# CS578- Speech Signal Processing
## Lecture 11: Hidden Markov Models, HMM

Yannis Stylianou

University of Crete, Computer Science Dept., Multimedia Informatics Lab
yannis@csd.uoc.gr

Univ. of Crete

# OUTLINE

- a *Markov chain* or *process* is a sequence of events, usually called *states*, $Q = \{q_1, \cdots q_K\}$), the probability of each of which is dependent only on the event immediately preceding it.

- a *Hidden Markov Model* (HMM) represents stochastic sequences as Markov chains where the states are not directly observed, but are associated with a probability density function (pdf)

# FORMULAS AND DEFINITIONS

$Q = \{ \cdots \}$

- a *Markov chain* or *process* is a sequence of events, usually called *states*, $Q = \{q_1, \cdots q_K\}$), the probability of each of which is dependent only on the event immediately preceding it.

- a *Hidden Markov Model* (HMM) represents stochastic sequences as Markov chains where the states are not directly observed, but are associated with a probability density function (pdf)

- The generation of a random sequence in HMM is the result of a random walk in the chain (i.e. the browsing of a random sequence of states $Q = \{q_1, \cdots q_K\}$) and of a draw (called an *emission*) at each visit of a state.

- In pattern recognition (and speech recognition) with HMMs, we are interested to associate a sequence of states $Q = \{q_1, \cdots q_K\}$ to a sequence of observations $X = \{x_1, \cdots x_K\}$).

- The true sequence of states is therefore *hidden* by a first layer of stochastic processes.

- The generation of a random sequence in HMM is the result of a random walk in the chain (i.e. the browsing of a random sequence of states $Q = \{q_1, \cdots q_K\}$) and of a draw (called an *emission*) at each visit of a state.
- In pattern recognition (and speech recognition) with HMMs, we are interested to associate a sequence of states $Q = \{q_1, \cdots q_K\}$ to a sequence of observations $X = \{x_1, \cdots x_K\}$).
- The true sequence of states is therefore *hidden* by a first layer of stochastic processes.

# FORMULAS AND DEFINITIONS

- The generation of a random sequence in HMM is the result of a random walk in the chain (i.e. the browsing of a random sequence of states $Q = \{q_1, \cdots q_K\}$) and of a draw (called an *emission*) at each visit of a state.

- In pattern recognition (and speech recognition) with HMMs, we are interested to associate a sequence of states $Q = \{q_1, \cdots q_K\}$ to a sequence of observations $X = \{x_1, \cdots x_K\}$.

- The true sequence of states is therefore *hidden* by a first layer of stochastic processes.

# HMM TERMINOLOGY

- *Emission probabilities*: are the pdfs (usually Gaussians) that characterize each state $q_i$, i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$.

- *Transition probabilities*: are the probability to go from a state $i$ to a state $j$, i.e. $P(q_j|q_i)$. They are stored in matrices where each term $a_{ij}$ denotes a probability $P(q_j|q_i)$.

- *Non-emitting initial and final states*: For a finite length random sequence, two additional states are used in order to model the "start" or "end" events. These states are not associated with some emission probabilities.

- *Initial state distribution $P(I|q_j)$*: Transitions starting from the initial state.

- *Final-absorbent state*: The final state usually has only one non-null transition that loops onto itself with a probability of 1

# HMM TERMINOLOGY

- *Emission probabilities* : are the pdfs (usually Gaussians) that characterize each state $q_i$, i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$.

- *Transition probabilities* : are the probability to go from a state $i$ to a state $j$, i.e. $P(q_j|q_i)$. They are stored in matrices where each term $a_{ij}$ denotes a probability $P(q_j|q_i)$.

- *Non-emitting initial and final states* : For a finite length random sequence, two additional states are used in order to model the "start" or "end" events. These states are not associated with some emission probabilities.

- *Initial state distribution* $P(I|q_j)$ : Transitions starting from the initial state.

- *Final-absorbent state* : The final state usually has only one non-null transition that loops onto itself with a probability of 1

# HMM TERMINOLOGY

- *Emission probabilities* : are the pdfs (usually Gaussians) that characterize each state $q_i$, i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$.

- *Transition probabilities* : are the probability to go from a state $i$ to a state $j$, i.e. $P(q_j|q_i)$. They are stored in matrices where each term $a_{ij}$ denotes a probability $P(q_j|q_i)$.

- *Non-emitting initial and final states* : For a finite length random sequence, two additional states are used in order to model the "start" or "end" events. These states are not associated with some emission probabilities.

- *Initial state distribution $P(I|q_j)$* : Transitions starting from the initial state.

- *Final-absorbent state* : The final state usually has only one non-null transition that loops onto itself with a probability of 1
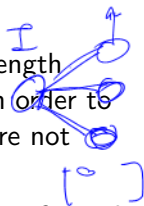
# HMM TERMINOLOGY

- *Emission probabilities*: are the pdfs (usually Gaussians) that characterize each state $q_i$, i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$.
- *Transition probabilities*: are the probability to go from a state $i$ to a state $j$, i.e. $P(q_j|q_i)$. They are stored in matrices where each term $a_{ij}$ denotes a probability $P(q_j|q_i)$.
- *Non-emitting initial and final states*: For a finite length random sequence, two additional states are used in order to model the "start" or "end" events. These states are not associated with some emission probabilities.
- *Initial state distribution $P(I|q_j)$*: Transitions starting from the initial state.
- *Final-absorbent state*: The final state usually has only one non-null transition that loops onto itself with a probability of 1

# HMM TERMINOLOGY

- *Emission probabilities*: are the pdfs (usually Gaussians) that characterize each state $q_i$, i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$.

- *Transition probabilities*: are the probability to go from a state $i$ to a state $j$, i.e. $P(q_j|q_i)$. They are stored in matrices where each term $a_{ij}$ denotes a probability $P(q_j|q_i)$.

- *Non-emitting initial and final states*: For a finite length random sequence, two additional states are used in order to model the "start" or "end" events. These states are not associated with some emission probabilities.

- *Initial state distribution* $P(I|q_j)$: Transitions starting from the initial state.

- *Final-absorbent state*: The final state usually has only one non-null transition that loops onto itself with a probability of 1

# HMM TERMINOLOGY

- *Ergodic HMM*: an HMM allowing for transitions from any emitting state to any other emitting state
- *Left-right HMM*: an HMM where the transitions only go from one state to itself or to a unique follower.
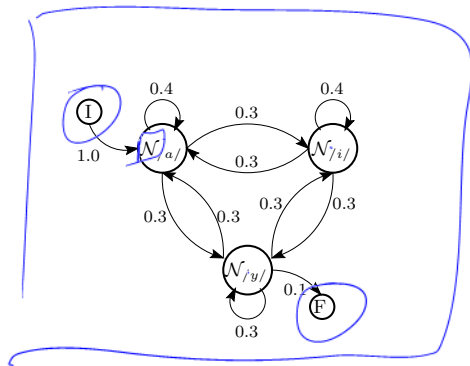
# HMM TERMINOLOGY



- *Ergodic HMM*: an HMM allowing for transitions from any emitting state to any other emitting state
- *Left-right HMM*: an HMM where the transitions only go from one state to itself or to a unique follower.
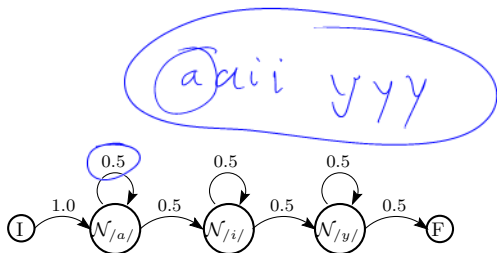
# HMM examples

**HMM1:**



**Transition matrix**

$$
\begin{bmatrix}
0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\
0.0 & 0.4 & 0.3 & 0.3 & 0.0 \\
0.0 & 0.3 & 0.4 & 0.3 & 0.0 \\
0.0 & 0.3 & 0.3 & 0.3 & 0.1 \\
0.0 & 0.0 & 0.0 & 0.0 & 1.0
\end{bmatrix}
$$

$a_{i,j}$

# HMM examples

**HMM2:**



**Transition matrix**

$$
\begin{bmatrix}
0.0 & \mathbf{1.0} & 0.0 & 0.0 & 0.0 \\
0.0 & \mathbf{0.5} & \mathbf{0.5} & 0.0 & 0.0 \\
0.0 & 0.0 & \mathbf{0.5} & \mathbf{0.5} & 0.0 \\
0.0 & 0.0 & 0.0 & \mathbf{0.5} & \mathbf{0.5} \\
0.0 & 0.0 & 0.0 & 0.0 & \mathbf{1.0}
\end{bmatrix}
$$

# HMM EXAMPLES

**HMM3:**



**Transition matrix**

$$
\begin{bmatrix}
0.0 & \mathbf{1.0} & 0.0 & 0.0 & 0.0 \\
0.0 & 0.95 & 0.05 & 0.0 & 0.0 \\
0.0 & 0.0 & 0.95 & 0.05 & 0.0 \\
0.0 & 0.0 & 0.0 & 0.95 & 0.05 \\
0.0 & 0.0 & 0.0 & 0.0 & \mathbf{1.0}
\end{bmatrix}
$$

# HMM MODEL: $\Theta$

$$\Theta = \left\{ a_{ij}, \mu_i, \Sigma_i \right\} \qquad \Theta = \left\{ a_{ij}, p_k^i, \mu_k^i, \Sigma_k^i \right\}$$

In the case of HMMs with Gaussian emission probabilities, the parameter set $\Theta$ comprises :

- the transition probabilities $a_{ij}$;
- the parameters of the Gaussian densities characterizing each state, i.e. the means $\mu_i$ and the variances $\Sigma_i$.

The initial state distribution is sometimes modeled as an additional parameter instead of being represented in the transition matrix.

# Size of an HMM model: Ergodic and Gaussian case

In the case of an ergodic HMM with $N$ emitting states and Gaussian emission probabilities, we have :

- $(N-2) \times (N-2)$ transitions, plus $(N-2)$ initial state probabilities and $(N-2)$ probabilities to go to the final state;

- $(N-2)$ emitting states where each pdf is characterized by a $D$ dimensional mean and a $D \times D$ covariance matrix.

Hence, in this case, the total number of parameters is $(N-2) \times (N + D \times (D+1))$.

# OUTLINE

Likelihood of a sequence given a HMM:

$$p(\underline{X}|\Theta),$$

i.e. the likelihood of an observation sequence given a model.

- Assume a state sequence $Q = \{q_1, \cdots, q_T\}$
- Probability of a state sequence:

$$P(Q|\Theta) = \prod_{t=1}^{T-1} a_{t,t+1} = a_{1,2} \cdot a_{2,3} \cdots a_{T-1,T}$$

$$X = \{ X_1, X_2, \cdots, X_T \}$$

- Assume a state sequence $Q = \{ q_1, \cdots, q_T \}$
- *Probability of a state sequence* :

$$P(Q|\Theta) = \prod_{t=1}^{T-1} a_{t,t+1} = a_{1,2} \cdot a_{2,3} \cdots a_{T-1,T}$$

$Q$

# LIKELIHOOD OF AN OBSERVATION SEQUENCE GIVEN A STATE SEQUENCE

- Assume an observation sequence $X = \{x_1, x_2, \cdots, x_T\}$ and a state sequence $Q = \{q_1, \cdots, q_T\}$
- Likelihood of an observation sequence along a single path, $Q$, for an HMM, $\Theta$:

$$p(X|Q, \Theta) = \prod_{i=1}^{T} p(x_i|q_i, \Theta) = b_1(x_1) \cdot b_2(x_2) \cdots b_T(x_T)$$

# LIKELIHOOD OF AN OBSERVATION SEQUENCE GIVEN A STATE SEQUENCE

$P(Q|\theta)$

- Assume an observation sequence $X = \{x_1, x_2, \cdots, x_T\}$ and a state sequence $Q = \{q_1, \cdots, q_T\}$
- Likelihood of an observation sequence along a single path, $Q$, for an HMM, $\Theta$:

$$p(X|Q, \Theta) = \prod_{i=1}^{T} p(x_i|q_i, \Theta) = b_1(x_1) \cdot b_2(x_2) \cdots b_T(x_T)$$

# LIKELIHOODS

- *Joint likelihood of an observation sequence X and a path Q* : it consists in the probability that $X$ and $Q$ occur simultaneously, $p(X, Q|\Theta)$, and decomposes into a product of the two quantities defined previously :

$$p(X, Q|\Theta) = p(X|Q, \Theta)P(Q|\Theta) \qquad \text{(Bayes)}$$

- *Likelihood of a sequence with respect to a HMM* : the likelihood of an observation sequence $X = \{x_1, x_2, \cdots, x_T\}$ with respect to a Hidden Markov Model with parameters $\Theta$ expands as follows :

$$p(X|\Theta) = \sum_{\text{every possible } Q} p(X, Q|\Theta)$$

# LIKELIHOODS

- *Joint likelihood of an observation sequence X and a path Q* : it consists in the probability that $X$ and $Q$ occur simultaneously, $p(X, Q|\Theta)$, and decomposes into a product of the two quantities defined previously :

$$p(X, Q|\Theta) = p(X|Q, \Theta)P(Q|\Theta) \qquad \text{(Bayes)}$$

- *Likelihood of a sequence with respect to a HMM* : the likelihood of an observation sequence $X = \{x_1, x_2, \cdots, x_T\}$ with respect to a Hidden Markov Model with parameters $\Theta$ expands as follows :

$$p(X|\Theta) = \sum_{\text{every possible } Q} p(X, Q|\Theta)$$

# FORWARD RECURSION

- There is a recursive way to compute $p(X|\Theta)$: Forward Recursion (FR)

- In FR, we define a *forward* variable:

$$p_t(i) = p(x_1, x_2, \cdots x_t, q^t = q_i | \Theta)$$

i.e. $p_t(i)$ is the probability of having observed the partial sequence $\{x_1, x_2, \cdots, x_t\}$ *and* being in the state $i$ at time $t$, given parameters $\Theta$.

# Forward Recursion

- There is a recursive way to compute $p(X|\Theta)$: Forward Recursion (FR)

- In FR, we define a *forward* variable:

$$p_t(i) = p(x_1, x_2, \cdots x_t, q^t = q_i | \Theta)$$

i.e. $p_t(i)$ is the probability of having observed the partial sequence $\{x_1, x_2, \cdots, x_t\}$ *and* being in the state $i$ at time $t$, given parameters $\Theta$.

Assume $N$ states with $N-2$ emitting states.

- **Initialization:**

$$p_1(j) = a_{1j} b_j(x_1)$$

  with $2 \leq j \leq N-1$

- **Recursion:**

$$p_t(j) = \left[ \sum_{i=2}^{N-1} p_{t-1}(i) \cdot a_{ij} \right] b_j(x_t),$$

  with $2 \leq t \leq T$ and $2 \leq j \leq N-1$

- **Termination:**

$$p(X|\Theta) = \left[ \sum_{i=2}^{N-1} p_T(i) \cdot a_{iN} \right]$$

# Computation of the forward variable

Assume $N$ states with $N - 2$ emitting states.

- **Initialization:**

$$p_1(j) = a_{1j}b_j(x_1)$$

with $2 \leq j \leq N - 1$

- **Recursion:**

$$p_t(j) = \left[\sum_{i=2}^{N-1} p_{t-1}(i) \cdot a_{ij}\right] b_j(x_t),$$

with $2 \leq t \leq T$ and $2 \leq j \leq N - 1$

- Termination:

$$p(X|\Theta) = \left[\sum_{i=2}^{N-1} p_T(i) \cdot a_{iN}\right]$$

# COMPUTATION OF THE FORWARD VARIABLE

Assume $N$ states with $N - 2$ emitting states.

- **Initialization:**

$$p_1(j) = a_{1j} b_j(x_1)$$

  with $2 \leq j \leq N - 1$

- **Recursion:**

$$p_t(j) = \left[ \sum_{i=2}^{N-1} p_{t-1}(i) \cdot a_{ij} \right] b_j(x_t),$$

  with $2 \leq t \leq T$ and $2 \leq j \leq N - 1$

- **Termination:**

$$p(X|\Theta) = \left[ \sum_{i=2}^{N-1} p_T(i) \cdot a_{iN} \right]$$

# BAYESIAN CLASSIFICATION

- Assume that there are many HMMs, $\Theta_i$, $i = 1, \cdots, M$
- Given the likelihood $p(X|\Theta_i)$ computed using the forward recursion algorithm, we can compute the probability of $\Theta_i$, using Bayes' rule:

$$
\begin{aligned}
P(\Theta_i|X) &= \frac{p(X|\Theta_i)P(\Theta_i)}{P(X|\Theta)} \\
&\propto p(X|\Theta_i)P(\Theta_i)
\end{aligned}
$$

- Other solution: Maximum likelihood.

# BAYESIAN CLASSIFICATION

- Assume that there are many HMMs, $\Theta_i$, $i = 1, \cdots, M$
- Given the likelihood $p(X|\Theta_i)$ computed using the forward recursion algorithm, we can compute the probability of $\Theta_i$, using Bayes' rule:

$$P(\Theta_i|X) = \frac{p(X|\Theta_i)P(\Theta_i)}{P(X|\Theta)}$$
$$\propto p(X|\Theta_i)P(\Theta_i)$$

- Other solution: Maximum likelihood.

- Assume that there are many HMMs, $\Theta_i$, $i = 1, \cdots, M$
- Given the likelihood $p(X|\Theta_i)$ computed using the forward recursion algorithm, we can compute the probability of $\Theta_i$, using Bayes' rule:

$$P(\Theta_i|X) = \frac{p(X|\Theta_i)P(\Theta_i)}{P(X|\Theta)}$$
$$\propto p(X|\Theta_i)P(\Theta_i)$$

- Other solution: Maximum likelihood.

# Outline

# Definitions

- *Highest* likelihood $\delta_t(i)$ along a *single* path among all the paths ending in state $i$ at time $t$:

$$\delta_t(i) = \max_{q_1, q_2, \cdots, q_{t-1}} p(q_1, q_2, \cdots, q_{t-1}, q^t = q_i, x_1, x_2, \cdots x_t | \Theta)$$

- *Buffer* $\psi_t(i)$ which allows to keep track of the "best path" ending in state $i$ at time $t$:

$$\psi_t(i) = \underset{q_1, q_2, \cdots, q_{t-1}}{\arg\max} \; p(q_1, q_2, \cdots, q_{t-1}, q^t = q_i, x_1, x_2, \cdots x_t | \Theta)$$

- *Highest* likelihood $\delta_t(i)$ along a *single* path among all the paths ending in state $i$ at time $t$:

$$\delta_t(i) = \max_{q_1, q_2, \cdots, q_{t-1}} p(q_1, q_2, \cdots, q_{t-1}, q^t = q_i, x_1, x_2, \cdots x_t | \Theta)$$

- *Buffer* $\psi_t(i)$ which allows to keep track of the "best path" ending in state $i$ at time $t$:

$$\psi_t(i) = \operatorname*{argmax}_{q_1, q_2, \cdots, q_{t-1}} p(q_1, q_2, \cdots, q_{t-1}, q^t = q_i, x_1, x_2, \cdots x_t | \Theta)$$

# VITERBI ALGORITHM

**1** **Initialization :**

$$\delta_1(i) = a_{1i} \cdot b_i(x_1), \quad 2 \leq i \leq N-1$$
$$\psi_1(i) = 0$$

**2** **Recursion :**

$$\delta_{t+1}(j) = \max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{ij}] \cdot b_j(x_{t+1}), \quad \begin{array}{l} 1 \leq t \leq T-1 \\ 2 \leq j \leq N-1 \end{array}$$

$$\psi_{t+1}(j) = \arg\max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{ij}], \quad \begin{array}{l} 1 \leq t \leq T-1 \\ 2 \leq j \leq N-1 \end{array}$$

**3** **Termination :**

$$p^*(X|\Theta) = \max_{2 \leq i \leq N-1} [\delta_T(i) \cdot a_{iN}]$$
$$q_T^* = \arg\max_{2 \leq i \leq N-1} [\delta_T(i) \cdot a_{iN}]$$

**4** **Backtracking :**

$$Q^* = \{q_1^*, \cdots, q_T^*\} \quad \text{so that} \quad q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2,$$

# OUTLINE

# ACKNOWLEDGMENTS