

NLP 201: Sequence Models

Jeffrey Flanigan

University of California Santa Cruz
jmflanig@ucsc.edu

Slides and figures adapted from David Bamman, Diyi Yang, and Noah Smith

Fall 2023

Plan

- Tagging (sequence labeling)
- Hidden Markov Models (HMMs)
- Viterbi algorithm
- Applications (BIO tagging, etc)
- Semi-rings
- Weighted automata

Sequence Labeling (Tagging)

The process of assigning a tag to each word in an input text.

- **Input:** a sequence of tokenized words
- **Output:** a sequence of tags, one per token

Sequence Labeling (Tagging)

Part-of-Speech (POS) tagging: assign a POS tag to each word.

We will use POS tagging as an example, but these methods can be used for any task that can be formulated as sequence labeling, such as named entity recognition (NER), supertagging, noun phrase chunking, etc.

Why is POS Tagging Useful?

- First step of a vast number of practical tasks
- Parsing
 - Need to know if a word is an N or V before you can parse
- Information extraction
 - Finding names, relations, etc.
- Speech synthesis/recognition
 - Object obJECT
 - OVERflow overFLOW
 - DIScount disCOUNT
 - CONtent conTENT
- Machine Translation

Tagging Is A Disambiguation Task

- Words often have more than one POS: *back*
 - The *back* door = JJ
 - On my *back* = NN
 - Win the voters *back* = RB
 - Promised to *back* the bill = VB

Let's Work on POS Tagging

- Mrs Shaefer never got around to joining
- All we do is go around the corner
- Chateau Petrus costs around 250

Let's Work on POS Tagging

- Mrs/NNP Shaefer/NNP never/RB got/VBD around/RP to/TO joining/VBG
- All/DT we/PRP do/VB is/VBZ go/VB around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ around/RB 250/CD

POS Tags in the Penn Treebank

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	's	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRP\$	possess. pronoun	<i>your, one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WP\$	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative adverb	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	\$
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	#
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &</i>	"	left quote	' or "
LS	list item marker	<i>1, 2, One</i>	TO	"to"	<i>to</i>	"	right quote	' or "
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(left paren	[, (, {, <
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>)	right paren],), }, >
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	,
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	. ! ?
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	: ; ... - -

Figure 8.1 Penn Treebank part-of-speech tags (including punctuation).

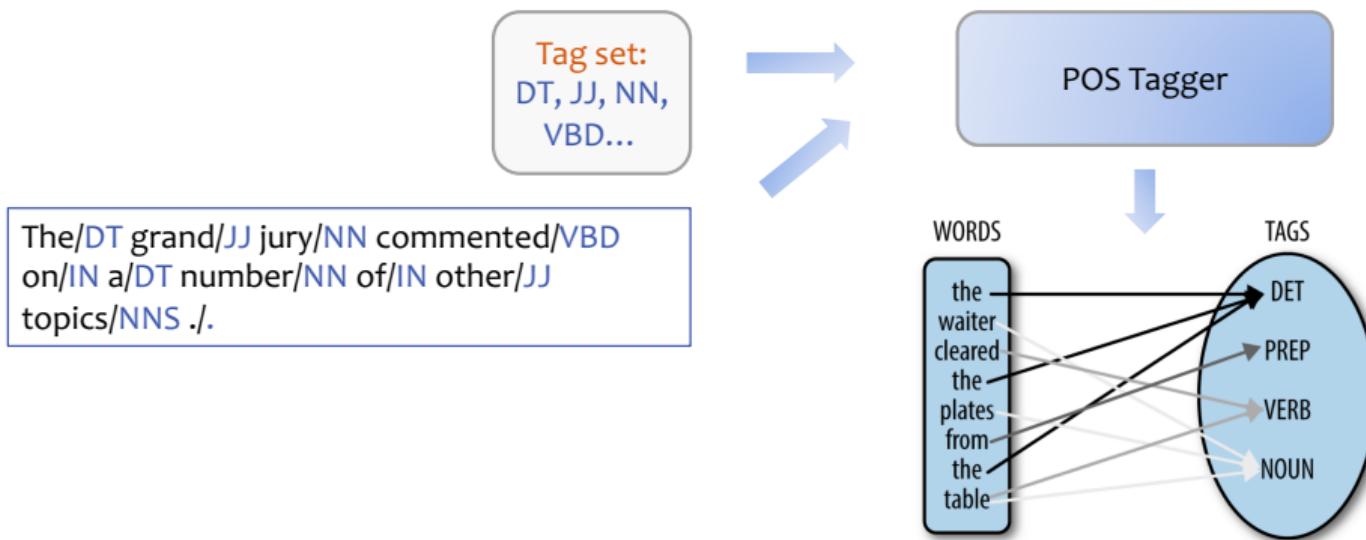
Source: SPLP 3 Ch 8

Deciding On the Correct POS Can Be Difficult for People

- Mrs/NNP Shaefer/NNP never/RB got/VBD **around/RP**
to/TO joining/VBG
- All/DT we/PRP do/VB is/VBZ go/VB **around/IN** the/DT
corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ **around/RB** 250/CD

Building a POS Tagger

- Supervised Learning
 - Assume linguistics have annotated several examples



Statistical POS Tagging

- What is the **most likely sequence of tags** for the given **sequence of words w**

$$\begin{aligned}\operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}|\mathbf{w}) &= \operatorname{argmax}_{\mathbf{t}} \frac{P(\mathbf{t}, \mathbf{w})}{P(\mathbf{w})} \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}, \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t})P(\mathbf{w}|\mathbf{t})\end{aligned}$$

Statistical POS Tagging

- What is the **most likely sequence of tags** for the given **sequence of words w**

$$\begin{aligned} P(\text{DT JJ NN} \mid \text{a smart dog}) &= P(\text{DD JJ NN a smart dog}) / P(\text{a smart dog}) \\ &\propto P(\text{DD JJ NN a smart dog}) \\ &= P(\text{DD JJ NN}) P(\text{a smart dog} \mid \text{DD JJ NN}) \end{aligned}$$

Transition Probability

- Joint probability $P(\mathbf{t}, \mathbf{w}) = P(\mathbf{t})P(\mathbf{w}|\mathbf{t})$
- $P(\mathbf{t}) = P(t_1, t_2, \dots, t_n)$
$$= P(t_1)P(t_2 | t_1)P(t_3 | t_2, t_1) \dots P(t_n | t_1 \dots t_{n-1})$$
$$\sim P(t_1)P(t_2 | t_1)P(t_3 | t_2) \dots P(t_n | t_{n-1})$$
$$= \prod_{i=1}^n P(t_i | t_{i-1})$$

Markov Assumption

- Bigram model over POS tags!
(similarly, we can define a n-gram model over POS tags, usually we called high-order HMM)

Emission Probability

- Joint probability $P(\mathbf{t}, \mathbf{w}) = P(\mathbf{t})P(\mathbf{w}|\mathbf{t})$
- **Assume** words only depend on their POS-tag
- $P(\mathbf{w}|\mathbf{t}) \sim P(w_1 | t_1)P(w_2 | t_2) \dots P(w_n | t_n)$
 $= \prod_{i=1}^n P(w_i | t_i)$

i.e., $P(\text{a smart dog} | \text{ DD JJ NN })$

$$= P(\text{a} | \text{DD}) P(\text{smart} | \text{JJ}) P(\text{dog} | \text{NN})$$

Independent Assumption

Putting them Together

- Two independent assumptions
 - Approximate $p(t)$ by a bi(or N)-gram model
 - **Markov assumption:** $P(t_i|t_1 \dots t_{i-1}) = P(t_i|t_{i-1})$
 - Assume each word depends only on its POS tag
 - **Output independence:** $P(w_i|t_1 \dots t_i, w_1 \dots w_{i-1} \dots w_T) = P(w_i|t_i)$

Putting them Together

- Joint probability $P(\mathbf{t}, \mathbf{w}) = P(\mathbf{t})P(\mathbf{w}|\mathbf{t})$

- $P(\mathbf{t}, \mathbf{w})$

$$\begin{aligned}&= P(t_1)P(t_2 | t_1)P(t_3 | t_2) \dots P(t_n | t_{n-1})P(w_1 | t_1)P(w_2 | t_2) \dots P(w_n | t_n) \\&= \prod_{i=1}^n P(w_i | t_i)P(t_i | t_{i-1})\end{aligned}$$

e.g., $P(\text{a smart dog , DD JJ NN })$

$$\begin{aligned}&= P(\text{a} | \text{DD}) P(\text{smart} | \text{JJ}) P(\text{dog} | \text{NN}) \\&\quad P(\text{DD} | \text{start}) P(\text{JJ} | \text{DD}) P(\text{NN} | \text{JJ})\end{aligned}$$

Board work: Hidden Markov Models (HMMs)

(5 min) On your own

- Draw the Bayesian network for this model

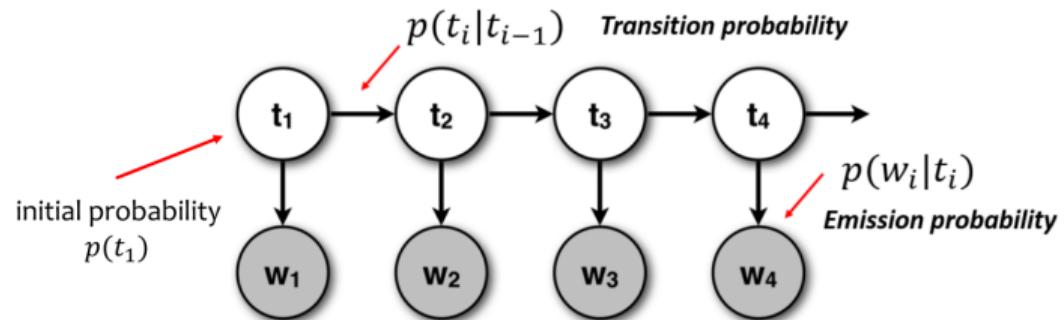
$$P(\mathbf{t}, \mathbf{w}) = \left(\prod_{t=1}^n P(w_i | t_i) P(t_i | t_{i-1}) \right) P(t_{n+1} | t_n)$$

t_0 is the <START> symbol, and t_{n+1} is the <STOP> symbol

- If you were given a corpus of tagged sentences, how would you estimate the conditional probability tables from the data?

(3 min) Discuss with a partner

HMM



Learning from Labeled Data

- We count how often we see $t_{i-1}t_i$ and w_it_i then normalize.

Learning the transition probabilities:

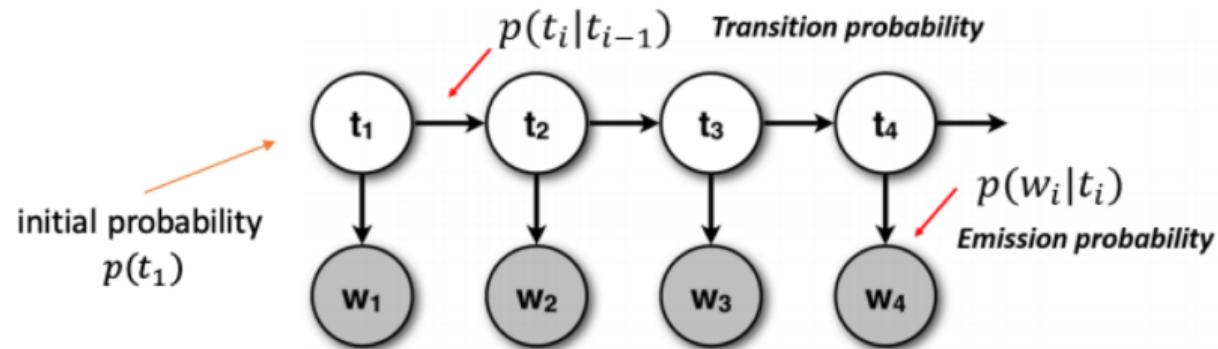
$$P(t_j|t_i) = \frac{C(t_i t_j)}{C(t_i)}$$

Learning the emission probabilities:

$$P(w_j|t_i) = \frac{C(w_j - t_i)}{C(t_i)}$$

Finding the Best Tag Sequence

What is the most likely sequence of tags for given sequence of words w



Viterbi algorithm

- Basic idea: if an optimal path through a sequence uses **label L at time T**, then it must have used an optimal path to get to label L at time T
- We can discard all non-optimal paths up to label L at time T

END							
DT							
NNP							
VB							
NN							
MD							
START							
	^	Janet	will	back	the	bill	\$

- At each time step t ending in label K, we find the max probability of any path that led to that state

END		
DT		v ₁ (DT)
NNP		v ₁ (NNP)
VB		v ₁ (VB)
NN		v ₁ (NN)
MD		v ₁ (MD)
START		

Janet

What's the HMM probability of ending in Janet = NNP?

$$P(y_t \mid y_{t-1})P(x_t \mid y_t)$$

$$P(\text{NNP} \mid \text{START})P(\text{Janet} \mid \text{NNP})$$

END		
DT		$v_1(DT)$
NNP		$v_1(NNP)$
VB		$v_1(VB)$
NN		$v_1(NN)$
MD		$v_1(MD)$
START		

Best path through time step 1
ending in tag y (trivially - best
path for all is just START)

Janet

$$v_1(y) = \max_{\mathbf{u} \in \mathcal{Y}} [P(y_t = y \mid \mathbf{y}_{t-1} = \mathbf{u}) P(x_t \mid y_t = y)]$$

END			
DT		$v_1(DT)$	$v_2(DT)$
NNP		$v_1(NNP)$	$v_2(NNP)$
VB		$v_1(VB)$	$v_2(VB)$
NN		$v_1(NN)$	$v_2(NN)$
MD		$v_1(MD)$	$v_2(MD)$
START			

Janet will

What's the **max** HMM probability of ending in will = MD?

First, what's the HMM probability of a single path
ending in will = MD?

END			
DT		v ₁ (DT)	v ₂ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)
VB		v ₁ (VB)	v ₂ (VB)
NN		v ₁ (NN)	v ₂ (NN)
MD		v ₁ (MD)	v ₂ (MD)
START			

Janet will

$$P(y_1 \mid START)P(x_1 \mid y_1) \times P(y_2 = \text{MD} \mid y_1)P(x_2 \mid y_2 = \text{MD})$$

END			
DT		$v_1(DT)$	$v_2(DT)$
NNP		$v_1(NNP)$	$v_2(NNP)$
VB		$v_1(VB)$	$v_2(VB)$
NN		$v_1(NN)$	$v_2(NN)$
MD		$v_1(MD)$	$v_2(MD)$
START			

Best path through time step 2
ending in tag MD

Janet will

$$P(DT \mid START) \times P(Janet \mid DT) \times P(y_t = MD \mid P(y_{t-1} = DT) \times P(will \mid y_t = MD))$$

$$P(NNP \mid START) \times P(Janet \mid NNP) \times P(y_t = MD \mid P(y_{t-1} = NNP) \times P(will \mid y_t = MD))$$

$$P(VB \mid START) \times P(Janet \mid VB) \times P(y_t = MD \mid P(y_{t-1} = VB) \times P(will \mid y_t = MD))$$

$$P(NN \mid START) \times P(Janet \mid NN) \times P(y_t = MD \mid P(y_{t-1} = NN) \times P(will \mid y_t = MD))$$

$$P(MD \mid START) \times P(Janet \mid MD) \times P(y_t = MD \mid P(y_{t-1} = MD) \times P(will \mid y_t = MD))$$

END			
DT		$v_1(DT)$	$v_2(DT)$
NNP		$v_1(NNP)$	$v_2(NNP)$
VB		$v_1(VB)$	$v_2(VB)$
NN		$v_1(NN)$	$v_2(NN)$
MD		$v_1(MD)$	$v_2(MD)$
START			

Janet will

Best path through time step 2
ending in tag MD

Let's say the best path ending $\text{will} = \text{MD}$ includes $\text{Janet} = \text{NNP}$. By definition, every other path ending in $\text{will} = \text{MD}$ has lower probability.

END			
DT		$v_1(DT)$	$v_2(DT)$
NNP		$v_1(NNP)$	$v_2(NNP)$
VB		$v_1(VB)$	$v_2(VB)$
NN		$v_1(NN)$	$v_2(NN)$
MD		$v_1(MD)$	$v_2(MD)$
START			

Best path through time step 2
ending in tag MD

Janet will

$$P(DT \mid START) \times P(Janet \mid DT) \times P(y_t = MD \mid P(y_{t-1} = DT) \times P(will \mid y_t = MD))$$

$$P(NNP \mid START) \times P(Janet \mid NNP) \times P(y_t = MD \mid P(y_{t-1} = NNP) \times P(will \mid y_t = MD))$$

$$P(VB \mid START) \times P(Janet \mid VB) \times P(y_t = MD \mid P(y_{t-1} = VB) \times P(will \mid y_t = MD))$$

$$P(NN \mid START) \times P(Janet \mid NN) \times P(y_t = MD \mid P(y_{t-1} = NN) \times P(will \mid y_t = MD))$$

$$P(MD \mid START) \times P(Janet \mid MD) \times P(y_t = MD \mid P(y_{t-1} = MD) \times P(will \mid y_t = MD))$$

To calculate this full probability, notice that we can reuse information we've already computed.

$$\frac{P(\text{DT} \mid \text{START}) \times P(\text{Janet} \mid \text{DT}) \times P(y_t = \text{MD} \mid P(y_{t-1} = \text{DT})) \times P(\text{will} \mid y_t = \text{MD})}{v_1(\text{DT})}$$

$$\frac{P(\text{NNP} \mid \text{START}) \times P(\text{Janet} \mid \text{NNP}) \times P(y_t = \text{MD} \mid P(y_{t-1} = \text{NNP})) \times P(\text{will} \mid y_t = \text{MD})}{v_1(\text{NNP})}$$

$$\frac{P(\text{VB} \mid \text{START}) \times P(\text{Janet} \mid \text{VB}) \times P(y_t = \text{MD} \mid P(y_{t-1} = \text{VB})) \times P(\text{will} \mid y_t = \text{MD})}{v_1(\text{VB})}$$

...

END			
DT		v ₁ (DT)	v ₂ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)
VB		v ₁ (VB)	v ₂ (VB)
NN		v ₁ (NN)	v ₂ (NN)
MD		v ₁ (MD)	v ₂ (MD)
START			

Janet will

$$v_t(y) = \max_{\textcolor{magenta}{u} \in \mathcal{Y}} [v_{t-1}(\textcolor{magenta}{u}) \times P(y_t = y \mid \textcolor{magenta}{y_{t-1}} = u)P(x_t \mid y_t = y)]$$

END				
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)
START				

Janet will back

25 paths ending in back = VB

END				
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)
START				

Janet will back

Let's say the best path ending in **back = VB** includes
will = MD.

END				
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)
START				

Janet will back

If the best path ending in **will = MD** includes
 Janet=NNP, we can forget all paths with Janet \neq NNP
 for any path including **will = MD** because we know
 they are less likely.

END					
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)	v ₄ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)	v ₄ (NNP)
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)	v ₄ (MD)
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)	v ₄ (NN)
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)	v ₄ (MD)
START					

Janet will back the

125 possible paths ending in the = DT, but we only need
 to consider 5 (best path ending in back = DT, back =
 NNP, back = VB, back = NN, back = MD)

END						
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)	v ₄ (DT)	v ₅ (DT)
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)	v ₄ (NNP)	v ₅ (NNP)
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)	v ₄ (MD)	v ₅ (MD)
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)	v ₄ (NN)	v ₅ (NN)
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)	v ₄ (MD)	v ₅ (MD)
START						

Janet will back the bill

END							$v_T(\text{END})$
DT		$v_1(\text{DT})$	$v_2(\text{DT})$	$v_3(\text{DT})$	$v_4(\text{DT})$	$v_5(\text{DT})$	
NNP		$v_1(\text{NNP})$	$v_2(\text{NNP})$	$v_3(\text{NNP})$	$v_4(\text{NNP})$	$v_5(\text{NNP})$	
VB		$v_1(\text{VB})$	$v_2(\text{VB})$	$v_3(\text{VB})$	$v_4(\text{MD})$	$v_5(\text{MD})$	
NN		$v_1(\text{NN})$	$v_2(\text{NN})$	$v_3(\text{NN})$	$v_4(\text{NN})$	$v_5(\text{NN})$	
MD		$v_1(\text{MD})$	$v_2(\text{MD})$	$v_3(\text{MD})$	$v_4(\text{MD})$	$v_5(\text{MD})$	
START							

Janet will back the bill

$v_T(\text{END})$ encodes the best path through the entire sequence

END							$v_T(\text{END})$
DT							
NNP							
VB							
NN							
MD							
START							

Janet will back the bill

For each timestep $t + \text{label}$, keep track of the max element from $t-1$ to reconstruct best path

Formal Algorithm

Let's formalize the Viterbi Algorithm

Hidden Markov Models (Formal)

- States $T = t_1, t_2 \dots t_N$ (N POS tags)
- Observations $W = w_1, w_2 \dots w_M$ (M words in a sentence)
 - Each observation is a symbol from a vocabulary V
- Transition probabilities
 - Transition probability matrix $A = \{a_{ij}\}$
- Observation likelihoods
 - Output probability matrix $B = \{b_i(w_t)\}$
- Initial probability vector π
 - $\pi_i = P(t_1 = i) \quad 1 \leq i \leq N$

The Components of An HMM Tagger

- An HMM has two components: A and B Probabilities
- A:
 - Tag transition probabilities $p(t_i|t_{i-1})$
 - The prob of a tag occurring given the previous tag
 - Compute the MLE by counting $p(t_i|t_{i-1}) = \frac{c(t_{i-1}, t_i)}{c(t_{i-1})}$
 - Example:
 - Modal verb **will** and a verb in the base form, a VB, like **race**

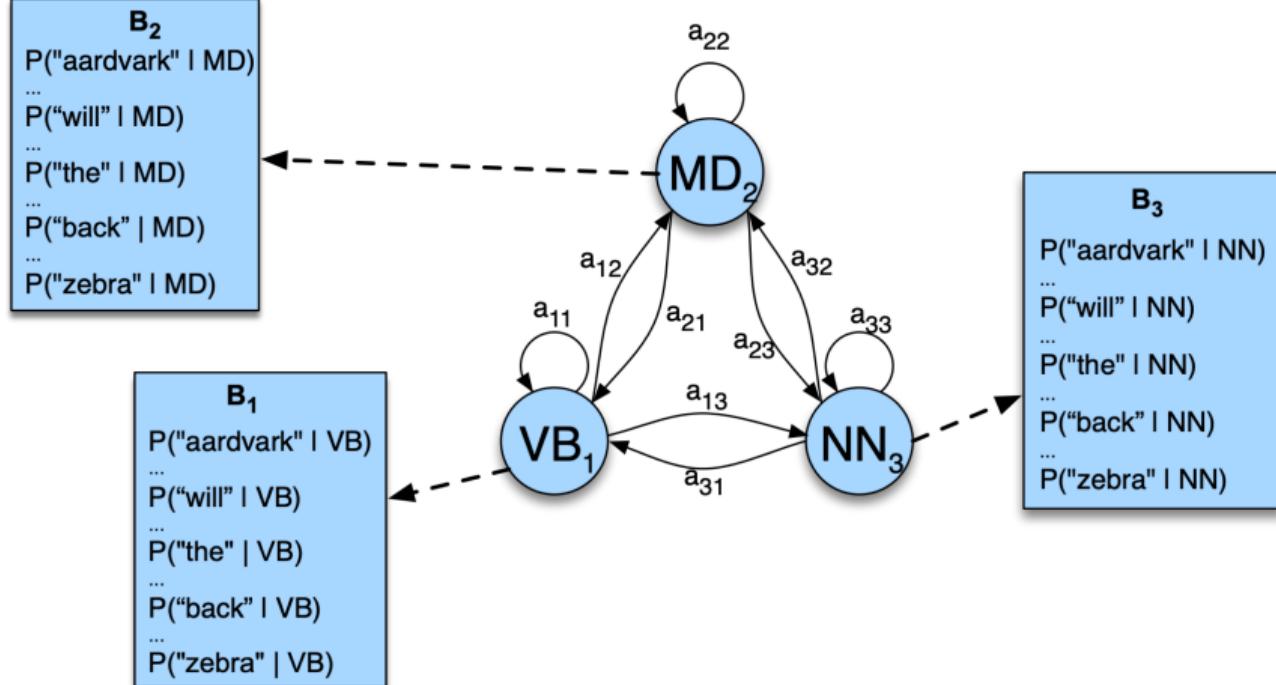
$$P(VB|MD) = \frac{C(MD, VB)}{C(MD)} = \frac{10471}{13124} = .80$$

The Components of An HMM Tagger

- An HMM has two components: A and B Probabilities
- B:
 - Emission probabilities $p(w_i|t_i)$
 - Given a tag (e.g., MD), that it will be associated with a given word (e.g., **will**)
 - Compute the MLE by counting $p(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$
 - Example:
 - Modal verb MD and the word **will**

$$P(\text{will}|MD) = \frac{C(MD, \text{will})}{C(MD)} = \frac{4046}{13124} = .31$$

The Components of An HMM Tagger



An illustration of the two parts of an HMM representation: the A transition probabilities used to compute the prior probability, and the B observation likelihoods that are associated with each state, one likelihood for each possible observation word.

Viterbi Algorithm

```
function VITERBI(observations of len T,state-graph of len N) returns best-path, path-prob
    create a path probability matrix viterbi[N,T]
    for each state s from 1 to N do           ; initialization step
        viterbi[s,1]← $\pi_s * b_s(o_1)$ 
        backpointer[s,1]←0
    for each time step t from 2 to T do       ; recursion step
        for each state s from 1 to N do
            viterbi[s,t]← $\max_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$ 
            backpointer[s,t]← $\operatorname{argmax}_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$ 
        bestpathprob← $\max_{s=1}^N viterbi[s,T]$            ; termination step
        bestpathpointer← $\operatorname{argmax}_{s=1}^N viterbi[s,T]$        ; termination step
    bestpath←the path starting at state bestpathpointer, that follows backpointer[] to states back in time
    return bestpath, bestpathprob
```

Figure A.9 Viterbi algorithm for finding optimal sequence of hidden states. Given an observation sequence and an HMM $\lambda = (A, B)$, the algorithm returns the state path through the HMM that assigns maximum likelihood to the observation sequence.

Parts of speech

- Parts of speech are categories of words defined **distributionally** by the morphological and syntactic contexts a word appears in.

Morphological distribution

POS often defined by distributional properties; verbs
= the class of words that each combine with the
same set of affixes

	-s	-ed	-ing
walk	walks	walked	walking
slice	slices	sliced	slicing
believe	believes	believed	believing
of	*ofs	*ofed	*ofing
red	*reds	*redded	*reding

Morphological distribution

We can look to the function of the affix (denoting past tense) to include irregular inflections.

	-s	-ed	-ing
walk	walks	walked	walking
sleep	sleeps	slept	sleeping
eat	eats	ate	eating
give	gives	gave	giving

Syntactic distribution

- Substitution test: if a word is replaced by another word, does the sentence remain **grammatical**?

Kim saw the

elephant

before we did

dog

idea

*of

*goes

Syntactic distribution

- These can often be too strict; some contexts admit substitutability for some pairs but not others.

Kim saw the

elephant

before we did

*Sandy

both nouns but
common vs. proper

Kim *arrived the

elephant

before we did

both verbs but
transitive vs. intransitive

Nouns	People, places, things, actions-made-nouns ("I like swimming"). Inflected for singular/plural
Verbs	Actions, processes. Inflected for tense, aspect, number, person
Adjectives	Properties, qualities. Usually modify nouns
Adverbs	Qualify the manner of verbs ("She ran <i>downhill extremely quickly yesterday</i> ")
Determiner	Mark the beginning of a noun phrase ("a dog")
Pronouns	Refer to a noun phrase (he, she, it)
Prepositions	Indicate spatial/temporal relationships (<i>on</i> the table)
Conjunctions	Conjoin two phrases, clauses, sentences (and, or)

Sequence labeling

$$x = \{x_1, \dots, x_n\}$$

$$y = \{y_1, \dots, y_n\}$$

- For a set of inputs x with n sequential time steps, one corresponding label y_i for each x_i

Named entity recognition

B-PERS I-PERS O O O O ORG

tim cook is the ceo of apple

3 or 4-class:

- person
- location
- organization
- (misc)

7-class:

- person
- location
- organization
- time
- money
- percent
- date

BIO tagging

For segmentation tasks, can use BIO tagging.

- B (with a label) denotes the start of a segment (with that label)
- I (with a label) denotes continuing the segment
- O denotes “outside” i.e. not a segment

Variations to this tagging scheme can improve performance.

A popular one: use B only if two spans are next to each other, otherwise default to I for the start.

Supersense tagging

O B-artifact I-artifact B-motion O B-time O O O B-group

The station wagons arrived at noon, a long shining line

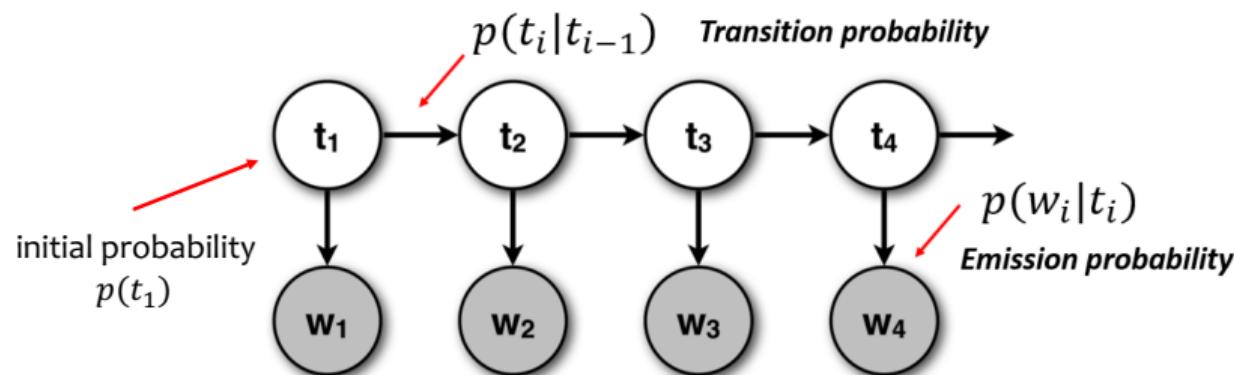
O B-motion O O B-location I-location

that coursed through the west campus.

1	person	7	cognition	13	attribute	19	quantity	25	plant
2	communication	8	possession	14	object	20	motive	26	relation
3	artifact	9	location	15	process	21	animal		
4	act	10	substance	16	Tops	22	body		
5	group	11	state	17	phenomenon	23	feeling		
6	food	12	time	18	event	24	shape		

Tagging

- What is the **most likely sequence of tags** for the given **sequence of words w**



Three Basic Problems for HMMs

- **Likelihood** of the input: How likely the sentence "I love cat" occurs
 - Forward algorithm
- **Decoding** (tagging) the input: POS tags of "I love cat" occurs
 - Viterbi algorithm
- **Estimation** (learning): How to learn the model?
 - Find the best model parameters
 - Case 1: supervised – tags are annotated (MLE)
 - Case 2: unsupervised -- only unannotated text (Forward-backward algorithm)

Last time: Viterbi Algorithm

END							v _T (END)
DT		v ₁ (DT)	v ₂ (DT)	v ₃ (DT)	v ₄ (DT)	v ₅ (DT)	
NNP		v ₁ (NNP)	v ₂ (NNP)	v ₃ (NNP)	v ₄ (NNP)	v ₅ (NNP)	
VB		v ₁ (VB)	v ₂ (VB)	v ₃ (VB)	v ₄ (MD)	v ₅ (MD)	
NN		v ₁ (NN)	v ₂ (NN)	v ₃ (NN)	v ₄ (NN)	v ₅ (NN)	
MD		v ₁ (MD)	v ₂ (MD)	v ₃ (MD)	v ₄ (MD)	v ₅ (MD)	
START							

Janet will back the bill

Each cell keeps track of the score of the best tag sequence ending in that tag at time i . Let's call this $\pi(i, y_i)$

Derivation: Viterbi Algorithm (Dynamic Programming)

$$\hat{y} = \operatorname{argmax}_{y_1 \dots y_n} p(x_1 \dots x_n, y_1 \dots y_n)$$

- Define $\pi(i, y_i)$ to be the max score of a sequence of length i ending in tag y_i

$$\begin{aligned}\pi(i, y_i) &= \max_{y_1 \dots y_{i-1}} p(x_1 \dots x_i, y_1 \dots y_i) \\ &= \max_{y_{i-1}} p(x_i | y_i) p(y_i | y_{i-1}) \max_{y_1 \dots y_{i-2}} p(x_1 \dots x_i, y_1 \dots y_{i-1}) \\ &= \max_{y_{i-1}} p(x_i | y_i) p(y_i | y_{i-1}) \pi(i-1, y_{i-1})\end{aligned}$$

- We now have an efficient algorithm. Start with $i = 0$ and work your way to the end of the sentence.

Sum over all tag sequences: Forward Algorithm

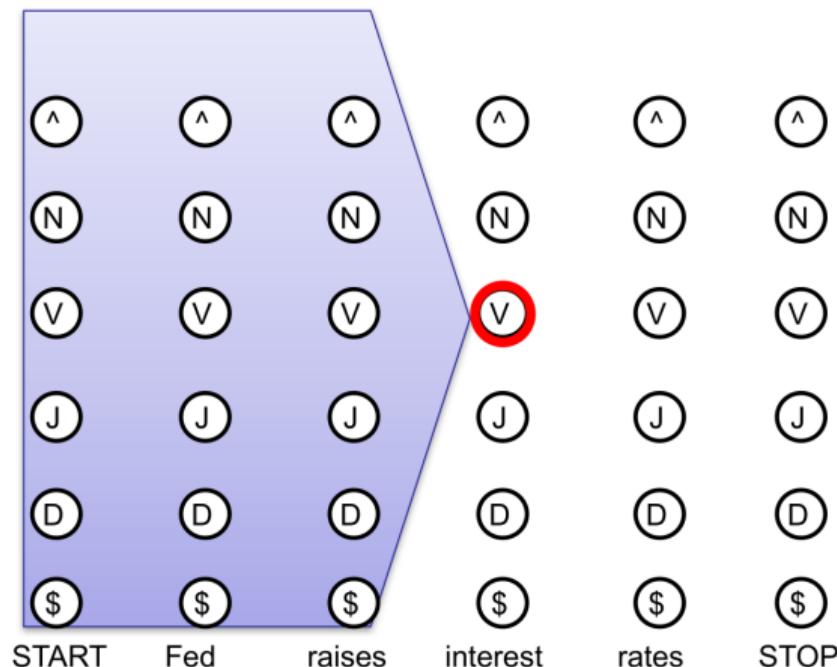
$$p(x_1 \dots x_n) = \sum_{y_1 \dots y_n} p(x_1 \dots x_n, y_1 \dots y_n)$$

- Define $\pi(i, y_i)$ to be the sum over all tag sequences of length i ending in tag y_i

$$\begin{aligned}\pi(i, y_i) &= \sum_{y_1 \dots y_{i-1}} p(x_1 \dots x_i, y_1 \dots y_i) \\ &= \sum_{y_{i-1}} p(x_i | y_i) p(y_i | y_{i-1}) \sum_{y_1 \dots y_{i-2}} p(x_1 \dots x_i, y_1 \dots y_{i-1}) \\ &= \sum_{y_{i-1}} p(x_i | y_i) p(y_i | y_{i-1}) \pi(i-1, y_{i-1})\end{aligned}$$

- Again, we have an efficient algorithm. Start with $i = 0$ and work your way to the end of the sentence.

Sum over all tag sequences: Forward Algorithm



Each cell keeps track of the score of the sum over all sequences ending in that tag at time i . We call this $\pi(i, y_i)$

Runtime of Viterbi and Forward algorithms

- Linear in sentence length n
- Polynomial in the number of possible tags $|K|$

$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

- Specifically:

$O(n|K|)$ entries in $\pi(i, y_i)$

$O(|K|)$ time to compute each $\pi(i, y_i)$

- Total runtime: $O(n|K|^2)$

- Q: Is this a practical algorithm?
- A: depends on $|K|$

Speeding Up Viterbi: Beam search

- **Beam search:** At each timestep, only keep the k best tags
Only need to check k previous tags at each timestep.
 k is called the “beam size”
- This is an **approximate algorithm** (no guarantee it will produce the Viterbi tag sequence)
- Runtime becomes $O(nk|T|)$
- If choose $k = 1$ we get a **greedy algorithm** (It never goes back to consider alternatives to choices it makes. It chooses the best tag sequence in a greedy fashion left to right)

Higher-Order HMMs

- Second order HMM: $P(\mathbf{t}, \mathbf{w}) = \prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1}, t_{i-2})$
- Nth-order HMM: $P(\mathbf{t}, \mathbf{w}) = \prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1}, t_{i-2}, \dots, t_{i-(N-1)})$
- Viterbi and Forward algorithms exist for higher-order HMMs
- Need to keep track of $N - 1$ tags in each Viterbi cell
- Runtime: $O(n|T|^N)$

Working in log space

- Similar to n-gram LMs, if we work with probabilities, we will get underflows from multiplying many small numbers
- $\log(a * b) = \log(a) + \log(b)$
- Viterbi algorithm
 - Use $\log(Pr(t_i|t_{i-1}))$ and $\log(Pr(w_i|t_i))$
 - Instead of multiply, add
 - The Viterbi array stores log probabilities
- Forward algorithm
 - Use $\log(Pr(t_i|t_{i-1}))$ and $\log(Pr(w_i|w_{i-1}))$
 - Instead of multiply, add
 - To add log probabilities, use **log-sum**:
If $a > b$ (otherwise, switch them):
$$\log(a + b) = \log(a(1 + b/a)) = \log(a) + \log(1 + b/a) = \log(a) + \log(1 + e^{\log(b) - \log(a)})$$

Viterbi Algorithm

```
function VITERBI(observations of len T,state-graph of len N) returns best-path, path-prob
    create a path probability matrix viterbi[N,T]
    for each state s from 1 to N do           ; initialization step
        viterbi[s,1] ←  $\pi_s * b_s(o_1)$ 
        backpointer[s,1] ← 0
    for each time step t from 2 to T do       ; recursion step
        for each state s from 1 to N do
            viterbi[s,t] ←  $\max_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$ 
            backpointer[s,t] ←  $\operatorname{argmax}_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$ 
        bestpathprob ←  $\max_{s=1}^N viterbi[s,T]$            ; termination step
        bestpathpointer ←  $\operatorname{argmax}_{s=1}^N viterbi[s,T]$        ; termination step
        bestpath ← the path starting at state bestpathpointer, that follows backpointer[] to states back in time
    return bestpath, bestpathprob
```

Figure A.9 Viterbi algorithm for finding optimal sequence of hidden states. Given an observation sequence and an HMM $\lambda = (A, B)$, the algorithm returns the state path through the HMM that assigns maximum likelihood to the observation sequence.

For details, see Jurafsky and Martin Appendix A.

Forward Algorithm

```
function FORWARD(observations of len  $T$ , state-graph of len  $N$ ) returns forward-prob
    create a probability matrix forward[ $N,T$ ]
    for each state  $s$  from 1 to  $N$  do                                ; initialization step
         $\text{forward}[s,1] \leftarrow \pi_s * b_s(o_1)$ 
    for each time step  $t$  from 2 to  $T$  do          ; recursion step
        for each state  $s$  from 1 to  $N$  do
             $\text{forward}[s,t] \leftarrow \sum_{s'=1}^N \text{forward}[s',t-1] * a_{s',s} * b_s(o_t)$ 
     $\text{forwardprob} \leftarrow \sum_{s=1}^N \text{forward}[s,T]$            ; termination step
    return forwardprob
```

Figure A.7 The forward algorithm, where $\text{forward}[s,t]$ represents $\alpha_t(s)$.

For details, see Jurafsky and Martin Appendix A.

Viterbi vs Forward

If we just want the score of the best sequence (not the actual sequence), then Viterbi and Forward are the same except:

- We have a *max* in Viterbi, and a *sum* in Forward
- Everything else is the same!
- **If we use semirings, we can use the same code to do either algorithm**

(Also easier to code, less likely to have bugs due to backpointers, get k-best outputs for free)

Definition: Semiring

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements A (so far real numbers)
- \otimes generalized times operation
- \oplus generalized plus operation
- 0_s semiring 0 (identity for \oplus : $a \oplus 0_s = a$ for all $a \in A$)
- 1_s semiring 1 (identity for \otimes : $a \otimes 1_s = a$ for
- \otimes, \oplus must satisfy the semiring axioms

Semiring Axioms

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- Satisfies commutative law (for \oplus), distributive law, associative law
- Missing: no inverse for addition (not true: all $a \in A$, exists b such that $a \oplus b = 0_s$)
Rings have additive inverses
- Missing: no inverse for multiplication (not true: all $a \in A$, exists b such that $a \otimes b = 1_s$)
Fields have multiplicative inverses and additive inverses

Example: Real Semiring

Used in Forward algorithm

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R}$: Real numbers
- $a \otimes b = a \times b$
- $a \oplus b = a + b$
- $0_s = 0$ (identity for \oplus : $a \oplus 0_s = a$ for all $a \in A$)
- $1_s = 1$ (identity for \otimes : $a \otimes 1_s = a$ for all $a \in A$)

Example: Viterbi Semiring

Used in Viterbi algorithm to get the max probability

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R}^+$: Real numbers ≥ 0
- $a \otimes b = a \times b$
- $a \oplus b = \max(a, b)$
- $0_s = 0$ (identity for \oplus : $a \oplus 0_s = a$ for all $a \in A$)
- $1_s = 1$ (identity for \otimes : $a \otimes 1_s = a$ for all $a \in A$)

Board work: Viterbi Semiring

The elements in the semiring: $A = \mathbb{R}^+$

$a \otimes b = a \times b$ and $a \oplus b = \max(a, b)$

(5 min) On your own

- Verify semiring 0 and 1 ($0_s = 0$ and $1_s = 1$) satisfy $a \oplus 0_s = a$ for all $a \in A$ and $a \otimes 1_s = a$ for all $a \in A$
- Verify semiring distributive law: $a \oplus (b \otimes c) = (a \otimes b) \oplus (a \otimes c)$

(3 min) Discuss with a partner

Example: Log-Real Semiring

Used in Forward algorithm, in log space to avoid underflows

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R}$: Real numbers
- $a \otimes b = a + b$
- $a \oplus b = \log(e^a + e^b)$
- $0_s = -\infty$
- $1_s = 0$

Example: Max-Plus (Tropical) Semiring

Used in Viterbi algorithm, in log space to avoid underflows

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R}$: Real numbers
- $a \otimes b = a + b$
- $a \oplus b = \max(a, b)$
- $0_s = -\infty$
- $1_s = 0$

Example: Viterbi-Derivation Semiring

Use in Viterbi algorithm to get the max tag sequence and probability score

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R} \times T^*$: Real number and tag sequence
- $(a_v, a_s) \in A$ **value and sequence**
- $(a_v, a_s) \otimes (b_v, b_s) = (a_s \times b_s, a_s.b_s)$ **multiply the numbers, concatenate the tag sequence**
- $(a_v, a_s) \oplus (b_v, b_s) = (a_v, a_s) \text{ if } a_v > b_v, \text{ else } (b_v, b_s)$
- $0_s = (-\infty, \epsilon)$
- $1_s = (1, \epsilon)$

Example: Viterbi-Derivation Semiring, Tropical Version

Use in Viterbi algorithm to get the max tag sequence and log-probability score

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = \mathbb{R} \times T^*$: Real number and tag sequence
- $(a_v, a_s) \in A$ **value and sequence**
- $(a_v, a_s) \otimes (b_v, b_s) = (a_s + b_s, a_s.b_s)$ **add the numbers, concatenate the tag sequence**
- $(a_v, a_s) \oplus (b_v, b_s) = (a_v, a_s) \text{ if } a_v > b_v, \text{ else } (b_v, b_s)$
- $0_s = (-\infty, \epsilon)$
- $1_s = (0, \epsilon)$

Example: k-Best Semiring

Use to get the k best tag sequences and their probability scores

$$S = \langle A, \otimes, \oplus, 0_s, 1_s \rangle$$

- The elements $A = (\mathbb{R} \times T^*)^k$: k Real numbers and tag sequences
- $((a_v^1, a_s^1) \dots (a_v^k, a_s^k)) \in A$
- \otimes : compute all combinations from both lists by multiplying the numbers and concatenating the tag sequence. Keep only the k best
- \oplus : merge the two lists by keeping the k best
- $0_s = (-\infty, \epsilon)^k$
- $1_s = (1, \epsilon)^k$

Implementing a Semiring Viterbi/Forward algorithm

In the code:

- Create a class for your semiring elements (for example, Viterbi-derivation semiring)
 - Implement semiring plus and times functions
 - Implement constructors for transition, emission, starting probabilities
 - Implement constructors for semiring 0 and 1
- Implement Forward algorithm, using semiring plus and times
- To get the Forward algorithm in log space, use Log-Real semiring
- To get the Viterbi algorithm in log space, use the Viterbi-derivation semiring,
Tropical version

Forward Algorithm

```
function FORWARD(observations of len  $T$ , state-graph of len  $N$ ) returns forward-prob
    create a probability matrix forward[ $N, T$ ]
    for each state  $s$  from 1 to  $N$  do ; initialization step
         $\text{forward}[s, 1] \leftarrow \pi_s * b_s(o_1)$ 
    for each time step  $t$  from 2 to  $T$  do ; recursion step
        for each state  $s$  from 1 to  $N$  do
             $\text{forward}[s, t] \leftarrow \sum_{s'=1}^N \text{forward}[s', t-1] * a_{s', s} * b_s(o_t)$ 
     $\text{forwardprob} \leftarrow \sum_{s=1}^N \text{forward}[s, T]$  ; termination step
    return forwardprob
```

Figure A.7 The forward algorithm, where $\text{forward}[s, t]$ represents $\alpha_t(s)$.

For details, see Jurafsky and Martin Appendix A.

Weighted Automata

Many of the models (n-gram language models, HMMs) we've covered so far can be represented as WFSAs and WFSTs

Weighted Automata

- Weighted automata can be understood as transducing strings to weights
- Weighted transducers can be understood as transducing pairs of strings to weights

Weighted Finite-State Automaton

Element	Definition
Q	Finite set of states
Σ	Finite vocabulary
$I \subseteq Q$	Set of initial states
$F \subseteq Q$	Set of final states
$E \subseteq (Q \times (\Sigma \cup \{\epsilon\}) \times Q)$	Set of transitions (edges)
$\lambda : I \rightarrow \mathbb{R}_{\geq 0}$	Initial weights
$\rho : F \rightarrow \mathbb{R}_{\geq 0}$	Final weights
$w : E \rightarrow \mathbb{R}_{\geq 0}$	Transition weights

Paths in WFSAs

- Define the function $p(e)$ to be the **previous state** of an edge e in E
- Define the function $n(e)$ to pick out the **next state** of an edge e
- A **path** π in E^* is a sequence of transitions $e_1e_2\dots e_\ell$ such that $n(e_i)=p(e_{i+1})$ for all $i\in[1,\ell)$
- We overload p and n to be defined on paths
 $p(\pi)=p(e_1) \quad n(\pi)=p(e_\ell)$

Weight of a Path

- We generalize the **transition weight** of the path as the **product** of all transitions

$$w(\pi) = w(e_1) \times w(e_2) \times \cdots \times w(e_\ell)$$

$$= \prod_{i=1}^{\ell} w(e_i)$$

- The **output weight** (or **score**) of the path is defined to be

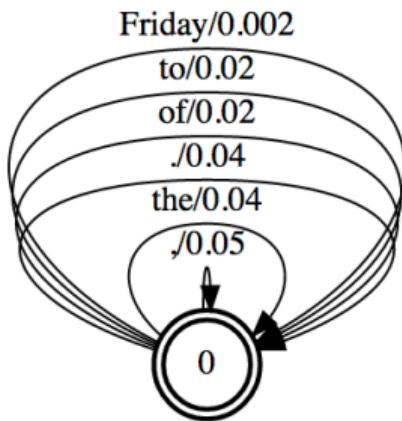
$$score(\pi) = \lambda(p(\pi)) \times w(\pi) \times \rho(n(\pi))$$

Deterministic and Nondeterministic

- Deterministic WFSA
 - There is one initial state q_0
 - For each q in Q and s in Σ , there is at most one r in Q such that $(q, s, r) \in E$.
 - There is no $(q, r) \in (Q \times Q)$ s.t. $(q, \varepsilon, r) \in E$
- Can we determinize WFSAs?
 - Sometimes, but not always
- Can we minimize (deterministic) WFSAs?
 - Yes. There might be more than one way to define “minimization.”

Language Models can be represented as a WFSA

The Simplest Language Model



λ

1

ρ

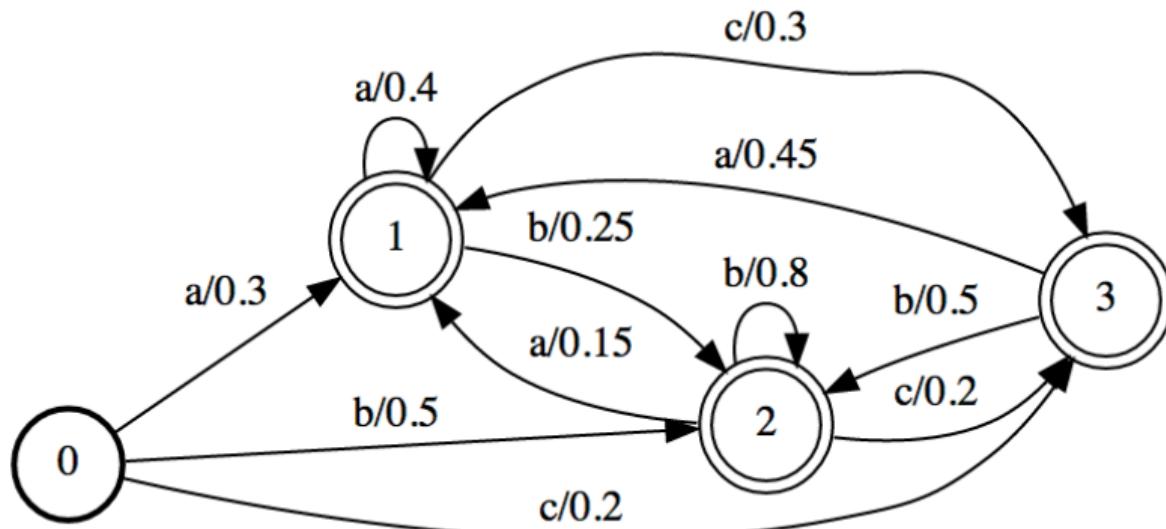
0.04

word	transition weight (w)
,	0.05
the	0.04
.	0.04
of	0.02
to	0.02
...	...
Friday	0.002

What Else Might We Do?

- We can encode more information in the states.
 - More states: more information.
 - Often used to encode history and increase the “memory” of the process.
 - Example: create a distinct state for each (n-1) word history ... this is an **n-gram** language model.

A Simple Bigram Model on {a, b, c}



WFSA

weighted languages
disambiguation

WFST

weighted string-to-string mappings



FSA

regular languages



FST

regular *relations*
string-to-string mappings

Beyond WFSAs: WFSTs

- Weighted finite-state *transducers*: combine the two-tape idea from last time with the weighted idea from today.
 - Very general framework; key operation is **weighted composition**.
 - Best-path algorithm variations: best path and output given input, best path given input and output, ...
- Weights do not have to be numbers.
 - Boolean weights = traditional FSAs/FSTs.
 - Real weights = the case we explored so far.
 - Many other *semirings* (sets of possible weight values and operations for combining weights)

WFSAs and WFSTs: Algorithms

- OpenFST documentation contains high-level view of most algorithms you think you want, and pointers.
- The set of people developing general WFST algorithms largely overlaps with the OpenFST team.
 - Mehryar Mohri, Cyril Allauzen, Michael Reilly
- This is an active area of research!

WFST Examples

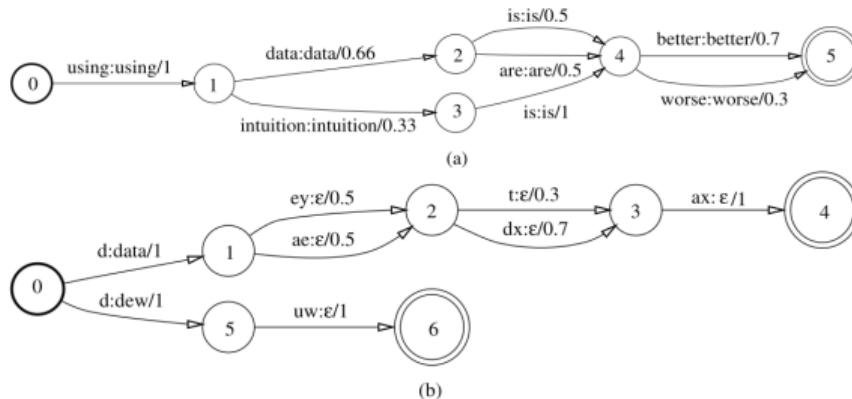


Figure 2: Weighted finite-state transducer examples.

Board Work

HMMs are can be represented as WFSTs: the input is the word sequence, the output is the tag sequence.

The easiest way: one WFST for emission probabilities $Pr(w_i|t_i)$ and one WFST (actually a WFSA as an WFST) for transition probabilities $Pr(t_i|t_{i-1})$.

These two WFSTs are composed to get the HMM

On your own (5 min)

For an HMM with two tags U, V and vocabulary a, b :

- Construct the emission WFST
- Construct the transition WFST

Discuss with a partner (3 min)

Why WFSAs / WFSTs

- Can use off the shelf software to find best path (Viterbi algorithm) and sum over paths (Forward algorithm)
- Example application: speech recognition (WFST for acoustic model, and WFST for n-gram LM)
- Nowadays WFSAs used to analyse seq2seq neural networks (theory)