



Instructor: Preethi Jyothi

HMMs and WFSTs

Lecture 4

Recap: HMMs for Acoustic Modeling



1. What is the forward algorithm? What is it used to compute?

tion sequence *O*, determine the likelihood $P(O|\lambda)$.



 $Q = q_1 q_2 q_3 \dots q_T.$



3. What is the Baum-Welch algorithm? What does it compute?

What are the three fundamental problems related to HMMs?

- **Computing Likelihood:** Given an HMM $\lambda = (A, B)$ and an observa-
- 2. What is the Viterbi algorithm? What is it used to compute?
 - **Decoding**: Given as input an HMM $\lambda = (A, B)$ and a sequence of observations $O = o_1, o_2, ..., o_T$, find the most probable sequence of states

Learning: Given an observation sequence O and the set of possible states in the HMM, learn the HMM parameters A and B.

Monophone HMMs? Not good enough. Need Triphones.

- A phone is affected by its phonetic context.
 - E.g. Coarticulation: Production of a speech sound is affected by adjacent speech sounds. "soon" vs. "seat". "ten" vs. "tenth".
- For modelling, use phones in context instead of monophones. E.g. diphones or triphones.
- Triphones are commonly used in ASR systems. Phone p with left context I and right context r is written as "I-p+r"

^{- &}quot;hello world" \rightarrow sil-h+eh h-eh+l eh-l+ow l-ow+w ow-w-er w-er+l er-l+d l-d+sil

Training



Estimate AM parameters $\theta = \{A_{jk}, (\mu_{jm}, \Sigma_{jm}, C_{jm})\}$ over all triphone states. Use Baum-Welch.

HMM of *i*th training utterance determined by using a word-to-triphone mapping applied to $W_1^l, \ldots, W_{\ell i}^l$

Estimate LM parameters β using $\mathbf{w}^1, \ldots, \mathbf{w}^N$

Overall Summary



 $W^* = \arg\max_{W} P(W|O) = \arg\max_{W} P_{\theta}(O|W)P_{\beta}(W)$

Compute using Viterbi algorithm

Search all possible state sequences arising from all word sequences most likely to have generated **O**

Computationally infeasible for continuous speech! WFSTs to the rescue, to make this more tractable.





What are Weighted Finite State Transducers (WFSTs)?

(Weighted) Automaton



- Accepts a subset of strings (over an alphabet), and rejects the rest Mathematically, specified by $L \subseteq \Sigma^*$ or equivalently $f : \Sigma^* \rightarrow \{0,1\}$ •
- Weighted: outputs a "weight" as well (e.g., probability)
 - $f: \Sigma^* \to W$
- Transducer: outputs another string (over possibly another alphabet)

•
$$f: \Sigma^* \times \Delta^* \to W$$



(Weighted) Finite State Automaton



- Functions that can be implemented using a machine which:
 - reads the string one symbol at a time
 - has a fixed amount of memory: so, at any moment, the machine can be in only one of finitely many states, irrespective of the length of the input string
- Allows efficient algorithms to reason about the machine
 - e.g., output string with maximum weight for input $\alpha\beta\gamma$





- Powerful enough to (reasonably) model processes in language, speech, computational biology and other machine learning applications
 - WFSTs, e.g., speech recognition systems
- the models and to make inferences

Why WFSTs?

Simpler WFSTs can be combined to create complex

If using WFST models, efficient algorithms available to train

Toolkits that don't have domain specific dependencies

Structure: Finite State Transducer (FST)



Elements of an FST

- States
- Start state (0) •

FST maps input strings to output strings

- Arcs (transitions)
- Input symbols (from alphabet Σ)
- Final states (1 & 2) \cdot Output symbols (from alphabet Δ)



- A successful "path" → Sequence of transitions from the start state to any final state
- Input label of a path → Concatenation of input labels on arcs.
 Similarly for output label of a path.



- Finite state acceptors (FSAs)
 - a label
 - FSA accepts a set of strings, $L \subseteq \Sigma^*$

- Finite state transducers (FSTs)
 - Each transition has a source & destination state, • an input label and an output label
 - FST represents a relation, $R \subseteq \Sigma^* \times \Delta^*$

FSAs and FSTs

Each transition has a source & destination state, and

FSA can be thought of as a special kind of FST



Example of an FSA

Accepts strings {c, a, ab}

Equivalent FST

Accepts strings {c, a, ab} and outputs identical strings {c, a, ab}

Barking dog FST



 $Σ = { yelp, bark }, Δ = { a, ..., }$

yelp $\rightarrow y i \rho$. bark $\rightarrow \mu$

Special symbol, ε (epsilon) : allows to make a move without consuming an input symbol

or without producing an output symbol



- sequence to an output sequence
- How are the weights accumulated along a path?

 "Weights" can be probabilities, negative log-likelihoods, or any cost function representing the cost incurred in mapping an input

Weighted Path: Probabilistic FST



 $= w(e_1) \times w(e_2) = 0.25 \times 1.0 = 0.25$

$$w(e) = \Pr[e \text{ taken } | \text{ state=0, in-symbol=} \alpha]$$

Weighted Path: Probabilistic FST



- $T(\alpha\beta,ab) = Pr[output=ab, accepts | input=\alpha\beta, start]$
- More generally, $T(x,y) = \Sigma_{\pi \in P(x,y)} \prod_{e \in \pi} w(e)$

$$w(e) = Pr[e \text{ taken } | \text{ state}=0, \text{ in-symbol}=\alpha]$$

= $\Pr[\pi_1 \mid \text{input} = \alpha\beta, \text{ start }] + \Pr[\pi_2 \mid \text{input} = \alpha\beta, \text{ start }]$

where P(x,y) is the set of all accepting paths with input x and output y

Weighted Path



- But not all WFSTs are probabilistic FSTs •
- But helpful to retain some basic algebraic properties of weights: abstracted as semirings

Weight is often a "score" and maybe accumulated differently

Semirings

and \otimes , along with their identity values $\overline{0}$ and $\overline{1}$

- A semiring is a set of values associated with two operations \oplus
 - Weight assigned to an input/output pair $T(x,y) = \bigoplus_{\pi \in P(x,y)} \bigotimes_{e \in \pi} w(e)$
- where P(x,y) is the set of all accepting paths with input x, output y
 - (generalizing the weight function for a probabilistic FST)

Semirings

Some popular semirings [M02]

SEMIRING	SET	\oplus	\bigotimes	$\overline{0}$	1	
Boolean	{F,T}	V	\wedge	F	Τ	Is there
Real	\mathbb{R}_+	+	X	0	1	
Log	ℝ ∪ {-∞, +∞ }	⊕log	₽	+∞	0	Log Dr[1/1/
Tropical	ℝ∪ {-∞, +∞}	min	÷	+∞	0	"Viterbi Approx." of -log Pr[v x]
Onarata	nr Au dofinad as v	θ. v	— _Io	a (o-x	т о-/	

Operator \bigoplus_{\log} defined as: $x \oplus_{\log} y = -\log(e^{-x} + e^{-y})$

[M02] Mohri, M. Semiring frameworks and algorithms for shortest-distance problems, Journal of Automata, Languages and Combinatorics, 7(3):321-350, 2002.



- Weight of a path π is the \otimes -product of all the transitions in π • $w(\pi)$: (0.5 \otimes 1.0) = 0.5 + 1.0 = 1.5
- Weight of a sequence " $\underline{x},\underline{y}$ " is the \oplus -sum of all paths labeled with " $\underline{x},\underline{y}$ " W((an), $(a \ B)$) = (1.5 $\oplus \overline{0}$) = min(1.5, ∞) = 1.5



• W((an), (a n)) = ?Path 1: $(0.5 \otimes 1.0) = 1.5$ Path 2: $(0.3 \otimes 0.1) = 0.4$ Weight of "((an), (a n))" = (1.5 \oplus 0.4) = 0.4

Weight of a sequence " $\underline{x},\underline{y}$ " is the \oplus -sum of all paths labeled with " $\underline{x},\underline{y}$ "

Shortest Path

• Recall $T(x,y) = \bigoplus_{\pi \in P(x,y)} w(\pi)$

- P(x,y) : Shortest Path

where $P(x,y) = \text{set of paths with input/output } (x,y); w(\pi) = \bigotimes_{e \in \pi} w(e)$

• In the tropical semiring \bigoplus is min. T(x,y) associated with a single path in

• Can be found using Dijkstra's algorithm : $\Theta(|E| + |Q| \cdot \log|Q|)$ time

 $\alpha: A/0.25$ $\alpha: a/0.25$ $\alpha: a/0.5$ $\alpha: a/0.$

ß:B/1.0

 $T(``\alpha", ``a") = ?$ $T(``\alpha\alpha", ``aa") = ?$

Shortest Path



Inversion



This operation comes in handy, especially during composition!

Swap the input and output labels in each transition

Weights (if they exist) are retained on the arcs

Projection





Project onto the input or output alphabet

a:a/0.5

Basic FST Operations (Rational Operations)

The set of weighted transducers are closed under the following operations [Mohri '02]:

- 1. Sum or Union: $(T_1 \oplus T_2)(\underline{x}, \underline{y}) = T_1(\underline{x}, \underline{y}) \oplus T_2(\underline{x}, \underline{y})$
- 2. Product or Concatenation: (T
- 3. Kleene-closure: $T^*(x, y) = \bigoplus_{i=1}^{\infty}$ n=0

$$\Gamma_1 \otimes T_2(\mathbf{x}, \mathbf{y}) = \bigoplus_{\substack{\mathbf{x} = \mathbf{x}_1 \mathbf{x}_2 \\ \mathbf{y} = \mathbf{y}_{1y_2}}} T_1(\mathbf{x}_1, \mathbf{y}_1) \otimes T_2(\mathbf{x}_2, \mathbf{y}_2)$$

Basic FST Operations (Rational Operations)

The set of weighted transducers are closed under the following operations [Mohri '02]:

1. Sum or Union: $(T_1 \oplus T_2)(\underline{x}, \underline{y}) = T_1(\underline{x}, \underline{y}) \oplus T_2(\underline{x}, \underline{y})$

3. Kleene-closure: $T^*(x, y) = \int_{-\infty}^{\infty} T^n(x, y)$



Example: Recall Barking Dog



Example: Union



Animal farm!

Basic FST Operations (Rational Operations)

The set of weighted transducers are closed under the following operations [Mohri '02]:

- 1. Sum or Union: $(T_1 \oplus T_2)(\underline{x}, \underline{y}) = T_1(\underline{x}, \underline{y}) \oplus T_2(\underline{x}, \underline{y})$

3. Kleene-closure: $T^*(x, y) = \bigoplus T^n(x, y)$

2. Product or Concatenation: $(T_1 \otimes T_2)(x, y) = \bigoplus T_1(x_1, y_1) \otimes T_2(x_2, y_2)$ $X = X_1 X_2$ **y=y**_{1y2} 00

Example: Concatenation

Suppose the last "baa" in a bleat should be followed by one or more a's

(e.g., "baabaa" is not OK, but "baaa" and "baabaaaaa" are)

Basic FST Operations (Rational Operations)

The set of weighted transducers are closed under the following operations [Mohri '02]:

- 1. Sum or Union: $(T_1 \oplus T_2)(\underline{x}, \underline{y}) = T_1(\underline{x}, \underline{y}) \oplus T_2(\underline{x}, \underline{y})$
- 3. Kleene-closure: $T^*(x, y) = \bigoplus T^n(x, y)$ n=0

2. Product or Concatenation: $(T_1 \otimes T_2)(x, y) = \bigoplus T_1(x_1, y_1) \otimes T_2(x_2, y_2)$ $X = X_1 X_2$ **y=y**1y2

Example: Closure

Animal farm: allow arbitrarily long sequence of sounds!

bark moo yelp bleat $\rightarrow \omega \circ \sigma$

Acoustic Model WFST

Composition

 If T₁ transduces x to z, and transduces x to y

$$(T_1 \circ T_2)(x, y) = \bigoplus_{Z} T_1(x, y)$$

• If T_1 transduces x to z, and T_2 transduces z to y, then $T_1 \circ T_2$

z) T₂ (z, y)

Composition: Construction

 $M_1 \circ M_2$

 $M_1 \circ M_2$

Composition

Composition: Example 1

6:B

 $M_1 \circ M_2$

Composition: Example 2

