

NLP

Module 2 - Word Level Analysis

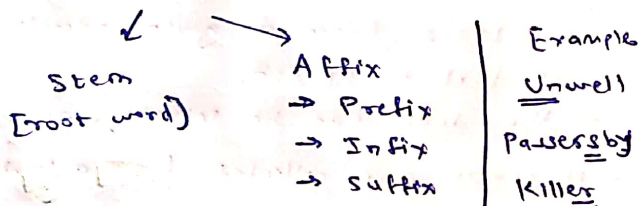
Morphology Analysis

- Morph → Form / shape
ology → study of
- Morphology is the study of the way words are formed from morphemes.
- Morphemes → smallest/meaning-bearing units.
- Example: Cats
morpheme cat and
morpheme -s

Root morphological variant.

- ① Walk - walks, walked, walking
- ② Noise - Noisy
- ③ atom - Atomic
- ④ order - reorder, orderly

① Morphemes



② Free and Bound morphemes.

→ Free morphemes are independent words, such as camera, pen, boat.

③ Free morpheme

- ① Lexical morpheme ← Picture words, ex → yellow
↳ Unlimited
- ② Grammatical morpheme
↳ Limited

④ Bound morpheme

- Bound morphemes does not have meaning unless attached with free morphemes.
- eg. -ing

- 2 Types.

① Inflection morpheme

Cat + s = <u>Cats</u>	It could be Suffix or Infix. [No Prefix]
↓ Noun	
↓ Noun	

② Derivation morpheme

danger + ou = dangerous
↓ Noun
↓ Adjective

FST Properties

- ① Union
- ② Inversion
- ③ Composition

Inflection v/s Derivation

Inflection

- ① It is a morphological process that adopts existing words so that they function effectively in sentences without changing pos of base morpheme.
- ② They close off the word.
eg. Plays.
- ③ They cannot be found in dictionaries
- ④ Inflection is relevant to syntax.
- ⑤ It is obligatory
- ⑥ Express the same concept as base

Derivation

- ① It is concerned with the way morphemes are connected to existing lexical forms as affixes.
- ② They never close off the word.
eg. Playful.
- ③ They can be found in dictionaries.
- ④ Derivation is irrelevant to syntax.
- ⑤ It is optional
- ⑥ Expresses a new concept.
- ⑦ Derivation is semantically irregular
- ⑧ It is expressed closed to the root
- ⑨ Meanings are less relevant to the meaning of the base
- ⑩ It expresses a relatively abstract meaning.

- ⑪ Can be prefix and suffix.

⑫ Example:

Cat + s = <u>Cats</u>
↓ Noun

danger + ou = dangerous
↓ Noun
↓ Adjective

- FST is a multi-function device

① Translator

- It reads one string on one tape and outputs another string

② Recognizer

- It takes a pair of strings as two tapes and accepts/rejects based on their match

③ Generator

- It outputs a pair of strings on two tapes, along with yes/no result based on matching.

④ Relator

- It computes the relationship between two sets of strings available on two tapes

Stemming

- ① Stemming is faster because it chops words without knowing the context. the word is given sentence.
- ② It is a rule based approach.
- ③ Accuracy is less.
- ④ When we convert any word into root from then stemming may create the non-existence meaning of word.
- ⑤ Stemming is preferred when the meaning of the word is not important of analysis.
Example: Spam Detection.
- ⑥ Example:
"studies" \Rightarrow "studi"

Lemmatization

- ① Lemmatization is slower as compared to stemming but it knows the context of the word before proceeding.
- ② It is a dictionary based approach.
- ③ Accuracy is more.
- ④ Lemmatization always gives the dictionary meaning word while converting into root form.
- ⑤ Lemmatization would be recommended when the meaning of the word is important for analysis.
Example: Question Answer.
- ⑥ Example:
"studied" \Rightarrow "study"

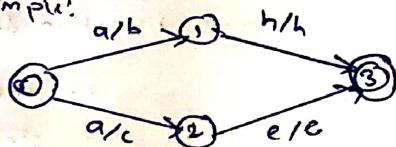
Regular Expression

- It is a language used for specifying text search string.
- Also called as Regex.
- Example:
 - ① `/[abcd]/` - will match a, b, c & d.
 - ② `/[3-6]/` - Specifies any one of the digit 3, 4, 5, 6.
 - ③ `/[A-Z]/` - Not an upper case letter.
[1] - Not.
 - ④ `*` or `+` - Occurrences

Finite Automata / Finite State Automata

Finite state Transducers.

- It is finite state machines with two steps.
 - ① Input tape
 - ② Output tape
- Unlike finite state automation which only has one tape.
- Example:



- FSA represents a set of strings.
eg. { walk, walks, walked }
- FST represents a set of pairs of strings. (input, output pair).
- FST is a multifunction device
 - ① Translator
 - ② Recognizer
 - ③ Generator
 - ④ Relator

N-Gram

- It is a continuous sequence of n items from a given sample of text or speech.
- The items can be letters, words.
- N-grams are collected from a text or speech corpus.

N-Gram language model.

- It predicts the probability of a given N-gram within any sequence of words in the language.
- A good N-gram model can predict the next word in the sentence.
- Example:
 - ① Unigram
("This", "article", "is", "on", "NLP")
 - ② Bigram
("The article", "article is", "is on", "on NLP")

N-gram for spelling correction.

- N-gram is used without a dictionary, this way employs to find in which position in the incorrect word the error occurs.
- If there is a special way to change the incorrect word so that it contains only correct n-gram, there is a correction.