WFSTs in ASR &

Basics of Speech Production





Instructor: Preethi Jyothi

Lecture 6



Determinization/Minimization: Recap

- A (W)FST is deterministic if:
 - Unique start state
 - •
 - No epsilon input labels •
- number of states (and transitions)
 - automaton
- and the weight semiring (allowing for weight pushing)

No two transitions from a state share the same input label

Minimization finds an equivalent deterministic FST with the least

 For a deterministic weighted automaton, weight pushing + (unweighted) automata minimization leads to a minimal weighted

Guaranteed to yield a deterministic/minimized WFSA under some technical conditions characterising the automata (e.g. twins property)

Example: Dictionary WFST



Determinized Dictionary WFST



Minimized Dictionary WFST









C⁻¹: Arc labels: "monophone : phone / left-context_right-context"

Figure reproduced from "Weighted Finite State Transducers in Speech Recognition", Mohri et al., 2002









Constructing the Decoding Graph



- Construct decoding search graph using $H \circ C \circ L \circ G$ that maps acoustic states to word sequences
- Carefully construct D using optimization algorithms:
 - $D = min(det(H \circ det(C \circ det(L \circ G))))$
- Decode test utterance O by aligning acceptor X (corresponding to O) with $H \circ C \circ L \circ G$:
 - $W = out[\pi]$

Decoding graph, $D = H \circ C \circ L \circ G$

```
W^* = \arg\min \mathbf{X} \circ \mathbf{H} \circ \mathbf{C} \circ \mathbf{L} \circ \mathbf{G}
```

where π is a path in the composed FST, out[π] is the output label sequence of π

"Weighted Finite State Transducers in Speech Recognition", Mohri et al., Computer Speech & Language, 2002

Constructing the Decoding Graph



"Weighted Finite State Transducers in Speech Recognition", Mohri et al., Computer Speech & Language, 2002

Constructing the Decoding Graph



- Each f_k maps to a distinct triphone HMM state j
- Weights of arcs in the ith chain link correspond to observation probabilities $b_i(o_i)$
- X is a very large FST which is never explicitly constructed!
- H o C o L o G is typically traversed dynamically (search algorithms will be covered later in the semester)

Impact of WFST Optimizations 40K NAB Evaluation Set '95 (83% word accuracy)



Tables from http://www.openfst.org/twiki/pub/FST/FstHltTutorial/tutorial_part3.pdf

states	transitions
1,339,664	$3,\!926,\!010$
8,606,729	$11,\!406,\!721$
$7,\!082,\!404$	9,836,629
$7,\!273,\!035$	$10,\!201,\!269$
$18,\!317,\!359$	$21,\!237,\!992$

x real-time
12.5
1.2
1.0
0.7

Toolkits to work with finite-state machines

- AT&T FSM Library (no longer supported) implement.shtml
- RWTH FSA Toolkit
- Carmel https://www.isi.edu/licensed-sw/carmel/
- MIT FST Toolkit ullethttp://people.csail.mit.edu/ilh/fst/
- OpenFST Toolkit (actively supported) http://www.openfst.org/twiki/bin/view/FST/WebHome

http://www3.cs.stonybrook.edu/~algorith/implement/fsm/

https://www-i6.informatik.rwth-aachen.de/~kanthak/fsa.html

Brief Introduction to the OpenFST Toolkit

Quick Intro to OpenFst (www.openfst.org)



Quick Intro to OpenFst (www.openfst.org)



a:a/0.5



an	а	0.5
<eps></eps>	n	1.0
a	а	0.5

Compiling & Printing FSTs

The text FSTs need to be "compiled" into binary objects before further use with **OpenFst utilities**

• Command used to compile:

fstcompile --isymbols=in.txt --osymbols=out.txt A.txt A.fst

•

fstprint --isymbols=in.txt --osymbols=out.txt A.fst A.txt

Get back the text FST using a print command with the binary file:

Composing FSTs

The text FSTs need to be "compiled" into binary objects before further use with OpenFst utilities

• Command used to compose:

fstcompose A.fst B.fst AB.fst

• appropriately sorted before composition

> fstarcsort --sort_type=olabel A.fst |\ fstcompose - B.fst AB.fst

OpenFST requirement: One or both of the input FSTs should be

Drawing FSTs

Small FSTs can be visualized easily using the draw tool:

fstdraw --isymbols=in.txt --osymbols=out.txt A.fst |\ dot -Tpdf > A.pdf



FSTs can get very large!



Basics of Speech Production

Speech Production





Schematic representation of the vocal organs

Schematic from L.Rabiner and B.-H.Juang , Fundamentals of speech recognition, 1993 Figure from <u>http://www.phon.ucl.ac.uk/courses/spsci/iss/week6.php</u>

Sound units

- Phones are acoustically distinct units of speech
- Phonemes are abstract linguistic units that impart different meanings in a given language
 - Minimal pair: <u>p</u>an vs. <u>b</u>an
- Allophones are different acoustic realisations of the same phoneme
- Phonetics is the study of speech sounds and how they're produced
- Phonology is the study of patterns of sounds in different languages

 Sounds produced with no obstruction to the flow of air through the vocal tract





Vowels

Image from https://en.wikipedia.org/wiki/File:IPA_vowel_chart_2005.png

Formants of vowels

- Formants are resonance frequencies of the vocal tract (denoted by F1, F2, etc.)
- F0 denotes the fundamental frequency of the periodic source (vibrating vocal folds)
- Formant locations specify certain vowel characteristics

Vowel	F1(Hz)	F2(Hz)	F3(Hz)
i	280	2620	3380
Ι	360	2220	2960
e	600	2060	2840
æ	800	1760	2500
Λ	760	1320	2500
a:	740	1180	2640
D	560	920	2560
3:	480	760	2620
U	380	940	2300
u	320	920	2200
31	560	1480	2520

Adult male formant frequencies in Hertz collected by J.C.Wells around 1960.

Spectrogram

- time, with amplitude of the frequency components expressed as a heat map
- Spectrograms of certain vowels:
- formants/pitch curves, etc.)

Spectrogram is a sequence of spectra stacked together in

http://www.phon.ucl.ac.uk/courses/spsci/iss/week5.php

 Praat (<u>http://www.fon.hum.uva.nl/praat/</u>) is a good toolkit to analyse speech signals (plot spectrograms, generate

Consonants (voicing/place/manner)

 "Consonants are made by restricting or blocking the airflow in some way, and may be voiced or unvoiced." (J&M, Ch. 7)

Voiced/Unvoiced Sounds

- Sounds made with vocal cords vibrating: voiced
 - E.g. /g/, /d/, etc.
 - All English vowel sounds are voiced
- - E.g. /k/, /t/, etc.

Sounds made without vocal cord vibration: voiceless

Consonants (voicing/place/manner)

- "Consonants are made by restricting or blocking the airflow in some way, and may be voiced or unvoiced." (J&M, Ch. 7)
- Consonants can be labeled depending on
 - where the constriction is made
 - how the constriction is made

Place of articulation



- Bilabial (both lips)
 [b],[p],[m], etc.
- Labiodental (with lower lip and upper teeth)
 [f], [v], etc.
- Interdental (tip of tongue between teeth)
 [θ] (thought), [δ] (this)

Manner of articulation



- Plosive/Stop (airflow completely blocked followed by a release) [p],[g],[t],etc.
- Fricative (constricted airflow) [f], [s], [th], etc.
- Affricate (stop + fricative) [ch], [jh], etc.
- Nasal (lowering velum) [n], [m], etc.

See realtime MRI productions of vowels and consonants here: <u>http://sail.usc.edu/span/rtmri_ipa/je_2015.html</u>