

Acoustic Modeling for Speech Recognition

Berlin Chen 2003

References:

1. X. Huang et. al., Spoken Language Processing, Chapter 8
2. The HTK Book (for HTK Version 3.2)

Introduction

- For the given acoustic observation $X = x_1, x_2, \dots, x_n$, the goal of speech recognition is to find out the corresponding word sequence $W = w_1, w_2, \dots, w_m$ that has the maximum posterior probability $P(W | X)$

$$\begin{aligned} \hat{W} &= \arg \max_w P(W | X) \\ &= \arg \max_w \frac{P(W)P(X | W)}{P(X)} \\ &= \arg \max_w P(W)P(X | W) \end{aligned}$$

$$W = w_1, w_2, \dots, w_i, \dots, w_m$$

where $w_i \in V : \{v_1, v_2, \dots, v_N\}$

Language Modeling


Acoustic Modeling

Possible variations

domain, topic, style, etc.

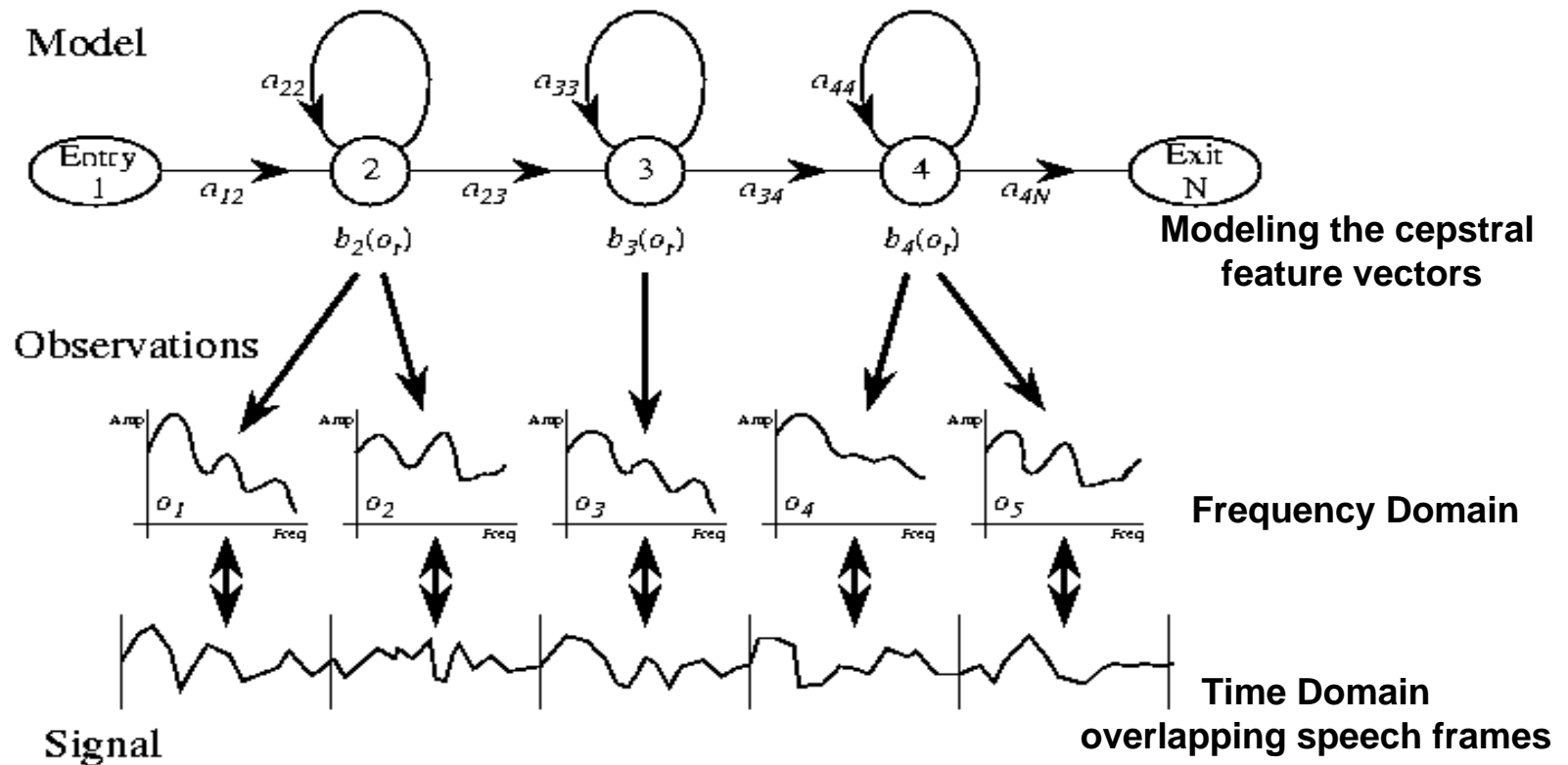
and

speaker, pronunciation, environment, context, etc.

To be discussed later on ! 

Review: HMM Modeling

- Acoustic Modeling using HMMs

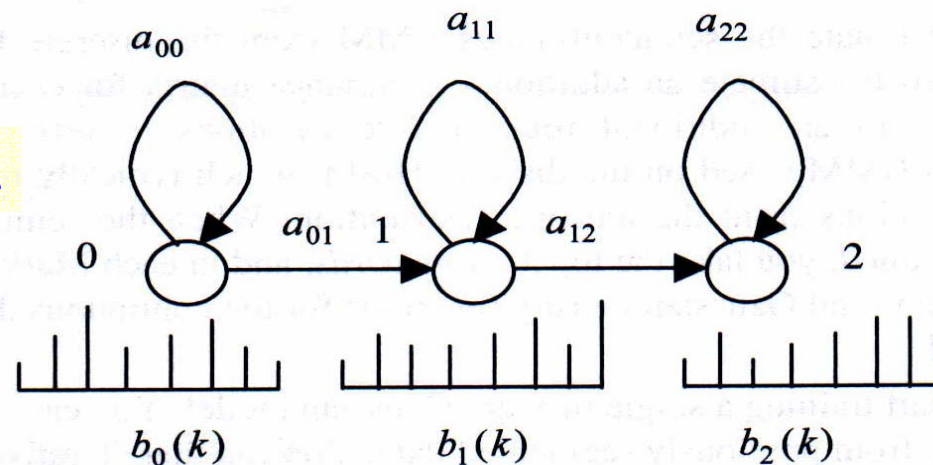


- Three types of HMM state output probabilities are used

Review: HMM Modeling

- Discrete HMM (DHMM): $b_j(\mathbf{v}_k) = P(\mathbf{o}_t = \mathbf{v}_k | s_t = j)$
 - The observations are quantized into a number of symbols
 - The symbols are normally generated by a vector quantizer

A left-to-right HMM



- With multiple codebooks

$$b_j(\mathbf{v}_k) = \sum_{m=1}^M c_{jm} p(\mathbf{o}_t = \mathbf{v}_k | m, s_t = j) \quad \sum_{m=1}^M c_{jm} = 1$$

codebook index

Review: HMM Modeling

- Continuous HMM (CHMM)
 - The state observation distribution of HMM is modeled by multivariate Gaussian mixture density functions (M mixtures)

$$b_j(\mathbf{o}_t) = \sum_{m=1}^M c_{jm} b_{jm}(\mathbf{o}_t)$$

$$= \sum_{m=1}^M c_{jm} N(\mathbf{o}_t; \boldsymbol{\mu}_{jm}, \boldsymbol{\Sigma}_{jm}) = \sum_{m=1}^M c_{jm} \left(\frac{1}{(\sqrt{2\pi})^{L/2} |\boldsymbol{\Sigma}_{jm}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{o}_t - \boldsymbol{\mu}_{jm})^T \boldsymbol{\Sigma}_{jm}^{-1} (\mathbf{o}_t - \boldsymbol{\mu}_{jm})\right) \right), \quad \sum_{m=1}^M c_{jm} = 1$$

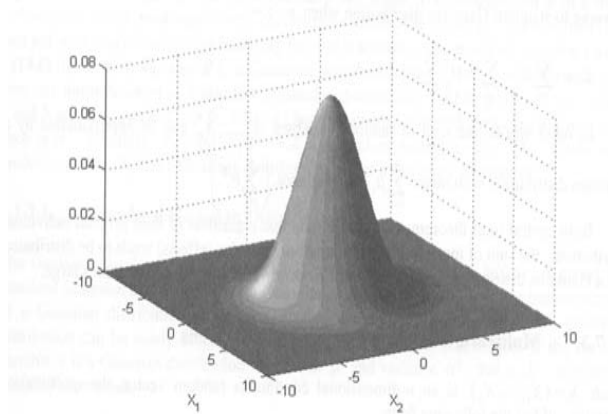


Figure 3.12 A two-dimensional multivariate Gaussian distribution with independent random variables x_1 and x_2 that have the same variance.

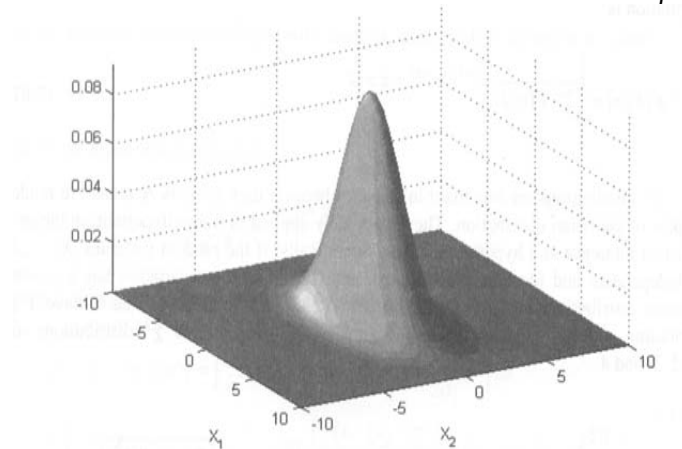
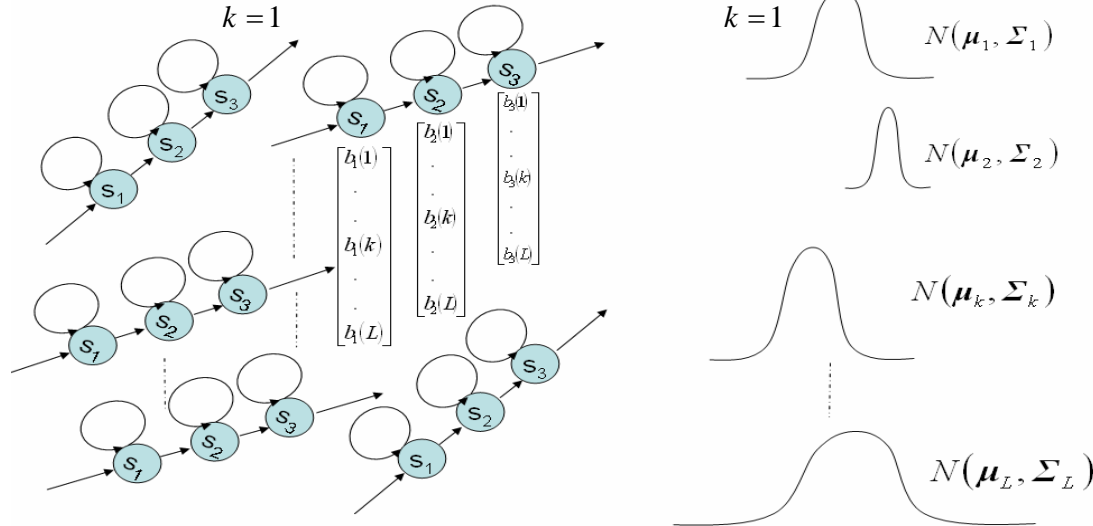


Figure 3.13 Another two-dimensional multivariate Gaussian distribution with independent random variable x_1 and x_2 which have different variances.

Review: HMM Modeling

- Semicontinuous or tied-mixture HMM (SCHMM)
 - The HMM state mixture density functions are tied together across all the models to form a set of shared kernels (shared Gaussians)

$$b_j(\mathbf{o}) = \sum_{k=1}^L b_j(k) f(\mathbf{o} | v_k) = \sum_{k=1}^L b_j(k) N(\mathbf{o}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



- With multiple codebooks

$$b_j(\mathbf{o}) = \sum_{m=1}^M c_m \sum_{k=1}^L b_{jm}(k) f(\mathbf{o} | v_{m,k}) = \sum_{m=1}^M c_m \sum_{k=1}^L b_{jm}(k) N(\mathbf{o}; \boldsymbol{\mu}_{m,k}, \boldsymbol{\Sigma}_{m,k})$$

Review: HMM Modeling

- Comparison of Recognition Performance

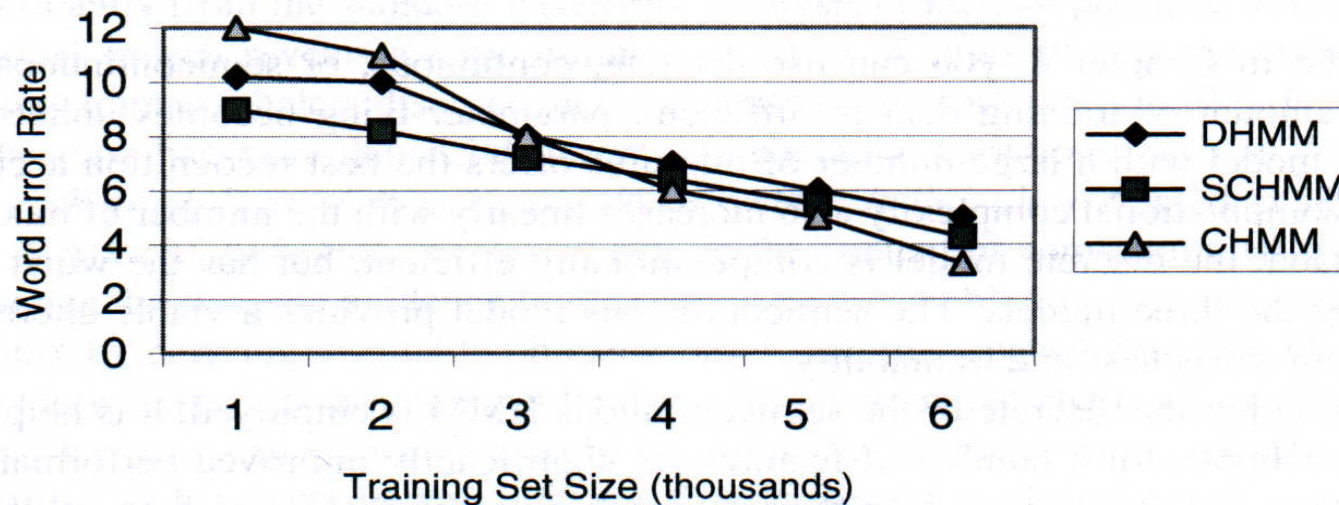


Figure 9.8 Continuous speaker-independent word recognition error rates of the discrete HMM (DHMM), SCHMM, and the continuous HMM (CHMM) with respect to the training set sizes (thousands of training sentences). Both the DHMM and SCHMM have multiple codebooks. The CHMM has 20 mixture diagonal Gaussian density functions.

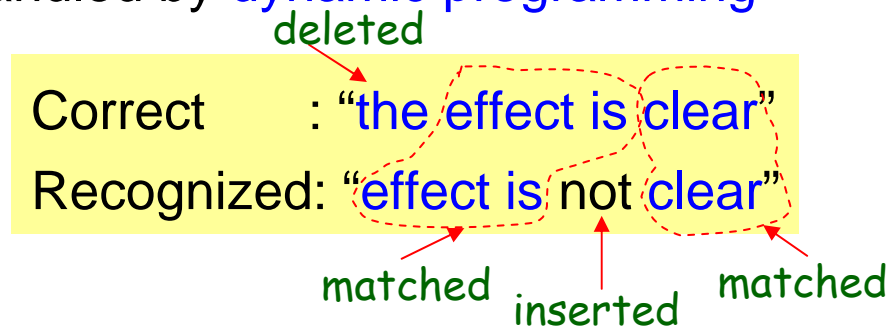
Measures of Speech Recognition Performance

- Evaluating the performance of speech recognition systems is critical, and the **Word Recognition Error Rate (WER)** is one of the most important measures
- There are typically three types of word recognition errors
 - **Substitution**
 - An incorrect word was substituted for the correct word
 - **Deletion**
 - A correct word was omitted in the recognized sentence
 - **Insertion**
 - An extra word was added in the recognized sentence
- How to determine the minimum error rate?

Measures of Speech Recognition Performance

- Calculate the WER by aligning the correct word string against the recognized word string
 - A **maximum substring matching** problem
 - Can be handled by **dynamic programming**

• Example:



- Error analysis: **one deletion and one insertion**
- Measures: **word error rate (WER)**, **word correction rate (WCR)**, **word accuracy rate (WAR)**

Might be higher than 100%

$$\text{Word Error Rate} = 100\% \frac{\text{Sub. + Del. + Ins. words}}{\text{No. of words in the correct sentence}} = \frac{2}{4} = 50\%$$

$$\text{Word Correction Rate} = 100\% \frac{\text{Matched words}}{\text{No. of words in the correct sentence}} = \frac{3}{4} = 75\%$$

$$\text{Word Accuracy Rate} = 100\% \frac{\text{Matched - Ins. words}}{\text{No. of words in the correct sentence}} = \frac{3-1}{4} = 50\%$$

WER+
WAR
=100%

Might be negative

Measures of Speech Recognition Performance

- A Dynamic Programming Algorithm (Textbook)

ALGORITHM 9.1: ALGORITHM TO MEASURE THE WORD ERROR RATE

Step 1: Initialization $R[0,0] = 0$ $R[i, j] = \infty$ if $(i < 0)$ or $(j < 0)$ $B[0,0] = 0$

Step 2: Iteration

for $i = 1, \dots, n$ { //denotes for the word length of the correct/reference sentence

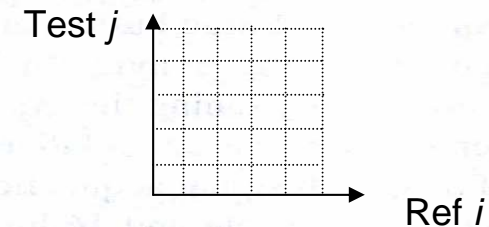
for $j = 1, \dots, m$ //denotes for the word length of the recognized/test sentence

minimum word
error alignment
at the a grid $[i,j]$

$$R[i, j] = \min \begin{bmatrix} R[i-1, j] + 1 \text{ (deletion)} \\ R[i-1, j-1] \text{ (match) /hit} \\ R[i-1, j-1] + 1 \text{ (substitution)} \\ R[i, j-1] + 1 \text{ (insertion)} \end{bmatrix}$$

kinds of
alignment

$$B[i, j] = \begin{cases} 1 & \text{if deletion} \\ 2 & \text{if insertion} \\ 3 & \text{if match /hit} \\ 4 & \text{if substitution} \end{cases} \} \}$$



Step 3: Backtracking and termination

$$\text{word error rate} = 100\% \times \frac{R(n, m)}{n}$$

$$\text{optimal backward path} = (s_1, s_2, \dots, 0)$$

$$\text{where } s_1 = B[n, m], s_t = \begin{bmatrix} B[i-1, j] & \text{if } s_{t-1} = 1 \\ B[i, j-1] & \text{if } s_{t-1} = 2 \\ B[i-1, j-1] & \text{if } s_{t-1} = 3 \text{ or } 4 \end{bmatrix} \text{ for } t = 2, \dots \text{ until } s_t = 0$$

Measures of Speech Recognition Performance

- Algorithm (by Berlin Chen)

```

Step 1: Initialization :
G[0][0] = 0;
for i = 1,..., n { //test
    G[i][0] = G[i - 1][0] + 1;
    B[i][0] = 1; //Insertion
} (Horizontal Direction)
for j = 1,..., m { //reference
    G[0][j] = G[0][j - 1] + 1;
    B[0][j] = 2; // Deletion
} (Vertical Direction)
    
```



```

Step 2: Iteration:
for i = 1,...,n { //test
    for j = 1,...,m { //reference
        G[i][j] = min [
            G[i - 1][j] + 1 (Insertion)
            G[i][j - 1] + 1 (Delection)
            G[i - 1][j - 1] + 1 (if LR[i] != LT[i], Substitution)
            G[i - 1][j - 1] (if LR[i] = LT[i], Match)
        ]
        B[i][j] {
            1; //Insertion, (Horizontal Direction)
            2; //Deletion, (Vertical Direction)
            3; //Substitution (Diagonal Direction)
            4; //match (Diagonal Direction)
        }
    } //for j, reference
} //for i, test
    
```



Step 3: Measure and Backtrace :

$$\text{Word Error Rate} = 100\% \times \frac{G[n][m]}{m}$$

$$\text{Word Accuracy Rate} = 100\% - \text{Word Error Rate}$$

Optimal backtrace path = (B[n][m] → → B[0][0])

if B[i][j] = 1 print " LT[i]" ; //Insertion, then go left

else if B[i][j] = 2 print "LR[j] " ; //Deletion, then go down

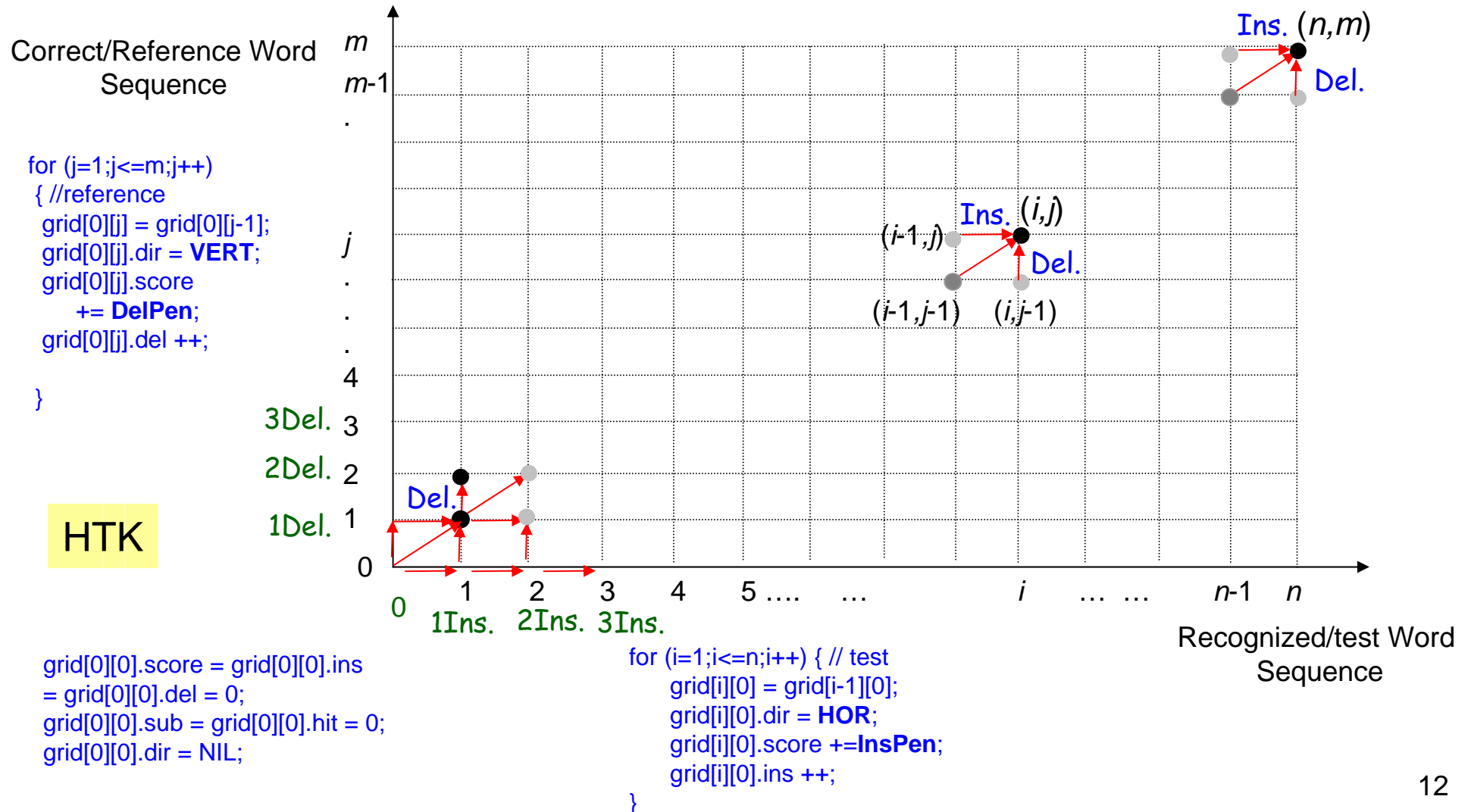
else print "LR[j] LR[i]" ; //Hit/Match or Substitution, then go down diagonally

Note: the penalties for substitution, deletion and insertion errors are all set to be 1 here

Measures of Speech Recognition Performance

- A Dynamic Programming Algorithm

– Initialization



Measures of Speech Recognition Performance

- Program

```

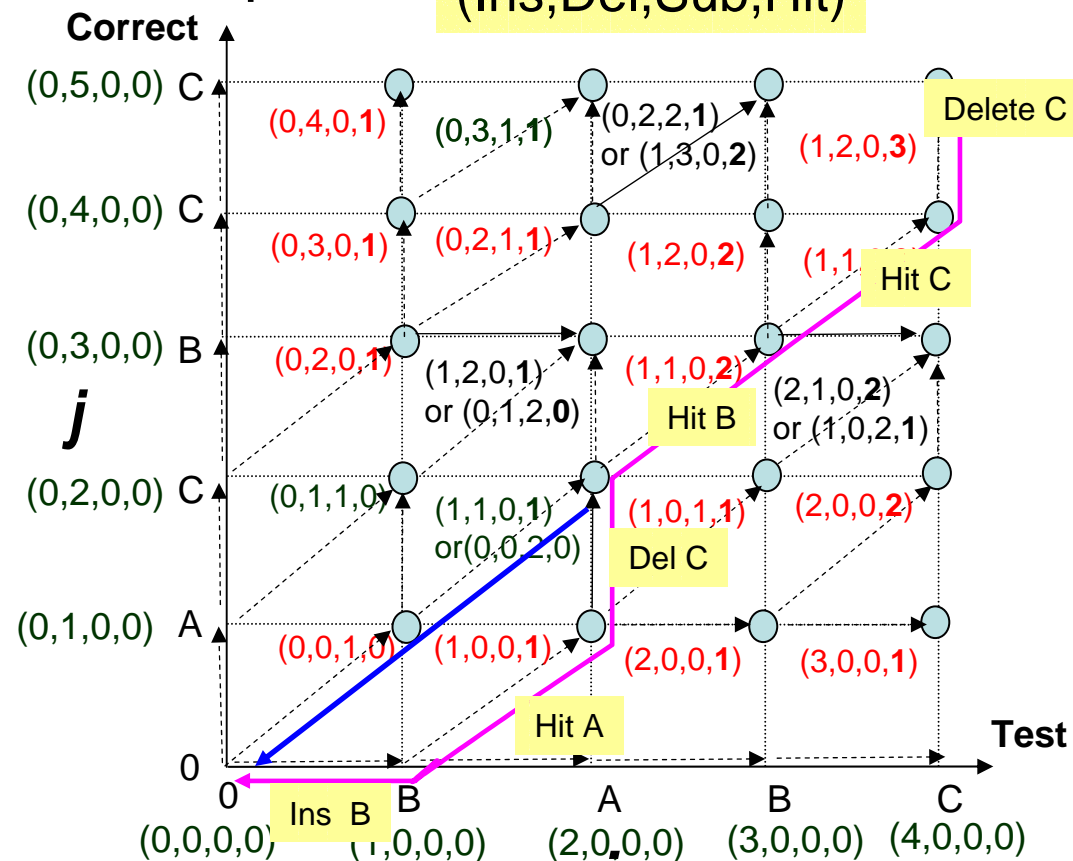
for (i=1;i<=n;i++) //test
{ gridi = grid[i]; gridi1 = grid[i-1];
  for (j=1;j<=m;j++) //reference
  {   h = gridi1[j].score +insPen;
      d = gridi1[j-1].score;
      if (lRef[j] != lTest[i])
        d += subPen;
      v = gridi[j-1].score + delPen;
      if (d<=h && d<=v) { /* DIAG = hit or sub */
        gridi[j] = gridi1[j-1];
        gridi[j].score = d;
        gridi[j].dir = DIAG;
        if (lRef[j] == lTest[i]) ++gridi[j].hit;
        else ++gridi[j].sub;
      }
      else if (h<v) { /* HOR = ins */
        gridi[j] = gridi1[j];
        gridi[j].score = h;
        gridi[j].dir = HOR;
        ++gridi[j].ins;
      }
      else { /* VERT = del */
        gridi[j] = gridi[j-1];
        gridi[j].score = v;
        gridi[j].dir = VERT;
        ++gridi[j].del;
      }
    }
  }
} /* for i */

```

HTK

- Example 1

(Ins,Del,Sub,Hit)



Alignment 1: WER= 60%

Correct: (A) (C) (B) (C) (C)

Test: (B) (A) (B) (C) (C)

Ins B Hit A Del C Hit b Hit c Del c

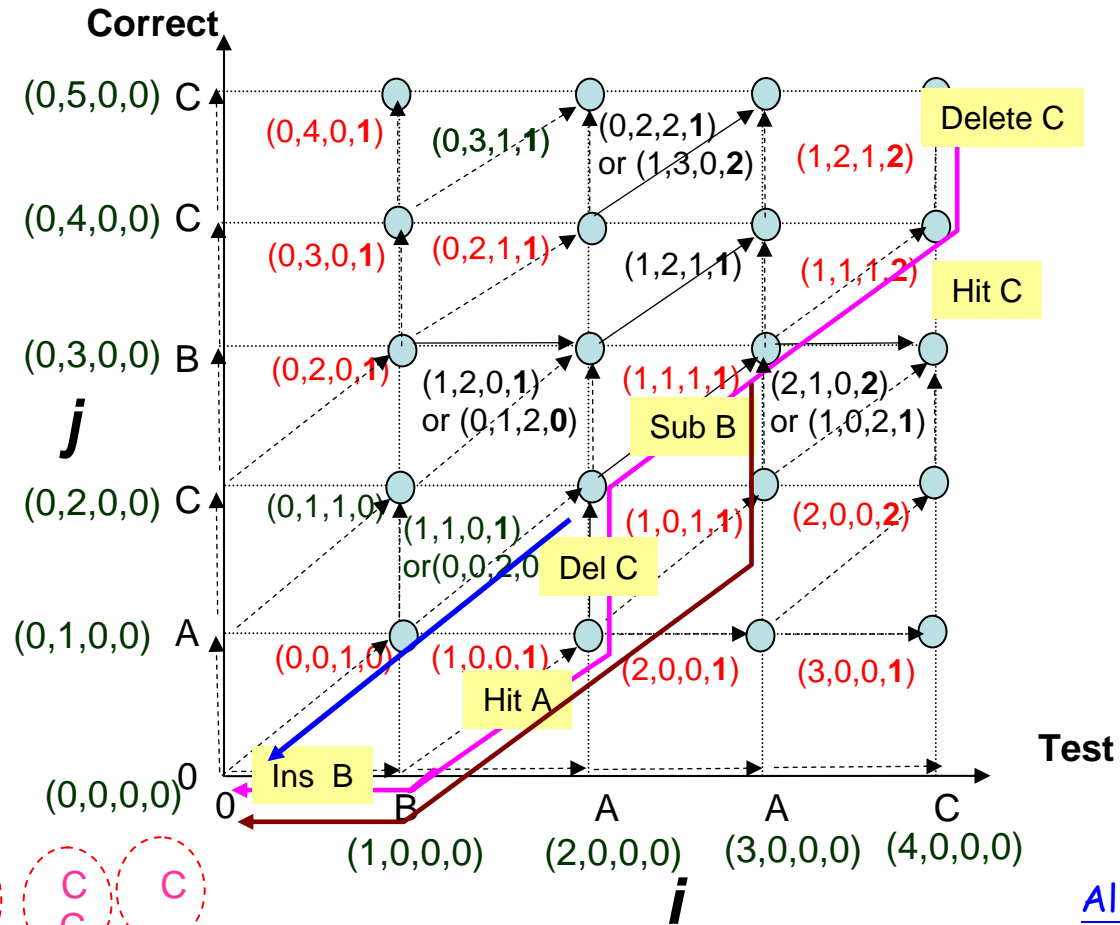
Still have an Other optimal alignment !

Measures of Speech Recognition Performance

- Example 2

Note: the penalties for substitution, deletion and insertion errors are all set to be 1 here

(Ins, Del, Sub, Hit)



Alignment 1: WER= 80%

Correct: A C B C C
 Test: B A A C C

Ins B Hit A Del C Sub B Hit c Del c

Correct: A C B C C
 Test: B A A C C

Hit A Del C Sub B Hit c Del c

Alignment 2:
WER=80%

Alignment 3:
WER=80%

Correct: A C B C C
 Test: B A A C C

Ins B Hit A Sub C Del B Hit c Del c

Measures of Speech Recognition Performance

- Two common settings of different penalties for substitution, deletion, and insertion errors

```
/* HTK error penalties */
```

```
subPen = 10;
```

```
delPen = 7;
```

```
insPen = 7;
```

```
/* NIST error penalties*/
```

```
subPenNIST = 4;
```

```
delPenNIST = 3;
```

```
insPenNIST = 3;
```

Choice of Appropriate Units for HMMs

- Issues for HMM Modeling units
 - **Accurate**: accurately represent the acoustic realization that appears in different contexts
 - **Trainable**: have enough data to estimate the parameters of the unit (or HMM model)
 - **Generalizable**: any new word can be derived from a predefined unit inventory for task-independent speech recognition

Choice of Appropriate Units for HMMs

- Comparison of different units

- **Word:**

- Semantic meaning, capturing within-word coarticulation, can be accurately trained for small-vocabulary speech recognition, but not generalizable

- **Phone:**

- More trainable and generalizable, but less accurate
 - There are only about 50 context-independent phones in English and 30 in Mandarin Chinese
 - Drawbacks: the realization of a phoneme is strongly affected by neighboring phonemes (e.g., /t s/ and /t r/)

- **Syllable:**

- A compromise between the word and phonetic models. Syllables are larger than phone
 - There only about 1,300 tone-dependent syllables in Chinese and 50 in Japanese. However, there are over 30,000 in English

subword

Choice of Appropriate Units for HMMs

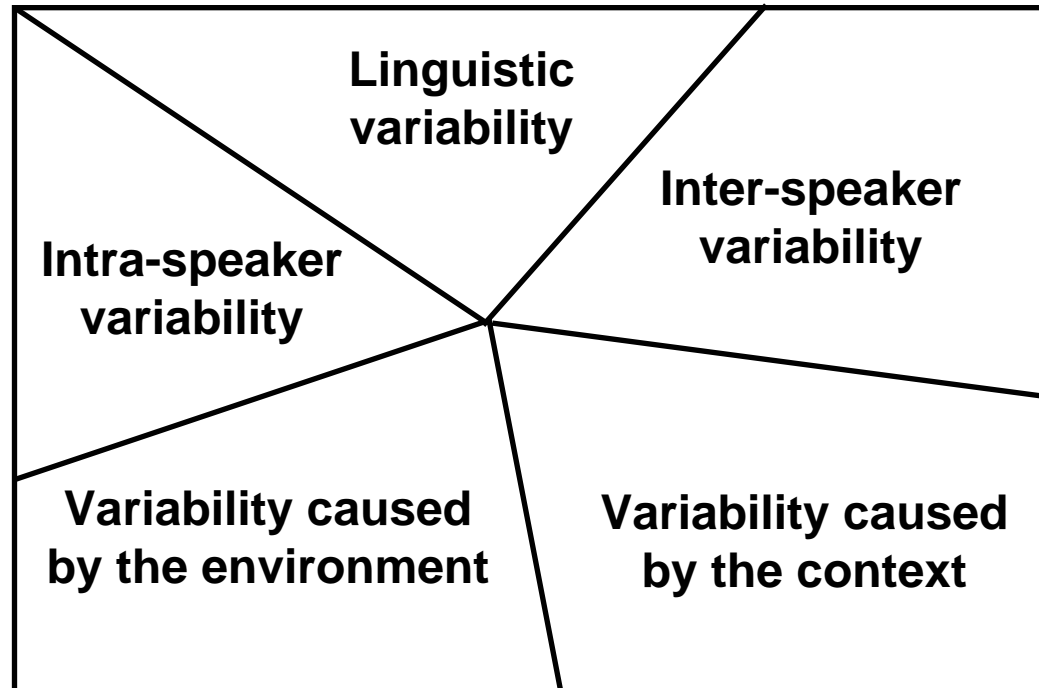
- Phonetic Structure of Mandarin Syllables

Syllables (1,345)		
Base-syllables (408)		Tones (4+1)
INITIAL's (21)	FINAL's (37)	
Phone-like Units/Phones (33)		

Variability in the Speech Signals

Pronunciation
Variation

Speaker-independency
Speaker-adaptation
Speaker-dependency



Robustness
Enhancement

Context-Dependent
Acoustic Modeling

Variability in the Speech Signals

- Context Variability

- Context variability at word/sentence level

- E.g., “Mr. Wright should write to Ms. Wright right away about his Ford or four door Honda”
 - Same pronunciation but different meaning (*Wright* , *write* , *right*)
 - Phonetically identical and semantically relevant (*Ford or*, *four door*)

Pause or intonation information is needed

- Context variability at phonetic level

- The acoustic realization of phoneme /ee/ for word *peat* and *wheel* depends on its left and right context

the effect is more important in fast speech or spontaneous conversations, since many phonemes are not fully realized!

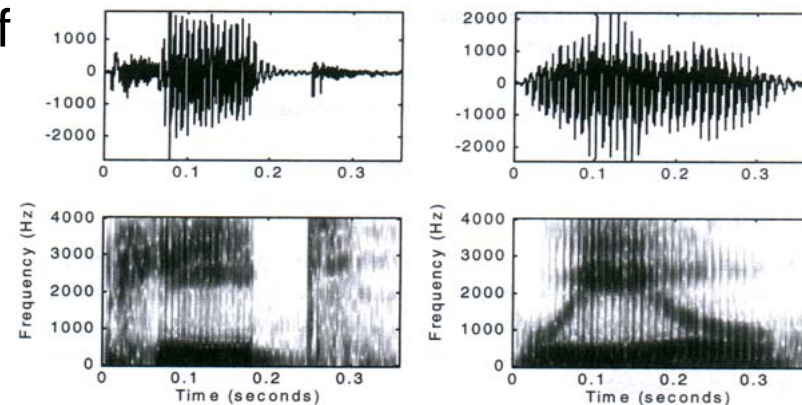


Figure 9.1 Waveforms and spectrograms for words *peat* (left) and *wheel* (right). The phoneme /ee/ is illustrated with two different left and right contexts. This illustrates that different contexts may have different effects on a phone.

Variability in the Speech Signals

- Style Variability (also including intra-speaker and linguistic variability)
 - Isolated speech recognition
 - Users have to pause between each word (a clear boundary between words)
 - Errors such as “*Ford or*” and “*four door*” can be eliminated
 - But unnatural to most people
 - Continuous speech recognition
 - Causal, spontaneous, and conversational
 - Higher speaking rate and co-articulation effects
 - Emotional changes also introduce more significantly variations

TABLE I
SPEAKING RATES OF THE BROADCAST NEWS SPEECH USED IN THIS STUDY

Data Types	Training					Testing
	PRS	UFO	VOH	VOT	Average	BCC
Average Speaking Rates (characters/sec)	4.9	6.2	4.8	5.7	5.4	5.8

Statistics of the speaking rates of the broadcast new speech collected in Taiwan

Variability in the Speech Signals

- Speaker Variability

- Interspeaker

- Vocal tract size, length and width of the neck and a range of physical characteristics
- E.g., gender, age, dialect, health, education, and personal style

- Intraspeaker

- The same speaker is often unable to precisely produce the same utterance
- The shape of the vocal tract movement and rate of delivery may vary from utterance to utterance

- Issues for acoustic modeling

- Speaker-dependent (SD), speaker-independent (SI) and speaker-adaptive (SA) modeling
- Typically an SD system can reduce WER by more than 30% as compared with a comparable SI one

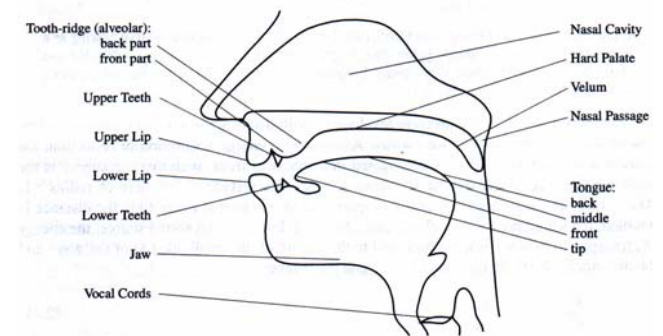


Figure 2.4 A schematic diagram of the human speech production apparatus.

Variability in the Speech Signals

- Environment Variability
 - The world we live in is full of sounds of varying loudness from different sources
 - Speech recognition in hands-free or mobile environments remain one of the most severe challenges
 - The spectrum of noises varies significantly
 - Noise may also be present from the input device itself, such as microphone and A/D interface noises
 - We can reduce the error rates by using multi-style training or adaptive techniques
 - Environment variability remains as one of the most severe challenges facing today's state-of-the-art speech systems

Context Dependency

- Review: Phone and Phoneme

- In speech science, the term **phoneme** is used to denote any of the minimal units of speech sound in a language that can serve to distinguish one word from another
- The term **phone** is used to denote a phoneme's acoustic realization
- E.g., English phoneme /t/ has two very different acoustic realizations in the word *sat* and *meter*
 - We have better treat them as **two different phones** when building a spoken language system

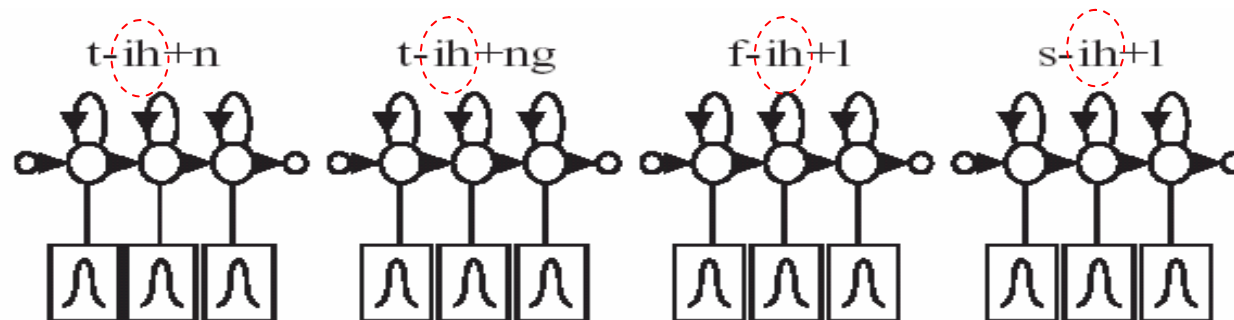
Context Dependency

- Why Context Dependency
 - If we make unit context dependent, we can significantly improve the recognition accuracy, provided there are enough training data for parameter estimation
 - A context usually refers to the immediate left and/or right neighboring phones
 - Context-dependent (CD) phonemes have been widely used for LVCSR systems

Context Dependency

allophones: different realizations of a phoneme is called allophones
→ **Triphones** are examples of allophones

- Triphone (Intra-word triphone)
 - A triphone model is a phonetic model that takes into consideration both the left and right neighboring phones
 - It captures the most important coarticulatory effects
 - Two phones having the same identity but different left and right context are considered different triphones
 - **Challenging issue:** Need to balance trainability and accuracy with a number of parameter-sharing techniques



Context Dependency

- Modeling inter-word context-dependent phone (like triphones) is complicated
 - Although the juncture effect on word boundaries is one of the most serious coarticulation phenomena in continuous speech recognition
 - E.g., speech /*s p iy ch*/ → /*s*/ and /*ch*/ are depending on the preceding and following words in actual sentences
 - Should be taken into consideration with the decoding/search scheme adopted
- Even with the same left/right context, a phone may have significant different realizations at different word positions
 - E.g., that rock /*t*/ → extinct! , theatrical /*t*/ → /*ch*/

Context Dependency

- Stress information for context dependency
 - **Word-level stress** (free stress)
 - The stress information: longer duration, higher pitch and more intensity for stressed vowels
 - E.g., *import* (n) vs. *import* (v), *content* (n) vs. *content* (v)

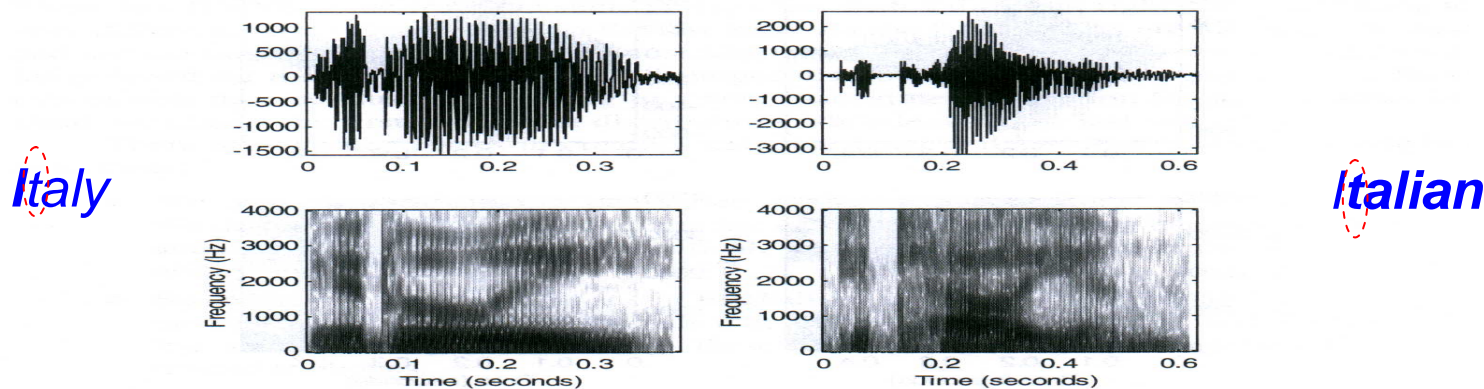


Figure 9.3 The importance of stress is illustrated in *Italy* vs. *Italian* for phone /t/. The realizations are quite different, even though they share the same left and right context.

- **Sentence-level stress** (including **contractive** and **emphatic** stress)
 - Sentence-level stress is very hard to model without incorporate semantic and pragmatic knowledge
 - Contractive: e.g., “I said **import** records not **export**”
 - Emphatic: e.g., “I **did** have dinner”

Clustered Acoustic-Phonetic Units

- Triphone modeling assumes that every triphone context is different. Actually, many phones have similar effects on the neighboring phones
 - /b/ and /p/ (labial stops) (or, /r/ and /w/ (liquids)) have similar effects on the following vowel
- It is desirable to find instances of similar contexts and merge them
 - A much more manageable number of models that can be better trained

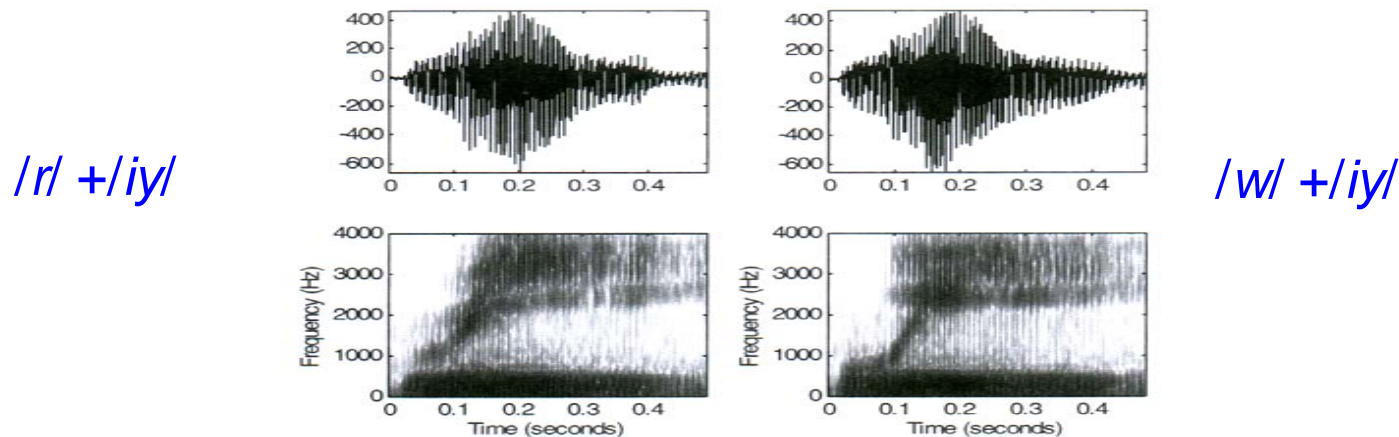
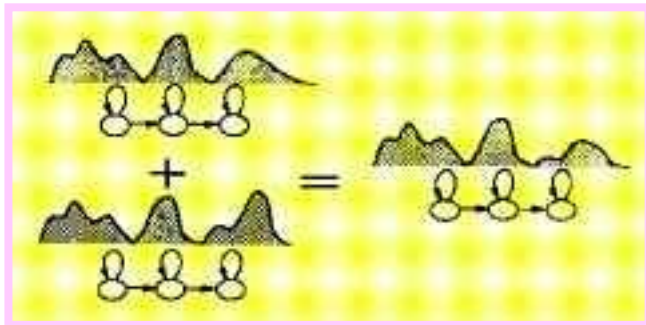


Figure 9.4 The spectrograms for the phoneme /iy/ with two different *left-contexts* are illustrated. Note that /r/ and /w/ have similar effects on /iy/. This illustrates that different left-contexts may have similar effects on a phone.

Clustered Acoustic-Phonetic Units

- Model-based clustering



- State-based clustering (state-tying)
 - Keep the dissimilar states of two models apart while the other corresponding states are merged

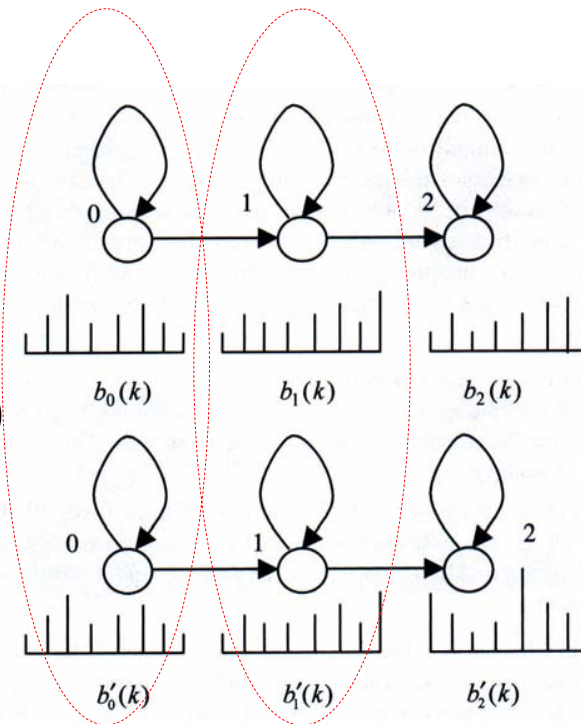
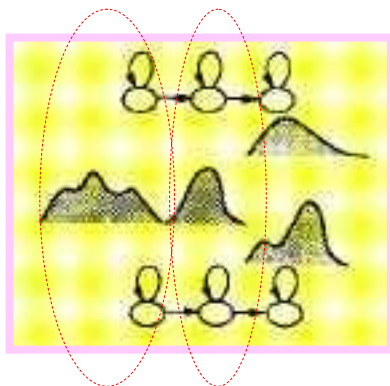
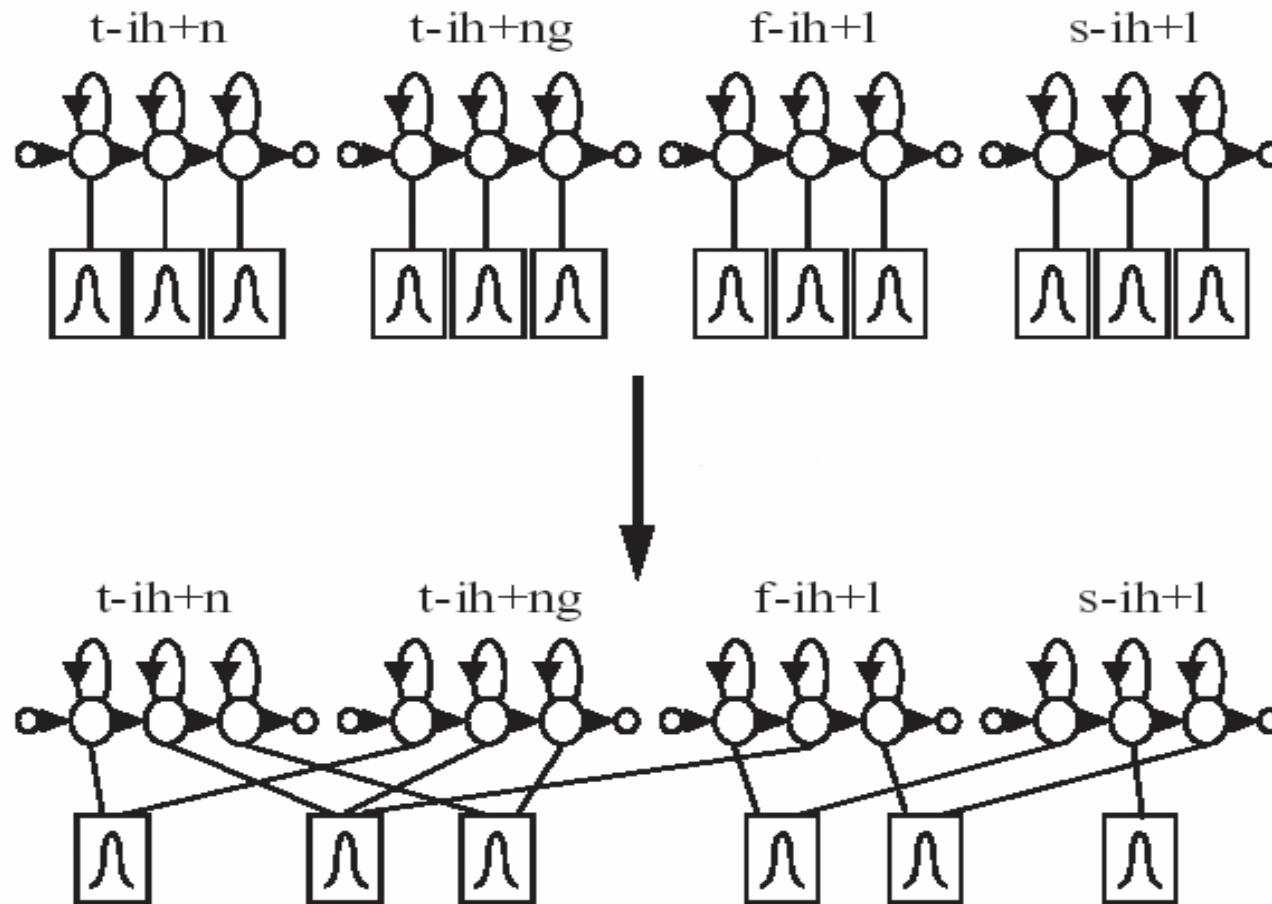


Figure 9.5 State-based vs. model-based clustering. These two models are very similar, as both the first and the second output distributions are almost identical. The key difference is the output distribution of the third state. If we measure the overall model similarity, which is often based on the accumulative output distribution similarities of all states, these two models may be clustered, leading to a very inaccurate distribution for the last state. If we cluster output distributions at state level, we can cluster the first two output distributions while leaving the last ones intact, leading to more accurate representations.

Clustered Acoustic-Phonetic Units

- State-tying of triphones

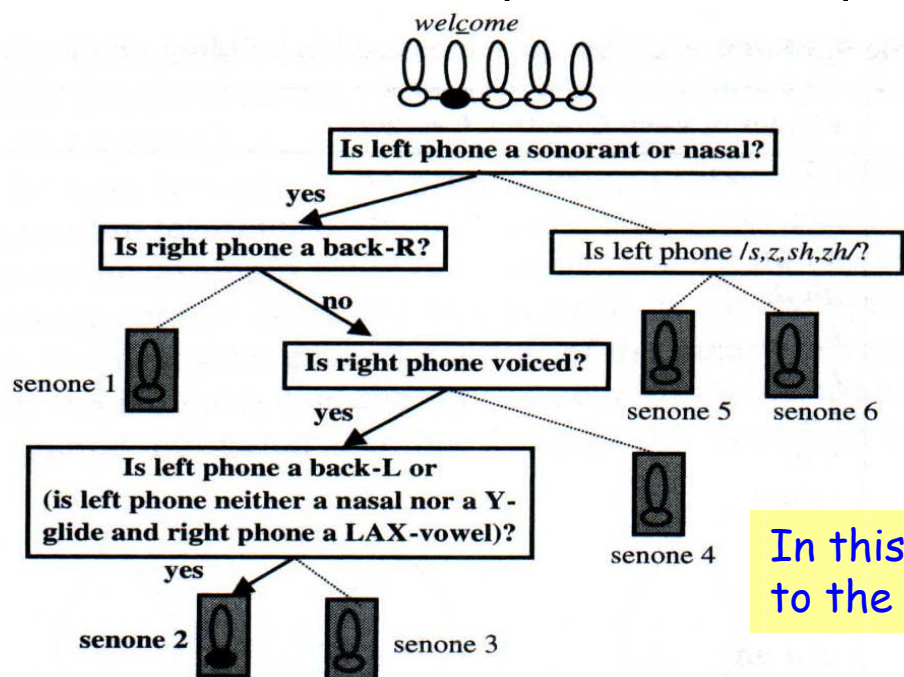


Clustered Acoustic-Phonetic Units

- Two key issues for CD phonetic or subphonetic modeling
 - Tying the phones with similar contexts to improve trainability and efficiency
 - Enable better parameter sharing and smoothing
 - Mapping the unseen triphones (in the test) into appropriately trained triphones is important
 - Because the possible of triphones could be very large
 - E.g., English has over 100,000 triphones

Clustered Acoustic-Phonetic Units

- Microsoft's approach - State-based clustering
 - Generate clustering to the state-dependent output distributions across different phonetic models
 - Each cluster represents a set of similar HMM states and is called **senone**
 - A subword model is composed of a sequence of senons



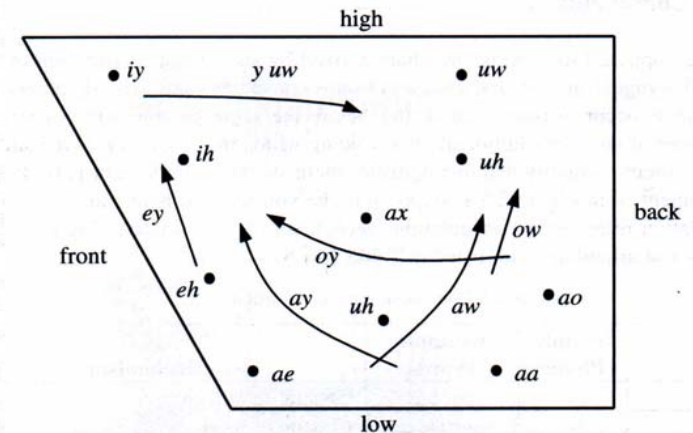
In this example, the tree can be applied to the second state of any /k/ triphone

Figure 9.6 A decision tree for classifying the second state of K -triphone HMMs [48].

Clustered Acoustic-Phonetic Units

- Some example questions used in building senone trees

Questions	Phones in Each Question Category
<i>Aspseg</i>	<i>hh</i>
<i>Sil</i>	<i>sil</i>
<i>Alvstp</i>	<i>d t</i>
<i>Dental</i>	<i>dh th</i>
<i>Labstp</i>	<i>b p</i>
<i>Liquid</i>	<i>l r</i>
<i>Lw</i>	<i>l w</i>
<i>S/Sh</i>	<i>s sh</i>
<i>Sylbic</i>	<i>er axr</i>
<i>Velstp</i>	<i>g k</i>
<i>Affric</i>	<i>ch jh</i>
<i>Lqgl-B</i>	<i>l r w</i>
<i>Nasal</i>	<i>m n ng</i>
<i>Retro</i>	<i>r er axr</i>
<i>Schwa</i>	<i>ax ix axr</i>
<i>Velar</i>	<i>ng g k</i>
<i>Fric2</i>	<i>th s sh f</i>
<i>Fric3</i>	<i>dh z zh v</i>
<i>Lqgl</i>	<i>l r w y</i>
<i>S/Z/Sh/Zh</i>	<i>s z sh zh</i>
<i>Wglide</i>	<i>uw aw ow w</i>
<i>Labial</i>	<i>w m b p v</i>
<i>Palatl</i>	<i>y ch jh sh zh</i>
<i>Yglide</i>	<i>iy ay ey oy y</i>
<i>High</i>	<i>ih ix iy uh uw y</i>
<i>Lax</i>	<i>eh ih ix uh ah ax</i>
<i>Low</i>	<i>ae aa ao aw ay oy</i>
<i>Orstp2</i>	<i>p t k</i>
<i>Orstp3</i>	<i>b d g</i>
<i>Alvelr</i>	<i>n d t s z</i>
<i>Diph</i>	<i>uw aw ay ey iy ow oy</i>
<i>Fric1</i>	<i>dh th s sh z zh v f</i>
<i>Round</i>	<i>uh ao uw ow oy w axr er</i>
<i>Frnt-R</i>	<i>ae eh ih ix iy ey ah ax y aw</i>
<i>Tense</i>	<i>iy ey ae uw ow aa ao ay oy aw</i>
<i>Back-L</i>	<i>uh ao uw ow aa er axr l r w aw</i>
<i>Frnt-L</i>	<i>ae eh ih ix iy ey ah ax y oy ay</i>
<i>Back-R</i>	<i>uh ao uw ow aa er axr oy l r w ay</i>
<i>Orstp1</i>	<i>b d g p t k ch jh</i>
<i>Vowel</i>	<i>ae eh ih ix iy uh ah ax aa ao uw aw ay ey ow oy er axr</i>
<i>Son</i>	<i>ae eh ih ix iy ey ah ax oy ay uh ao uw ow aa er axr aw l r w y</i>
<i>Voiced</i>	<i>ae eh ih ix iy uh ah ax aa ao uw aw ay ey ow oy l r w y er axr m n ng jh b d dh g v z zh</i>



Clustered Acoustic-Phonetic Units

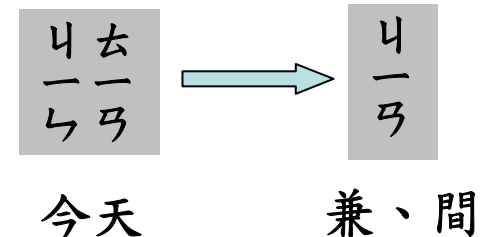
- Comparison of recognition performance for different acoustic modeling

Table 9.4 Relative error reductions for different modeling units. The reduction is relative to that of the preceding row.

Units	Relative Error Reductions
Context-independent phone	Baseline
Context-dependent phone	+25%
Clustered triphone	+15%
Senone	+24%

Pronunciation Variation

- We need to provide alternative pronunciations for words that may have very different pronunciations
 - In continuous speech recognition, we must handle the modification of **interword pronunciations** and **reduced sounds**
- Variation kinds
 - **Co-articulation** (Assimilation)
 - “did you” /*d ih jh y ah*/, “set you” /*s eh ch er*/
 - Assimilation: a change in a segment to make it more like a neighboring segment
 - **Deletion**
 - /t/ and /d/ are often deleted before a consonant
- Variation can be drawn between
 - Inter-speaker variation (social)
 - Intra-speaker variation (stylistic)



Pronunciation Variation

- Pronunciation Network (a probabilistic finite state machine)

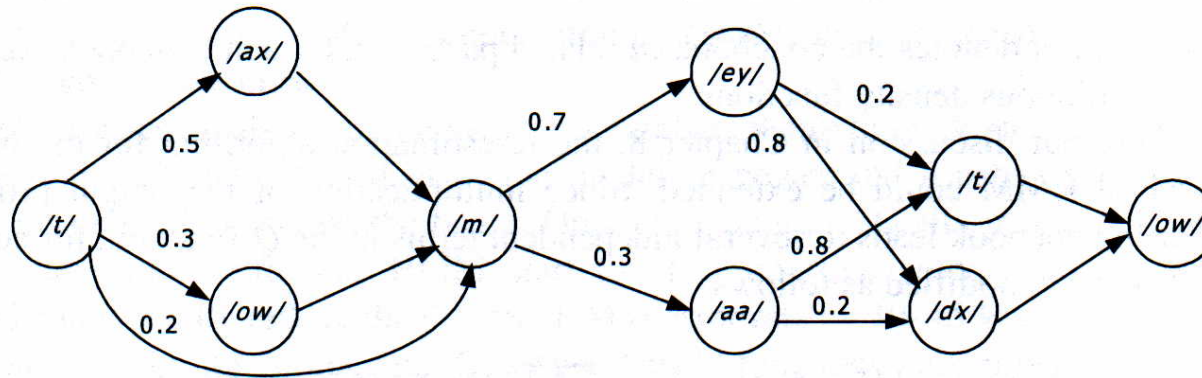


Figure 9.7 A possible pronunciation network for word *tomato*. The vowel */ey/* is more likely to flap, thereby having a higher transition probability into */dx/*.

- Examples:
 - E. g., word “*that*” appears 328 times in one corpus, with 117 different tokens of the 328 times (only 11% of the tokens are most frequent)
Greenberg, 1998
 - Cheating experiments show big performance improvements achieved if the tuned pronunciations were applied to those in test data (e.g. Switchboard WER goes from 40% to 8%)

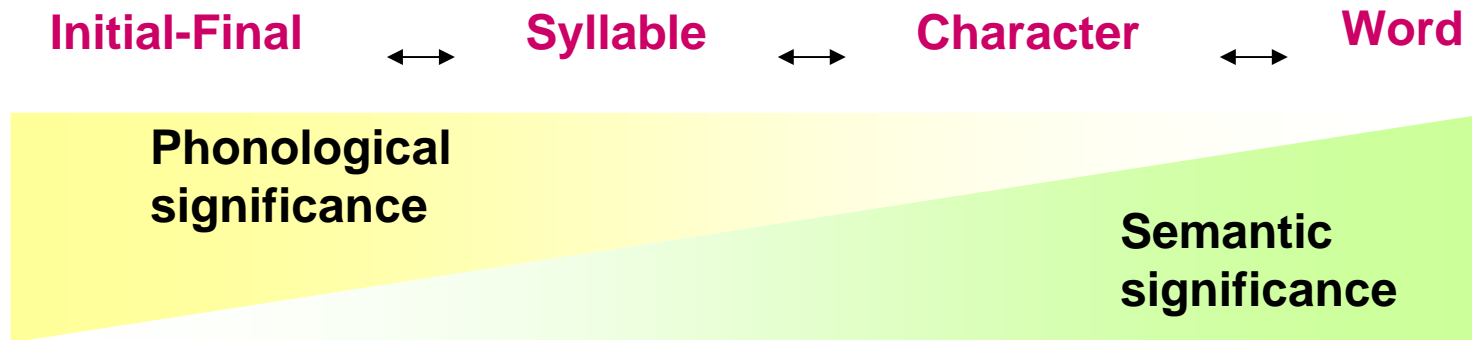
McAllaster et al., 1998

Pronunciation Variation

- Adaptation of Pronunciations
 - Dialect-specific pronunciations
 - Native vs. non-native pronunciations
 - Rate-specific pronunciations
- Side Effect
 - Adding more and more variants to the pronunciation lexicon increases size and confusion of the vocabulary
 - Lead to increased ASR WER

Characteristics of Mandarin Chinese

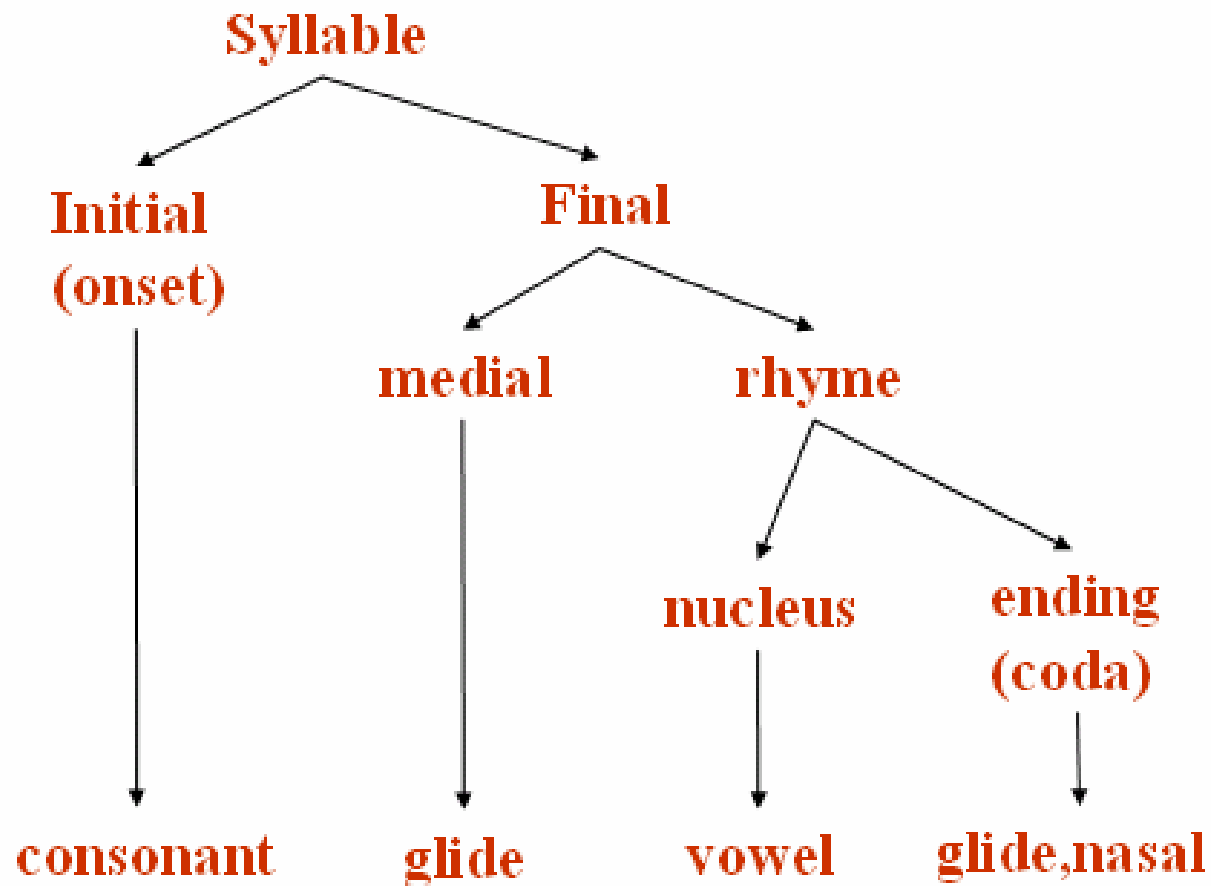
- **Four levels of linguistic units**



- **A monosyllabic-structure language**
 - All characters are monosyllabic
- **Most characters are morphemes (詞素)**
- **A word is composed of one to several characters**
- **Homophones**
 - Different characters sharing the same syllable

Characteristics of Mandarin Chinese

- Chinese syllable structure



Characteristics of Mandarin Chinese

- Sub-syllable HMM Modeling
 - INITIALS

Context -Independent INITIALS
b(ㄅ), p(ㄆ), m(ㄇ), f(ㄈ), d(ㄉ), t(ㄊ), n(ㄋ), l(ㄌ), g(ㄍ), k(ㄎ), h(ㄏ), j(ㄐ), ch(ㄑ), sh(ㄒ), j(ㄓ), ch(ㄔ), sh(ㄕ), r(ㄖ), tz(ㄗ), ts(ㄘ), s(ㄙ)

Table 2.2.1a 21 context-independent INITIALS.

Sub-Syllable HMM Modeling

- Sub-syllable HMM Modeling

- FINALs

Cluster	FINALs
1.(empty)	empty
2.(a)	a(ㄚ), ai(ㄞ), au(ㄨ), an(ㄢ), ang(ㄤ)
3.(o)	o(ㄛ), ou(ㄨ)
4.(e)	e(ㄝ), en(ㄣ), eng(ㄤ), er(ㄦ)
5.(i)	i(ㄟ), ia(ㄟㄚ), ie(ㄟㄝ), iai(ㄟㄞ), iau(ㄟㄨ), ian(ㄟㄢ), in(ㄟㄣ), ing(ㄟㄤ), iang(ㄟㄤ), iou(ㄟㄨ) , io (ㄟㄛ, e.g., for 唷 was ignored here)
6.(u)	u(ㄨ), ua(ㄨㄚ), uo(ㄨㄛ), uai(ㄨㄞ), uei(ㄨㄝ), uan(ㄨㄢ), uen(ㄨㄣ), ueng(ㄨㄤ), uang(ㄨㄤ)
7.(iu)	iu(ㄟㄨ), iue(ㄟㄨㄝ), iuan(ㄟㄨㄢ), iun(ㄟㄨㄣ), iung(ㄟㄨㄤ)
8.(E)	ei(ㄝ)

Table 2.2.1b 37 FINALs, which can be divided into 8 clusters according to the beginning phone.

Classification and Regression Trees (CART)

- CART are **binary decision trees**, with splitting questions attached to each node
 - Act like a rule-based system where the classification carried out by a sequence of decision rules
- CART provides an easy representation that interprets and predicates the structure of a set of data
 - Handle data with high dimensionality, mixed data type and nonstandard data structure
- CART also provides an automatic and data-driven framework to construct the decision process based on objective criteria, not subjective criteria
 - E.g., the choice and order of rules
- CART is a kind of clustering/classification algorithms

Classification and Regression Trees

- Example: height classification
 - Assign a person to one of the following five height classes

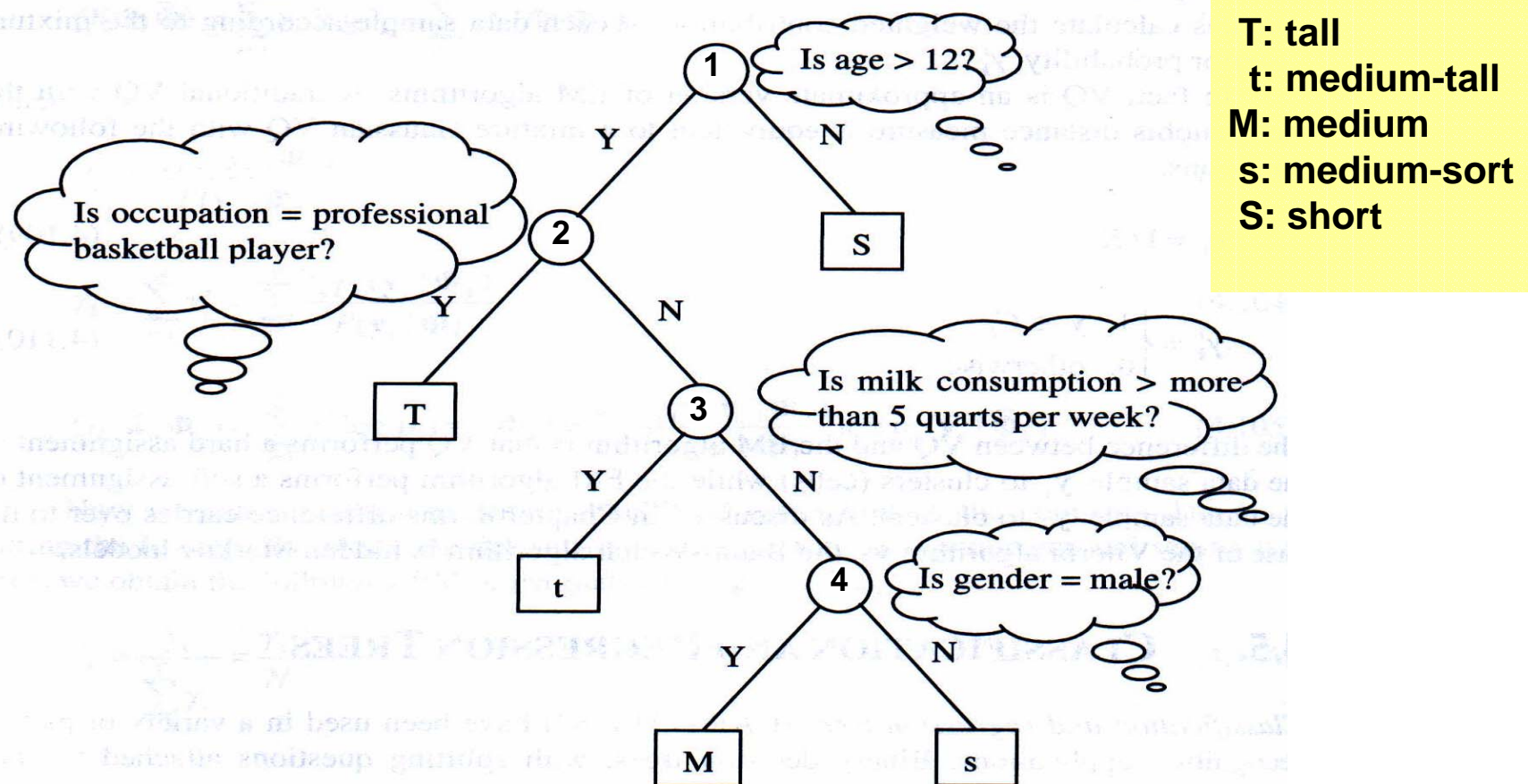


Figure 4.14 A binary tree structure for height classification.

Classification and Regression Trees

- Example: height classification (cont.)
 - Can easily predict the height class for any new person with all the measured data (age, occupation, milk-drinking, etc.) **but no height information**, by traversing the binary tree (based on a set of questions)
 - “No”: right branch, “Yes” left branch
 - When reaching a leaf node, we can use its attached label as the height class for the new person
 - **Also can use the average height in the leaf node to predict the height of the new person**

CART Construction using Training Samples

- **Steps**

1. First, find a set of questions regarding the measured variable
 - E.g., “Is age>12?”, “Is gender=male?”, etc.
2. Then, place all the training samples in the root of the initial tree
3. Choose the **best question** from the question set to split the root into two nodes (**need some measurement !**)
4. Recursively split the most promising node with the **best question** until the right-sized tree is obtained

How to choose the best question?

- E.g., reduce the uncertainty of the event being decided upon
i.e., **find the question which gives the greatest entropy reduction**

CART Construction using Training Samples

- **Splitting Criteria (for discrete pdf)**
 - How to find the best question for a node split ?
 - I.e., find the best split for the data samples of the node
 - Assume training samples have a probability (density) function $P(\omega|t)$ at each node t
 - E.g. $P(\omega_i|t)$ is the percentage of data samples for class i at a node t and $\sum_i P(\omega_i|t) = 1$

CART Construction using Training Samples

- **Splitting Criteria (for discrete pdf)**

- Define the weighted entropy for any tree node t

$$\bar{H}_t(Y) = H_t(Y)P(t)$$

$$H_t(Y) = -\sum_i P(\omega_i|t)\log P(\omega_i|t), \text{ Entropy: average amount of information}$$

- Y is the random variable for classification decision
- $P(t)$ is the prior probability of visiting node t (ratio of numbers of samples in a node t and the total number of samples)

CART Construction using Training Samples

- **Splitting Criteria (for discrete pdf)**

- Entropy reduction for a question q to split a node t into nodes l and r
 - Pick the question with the greatest entropy reduction

$$\Delta \bar{H}_t(q) = \bar{H}_t(Y) - (\bar{H}_l(Y) + \bar{H}_r(Y)) = \bar{H}_t(Y) - \bar{H}_t(Y|q)$$

$$q^* = \arg \max_q [\Delta \bar{H}_t(q)]$$

Review: Fundamentals in Information Theory

- Three interpretations for quantity of information

1. The amount of **uncertainty** before seeing an event
2. The amount of **surprise** when seeing an event
3. The amount of **information** after seeing an event

- The definition of information:

$$I(x_i) = \log \frac{1}{P(x_i)} = -\log P(x_i)$$

- $P(x_i)$ the probability of an event x_i

- Entropy: the average amount of information

$$H(X) = E[I(X)]_X = E[-\log P(x_i)]_X = \sum_{x_i} -P(x_i) \cdot \log P(x_i)$$

where $S = \{x_1, x_2, \dots, x_i, \dots\}$

- Have maximum value when the probability (mass) function is a uniform distribution

CART Construction using Training Samples

- **Splitting Criteria (for discrete pdf)**

- **Example**

$$X:\{x_i=1, 1, 3, 3, 8, 8, 9, 9\}$$

$$P(x=1)=1/4$$

$$P(x=3)=1/4$$

$$P(x=8)=1/4$$

$$P(x=9)=1/4$$

$$H=-4*(1/4)\log_2(1/4)=2$$

① $Y:\{y_i=1, 1\}$
 $P(y=1)=1$

$Z:\{z_i=3, 3, 8, 8, 9, 9\}$

$P(z=3)=1/3$

$P(z=8)=1/3$

$P(z=9)=1/3$

$H_l = -1*(1)\log_2(1)=0;$

$\bar{H}_l = H_l \cdot P(\text{Node}_l) = 0 \cdot 1/4 = 0$

$H_r = -3*(1/3)\log_2(1/3)=1.6;$

$\bar{H}_r = H_r \cdot P(\text{Node}_r) = 1.6 \cdot 3/4 = 1.2$

$\bar{H} = \bar{H}_l + \bar{H}_r = 1.2$

② $Y:\{y_i=1, 1, 3, 3\}$
 $P(y=1)=1/2$
 $P(y=3)=1/2$

$Z:\{z_i=8, 8, 9, 9\}$

$P(z=8)=1/2$

$P(z=9)=1/2$

$H_l = -2*(1/2)\log_2(1/2)=1;$

$H_r = -2*(1/2)\log_2(1/2)=1;$

$\bar{H}_l = H_l \cdot P(\text{Node}_l) = 1 \cdot 1/2 = 1/2$

$\bar{H}_r = H_r \cdot P(\text{Node}_r) = 1 \cdot 1/2 = 1/2$

$\bar{H} = \bar{H}_l + \bar{H}_r = 1.0$

CART Construction using Training Samples

- **Entropy for a tree**

- the sum of weighted entropies for all terminal nodes

$$\bar{H}(T) = \sum_{t \text{ is terminal}} \bar{H}_t(Y)$$

- It can be show that the above tree-growing (splitting) procedure repeatedly reduces the entropy of the tree
- The resulting tree has a better classification power

CART Construction using Training Samples

- **Splitting Criteria** (for continuous pdf)

- the likelihood gain is often used instead of the entropy measure
- Suppose one split divides the data X into two groups X_1 and X_2 , which can be respectively represented as two Gaussian distributions $N_1(\mu_1, \Sigma_1)$ and $N_2(\mu_2, \Sigma_2)$

$$L_1(X_1 | N_1) = \log \prod_{x_1} N(x_1; \mu_1, \Sigma_1)$$

$$L_2(X_2 | N_2) = \log \prod_{x_2} N(x_2; \mu_2, \Sigma_2)$$

See textbook P. 179-180 and complete the derivation
Due 12/9

Log likelihood gain at node t

$$\Delta \bar{L}_t(q) = L_1(X_1 | N_1) + L_2(X_2 | N_2) - L_X(X | N)$$

= ...

$$= (a + b) \log |\Sigma| - a \log |\Sigma_1| - b \log |\Sigma_2|$$

a, b are the sample counts for X_1 and X_2