

Artificial Intelligence :-

"Man" made + "Thinking Power".

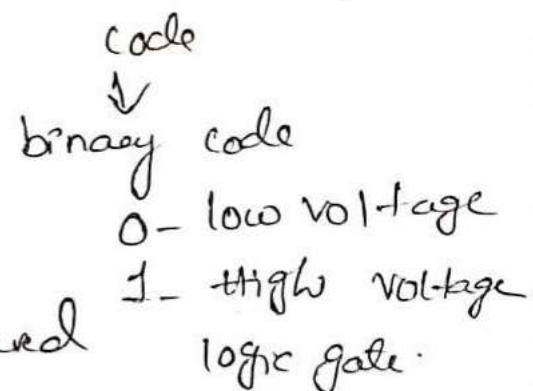
It is branch of computer science by which we can create intelligent machine, which can behave like a human, think like humans & able to make decision.

→ AI do not need to Pre Program machine to do some work.

→ You have to create a machine with programmed algorithm, which can work with own intelligence.

NEP :-

^{machine}
Computer (AI) knows only
binary code (0101...), they don't
know the human language that achieved
by using NLP



→ why we study natural language Processing?

→ NLP stands for natural language Processing.

→ It is a subfield of AI.

→ It deals with the interaction b/w humans & computers in natural language.

→ It uses computational techniques to process & analyse the data.

Applications:-

→ Question Answering.

Ex:- Alexa, Siri &

→ Spam detection.

Ex:- If you get a mail, that mail is danger. Spam.

→ Machine translation.

Ex:- convert to one language to another.

→ Speech Recognition.

Ex:- google assist, If you have YouTube, that video has no subtitle it will click **CC** it shows text.

→ Chatbot  chat box in website.

→ Sentiment analysis.

Ex:- Sad, happy, crying.

Ex:- I am enjoying right now,

Phases of NLP :- 5 phases in NLP

1.

Lexical analysis (or)
Morphological analysis

Lexeme :— [you giving sentence in
sub word / words] → Converting the given sentence
into stream of lexemes / morpho.
2.

Syntactic analysis

It is used to check grammar,
word arrangement and relationships
→ b/w words.
3.

Semantic analysis
(meaning of words / sentence /
phrases)

It is used to find the meaning
of words, phrases and sentences.
4.

Discourse integration
(group of sentence / paragraphs
which are integrated together
so that each sentence depends
on the previous sentence)

Ex:- (Ramu is a
boy) standard, (He
is a good boy) child
5.

Pragmatic Analysis

It first understands
what is stated & then
understands the intended
meaning of the sentence

①

words and their components:-

words - words are nothing but a meaningful unit of a sentence, by using the meaningful words only we construct the meaningful sentence.

→ we have 4 components.

1. tokens
2. lexemes
3. morphology
4. typology

tokens:-

If refers to a sequence of characters constitute that are treated as a single entity.

→ Tokens are words that are created by dividing the text into smaller units.

→ Process:- Tokenization. (By using tokenization we obtained the tokens)

(1) word tokens

(2) character tokens.

(3) subword tokens.

Ex:- I | love | reading | news papers.

↳ words tokens

↳ reading

↳ character tokens

news | paper

↳

↳ subword tokens

2. Lemmas:-

- They are the base (or) canonical form of words
- ex:- The base word of running, ran is 'run' (base word)
inflection/variation
- ex2:- largest, larger is large (base word).
- fundamental unit of meaning associated with a word, regardless of the inflection variations.

Process:- Lemmatization.

3. morphemes:-

- They are the smallest meaningful unit.
- Many words are formed by combining more than one morpheme.

Types:-

1. Free morpheme :- ex:- agree
 2. Bound morpheme :- ex:- dis^{agreement}
- ex:- di^{agreement}.

Process:- morphological process

4. Typology:-

- Typology refers to the classification or categorization of language based on their structural or grammatical features.
- Aim:- To identify common pattern & variation across language to understand their properties or relationships.

1. Isolating (or) Analytical language:-

- Have no (or) relatively few words that are composed of more than one morpheme.

ex:- Chinese, Thai

(In this kind of languages ~~there~~ there are very less words which are formed by one (or) more morphemes).

2. Synthetic languages:-

- It is opposite of 1st language.
- Combines more morphemes in one word.

3. Agglutinative language:-

morphemes associated with only a single function at a time. (Complex words to explain a simple topic).

ex:- Korean, Japanese, Tamil, Finnish.

4. Fusional languages:-

ex:- bat \leftarrow ^{animal} cricket bat.

→ feature (meaning) — Per-morpheme variation higher than one.

ex:- Arabic, Sanskrit, German etc.

①

Issues and challenges:-

There are 3 issues and challenges.

1. Irregularity
2. Ambiguity
3. Productivity

Irregularity:-

The phenomenon where certain words (~~or~~) words and forms does not follow regular patterns (~~or~~) rules involving ~~or~~ other morphology (~~or~~) syntax.

→ It is a challenge for algorithms which follow particular patterns.

(i) Irregular verbs & nouns:-

which does not follow standard pattern (~~or~~) inflection.

ex:-	Past	Present	Future	
	choose	chose	chosen	(choose)
	bite	bit	bitten	(bite)
	went			(go) Irregular words

(ii) Exceptional Inflection:-

Comparative and superlative adjectives.

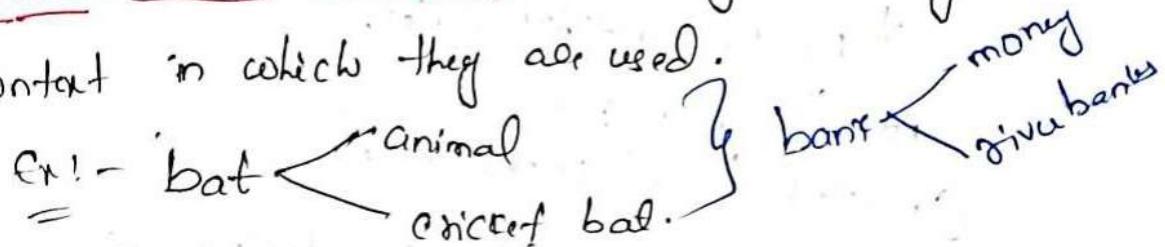
ex:- big, bigger, biggest

dark, darker, darkest

good, better, best Irregular words.

- ② Ambiguity:- (The word which is having multiple meanings; same spelling).
- word forms that can be understood in multiple ways out of context.
- word forms that look same but have distinct function (or) meaning (Homonyms).
- Ambiguity arises in morphological Processing and language Processing.

(a) word sense ambiguity:- Meaning depending on the context in which they are used.



(b) Parts of speech ambiguity:- Different parts of speech on their usage.

Ex:- I run (Verb)

He went for a run (noun).

(c) Structural Ambiguity:- multiple valid syntactic structure.

(d) Referential Ambiguity:-

Referring a Person by he/she
girl:- she/her, boy:- he/him.

③ Productivity:-

Ability to generate new words (or word forms) using Productivity rules. (It generates some new words by taking some unknown words).

e.g. According to wikiPedia,

- goog~~l~~ means 1 followed by 100 'g's
↓

From this word, an unknown word

[google] is generated.

→ By Productivity rules.

new words are
generated based on that word
i.e., googling, googlisch, googleology etc.

→ Names of People, organisation, location have unique structure where the Productivity issues arises.

① Morphological Models:-

- The morphological models used to analyze the structure and also formation of the words.
- It also helps us to identify the morphology in a particular sentence.

5 morphological models.

1. Dictionary lookup.
2. finite-state morphology.
3. Unification-based morphology.
4. Functional morphology.
5. Morphology Induction

(a) Dictionary lookup:- this used to analyse the structure of a particular word.

↓
word → base form (or) Canonical form
↓ e.g. - beautiful → beauty

Search in dictionary

↓
returning these information.

→ If you don't found the word in dictionary (or)
you see more topics in that particular word.

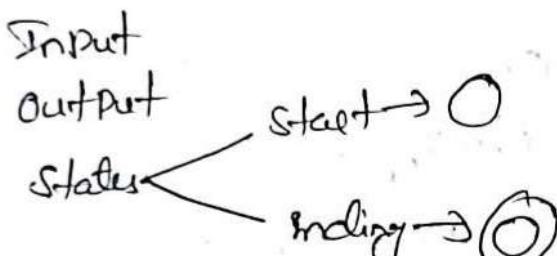
It will combine the another algorithm and shown
the result word.

(b) Finite-state morphology:-

→ Based on finite-automation and formal language theory.

→ used for generation & recognition tasks.

formal language theory:-

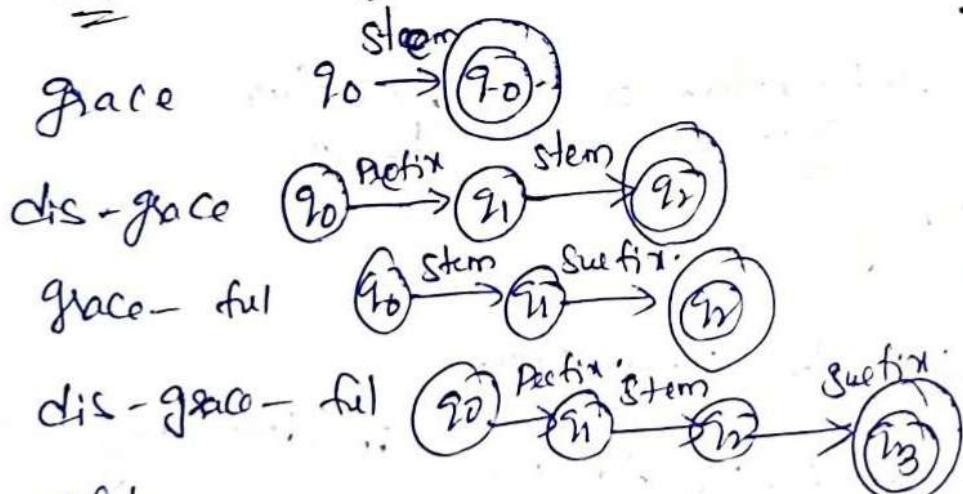


Transitions:— (→)

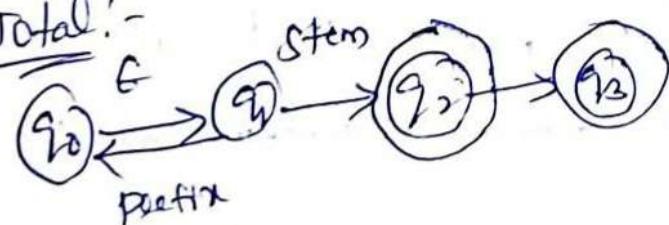
Finite-state transitions:-

This process is represented by FST's

Ex:- word-glace



Total:-



(c) Unification-based morphology:-

One word ^{which} is formed by combining these different words. That formed word should be meaningful.

$$\boxed{\text{word} = \text{word 1} + \text{word 2} + \text{word 3.}}$$

meaningful.

Eg:- window = word 1 + word 2 + word 3.

= DIY = Do + it + yourself.

HDTV = High + definition Television.

(d) Functional-based morphology:-

→ used to analyze the role of a morpheme in a word. (which is helpful in forming another word).

→ and also analyze their contribution in overall grammatical structure of a sentence.

Eg:- walk, walking,
= walk.
↓ word.

(e) morphological Induction:-

→ It is used to discover the pattern of a particular word (or) structures without helping in a any human being. The system only automatically detects.

→ By ~~detecting~~ preparing classical morphology

Finding the structure of documents:-

→ Extracting the structure of documents helps Natural language Processing tasks like Parsing, machine translation etc.

→ How?

- By chunking (Parsing) the input text (or) speech to blocks → Segmentation.
- Different approaches for different languages.

For example:-

Chinese documents → 1. Segment characters sequence into words

Morphologically rich language → 1. PreProcessing step
(Determines the tokens)
2. Apply algorithms.

→ In order to chunk the text, Segmentation is used.

Segmentation:- (Given text into some sentences).

→ Decides whether to mark a boundary b/w tokens as sentence boundary (or) topic

2 types.

- (i) Sentence-boundary detection.
- (ii) Token-boundary detection.

→ (i) Sentence-boundary detection:-

→ Also called Sentence Segmentation.

→ The process of segmenting sequence of words into units.

Sentence can be identified by:-

In written English:-

→ Beginning of the sentence (uppercase).

→ End of sentence (? , . , !)

But sometimes.

→ Capitalizations are ~~commonly~~ used for nouns and abbreviations.

→ Punctuation marks are used inside the sentences.

Example:-

I talked to Dr. Smith and my house is on mountain Dr.

(i) Optical character Recognition:- Segments the I/p text into sentence. These systems confuses with Periods and Commas.

→ It result in meaningful sentence.

Speaker character Recognition:-

(iii) Automatic Speech Recognition:- Segments the speech into sentence. These systems confuse with Punctuation marks (?, !, .)

→ Code switching Problem arises when segmentation done for other languages.

Example:- Spanish uses inverted question marks & exclaiming marks (¡ and ?)

→ Conventional rules-based Patterns are used to identify Potential ends of Sentence.

(ii) Topic boundary detection:-

→ It is also called as discourse segmentation/text Segmentation.

→ The process of dividing the speech into homogenous blocks is called topic Segmentation.

Applications:-

1) Text extraction and Retrieval.

2) Text Summarization.

Identification:-

(i) Text Segmentation clues:-

→ by following the headlines.

→ By Paragraph breaks.

(ii) Speech Segmentation clues:-

→ Pause duration

→ Speaker changes.

Methods:-

1) Generative Sequence classification method:-

Observation :- words, punctuation. (? , !)

Labels :- Sentence boundary, topic boundary.

→ These methods have models which are not only capable of predicting the most likely label but also generating the sequence itself based on learned probabilities.

One example of this method is Hidden marker.

How HMM ~~works~~ works?

- 1) Learn from data about the observations and its corresponding hidden state (POS).
- 2) Predict the labels (or) sequence generation is done.
- 3) classify the new sequence.

Example:-

Sentence 1:- 'I / PRON Love / VERB coding / Noun'.

Sentence 2:- "the / DET quick / ADJ brown / ADJ, for
Noun jumps / VERB".

Sentence 3:- 'She / PRON sells / VERB sandwich / Noun'.

→ Based on these 3 sentences sequence generation and classification is done.

2) Discriminative Local classification method:-

Local features:- word Identifies, Prefix, Suffixes & nearby Pos (Part of speech).

→ These approaches focuses on making decision based on local information, typically without considering the entire sequence of data.

→ These methods directly model the relationship b/w input features and output labels.

Example of these method be maximum Entropy model (or) support vector machine.

Applications:-

- 1) Parts-of-Speech Tagging. (It Predict the ~~Particular~~ at the Particular Segments)
- 2) Speech Recognition. (Ex:- YouTube)
- 3) Named Entity Recognition. (Person name, location name etc).

Complexity of the approaches:-

- * Whether comes to Complexity / one approach may be better than other.
- * The complexity can be measured by rating their training and Prediction.
- * ALSO measure their Performance.

Performances of the APProaches:-

$$\text{Error \%} = \frac{\text{no. of errors}}{\text{no. of examples}} \cdot$$

$$\text{Precision} = \frac{\text{total no. of correct Positive Predictions}}{\text{total no. of Positive Predictions.}}$$

$$\text{Recall} = \frac{\text{total no. of all correct Positive Prediction}}{\text{total no. of Positive instances.}}$$

$$f1\text{ score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall.}}$$