

UNIT - 1

④ NLP :- Natural Language processing.

Natural Language processing is the subfield of Artificial Intelligence (AI) focused on enabling machines to understand, interpret, and generate human language. chatbots, search engines, translation apps. (like Google translation).

Core challenges of NLP :-

- * Ambiguity :-
 - * Lexical Ambiguity :- words with multiple meanings (bank as a financial institution vs. riverbank).
 - * Syntactic Ambiguity :- sentences with multiple valid grammatical interpretations.
(I saw the man with the telescope.)
 - * Semantic Ambiguity :- sentences whose meaning changes based on context
(The chicken is ready to eat.)
 - * Contextual understandings :- difficulty in understanding context-dependent meanings, especially in long conversations (or) documents

Pragmatics: Recognizing implied meanings.
Sarcasm, and humor remains challenging
for NLP models.

- * Low-Resource Language:
NLP systems often lack sufficient
training data for Low-resource.
- * Ethical and Bias Concerns:
 - * Models inherit biases from training
datasets, potentially leading to unfair
(or) prejudiced outputs.

Future of NLP!

- 1) Zero-shot and Few-shot Learning:
 - * Models capable of performing tasks with minimal (or) no labeled data.
- 2) Cross-Lingual NLP:
 - * seamless functioning of models across multiple languages without retraining.
- 3) Bias mitigation:
 - * strategies for identifying and
addressing dataset biases.

- 4) Integration with other Modalities;
→ combining NLP with computer vision and audio processing for enhanced applications.

Tab
→

⑩ Language Modelling (LM) VIM -

- * Language Modelling is the process of predicting the next word in a sequence of given a set of preceding words.
- It is the fundamental task in NLP.
- serves as back bone for applications like speech recognition, machine translation and text prediction.

Applications of Language Modelling (LM).

- * ① speech Recognition : converting spoken words into a text.
- * ② Machine Translation : Generating accurate translations.
- * ③ chat bots : creating human like response
- * ④ spell checking : Identifying and correcting errors & grammarly

Types of Language Models

1. Grammer Based models.

2. statistical Models.

* Grammer Based models.

* A Grammer Based model Follows prede
fined grammer rules. These rules tell the
model how to structure sentences correctly.

sentence structure.

subject + verb + object.

Ex: The cat eats fish ✓

she loves chocolate ✓

Advantages :

① * High Accuracy in Grammer.

* since it follows strict rules, it
ensures that sentences are grammatically
correct.

Ex: She eats an apple. ✓

She eating apple ✗

*② Good for structured Languages

- * works well for language with clear grammar rules.

Ex: programming Languages. [c, c++, java, python]

*③ useful for Language Learning & grammar checking.

- * Helps Learners by correcting grammar mistakes.

Ex: Grammarly uses grammar rules to suggest correct sentence structure.

*④ predictable and reliable.

- * Since rules are predefined, it always follows them reducing random errors.

Ex: In legal documents, where exact wording is important, it ensures correct net.

Limitations of Grammer - Based Models.

* ① Too Rigid (not Flexible)

- * Language is always evolving and fixed rules don't always work for informal speech.

Ex :- Gonna grab a coffee.

(casual but correct in real life,
yet a grammar model may flag it
as wrong.)

* ② Needs manual Rule creation.

- * Humans must write and update grammar. rules manually, which is time consuming.

Ex :- Adding new slang ^(ex) cool phrases requires extra work.

* ③ Cannot Handle complex sentences well

- * some sentences follow the rule but still sound awkward.

Ex: The dog which barks loudly is outside.
(Technically correct but unnatural - better)

as1 " The loudly barking dog is outside.

Grammatically correct The loudly barking dog is outside.

* ④ struggles with meaning in context.

- * It only checks if the structure is correct but doesn't understand the actual meaning.

* The cat ate the banana.

(Grammatically correct but makes little sense.)

* Parsing Techniques in Grammars - Based LMs.

- * Parsing in Grammar - Based models (LMs) refers to analyzing a sentence's structure based on grammatical rules.

* Top-down parsing

- * start with main rule (e.g. sentence) and breaks it into smaller parts (subject, verb, object etc.).

Ex: If the rule is Sentence \rightarrow Noun phrase + Verb phrase,

i.e first assumes this structure and tries to match the words.

* Bottom-up parsing

- * Begins with individual words and builds up to larger structure

Ex: It first identifies 'dog' as a noun 'runs' as a verb then recognizes the phrase "The dog runs".

CYK Algorithm (Cocke-Younger-Kasami).

- * A dynamic programming algorithm used for parsing context-free grammars efficiently.
- * Advantages: guarantees optimal parsing for CFG's.

* Parsing Techniques in Grammar-Based LMs

| Aspect | Grammar-Based LM | Statistical LM |
|--------------------|-----------------------|---------------------------|
| Approach | Rule-based | Data-driven |
| Interpretability | Highly interpretable. | Less interpretable. |
| Flexibility | Rigid structures | Adaptive to data. |
| Handling ambiguity | Limited | Better ambiguity handling |
| Scalability | Difficult to scale. | Scales well. |

Applications of Grammar-based LM

1. Voice-controlled system :- Parsing restricted voice commands ("Turn off the fan").
2. Text-to-speech system :- ensuring grammatically accurate speech synthesis.
3. Education Tools :- Grammar-checking software to guide language learners.
4. Compiler Design :- validating syntax in programming languages.

statistical Language modeling (SLM)

- + SLM is a way of predicting the next word in a sentence based on probability.
- * It helps computers understand and generate human like text.
- * Unigram model (one-word at a time).
 - * Assumes that each word is independent of the previous ones.

Ex: I love NLP

I, love NLP.

* Bigram model

- * considers one preceding word to predict the next word.

(word + next word)

I love pizza,

I love pasta-

user-type I it check most common words comes after "i" most e.g. love,

③ Trigram model.

- * Looks at groups of three words (word + next two words).

Ex: If a chatbot sees "I love pizza" many times, it learns that after "I love", "pizza" is likely.

- ↳ Better than unigram model.

Neural Language models (NLMs)

- * uses neural networks (complex mathematical functions) to understand entire sentences instead of just word pairs (or) triplets.

Ex: If a user type "I want to eat", a neural model (like GPT or BERT). can predict "a burger" (or) some food based on deep context

probabilistic approaches in statistical language modeling.

PASLM help predict the next word in a sequence in sentence using probability. These approaches assign a likelihood to each possible word based on past observation.

* smoothing Techniques:

To handle zero probabilities (unseen word sequences), smoothing techniques are applied.

1. Add-one (Laplace) smoothing:

Adds 1 to all word counts.

2. Good-Turing Discounting:

Adjusts probabilities for unseen events.

3. Kneser-Ney smoothing:

Redistributes probabilities effectively across N-Gram levels

Backoff : If a higher-order N-gram probability is unavailable, fall back to lower-order probabilities.

Interpolation : combine probabilities from different N-gram levels for better predictions.

Evaluation Metrics for Language models

When we build a language model (LM), we need to check how good it is at predicting text. Evaluation metrics help measure the model's performance.

* Perplexity (PPL)

* Measures how well the model predicts a sentence.

* Lower is better (a perfect model has low perplexity).

Cross-Entropy

- * Measures the average bits required to encode a message.

challenges in statistical LM.

* Data Sparsity:

- * some words combination appear very rarely, so the model doesn't learn them well.

* Out of Vocabulary:

- * If the model has never seen a word before, it doesn't know how to handle it.

* Eg: context Length

- * Traditional models (like n-grams) only look at a few previous words and can't understand long sentences properly.

computational overhead :-

More Complexity , more Resources .

Ambiguity :- Multiple meanings are hard.

Applications in Statistical LM.

- * Search Engines :- Enhancing query relevance.
- * Text Generation :- creating context aware text content .

Finite - state Automata (FSA)

A Finite state Automaton (FSA) is like a simple machine that moves through different states based on inputs. It is used to recognize patterns in text, speech, and programming.

Key components of FSA

1. States → A finite set of states.
2. Alphabet → A finite set of input symbols.
3. start state → The initial state.
4. Accepting states → states that signify valid string recognition.
5. Transition Function : Rules defining transitions between states.

English Morphology

* Morphology is the study of how words are formed and structured in a language.

* In English, words are made up of

smaller meaningful parts called morphemes.

cut, table

to, she, the

lexical

functional

derivational unhappy

inflectional in walking.

Morphemes

free

Bound.

Free : Can stand alone as independent

words.

Ex: cat, run, tree.

Bound : Cannot stand alone, they need

to be attached to another word.

Ex: plural cats, dogs.

negation unhappy, unclear.

1) Lexical words - Meaningful words.

These words carry the main meaning of a sentence.

Ex: nouns → (dog, car, happiness)

verb → (run, eat, jump)

Adjectives → (beautiful, big, fast)

Adverbs → (quickly, softly, yesterday)

The dog runs quickly in the park.

carry the main meaning → Lexical words.

2) Functional words - Grammatical helpers.

* These words don't carry meaning but help connect and structure sentences.

* It shows the relationships between lexical words.

The dog is running in the park

is in the → Functional words

* Inflectional Morphology (^{changes grammar}_{not meaning}).

changes: tense, number, (or) comparison
but doesn't create a new word.

Ex: play → played. (past tense)

big → bigger (comparison)

dog → dogs (plural)

Derivational (creates new words).

* changes the meaning (or) word type
(noun → verb, verb → adjective).

Ex: "happy" → "unhappy" (meaning changed).

"teach" → "teacher" (verb → noun)

"act" → "action" (verb → noun)

Morphological Parsing

* Morphological parsing is the process of analyzing words into their morphemes to extract structure and meaning.

steps in morphological

1. Tokenization : Breaking text into individual words.
2. morphological segmentation : splitting words into morphemes.
3. Morpheme Tagging : Identifying each morpheme's role. (e.g. root, prefix, suffix)
4. Reconstructing Meaning : combining the meanings of morphemes.

Tools for Morphological parsing

- * Porter stemmer :- Reduces words to their root form.
- * Lancaster stemmer :- Aggressive stemming algorithm.
- 3 spaCy :- NLP library for morphological analysis.

challenges in morphological \Rightarrow analysis

- * Ambiguity :- Words with multiple meanings
- * complex morphology :- Highly inflected languages (e.g. Finnish, Turkish)
- * compound words :- Difficulty breaking down compound words (e.g. notebook)
- * Irregular forms :- Handling irregular verbs and plurals (go \rightarrow went.)

Low - Resource Language : Limited linguist
ic resources for less common languages

Transducers for Lexicon and Rules

1. Transducer : A Transducer is a type of finite-state machine. that takes input and produces output.

Ex Input : "play" → output "played".
Input : "jump" → output "jumped".

Lexicon Rules ?

Lexicon rules define how words change.
(like verb tense, plural forms,).

* Rule Based Transduction

* If a word ends in "y", replace it with "ied" (e.g., carry - carried)

* If a word ends in "e" add "d"
(like → "liked")

Finite state Transducer (FST)

- * similar to a Finite state Automatos (FSA) but with two layers.
 - * one layer read input (e.g "run").
 - * one layer gives output (e.g "running").
- + it follows predefined rules to transform one word into another.

Input. "play" → output "played"

Input "jump" → output "jumped"

Applications of FST

- ① Morphological Analysis (word structure processing).

- * Helps break words into roots and suffix.

Input = running → output runing
(root + suffix).

used in :- spell checkers, grammars

tools.

Speech Recognition (converting sound to text)

- * converts spoken words into text format

Ex : Input : voice command "call mom"

Output : written txt "call mom!"

used in : Siri, Alexa, Google Assistant

ant.

Machine Translation (language conversion).

- * helps convert words between different languages

Ex : dog (English) → output "perro" (Spanish)

used in : Google Translate.

Text-to-Speech systems (TTS)

- * converts written text into spoken words.

Ex : "Hello, how are you?" → Input.

Output : spoken words in AI voice.

used in : Audiobooks, Screen readers.