

## Internal Assessment 1 solution

|       |                             |           |          |            |           |            |         |  |
|-------|-----------------------------|-----------|----------|------------|-----------|------------|---------|--|
| Sub:  | Natural Language Processing |           |          |            | Sub Code: | 15CS741    | Branch: |  |
| Date: | 20.09.18                    | Duration: | 90 min's | Max Marks: | 50        | Sem / Sec: | VII/C   |  |

Answer any FIVE FULL Questions

MARKS  
[04]

1 (a) What is NLP? Explain the challenges and applications of NLP.

⇒ NLP is concerned with development of computational models of aspects of human lang. processing. 2 main reasons for development:

- i) Develop automated tools for language processing.
- ii) Gain a better understanding of human comm.

Challenges of NLP

- No. of factors that make NLP difficult.
  - Relate to problems of representation + interpretation
  - Lang. computing ⇒ precise representation of combi-N.L. ⇒ highly ambiguous, vague ⇒ achieving such representation can be difficult.
- ⇒ Inability to capture all the required knowledge impossible to<sup>||</sup> embody all sources of knowledge that humans use.  
not possible<sup>||</sup> to write procedure that imitate lang. as done by humans.

⇒ greatest difficulty → identifying its semantics.  
⇒ Principle of Compositional Semantics  
considers the meaning<sup>||</sup> of a sentence to be a composition of the meaning of words.

⇒ words alone do not make a sentence.

eg: I do not like icecream  
do not I icecream like. x

(or)

Kabir and Ayan are Married } both give different  
Kabir + Suba " " } meaning even if  
has same structure.

- ⇒ Only words along with their syntactic + semantic relation give meaning to a sentence.
- ⇒ Inability to capture all the required knowledge impossible to<sup>||</sup> embody all sources of knowledge that humans use.  
not possible<sup>||</sup> to write procedure that imitate lang. as done by humans.

⇒ greatest difficulty → identifying its semantics.  
⇒ Principle of Compositional Semantics  
considers the meaning<sup>||</sup> of a sentence to be a composition of the meaning of words.

⇒ words alone do not make a sentence.

eg: I do not like icecream  
do not I icecream like. x

(or)

Kabir and Ayan are Married } both give different  
Kabir + Suba " " } meaning even if  
has same structure.

- ⇒ Only words along with their syntactic + semantic relation give meaning to a sentence.

- ① Idioms, Metaphor + ellipses  $\Rightarrow$  more complexity  
to identify the meaning of the written text.  
 $\Rightarrow$  eg. oldman finally kicked the bucket.
- ② Quantifier Scoping - scope of quantifier is not clear
- ③ Ambiguity of natural languages.
- ④ Incorporating Contextual + world language knowledge poses the greatest difficulty.

NLP applications:

1. Machine Translation
2. Speech Recognition
3. Speech synthesis
4. Natural language interfaces to databases
5. Information retrieval
6. Information extraction
7. Question Answering.

(b) Explain the stages of transformational grammar for the sentence "The boy hit the girl".

[06]

$\Rightarrow$  transformational grammar has 3 components

- ① Phrase structure
- ② Transformational rules.
- ③ Morphophonemic rules.  $\rightarrow$  match each sentence represent to a string of phonemes.

S  $\rightarrow$  NP + VP  
 VP  $\rightarrow$  V + NP  
 NP  $\rightarrow$  Det + Noun  
 V  $\rightarrow$  Aux + Verb  
 Det  $\rightarrow$  the, a, an, ...  
 Verb  $\rightarrow$  Catch, write, eat, ...  
 Noun  $\rightarrow$  police, snatcher  
 Aux  $\rightarrow$  will, is, can, ...

S  $\Rightarrow$  Sentence  
 NP  $\Rightarrow$  Noun phrase  
 VP  $\Rightarrow$  Verb phrase  
 Det  $\Rightarrow$  Determiner

2 (a) Explain various grammar based language models.

[03]

1. Generative grammar
2. Hierarchical grammar
3. Government and binding
4. Lexical Functional grammar
5. Paninian framework

(b) Explain the types of statistical language model with an example.

[07]

Statistical Language model

- A statistical language model is a probability distribution  $P(S)$  over all possible words sequences.

(i) n-gram model:

- The goal is to estimate the probability of sentence.

- By decomposing sentence probability into a product of conditional probabilities using the chain rule as follows:

$$P(S) = P(w_1, w_2, w_3, \dots, w_n)$$

$$= P(w_1) P(w_2/w_1) P(w_3/w_1, w_2) P(w_4/w_1, w_2, w_3) \dots P(w_n/w_1, w_2, w_3, \dots, w_{n-1})$$

$$= \prod_{i=1}^n P(w_i/h_i)$$

$h_i$  is history of word  $w_i$  defined as  
 $w_1, w_2, \dots, w_{i-1}$

An n-gram model simplifies the task by approximating the probability of a word  $w_i$  on all the previous words by the conditional probability on previous  $n-1$  words only.

$$P(w_i/h_i) \approx P(w_i/w_{i-n+1}, \dots, w_{i-1})$$

Thus, an ngram model calculates  $P(w_i/h_i)$  by modelling language as Markov model of order  $n-1$  by looking at previous  $n-1$  words only.

A model that limits the history to the previous one word only, is termed a bi-gram ( $n=2$ ) model.



using bi-gram and tri-gram estimate, the probability of sentence can be calculated as,

$$P(s) \approx \prod_{i=1}^n P(w_i | w_{i-1}) \text{ and}$$

$$P(s) \approx \prod_{i=1}^n P(w_i | w_{i-2} \cdot w_{i-1})$$

eg: the bi-gram approximation of  $P(\text{east} | \text{The Arabian Knights are fairy tales of the})$  is

$$P(\text{east} | \text{the})$$

whereas a tri-gram approximation is

$$P(\text{east} | \text{of the})$$

Count a particular n-gram in the training corpus + divide it by the sum of all n-grams that share the same prefix.

$$P(w_i | w_{i-n+1} \dots w_{i-1}) = \frac{c(w_{i-n+1} \dots w_{i-1}, w_i)}{\sum_w c(w_{i-n+1} \dots w_{i-1}, w)}$$

The sum of all n-grams that share first n-1 words is equal to the count of the common prefix  $w_{i-n+1} \dots w_{i-1}$

$$P(w_i | w_{i-n+1} \dots w_{i-1}) = \frac{c(w_{i-n+1} \dots w_{i-1}, w_i)}{c(w_{i-n+1} \dots w_{i-1})}$$

The model parameter we get using these estimates, maximizes the probability of the training set  $T$  on the model  $M$ . i.e.  $P(T|M)$ .

eg:

training set -

The Arabian Knights

These are the fairy tales of the east

The stories of the Arabian Knights are translated in many languages.

Bi-gram model:

$$P(\text{the} | \langle s, \rangle) = 0.67 \quad P(\text{Arabian} | \text{the}) = 0.4$$

$$P(\text{Knights} | \text{Arabian}) = 1.0$$

$$P(\text{are} | \text{there}) = 1.0 \quad P(\text{the} | \text{are}) = 0.5$$

$$P(\text{fairy} | \text{the}) = 0.2$$

$$P(\text{tales} | \text{fairy}) = 1.0 \quad P(\text{of} | \text{tales}) = 1.0$$

$$P(\text{the} | \text{of}) = 1.0$$

$$P(\text{east} | \text{the}) = 0.2 \quad P(\text{stories} | \text{the}) = 0.2$$

$$P(\text{of} | \text{stories}) = 1.0$$

$$P(\text{are} | \text{Knights}) = 1.0 \quad P(\text{translated} | \text{are}) = 0.5$$

$$P(\text{in} | \text{translated}) = 1.0$$

$$P(\text{many} | \text{in}) = 1.0$$

$$P(\text{languages} | \text{many}) = 1.0$$

Test Sentence (S): The Arabian Knights are the  
fairy tales of the East.

$$P(\text{The} | \langle s \rangle) \times P(\text{Arabian} | \text{the}) \times P(\text{Knights} | \text{Arabian}) \times P(\text{are} | \text{Knights}) \times P(\text{the} | \text{are}) \times P(\text{fairy} | \text{the}) \times P(\text{tales} | \text{fairy}) \times P(\text{of} | \text{tales}) \times P(\text{the} | \text{of}) \times P(\text{East} | \text{the})$$

$$= 0.67 \times 0.5 \times 1.0 \times 1.0 \times 0.5 \times 0.2 \times 1.0 \times 1.0 \times 1.0 \times 0.2$$

$$= 0.0067$$

### Adds-one Smoothing

It adds a value of one to each n-gram frequency before normalizing them into probabilities. Thus the conditional probability becomes:

$$P(w_r | w_{r-n+1} \dots w_{r-1}) = \frac{c(w_{r-n+1} \dots w_{r-1}, w_r)}{c(w_{r-n+1} \dots w_{r-1}) + V}$$

✓ if the vocabulary size (V) size of the set of all the words being considered.

### Good-Turing Smoothing

- adjusts the frequency  $f$  of an n-gram using the count of n-grams having a frequency of occurrence  $f+1$

$$f^* = \frac{(f+1) n_{f+1}}{n_f}$$

$n_f$  - no. of n-gram that occur exactly  $f$  times in the training corpus.

eg: occur 4 times is 25,108 + The no. of n-gram that occur 5 times is 20,542, smoothed count for 5 will be

$$\frac{20,542 \times 5}{25,108} = 4.09$$

### Caching Technique

- The frequency of n-gram is not uniform across the text segments or corpus.

- Certain words occur more frequently in certain segments and rarely in other. eg: the frequency of the word 'n-gram' is high, whereas it occurs rarely.

### Caching Technique

- The frequency of n-gram is not uniform across the text segments or corpus.

- Certain words occur more frequently in certain segments and rarely in other. eg: the frequency of the word 'n-gram' is high, whereas it occurs rarely. The basic n-gram model ignores this sort of variation of n-gram frequency.

The cache model combines the most recent n-gram frequency with the std n-gram model to improve its performance locally.

3. Explain briefly about Government and binding grammar.

[10]

## 10) Government and binding

Chomsky - 1985

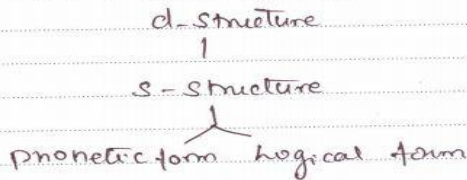
The structure of a language is understood at the level of its meaning. The sentences are gr. at the syntactical level and the transformation from meaning to syntax or vice versa is not well understood.

Transformational grammar has two level of existence of sentences

- (i) Surface level
- (ii) Deep root level

Government and binding theories renamed as:

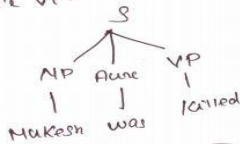
- (i) S-level
- (ii) d-level and identified 2 more levels.
- (iii) phonetic form
- (iv) logical form



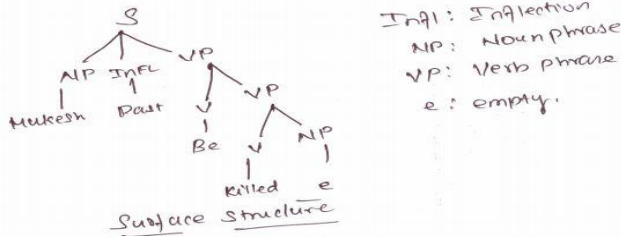
### Different levels of representation in GB.

Ex: Makesh was killed

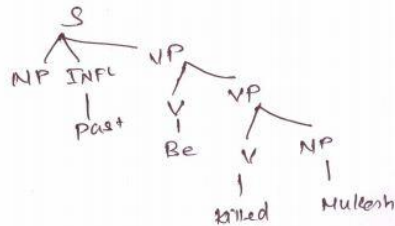
i) In transformational grammar, this can be represented as S-NP Aux VP.



ii) In GB, the S-structure and d-structure:



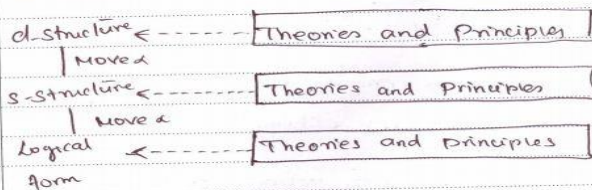
INFL: Inflection  
NP: Noun phrase  
VP: Verb phrase  
e: empty.



### Deep Structure

#### Components of GB:

- Comprises a set of theories that map the structures from d-structure to s-structure and to logical form.
- Transformational rule: 'Move  $\alpha$ ' is applied to d-structure level and s-structure level.



#### Components of GB

GB  $\Rightarrow$   $\alpha$  Series of modules that contain constraints and principles applied at various levels of representation and transformation.



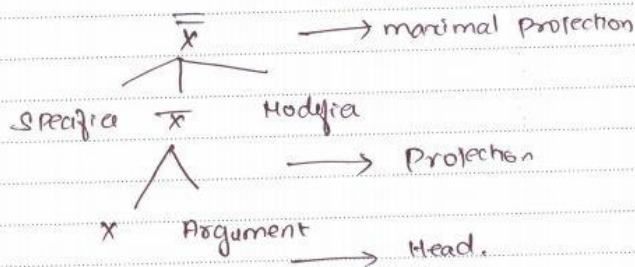
## General characteristics of GB:

- \* X-bar theory
- \* Projection principle
- \*  $\theta$ -theory
- \*  $\theta$ -criterion
- \* C-Command and government
- \* Case Study or theory
- \* Empty category principle (ECP)
- \* Binding theory

## X-Theory:

- Central concept in GB
  - X theory defines both phrase structure and the sentence structure with separate set of rules as maximal projections.
  - Entities defined become language independent
- Thus, noun phrase, verb phrase, adjective phrase, prepositional phrase are maximal projections of nouns, verb, adjective and preposition.

$$X = \{N, V, A, P\}$$



## Sub-categorization:

It is to be noted that GB does not consider traditional phrase structure as an appropriate device for defining language constructs.

It places the burden of ascertaining well formedness to sub-categorization frames of heads.

Any maximal projection can be argument of a head, but sub-categorization is used as a filter to permit various heads to select a certain subset of the range of maximal projection.

eg: we know that the verb 'eat' and sub-categorize for NP, whereas the verb 'sleep' cannot. Hence 'ate food' is well formed but the sentence 'slept the bed' not.

## Projection principle

- place a constraint on the three syntactic representations and their mapping from one to the other.

- The principle states that representations at all syntactic levels (rel d-level, S-level, LF-level) are projections from the lexicon.

## Theta Theory ( $\theta$ -theory) or The theory of thematic relations:

Sub-categorization only places a restriction on syntactic categories which a head can accept. GB

Put another restriction on the lexical heads through which it assign certain roles to its arguments. These roles are preassigned and cannot be violated at any syntactical level as per the Projection principle. These role assignments are called theta roles and are related to 'semantic-selection'.

### Theta-role and theta-criterion:

There are certain thematic roles from which a head can select. These are called  $\theta$ -roles and they are mentioned in the lexicon, eg: The Verb 'eat' can take arguments with  $\theta$ -roles. (Agent - theme)

### C-command and government:

As 'government' is a special case of 'C-command' we will first define C-command.

C-command: defines the scope of maximal projection. It's a basic mechanism through which many constraints are defined on Move  $\alpha$ .

If any word or phrase falls within the scope and is dominated by a maximal projection. Now if there are two structures  $\alpha$  and  $\beta$  related in such a way that 'every maximal projection dominating  $\alpha$  dominates  $\beta$ '.

### Government:

- $\alpha$  governs  $\beta$  iff
- $\alpha$  C-commands  $\beta$

$\alpha$  is an X (head, eg: Noun, Verb, Preposition, adjective,  $\neq$  inflection), and every maximal projection dominating  $\beta$  dominates  $\alpha$ .

### Movement, Empty Category and Co-indexing

In GB, Move  $\alpha$  is described as 'move anything anywhere', though it provides restriction for valid movements:

### Binding theory:

Binding is defined by Sells (1982)

- $\alpha$  binds  $\beta$  iff
- $\alpha$  C-commands  $\beta$  and
- $\alpha$  and  $\beta$  are co-indexed

[e<sub>i</sub> INFL kill Mukesh]  
[Mukesh<sub>i</sub> was killed (by e<sub>i</sub>)]  
Mukesh<sub>i</sub> was killed.

Empty clause (e<sub>i</sub>) and Mukesh (NP<sub>i</sub>) are bound. This theory gives a relationship b/w NPs

Now, binding theory can be given as follows:

- An anaphor (t<sub>e</sub>) is bound in its governing category.
- A pronominal (t<sub>p</sub>) is free in its governing category.
- An R-expression (- <sub>$\alpha$</sub> , -<sub>P</sub>) is free.

### Empty Category Principle (ECP)

'Proper government':

- $\alpha$  properly governs  $\beta$  iff
- $\alpha$  governs  $\beta$  and  $\alpha$  is lexical (ie. N, V, A or P) or
- $\alpha$  locally  $\bar{A}$  binds  $\beta$ .

The ECP says 'A trace must be properly governed'.

This principle justifies the creation of empty categories during NP-trace and wh-trace and also explains the subject/object asymmetries to some extent.

- What<sub>i</sub> do you think that Mukesh ate e<sub>i</sub>?
- What<sub>i</sub> do you think Mukesh ate e<sub>i</sub>?

### Bounding and control Theory:

There are many other types of constraints on movement.

The long distance movement for complement clause can be explained by bounding theory if NP and S are taken to be bounding nodes. The theory says that the application of movement may not cross more than one bounding node. The theory of control involves syntax, semantics and pragmatics.

### Case Theory and Case filter:

In GB, Case Theory deals with the distribution of NPs and mentions that each NP must be assigned a Case.

4(a) Explain Lexical functional grammar (LFG) model for the sentence "She saw stars in the sky".

[06]

### Lexical functional Grammar model

LFG represents sentences at two syntactic levels - constituent structure (c-structure) and functional structure (f-structure).

Based on Wood's Augmented Transition Networks (1970), which used phrase structure trees to represent the surface structure of sentences and the underlying predicate-argument structure.

### C-structure and f-structure in LFG:

An LFG is aimed at providing exact computational algorithms. It provides well defined objects called constituent structure (c-structure) and functional structure (f-structure).

The c-structure is derived from the usual phrase and sentence structure syntax, as in CFG. However the grammatical-functional role cannot be derived directly from phrase and sentence structure, functional specifications are annotated on the nodes of c-structure, which applied on sentences, results in f-structure, which is the final product.



which encodes the information obtained from phrase +  
Sentence Structure rules and functional specification.

eg: She saw stars in the sky

CFG rules to handle this sentence are:

$S \rightarrow NP VP$

$VP \rightarrow V \{NP\} \{NP\} PP^* \{S'\}$

$PP \rightarrow P NP$

$NP \rightarrow Det N \{PP\}$

$S' \rightarrow Comp S$

where

S: Sentence V: verb P: Preposition N: noun

S': clause Comp: Complement  $\{ \}$  optional

\*: phrase can appear any number of times including blank.

When annotated with functional specifications,

the rules become:

Rule 1:  $S \rightarrow NP VP$   
 $\uparrow \text{e.i.w.} = \downarrow \uparrow = \downarrow$

Rule 2:  $VP \rightarrow V \{NP\} \{NP\} PP^* \{S'\}$   
 $\uparrow \text{obj} = \downarrow \uparrow \text{obj} 2 = \downarrow \uparrow (\downarrow \text{case}) = \downarrow \uparrow \text{comp} = \downarrow$

Rule 3:  $PP \rightarrow P NP$   
 $\uparrow \text{obj} = \downarrow$

Rule 4:  $NP \rightarrow \{Det\} N \{PP\}$   
 $\uparrow \text{Adjunct} = \downarrow$

Rule 5:  $S' \rightarrow Comp S$   
 $\uparrow = \downarrow$

$\uparrow$  - f-structure of the mother node that it is on the left hand of the rule.

$\downarrow$  - f-structure of the node under which it is denoted.

Hence in Rule 1, ( $\uparrow \text{Subj} = \downarrow$ ) indicates that the f-structure of the first NP goes to the f-structure of the subject of the sentence.

while ( $\uparrow = \downarrow$ ) indicates that the f-structure of the VP node goes directly to the f-structure of the sentence VP.

Similarly in Rule 2, the f-structure of VP is defined by the lexical item V, the two optional NPs, any no. of PP's and the optional clause (S').

She saw stars.

She N ( $\uparrow \text{pred}$ ) = 'PRO'

( $\uparrow \text{pers}$ ) = 3

( $\uparrow \text{Num}$ ) = SG

( $\uparrow \text{gen}$ ) = FEM

( $\uparrow \text{case}$ ) = NOM

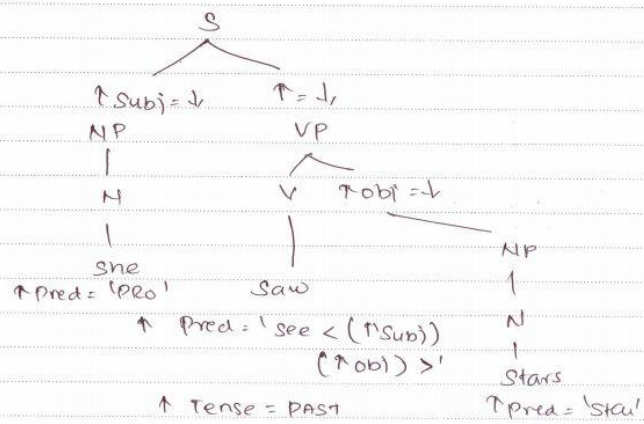
Saw V ( $\uparrow \text{Pred}$  = 'saw < ( $\uparrow \text{Subj}$ ) ( $\uparrow \text{Obj}$ ) >')

( $\uparrow \text{Tense}$  = PAST)

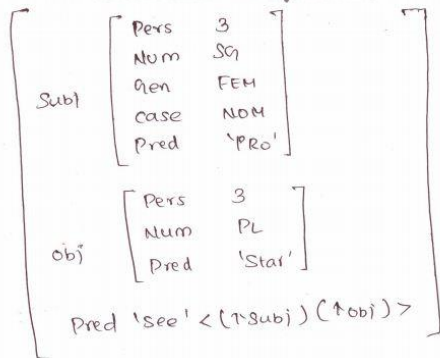
Stars N ( $\uparrow \text{Pred}$  = 'Star'

$\uparrow \text{Pers}$  = 3

$\uparrow \text{Num}$  = PL



C - Structure of sentence.



LFQ imposes three conditions on  $\uparrow$ -Structure:

Consistency: In a given  $\uparrow$ -structure, a particular attribute may have at most one value. Hence, while unifying two  $\uparrow$ -structures if the attribute Num has value SG in one and PL in the other, it will be rejected.

LFQ imposes three conditions on  $\uparrow$ -Structure:

Consistency: In a given  $\uparrow$ -structure, a particular attribute may have at most one value. Hence, while unifying two  $\uparrow$ -structures if the attribute Num has value SG in one and PL in the other, it will be rejected.

Completeness: A function is called governable if it appears within the Pred value of some lexical form. eg: Subj, Obj, Adjunct. Adjunct is not a governable function.

Coherence: coherence maps the completeness property in the reverse direction. It requires that all governable functions of an  $\uparrow$ -structure, and all its subsidiary  $\uparrow$ -structures must be governed by their respective Predicals.

(b) Explain in detail about Paninian framework and its layered representation.

[04]

Paninian Framework

Paninian Grammar-based model was written by Panini in 500BC in Sanskrit, the framework can be used for other Indian languages and possibly some Asian languages.

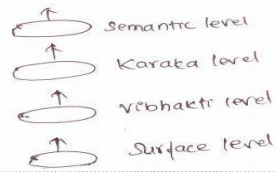
Unlike English, Asian languages are SOV (Subject-Object-Verb) ordered and inflectionally rich. The inflections provide important syntactic and semantic cues for language analysis.



### Layered representation in PG

The GB Theory represents three syntactic levels: deep structure, surface structure and logical form, where the LF is nearer to semantics.

Unlike GB, Paninian Grammar framework is said to be syntactico-semantic that is one can go from surface layer to deep semantics by passing through intermediate layers.



### Karaka Theory:

- Central theme of PG framework
- Karaka relations are assigned based on the roles played by various participants in the main activity.
- These roles are reflected in the case markers + post position markers.

### Issues in Paninian grammar:

Two problems:

- i) Computational implementation of PG and
- ii) Adaptation of PG to Indian + other similar languages.

5. Explain the stages of morphological parsing using FST.

[10]

### Morphological Parsing:

- Morphology is a subdiscipline of linguistics.
- Studies word structure and the formation of words from smaller units (morphemes)
- goal to discover the morphemes that build a given word.
- morphemes => smallest meaning-bearing units in a language.

- morphemes  $\left\{ \begin{array}{l} \text{stems} \\ \text{affixes} \end{array} \right.$

→ Stem - main morpheme i.e) The morpheme that contains the central meaning.

→ Affixes - modify the meaning given by the stem

- Affixes  $\div$
- i) Prefix
  - ii) Suffix
  - iii) Infix
  - iv) Circumfix

3 main ways of word formation:

- inflection
- derivation
- Compounding

Inflection → a root word is combined with a grammatical morpheme to yield a word of the same class as the original stem.

Derivation → combines a word stem with a grammatical morpheme to yield a word belonging to a different class.

eg: formation of the noun computation from the verb compute.

⇒ the formation of a noun from a verb or adjective is called nominalization.

⇒ Compounding is the process of merging two or more words to form a new word.

eg: personal computer, desktop, overlook.

- A morphological parser uses the following information:

Sources:

i) Lexicon:

- lexicon lists stems and affixes together with basic information about them.

ii) Morphotactics:

- there exists certain ordering among the morphemes that form a word.

- can't be arranged arbitrarily.

eg: restlessness is a valid word

rest-ness-less is invalid.

- This morphotactics deals with the ordering of morphemes, it describes the way morphemes are arranged or touch each other.

iii) Orthographic:

- spelling rules that specify the changes that occur when two given morphemes combine.

eg: y → ier. spelling rule changes

easy → easier. not to easier.

- morphological analysis can be avoided if an exhaustive lexicon is available that lists features for all the word-forms of all the roots.

| Word form           | Category | Root   | Gender    | Number   | Person          |
|---------------------|----------|--------|-----------|----------|-----------------|
| ghodhaa             | noun     | ghodaa | masculine | singular | 3 <sup>rd</sup> |
| ghodhi <sup>o</sup> | do       | do     | feminine  | do       | do              |
| ghodhon             | do       | do     | masculine | plural   | do              |
| ghodhe              | do       | do     | do        | do       | do              |



But this approach has several limitations:

- ① It puts a heavy demand on memory.
  - have to list every form of the word which results in a large number of redundant entries in the lexicon.
- ② Exhaustive lexicon fails to show the relationship between different roots having similar word forms.
  - approach fails to capture linguistic generalization, which is essential to develop a SLM capable of understanding unknown words.
- ③ Morphologically complex languages like Turkish, the no. of possible word forms may be theoretically infinite.
  - Not possible to list all possible word forms in these languages.

- Stemming algorithms work in two steps:

i) Suffix removal: this step removes predefined endings from words.

ii) Recoding: this step adds predefined endings to the output of the first step.

- These two steps can be performed sequentially as in Lovin's stemmer or in Porter's stemmer.

eg: porter's stemmer makes use of the following transformation rule!

rotational  $\rightarrow$  ate

to transform word such as 'rotational' into 'rotate'.

- In this model, a word is represented in

- i) lexical level form  $\rightarrow$  represents the concatenation of its constituent morphemes.
- ii) surface level form.

$\downarrow$   
represents the actual spelling of the word

$\rightarrow$  morphological parsing is a mapping from the surface level into morpheme and feature sequences on the lexical level.

- surface form

'playing'  $\rightarrow$  play + V + PP

- lexical form

$\Rightarrow$  stem 'play' followed by the morphological

information +V, +PP

Surface level: playing

Lexical level: play + VPP

$\Rightarrow$  Surface form 'books'

lexical form  $\frac{\text{book} + \text{N} + \text{PL}}$

$\downarrow$   
morphological information

$\downarrow$   
Surface level form is plural noun.

This model is usually implemented with a kind of finite state automata  $\Rightarrow$  finite state transducer (FST)

- Transducer maps a set of symbols to another.

FST - 2 state automaton

- recognizes or generates a pair of strings.  
 - Passes over the input string by consuming the input symbols on the tape it traverses and consists it to the output string in the form of symbols.

- Finite state transducer is a 6 tuple

$(\Sigma_1, \Sigma_2, Q, \delta, S, F)$

$Q$  - is a set of states

$S$  - is the initial state

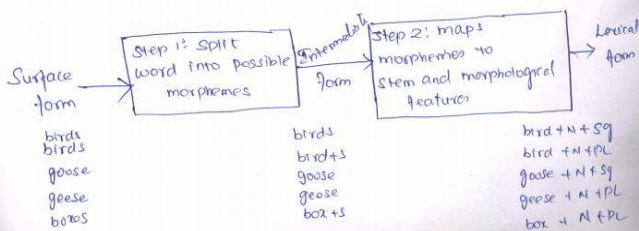
$F \subseteq Q$  - set of final states

$\Sigma_1$  - input alphabet

$\Sigma_2$  - output alphabet

$\delta$  - function mapping

### Two Step morphological parser



### FST



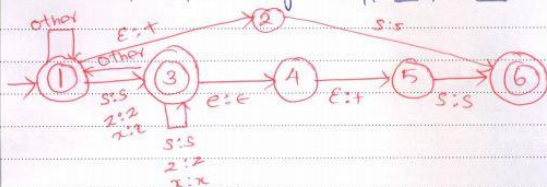
$\Rightarrow$  FST represents the information that the comparative form of adjective 'less' is 'lesser', 'e' here is empty string.

$\Rightarrow$  The automaton is inherently bi-directional. The same transducer can be used for analysis  $\rightarrow$  surface up, for generation  $\rightarrow$  lexical down.

$\rightarrow$  The plural form of regular nouns usually end with -s or -es.

However a word ending in 's' need not be plural. There are number of singular words ending in 's' like miss, ass.

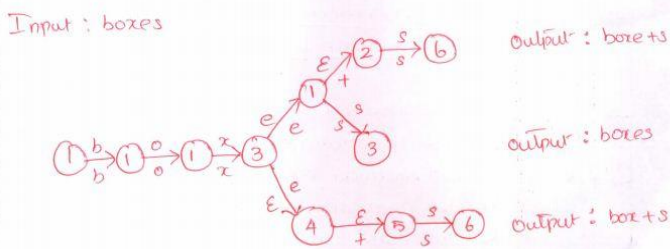
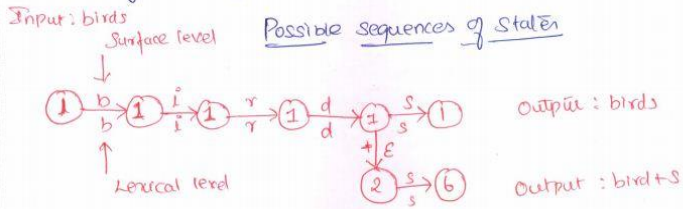
$\rightarrow$  One of the required translations is the deletion of the 'e' when introducing a morpheme boundary. This deletion is usually required for words ending in 'xes, ses, zes' eg: (Supplies, boxes).



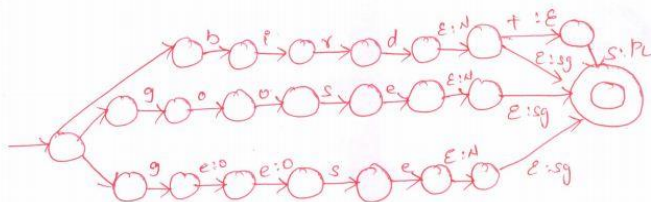
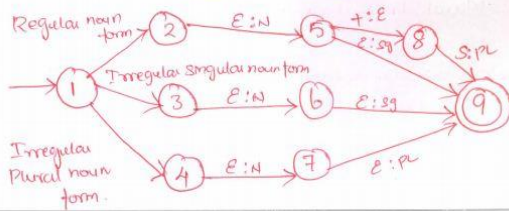
Simplified FST, mapping English nouns to the intermediate form.



- Next step is to develop a transducer, that does the mapping from the intermediate level to the lexical level.
- The input to transducer has one of the following forms:
  - \* Regular noun stem, eg: bird, cat
  - \* Regular noun stem + s, eg: bird+s
  - \* Singular Irregular noun stem, eg: goose
  - \* Plural Irregular noun stem, eg: geese
- First case, the transducer has to map all symbols of the stem to themselves and then output N and s.
- Second case, it has to map all symbols of the stem



Transducer for step 2



A transducer mapping nouns to their stem and morphological features

6. Describe the concept of spelling error detection and correction by computing minimum edit distance algorithm.

## Spelling Error Detection & Correction:

- In computer-based information slm, errors of typing and spelling causes variation b/w strings.

- These errors are investigated and that are:

\* Single character omission

\* Insertion

\* Substitution

\* Reversal  $\Rightarrow$  most common typing mistake.

- Dameanu (1964) reported that over 80% of the typing errors were single-error misspellings:

i) Substitution of a single letter

ii) Omission of a " "

iii) Insertion of " " "

iv) Transposition of two adjacent letters.

Single character omission occurs when a single character is missed (deleted) eg: when 'concept' is accidentally typed as 'concp'.

Insertion error refers to the presence of an extra character in a word. eg: when 'error' is misspelled as 'erron'.

Substitution error occurs when a wrong letter is typed in place of the right one, eg: 'error' where 'r' appears in place of 'o'.

Reversal refers to a situation in which the sequence of characters is reversed. eg: 'aer' instead of 'are', also called as transposition.

$\Rightarrow$  Optical character Recognition (OCR) and other automatic reading device introduce errors of substitution, deletion, and insertion but not of reversal.

$\Rightarrow$  Spelling errors are mainly due to phonetic and it is one of two distinct categories:

① Non word errors

② Real word errors.

$\rightarrow$  When an error results in a word that does not appear in a lexicon or valid orthographic word form, it's termed as non word error.

$\rightarrow$  Most of the research found the detection of non word error:

$\rightarrow$  2 main techniques  $\left\{ \begin{array}{l} \text{n-gram analysis} \\ \text{dictionary lookup.} \end{array} \right.$

$\rightarrow$  Now it's solved problem.



→ Real word error results in actual words of the language.

→ Occurs due to typographical mistakes or Spelling errors.

eg: substituting the spelling of a homophone or near homophone. Such as Peace for piece  
Meat for meet.

→ Real world errors cause:

- i) Local syntactic errors
- ii) Global syntactic errors
- iii) Semantic errors
- iv) Discourse errors or pragmatic levels.

→ To decide a word is wrong is only by contextual information.

⇒ Spelling correction:

- Detecting
- Correcting errors.

⇒ Spelling correction:

- Detecting
- Correcting errors.

Error detection:

- process of finding misspelled words.

Error correction:

- process of suggesting correct words to a misspelled word.

⇒ 2 ways:

- 1) Isolated error detection and correction
- 2) Context-dependent error detection + correction.

Isolated error detection and correction:

- Each word is checked separately, independent

② Context dependent error detection and correction:

- utilize the context of a word to detect and correct error.
- It requires grammatical analysis and is thus more complex and language dependent.

- Even in context the list of candidate words must first be obtained using an isolated word before making a selection depending on the context.

## Spelling error Correction algorithm:

### ① Minimum edit distance:

the minimum edit distance between two strings is the minimum number of operations required to transform one string into another.

### ② Similarity key techniques:

the idea is to change a given string into a key such that similar strings will change into the same key. SOUNDEX (Knudsen & Russett 1918) is an example which is used in phonetic spelling correction appln.

### ③ n-gram based techniques:

- used for both nonword and real word error detection. eg. in English bi-grams, tri-grams of letters never occur. ex. trigram - q.st and bigram - q.d.

## Minimum edit distance

- Number of insertions, deletion and substitutions required to change one string into another.

- the distance between two strings is the minimum edit distance.

eg. MED blo tutor and tumor is 2

→ substitute 'm' for 't'

→ Insert 'u' before 'r'

- Edit distance between two strings can be represented as binary function, ed which maps two strings to their edit distance. ed is symmetric. For any two strings s and t,  $ed(s,t)$  is always equal to  $ed(t,s)$ .

Input: Two strings, X and Y

Output: The minimum edit distance between X and Y

$m \leftarrow \text{length}(X)$

$n \leftarrow \text{length}(Y)$

for  $i=0$  to  $m$  do

$dist[i,0] \leftarrow i$

for  $j=0$  to  $n$  do

$dist[0,j] \leftarrow j$

for  $i=0$  to  $m$  do

for  $j=0$  to  $n$  do

$$dist[i,j] = \min \left\{ \begin{array}{l} dist[i-1,j] + \text{insert\_cost}, \\ dist[i-1,j-1] + \text{subst\_cost}(X_i, Y_j), \\ dist[i,j-1] + \text{delet\_cost} \end{array} \right\}$$

Minimum edit distance algorithm.

| # | t | u | m | o | u | r |   |
|---|---|---|---|---|---|---|---|
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| t | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| u | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| t | 3 | 2 | 1 | 1 | 2 | 3 | 4 |
| o | 4 | 3 | 2 | 2 | 1 | 2 | 3 |
| r | 5 | 4 | 3 | 3 | 2 | 2 | 2 |

Computing minimum edit distance

7(a) What is Regular Expression and write down the regexes for email address.

[6]

### Regular Expressions:

- RE (or) Regexes.
- Pattern matching std for string parsing and replacement.
- RE can be used to \* parse dates \* urls and email addresses \* log files \* configuration files \* command line switches \* programming scripts.

$\backslash n$  - matches a newline character

$\backslash t$  - matches a tab

$\backslash d$  - matches a digit [0-9]

$\backslash D$  - matches a nondigit

$\backslash w$  - matches an alphanumeric character

$\backslash W$  - matches a non alphanumeric character

$\backslash s$  - matches a whitespace character

$\backslash S$  - matches a non whitespace character

$\backslash$  - escape special characters

$\backslash .$  - matches a dot

$\backslash *$  - matches a \*

$\backslash \backslash$  - matches a backslash

-> email address consist of a non empty sequence of characters followed by 'at' symbol, @, followed by non empty sequence of characters ending with pattern like .xx, .xxx, .xxxx, etc.

The RE for an email address is:

$$^{[A-Za-z0-9\.\-]+}@[A-Za-z0-9\.\-]+[A-Za-z0-9][A-Za-z0-9]\$$$
$$^{[A-Za-z0-9\.\-]+}$$

↓

match a positive no. of acceptable characters at the start of the string.

@ => matches the @ sign

$$[A-Za-z0-9\.\-]+$$

↓

match any domain name, including a dot

$$[A-Za-z0-9][A-Za-z0-9]\$$$

↓

match two acceptable characters but not a dot. This ensures that the email address ends with .xx, .xxx, .xxxx, etc.

(b) Describe in detail about Finite State Automata with an example.



## Finite State Automata

Some games that we have played that fits the following description:

- Pieces are set up on a playing board
- Dice are thrown or a wheel is spun and a number is generated at random.
- Based on the number appearing on the die, the pieces on the board are rearranged specified by the rules of the game.
- Here, all the possible positions of the piece on the board and call them states.
- State begins  $\rightarrow$  initial state
- winning positions  $\rightarrow$  final state.

$\rightarrow$  A finite automaton has the following properties.

- A finite set of states, one of which is designated the initial or start state, and one or more of which are designated as final state.
- A finite alphabet set,  $\Sigma$  consisting of input symbols.
- A finite set of transitions that specify for each state and each symbol of I/P alphabet, the state to which it next goes.

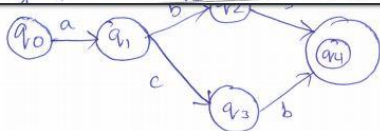
A finite automaton  $\left\{ \begin{array}{l} \text{deterministic} \\ \text{Non deterministic} \end{array} \right.$

$\downarrow$   
more than one transition out of a state is possible for the same input symbol.

### Example:

Suppose  $\Sigma = \{a, b, c\}$ , the set of states =  $\{q_0, q_1, q_2, q_3, q_4\}$  with  $q_0$  being the start state and  $q_4$  the final state the transition are:

- From state  $q_0$  and with input  $a$ , go to state  $q_1$
- " "  $q_1$  " " "  $b$ , go to state  $q_2$
- " "  $q_1$  " " "  $c$ , " " "  $q_3$
- " "  $q_2$  " " "  $b$ , " " "  $q_4$



Nodes  $\rightarrow$  states, arcs  $\rightarrow$  transitions.

$\odot \rightarrow$  final state.

$\Rightarrow$  There is exactly one transition leading out of each state hence it's deterministic automaton.

$\Rightarrow$  FSA are used in:

- \* Areas including linguistics
- \* Electrical engineering
- \* Computer Science
- \* Mathematics and logic

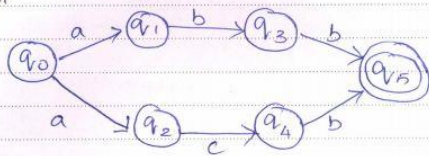
### ① Deterministic finite state automaton

- 5 tuple  $(Q, \Sigma, \delta, S, F)$

- $Q \rightarrow$  set of states
- $F \subseteq Q \rightarrow$  set of final states
- $\Sigma \rightarrow$  alphabet
- $\delta \rightarrow$  transition function.
- $S \rightarrow$  start state

→ The transition function of a non-deterministic finite state automaton (NFA) maps  $Q \times (\Sigma \cup \{\epsilon\})$  to a subset of the power set of  $Q$ .

→ for each state, there can be more than one transition on a given symbol, each leading to a different state.



→ A path is a sequence of transitions beginning with the start state, a path leading to one of the final states is a successful path.

→ FSAs encode regular language.

→ The language that an FSA encodes is the set of strings that can be formed by concatenating the symbols string accepts.

→ Language defined by this automaton can be regular expression /ablabcb/

→ All transition rules listing is inconvenient, so represent automaton as a state transition table.

Rows → States

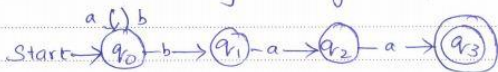
Columns → input

$\phi$  → missing transition

| State        | a      | b      | c      |
|--------------|--------|--------|--------|
| Start: $q_0$ | $q_1$  | $\phi$ | $\phi$ |
| $q_1$        | $\phi$ | $q_2$  | $\phi$ |
| $q_2$        | $\phi$ | $q_4$  | $\phi$ |
| $q_3$        | $\phi$ | $q_4$  | $\phi$ |
| Final: $q_4$ | $\phi$ | $\phi$ | $\phi$ |

→ strings containing

NFA implementing this regular expression:



/ (a|b)\*baa /

State transition table:

| State        | a         | b              |
|--------------|-----------|----------------|
| Start: $q_0$ | { $q_0$ } | { $q_0, q_1$ } |
| $q_1$        | { $q_2$ } | $\phi$         |
| $q_2$        | { $q_3$ } | $\phi$         |
| Final: $q_3$ | $\phi$    | $\phi$         |

Two automata that define the same language are said to be equivalent.

DFA:

