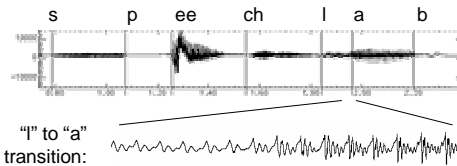


Speech Recognition: Acoustic Waves

- Human speech generates a wave
 - like a loudspeaker moving
- A wave for the words "speech lab" looks like:



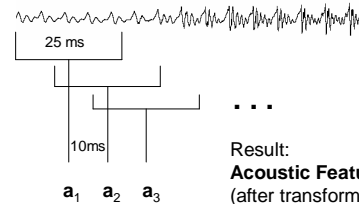
Graphs from Simon Amfield's web tutorial on speech, Sheffield:
<http://www.psyc.leeds.ac.uk/research/cogn/speech/tutorial/>

June 1999

1

Acoustic Sampling

- 10 ms frame (ms = millisecond = 1/1000 second)
- ~25 ms window around frame [wide band] to allow/smooth signal processing – it let's you see formants



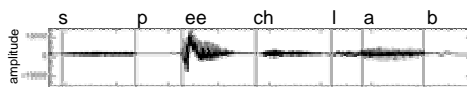
Result:
Acoustic Feature Vectors
 (after transformation,
 numbers in roughly \mathbb{R}^{14})

June 1999

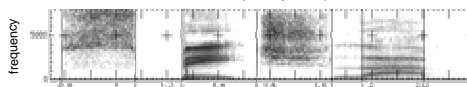
2

Spectral Analysis

- Frequency gives pitch; amplitude gives volume
 - sampling at ~8 kHz phone, ~16 kHz mic (kHz=1000 cycles/sec)



- Fourier transform of wave displayed as a spectrogram
 - darkness indicates energy at each frequency
 - hundreds to thousands of frequency samples



June 1999

3

The Speech Recognition Problem

- The **Recognition Problem: Noisy channel model**
 - We started out with English words, they were encoded as an audio signal, and we now wish to decode.
 - Find most likely sequence \mathbf{w} of "words" given the sequence of acoustic observation vectors \mathbf{a}

- Use Bayes' law to create a **generative model** and then decode
- $\text{ArgMax}_{\mathbf{w}} P(\mathbf{w}|\mathbf{a}) = \text{ArgMax}_{\mathbf{w}} P(\mathbf{a}|\mathbf{w}) P(\mathbf{w}) / P(\mathbf{a})$
 $= \text{ArgMax}_{\mathbf{w}} P(\mathbf{a}|\mathbf{w}) P(\mathbf{w})$

- Acoustic Model:** $P(\mathbf{a}|\mathbf{w})$

- Language Model:** $P(\mathbf{w})$ ←

A probabilistic theory of a language

June 1999

4

Probabilistic Language Models

- Assign probability $P(\mathbf{w})$ to word sequence $\mathbf{w} = w_1, w_2, \dots, w_k$
- Can't directly compute probability of long sequence – one needs to decompose it
- Chain rule provides a **history-based model**:

$$P(w_1, w_2, \dots, w_k) = P(w_1) P(w_2|w_1) P(w_3|w_1, w_2) \dots P(w_k|w_1, \dots, w_{k-1})$$
- Cluster** histories to reduce number of parameters
- E.g., just based on the last word (1st order Markov model):

$$P(w_1, w_2, \dots, w_k) = P(w_1|<s>) P(w_2|w_1) P(w_3|w_2) \dots P(w_k|w_{k-1})$$
- How do we estimate these probabilities?
 - We count word sequences in corpora
 - We "smooth" probabilities so as to allow unseen sequences

June 1999

5

N-gram Language Modeling

- n -gram assumption clusters based on last $n-1$ words
 - $P(w_i|w_1, \dots, w_{i-1}) \approx P(w_i|w_{i-n+1}, \dots, w_{i-1})$
 - unigrams $\sim P(w_i)$
 - bigrams $\sim P(w_i|w_{i-1})$
 - trigrams $\sim P(w_i|w_{i-2}, w_{i-1})$
- Trigrams often interpolated with bigram and unigram:

$$\hat{P}(w_3 | w_1, w_2) = \lambda_3 \frac{F(w_3 | w_1, w_2)}{\sum_k F(w_k | w_1, w_2)} + \lambda_2 \frac{F(w_3 | w_2)}{\sum_k F(w_k | w_2)} + \lambda_1 \frac{F(w_3)}{\sum_k F(w_k)}$$

- the λ_i typically estimated by maximum likelihood estimation on held out data ($F(\cdot, \cdot)$ are relative frequencies)
- many other interpolations exist (another standard is a non-linear **backoff**)

June 1999

6