

7 Jan '25

11

NATURAL LANGUAGE PROCESSING

- subfield of AI that uses ML to understand & recognise human language.
- enables digital system to understand / recog and also generate text & speech by using computational linguistics and statistics / ML / DL.
 - ↳ techniques / rule to understand language.

• Generative AI

- LLM / Img recog or ability to understand prompts (LSTM)

• Benefits

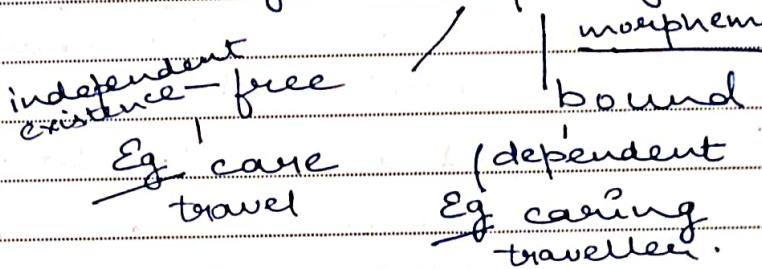
- powerful content generation
- automation of repetitive tasks (+ customer support, doc. classification)
- enhanced speech search
- improved data analysis: recommendation / analysis

• Levels of NLP.

(1b)

(1c)

- 1) Lexical Analysis - Phonological / Morphological
- 2) Syntactic
- 3) Semantic
- 4) Discourse Integration
- 5) Pragmatic.



1 a) Tokenization

→ Stop words

identify common
insignificant words

Eg: am, is, but, etc.

~~break word to manageable chunks~~

~~Lemmatization~~

scan / segregate text into sentence / para / word

~~big short karma~~ → manageably token. → to understand text at the word level.

better.

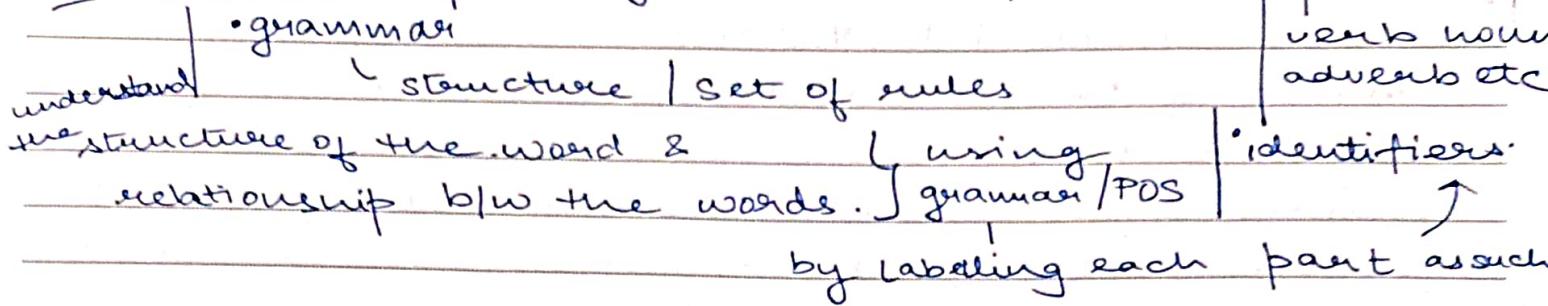
~~takes Lemmatization | Stemming~~ → does not check relevance.

~~Caring~~
~~break karke.~~

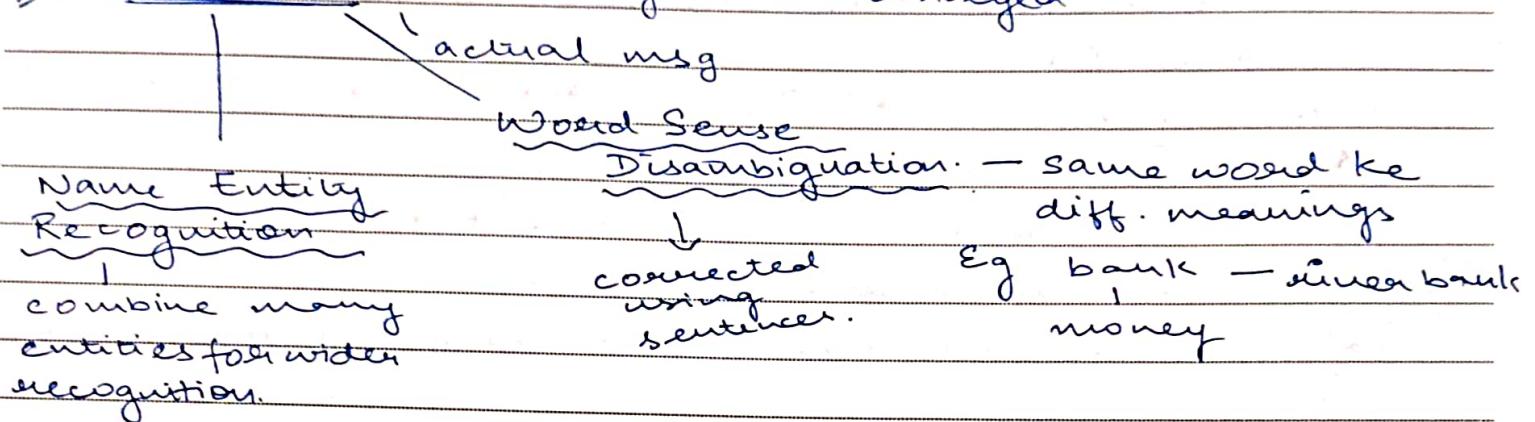
~~care~~ banao.

↳ Caring se break karke ~~Car~~
closest root value. ↓ irrelevant.

2 Syntactic : - parsing ; - Part of Speech tags (POS)



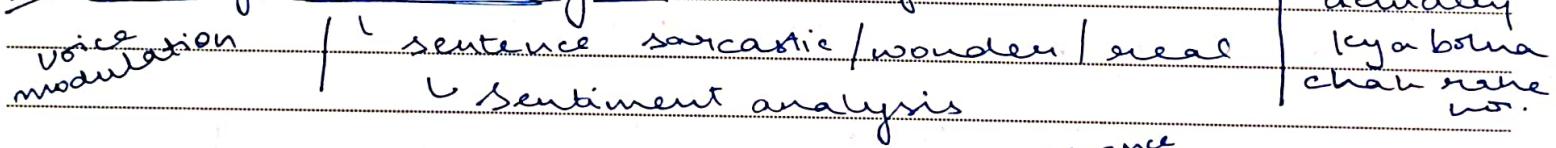
3 Semantic : meaning to be conveyed



4 Discourse Integration => Sequence.

- each sentence has some significance / relevance
 - Contextual Reference → same type ki cheez ki.
 - Anaphora Resolution
 - persons type ka
 - tere sentence mei NLP bola.
 - baadmehi, this subject boldo
- NLP = DS ✓
NLP = history ✗

5 Pragmatic Analysis - change in tone.



Feature Extraction

1) Term frequency (TF)

No. of total words in doc

reduce dimension.

2) Inverse Doc. frequency (IDF) — rarity of word in corpus

3) TF - IDF — statistical measure.

$$L \boxed{TF \times IDF}$$

significance of a word DOMS and rarity of word in doc

Corpus - D_1 : Natural Lang
 D_2 : Lang Processing
 D_3 : Natural Lang. Processing.] combination

Vocabulary - $\#1$: Natural] ek ek karte
 $\#2$: Lang]
 $\#3$: Processing]

$D_n \rightarrow$ Documents

corpus is a collection of documents

Term frequency	D_1	D_2	D_3
of Natural	$1/2$	0	$1/3$
Lang	$1/2$	$1/2$	$1/3$
Processing	0	$1/2$	$1/3$

IDF : entire corpus : same documents

$$IDF = \log \left(\frac{\text{Total no. of doc in corpus.}}{\text{Total no. of appearance of a word across all doc}} \right)$$

$$\text{Natural} - \log \left(\frac{3}{2} \right)$$

$$\text{Language} - \log \left(\frac{3}{3} \right)$$

$$\text{Processing} - \log \left(\frac{3}{2} \right)$$

3 doc. { Total D_1, D_2
 D_3

9/01/25

TF-IDF

	D_1	D_2	D_3
of Natural	$\log \left(\frac{3}{2} \right) \times \frac{1}{2}$	$\log \left(\frac{3}{2} \right) \times 0$	$\log \left(\frac{3}{2} \right) \times \frac{1}{3}$
Language	$\log \left(\frac{3}{3} \right) \times \frac{1}{2}$	$\log \left(\frac{3}{3} \right) \times \frac{1}{2}$	$\log \left(\frac{3}{3} \right) \times \frac{1}{3}$
Processing	$\log \left(\frac{3}{2} \right) \times 0$	$\log \left(\frac{3}{2} \right) \times \frac{1}{2}$	$\log \left(\frac{3}{2} \right) \times \frac{1}{3}$

DOMS

N-Grams - how many words together: contiguous sequence
 ↳ $N=4, 5$: most effective.

Eg NLP

$$N=1 \Rightarrow \{N, L, P\}$$

$$N=2 \Rightarrow \{NL, LP\}$$

$$N=3 \Rightarrow \{NLP\}$$

Eg Find the probability of $\langle S \rangle I I \text{ am } \langle /S \rangle$ in the given corpus

1: $\langle S \rangle I$ am a human $\langle /S \rangle$

2: $\langle S \rangle I$ am not a robot $\langle /S \rangle$

3: $\langle S \rangle I$ am not $\langle /S \rangle$

Prob of $(\frac{\langle S \rangle \text{ and } I \text{ together}}{\text{just } \langle S \rangle})$

$$\Rightarrow P(I | \langle S \rangle) \times P(I | I) \times P(\text{am} | I) \times P(I \text{ and } \text{am} \text{ tog} | \langle S \rangle)$$

Eg $\langle S \rangle I$

$P(\frac{I \text{ and } I \text{ tog}}{\text{just } I})$

$P(\langle /S \rangle / \text{am})$

=> Count

$$\frac{3}{3} \times \frac{0}{3} \times \frac{3}{3} \times \frac{0}{3} = 0$$

Q Predict: not _____?

$$\Rightarrow P(\text{not } / \text{not})$$

$$P(I | \langle S \rangle) \Rightarrow \frac{\langle S \rangle \text{ and } I \text{ tog}}{\text{just } \langle S \rangle} = \frac{3}{3} = 1$$

$$P(\text{am} | I) \Rightarrow \frac{I \text{ and am}}{\text{just } I} = \frac{3}{3} = 1$$

$$P(a | \text{am}) \Rightarrow \frac{\text{am and a}}{\text{am}} = \frac{1}{3}$$

$$P(\text{human} | a) \Rightarrow \frac{a \text{ and human}}{a} = \frac{1}{2}$$

$$P(\text{not} | \text{am}) \Rightarrow \frac{\text{am and not}}{\text{am}} = \frac{2}{3}$$

$$* P(a | \text{not}) \Rightarrow \frac{\text{not and a}}{\text{not}} = \frac{1}{2}$$

$$P(\text{robot} | a) \Rightarrow \frac{a \text{ and robot}}{a} = \frac{1}{2}$$

$$P(\langle /S \rangle | \text{human}) \Rightarrow \frac{\text{human } \langle /S \rangle}{\text{human}} = \frac{1}{1} = 1$$

$$P(\langle /S \rangle | \text{robot}) \Rightarrow \frac{\text{robot } \langle /S \rangle}{\text{robot}} = \frac{1}{1} = 1$$

$$* P(\langle /S \rangle | \text{not}) \Rightarrow \frac{\text{not } \langle /S \rangle}{\text{not}} = \frac{1}{2} = \frac{1}{2}$$

$$P(a | \text{not}) = P(\langle /S \rangle | \text{not})$$

end of sentence

DOMS

not considered

against occurrence
 along path
 along prob
 jiski vali
 syada consider

So, ans: a

13 Jan 2025, Transformation based learning by Rule + Stochastic Syntactic level of NLP - noun verbs | etc.

Parts of Speech Tagging. ① Dissemination of the Word Class

- I am about to book a flight.

- I like to read book.

open

closed

② Basis

Rule based

Stochastic based

Noun

Adverb

Adjective

Verb

Preposition

Pronoun

Determinant

→ fixed

→ appear

frequently
in a sentence

→ shorten words

Word frequency measurement

Tag sequence probability

Rule based:

- When a word is preceded by a determinant and followed by a noun: the word is an adjective
- higher time complexity, inefficient, inaccurate

Stochastic based: better based. I like to play →]
better based. I like to play →]
acket

Word frequency

counting mere occurrence of word

8 → verb 80%

2 → noun 20%

Ten $\frac{8}{10}$ means ki out of 10 sentences, 8 mei (play) is a verb, $\Rightarrow 80\%$. 2 mei (play) is a noun $\Rightarrow 20\%$.

PVP 2/V

syada
reliable

10

play → V

Tag Sequence

look for patterns
 $w_1 \ w_2 \ w_3 \ w_4$

Identifying
a pattern

ki agar most sentences mei pronoun - verb - prepo - verb
hai, ton (PVPP) DOMS
ke bare fine V ni aayega.

~~15 Jan 2025~~

decision making factors /

Transformation Based Learning

Rule +

Stochastic

Training Data

- i/p Tagged corpus
- list of dictionary words with their most frequent tags.

o/p

- sequence of transformation rules
- accuracy

↓
change the tags

not used
a lot ← Underflow condition → negligible data

~~20 Jan 2025~~

Smoothing

used when scarce data
or missing info

(1) Log based

(2) Laplace

(3) Add - k / Lidstone / Additive

→ redistribute vocabulary massed

Training : Natural Language Processing

$$P(\cdot) = \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}$$

Testing : Natural Understanding

$$P(\cdot) = \frac{1}{2} \times \frac{0}{2} = \underline{0}$$

$$\text{Vocabulary} = \{N, L, P, U\}$$

dis
adv
ant
age

toyota yes

This is why
we do
Smoothing

still probability
gives 0.
DOMS

"Understanding" is
in vocabulary

Laplace Smoothing:

~~Explanation~~: add 1 to numerator

~~Explanation~~: add total no. of words to denominator

$$\frac{1}{2} \times \frac{0}{2} \text{ becomes } \frac{1+1}{2+4} \times \frac{0+1}{2+4}$$

How to actually solve:

~~Final Rule~~: $P(N) \times P(U)$

↓ ↓
total occurrences + 1
of N is

no. of words no. of
in testing & + words in
training vocabulary.

} final / count the
word in training.
if (training = true)
check testing
else 0

So,

Training : Natural Lang Processing

Testing : Natural Understanding

Vocab = {N, L, P, U}

$$P(N) \times P(U)$$

$$\Rightarrow \left(\frac{2+1}{5+4} \right) \times \left(\frac{0+1}{5+4} \right)$$

$$= \frac{3}{9} \times \frac{1}{9} = \frac{3}{81} = \frac{1}{27}$$

} if word is
in testing
training,
check in
testing

if word not
in training
then sehra 0

22/01/25

more complex than Laplace

Add K / Lidstone

keeps on changing

NATURAL LANGUAGE PROCESSING

count(N)

$$\textcircled{1} + 0.5$$

$$1 + 0.5$$

$$1 + \textcircled{0.5}$$

$$\textcircled{3} + (0.5 \times 3)$$

$$3 + (0.5 \times 3)$$

$$3 + (0.5 \times \textcircled{3})$$

Total words

3 vocabulary

\propto or K

Good Turing Smoothing

Count

N

$c = 2$

Natural

2

$\frac{2}{2} \{ N_2 \}$

$[N, L]$

Language

2

$\frac{2}{2}$

Processing

1

$\frac{1}{N_1}$

Understanding

0

$N_0 = 9$

$P = 1 \left[\frac{2}{P} \right]$

5

$N_2 :$

$[N, L]$

no. of elem 2
jisme N, L ke
'N' column mei
2 aajayega

N_1

$[P]$

no. of elem = 1
jisme P ke
'N' column mei
1 aajayega

No

$[U]$

no. of elem
= 0

$$C^* = (C+1)^* \frac{N_{c+1}}{N_c}$$

$$P^*() = \frac{C^*}{N} \quad \text{total no. of words}$$

$$P(\text{museum}) = \frac{N_1}{N^*}$$

DOMS

LL

	N_2 [N, L]	N_1 [P]	N_0 [U]
Count	$\frac{2}{5}$	$\frac{1}{5}$	0

$$C^* = \frac{(2+1)^* N_3}{N_2} \quad | \quad C^* = (1+1)^* N_2 \quad | \quad C^* = (1)^* \left(\frac{N_1}{N_{10}} \right)$$

\downarrow

$\because \text{no } N_3$

$\Rightarrow \text{can't calculate}$

$$: 2^* \frac{2}{1}$$

\downarrow

$: 4$

$$\cancel{C^* = (1)^* \left(\frac{N_1}{N_{10}} \right)}$$

$$\cancel{= 1^* 1}$$

unseen probability

$$P^* = \frac{C^*}{N} \times \quad | \quad P^* = \frac{4}{5} \quad | \quad P^* = \frac{1}{5}$$

can't calculate

\downarrow

$\text{total no. of words}$

\downarrow

$N_1 = 1$

\downarrow

from $P(\text{unseen})$ vala formula.

agar 2 unseen words, then divide by 2

23/01/25 Regular Expression.

- ab^* → b is n no. of times
- $abc?$ → c can occur 0 or 1 times
- $abc\{2,3\}$ → c can occur max 2 times
- $abc\{2,3\}$ → c's range starts with 2
- abc^+ → 1 or n no. of c
- $^a a$ → string starts with a
- $b\$$ → string ends with b
- a^b → a and b
- $^A (A - Z)$ → Mean A acts as negation,
 $A - Z$ range is not at all present
(a - z) present -

Example questions

- Phone Number (starts with 8 or 9)

$[8\ 9]\ [0 - 9]\{9\}$

Phone number in format

XX - XX - XXXX

$[8\ 9]\{1\} [0 - 9]\{1\} [\backslash -]$

$[0 - 9]\{2\} [\backslash \backslash -]$

$[0 - 9]\{4\}$ ↓
escape character

(\) ke baad jo likhoge
as it is print ho Jayega.

Email

$[a - z\ A - Z\ 0 - 9\ \backslash -\ \backslash +\ \backslash \dots]$

$+ [\backslash @]\ [a - z] + [\backslash \cdot]$

$[a - z]\{2, 3\}$

27/01/25

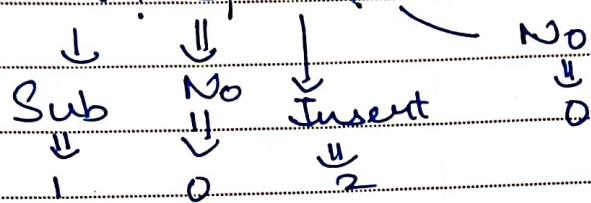
EDIT DISTANCE minimize the cost.

3 operations

- Insert → 2
- Delete → 2
- Substitution → 1

$S_1 : a | b | c$

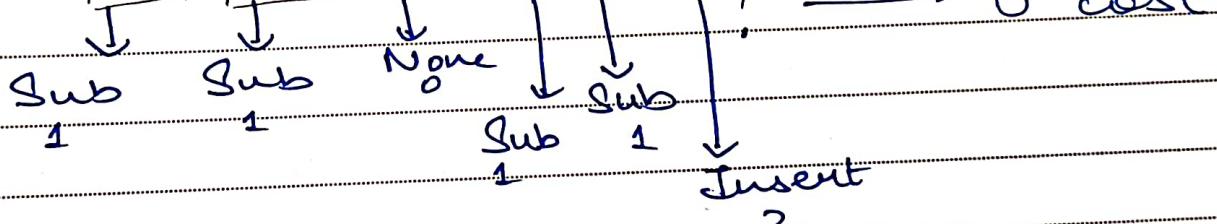
$S_2 : d | b | e | c$



⇒ 3 cost

$S_1 : I | N | T | E | N | C | T | I | O | N$

$S_2 : E | X | T | R | A | C | T | I | O | N$



⇒ Cost = 6