## Module 2: Word Level Analysis

**Q1** Write a short note on morphological parsing and Morphology Analysis.

=> Morphological parsing is the task of recognizing the morphemes inside a word.

- Morphemes are the minimal meaning - bearing unit in a language.

- Example: Mangoes

Mango     es  ⟶ ①     Here there are two Morphemes

= Morphemes can be Stems (Root word) or an Affix

- Now this Affix is divided into three parts An affix can be prefix (eg reform) or Suffix (eg loved) or infix (passersby).

- So here Mango is a Stem and es is a suffix because it is attached after the main word

= Following are the requirments for building a morphological Parser

1. Lexicon: It includes the list of stems and affixes along with the basic information about them.
      eg Stem is a noun Stem or a verb Stem.

Morphotactics:

## Morphotactics:

Morphotactics has a set of rules by the help of which it decides the ordering of words.

example: Use able ness

Useableness ✓

able use ness

ableuseness ✗

## Orthographic rules:

These spelling rules are used to model the changes occuring in a word.

example: lady + s = ladys ✗

lady + s = ladies ✓

- The study of formation of words is called morphology.

- Some words are self sufficient they have there own meaning example - camera, pen

1. Some words if devided has there own meanings example: Showcase = Show & case both words have there own meanings.

- Some words if combine dont have any meaning but if they are combined with a word it becomes a meaningful word.

example: ing has no meaning but if combined with love its loving which have there meaning.

- So basically there are different words which if used in a right way we can get a meaningful word.

- So analysis is Studying in detail and Analysis of Morphology is Morphology analysis.

**Q2** Write a short note on Inflectional and Derivational morphology.

=> Before ⎯ Inflectional and derivational morphology we need to understand what are morphemes. Morpheme is a word or a part of a word that has meaning and a morpheme cannot be devided further into meaningfull units. example of morpheme : cat

If we try to devide morpheme more it will be a meaningless result.

There are two types of morphemes
① Free Morphemes ② Bound morphemes

① Free Morpheme : Free morpheme is a morpheme which has its own meaning or it has its complete meaning example : fan, camera etc.

― Free morphemes are of two types lexical morphemes and grametical morphemes

• Lexical morphemes are the picture words they are noun, adjective, verbs, adverbs example : black, yello, chair
Every year new and new lexical morphemes are added in a language.

― Gramatical morphemes are grammar words which are limited in each and every language.

These morphemes dont change frequently like Lexical morphemes they are preposition, conjunctions, etc

② Bound Morphemes: These morphemes are of two types Inflectional and Derivational. But before going further we must know what is bound morphemes.

Bound morphemes are those morphemes whose meaning is not complete in themself. And that is the reason why they depend on the free morphemes for meaning

Now Affeeres are bound morphemes and Affieres are of three types prefix, Postfix, suffix

prefix →     Because

Soffix → loveable

Infix → passersby

Inflectional morpheme: Inflectional morpheme is one which when attached to a root word dosent change its class

$$\underset{Noun}{\underline{Book}} + S = \underset{Noun}{\underline{Books}}$$

Inflectional morphemes are Infixes and Suffixes and cant Be prefix

## Derivational Morphemes:

Derivational Morphemes are one which when added to a word changes its class.

$$\underline{Teach} + \underline{er} = \underline{Teacher}$$
$$\text{Verb} \qquad\qquad \text{Noun}$$

Derivational morphemes are of two types class changing which was the above one and class maintaining which when added to word changes the word but cant change the class.

$$\underline{Child} + \underline{hood} = \underline{childhood}$$

Common noun          abstract noun but the class is same
Which is Noun.

**Q3** Design a finite state automata (FSA) for boat!

=>. Q : finite set of states
$q_0, q_1, q_2, q_3, q_4$

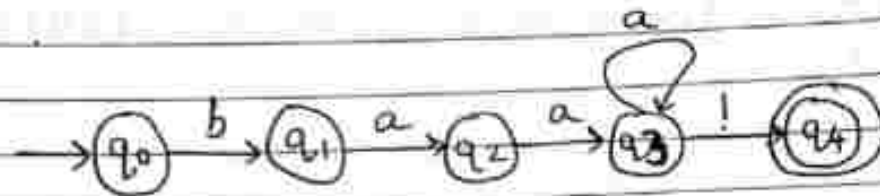$\Sigma$ : set of input alphabets
: $\{a, b, !\}$

$q_0$ : Start State.

F : Set of final states
: $F \subseteq Q$

$\delta(q, i)$ : the transition function $\overset{or}{,}$ transition matrix between states. Given a state $q \in Q$ and input symbol $i \in \Sigma$, $\delta(q, i)$ returns a new state $q' \in Q$. $\delta$ is thus a relation from $Q \times \Sigma$ to $Q$;

Transition table

| Input<br>States | a | b | ! |
|---|---|---|---|
| $q_0$ | $\phi$ | $q_1$ | $\phi$ |
| $q_1$ | $q_2$ | $\phi$ | $\phi$ |
| $q_2$ | $q_3$ | $\phi$ | $\phi$ |
| $q_3$ | $q_3$ | $\phi$ | $q_4$ |
| $q_4$ | $\phi$ | $\phi$ | $\phi$ |

**Q4** Design a finite State Automata for divisibility by 5 teeter for binary number.

Q => finite set of state

q$ => Start state ('Might or Might not include yours)

$q_0$ => rem 0 State

$q_1$ => rem 1 State

$q_2$ => rem 2 State

$q_3$ => rem 3 State

$q_4$ => rem 4 State

rem = reminder.

The question can also be written as design a FSA to check wether given binary no is divisible by 5 or not

$\xi$ => Set of input alphabets = {0,1}

$\delta$ = Transition function

$\delta : Q \times \xi \to Q$

$q_0$ = Start state = q$

F = final state

$q_0, F \subseteq Q$

[or you can write like previous Numerical]

Scanned with OKEN Scanner

| $\delta$ | 0 | 1 |
|---|---|---|
| $q_5$ | $q_0$ | $q_1$ |
| * $q_0$ | $q_0$ | $q_1$ |
| $q_1$ | $q_2$ | $q_3$ |
| $q_2$ | $q_4$ | $q_0$ |
| $q_3$ | $q_1$ | $q_2$ |
| $q_4$ | $q_3$ | $q_4$ |



**Q5-** Design a DFA of a string that should end with 100

=> $M = \{ Q, \Sigma, \delta, q_0, F \}$

$q_0$ = initial state
$q_1$ = String ending with 1
$q_2$ = String ending with 10
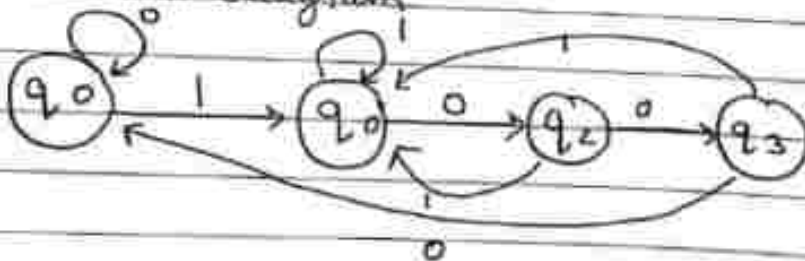$q_3$ = String ending with 100

$\Sigma = \{ 0, 1 \}$
$q_0 = \{ q_0 \}$
F = final state = $q_3$

## Transition Table

|     | 0   | 1   |
| --- | --- | --- |
| $q_0$ | $q_0$ | $q_1$ |
| $q_1$ | $q_2$ | $q_1$ |
| $q_2$ | $q_3$ | $q_1$ |
| $q_3$ | $q_0$ | $q_1$ |

## Transition diagram

**Q6** Differentiate between Inflectional & Derivational morphology.

| | Inflectional Morphology | Derivational Morphology |
|---|---|---|
| 1. | It is a morphological process that adapts existing words so that they function efficisively in sentences without changing PoS of base morpheme m | It is concerned with the way morphemes are connected to existing lexical forms as affixes. |
| 2. | Regular: It is more Regular | It is very less regular |
| 3. | Use: Can only be Suffix or infix and not prefix | Can be both prefix & Suffix |
| 4. | Change in : Never change Part of the gramatical Speech category or PoS | It can change the gramatical category or PoS |
| 5. | Example: $\underset{Noun}{Cat} + S = \underset{Noun}{cats}$ | $\underset{Noun}{danger} + ous = \underset{Adjective}{dangerous}$ |

**Q7** Write a short note on language model

**Q1. Write a short note on language model**

=> The goal of a language model is to assign probability to a sentence.

With this is also decides which sentence is more accurate at the moment.

example: She is a tall girl is more accurate than she is a long girl

Statistical language Modelling, or language Modelling is the development of probabilistic models that are able to predict the next word in the sequence given the words that ~~predict~~ precede it. It is a probability distribution over sequences of words.

Given such a sequence, say of length $m$, it assigns a probability $P(w_1, \ldots, w_m)$ to the whole sequence.

The goal of probabilistic language modelling is to calculate the probability of a sentence of sequence of words: $P(w) = P(w_1, w_2, w_3 \ldots w_n)$ and can be used to find the probability of the next word in the sequence:
$P(w_5 / w_1, w_2, w_3, w_4)$
A model that computes either of these is called language model

Method of calculating Probability:
Conditional probability:
let A and B be two events with $P(B) \neq 0$, the
conditional probability of A given B is:

$$P(A/B) = \frac{P(A,B)}{P(B)}$$

$$P(x_1, x_2, \ldots x_n) = P(x_1) P(x_2/x_1) \ldots P(x_n/x_1 \ldots x_{n-1})$$

for example: P("its water is so transparent") =
$\quad$ P(its) * P(water/its) * P(is/its water) *
P(so/its water is) * P(transparent/its water is so)

we can estimate this by simply counting and
dividing the results.

$$P(\text{transparent / its water is so}) = \frac{\text{count (its water is so transparent)}}{\text{count (its water is so)}}$$

Markov Property : It says that the probability
of the next word can be estimated given
only the previous K number of words.

for example, if $k=1$:
P (transparent / its water is so) $\approx$ P( transparent /
so )
or if $K=2$:
P(transparent / its water is so) $\approx$ P( transparent /
is so)

general equation for the Markov Assumption, $K = i$:

$$P(w_i | w_1, w_2 \ldots w_{i-1}) \approx P(w_i | w_{i-K} \ldots w_{i-1})$$

**Q8** Write a short note on N-Gram model.

=> Note: This answer is in continuation with the previous answer. You have to decide the length depending the marks alloted to Question.

- The Simplest case of markov model is a unigram model, In this model we we simply estimate the probability of the whole sequence of words by the product of probabilities of individual words - unigrams. and if we generated sentences by randomly picking words,           It would be

    Sixth, the, rupees, abduction

    It would be just a random sequence of words

$$P(w_1, w_2 \ldots w_n) \approx \prod_i P(w_i)$$

- Slightly more intelligent is the bigram model where we condition on the single previous word.

$$P(w_i | w_1, w_2 \ldots w_{i-1}) \approx P(w_i | w_{i-1})$$

We can extend this to trigrams, 4-grams, 5-grams.

But in general this is an insufficient model of language.

○ because language has long distance dependencies

example:

"The computer which I had just put into the machine room on the fifth floor crashed."

And if we say predict the next word after floor so it is very unlikely to predict crash but if we compare with or bring the main subject "computer" in the picture then we are more likely to guess crashed or predict crash as a next word.

Q9 corpus:
  &lt;s&gt; I am a human &lt;/s&gt;
  &lt;s&gt; I am not a stone &lt;/s&gt;
  &lt;s&gt; I I live in Mumbai &lt;/s&gt;
  check the probability of &lt;s&gt; I I am not &lt;/s&gt;
  using bigram

⇒ $P(\,I\,I\,am\,not\,)$
= $P(I\,|\,\langle s\rangle)\,P(I\,|\,I)\,P(am\,|\,I)\,P(not\,|\,am)$
  $P(\langle /s\rangle\,|\,not\,)$

= $\dfrac{C(\langle s\rangle\,|\,I)}{C(\langle s\rangle)}\ \dfrac{C(I\,|\,I)}{C(I)}\ \dfrac{C(I\,|\,am)}{I(I)}\ \dfrac{C(am\,|\,not)}{C(am)}$

  $\dfrac{C(not\,|\,\langle /s\rangle)}{C(not)}$

= $\dfrac{3}{3}\times\dfrac{1}{4}\times\dfrac{2}{4}\times\dfrac{1}{2}\times\dfrac{0}{1}$

= $0$

**Q10** consider following training data.

$\langle S \rangle$ I am Jack $\langle /S \rangle$

$\langle S \rangle$ Jack I am $\langle /S \rangle$

$\langle S \rangle$ Jack I like $\langle /S \rangle$

$\langle S \rangle$ Jack I do like $\langle /S \rangle$

$\langle S \rangle$ do I like Jack $\langle /S \rangle$

Assume that we use a biagram language model based on above data.

what is most probable next word predicted by the model?

1. $\langle S \rangle$ Jack....     2. $\langle S \rangle$ Jack I do....

3. $\langle S \rangle$ Jack I am Jack.....

4. $\langle S \rangle$ do I like....

$\Rightarrow$ 

$P(I|\langle s \rangle) = C(\langle s \rangle|I)/C(\langle s \rangle) = 1/5$

$P(Jack|\langle s \rangle) = (\langle s \rangle|Jack)/C(\langle s \rangle) = 3/5$

$P(do|\langle s \rangle) = C(\langle s \rangle|do)/C(\langle s \rangle) = 1/5$

$P(am|I) = C(I|am)/c(I) = 2/5$

$P(like|I) = C(I|like)/C(I) = 2/5$

$P(do|I) = C(I|do)/C(I) = 1/5$

$P(\langle /S \rangle|Jack) = C(Jack|\langle /S \rangle)/C(Jack) = 2/5$

$P(\langle /S \rangle|like) = C(like|\langle /S \rangle)/C(like) = 2/3$

$P(\langle /S \rangle|am) = C(am|\langle /S \rangle)/C(am) = 1/2$

$P(I|Jack) = C(Jack|I)/c(Jack) = 3/5$

$P(like|do) = C(do|like)/c(do) = 1/2$

$P(I|do) = C(do|I)/C(do) = 1/2$

$P(Jack|like) = C(like|Jack)/c(like) = 1/3$

$P(Jack|am) = C(am|Jack)/c(am) = 1/2$

1. \<S\> Jack I    2. \<S\> Jack I do like A I

3. \<S\> Jack I am Jack I

4. \<S\> do g like \</S\>