

Hidden Markov Models: lecture 8

Application to speech recognition

Xavier Didelot

Introduction to speech recognition

- ▶ Understanding spoken language understanding is a difficult task, and it is remarkable that humans do as well at it as we do
- ▶ The goal of automatic speech recognition research is to address this problem
- ▶ Hidden Markov Models have a long successful history of being used in this field computationally by building systems that map from an acoustic signal to a string of words
- ▶ Comparison of error rates:

Task	Vocabulary	Machine	Human
Digits	11	0.5	0.009
Wall Street Journal clean	5k	3	0.9
Wall Street Journal noisy	5k	9	1.1
Broadcast new	65k	20	4

Prior probability of a word

- ▶ Even before we hear the beginning of a word, not all words are equally likely...
- ▶ An N-gram probability is the conditional probability of a word given the previous N-1 words
- ▶ N-gram probabilities can be computed by simply counting in a corpus and normalizing (the Maximum Likelihood Estimate) or they can be computed by more sophisticated algorithms
- ▶ The advantage of N-grams is that they take advantage of lots of rich lexical knowledge
- ▶ A disadvantage for some purposes is that they are very dependent on the corpus they were trained on

Bigram

Bigram for eight of the words (out of $V = 1446$) in a corpus of 9332 sentences.

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Random sentences from Shakespeare's work

Unigram

To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have
Every enter now severally so, let
Hill he late speaks; or! a more to leg less first you enter
Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like

Bigram

What means, sir. I confess she? then all sorts, he is trim, captain.
Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.
What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?

Trigram

Sweet prince, Falstaff shall die. Harry of Monmouth's grave.
This shall forbid it should be branded, if renown made it empty.
Indeed the duke; and had a very good friend.
Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

Quadrigram

King Henry.What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;
Will you not tell me who I am?
It cannot be but so.
Indeed the short and the long. Marry, 'tis a noble Lepidus.

Random sentences from Wall Street Journal

Unigram

Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

Bigram

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

Trigram

They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

Phones

- ▶ We can represent the pronunciation of words in terms of units called phones
- ▶ The standard system for representing phones is the International Phonetic Alphabet or IPA

International Phonetic Alphabet (IPA) ˌɪntəˈnæʃnəl fəˈnetɪk ˈælfəbet

Consonants (pulmonic)

	Bilabial	Labio-dental	Dental	Alveolar	Post-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ		
Trill	ʙ			r					ʀ		
Tap or flap		ⱱ		ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral fricative				ɬ ɮ							
Approximant		ʋ		ɹ		ɻ	j	ɰ			
Lateral approximant				l		ɭ	ʎ	ʟ			

- ▶ The most common computational system for transcription of English is the ARPAbet, which conveniently uses ASCII symbols

Phones for consonants

ARPAbet Symbol	IPA Symbol	Word	ARPAbet Transcription
[p]	[p]	<u>p</u> arsley	[p aa r s l iy]
[t]	[t]	<u>t</u> ea	[t iy]
[k]	[k]	<u>c</u> ook	[k uh k]
[b]	[b]	<u>b</u> ay	[b ey]
[d]	[d]	<u>d</u> ill	[d ih l]
[g]	[g]	<u>g</u> arlic	[g aa r l ix k]
[m]	[m]	<u>m</u> int	[m ih n t]
[n]	[n]	<u>n</u> utmeg	[n ah t m eh g]
[ŋg]	[ŋ]	bak <u>ing</u>	[b ey k ix ŋg]
[f]	[f]	<u>f</u> lour	[f l aw axr]
[v]	[v]	clo <u>v</u> e	[k l ow v]
[θ]	[θ]	<u>th</u> ick	[θ ih k]
[ðh]	[ð]	<u>th</u> ose	[ðh ow z]
[s]	[s]	<u>s</u> oup	[s uw p]
[z]	[z]	egg <u>s</u>	[eh g z]
[ʃh]	[ʃ]	squash <u>h</u>	[s k w aa sh]
[zh]	[ʒ]	ambros <u>ia</u>	[ae m b r ow zh ax]
[ch]	[tʃ]	<u>ch</u> erry	[ch eh r iy]
[jh]	[dʒ]	<u>j</u> ar	[jh aa r]
[l]	[l]	li <u>c</u> orice	[l ih k axr ix sh]
[w]	[w]	ki <u>w</u> i	[k iy w iy]
[r]	[r]	ri <u>c</u> e	[r ay s]
[y]	[j]	ye <u>ll</u> ow	[y eh l ow]
[h]	[h]	<u>h</u> oney	[h ah n iy]
Less commonly used phones and allophones			
[q]	[ʔ]	<u>uh</u> -oh	[q ah q ow]
[dx]	[ɾ]	bu <u>tt</u> er	[b ah dx axr]
[nx]	[ɹ̥]	win <u>n</u> er	[w ih nx axr]
[el]	[l̥]	ta <u>bl</u> e	[t ey b el]

Phones for vowels

ARPabet Symbol	IPA Symbol	Word	ARPabet Transcription
[iy]	[i]	lily	[l ih l iy]
[ih]	[ɪ]	lily	[l ih l iy]
[ey]	[eɪ]	da <u>isy</u>	[d ey z iy]
[eh]	[ɛ]	pe <u>n</u>	[p eh n]
[ae]	[æ]	a <u>ster</u>	[æ s t axr]
[aa]	[ɑ]	po <u>ppy</u>	[p aa p iy]
[ao]	[ɔ]	o <u>r</u> chid	[ao r k ix d]
[uh]	[ʊ]	wo <u>o</u> d	[w uh d]
[ow]	[oʊ]	lot <u>u</u> s	[l ow dx ax s]
[uw]	[u]	tu <u>l</u> ip	[t uw l ix p]
[ah]	[ʌ]	bu <u>tt</u> erc <u>u</u> p	[b ah dx axr k ah p]
[er]	[ɜ]	bi <u>r</u> d	[b er d]
[ay]	[aɪ]	i <u>r</u> is	[ay r ix s]
[aw]	[aʊ]	sun <u>fl</u> ow <u>e</u> r	[s ah n f l aw axr]
[oy]	[oɪ]	so <u>i</u> l	[s oy l]
Reduced and uncommon phones			
[ax]	[ə]	lot <u>u</u> s	[l ow dx ax s]
[axr]	[ɜ]	hea <u>th</u> er	[h eh dh axr]
[ix]	[i]	tu <u>l</u> ip	[t uw l ix p]
[ux]	[ʊ]	du <u>d</u> e ¹	[d ux d]

Words are not deterministic sequences of phones

- For example there are several pronunciations possible for “because” and “about”:

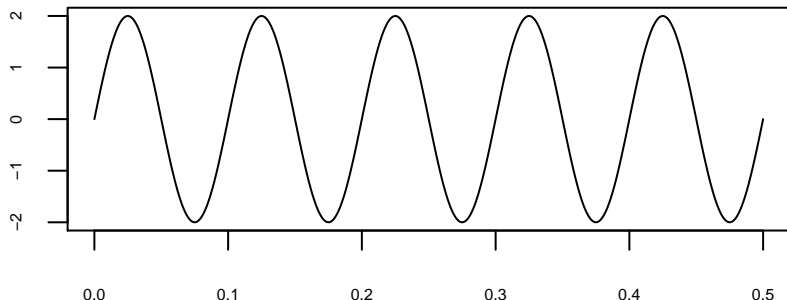
because				about			
ARPAbet	%	ARPAbet	%	ARPAbet	%	ARPAbet	%
b i y k a h z	27%	k s	2%	a x b a w	32%	b a e	3%
b i x k a h z	14%	k i x z	2%	a x b a w t	16%	b a w t	3%
k a h z	7%	k i h z	2%	b a w	9%	a x b a w d x	3%
k a x z	5%	b i y k a h z h	2%	i x b a w	8%	a x b a e	3%
b i x k a x z	4%	b i y k a h s	2%	i x b a w t	5%	b a a	3%
b i h k a h z	3%	b i y k a h	2%	i x b a e	4%	b a e d x	3%
b a x k a h z	3%	b i y k a a z	2%	a x b a e d x	3%	i x b a w d x	2%
k u h z	2%	a x z	2%	b a w d x	3%	i x b a a t	2%

Waves

- ▶ Acoustic analysis is based on trigonometric functions.

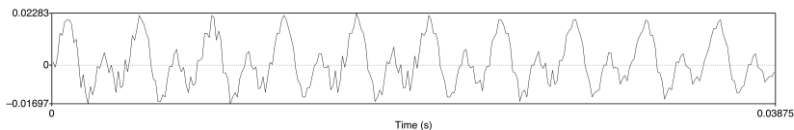
$$y = A\sin(2\pi ft)$$

- ▶ x is time measured in seconds
- ▶ y represents air pressure above and below the normal atmospheric pressure
- ▶ A is the signal amplitude
- ▶ f is the signal frequency
- ▶ For example if $A = 2$ and $f = 10$ Hertz:



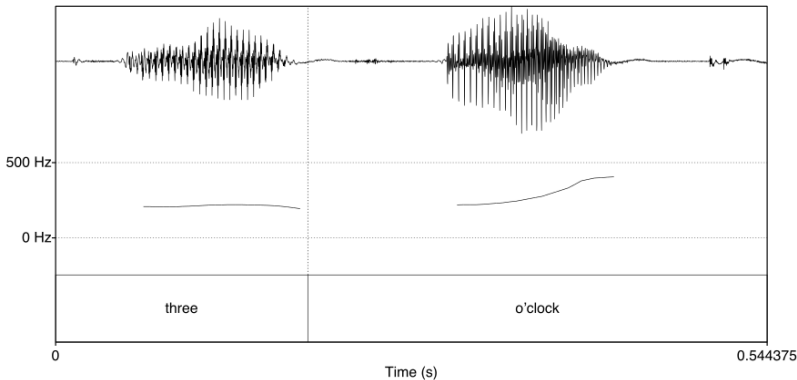
Real sounds

- ▶ Real spoken sounds are complex series of change in air pressure
- ▶ These changes in air pressure are caused by the specific way that air passes through the glottis and out the oral or nasal cavities
- ▶ For example the vowel [iy] from someone saying “She just had a baby”



Fundamental frequency

- ▶ Real sounds are more complex than the simple sine, but the signal still looks periodic with a frequency which is called the fundamental frequency
- ▶ Perceived pitch is (roughly) proportional to the fundamental frequency
- ▶ The fundamental frequency varies over time:

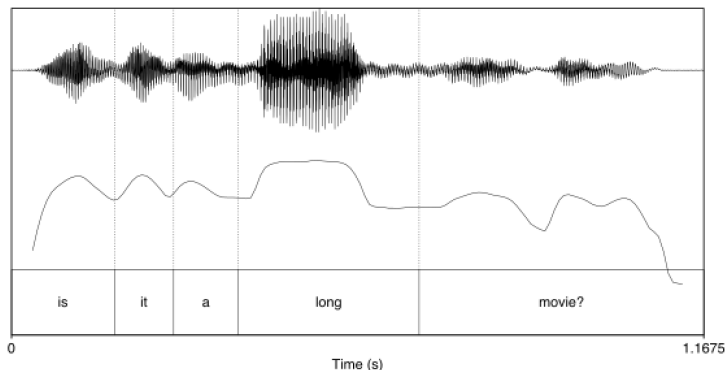


Power and intensity over time

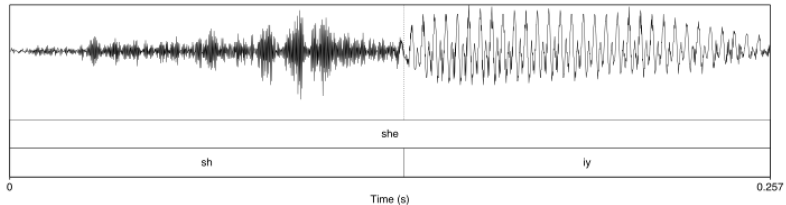
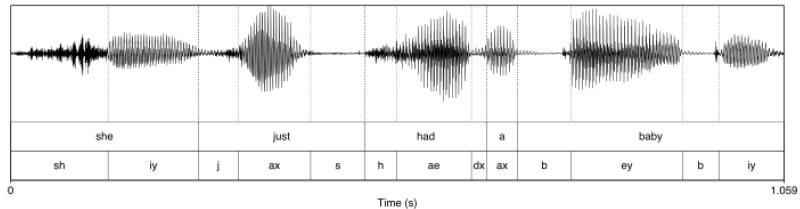
- ▶ The power and intensity of the signal are defined as:

$$P = \frac{1}{N} \sum_{i=1}^N x_i^2 \text{ and } I = 10 \log_{10}(P/P_0)$$

- ▶ Perceived loudness is (roughly) proportional to power
- ▶ Variation of intensity over time:

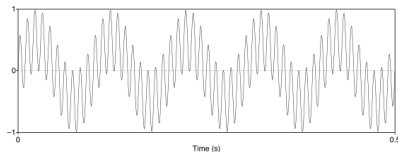


She just had a baby

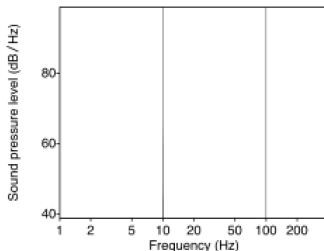


Fourier analysis and spectrum

- ▶ Idea: decompose complex waves into a sum of many sine waves of different frequencies
- ▶ For example the sum of two sine waveforms, one of frequency 10 Hz and one of frequency 100 Hz, both with amplitude 1:

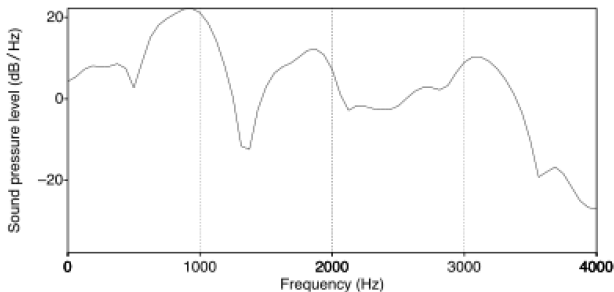
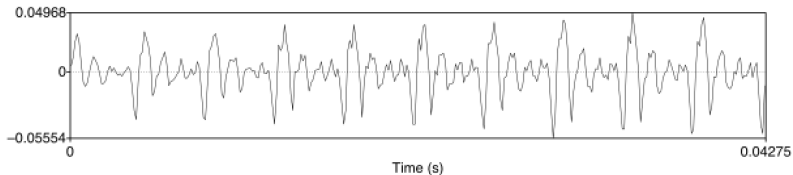


- ▶ We can represent the two component frequencies with a spectrum.



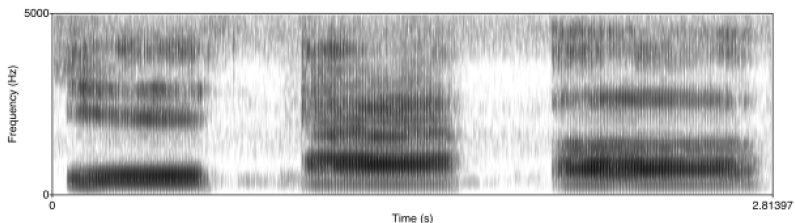
Discrete Fourier Transform

- ▶ Discrete Fourier Transform (DFT) computes a spectrum
- ▶ For example the vowel [ae] from the word “had”



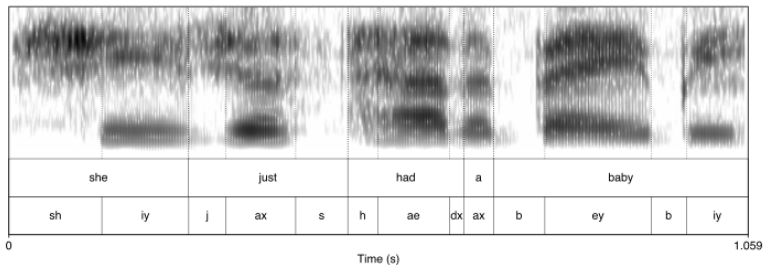
Spectrogram

- ▶ A spectrogram shows how the spectrum changes over time
- ▶ The x-axis is time
- ▶ The y-axis is the frequency shown on the x-axis of the spectrum
- ▶ The darkness represents the amplitude shown on the y-axis of the spectrum
- ▶ For example [ih], [ae], and [uh]

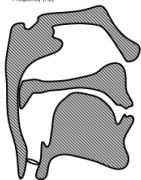
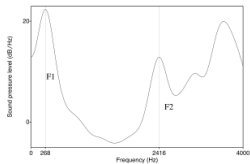


Formants

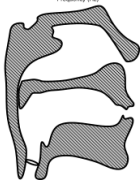
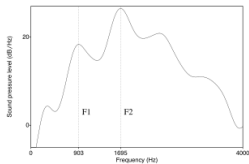
- ▶ Black bars on spectrograms correspond to peaks in spectrums and are called formants
- ▶ “She just had a baby”



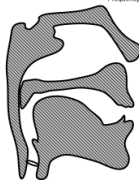
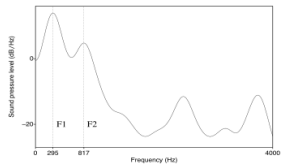
Comparing three phones



[iy] (tea)

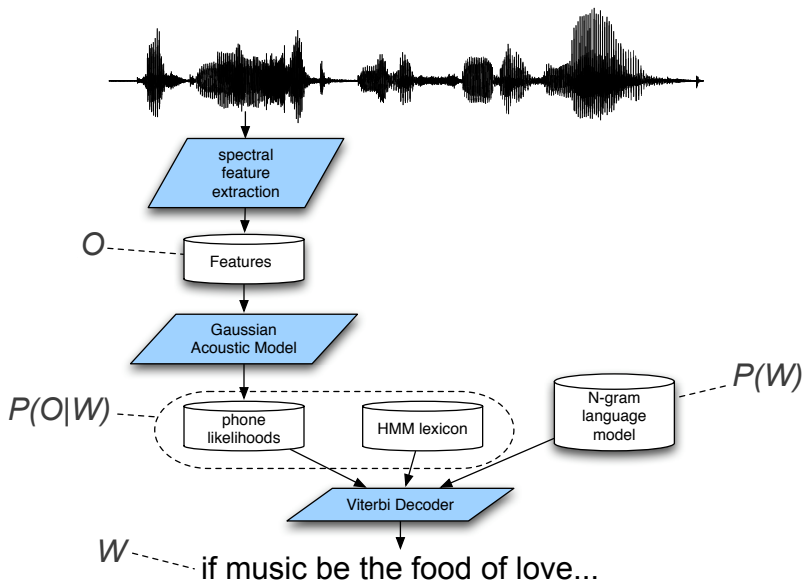


[ae] (cat)



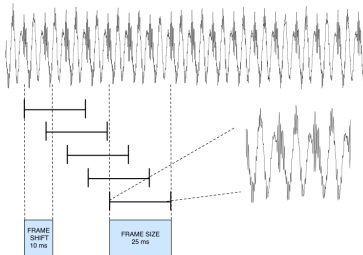
[uw] (moo)

Schematic for speech recognition of a sentence

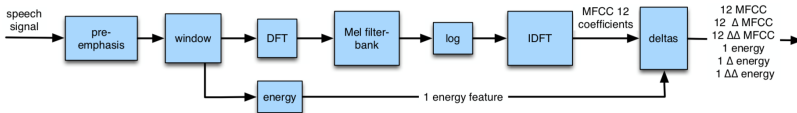


Features extraction

- ▶ Every $\sim 10\text{ms}$ a spectrum is computed based on the next $\sim 25\text{ms}$

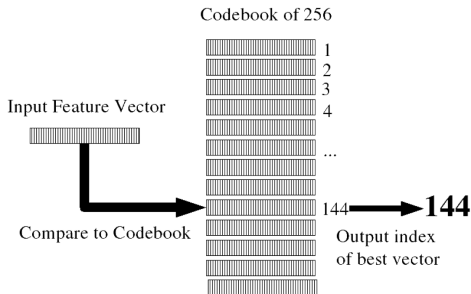


- ▶ From each spectrum, 39 features are computed (MFCC)



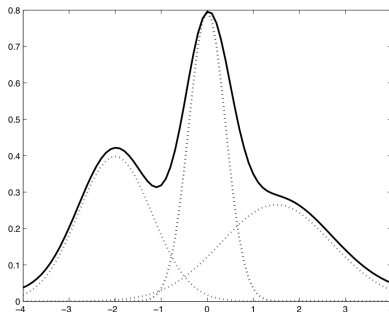
Vector quantization

- ▶ Old method, no longer used
- ▶ Compare feature vector to codebook
- ▶ Given an observation, select nearest vector from codebook



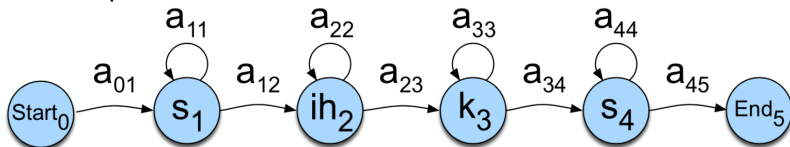
Gaussian acoustic model

- ▶ Consider that each feature is Normally distributed
- ▶ Learn parameters from training set
- ▶ Multivariate Gaussian: consider covariance between features
- ▶ Even better: mixture of Gaussians
- ▶ GMM: Gaussian Mixture Model
- ▶ Training similar to Baum-Welch algorithm



Building HMMs for speech recognition

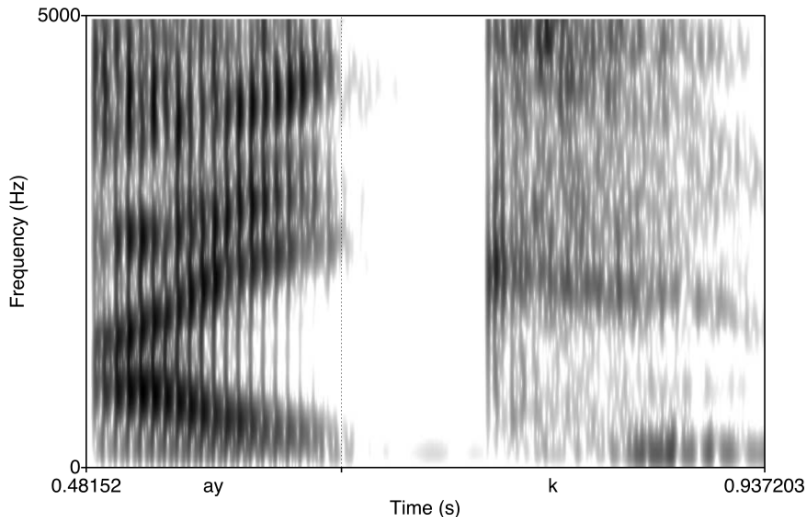
- ▶ Since words are made of phones, we consider a HMM in which each state is a phone
- ▶ For example for the word “six”:



- ▶ Note the self-loops, which allow a single phone to repeat so as to cover a variable amount of the acoustic input
- ▶ Phone durations vary hugely, dependent on the phone identify, the speaker's rate of speech, the phonetic context, and the level of prosodic prominence of the word

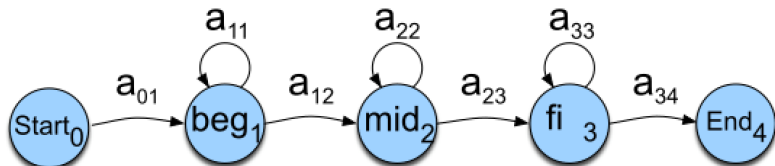
Spectrum variation within a phone

- ▶ The spectrum of a phone is not always constant, and the variation is sometimes really important
- ▶ For example the word “Ike” pronounced [ay k]:

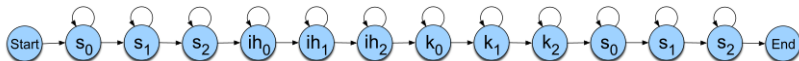


Spectrum variation within a phone

- To account for this, we consider that a phone is made of three states:



- For the word “six”, this gives:

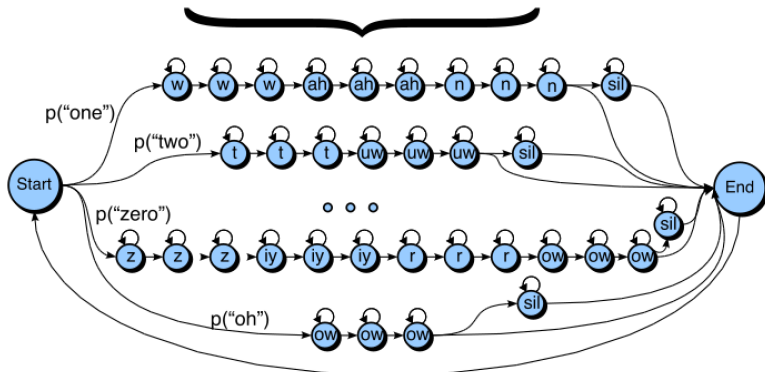
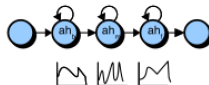


HMM for the digit recognition task

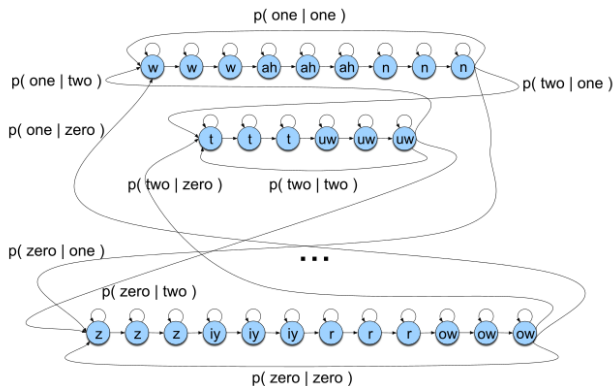
Lexicon

one	w ah n
two	t uw
three	th r iy
four	f ao r
five	f ay v
six	s ih k s
seven	s eh v ax n
eight	ey t
nine	n ay n
zero	z iy r ow
oh	ow

Phone HMM



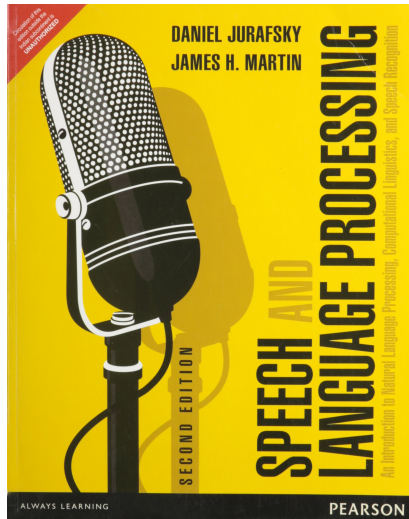
Bigram grammar network for the digit recognition task



HMM algorithms

- ▶ Now that we have a well defined HMM, we can use the algorithms we have learnt in previous lectures
- ▶ Training using the Baum-Welch algorithm on a large number of known words
- ▶ Forward algorithm to calculate $p(O|W)$ for several words
- ▶ Viterbi algorithm for global decoding

To find out more...



Conclusions

- ▶ The input to a speech recognizer is a series of acoustic waves
- ▶ Waveform, spectrum and spectrogram are useful visualization tools
- ▶ In the first step in speech recognition, sound waves are sampled, quantized, and converted into spectral feature vectors (often MFCC)
- ▶ GMM acoustic models are used to estimate the phonetic likelihoods of these feature vectors for each frame
- ▶ HMM models combine phone (or parts of phones) into words
- ▶ Can also account for independence of successive words via N-gram
- ▶ Efficiency of HMM algorithms is key to application in automatic speech recognition