

AUTOMATIC SPEECH RECOGNITION

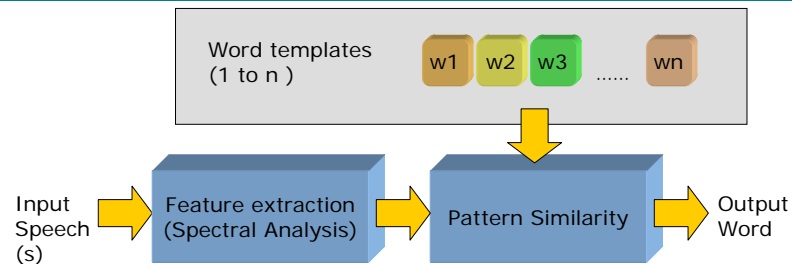
Lecture 8

Probability Review Overview of Hidden Markov Models

This Lecture

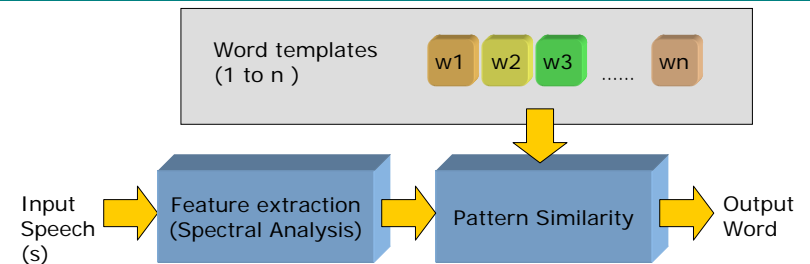
- Motivation
- Review of Probability
- Examples (somewhat) related to ASR
- Overview of HMMs

Word-based Template Matching Summary



- 1 Do spectral analysis to convert the input word to feature vectors
- 2 Compare the feature vectors for all the frames with each of the reference templates. Use DP matching to time align the patterns
- 3 Choose the output word which corresponds to the reference template with the minimum distance

How likely?



Suppose the distance between S and W is D_i .
 W yielding the minimal D_i should be picked as the output word.

How likely that the input speech is W ?




To indicate how likely...

We need
PROBABILITY
!!!



A broad picture

Suppose we want to classify between two words, W1 and W2.

Feature vector: 1-dimensional ... Let's say total Energy.  (Obviously, a bad one)

Training Phase

Learn the value distribution of E from many examples of W1
→ building a model of E for W1

Learn the value distribution of E from many examples of W2
→ building a model of E for W2

Recognizing (Decoding) Phase

Compare E_x , the total energy of the word (X) to be recognized, with the two distributions.

Calculate the probabilities of X being W1 and X being W2. (Given that X must be either one of them)



Important Concepts

- Probability
- Conditional Probability
- Total probability theorem
- Bayes' Rule
- Random variable (R.V.)
 - discrete R.V.
 - continuous R.V.
- Expected Value and Variance
- Gaussian Random Variable
- Joint PDF



Probability

- A : any event
- $P(A)$: probability that the event A happens
- $0 \leq P(A) \leq 1$



Conditional Probability

- $P(A|B)$ = probability of A, given that B has occurred.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

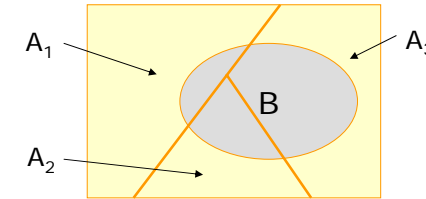
- Adjusting the universe to B

e.g. After finishing an ASR homework, a student has gone crazy and yell out one of the four cardinal vowels randomly.
 $P(\text{the vowel is /i/}) = ?$
 Given that he/she yells out a front vowel,
 $P(\text{/i/} \mid \text{front vowel}) = ?$



Total Probability Theorem

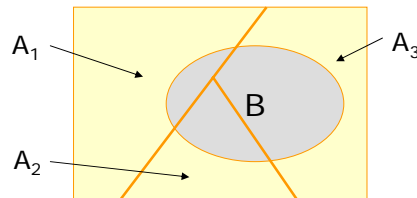
- Divide the universe into smaller partitions



$$P(B) = P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)$$



Bayes' Rule



$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{\sum_j P(A_j)P(B | A_j)}$$

$P(A_i)$: "Prior" probability



Random Variable

- X : Random variable (R.V.)
- x : experimental value of the R.V. X

X	x	
Type of a vowel V	$\{/a/, /i/, \dots\}$	} Discrete
# of syllable in a word W	$\{1, 2, 3, 4, 5, \dots\}$	
The first formant frequency	$(0, F_s/2)$	} Continuous
\log_{10} Energy of a signal	$(-\infty, \infty)$	
Probability of event A	$[0, 1]$	



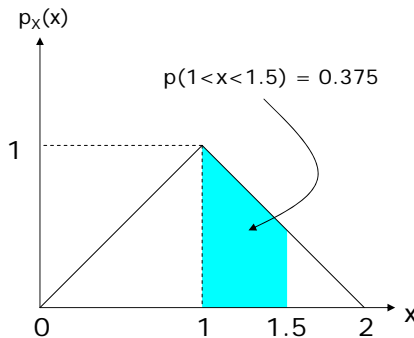
Probability Density Function (PDF)

- $p_X(x)$ is a density of probability measure on the event space for random variable X

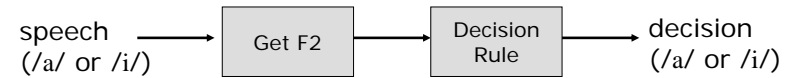
$$P(a < x \leq b) = \int_a^b p(x) dx$$

$$P(-\infty < x \leq \infty) = \int_{-\infty}^{\infty} p(x) dx = 1$$

$$P(x = a) = 0$$



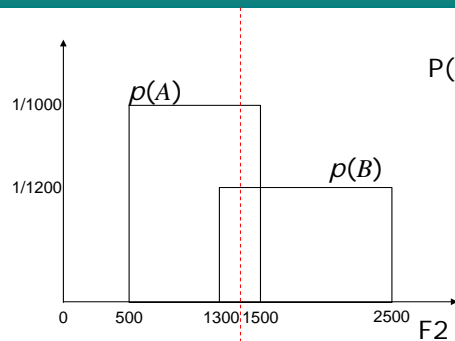
Example



- Let
 - $A = F2$ of /a/
 - $B = F2$ of /i/
- Suppose:
 - A uniformly distributes between 500 and 1500 Hz
 - B uniformly distributes between 1300 and 2500 Hz
- We want to classify an incoming vowel by the following rule:
 - if ($F2 > 1400$ Hz) choose /i/; else choose /a/;
- If both vowels are equally likely to be produced, what is $P(\text{error})$?



Example



$$\begin{aligned}
 P(\text{err}) &= P(\text{choose /a/} \mid /i/)P(/i/) + P(\text{choose /i/} \mid /a/)P(/a/) \\
 &= P(F2 < 1400 \mid /i/) \times 0.5 + P(F2 > 1400 \mid /a/) \times 0.5 \\
 &= P(B < 1400 \mid /i/) \times 0.5 + P(A > 1400 \mid /a/) \times 0.5 \\
 &= 100 \times (1/1200) \times 0.5 + 100 \times (1/1000) \times 0.5 \\
 &= 0.092
 \end{aligned}$$



Expected Value and Variance

- Expected value

$$E[x] = \int_{-\infty}^{\infty} xp(x) dx$$

$$E[g(x)] = \int_{-\infty}^{\infty} g(x)p(x) dx$$

- Variance (σ^2) (Standard Deviation = σ)

$$Var[x] = E[(x - E[x])^2] = \sigma^2 = \int_{-\infty}^{\infty} (x - E[x])^2 p(x) dx$$

$$E[(x - E[x])^2] = E[x^2] - (E[x])^2$$



Expected Value

- $E[a] = a$; a is a constant.
- $E[aX+b] = aE[X]+b$
- $E[X+Y] = E[X]+E[Y]$
- $\text{Var}[a] = 0$
- $\text{Var}[aX+b] = a^2\text{Var}[X]$

$$E[x | A] = \int_{-\infty}^{\infty} xp(x | A)dx$$

$$E[g(x) | A] = \int_{-\infty}^{\infty} g(x)p(x | A)dx$$

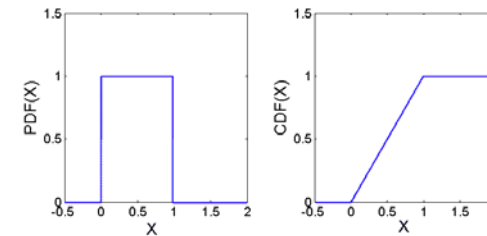
Conditional Expected Value



Cumulative Distribution Function

- The Cumulative Distribution Function (CDF) $p_{x \leq}$ provides the probability $P(X \leq x_0)$ of the event that the value of an R.V. X does not exceed x_0

$$p_{x \leq}(x_0) = \int_{-\infty}^{x_0} p(x)dx$$



Gaussian (Normal) Random Variable

- X is normal (Gaussian): $X \sim N(\mu, \sigma^2)$

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

$$E[x] = \mu$$

$$\text{Var}[x] = \sigma^2$$

- X is Standard normal (Standard Gaussian):
 $X \sim N(0,1)$ when $\mu=0, \sigma^2=1$

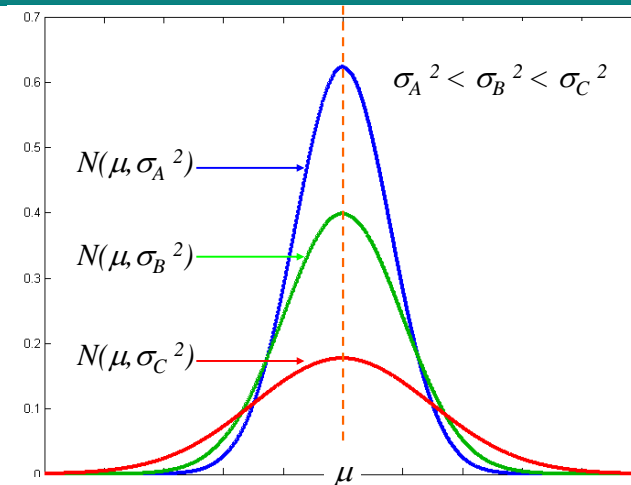
$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-(x-\mu)^2/2}$$

$$E[x] = 0$$

$$\text{Var}[x] = 1$$



Normal PDF





Standard Normal Variable

- Normality is preserved by linear transformation. Calculation involving the normal variable is usually done in terms of standard normal.
- Let $Y=aX+b$,
if $X \sim N(\mu, \sigma^2) \rightarrow Y \sim N(a\mu+b, a^2\sigma^2)$
- Let $Z=(X-\mu)/\sigma$,
if $X \sim N(\mu, \sigma^2) \rightarrow Z \sim N(0,1)$: Standard Normal



Standard Normal CDF

- X is a normal R.V.
- If we want to find $P(X \leq x_o)$, convert to a standard normal Z and then use Standard Normal CDF. (tabulated in the form of Φ)

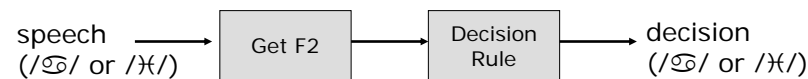
$$P(X \leq x_o) = P(\sigma Z + \mu \leq x_o) = P(Z \leq \frac{x_o - \mu}{\sigma}) = \Phi(\frac{x_o - \mu}{\sigma})$$

- $\Phi(z_o)$ is the total area under the standard normal PDF curve from $-\infty$ to z_o

$$\Phi(z_o) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{z_o} e^{-\frac{t^2}{2}} dt$$



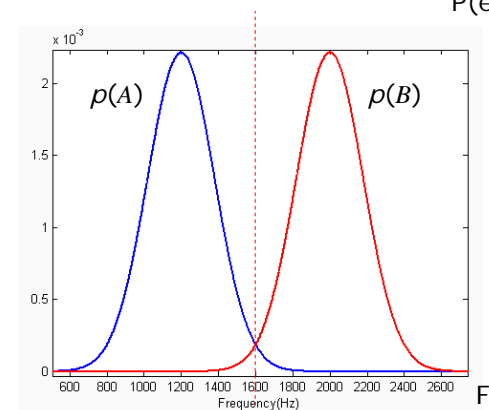
Example



- Let
 - $A = F2$ of $/\text{ɔ̃}/$
 - $B = F2$ of $/\text{ɤ}/$
- Suppose:
 - A is normal with mean 1200Hz and SD 180
 - B is normal with mean 2000 Hz and SD 180
- We want to classify an incoming vowel by the following rule:
 - if $(F2 > 1600\text{Hz})$ choose $/\text{ɤ}/$; else choose $/\text{ɔ̃}/$;
- If both vowels are equally likely to be produced, what is $P(\text{error})$?



Example



$$\begin{aligned}
 P(\text{err}) &= P(\text{choose } /a/ \mid /i/)P(/i/) \\
 &\quad + P(\text{choose } /i/ \mid /a/)P(/a/) \\
 &= P(F2 < 1600 \mid /i/) \times 0.5 \\
 &\quad + P(F2 > 1600 \mid /a/) \times 0.5 \\
 &= (1 - \Phi((1600 - 1200)/180)) \times 0.5 \\
 &\quad + (1 - \Phi((2000 - 1600)/180)) \times 0.5 \\
 &= (1 - \Phi(2.2)) \times 0.5 \\
 &\quad + (1 - \Phi(2.2)) \times 0.5 \\
 &= 1 - \Phi(2.2) = 0.0139
 \end{aligned}$$

Joint PDF

- Since we hardly use 1-dimensional feature vector, we mostly use joint PDFs to model the value distributions of multiple random variables.

Joint PDF

- $p_{XY}(x,y)$ is a joint density of probability measure on the event space for random variables X and Y

$$P(A) = \int \int_A p_{XY}(x, y) dx dy$$

$$p_X(x) = \int_{-\infty}^{\infty} p_{XY}(x, y) dy$$

$$p_Y(y) = \int_{-\infty}^{\infty} p_{XY}(x, y) dx$$

Independence

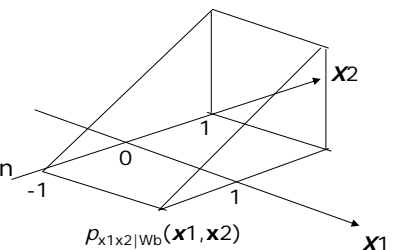
- When $p(X|Y) = p(X)$, it is said that the R.V. X is independent of Y . This leads to $p(XY) = p(X)p(Y)$
- E.g.:
 - A speaker says a phrase with two words
 - X is the duration in ms. of the first word.
 - Y is the duration in ms. of the second word.
 - if we do not know the total duration, knowing Y does not affect our knowledge of X
 - $p(X|Y) = p(X)$
 - Suppose we want to find $P(X>500, Y<500) \rightarrow$ multiply $P(X>500)$ and $P(Y<600)$ directly.

Example



- Vocab: W_a, W_b
- Feature vector: $(\mathbf{X1} = x_1, \mathbf{X2} = x_2)$
- Decision rule: Pick W_i giving maximal $p_{\mathbf{x1}, \mathbf{x2} | W_i}(x_1, x_2 | W_i)$

- Models: (from training data)
 - $p_{\mathbf{x1} | W_a} \sim$ Uniform over $(-1, 1)$
 - $p_{\mathbf{x2} | W_a} \sim$ Uniform over $(-1, 1)$
(Given W_a , $\mathbf{X1}$ and $\mathbf{X2}$ are independent.)
 - $p_{\mathbf{x1}, \mathbf{x2} | W_b}(\mathbf{X1}, \mathbf{X2} | W_b) \sim$ as shown





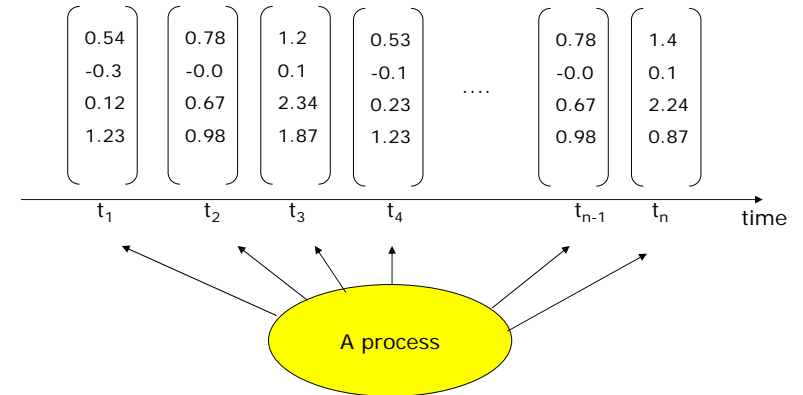
Example

- 4 test speech tokens:
 – (0.5,-0.3), (0.2,0,2), (-0.8,0.7), (-0.5,-0.5)
- What should the decisions be?
- What if the decision rule is changed to:
 Pick W_i giving maximal $p_{W_i|x_1x_2}(W_i|x_1,x_2)$
 and we know that W_b occurs twice as often as W_a ?



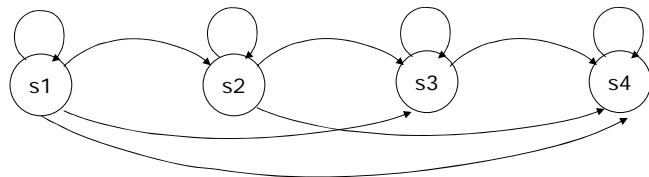
Hidden Markov Models

- HMMs are statistical models for processes that generate sequence of numbers.



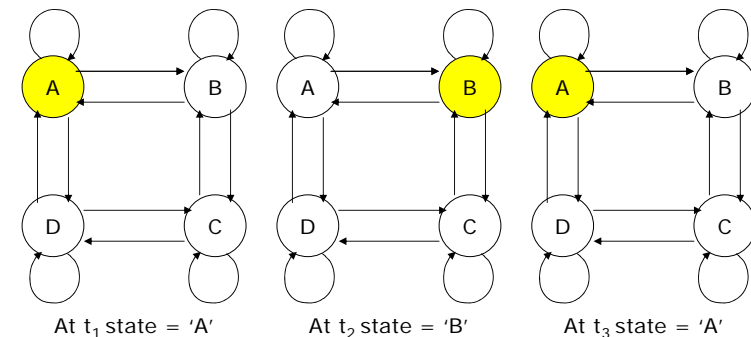
Hidden Markov Models

- Each process is described using a "Finite State Machine".



What is HMM?

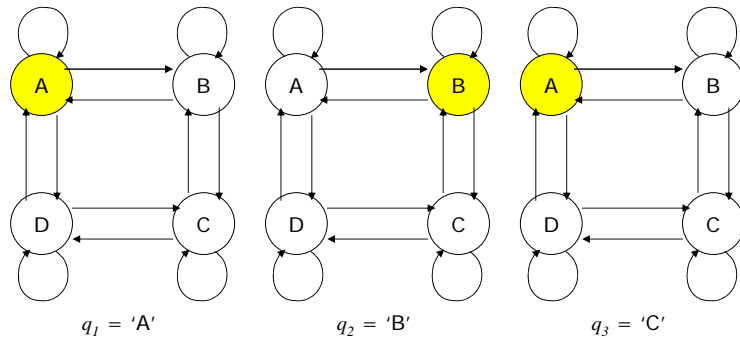
- The model assumes that underlying process can be in one of a number of states at any time instant.





What is HMM?

- Let q_i = the state that the process is in at time i

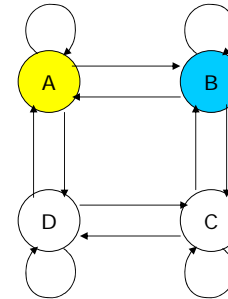


State sequence $Q = \{q_1, q_2, q_3, \dots, q_n\}$



What is HMM?

- The state that the process is in at any time instant is dependent only on the state it was in at the previous instant. (Markovian)



$$P(q_t | q_{t-1}) = P(q_t | q_{t-1}q_{t-2}q_{t-3}\dots q_0)$$



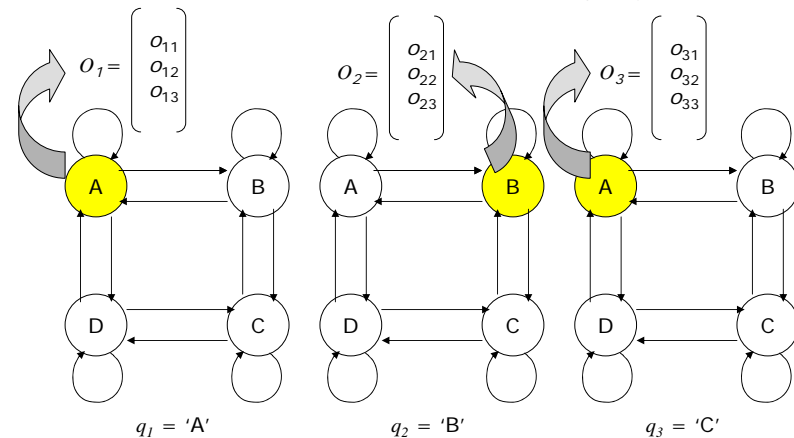
What is HMM?

- At each time instant the process generates an observation using probability distribution that is specific to the state that it is in.



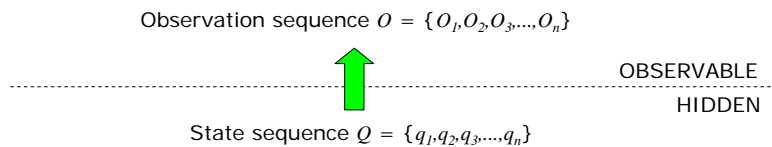
What is HMM?

Observation sequence $O = \{O_1, O_2, O_3, \dots, O_n\}$



What is HMM?

- The generated observations are all that we can see.
- The actual state sequences that happened are not directly observable.

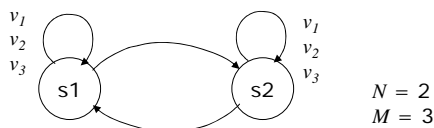


Discrete and Continuous Models

- HMMs can be classified into **discrete** models or **continuous** models
 - observable events assigned to each state are discrete, e.g. codewords after vector quantization
→ discrete model
 - observable events assigned to each state are continuous, e.g. many acoustic measurements without vector quantization
→ continuous model

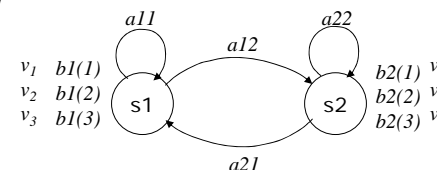
Elements of an HMM

- N : number of states in the model
 - states, $S = \{s_1, s_2, \dots, s_N\}$
 - state at time t , $q_t \in S$
- M : number of observation symbols (assume discrete observation)
 - observation symbols, $V = \{v_1, v_2, \dots, v_M\}$
 - observation at time t , $o_t \in V$



Elements of an HMM

- $A = \{a_{ij}\}$: state transition probability distribution
 - $a_{ij} = P(q_{t+1}=s_j | q_t=s_i); 1 \leq i, j \leq N$
- $B = \{b_j(k)\}$: observation symbol probability distribution in state j
 - $b_j(k) = P(v_k \text{ at } t | q_t=s_j); 1 \leq j \leq N, 1 \leq k \leq M$
- $\pi = \{\pi_i\}$: initial state distribution
 - $\pi_i = P(q_1=s_i); 1 \leq i \leq N$



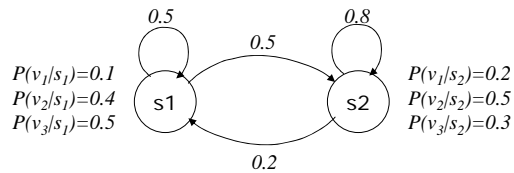


Example

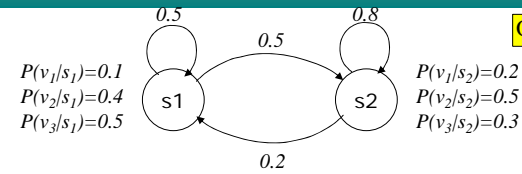
$$A = \begin{bmatrix} a11 & a12 \\ a21 & a22 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.2 & 0.8 \end{bmatrix}$$

$$B = \begin{bmatrix} b1(v1) & b1(v2) & b1(v3) \\ b2(v1) & b2(v2) & b2(v3) \end{bmatrix} = \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}$$

$$\pi = \{a01, a02\} = \{0.5, 0.5\}$$



Example

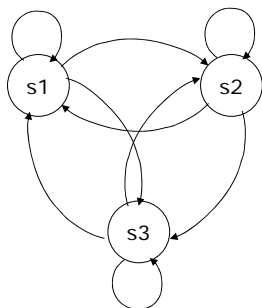


Observe $O = \{v2, v1, v3\}$, $P(O | A, B, \pi) = ?$

Possible state sequences, Q	$P(O A, B, \pi, q)$
{s1, s1, s1}	$0.5 \times 0.4 \times 0.5 \times 0.1 \times 0.5 \times 0.5$
{s1, s1, s2}	$0.5 \times 0.4 \times 0.5 \times 0.1 \times 0.5 \times 0.3$
{s1, s2, s1}	$0.5 \times 0.4 \times 0.5 \times 0.2 \times 0.2 \times 0.5$
{s1, s2, s2}	$0.5 \times 0.4 \times 0.5 \times 0.2 \times 0.8 \times 0.3$
{s2, s1, s1}	$0.5 \times 0.5 \times 0.2 \times 0.1 \times 0.5 \times 0.5$
{s2, s1, s2}	$0.5 \times 0.5 \times 0.2 \times 0.1 \times 0.5 \times 0.3$
{s2, s2, s1}	$0.5 \times 0.5 \times 0.8 \times 0.2 \times 0.2 \times 0.5$
{s2, s2, s2}	$0.5 \times 0.5 \times 0.8 \times 0.2 \times 0.8 \times 0.3$



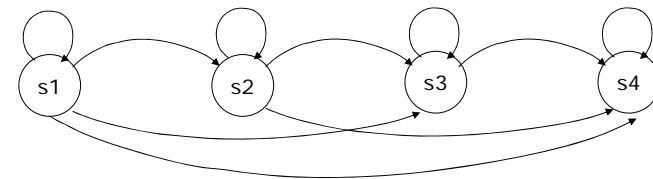
Ergodic Model



- An HMM is called ergodic or fully-connected model when every state of the model can be reached (in a single step) from every other state of the model.



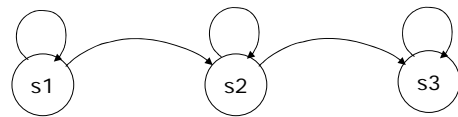
Left-to-Right Model



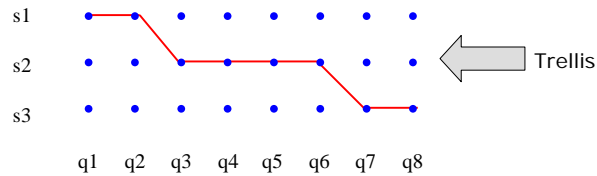
- The underlying state sequence associated with the model has the property that, as time increases, the state index increases.
- The left-to-right model exhibits the desirable property of being readily able to model speech whose properties change over time in a successive manner.



Representing State Diagram by Trellis

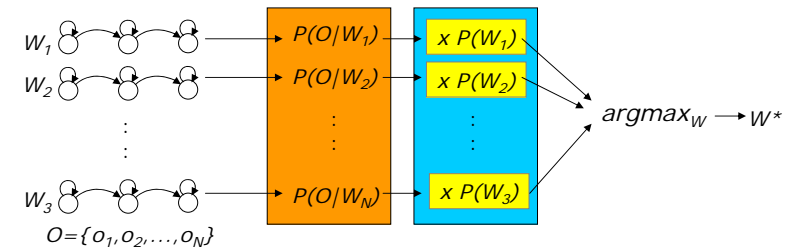


$$Q = \{ s_1, s_1, s_2, s_2, s_2, s_2, s_3, s_3 \}$$

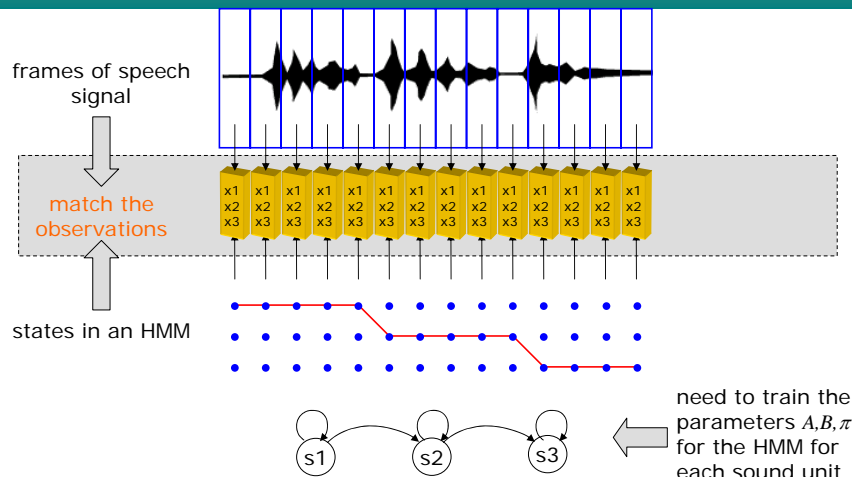


Modeling words with HMMs

- We seek to find $W^* = \text{argmax}_W(P(W/O))$
- $W^* = \text{argmax}_W(P(O/W)P(W)/P(O))$
- W is an element of set $\{W_1, W_2, W_3, \dots, W_N\}$
- So, we compare $P(O/W_i), i=1, 2, \dots, N_i$ and pick W_i that gives maximal $P(O/W_i)$.



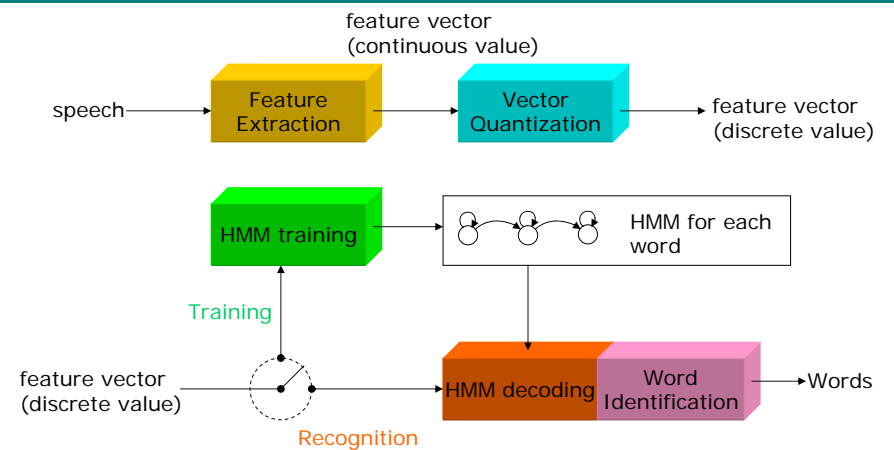
Acoustic Modeling using HMM



need to train the parameters A, B, π for the HMM for each sound unit



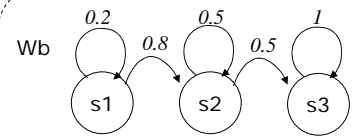
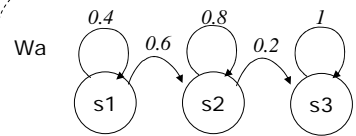
Speech Recognizer Based on discrete HMMs





Example

- 1D feature vector consisting of symbols {x,y,z}
- Vocab: {Wa,Wb}
- Observations: O = {x, x, z, y}
- Models: (already obtained from training)



$P(x/s_1)=0.2$	$P(x/s_2)=0.8$	$P(x/s_3)=0.3$	$P(x/s_1)=0.1$	$P(x/s_2)=0.3$	$P(x/s_3)=0.1$
$P(y/s_1)=0.3$	$P(y/s_2)=0.1$	$P(y/s_3)=0.2$	$P(y/s_1)=0.2$	$P(y/s_2)=0.3$	$P(y/s_3)=0.5$
$P(z/s_1)=0.5$	$P(z/s_2)=0.1$	$P(z/s_3)=0.5$	$P(z/s_1)=0.7$	$P(z/s_2)=0.4$	$P(z/s_3)=0.4$

- Decision rule: Pick W_i that maximize $P(O|W_i)$
- Constraints: 1st observation generated by s_1 . The last by s_3 .



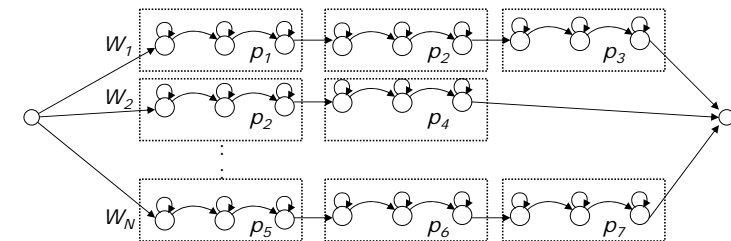
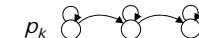
Modeling subword units with HMMS

- Suppose:

$$W_1 = p_1 p_2 p_3$$

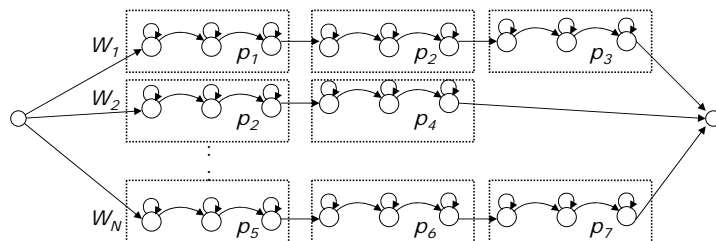
$$W_2 = p_2 p_4$$

$$W_N = p_5 p_6 p_7$$



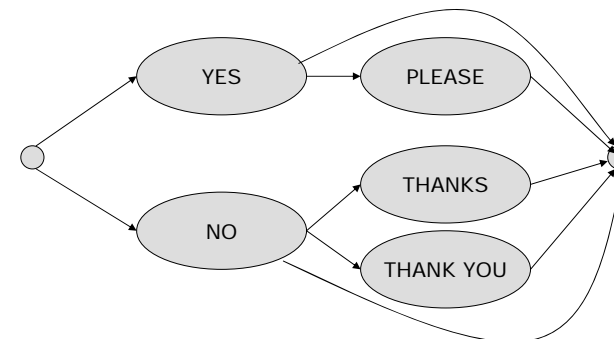
Modeling subword units with HMMS

- Find the best state sequence $Q = \{q_1, q_2, q_3, \dots, q_n\}$



Continuous Speech Recognizer

- Suppose a user can say:
"Yes" / "Yes, please" / "No" / "No. Thanks" / "No. Thank you"
- Find the best state sequence $Q = \{q_1, q_2, q_3, \dots, q_n\}$





Continuous Speech Recognizer

- Suppose the system need to recognize a bit string without knowing its length.

