### 8.0 Search Algorithms for Speech Recognition

#### References: 1. 12.1-12.5 of Huang, or

- 2. 7.2-7.6 of Becchetti, or
- 3. 5.1-5.7, 6.1-6.5 of Jelinek
- Progress in Dynamic Programming Search for LVCSR (Large Vocabulary Continuous Speech Recognition)", Proceedings of the IEEE, Aug 2000

## **Basic Approach for Large Vocabulary Speech Recognition**

#### A Simplified Block Diagram



Example Input Sentence

this is speech

Acoustic Models

(th-ih-s-ih-z-s-p-ih-ch)

• **Lexicon**  $(th-ih-s) \rightarrow this$ 

 $(ih-z) \rightarrow is$  $(s-p-iy-ch) \rightarrow speech$ 

• Language Model (this) – (is) – (speech)

P(this) P(is | this) P(speech | this is) $P(w_i|w_{i-1})$  $P(w_i|w_{i-1},w_{i-2})$  $P(w_i|w_{i-1},w_{i-2})$  $P(w_i|w_{i-1},w_{i-2})$ 

# **DTW and Dynamic Programming**

#### • Dynamic Time Warping (DTW)

- well accepted pre-HMM approach
- find an optimal path for matching two templates with different length
- good for small-vocabulary isolated-word recognition even today

#### • Test Template [y<sub>j</sub>, j=1,2,...N] and Reference Template [x<sub>i</sub>, i=1,2,...M]

- warping both templates to a common length L warping functions:  $f_x(m) = i$ ,  $f_v(m) = j$ , m = 1, 2, ... L
- endpoint constraints:  $f_x(1) = f_y(1) = 1$ ,  $f_x(L) = M$ ,  $f_y(L) = N$ monotonic constraints:  $f_x(m+1) \ge f_x(m)$ ,  $f_y(m+1) \ge f_y(m)$
- matching pair:  $x_{f_x(m)} \leftrightarrow y_{f_y(m)}$  for every m, m: index for matching pairs
- recursive relationship:

 $D(i, j) = \min_{(i', j')} \{D(i', j') + \overline{d}[(i', j'); (i, j)]\}, \quad D(i, j) : \text{accumulated minimum distance up to } (i, j)$ 

d[(i', j'); (i, j)]: additional distance extending (i', j') to (i, j)

d (i, j) = distance measure for  $x_i$  and  $y_j$ 

examples:  $\overline{d}[(i-1, j-1); (i, j)] = d(i, j)$ 

$$\overline{d}[(i-k, j-1); (i, j)] = \frac{1}{k} [d(i-k+1, j) + \dots + d(i-1, j) + d(i, j)]$$

- global constraints/local constraints
- lack of a good approach to train a good reference pattern

#### Dynamic Programming

- replacing the problem by a smaller sub-problem and formulating an iterative procedure

#### • Reference: 4.7 up to 4.7.3 of Rabiner and Juang

### DTW

#### progr<u>am</u>ming



## DTW





### **Basic Problem 2 for HMM** (P.20 of 4.0)

- Approach 2 —Viterbi Algorithm finding the single best sequence  $\overline{q}^* = q_1^* q_2^* \dots q_T^*$ 
  - Define a new variable  $\delta_t(i)$ 
    - $\delta_{t}(i) = \max_{q_{1}, q_{2}, \dots, q_{t-1}} P[q_{1}, q_{2}, \dots, q_{t-1}, q_{t} = i, o_{1}, o_{2}, \dots, o_{t} | \lambda]$ 
      - = the highest probability along a certain single path ending at state i at time t for the first t observations, given  $\lambda$
  - Induction

 $\delta_{t+1}(j) = \max_{i} \left[ \delta_t(i) a_{ij} \right] \bullet b_j(o_{t+1})$ 

- Backtracking  $\psi_{t}(j) = \arg \max_{1 \le i \le N} [\delta_{t-1}(i)a_{ij}]$ 

the best previous state at t–1 given at state j at time t keeping track of the best previous state for each j and t

## Viterbi Algorithm (P.21 of 4.0)



 $\delta_{t}(i) = \max_{q_{1},q_{2},...,q_{t-1}} P[q_{1},q_{2},...,q_{t-1}, q_{t} = i, o_{1},o_{2},...,o_{t} |\lambda]$ 

## Viterbi Algorithm (P.22 of 4.0)



## **Continuous Speech Recognition Example: Digit String Recognition— One-stage Search**

- Unknown Number of Digits
- No Lexicon/Language Model Constraints
- Search over a 3-dim Grid



## **Recognition Errors**



$$\frac{T - D - S - I}{T} \times 100\% = \text{Accuracy}$$

## **Continuous Speech Recognition Example: Digit String Recognition — Level-Building**

- Known Number of Digits
- No Lexicon/Language Model Constraints
- Higher Computation Complexity, No Deletion/Insertion





- number of levels = number of digits in an utterance
- automatic transition from the last state of the previous model to the first state of the next model

## **Time (Frame)- Synchronous Viterbi Search for Large-Vocabulary Continuous Speech Recognition**

#### •MAP Principle

 $W^* = \mathop{\operatorname{arg max}}_{W} [p(W|O)] = \mathop{\operatorname{arg max}}_{W} [\frac{p(O|W)p(W)}{p(O)}] = \mathop{\operatorname{arg max}}_{W} [p(O|W)p(W)]$   $p(O|W) = \sum_{all \ \overline{q}} p(O, \overline{q}|W), \ \overline{q}: a \ state \ sequence \ from \ from \ Language \ Model \ HMM$ 

#### •An Approximation

$$W^* = \mathop{\operatorname{arg\,max}}_{W} \left[ p(W) \sum_{a \parallel \overline{q}} p(O, \overline{q} | W) \right] \cong \mathop{\operatorname{arg\,max}}_{W} \left[ p(W) \cdot \mathop{\operatorname{max}}_{\overline{q}} p(O, \overline{q} | W) \right]$$

- the word sequence with the highest probability for the most probable state sequence usually has approximately the highest probability for all state sequences
- Viterbi search, a sub-optimal approach

#### •Viterbi Search—Dynamic Programming

- replacing the problem by a smaller sub-problem and formulating an iterative procedure
- time (frame)- synchronous: the best score at time t is updated from all states at time t-1
- •Tree Lexicon as the Basic Working Structure



- each arc is an HMM (phoneme, tri-phone, etc.)
- each leaf node is a word
- search processes for a segment of utterance through some common units for different words can be shared
- search space constrained by the lexicon
- the same tree copy reproduced at each leaf node in principle

#### **Basic Problem 2 for HMM** (P.24 of 4.0)

#### • Application Example of Viterbi Algorithm

- Isolated word recognition

$$\lambda_0 = (\mathbf{A}_0, \mathbf{B}_0, \boldsymbol{\pi}_0)$$
$$\lambda_1 = (\mathbf{A}_1, \mathbf{B}_1, \boldsymbol{\pi}_1)$$
$$\vdots$$
$$\lambda_n = (\mathbf{A}_n, \mathbf{B}_n, \boldsymbol{\pi}_n)$$

observation

 $\overline{O} = (o_1, o_2, \dots o_T)$   $k^* = \arg \max_{1 \le i \le n} P[\overline{O} \mid \lambda_i] \approx \arg \max_{1 \le i \le n} [P^* \mid \lambda_i]$   $\widehat{\Box}$ Basic Problem 1
Basic Problem 1
Forward Algorithm
(for all paths)
(for a single best path)

-The model with the highest probability for the most probable path usually also has the highest probability for all possible paths.

### **Tree Lexicon**



 $o_1 o_2 \dots o_t \dots o_T$ 



## **Time (Frame)- Synchronous Viterbi Search for Large –Vocabulary Continuous Speech Recognition**

#### Define Key Parameters

D (t,  $q_t$ , w) : objective function for the best partial path ending at time t in state  $q_t$  for the word w

 $h(t, q_t, w)$ : backtrack pointer for the previous state at the pervious time when the best partial path ends at time t in state  $q_t$  for the word w

#### • Intra-word Transition—HMM only, no Language Model

$$D(t, q_t, w) =_{q_{t-1}}^{\max} [d(o_t, q_t | q_{t-1}, w) + D(t-1, q_{t-1}, w)]$$
  
$$d(o_t, q_t | q_{t-1}, w) = \log p(o_t | q_t, w) + \log p(q_t | q_{t-1}, w)$$
  
$$\overline{q}(t, q_t, w) =_{q_{t-1}}^{\arg \max} [d(o_t, q_t | q_{t-1}, w) + D(t-1, q_{t-1}, w)]$$

 $h(t,q_t,w) = \overline{q}(t,q_t,w)$ 

#### • Inter-word Transition—Language Model only, no HMM (bi-gram as an example)

 $D(t,Q,w) = \int_{v}^{max} [\log p(v|u) + D(t,q_f(v),v)]$ 

*u* : the word before *v* 

Q:a pseudo initial state for the word w

 $q_f(v)$ : the final state for the word v

 $\overline{v}: \sup_{v} \max_{v} \left[\log p(v|u) + D(t, q_f(v), v)\right]$ 

 $\mathbf{h}(t, Q, w) = q_f(\overline{v})$ 

## **Time Synchronous Viterbi Search**



O.

18

## Viterbi Algorithm (P.21 of 4.0)



 $\delta_{t}(i) = \max_{q_{1},q_{2},...,q_{t-1}} P[q_{1},q_{2},...q_{t-1}, q_{t} = i, o_{1},o_{2},...,o_{t} |\lambda]$ 



## **Time (Frame)- Synchronous Viterbi Search for Large-Vocabulary Continuous Speech Recognition**

#### Beam Search

- at each time t only a subset of promising paths are kept
- example 1: define a beam width L (i.e. keeping only L paths at each time) example 2: define a threshold Th (i.e. all paths with D<  $D_{max,t}$ -Th are deleted)
- very helpful in reducing the search space
- Two-pass Search (or Multi-pass Search)



- use less knowledge or less constraints (e.g. acoustic model with less context dependency or language model with lower order) in the first stage, while more knowledge or more constraints in rescoring in the second path
- search space significantly reduced by decoupling the complicated search process into simpler processes
- N-best List and Word Graph (Lattice)



- similarly constructed with dynamic programming iterations

## **Some Search Algorithm Fundamentals**

• An Example – a city traveling problem



• Search Tree(Graph)



- Heuristic Search
  - Best-first Search
  - based on some knowledge, or "heuristic information"  $f(n) = g(n) + h^*(n)$ 
    - $f(n) = g(n) + h^*(n)$

g(n): distance up to node n

h\*(n): heuristic estimate for the remaining distance up to G

heuristic pruning

- S: starting G: goal
- to find the minimum distance path

### Blind Search Algorithms

- Depth-first Search: pick up an arbitrary alternative and proceed
- Breath-first Search: consider all nodes on the same level before going to the next level

11.5

11.8

- no sense about where the goal is

h<sup>\*</sup>(n): straight-line distance

## **Heuristic Search: Another Example**

• Problem: Find a path with the highest score from root node "A" to some leaf node (one of "L1","L2","L3","L4")

 $f(n) = g(n) + h^*(n)$ 

g(n): score from root node to node n



#### List of Candidate Steps

Тор	Candidate List			
A(15)	A(15)			
C(15)	C(15), B(13), D(7)			
G(14)	G (14), B(13), F(9), D(7)			
B(13)	B(13), L3(12), F(9), D(7)			
L3(12)	L3 (12), E(11), F(9), D(7)			

Node	<u>g(n)</u>	<u>h*(n)</u>	<u><b>f</b>(n)</u>	
Α	0	15	15	
В	4	9	13	
С	3	12	15	
D	2	5	7	
$\mathbf{E}$	7	4	11	
$\mathbf{F}$	7	2	9	
G	11	3	14	
L1	9	0	9	
L2	8	0	8	
L3	12	0	12	
L4	5	0	5	

# **A\* Search and Speech Recognition**

## Admissibility

- a search algorithm is admissible if it is guaranteed that the first solution found is optimal, if one exists (for example, beam search is NOT admissible)

### • It can be shown

- the heuristic search is admissible if

 $h^*(n) \ge h(n)$  for all n with a highest-score problem

 $-A^*$  search when the above is satisfied

### Procedure

 Keep a list of next-step candidates, and proceed with the one with the highest f(n) (for a highest-score problem)

## • A\* search in Speech Recognition

- example 1: use of weak constraints in the first pass to generate heuristic estimates in multi-pass search
- example 2: estimated average score per frame as the heuristic information

 $s_f = \left[\log P(\overline{o}_{i,j} | \overline{q}_{i,j})\right] / (j - i + 1)$ 

 $\overline{o}_{i,j}$ : observations from frame i to j,  $\overline{q}_{i,j}$ : state sequences from frame i to j estimated with many (i, j) pairs from training data

 $h^*(n)$  obtained from Max  $[s_f]$ , Ave $[s_f]$ , Min  $[s_f]$  and (T - t)