



January 17, 2020 Professor Marie Meteer Brandeis University

+ Today

- The speech problem
- Units of speech
- Hidden Markov Models
- Phonetic HMMs
- Recognition architecture
 - Training
 - Decoding

+ 1985: First significant breakthrough in Speech Recognition

- Hidden Markoff Models
 - Mathematical framework
 - Ability to model time and spectral variability simultaneously
 - Ability to automatically estimate parameters given data
 - Not longer need to hand segment into phonemes
 - Segmentation and modeling done in one step
 - Data driven \rightarrow Standard scientific procedures
 - Empirical!

A Decade of Progress in Speech Recognition



+ The speech problem





Sequence of discrete entities

ih	t	S	iy	z	iy	t	u	r	eh	k	0	g	n	iy	z	s	р	iy	ch
It's easy to				recognize							speech								

• Or did she say

It's easy to wreck a nice beach?

+ Challenge: Variability

- Linguistic
 - Can say many different things
 - Phonetics, phonology, syntax, semantic, discourse

Speaker

- Physical characteristics of speaker
- Co-articulation (mouth has to transition between sounds
- Native language/dialect

Channel

- Background noise
- Transmission channel (microphone/telephone quality)

6

+ A look at the speech sounds

Grey whales



Words	Grey			Whales	S	meov	v		
Phonemes	g	r	еу	w	еу	1	z	?	
Triphones	- g r	g r ey	r ey w	ey W y	weyı	ey I z	Z -		





Search through space of all possible sentences.

Pick the one that is most probable given the waveform.

+ Design Intuition

- Build a statistical model of the speech sounds-to-words
 - Collect lots and lots of speech, and transcribe all the words.
 - Train the model on the labeled speech
- Create a "search space" of all the possible word combinations
 - Write a grammar
 - "Learn" likely sequences from lots of text
- Paradigm: Supervised Machine Learning + Search
 - Use the model to find the best sequence of words given the input sounds



+ The Noisy Channel Model (Preview)

- What is the most likely sentence out of all sentences in the language L given some acoustic input O?
- Treat acoustic input O as sequence of individual observations
 - $O = o_1, o_2, o_3, \dots, o_t$
- Define a sentence as a sequence of words:
 - $W = w_1, w_2, w_3, \dots, w_n$

+ Noisy Channel Model

Probabilistic implication: Pick the highest prob S:

$$\hat{W} = \underset{W \in L}{\operatorname{argmax}} P(W \mid O)$$

• We can use Bayes rule to rewrite this:

$$\hat{W} = \underset{W \in L}{\operatorname{argmax}} \frac{P(O | W)P(W)}{P(O)}$$

Since denominator is the same for each candidate sentence W, we can ignore it for the argmax:

$$\hat{W} = \underset{W \in L}{\operatorname{arg\,max}} P(O | W) P(W)$$

+ The noisy channel model

 Ignoring the denominator leaves us with two factors: P(Source) and P(Signal|Source)



+ Speech Recognition Architecture



+ How is this connected to language?

- The Dictionary
- Hand built knowledge source that ties words to sounds
- Sequence of words \rightarrow
 - sequence of phonemes \rightarrow
 - HMM states



16

+ From Phonetics to the Dictionary

- Dictionaries spell out the phonemes for each word
- Variations in pronunciation come from
 - Slight differences in shape of the vocal tract
 - Slight variations in articulation (accents)
 - Co-articulation with neighboring phonemes
 - Predictable variation from the position of phoneme in the word (e.g the "ps" in "tap" and "pat" are different)

Some words have multiple pronunciations >RECORD R-EH-K-AXR-D R-IX-K-AO-R-D

+ Base + Domain Dictionary

>ABOUND AX-B-AW-N-D >ABOUNDED AX-B-AW-N-D-IX-D >ABOUNDING AX-B-AW-N-D-IX-NX >ABOUNDS AX-B-AW-N-D-7 >ABOUT AX-B-AW-T >ATIS >ABOUT'S AX-B-AW-T-S >ABOVE AX-B-AH-V >ABOVEBOARD AX-B-AH-V-B-AO-R-D >ABPI ANAI P AE-B-P-L-AX-N-AE-L-P >ABRA AA-B-R-AX >ABRACADABRA AF-B-R-AX-K-AX-D-AF-B-R-AX >ABRAHAM EY-B-R-AX-HH-AE-M

>AIRFORCE
>APPROACHING
AX-P-R-OW-CH-IX-N
AX-P-R-OW-CH-IX-NX
P-R-OW-CH-IX-N
>ARAC
EY-R-AE-K
>ATIS
EY-T-IX-S
>AZIMUTH
AE-Z-M-EH-TH
>BLACKCAT
B-L-AE-K-AE-T
>CAIRNS
K-AE-R-IX-N-Z

© MM Consulting 2011

1/10/11

+ Architecture: Five easy pieces

- Feature extraction
- HMMs, Lexicons, and Pronunciation
- Acoustic Modeling
- Language Modeling
- Decoding

+ Speech Recognition Architecture



+ Speech Recognition Pipeline



+ Speech Recognition Knowledge Sources

Model of speech, pauses, coughs, and other sounds

List of all the words and their pronunciations, the "phonetic spelling"

Models the relationship between the sounds and the phonemes. Specific to a language (English or Spanish) and a channel (telephony or broadcast)

Grammar: all possible sentences

Statistical language model (SLM): Captures the likelihood of sequences of words

Thousands of hours of transcribed speech



Speech/

non-speech



CS 224S Winter 2007

1/23/20



© MM Consulting 2015

24

+ Embedded Training



CS 224S Winter 2007

1/23/20

+ Statistical Language Model (SLM)

- Captures how words are used in a particular domain
 - Specific to dialect ("in the hospital" vs. "in hospital")
 - Specific to domain (frequency of different words)
 - Disambiguate homophones ("disc" vs. "disk)
- But allows any word order
- Trigram model
 - Probability of a word given the previous two words

What is the missing word?

"I love ____"

NOW what is the missing word?

"My favorite TV show was I love _____"

© MM Consulting 2011

Language (or Domain) Model: Statistical Grammar

- Word order is guided by statistics
- Data is used to count the likelihood a word given the previous
 - "in the hospital"
 - "the in hospital"
 - "the hospital in"
- Creates a lattice that is searched
- Training data must be relevant!



Bigram

+ Back to HMMs

Why Markov model?

Markov chain models sequences

Transitions in the chain are \rightarrow probabilistic

- → Speech is sequential can be modeled as a sequence of states specified by the dictionary
 - Probabilities model uncertainly well

Why Hidden Markov model?

- Output symbols are probabilistic distribution over all labels
- The actual sequence of states for a particular output is "Hidden"
- There is one sequence that is the most probable to generate the output symbols

+ Hidden Markov Models formally

- States $Q = q_1, q_2...q_{N_i}$
- Observations O= o₁, o₂...o_{N;}
 Each observation is a symbol from a vocabulary V = {v₁, v₂,...v_V}
- Transition probabilities
 - Transition probability matrix $A = \{a_{ij}\}$
- Observation likelihoods
 Output probability matrix B={b_i(k)}
- Special initial probability vector π

$$a_{ij} = P(q_t = j | q_{t-1} = i) \quad 1 \le i, j \le N$$

$$b_i(k) = P(X_t = o_k \mid q_t = i)$$

$$\pi_i = P(q_1 = i) \quad 1 \le i \le N$$

1/23/20

29

+ Different types of HMM structure



Bakis = left-to-right (allowing skips)



Ergodic = fully-connected

+ HMMs for speech

Dictionary SIX SIHKS

State sequence for every word a_{11} a_{22} a_{33} a_{44} a_{45} a_{45}

Each phone has 3 subphones



End

+ HMM for digit recognition task





Search through space of all possible sentences
 Defined by the HMM

Pick the one that is most probable given the waveform. Based on the transition and output probabilities in the HMM

+ The Noisy Channel Model

- What is the most likely sentence out of all sentences in the language L given some acoustic input O?
- Treat acoustic input O as sequence of individual observations
 - $O = o_1, o_2, o_3, \dots, o_t$
- Define a sentence as a sequence of words:
 - $W = w_1, w_2, w_3, \dots, w_n$

+ The Evaluation (forward) problem for speech

The observation sequence O is a series of MFCC vectors

- The hidden states W are the phones and words
- For a given phone/word string W, our job is to evaluate P(O|W)
- Intuition: how likely is the input to have been generated by just that word string W

+ HMM for speech: Consider all different paths!

- fayayayayvvvv
- ffay ay ay ay v v v
- ffffay ay ay ay v
- ffay ay ay ay ay ay v
- ffay ay ay ay ay ay ay ay v
- ffayvvvvvv



+ The forward lattice for "five"

Computes all possible paths





+ Forward trellis for "five"

	V	0.0000	0.0000	0.0080	0.0093	0.0114	0.0070	0.0069	0.0036	0.0021	0.0011
ates	AY	0.0000	0.0400	0.0540	0.0664	0.0355	0.0160	0.0051	0.0016	0.0004	0.0001
	F	0.8000	0.3200	0.1120	0.0224	0.0045	0.0009	0.0002	0.0001	0.0000	0.0000
Sti		1	2	3	4	5	6	7	8	9	10
-	time	е									

"Emission" probabilities (observation likelihood for the observation *o* at each frame)

f	0.8	0.8	0.7	0.4	0.4	0.4	0.5	0.5	0.5	0.5
ay	0.1	0.1	0.3	0.8	0.8	0.8	0.8	0.6	0.5	0.4
v	0.6	0.6	0.4	0.3	0.3	0.3	0.3	0.6	0.8	0.9
р	0.4	0.4	0.2	0.1	0.1	0.1	0.1	0.1	0.3	0.3
iy	0.1	0.1	0.3	0.6	0.6	0.6	0.5	0.5	0.5	0.4

"Transition" probabilities

	F	AY	V
F	0.5	0.5	0.0
AY	0.0	0.5	0.5
V	0.0	0.0	0.5



States AX E 0.0000 0.064 0.0168 0.002688 0.0010752 0.00086016 0.000344064 0.000137626 0.00672 0.04 0.048 0.0448 0.0000 0.01792 0.007168 0.0028672 0.00021504 6.88128E-05 0.00086016 0.8000 0.32 0.112 0.0224 0.00896 0.003584 0.001792 0.0007168 0.00021504 8.6016e-05 3 5 7 9 2 4 8 10 1 6

time

"Emission" (observation) probabilities

+ Viterbi trellis for "five"

f	0.8	0.8	0.7	0.4	0.4	0.4	0.5	0.5	0.5	0.5
ay	0.1	0.1	0.3	0.8	0.8	0.8	0.8	0.6	0.5	0.4
v	0.6	0.6	0.4	0.3	0.3	0.3	0.3	0.6	0.8	0.9
р	0.4	0.4	0.2	0.1	0.1	0.1	0.1	0.1	0.3	0.3
iy	0.1	0.1	0.3	0.6	0.6	0.6	0.5	0.5	0.5	0.4

Transition probabilities

		F	AY	V
	F	0.5	0.5	0.0
	AY	0.0	0.5	0.5
Thanks			10.0	0,5
THATKS		III JULAISK	y for the	se shues



+ Viterbi backtrace



Thanks to Dan Jurafsky for these slides



+ "Output symbols": the great hack

- Markov models are "generative" models: Find the most likely sequence that generates the output symbols
- "Output symbols" = Observations
 - What we normally think of as input