

Sequence Models - I

Spring 2022



COS 484

Why model sequences?



Part of Speech tagging

Why model sequences?



Named Entity recognition

Why model sequences?

Brazil ranks number 5 in the list of countries by population.

> The term **"Ibu Negara**" (Lady/Mother of the State) is used for wife of the President of Indonesia.

Game of Thrones is an adaptation of A Song of Ice and Fire, George R. R. Martin's series of fantasy novels. It ranks fourth among the IMDB Top Rated TV Shows

THE COUNTRIES WITH THE LARGEST POPULATION

THE COUNTRY'S' FIRST LADIES

Brigitte Macron Melania Trump Iriana Widodo

IMDB TOP RATED TV SHOWS

- **1** Planet Earth II (2016) 9.6. **2** Band of Brothers (2001) 9.5. 3 Planet Earth (2006) 9.5.

- **4** Game of Thrones (2011) 9.4.
- 5 Breaking Bad (2008) 9.4.

- China **1** 1,388,232,693
- India 2 1,342,512,706
- Unites States **3** 326,474,013
 - Indonesia **4** 263,510,146
 - Brasil **5** 174,315,386
- Spouse: Emmanuel Macron, President of France (2017)
- Spouse: Donald J. Trump, U.S. President (2017-)
- Spouse: Joko Widodo, President of Indonesia (2014) - Also known as: "Ibu Negara" (Lady/Mother of the State)

Information Extraction

- Hidden markov models (HMM)
- Viterbi algorithm

Overview

What are part of speech tags?



1. The/DT cat/NN sat/VBD on/IN the/DT mat/NN

2. Princeton/NNP is/VBZ in/IN New/NNP Jersey/NNP

3. The/DT old/NN man/VB the/DT boat/NN

- Word classes or syntactic categories
 - Reveal useful information about a word (and its neighbors!)

- Different words have different functions lacksquare
- Can be roughly divided into two classes
- Closed class: fixed membership, function words
 - e.g. prepositions (*in*, *on*, *of*), determiners (*the*, *a*)
- **Open class**: New words get added frequently
 - e.g. nouns (Twitter, Facebook), verbs (google), adjectives, adverbs

Parts of Speech



• How many part of speech tags do you think English has?

A. < 10

- B. 10 20
- C. 20 40

D. > 40

Parts of Speech







Penn Tree Bank tagset

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

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

Other corpora: Brown, WSJ, Switchboard

45 tags

(Marcus et al., 1993)

Part of Speech Tagging

- Tag each word in a sentence with its part of speech
- Disambiguation task: each word might have different senses/ functions
 - The/DT man/NN bought/VBD a/DT boat/NN
 - The/DT old/NN man/VB the/DT boat/NN

Types:	W	SJ	Brov	wn	
Unambiguous (1 tag	() 44,432	(86%)	45,799	(85%)	
Ambiguous (2+ ta	ngs) 7,025	(14%)	8,050	(15%)	
Tokens:					
Unambiguous (1 tag	s) 577,421	(45%)	384,349	(33%)	
Ambiguous (2+ ta	ngs) 711,780	(55%)	786,646	(67%)	
Figure 8.2 Tag ambiguity for word types in Brown and WSJ, using Treebank-3 (45-tag)					
tagging. Punctuation were treated a	as words, and words w	vere kept i	n their orig	ginal case.	



Same word, different tags



Part of Speech Tagging

- Tag each word with its part of speech
- Disambiguation task: each word might have different senses/ functions
 - The/DT man/NN bought/VBD a/DT boat/NN
 - The/DT old/NN man/VB the/DT boat/NN

earnings growth took a **back/JJ** seat a small building in the **back/NN** a clear majority of senators **back/VBP** the bill Dave began to **back/VB** toward the door enable the country to buy **back/RP** about debt I was twenty-one **back/RB** then



Same word, different tags

Some words have many functions!



A simple baseline

- Many words might be easy to tag
- Most frequent class: Assign each token (word) to the class it occurred most in the training set. (e.g. man/NN)
- Accurately tags 92.34% of word tokens on Wall Street Journal (WSJ)!
- How accurate do you think this baseline would be at tagging words? A) <50%
- B) Average English sentence ~ 14 words C) 75-90%
- D) \rightarrow Sentence level accuracies: $0.92^{14} = 31\%$ vs $0.97^{14} = 65\%$
- POS tagging not solved yet!



Some observations

- The function (or POS) of a word depends on its context
 - The/DT old/NN man/VB the/DT boat/NN
 - The/DT old/JJ man/NN bought/VBD the/DT boat/NN
- Certain POS combinations are extremely unlikely
 - <JJ, DT> ("good the") or <DT, IN> ("the in")
- Better to make decisions on entire sentences instead of individual words (Sequence modeling!)

Hidden Markov Models



- Model probabilities of **sequences** of variables
- Each state can take one of K values (can assume {1, 2, ..., K} for simplicity)
- Markov assumption: $P(s_t | s_{< t}) \approx P(s_t | s_{t-1})$

Markov chains



Where have we seen this before? Language models!



The/DT cat/NN sat/VBD on/IN the/DT mat/NN

Markov chains can help us model entire sentences.





BUT we don't normally see sequences of POS tags appearing in text.

Markov chains

The/?? cat/?? sat/?? on/?? the/?? mat/??

Hidden Markov Model (HMM)



- We don't normally see sequences of POS tags in text
- However, we do observe the words!
- The HMM allows us to jointly reason over both hidden and observed events.
 - Assume that each position has a tag that generates a word (Generative model)

The/?? cat/?? sat/?? on/?? the/?? mat/??

Components of an HMM



- 1. Set of states $S = \{1, 2, ..., K\}$ and set of observations O
- 2. Initial state probability distribution $\pi(s_1)$
- 3. Transition probabilities $P(s_{t+1} | s_t)$ (OR $\theta_{s_t \rightarrow s_{t+1}}$)
- 4. Emission probabilities $P(o_t | s_t)$ (OR $\phi_{s_t \rightarrow o_t}$)





Markov assumption: 1.

 $P(s_{t+1} | s_1, \dots, s_t) \approx P(s_{t+1} | s_t)$

Output independence: 2.

 $P(o_t | s_1, \ldots, s_t) \approx P(o_t | s_t)$

Assumptions



Depends on language! Which do you think is a stronger 1) assumes POS tag sequences do not have very strong priors/ A) Markov assumption long-range dependencies B) Output independence 2) assumes neighboring tags don't affect current word

Sequence likelihood



(joint likelihood of seeing both sequences)



Sequence likelihood



 $P(S, 0) = P(S_1, S_2 \cdots)$ = $T(S_1) P(0)$

$$S_{n}, \partial_{1}, \partial_{2} \dots \partial_{n}$$

$$|s_{1}\rangle = 2^{n} P(s_{1}, 0; |s_{1}\rangle)$$

Sequence likelihood



1=2

Transition

Emission

Example: Sequence likelihood



Dummy start state

J		<i>S</i> _{<i>t</i>+1}	<i>S</i> _{<i>t</i>+1}		
		DT	NN		
	Ø	0.8	0.2		
s _t	DT	0.2	0.8		
	NN	0.3	0.7		



 O_t

the	cat
0.9	0.1
0.5	0.5

What is the joint probability *P*(the cat, DT NN)?

A) (0.8 * 0.8) * (0.9 * 0.5)B) (0.2 * 0.8) * (0.9 * 0.5)C) (0.3 * 0.7) * (0.5 * 0.5)

Ans: A





Learning

Training set:

1 Pierre/NNP Vinken/NNP ,/, 61/CD years/NNS old/JJ ,/ join/VB the/DT board/NN as/IN a/DT nonexecutive/JJ di Nov./NNP 29/CD ./.

2 Mr./NNP Vinken/NNP is/VBZ chairman/NN of/IN Elsev N.V./NNP ,/, the/DT Dutch/NNP publishing/VBG group/ 3 Rudolph/NNP Agnew/NNP ,/, 55/CD years/NNS old/JJ chairman/NN of/IN Consolidated/NNP Gold/NNP Fields/N ,/, was/VBD named/VBN a/DT nonexecutive/JJ director/ this/DT British/JJ industrial/JJ conglomerate/NN ./.

38,219 It/PRP is/VBZ also/RB pulling/VBG 20/CD peopl of/IN Puerto/NNP Rico/NNP ,/, who/WP were/VBD help Huricane/NNP Hugo/NNP victims/NNS ,/, and/CC sendin them/PRP to/TO San/NNP Francisco/NNP instead/RB ./

Maximum likelihood estimate:

 $P(s_i | s_j) = \frac{Count(s_j, s_i)}{Count(s_j)}$

 $P(o \mid s) = \frac{Count(s, o)}{Count(s)}$



1. the/DT cat/NN sat/VBD on/IN the/DT mat/NN

2. Princeton/NNP is/VBZ in/IN New/NNP Jersey/NNP

3. the/DT old/NN man/VB the/DT boats/NNS

$$P(NN|DT) = \frac{3}{4}$$

$$P(cat \,|\, NN) = \frac{1}{3}$$

Example

• Maximum likelihood estimate:

 $P(s_i | s_j) = \frac{Count(s_j, s_i)}{Count(s_j)}$

 $P(o \mid s) = \frac{Count(s, o)}{Count(s)}$



Task: Find the most probable sequence of states (s_1, s_2, \ldots, s_n) given the observations (o_1, o_2, \ldots, o_n) $\hat{S} = ahgmax P(S|O) = ahgmax P(S) P(O|S)$ S = by S P(O) = by S P(O) S P(O)





Task: Find the most probable sequence of states (s_1, s_2, \ldots, s_n) given the observations (o_1, o_2, \ldots, o_n) $\hat{S} = ahgmax P(S|0) = ahgmax P(S) P(0|S)$ $S = begin{pmatrix} S = b$ = algmax P(S) P(0|S)





S = algmax P(S) P(O|S)= $argmax \prod P(S; |S_{i-1}) P(O; |S_{i})$ 121 How can we maximize this? Emission Transition Search over all state sequences?

Task: Find the most probable sequence of states $\langle s_1, s_2, \ldots, s_n \rangle$ given the observations $\langle o_1, o_2, \ldots, o_n \rangle$



Greedy decoding



Decode/reveal one state at a time

argmax
$$\pi(S = S) P(\pi | S)$$

S = ` DT'



Greedy decoding





argmax $P(s_2=s|DT)P(cat|s)$ s = 'NN'







• Not guaranteed to produce the overall optimal sequence

• Local decisions

Greedy decoding

max
$$P(S|S_t) P(O_{t+1}|S)$$

Viterbi decoding

- Use dynamic programming!
- Maintain some extra data structures
- Probability lattice, M[T, K] and backtracking matrix, B[T, K]
 - T: Number of time steps
 - *K* : Number of states
- B[i, j] is the tag at time i-1 in the most probable sequence ending with tag j at time i

• M[i, j] stores joint probability of most probable sequence of states ending with state j at time i

Viterbi decoding



 $M[1,DT] = \pi(DT) P(\text{the} | DT)$

 $M[1,NN] = \pi(NN) P(\text{the}|NN)$

Initialize the table

 $M[1, VBD] = \pi(VBD) P(\text{the} | VBD)$

 $M[1,IN] = \pi(IN) P(\text{the} | IN)$

Forward







 $M[2,NN] = \max M[1,k] P(NN|k) P(\operatorname{cat}|NN)$

 $M[2, VBD] = \max M[1, k] P(VBD | k) P(\operatorname{cat} | VBD)$

 $M[2,IN] = \max M[1,k] P(IN|k) P(\operatorname{cat}|IN)$



Viterbi decoding



k



What is the time complexity of this algorithm?



- A) O(n)
- *B) O*(*nK*)
- C) $O(nK^2)$
- D) $O(n^2K)$

n = number of timesteps K = number of states







If K (number of possible hidden states) is too large, Viterbi is too expensive!



• If K (number of states) is too large, Viterbi is too expensive!

Observation: Many paths have very low likelihood!

- Keep a fixed number of hypotheses at each point
 - Beam width, β

• If K (number of states) is too large, Viterbi is too expensive!

• Keep a fixed number of hypotheses at each point

 $\beta = 2$



Beam Search

• Keep a fixed number of hypotheses at each point

 $\beta = 2$



score = -16.5score = -6.5score = -3.0score = -22.1score = -0.5score = -13.5score = -32.0score = -20.3

Accumulated scores

Step 1: Expand all partial sequences in current beam

• Keep a fixed number of hypotheses at each point



Step 2: Prune set back to top β sequences (sort and select)

... and Repeat!

Keep a fixed number of hypotheses at each point

k



What is the time complexity of this algorithm?

- n = number of timesteps
- K = number of states
- = beam width B

- Keep a fixed number of hypotheses at each point
 - Beam width, β
- Trade-off (some) accuracy for computational savings

• If K (number of states) is too large, Viterbi is too expensive!

Beyond bigrams (Advanced)

- Trigram HMM: $P(s_{t+1} | s_1, s_2)$



Pros?

• Real-world HMM taggers have more relaxed assumptions

$$S_2, \ldots, S_t) \approx P(s_{t+1} | s_{t-1}, s_t)$$

https://forms.gle/byRfQJ5WsdYMYKCa6

Give us feedback!