

01-11-23

UNIT-I

Wednesday

Words and their components

1. Tokens

2. Lexemes

3. Morphemes

4. Typology

① Tokenization: It is a way of separating a piece of text into smaller units called Tokens.

* Token may be either words, characters or subwords.

Ex: NLP is sub branch of AI

Tokens = 6.

② Lexemes:

The set of alternate forms that can express a word are called lexical items.

Ex: The lexem "play" can have many forms such as playing, play, plays, etc.

③ Morphemes:

* It is smallest unit of a word provides a specific meaning to the string of letters.

* There are two types of morphemes

(i) Free morpheme:

un+happy = unhappy

play+ing = playing

(ii) Bound morphemes:

Dis+place+ment = Displacement

Un+employ+ment = unemployment

im+prison+ment = imprisonment



④ Typology:-

- * Typology is a study, the ways in which the languages of the world vary in their patterns
 - * It is concerned with discovery what grammatical patterns are common to many languages; they are specified according to three criteria,
- (i) Genealogical families
 - (ii) Areal families
 - (iii) Geographical families
- (i) Genealogical families:
This classification indicates the historical connections between the languages and it uses the historical and linguistic criteria as a bases.

Ex: In Europe the Basque language is called language isolates, as it cannot relate to any other language.

(ii) Areal families:
These kinds of languages are creating language unions such as Balkan language union, encompassing Macedonian, Bulgarian, Serbian, Albanian languages.

(iii) Geographical families:
These covers the languages that are close by and have developed similar characteristics in terms of structure.

03/11/23 Issues and challenges
Friday

① Irregularity

② Ambiguity

③ Productivity

① Irregularity

Irregularity means existence of words and structures that are not described appropriately by prototypical linguistic (model) forms.

Ex:- eat ate eats eating eaten
 ↓ ↓ ↓ (part of speech)
 v1 (past) part(v2) eat+_s eating eaten

① watch watches watching watched
 ↓ ↓ (part of speech)
 v1 watching watched

② write wrote writing written
 ↓ ↓ (part of speech)
 v2 writing written

② Ambiguity

Types of Ambiguity:
 1. Lexical (semantic) [meaning of the word]
 2. Syntactical (structural)

- ① Lexical (semantic) [meaning of the word]
- ② Syntactical (structural)
- ③ Anaphoric [reference] [referring from another sentence]
- ④ Pragmatic

① Lexical Ambiguity:

The ambiguity of a single word is called Lexical Ambiguity (or) when

Lexical Ambiguity happens when the word can be themselves interpreted when the meaning of the word can be themselves misinterpreted (or)

Semantic ambiguity happens when the sentence contains an ambiguous words or phrases

Ex: watch (having 2 meanings).

① wrist watch

② watching movie

Ex: 2

Live - (Live match)

Live - (existing)

Ex: 3 plant - (planting trees)

plant - (solar plant, steel plant, water plant)

bat - (cricket bat)

bat - (mammal)

The car hit the pole while it was moving.

→ The sentence is having two meanings:

- ① the car, while moving hit the pole
- ② the car hit the pole, while pole was moving.

③ Syntactical Ambiguity:

This type of ambiguity occurs when sentence is parsed in different ways

Ex: The man saw the girl with the telescope.

→ It is ambiguous whether the man saw the girl carrying telescope or he saw her through his telescope.

Ex: This book costs a large amount

Ex: The chicken is ready to eat

Ex: Ram bought green shirts & shoes.

→ It is ambiguous to who do we refer shirts and shoes.



③ Anaporic Ambiguity:

This kind of ambiguity arises due to the use of

Anapora entities

Ex: The horse ran up the hill. It was steep. It got tired.



Here, the anaporic reference "It is ambiguous"

Ex: Jyoshna invited Sudeshna for a visit, but she told her, she had to go to work.

Ex: I am allergic to tomatoes. Also fish

④ Pragmatic Ambiguity: This type of ambiguity occurs when phrases give multiple interpretation.

Ex: I like you too.

multiple interpretations are I like you just you like me or someone else.

(or) I like you someone else that.

⑤ Pragmatic Ambiguity:

Ex: Mother : what you want to eat

son : If ice cream is good this time of year

Ex: I tried to reach her and on the phone, but she didn't answer.

Productivity:

- ① Set objectives
- ② Boost staff morale
- ③ Better communication
- ④ Learning and development
- ⑤ changing behavior

06/11/23

Monday

NLP Tools and techniques

- ① MonkeyLearn (diagnosis, call centers)
- ② IBM Watson
- ③ Google Cloud NLP API
- ④ NLTK (many methods & fns are implemented)
- ⑤ spaCy
- ⑥ pytorch

NLP techniques

- ① Lemmatization
- ② Lexicon
- ③ NER
- ④ N-gram

① Lemmatization:

Lemmatization is a process of normalize a word and remove its inflection form to arrive at the base word.

Ex: base word for gone or went is go

for better & best → good



Scanned with OKEN Scanner

write = wrote · written

base word

② Lexicon:
Lexicon is the vocabulary of a language or subject along with its usage.

Ex:- medical terminology, computer [default, web, shutdown, application]

③ NER: [name, Named entity recognition]

NER is a popular technique used in information extraction to identify and segment the named entities and classify or categorize them under various pre-defined clauses. En Real world objects (people, places, organization, name)

Delhi is ~~the~~ the capital of India

sundar pichhai is the ceo of Google

Name of the person

sudeshna ~~is~~ belongs to Kadapa

Name

④ N-gram

N-gram is a contiguous sequence of N tokens from a given text

* N-gram based models credit the most probable words that might follow the entered text sequence

En Applications: including machine translation, auto-complete

Corpus:

- * Corpus in NLP means collection of text or other digital data across languages

Ex: Wikipedia

10/11/23

Friday

Morphological models (study of internal structure of the word)

- ① Dictionary Lookup (collection of meanings, Nouns, adverbs, adj, pron)
- ② Finite state morphology (having nodes & edges states arcs)
- ③ Unification based morphology
- ④ Functional morphology
- ⑤ Morphology Induction

Functional

Number marker - Three catz

possessive marker - John's book

inflection - she walks

* Morphological model is the study of internal structure of word.

① Dictionary Lookup:

- * Dictionary Lookup is the fundamental concept in NLP

* That involves searching for a word or phrase in a dictionary or lexicon to retrieve its definition.

translation, parts of the speech or reliable relevant information

It is often used to extract meaning of the texts

It is often used to extract meaning of the texts

a) lexicon / dictionary:

- * In NLP, a dictionary or lexicon is a collection of words or phrases along with associated information, such as meaning, part of the speech, pronunciation, translation, and more.
- * Dictionaries can be simple lookup tables or complex databases.

Purpose of dictionary lookup:

- * Dictionary lookup is used to provide context and understanding to words, or phrases in the text.
- * It can be used for various tasks like word sense disambiguation:
 - ① Identifying a correct meaning of the word based on the context.
 - ② Language translation: - converting words or phrases from one language to another.
 - ③ Parts of speech: Assigning the appropriate part of the speech to word like noun, verb, etc.

Information retrieval

- * Retrieving additional information such as synonyms, antonyms, or definitions.

Named Entity Recognition: Identifying named entities like names of people, places or organizations

Ex: Let consider a word bank and perform dictionary lookup



word: Bank

Part of speech: Noun

def: Financial Institute

Synonyms: FI, credit union

Antonyms: borrowers, debtor

word: bank

pos: Noun

def: side of river, lake

Synonyms: shore, edge

Antonyms: intellect, foolish

word: Bank

pos: verb

def: to deposit money in F.I

Synonyms: deposit, invest

Antonyms: withdraw, disburse

11/11/23

Saturday

② Finite state morphology:

Finite state morphology is a computational

linguistic approaches used to model and

analyze the morphology of a languages in NLP

* It focus on the use of finite machines [FSM]

and finite state transducer [PST] to represent

and manipulate the structure and variation of

words in a language

FSM:

FSM's are computational models that consists

of states, transitions and input symbols.

In the context of morphology state represents

linguistic forms or morphemes

the state transition

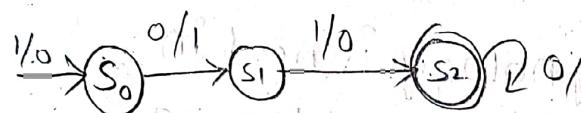
the input symbol

the output symbol



Scanned with OKEN Scanner

- Transitions represents transformations between these forms and input symbols represents the letters or sounds of language



FST:

- * FST's are the finite state machines that can map input symbols to output symbols

- * They are commonly used in finite state morphology to represent relationship between the surface form and root form

Ex: playing = play + ing
 root surface

Purpose of Finite state Morphology

- * It is used for various linguistic tasks including

- stemming, lemmatization, inflection and understanding the internal structure of words
- understanding the internal structure of words
- It is particularly valuable for languages with complex inflectional processes, where words can change to base form according to their grammatical features

Ex:- noun ^{non-} pluralization

- (i) singular nouns ends in -s, -o, -x and to pluralize them we add -es

Ex: bus, box, mango

- (ii) singular noun ending in a consonant followed by vowel changes to -ies

Ex: baby → babies Monkey → Monkeys

③ Unification based morphology

- * It employs unification as a key concept which is a process of merging or matching linguistic features and constraints to derive word forms
- * Unification based morphology is used in computational linguistics, formal language grammar.
- * Unification linguistics
- * Unification is the process of merging linguistic features and constraints from different sources to create an unified representation of a linguistic object such as word.
- * In this approach features and constraints are represented using feature structures which can be taught of attribute ~~phase~~ pairs.

Features and constraints :

- * Features and constraints describes various aspect of a words such as its lexical form, grammatical category and so on.
- * Constraints specify the rules and restrictions on how features can be combined.

Purpose of unification based morphology

- * The goal of unification based morphology is to model the structure and formation of words.

Ex: Features and constraints for noun pluralization



1) Features:

Lexical Form(LF): The base form of the noun is cat.

Number(Num): The grammatical number of a noun

[singular(+)plural] is specified.

Constraints:

① constraint-1: If the number is singular and
lexical form ends with 's, 'o' etc.

constraint-2: If the number is singular, lexical
form ends with 'y'.

Rule: adjective + er = comparative adjective
Ex: fast + er = faster

4) Functional morphology:

Morphemes are analyzed based on their roles within sentences such as how they contribute to meaning, syntax and structures.

This approach considers the functional aspects of semantic, pragmatic aspects of morphemes.

Ex: (i) Number marker - ("Three cats")
the plural s indicates three numbers of cats.

(ii) Possessive marker - ("John's book")
In the phrase - ("John's book")
s serves as the book belongs to John.

(iii) Verbal Inflection - ("she walks")
In the sentence - ("she walks")
s indicates the third person singular

present tense

Morphology Induction

- * Aims to automatically identify and extract morphological patterns or units from a corpus of text, without prior knowledge of the language morphology. This process involves discovering how words are formed, segmented, and inflected in a language.

① Unsupervised learning: Morphology induction

typically uses unsupervised machine learning approaches, such as statistical models or algorithms, to analyze a large corpus of text and identify patterns within it.

② Identifying morphemes:

The goal is to identify morphemes, which are the smallest meaningful units in a language. Morphemes can include prefixes, suffixes, roots and other components that contribute to word formation and meaning.

③ clustering and segmentation

It involves clustering similar words together and segmenting them into their constituent morphemes based on statistical or distributional patterns. This process often involves grouping words that share similar prefixes, suffixes or root.



Ex:

Unhappiness, rebuild and undo

In morphological induction process, the algorithm identify patterns such as,

un - is a prefix indicates negation

-ness - is a suffix indicates quality or state

re - is a prefix that indicates doing something again

build and do are root words

the algorithm might group these words

1. clusters 1 : unhappiness, unhappy, unhappily

Prefix : un -

Suffix : -ness, -ly

cluster 2 : rebuild, rebuilding, rebuilt

Pre :- re -

Root : build

cluster 3 : undo, undid, undoing

Pre :- un

Root :- do

cluster 4 : happy, happier, happiest

Pre : happy

Root :- -ier, -iest