

Natural Language Processing (NLP) (Day-1) 12-08-2023

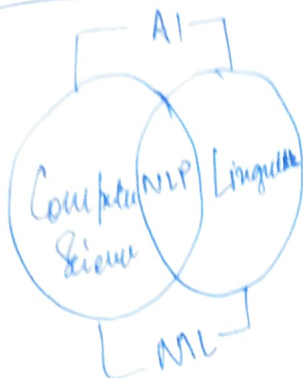
Meenakshi.

UG - VIT

PG -

PHD - Research in ML, DL

CV, NLP Advancing ML



→ NLP deals with only text data, emotion detection is not considered.

→ Phases in NLP.

* Lexical Analysis → Converting the sentences into tokens like individual words.
→ Morphology.
→ Formulation of words.

* Syntactic Analysis → whether the language is following proper rules of the language grammar.

* Semantic Analysis → Understanding the meaning.

* Discourse Analysis → Analysing text in the paragraph level.

* Pragmatic Analysis → Let the system understand the general day to day conversation.
ex:- I drink Taj → Taj is in context denoting Taj tea not Taj mahal, Taj hotel etc.

Applications:-

→ Some of the applications.

→ Check Credit Analysis, Worthiness, Language Translation, Sentiment Analysis, Customer Support, Work Routing, Identify Similar legal cases. (Document similarity), Information Retrieval, Information Extraction, Q&A Chatbot, Optical Character Recognition.
Text categorization (a classification), Word prediction (Auto completion), Speech recognition (Alexa), Machine Translation (Google Translate).

→ Ambiguity (Challenges of NLP) + environment

word line.

ex. - Bank, Can, Bat.

* Sentence Incl.

Ques: The man saw the girl with the telescope.

- We saw his duck.

- * Pragmatic level.

Ques. Tourist (checking out of the hotel): Waiter, go upstairs to my room and see if my sandals are there; do not be late; I have to catch ~~the~~ train in 45 minutes.

Qa. Waiter (nearly upstairs coming back downstairs): Yes sir, they are there.

* Discourse level.

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

→ Features in lexical programming -

→ NLP is based on:

* Prob. & Stats.

* ML

1. Linguistics -

★ Common Cause -

→ Standard terms.

- * Corpus: A body of text samples.

* Document: A text sample.

• Vocabulary: A list of words used in the corpus.

- Vocabulary: A list of words with their meanings.
- Language model: How the words are supposed to be organised.

- Language model: How the words are supposed to begin

→ NER → Named Entity Recognition.

Ex: Sports → Cricket
→ Football
→ Tennis } Entities.

→ NLP $\begin{cases} \rightarrow \text{NLU} \rightarrow \text{Nat. Lang. Understanding.} \\ \rightarrow \text{NLG} \rightarrow \text{Nat. Lang. Generation} \end{cases}$

→ GAN → Generative Adversarial Network.

→ Regular Expression ::

* Pattern matching.

* Purpose is to find the pattern within the corpus.
Regex.

Uses

→ Web-search,

→ Retrieval applications.

→ Word-processing, computation of frequencies.

→ Basic Regular Patterns.

→ Simple text search.

RE.

`/text/`

→ It is case-sensitive.

→ Disjunction of characters -
Syntax.

`/[word]/`

→ Range of characters.

Syntax `/[a-z]/`

→ Caret - Negation.

→ The caret ^ is the first symbol after the square bracket.

→ !

→ preceding character is ignored.

Ex color and colour.

color?r

↓
ignored.

→ Kleene Operator.

* The Kleene * means "Zero or more occurrences of the immediately previous character or regular expression.

→ The Kleene + means "One or more occurrences

→ Wild Card expressions.

→ Syntax /./

ex /beg.n/

If we need to find only "." we can use /\. / or

/[.]/

→ Anchor:

2 anchor expressions -

→ Caret - / ^ / → specifies start of the sentence / corpus -

→ Dollar - / \$ / → specifies end of the sentence.

Word Boundary

→ \b matches a word boundary. while \B matches a non-boundary.

→ How to specify for ex: guppy and guppies. use "|" operator.

ex: /guppy/

/guppy|ies/

→ Repetition.

ex: [a-z]{3}

→ Advanced operators.

\d (0-9)

\D (^0-9)

\w

\W

\s (\r \t \n \f)

\S (^ \s)

→ ex: questions.

Check the url like name@gmail.com.

Spell checking:

→ ex: piece
peace

} Two types of errors.

① word error. → when, peace is written

② Non-word error. instead of piece.

word error.

→ Replace.

→ Insert.

→ Removal.

ex:

tutor → misspell.

tumour → correct.

2 solutions.

1st } Replace t → m
1 Insertion u } optimal.

2nd } Removal t
1 Insertion m
1 Insertion u } General.

Dynamic.
Greedy
Backtracking

Minimum Edit Distance.

Ex. Misspelled a e c d g

Original a b c f g

 N a b c f g
N
a
e
c
d
g