

a) What is stemming? Explain Porter's algorithm in detail.

Stemming

- It is a process of reducing a word to its word stem that affixes to suffixes and prefixes or to the roots of words known as lemma
- Stemming is a part of linguistic studies in morphology and artificial intelligence, information extraction. Stemming and AI knowledge extract meaningful info from vast source like big data or the internet :: additional forms of a word related to a subject may need to be searched to get best results.

Porter's Stemming Algorithm

Consonant: a letter other than A, E, I, O, U and preceded by consonant

Vowel - Only other letter.

e.g. - Troubles
 $\overline{T} \overline{r} \overline{o} \overline{u} \overline{b} \overline{l} \overline{e} \overline{s}$
 $\overline{C} \overline{V} \overline{C} \overline{V} \overline{C} \overline{V} \overline{C}$.

The rules are of the form:

(conditions) $\rightarrow s_1 \rightarrow s_2$

Step 1 - SSEs \rightarrow ss coreses \rightarrow cores

I ES \rightarrow I

panies \rightarrow pari

ss \rightarrow ss

coreses \rightarrow cores

s \rightarrow e

cats \rightarrow cat

Step 2a : $(m > 1) EED \rightarrow E$ agreed \rightarrow agree
 $(\cancel{*} V^*) ED \rightarrow E$ plastered \rightarrow plaster
 $(\cancel{*} V^*) I \cancel{DNG} \rightarrow E$ monitoring \rightarrow monitor

Step 2b : AT \rightarrow AIE

BL \rightarrow BLE

$(\cancel{*} d f ! (\cancel{*} L o r \cancel{*} S o r \cancel{*} z)) \rightarrow$ single letter

$(m = 1 \& \cancel{*} o) \rightarrow E$

Step 3 $(\cancel{*} V^*) r \rightarrow i$

Step 4 - $(m > 0) ATIONAL \rightarrow ATE$

$(m > 0) IZATION \rightarrow IZE$

$(m > 0) BILITE \rightarrow BLE$

Step 5 : $(m > 0) ICATE \rightarrow IC$

$(m > 0) FUL \rightarrow \leftarrow$

$(m > 0) NFSS \rightarrow \leftarrow$

Step 6 : $(m > 0) ANCE \rightarrow \leftarrow$

$(m > 0) ENT \rightarrow \leftarrow$

$(m > 0) NE \rightarrow \leftarrow$

Step 7 : $(m > 1) E \rightarrow G$

$(m = 1 \& \cancel{f} ! \cancel{*} o) NESS \rightarrow \leftarrow$

2) Explain ambiguity in NLP with example.

- Ambiguity can be referred as the ability of having more than one meaning or being understood in more than one way.
- Natural languages are ambiguous, so continuous computers are not able to understand language the way people do.
- NLP is concerned with the development of computational models of aspects of human language processing.
- Ambiguity could be lexical, syntactic, semantic, pragmatic etc.
eg - I saw a man on a hill with telescope
Alternate meaning

1. There is a man on a hill, and I'm watching him with my telescope.
- 2) There's a man and he's on a hill that also has a telescope on it.

or what is NLP? Describe various stages involved in NLP process, with suitable

NLP :

It is a subfield of linguistics, CS, info engg and AI concerned with the interaction b/w computers and human languages in particular how to program computers to process and analyze large amount of

Teacher's Sign.: _____

natural language data

Challenges in NLP: involve speech recognition, natural language understanding and natural language generation

Steps in NLP:

1. Lexical Analysis

It involves identifying and analyzing the structure of words. Lexicon of a language means the collection of words and phrase in a language.

Lexical analysis is dividing the whole chunk of text into paragraphs, sentences and words.

2. Syntactic Analysis

It involves analysis of words in the sentence for grammar and arranging words in a manner that shows relationships.

3) Semantic Analysis

It draws the exact meaning or the dictionary meaning of the text. The text is checked for meaningfulness. It is done by mapping syntactic structures and objects in the texts domain. The semantic analyzer disregards sentence such as "hot" ice-cream

4. Discourse Integration:

The meaning of any sentence depends upon the meaning of sentence just before it. In addition, it also brings about the meaning of immediately succeeding sentence.

5. Pragmatic Analysis:

During this, what was said is re-interpreted on what it actually meant. It involves deriving those aspects of language which require real world knowledge.

Lexical Analysis



Syntactic Analysis



Semantic Analysis



Discourse Integration



Pragmatic Analysis

Q3

An Finite state Transducer

→ FS Automata

- An FST represents a set of strings
e.g. walk, walks, walked
- A recognizer function
 $\text{recognize}(\text{str}) \rightarrow \text{true or false}$

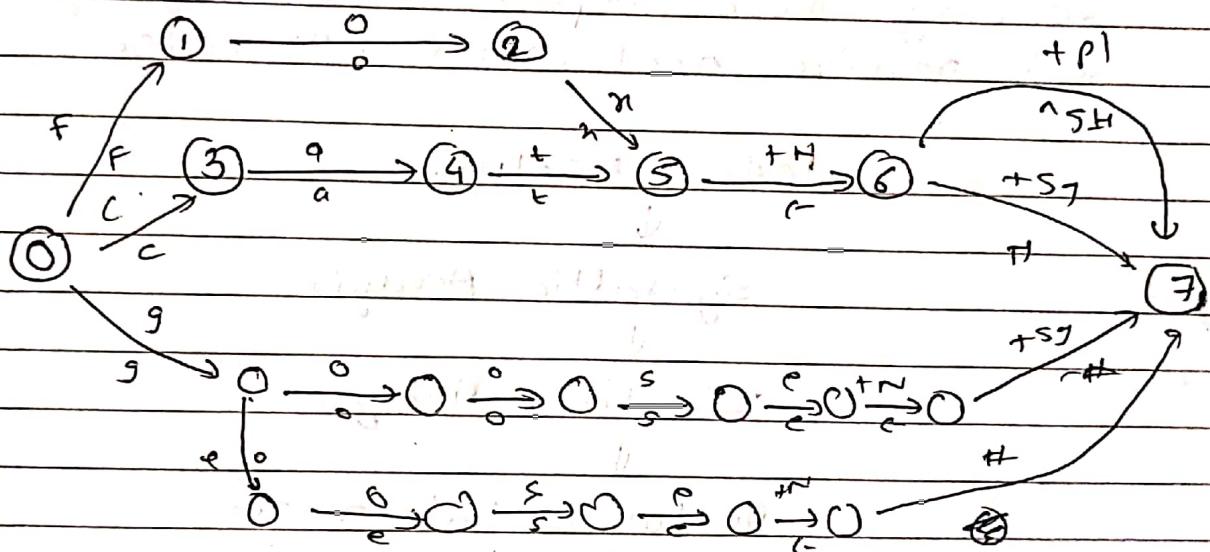
Teacher's Sign.: _____

- FST transducer.

> An FST represents a set of pairs of strings
(think of as ip, op, pairs)

can return multiple answers if ambiguous
eg if you don't have pos-tagged input

"walk" could be the verb "They walk to
the store versus the noun "I took a
walk."



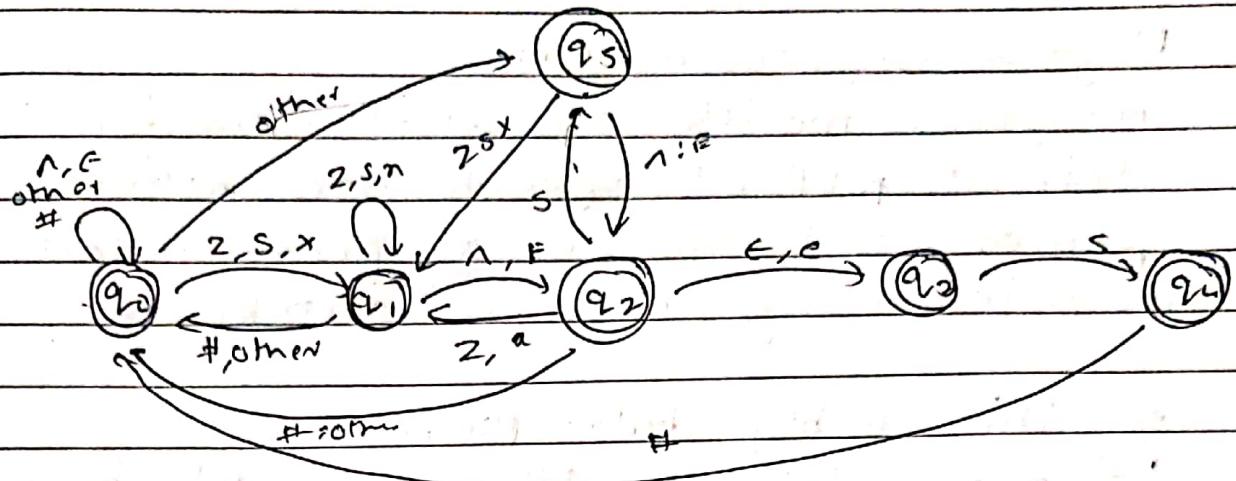
lexicon

|f|o| x |t| +n| +PI| | |

Intermediate

|f|o| x |n| s| #| |

Teacher's Sign.: _____



lexical

| f | o | x | i n | + p |

Intermediate

| f | o | x | i n | s | # |

Surface

| | f | o | x | E | s | T |

Q4 what is language models

write a note on program language model

language modelling is one of the most important parts of modern NLP. There are many sorts of application for language modelling like Machine translation, spell correction, speech recognition, summarization, Question Answering, sentiment analysis etc. Each of those tasks require use of language model.

Teacher's Sign.: _____

N -gram language model:

N -gram models are widely used in statistical NLP. In speech recognition, phonemes and sequences of phonemes are modelled using n -gram distribution.

e.g. This is big data AI book.

Unigram : This, is, big, data, AI, book = 6

Bigram : This is, is big, big data, data AI.

Trigram : This is a big, is big data, big data AI book = 4

The chain rule probability is:

$$P(w_1, \dots, w_n) = P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_1, w_2) \cdots P(w_n | w_1, w_2, \dots, w_{n-1})$$

The maximum likelihood estimate

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

e.g. <S> I am Sam </S>

<S> Sam I am </S>

<S> I don't like green eggs and ham </S>

$$P(I | <S>) = \frac{2}{3}$$

$$P(\text{Sam} | \text{am}) = \frac{1}{2}$$

$$P(\text{Sam} | <S>) = 1/3$$

$$P(\text{do} | I) = 1/3$$

$$P(\text{am} | I) = 2/3$$

$$P(</S> | \text{Sam}) = 1/2$$

Teacher's Sign.: _____

Q Explain various approaches to perform Part of speech (POS).

→ Part of ~~sp~~ speech tagging may be defined as the process of assigning one of the PoS to given word. In a simple words, POS tagging is a task of labelling each word in a sentence with its appropriate PoS.

Most of PoS tagging falls under Rule Based, stochastic and ~~and~~ transformation based tagging.

* Rule Based PoS tagging

One of the oldest techniques of tagging is rule-based PoS tagging. Rule-based tagger use dictionary or lexicon for getting possible tag for tagging each word. If the word has more than one possible tag, then rule-based tagger uses hand-written rules to identify the correct tag.

~~disambiguation~~ Disambiguation can also be performed in this by analyzing the linguistic feature of a word along with its preceding or next few words.

Rules may be either

- context : Pattern Rule
- Regex into FSA

Teacher's Sign.: _____

Stochastic Pos Tagging

- The model that includes frequency or probability can be called stochastic.

The simplest stochastic trigger applies the full approaches for Pos tagging.

Word frequency approach :

In this the stochastic tagger disambiguates the words based on the prob that word occurs with a particular tag.

The main issue is that it may yield inadmissible sequence of tags.

Tag sequence probability.

It is another approach where tagger calculates the prob of a given sequence of tags occurring. It is also called n-gram approach.

* Transformation Based Tagging:

It is also called Brill tagging. It is an instance of transition Based Learning which is a rule based algo for automatic tagging of pos. TBL allows to have linguistic knowledge in a readable form, transform one state to another using transition rules. If we see similarities b/w rule based and translation tagger then like rule-based, it is also based on rules.

Teacher's Sign.: _____

which specify what tags to be needed to assign specific words. It is machine learning technique in which rules are automatically included from data.

Teacher's Sign.: _____