CS626: Speech, Natural Language Processing and the Web

Part of Speech Tagging Pushpak Bhattacharyya Computer Science and Engineering Department IIT Bombay Week 2 of 1st August, 2022

What is NLP

NLP= Language + Computation

(due to ML)

= Linguistics + Probability



NLP Layers



Linguistic Strata vs. Languages

	Hindi	Swahili	Tamil	etc.
Sound: Phonetics, Phonology				
Structure: Morphology, Syntax				
Meaning: Semantic, Pragmatics				

Speech-NLP Stack vs. Languages

	Hindi	Swahili	Tamil	etc.
Sound: ASR, TTS				
Structure: MA, POS, NE, Chunker, Parser				
Meaning: SRL, Knowledge Nets, SA-EA-OM, QA, Summarizer				

76666466667.pug/200ak

Part of Speech Tagging

What is "Part of Speech"

- Words are divided into different kinds or classes, called Parts of Speech, according to their use; that is, according to the work they do in a sentence.
- The parts of speech are eight in number: 1. Noun. 2. Adjective. 3. Pronoun. 4. Verb. 5. Adverb. 6. Preposition. 7. Conjunction. 8. Interjection.





Understanding two methods of classification

- Syntactic
 - Important, e.g., for POS tagging
- Semantics
 - Important. E.g., for question answering
- Example: Adjectives
 - Syntactic: normal, comparative, superlative: good, better, best
 - Semantic: qualitative, quantitative: *tall man*, *forty horses*



Problem Statement

• Input: a sequence of words

Output: a sequence of labels of these words

Example

- Who is the prime minister of India?
 - Who- WP (who pronoun)
 - Is- VZ (auxiliary verb)
 - The- DT (determiner)
 - Prime-JJ (adjective)
 - Minister-NN (noun)
 - Of- IN (preposirtion)
 - India- NNP (proper noun)
 - ?- PUNC (punctuation)
- These POS tags are as per the Penn Treebank tagset

Motivation for POS tagging

Question Answering

Machine Translation

• Summarization

• Entailment

Information Extraction

Machine Translation

- "I bank on your moral support" to Hindi "main aapake naitik samarthan par nirbhar hUM"
- needs disambiguation of 'bank' (verb; POS tag VB) to 'nirbhar hUM'- as opposed to the noun 'bank' meaning the financial organization or the side of a water body.



Morphology

POS Annotation

 Who_WP is_VZ the_DT prime_JJ minister_NN of _IN India_NNP ?_PUNC

 Becomes the training data for ML based POS tagging 3 Generations of POS tagging techniques

- Rule Based POS Tagging
 - Rule based NLP is also called Model Driven NLP
- Statistical ML based POS Tagging (*Hidden Markov Model, Support Vector Machine*)
- Neural (Deep Learning) based POS Tagging

Necessity of POS Tagging (1/2)

• Command Center to Track Best Buses (Tol 30Jan21): POS ambiguity affects

 Elderly with young face increased covid 19 risk (Tol Oct 20)

Dependency Ambiguity



(it is reported by Maharastra Govt. that covid-19 cases have increased) root



(it is the Maharastra reports that have increased covid-19 cases!!!)

W :	^	Brown	foxes	jumped	over	the	fence	•
Т:	^	JJ	NNS	VBD	NN	DT	NN	•
		NN	VBS	JJ	IN		VB	
					JJ			
					RB			



A Brown

foxes jumped

over the

fence



Find the PATH with MAX Score.

What is the meaning of score?

Noisy Channel Model



Sequence *W* is transformed into sequence *T*



VV

HMM: Generative Model



This model is called Generative model. Here words are observed from tags as states. This is similar to HMM. 26562466pto157.pug/2100ak

Tag Set

Attach to each word a tag from
 Tag-Set

 Standard Tag-set : Penn Treebank (for English).

Penn POS TAG Set

1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential there
5.	FW	Foreign word
6.	IN	Preposition or subordinating conju
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun
19.	PRP\$	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative

A note on NN and NNS

 Why is there a different tag for plural-NNS?

- Answer: in English most nouns can act as verbs: I watched a *play_noun*; I *play_verb* cricket
- Also, the plural form of nouns coincide with 3rd person, singular number, present tense for verbs

Penn POS TAG Set (cntd)

22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	ТО	to
26.	UH	Interjection
27.	VB	Verb, base form
28.	VBD	Verb, past tense
29.	VBG	Verb, gerund or present participle
30.	VBN	Verb, past participle
31.	VBP	Verb, non-3rd person singular present
32.	VBZ	Verb, 3rd person singular present
33.	WDT	Wh-determiner
34.	WP	Wh-pronoun
35.	WP\$	Possessive wh-pronoun
36.	WRB	Wh-adverb





An Explanatory Example

Colored Ball choosing



Probability of transition to another Urn after picking a ball:

	U ₁	U ₂	U ₃
U_1	0.1	0.4	0.5
U_2	0.6	0.2	0.2
U ₃	0.3	0.4	0.3

Example (contd.)

G

0.5

0.4

0.1

B

0.2

0.5

0.3



Observation : RRGGBRGR

State Sequence : ??

Not so Easily Computable.

There are also initial probabilities of starting A particular urn: 3 probabilities

Diagrammatic representation (1/2)



Diagrammatic representation (2/2)



Classic problems with respect to HMM

 Given the observation sequence, find the possible state sequences- Viterbi
 Given the observation sequence, find its probability- forward/backward algorithm
 Given the observation sequence find the HMM prameters.- Baum-Welch algorithm

Illustration of Viterbi

- The "start" and "end" are important in a sequence.
- Subtrees get eliminated due to the Markov Assumption.

POS Tagset

- N(noun), V(verb), O(other) [simplified]
- ^ (start), . (end) [start & end states]

Illustration of Viterbi

<u>Lexicon</u>

people: N, V laugh: N, V

- •
- •
- •

Corpora for Training

$$\sum_{k=1}^{n} w_{11} t_{11} w_{12} t_{12} w_{13} t_{13} \dots w_{1k_{1}} t_{1k_{1}} t_{1k_{1}}$$

$$\sum_{k=1}^{n} w_{21} t_{21} w_{22} t_{22} w_{23} t_{23} \dots w_{2k_{2}} t_{2k_{2}} t_{2k_{2}}$$

$$\cdot$$

^ $w_{n1}_{n1} w_{n2}_{n2} t_{n2} w_{n3}_{n3} \dots w_{nk_{n}} t_{nk_{n}}$

Inference



Partial sequence graph

	^	Ν	V	0	•
٨	0	0.6	0.2	0.2	0
Ν	0	0.1	0.4	0.3	0.2
V	0	0.3	0.1	0.3	0.3
0	0	0.3	0.2	0.3	0.2
•	1	0	0	0	0

This transition table will change from language to language due to language divergences.

Lexical Probability Table

	E	people	laugh	•••	
٨	1	0	0	•••	0
N	0	1x10 ⁻³	1x10 ⁻⁵	•••	•••
V	0	1x10 ⁻⁶	1x10 ⁻³	•••	•••
0	0	0	1x10 ⁻⁹	•••	•••
•	1	0	0	0	0

Size of this table = # pos tags in tagset X vocabulary size

vocabulary size = # unique words in corpus

Inference

New Sentence: $^{\circ}$ people laugh . $\stackrel{\varepsilon}{\leftarrow}$ $\stackrel{N}{\longrightarrow}$ $\stackrel{N}{\longrightarrow}$ $\stackrel{N}{\longleftarrow}$ $\stackrel{V}{\longleftarrow}$ $\stackrel{V}{\longrightarrow}$ $\stackrel{V}{\rightarrow$

p(^ N N . | ^ people laugh .) = (0.6 x 0.1) x (0.1 x 1 x 10⁻³) x (0.2 x 1 x 10⁻⁵)

Computational Complexity

- If we have to get the probability of each sequence and then find maximum among them, we would run into exponential number of computations.
- If |s| = #states (tags + ^ + .) and |o| = length of sentence (words + ^ + .)

Then, #sequences = s^{|o|-2}

• But, a large number of partial computations can be reused using Dynamic Programming.

Dynamic Programming



Computational Complexity

- Retain only those N / V / O nodes which ends in the highest sequence probability.
- Now, complexity reduces from |s|^{|o|} to |s|.|o|
- Here, we followed the Markov assumption of order 1.

Back to urn example

Back to the Urn Example

- Here :
 - $\begin{array}{c|c} & S = \{U1, U2, U3\} \\ & V = \{R, G, B\} \end{array}$
- For observation: $- O = \{O_1 \dots O_n\}$
- And State sequence
 - $Q = \{q_1 \dots q_n\}$ B=
- **T** $\underset{\pi_i}{i} = P(q_1 = U_i)$

	U ₁	U ₂	U ₃
U_1	0.1	0.4	0.5
U ₂	0.6	0.2	0.2
U ₃	0.3	0.4	0.3
	R	G	В
U_1	R 0.3	G 0.5	B 0.2
U ₁ U ₂	R 0.3 0.1	G 0.5 0.4	B 0.2 0.5

Observations and states

- $S_i = U_1/U_2/U_3$; A particular state
- S: State sequence
- **O: Observation sequence**
- S* = "best" possible state (urn) sequence
- Goal: Maximize P(S*|O) by choosing "best" S



 Maximize P(S|O) where S is the State Sequence and O is the Observation Sequence

$$S^* = \arg\max_{S} (P(S \mid O))$$

False Start

 O_6 O_1 $O_2 O_3$ O_4 O_5 O_8 O_7 R G G B R G **OBS**: R R S₃ S_4 S_6 S_2 S₅ S_7 S₈ State: S₁

 $P(S \mid O) = P(S_{1-8} \mid O_{1-8})$ $P(S \mid O) = P(S_1 \mid O) \cdot P(S_2 \mid S_1, O) \cdot P(S_3 \mid S_{1-2}, O) \cdot \dots \cdot P(S_8 \mid S_{1-7}, O)$

By Markov Assumption (a state depends only on the previous state)

 $P(S \mid O) = P(S_1 \mid O).P(S_2 \mid S_1, O).P(S_3 \mid S_2, O)...P(S_8 \mid S_7, O)$

Baye's Theorem

$P(A \mid B) = P(A).P(B \mid A) / P(B)$

P(A) -: Prior P(B|A) -: Likelihood

 $\operatorname{arg\,max}_{S} P(S \mid O) = \operatorname{arg\,max}_{S} P(S) \cdot P(O \mid S)$

State Transitions Probability

$$P(S) = P(S_{1-8})$$

$$P(S) = P(S_1).P(S_2 | S_1).P(S_3 | S_{1-2}).P(S_4 | S_{1-3})...P(S_8 | S_{1-7})$$

By Markov Assumption (k=1)

 $P(S) = P(S_1).P(S_2 | S_1).P(S_3 | S_2).P(S_4 | S_3)...P(S_8 | S_7)$

Observation Sequence probability

 $P(O \mid S) = P(O_1 \mid S_{1-8}) \cdot P(O_2 \mid O_1, S_{1-8}) \cdot P(O_3 \mid O_{1-2}, S_{1-8}) \cdot \cdot \cdot P(O_8 \mid O_{1-7}, S_{1-8})$

Assumption that ball drawn depends only on the Urn chosen

 $P(O | S) = P(O_1 | S_1) . P(O_2 | S_2) . P(O_3 | S_3) ... P(O_8 | S_8)$

 $P(S \mid O) = P(S).P(O \mid S)$

 $P(S | O) = P(S_1).P(S_2 | S_1).P(S_3 | S_2).P(S_4 | S_3)...P(S_8 | S_7).$

 $P(O_1 | S_1).P(O_2 | S_2).P(O_3 | S_3)...P(O_8 | S_8)$

Grouping terms

O_0 (О ₁	O ₂	O_3	O_4	O_5	O_6	O ₇	O_8	
Obs: E	R	R	G	G	В	R	G	R	
State: S ₀ S	S ₁	S ₂	S ₃	S ₄	S ₅	S ₆	S ₇	S ₈	S_9

P(S).P(O|S)

 $= [P(O_0|S_0).P(S_1|S_0)].$ $[P(O_1|S_1). P(S_2|S_1)].$ $[P(O_2|S_2). P(S_3|S_2)].$ $[P(O_3|S_3).P(S_4|S_3)].$ $[P(O_4|S_4).P(S_5|S_4)].$ $[P(O_5|S_5).P(S_6|S_5)].$ $[P(O_6|S_6).P(S_7|S_6)].$ $[P(O_7|S_7).P(S_8|S_7)].$ $[P(O_8|S_8).P(S_9|S_8)].$ We introduce the states S_0 and S_9 as initial and final states respectively.

respectively. After S_8 the next state is S_9 with probability 1,

i.e., $P(S_{9}|S_{8})=1$

 O_0 is ε -transition

Introducing useful notation





Probabilistic FSM



The question here is:

"what is the most likely state sequence given the output sequence seen"

Developing the tree



Tree structure contd...



The problem being addressed by this tree is $S^* = \arg \max_{s} P(S \mid a_1 - a_2 - a_1 - a_2, \mu)$

a1-a2-a1-a2 is the output sequence and μ the model or the machine

POS tagging to be contd.