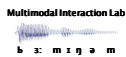


Word and Sub-Word HMMs



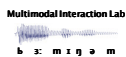
EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Objectives

- Word level HMMs
- Sub-word HMMs
- Context-sensitive sub-word HMMs
 - Biphone HMMs
 - Triphone HMMs
- Triphone HMM training issues
- Phoneme Decision Trees (PDTs)
- Notes: pp 43-46



EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Word Level HMMs

- Early systems (1980s) used word level HMMs
- I.e. each word modelled by a single, dedicated HMM (c.f. “zero” picture)
 - Advantages:
 - Good performance due to explicit modelling of word-dependent variability

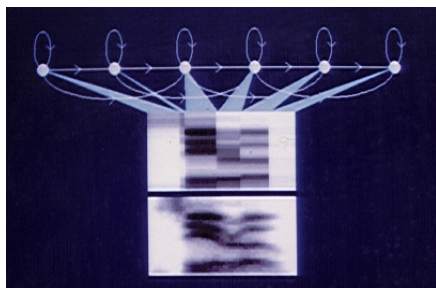


EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

6 state HMM of the digit 'zero'



Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



Word Level HMMs

- Disadvantages:
 - Many examples of each word required for training
 - Fails to exploit regularities in spoken language
- Word-level systems typically restricted to well-defined, demanding, small vocabulary applications

Multimodal Interaction Lab

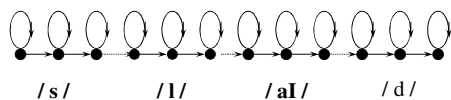
b s : m i g a m

EE4R Automatic Spoken Language Processing



Sub-Word Level HMMs

- Build HMMs for a complete set of sub-word 'building blocks'
- Construct word-level HMMs by concatenation of sub-word HMMs
- E.g. *slide* = / s l aɪ d /



Multimodal Interaction Lab

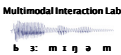
b s : m i g a m

EE4R Automatic Spoken Language Processing



Sub-Word Level HMMs

- Advantages
 - Able to exploit regularities in speech patterns
 - More efficient use of training data - e.g. in phoneme-based system “five” (/ f aI v I/) and “nine” (/n aI n I/) both contribute to /aI/ model.
 - Flexibility - acoustic models can be built **immediately** for words which did not occur in the training data

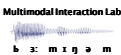


EE4R Automatic Spoken Language Processing



Phoneme-Level HMMs

- Why choose phonemes rather than any other sub-word unit?
- Disadvantages
 - Phonemes are defined in terms of the contrastive properties of speech sounds within a language - not their consistency with HMM assumptions!

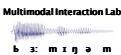


EE4R Automatic Spoken Language Processing



Advantages of Phoneme-HMMs

- Completeness & compactness - approx. 50 phonemes required to describe English.
- Well studied - potential for exploitation of ‘speech knowledge’ (e.g. pronunciation differences due to accent...)
- Availability of extensive phoneme-based pronunciation dictionaries



EE4R Automatic Spoken Language Processing



Context-Sensitivity

- Problem
 - Acoustic realization of a phoneme depends on the context in which it occurs
 - Think of your lip shape for the “k” sound in the words “book shop” and “thick”



EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Biphones and Triphones

- Solution
 - **Context-sensitive** phoneme-level HMMs
 - E.g.
 - ‘biphones’ : (k:_S) in “book shop”
 - ‘triphones’ : (k:u_S) in “book shop”
- Almost all systems use triphone HMMs



EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Triphones - problems

- Increased number of model parameters
 - Need more (well-chosen) training data
- Which triphone?
 - If a word in the application contains a triphone which was not in the training set, which triphone HMM should we use?



EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Number of parameters

- If there are 50 phones, the maximum number of triphone HMMs is $50^3=125,000$
- Most ruled out by **phonological** constraints – most phone triples never occur in speech
- But many are legal

Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Example: Model Parameters

- Each model has 3 emitting states
- Each state modelled as, say, a 10 component Gaussian mixture
- Each feature vector is 40 dimensional
- Hence number of parameters per model is:

$$3 \times (10 \times (40 + 40 + 1) + 9) = 2,457$$

Number of states Number of mixture components Mean vector Variance vector Mixture weight Transition probs

Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Acoustic model parameters

- So, even if we only have 1,000 acoustic models (instead of 125,000), total acoustic model parameters will be 2,457,000
- Too many to estimate with practical quantity of data
- Most common solution is HMM **parameter tying**
- **Different** HMMs share **same** parameters

Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM

Tied variance

- Variances are more costly to estimate than means
- Simple solution – divide set of all HMMs into classes, so that within a class all HMM state PDFs have same variance
- This is **tied variance**
- If **all** HMM state PDFs share the same variance, the variance is referred to as **grand variance**

Multimodal Interaction Lab

b s m i g a m

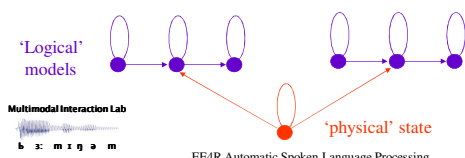
EE4R Automatic Spoken Language Processing



UNIVERSITY OF BIRMINGHAM

Phone decision trees

- Most common approach to general HMM tying is **decision tree clustering**
- Decision tree clustering can be applied to individual states or to whole HMMs – we'll consider states
- Basic idea is to use **knowledge** about which phones are likely to induce similar contextual effects



Multimodal Interaction Lab

b s m i g a m

EE4R Automatic Spoken Language Processing



UNIVERSITY OF BIRMINGHAM

Phonetic knowledge

- For example, we know that /f/ and /s/ are both unvoiced fricatives, produced in a similar manner
- Therefore we might **hypothesise** that, for example, an utterance of the vowel /e/ preceded by /f/ might be similar to one preceded by /s/
- This is the basic idea behind decision tree clustering

Multimodal Interaction Lab

b s m i g a m

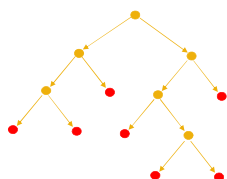
EE4R Automatic Spoken Language Processing



UNIVERSITY OF BIRMINGHAM

A phone decision tree for /e/

- A phone decision tree is just a **binary tree**, where each node of the tree is associated with:
 - A set of phones
 - A position (L or R)
- The root node of the tree corresponds to /e/
- The terminal nodes correspond to **significantly different** contextual variants of /e/



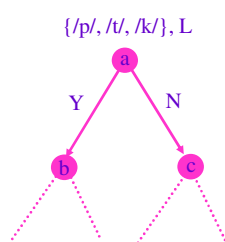
Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



A decision tree node (example)



- Want to choose a state for /e/ in a particular context
- At node (a), ask **question**: is the **Left** context one of the set $\{/p/, /t/, /k/\}$?
- If “yes” go to node (b), otherwise go to node (c)
- Continue until a terminal node is reached
- Choose associated HMM state

Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



Building a phone decision tree for /e/

- First choose a set of **questions**
 - Can be chosen using **phonetic knowledge**
 - ...plus pragmatics!
- Also need the set E of states which occur in a particular position in triphones for /e/
- Each question partitions E into two subsets
 - E_Y – states of /e/ for which answer to question is “Yes”
 - E_N – states of /e/ for which answer to question is “No”

Multimodal Interaction Lab

b s : m i g a m

EE4R Automatic Spoken Language Processing



Building a phone decision tree

- For each question Q , we can define a “quality measure” $g(Q)$
- $g(Q)$ is a measure of how similar to each other the states in E_Y are, and how similar to each other the states in E_N are
- Intuitively, $g(Q)$ is a measure of how **compact** or **‘homogeneous’** the sets E_Y and E_N are
- Choose the question Q for which $g(Q)$ is biggest

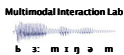


EE4R Automatic Spoken Language Processing



Building a phone decision tree

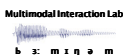
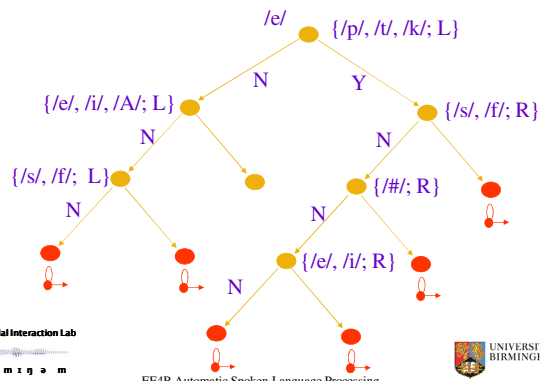
- States in E_Y (resp E_N) are assigned to the “Y” (resp “N”) successor nodes
- Whole process is repeated for each successor node
- Process stops when, for example, the number of states associated with a node reaches a minimum



EE4R Automatic Spoken Language Processing



Phone Decision Tree



EE4R Automatic Spoken Language Processing



Summary

- Word-level HMMs
- Sub-Word HMMs
- Phoneme-level HMMs
- Context-sensitivity
 - Biphones & triphones
- Triphone decision trees

Multimodal Interaction Lab



EE4R Automatic Spoken Language Processing



UNIVERSITY OF
BIRMINGHAM
