in speech recognition & language processing