

Hidden Markov Models

Sargur Srihari

srihari@cedar.buffalo.edu

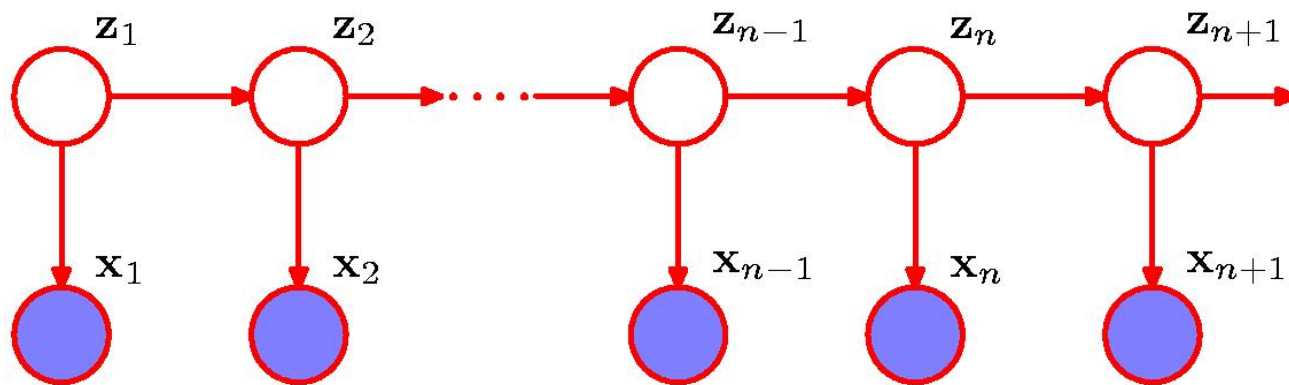
HMM Overview

1. What is an HMM?
2. State-space Representation
3. HMM Parameters
4. Generative View of HMM

- Ubiquitous tool for modeling time series data
 - Used in most all speech recognition systems
 - Computational molecular biology
 - Group amino acid sequences into proteins
- It is a BN for representing probability distributions over sequences of observations
- HMM has two defining properties:
 - Observation \mathbf{x}_t at time t was generated by some process whose state \mathbf{z}_t is hidden from the observer
 - Assumes that state at \mathbf{z}_t is dependent only on state \mathbf{z}_{t-1} and independent of all prior states (First order)
- Example:
 - \mathbf{z} are phoneme sequences, \mathbf{x} are acoustic observations

Graphical Model of an HMM

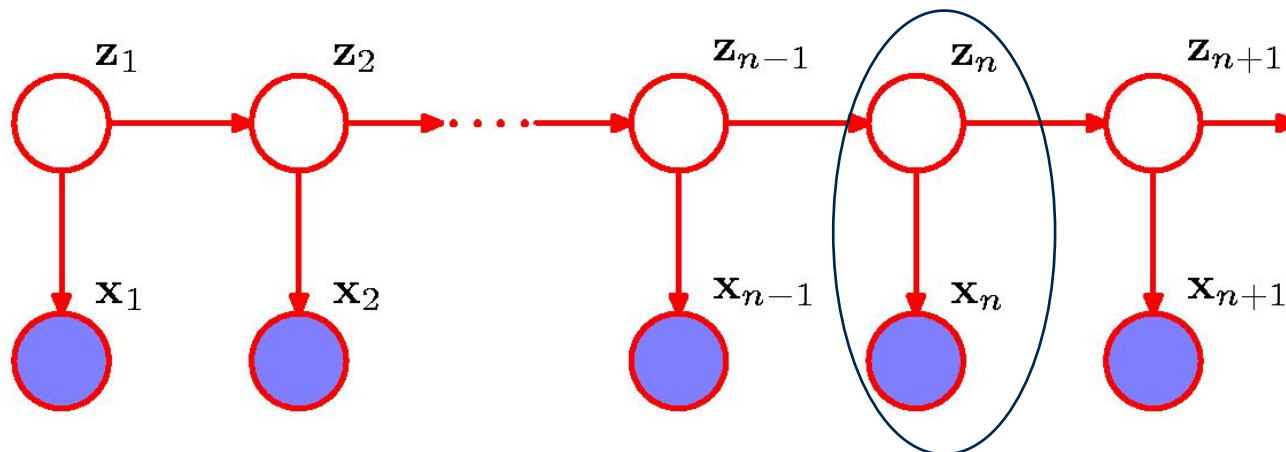
- Has the state space model shown below and latent variables are discrete



- Joint distribution has the form

$$p(x_1, \dots, x_N, z_1, \dots, z_N) = p(z_1) \left[\prod_{n=2}^N p(z_n | z_{n-1}) \right] \prod_{n=1}^N p(x_n | z_n)$$

Mixture Viewed as HMM



- A single time slice corresponds to a mixture distribution with component densities $p(x|z)$
- An extension of mixture model
 - Choice of mixture component depends on choice of mixture component for previous distribution
- Latent variables are multinomial variables z_n
 - That describe component responsible for generating x_n
- Can use *one-of-K* coding scheme

Transitional Probability Matrix

- Joint distribution:

$$p(x_1, \dots, x_N, z_1, \dots, z_n) = p(z_1) \left[\prod_{n=2}^N p(z_n | z_{n-1}) \right] \prod_{n=1}^N p(x_n | z_n)$$

- We allow z_n to depend on z_{n-1} via $p(z_n | z_{n-1})$

- Latent variables are 1-of- K , thus this conditional distribution is specified by a transition probability matrix A with $A_{jk} = p(z_{nk}=1 | z_{n-1,j}=1)$

- Conditional distribution is stated as

$$p(z_n | z_{n-1}, A) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{n,k}}$$

- Exponent $z_{n-1,j} z_{n,k}$, which is a product, is 0 or 1
- Hence product evaluates to a single A_{jk} for each setting of values of z_n , z_{n-1}
- Thus $p(z_n=3 | z_{n-1}=2) = A_{23}$

State of z_n State of z_{n-1} A matrix

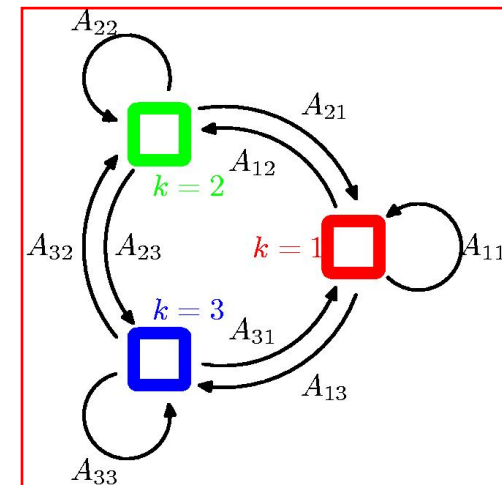
k	1	2	...	K
z_{nk}	0	1	...	0

j	1	2	...	K
$z_{n-1,j}$	1	0	...	0

	1	2	...	K
1				
2				
...			A_{jk}	
K				s

	$z_n=1$ $z_{n1}=1$ $z_{n2}=0$ $z_{n3}=0$	$z_n=2$ $z_{n1}=0$ $z_{n2}=1$ $z_{n3}=0$	$z_n=3$ $z_{n1}=0$ $z_{n2}=0$ $z_{n3}=1$
$z_{n-1}=1$ $z_{n-1,1}=1$ $z_{n-1,2}=0$ $z_{n-1,3}=0$	A_{11}	A_{12}	A_{13}
$z_{n-1}=2$ $z_{n-1,1}=0$ $z_{n-1,2}=1$ $z_{n-1,3}=0$	A_{21}	A_{22}	A_{23}
$z_{n-1}=3$ $z_{n-1,1}=0$ $z_{n-1,2}=0$ $z_{n-1,3}=1$	A_{31}	A_{32}	A_{33}

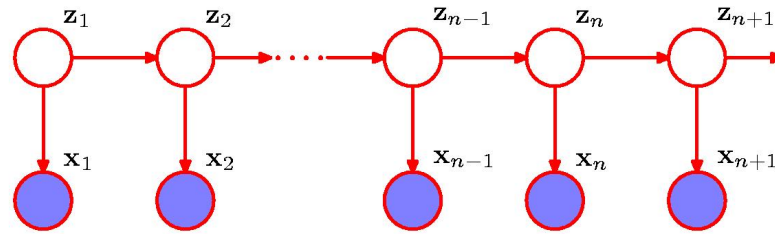
Transition diagram where latent variables have three possible states



- Not a graphical model since nodes are not separate variables but states of a single variable
- Here $K=3$

Matrix A has $K(K-1)$ independent parameters

Initial Variable Probabilities



- Initial latent node z_1 does not have a parent node
 - Its marginal distribution $p(z_1)$ is represented by a vector of probabilities π with elements $\pi_k = p(z_{1k}=1)$ so that

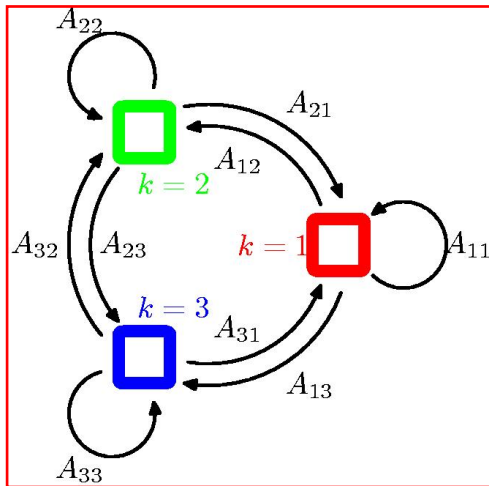
$$p(z_1 \mid \pi) = \prod_{k=1}^K \pi_k^{z_{1,k}} \text{ where } \sum_k \pi_k = 1$$

- Note that π is an HMM parameter
 - representing probabilities of each state for the first variable

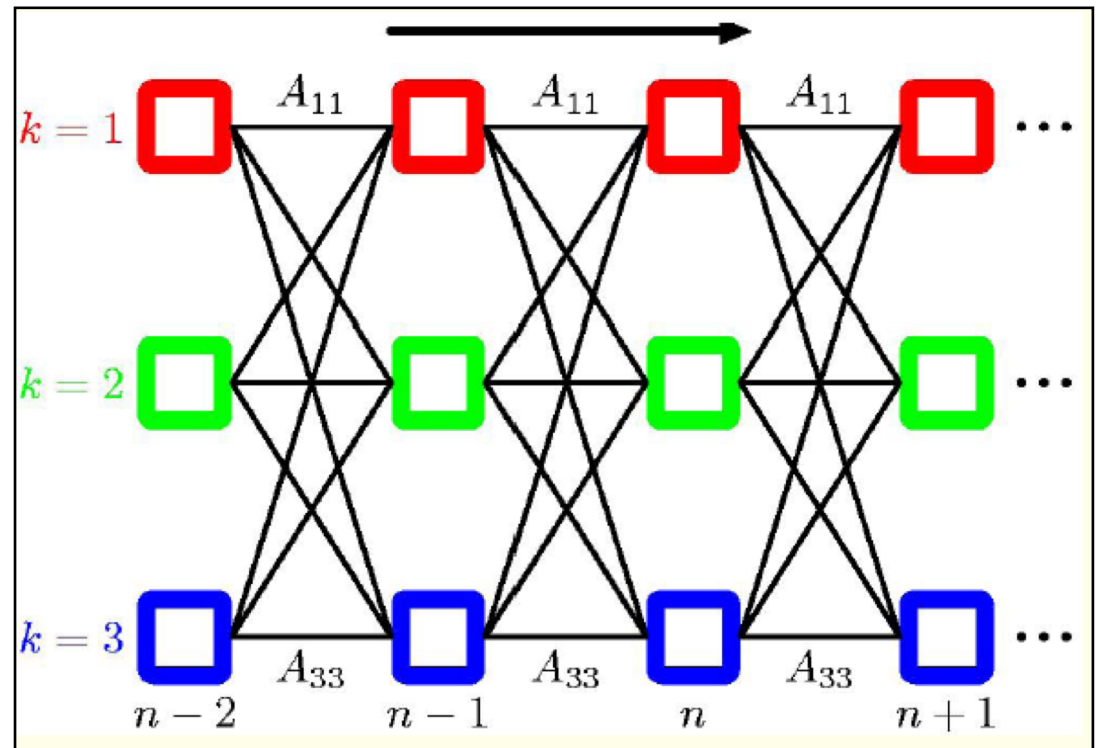
Unfolding State Transition over time

Lattice or trellis representation of the latent states

State Transition diagram



Each column corresponds to one of the z_n



Emission Probabilities $p(\mathbf{x}_n | \mathbf{z}_n)$

- Specification of probabilistic model is completed by defining conditional distributions of observed variables $p(\mathbf{x}_n | \mathbf{z}_n, \phi)$ where ϕ are parameters
 - These are known as *emission probabilities* which can be either continuous (by Gaussians) or discrete (by tables)
- Because \mathbf{x}_n is observed, distribution $p(\mathbf{x}_n | \mathbf{z}_n, \phi)$ consists of a table of K numbers corresponding to K states of binary vector \mathbf{z}_n
 - Emission probabilities can be represented as
$$p(\mathbf{x}_n | \mathbf{z}_n, \phi) = \prod_{k=1}^K (p(\mathbf{x}_n | \phi_k))^{z_{nk}}$$

Homogeneous Models

- All conditional distributions governing the latent variables share the same parameters \mathbf{A}
- All the emission distributions share the same parameters ϕ

Joint Distribution over Latent and Observed variables

- Joint can be expressed in terms of parameters:

$$p(X, Z | \theta) = p(z_1 | \pi) \left[\prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \prod_{m=1}^N p(x_m | z_m, \phi)$$

where $X = \{x_1, \dots, x_N\}$, $Z = \{z_1, \dots, z_N\}$, $\theta = \{\pi, A, \phi\}$

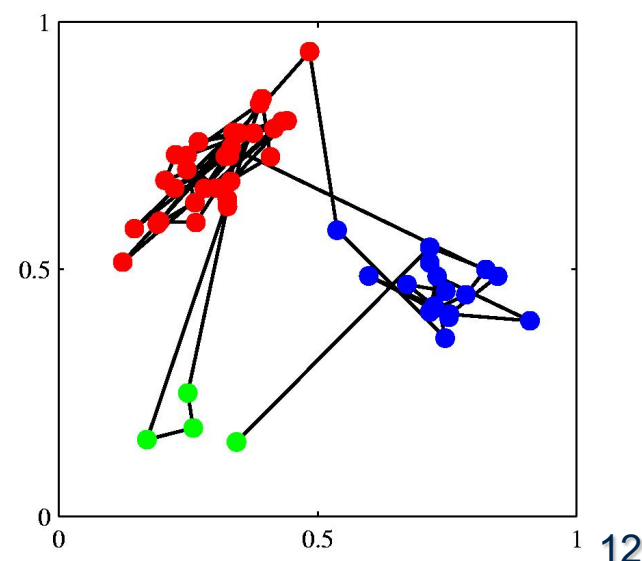
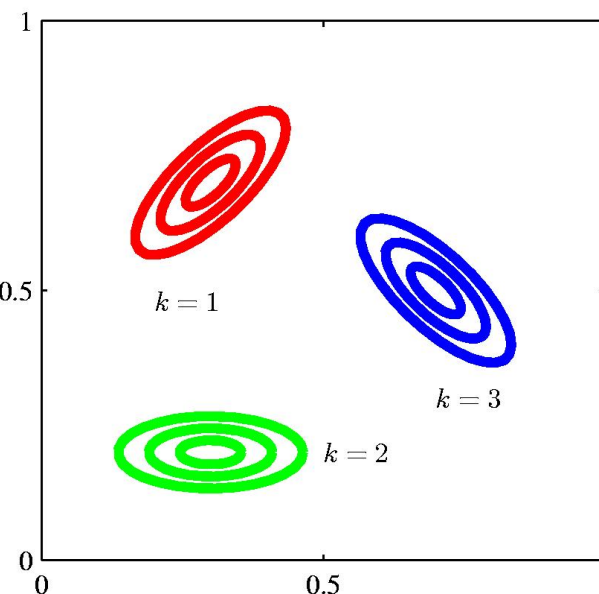
- Most discussion of HMM is independent of emission probabilities
 - Tractable for discrete tables, Gaussian, GMMs

HMM from a generative viewpoint

- Can get a better understanding of HMMs by considering it from a generative viewpoint
- First choose latent variable \mathbf{z}_1 with probabilities determined by probabilities π_k and then sample the corresponding observation \mathbf{x}_1
- Now we choose the state of the variable \mathbf{z}_2 according to the transition probabilities $p(\mathbf{z}_2|\mathbf{z}_1)$ according to the already instantiated value of \mathbf{z}_1
- This is an example of ancestral sampling for a directed PGM

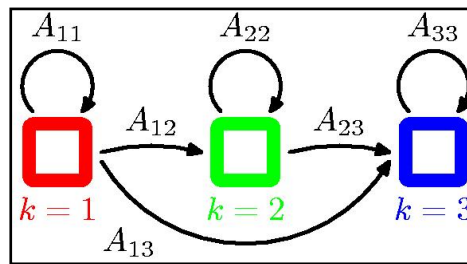
Example of sampling from a HMM

- Generative viewpoint
- Three states of a latent variable
- Gaussian emission model $p(\mathbf{x}|\mathbf{z})$
- Two-dimensional \mathbf{x}
- 50 data points generated
- Transition probabilities
 - 5% probability of making transition
 - 90% of remaining in same

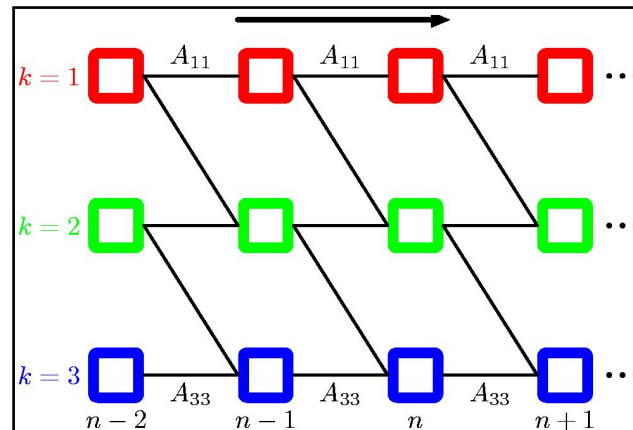


A variant of HMM

- Imposing restrictions on transition matrix A
 - Left-to right HMM
 - Setting the elements of $A_{jk}=0$ if $k < j$

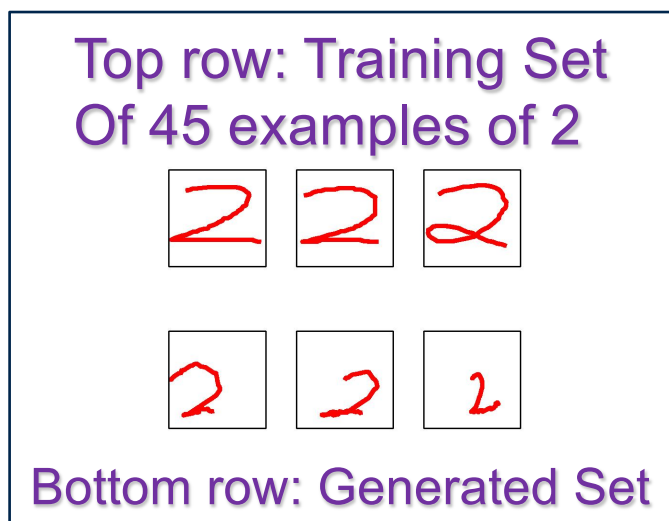


- Once a state has been vacated, it cannot be reentered
- Corresponding lattice



Left-to-Right HMM Applied to Digits

- Examples of on-line handwritten digits



- $K=16$ states: corresponding to a line segment of fixed length in one of 16 angles
- Emission probabilities: 16 x 16 table, associated with allowed angle values for each state
- Transition probabilities set to zero except for those that keep state index k the same or increment by one
- Parameters optimized by 25 iterations of EM

HMM invariance to time warping

- Time warping (compression and stretching)
- On-line handwriting recognition
 - A typical digit consists of two strokes
 - Arc starts at top left down to cusp/loop at bottom left
 - A straight sweep ending at bottom right
 - Relative sizes of the two sections vary and hence the location of the cusp/loop
 - HMM accommodates this by no of transitions to same state vs transitions to successive state
- Speech recognition
 - Warping of time axis is speed of speech
 - HMM accommodates such distortion