

Hidden Markov Model and its applications in Signal Processing

Kanru Hua

Univ. of Illinois at Urbana-Champaign

Apr. 27, 2016



What is Hidden Markov Model?

- A hidden Markov model (HMM) is a **statistical** Markov model in which the system being modeled is assumed to be a **Markov process** with **unobserved** (hidden) states.

– Wikipedia

- **Typical applications:**
 - Speech Recognition/Synthesis
 - Part-of-speech Tagging
 - Handwriting Recognition
 - Automatic Audio/Video Transcription
 - Signal Denoising

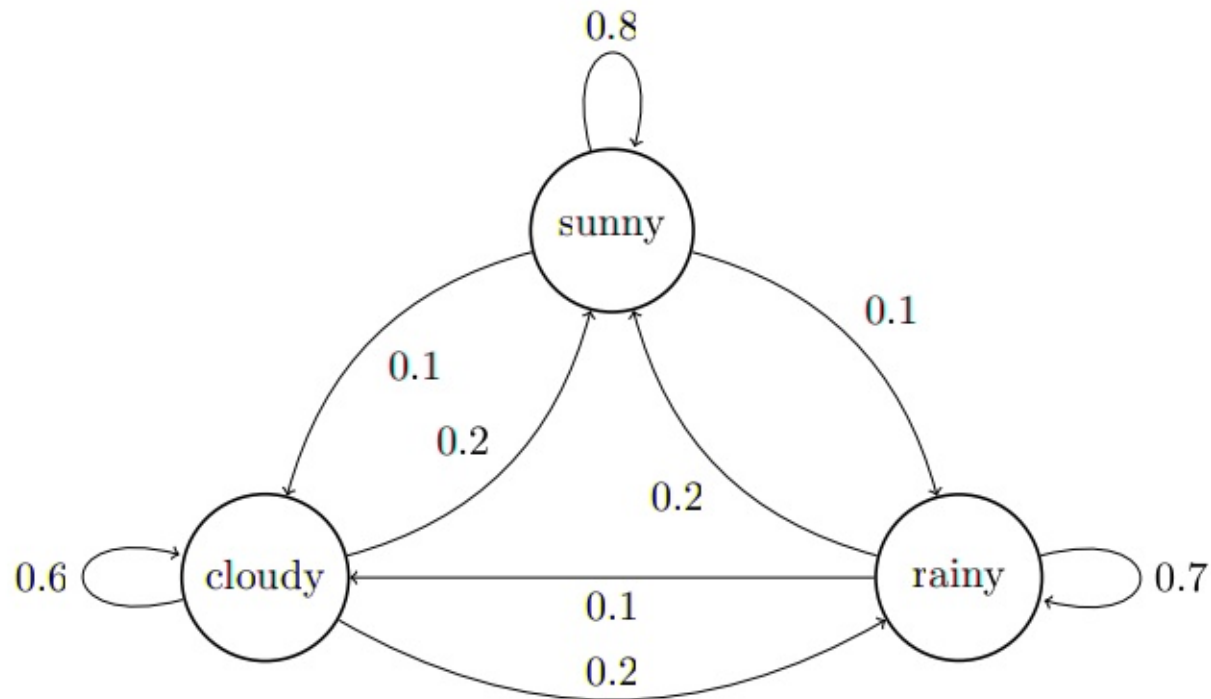


Outline

- **Markov Chain**
- **Hidden Markov Model**
 - Definition
 - Algorithms – Inference and Training
 - Run HMM in Matlab
- **Extensions to HMM**
 - Continuous (and Multivariate) Output
 - Multiple Observation Sequences
- **Applications**
 - Speech Recognition
 - Filtering/Denoising



Markov Chain



- **Random walk example:**

→ Sunny → Sunny → Cloudy → Rainy → Rainy → Sunny →
Sunny → Sunny → Rainy → Rainy → Cloudy → Sunny → ...



Markov Chain

- **Definition**

parameter set: $\lambda = \{\pi, \mathbf{A}\}$

initial probability: $\pi_j = P(s_1 = j | \lambda)$

transition probability: $a_{ij} = P(s_t = j | s_{t-1} = i, \lambda)$

- **Markov Property**

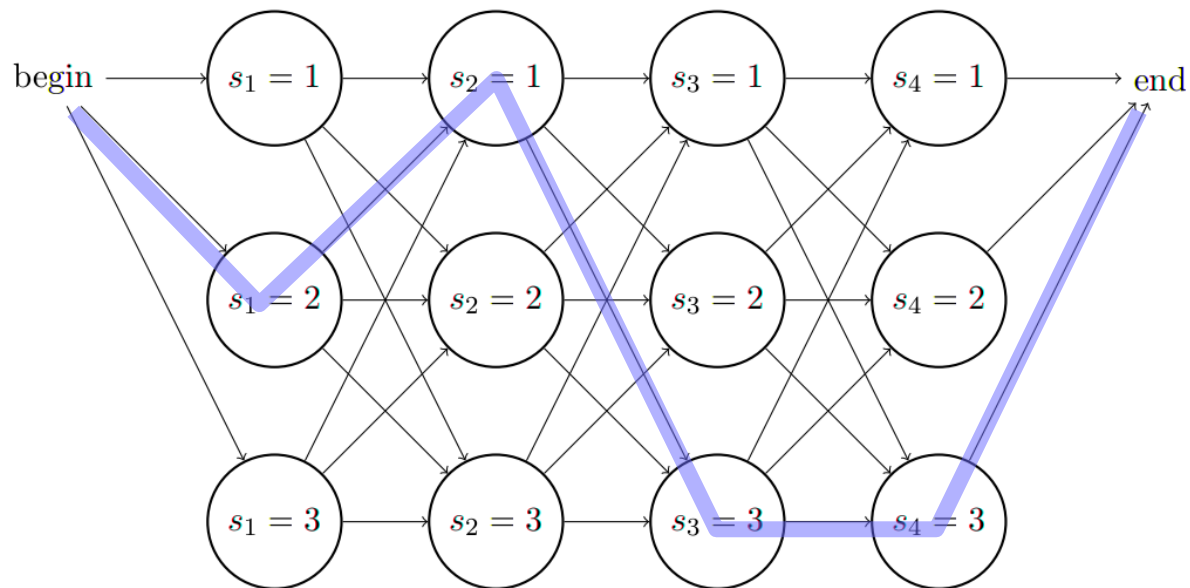
$$P(s_t = j | s_1, s_2, \dots, s_{t-2}, s_{t-1} = i, \lambda) = P(s_t = j | s_{t-1} = i, \lambda) = a_{ij}$$

Getting into current state only depends on the previous state but not any state before the previous state. (First-order Markov Chain)



Markov Chain

- A lateral view



- Total Probability

$$\begin{aligned} P(s_1=2, s_2=1, s_3=3, s_4=3|\lambda) &= P(s_1=2|\lambda)P(s_2=1|s_1=2,\lambda)P(s_3=3|s_2=1,\lambda)P(s_4=3|s_3=3,\lambda) \\ &= \pi_2 a_{2,1} a_{1,3} a_{3,3} \end{aligned}$$



Hidden Markov Model: an Example

On sunny days I usually go to ECEB by bike;

On cloudy days sometimes I walk.

On rainy days I either walk or take MTD, but I rarely bike.

- My behavior can be modelled by a HMM and we can
 - Given weather data and my transportation record, train a HMM.
 - Given a HMM, generate a weather sequence and corresponding transportation record.
 - Given HMM & my transportation record, estimate the most probable weather sequence.
 - Given HMM & my transportation record, estimate the probability of having {sunny, cloudy, rainy} weather on each day.
 - Given HMM & transportation records of a few people at different locations, pick the one that is most likely to be my record (in Champaign).



Hidden Markov Model

- Definition

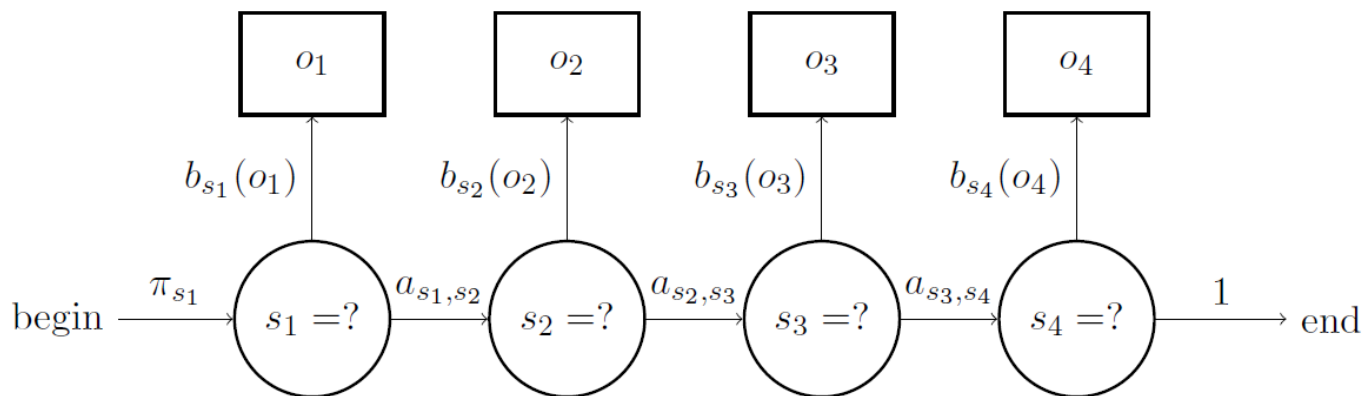
parameter set: $\lambda = \{\pi, \mathbf{A}, \mathbf{B}\}$

initial probability: $\pi_j = P(s_1 = j | \lambda)$

transition probability: $a_{ij} = P(s_t = j | s_{t-1} = i, \lambda)$

output probability (discrete case): $b_{jk} = P(o_t = k | s_t = j, \lambda)$

output probability density (continuous case): $b_j(o_t) = p(o_t | s_t = j, \lambda)$



Hidden Markov Model

- What can we do with HMM?
- Generation
 - Random walk
- Inference
 - Forward/Backward algorithm
 - Viterbi algorithm (most probable sequence)
 - Total probability and probability of a particular state
- Parameter estimation (training)
 - Initialization
 - Viterbi training algorithm
 - Baum-Welch algorithm



Extension – Continuous Output

- Use a continuous density distribution for $b_j(o_t)$
- e.g. normal distribution

$$b_j(o_t = x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- Or multivariate normal distribution with dimension k

$$b_j(o_t = x) = (2\pi)^{-\frac{k}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

- To reduce number of parameters we often assume diagonal covariance, which significantly simplifies the whole thing.



Extension – Continuous Output

- The inference algorithms are exactly same as the discrete case (except for function $b_j(o_t)$ itself, of course).
- Maximization step needs a little bit change:

$$\mu_j = \frac{\sum_{t=1}^T \gamma_t(j) o_t}{\sum_{t=1}^T \gamma_t(j)} \quad \Sigma_j = \frac{\sum_{t=1}^T \gamma_t(j) o_t o_t^T}{\sum_{t=1}^T \gamma_t(j)} - \mu_j \mu_j^T$$

- Gaussian Mixture Model as output density distribution

$$p(x \in \mathbf{R}^k | \mathbf{c}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{k}{2}} \sum_j c_j |\Sigma_j|^{-\frac{k}{2}} e^{-\frac{k}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}$$

$$g_{jk}(o_t) = p(o_t | s_t = j, m_t = k, \lambda) = p(o_t | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk})$$

$$b_j(o_t) = p(o_t | s_t = j, \lambda) = \sum_k c_{jk} g_{jk}(o_t)$$



Extension – Multiple Observation Sequences

- Training data set often contains more than one sample.
- In fact there are often hundreds or even thousands of samples.
- Use common sense, just add (average) bunch of things together!

$$\pi_i = \frac{1}{L} \sum_{l=0}^L \gamma_1^l(i)$$
$$a_{ij} = \frac{\sum_{l=0}^L \sum_{t=2}^T \gamma_t^l(i, j)}{\sum_{l=0}^L \sum_{t=2}^T \gamma_{t-1}^l(i)} \quad b_{ik} = \frac{\sum_{l=0}^L \sum_{t:o_t=k} \gamma_t^l(i)}{\sum_{l=0}^L \sum_{t=0}^T \gamma_t^l(i)}$$

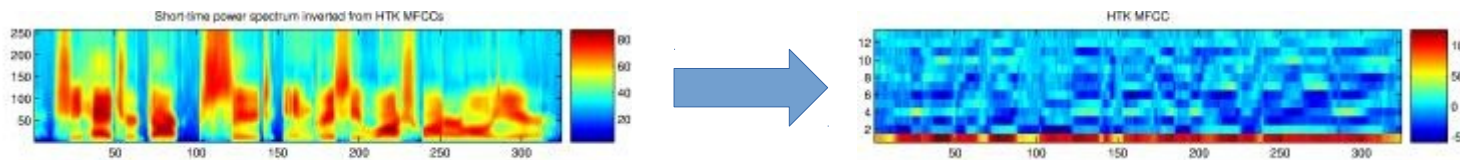
- To prove this you can change the “s” under summation in EM algorithm and rederive the whole thing again.



HMM for Speech Recognition

- **Speech feature**

MFCC (Mel-Frequency Cepstral Coefficient)



- **Output density function:**

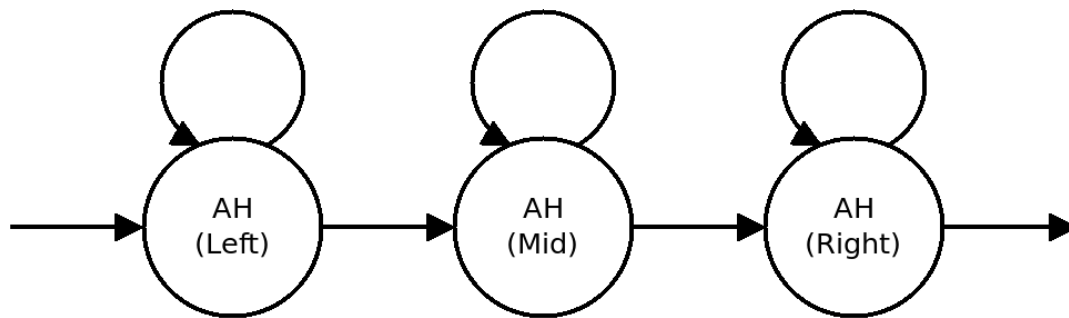
(Mixture of) multivariate normal distribution with diagonal covariance,
whose dimension = dimension of each feature vector

Some earlier approaches use discrete HMM with vector-quantized feature labels, output matrix is like a bunch of histograms.



HMM for Speech Recognition

- **Triphone Models**



	Left	Mid	Right	Next
Left	0.7	0.3	0	0
Mid	0	0.9	0.1	0
Right	0	0	0.7	0.3

“**LLLMMMMMMRRN**”

- **Forced Alignment**

Speech data often comes with phoneme transcription. Though exact state sequence is not always given, the model topology for each sample is known.

E.g. “Hello” → “HH AH L OW”

For this sample, HH(R) → AH(L) exists; AH(R) → L(L) exists; HH(R) → OW(L) doesn't exist.

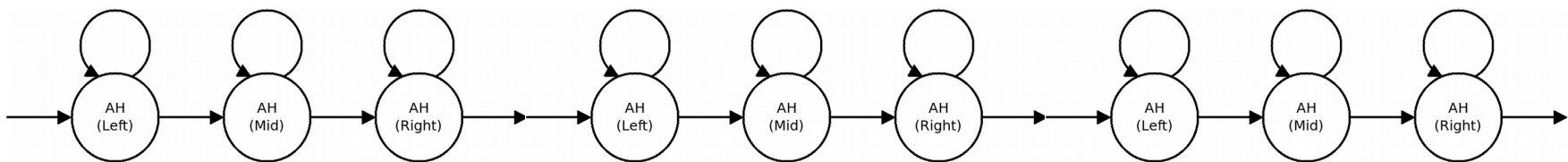


HMM for Speech Recognition

- **Forced Alignment (cont)**

We don't want the algorithm to consider impossible/confusing transitions (e.g. HH(R) \rightarrow HH(L) or HH(R) \rightarrow OW(L))

Then during inference, simply change the transition matrix for each training sample. Each state should only have *self transition* and *transition to the next state*. \rightarrow Left-to-Right HMM



Forced alignment – given phoneme transcription and speech, calculate the starting/ending time of each phoneme (using Viterbi algorithm).



HMM for Speech Recognition

- **State Tying**

What if a phoneme appears more than once in a sample?

e.g. “Hello, world.” → “HH AH **L** OW, W ER **L** D”, where we can find $L(R) \rightarrow OW(L)$ and $L(R) \rightarrow D(L)$

- Just create a copy of L(R, M, L) states, and name them L1(R, M, L)

HH AH **L** OW, W ER **L1** D

L and L1 share the same output distribution. During maximization step, statistics on L and L1 are collected and grouped together. This trick is called **state tying**.



HMM for Speech Recognition

- **Context Dependency**

So far we've only seen **context-independent (CI)** HMM – same L, M, R states are shared across the whole data set.

However pronunciation could vary with respect to context.

Full-context models: represent each occurrence of each phoneme by a unique set of L, M, R states.

- **Context-dependent (CD)** models: **cluster** full-context models into groups that share both similar context and similar output distribution, then **merge** & **tie** them together.
- We can use a decision tree for such clustering.

