

## Unit -1

finding the structure of words

- words and their components
- Issues and challenges
- Morphological models

finding the structure of Documents

- Introduction
- methods
- complexity of the approaches
- performances of the approaches
- features.

## Definition of NLP

Natural language processing is a cross disciplinary field of linguistic, computer science and NLP. It is related to interaction, b/w digitalized computing devices and human languages.

The NLP deals with designing and programming digital computational devices to process and analysis large amount of natural language data.

Natural language processing (NLP) is a field of AI that enables machine to understand, interpret, and generate human language through computational techniques.

## Applications of NLP

1. spam detection
2. machine translation
3. speech recognition
4. chat bot
5. sentiment analysis

## Phases of NLP:-

1. lexical analysis / morphological analysis
2. syntactic analysis
3. semantic analysis
4. discourse Integration
5. pragmatic analysis.

## 1. Lexical Analysis

- breaking down text into smaller units like paragraphs, sentences and words (tokenization)
- part learning techniques such as lemmatization, stop words removal, and correcting misspelled words.

- Identify morphemes, the smallest units of meaning in a word.

## 2. Syntactic Analysis (Parsing).

- understanding the structure of a sentence and its grammatical correctness
- Assigning parts-of-speech (POS) tags to each word (e.g., nouns, verbs)
- ensuring logical consistency and reducing ambiguity.

## 3. Semantic Analysis:-

- extracting the meaning of text.
- Understanding dictionary definitions and usage in context.
- Tasks like Named Entity Recognition (NER) and word sense Disambiguation.

(NER)

## 4. Discourse Integration

- comprehending relationships between the current sentence and previous sentences or larger context.

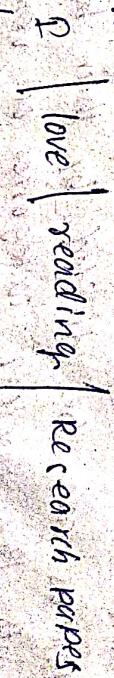
- Contextual references and anaphora resolution (e.g. understanding pronouns).

## 5. Pragmatic Analysis

- Interpreting the inferred meaning beyond literal content.
- Understanding intentions, implications, and underlying assumptions.
- Examples include figurative expressions and contextual interpretations.

## Words and their components

- words are nothing but a meaningful unit of sentences. by using this kind of words we can construct meaningful sentences
- ↓  
word or  
Token.



## 1. Tokens

Dividing sentence into words. The words are called tokens.

### 3 types of tokens

1. word token
2. character token
3. subword token.

## 2. Lemmas

Based on canonical forms of words

scanning → ran  
run → base word

### canonical form of word

In NL a word often denotes one linguistic item.

form in a given context and the concept behind the form and set of alternative forms. That can express it such soft core called lemmas or lexical items.

## 3. Morphemes

The morphemes is nothing but smallest meaningful unit of any language

### two types of morphems:-

1. free
2. Bound

free:- are nothing but this morphems are independent of another morphems

Bound:-  
Dependent on other morphems.

Ex:- disappointment

## 4. Syntactic categories

Refers to classification or categorization of language based on the structure of grammatical features.

It is used to identify common patterns or variations across the language <sup>and</sup> their properties or relationships.

Lexems can be divided by their behaviour into the lexical categories of verbs, noun, adjectives, conjunctions or other parts of speech.

word =  $w_1 + w_2 + \dots + w_n$

e.g. disagreement = dis + agreement

This model is used to analyse by using different unity of words by combining those words one particular word will be formed.

d. Functional based:-

It is also used to analyse the structure of a word from the given statement

5. Morphology Induction:-

which analyses the pattern,

## Tissues and challenges

→ 1) Irregularity

a) Ambiguity

b) Productivity

Irregularity :-

It means existence of such forms and structures, that or not describe appropriately by a prototypical linguistic model.

Some irregularities can be understood by redesigning the model, and improving its rules. Otherwise, lexical dependency irregularities can't be generalized.

1. Regular verbs & nouns

2. Exceptional inflection

3) Irregular verbs & nouns :- doesn't follow regular pattern.

generates new words by fusing word

generates new words by fusing  
existing words

4) Exceptional inflection:-

comparative & superlative

bige bigger biggest  
dark darker darkest  
good great best

2. Ambiguity :-

(multiple meanings)

word forms that can be understood multiple ways out of context

There are different types

1. word sense ambiguity (particular word differs in meaning)

2. parts of speech " of changing the context (particular word)

3. structural "

we can work in particular sentence in multiple ways.

3. Productivity :-

ability to generate new words or word forms using productivity rules

examples

Morphological models :- There are 5 types of morphological

1. Dictionary lookup

2. Finite state

3. Unification

4. Functional

5. Morphology induction

a. Dictionary lookup

b. It is used to analyse the particular word (or) specific word.

It is used to analyse the structure if a particular word.

Process of analyzing the structure of word :-

i. selecting word.

ii. convert the selected word

based form or canonical

iii. Search for a particular word in the

dictionary.

a. If it match with the base word  
or canonical word. It give  
info about particular word.

A dictionary is a data structure that  
directly enables the precompiled  
word analysis.

Dictionaries enumerate the set of  
association b/w word form &  
their descriptions.

b. finite state morphology

It is also used to analyse the  
structure of a word.  
But it uses the finite automata  
of formal language theory.

c. Unification

Unification is the key operation by  
which feature structures can be

## Finding structure of Document

Language is Structured and analyzing this structure helps various NLP tasks like parsing machine translation, and semantic role labeling.

Identifying sentence boundaries improves readability in speech recognition systems while topic segmentation determines where a topic starts and ends in a document.

### 1. Sentence Boundary detection

Identifying where sentences begin and end in a sequence of text.

challenges:-

- Abbreviations (e.g., "Dr", "Mr")
- Quoted sentences may contain multiple sentence boundaries

### 2. Topic Boundary detection.

Dividing text into meaningful topic-based segments.

Applications:-  
used in IR & text summarization

Date:

15 X 15  
4 IS 4

Finding the structure of Document

In human language words and  
sentences appear randomly  
but really have a structure

A combination of words form sentence  
meaningful grammatical units such  
as noun, verb, command etc.

If the sentence of document is  
entitled it makes further processing  
of text easy in NLP

The NLP tasks depends on the  
document structure or parsing,  
making translations, of semantic role  
and labeling in sentences.

Tree Parsing classification  
is the problem in NLP of deciding  
what is the initial position and

The result is also called a  
sentential segmentation, deals with  
segmenting.

or sentence of word token token  
After process of segmenting sentence  
of speech into units is nothing but  
sentence boundary detection

ex:-

I talked to Dr. Visha ~~about~~  
and I complain about health.

### No Boundary Detection

Topic Segmentation is the task of  
automatically dividing a spoken or  
text of speech into artificially  
large unit blot.

The first form boundary or the  
boundary of speech token  
includes other cases such as  
Speech provider, other cases such as  
Speaker duration of speech etc.

case

sentences

2 types

1. Generative sequence classification method.

2. Discriminative local classification method.

Generative :-

It refers to the techniques that utilizes generative models to classify sequences of text by first learning the underlying probability distribution of data across different classes.

It allows them to generate to new text examples similar to each class. And then

classify new input sequence, by comparing them. These generative random

→ These methods have models, which have not only predicting most likely labels or generating the sentence itself but also generating probabilities.

based on learned probabilities.  
→ It is capable of predicting the mostly likely label, it will decide the boundary of sentence, or boundary of topic.

Discriminative local classification method  
It will only consider local features of a particular sentence and then it will predict to particular sentence.

These approach focus on making decision based on local info.

→ we directly learn the decision boundary b/w different classes in the  $\text{D} \times \text{P}$  space rather modeling the entire data distribution.

Complexity of the approaches:-

The complexity can be measured by rating training & prediction and also measured by their performance of the system

Predicting means they will predicting new sentence based on training example

Performance of the approaches:-

We are going to calculate the performance

$$\text{error percentage} = \frac{\text{No. of errors}}{\text{No. of examples}}$$

precision =  $\frac{\text{Total no. of correct positive prediction}}{\text{Total no. of positive prediction}}$

recall =  $\frac{\text{Total no. of all correct positive prediction}}{\text{Total no. of positive instance}}$

f<sub>1</sub> score =  $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$

## Unit - II

### Syntax I:

- Parsing Natural Language
- Treebanks:
- A data-driven approach to syntax
- Representation of Syntax structure
- Parsing Algorithms.

## Parsing Natural language

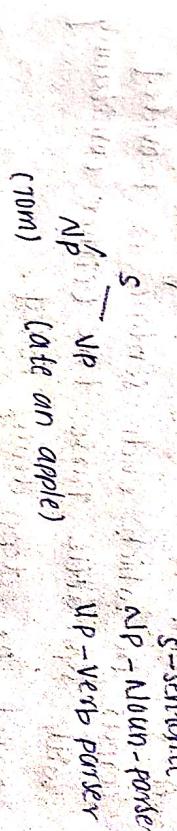
Parsing refers to analyzing a sentence to determine its grammatical structure based on a given formal grammar.

It is essential for applications like text-to-speech summarization, question answering, and more.

(01)

Parsing means analyzing a sentence to understand its structure - like breaking it into smaller parts (subject, verb, object, etc) so a computer can understand it.

Ex:- "Tom ate an apple"



→ Output of parsing is dt:

Parse tree or dependency trees

### Types of parsing structures

1. Constituent Tree (phrase structure)  
Breaks the sentence into nested phrases  
(like Noun phrases, Verb phrases)
2. Dependency tree  
Shows word-to-word relationships, like  
"ate" → main verb  
"Tom" → subject of "ate"  
"apple" → object of "ate"

What does a parser do

1. Takes a sentence
2. Uses grammar rules
3. Builds a tree (either dependency or constituency)
4. Helps applications like translation, voice commands etc.

Challenges in NL parsing:-

1. Ambiguity
2. Incomplete Grammar
3. Robustness
4. Performance.

Treebanks : A data-driven approach to syntax

→ treebank is a linguistically annotated corpus in which each sentence is provided with syntactic parse tree (either constituency based or dependency based).

→ Corpus means collection of something

Construction of a Treebank:-

- Built by linguists manually or through semi-automatic tools.
- Each sentence is:
  - pos tagger (part-of-speech tagging)
  - annotated with syntactic structure.
  - Requires annotation guidelines to maintain consistency across all sentences.

Parsing involves recovering syntactic information that is not explicitly present in a sentence.

- A parser needs knowledge (syntactic rules) beyond just the sentence to generate proper syntactic analysis.
- Traditionally, this knowledge is represented as a grammar (CFG), but natural language syntax is too complex to fully specify with explicit CFG rules.
- Beside rules, we must know which possible syntactic analysis is most plausible (probably) for a given sentence is the second knowledge problem.
- Treebank provides a data-driven soln to both these problems.
- A treebank is a corpus of sentence where each sentence has a complete syntactic analysis provided by human experts.
- These annotations represent the most probable analysis for each sentence.
- The treebank does not explicitly provide grammar rules; instead the syntactic analyses serve as direct knowledge for parsing.

- Treebanks enable parsers to:

- avoid explicitly encoding grammar rules

- use supervised machine learning on the annotated data to learn how to choose

- the best parse for new sentences.

- Two main forms of syntax representation in tree banks:

- Phrases structure trees (Constituency-based)

- Dependency graphs.

- Treebanks ensure consistent, high-quality annotation by following detailed guidelines during the human annotation process!

Ex:- John hit the ball

$s \rightarrow NP VP$

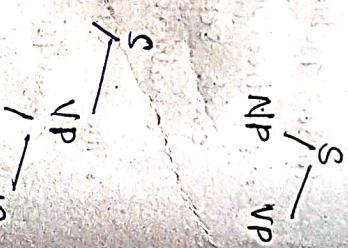
$NP \rightarrow N$

$NP \rightarrow DN$

$V \rightarrow hit$

$D \rightarrow the$

$N \rightarrow John | ball$



### Key:

- words in the sentence are the only nodes in the graph

- arcs show head  $\rightarrow$  dependent relations.

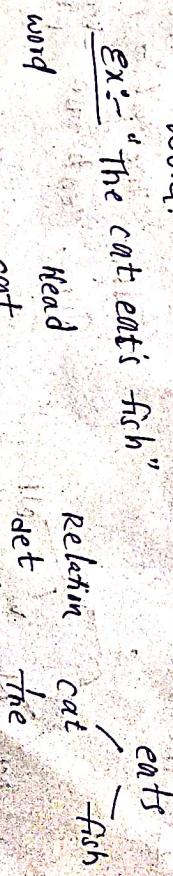
- One word acts as the root (under 0), representing the main action or verb

- It's a form of tree structure where each word (except the root) depends on exactly one head.

- It is asymmetric and directional

- The aim is to find the head of each word.

Ex:- "The cat eats fish"



words

cat  
eats  
fish

eat

the  
cat  
eats  
fish

eat

the

cat

eats

fish

## Representation of syntactic structure.

It can be represented in two main ways:

- Dependency Graphs

- Phrase structure trees (or) Constituency Trees

usage:-

- preferred for free word order languages like Czech, Turkish.

• used in Prague Dependency Treebank

- common in modern NLP application.

### Syntactic Analysis Using Phrase Structure Tree

Phrase structure tree (or constituency trees) break a sentence into nested phrases like Noun phrase (NP) and Verb phrase (VP) based on context-free grammar (CFG)

#### Key:-

- Derived from generative grammar
- shows hierarchical grouping of words into constituents.
- Helps handle long-distance dependencies like questions, gaps, and wh-movement.
- Each sentence is analyzed based on subject predicate structure.

e.g:- "The cat eats fish"

NP - "the cat" - Noun phrase  
NP - "eats fish" - verb phrase

NP - "cat eats fish" - S-sentence

Head-dependent relations

constituent phrases.

- Dependency Graphs are simple and effective for parsing languages with flexible word order. They model word-to-word relations directly.
- phrase structure trees are useful for analyzing nested grammatical structures and are better suited for languages with more fixed word order.
- Both methods help address the knowledge acquisition problem in syntax by providing structured syntactic annotations.
- Treebanks use these representations to create supervised data for training parsers in NLP.

• Dependency graph (Constituency tree)  
shows direct binary relations between words (head and dependent).

• phrase structure tree (constituency tree)  
shows hierarchical grouping of words into phrases (NP, VP, etc)

• Directed graph, with arcs as dependencies  
Tree structure with nested constituents.

• explicit main verb/root symbol is at index 0  
Implicit (format node is S-sentence)

• Head-dependent relations

free word order  
languages  
(Czech, French)

Parsing - Algorithms :-

fixed word order  
languages  
(English, French)

- Parsing algorithms helps in analyzing the  
Syntactic structure of the sentences.

single words connected  
via arts

- Parsing algorithm help in understanding of  
applying the grammatical rules, allowing  
systems to generate grammatically correct

less complex, more  
compact

- More detailed and  
hierarchical.

Machine learning parsing  
models, information  
extraction.

- 1. Recursive Descent Parser
- 2. Shift-reduce parser
- 3. Chart parser
- 4. Regex parser.

- Ans from head to  
branches (top-down)  
dependent  
(graph form)

- Understanding sentence  
structure, grammar  
teaching.

#### 1. Recursive Descent Parser

- It is a type of parsing techniques
- It is a top down parser and reads the  
input from left to right.

phrase labels  
(e.g., NP, VP, S)

relation tables  
(e.g. nsubj, obj, det)

- It works by recursively descending through the  
input string, checking each token against a  
set of production rules defined by the  
grammar of the language.
- The process continues until continues until the  
entire input string has been successfully  
parsed or until an error is encountered.

'3+4 \* (5-2)'

expression → term ('+' '-' ) term)\*

term → factor ('\*' '/' ) factor)\*

factor → Number | expression ()'

expression (3+4 \* (5-2))

term()

factor() → Parse the number 3

expression() → Parse '+' operator

term()

factor() → Parse the number '4'

expression()

term() → Parse '\*' operator

factor() → Parse '(' and ')'.

expression() → Parse the '-' operator

term()

factor() → Parse the number

'12' and '1'

Successfully parsed.

key :-

- Simple and intuitive
- Works well for unambiguous grammar
- Cannot handle left-recursion
- Commonly used in handwritten parsers.

Grammar :-

2. Shift-Reduce parser

- shift-reduce parser is used to reduce a string to the start symbol of a given grammar

- Two parts

1. stack :- shift-reduce parser uses a stack to hold / store the grammar

2. Input tape:

Shift-reduce parser uses input tape to hold the string.

- This parser performs two actions; that is why it is known as shift-reduce parser.

- At the shift action, the current symbol in the input string is pushed to the stack.

- At the reduction, the symbols will be replaced by the non-terminals.

Ex :- production rules

$T \rightarrow T + T$

$T \rightarrow T - T$

$T \rightarrow (T)$

$T \rightarrow C$

String :-

$C_1 - (C_2 + C_3)$

Input

$C_1 - (C_2 + C_3)$

$(C_2 + C_3) \$$

$(C_2 + C_3) \$$

$C_1 - (C_2 + C_3)$

$\$ T -$

$\$ \tau - c$

$C_2 + C_3 \$$

shift

$\$ \tau - C_2$

$+ C_3 \$$

shift

$\$ \tau - C \tau$

$+ C_3 \$$

reduce  
 $\tau \rightarrow C$

$\$ \tau - ( \tau +$

$C_3 ) \$$

shift

$\$ \tau - ( \tau + C_3$

$\$$

shift

$\$ \tau - ( \tau + \tau$

$\$$

reduce

$\$ \tau - ( \tau \del{C_3} + \tau$

$\$$

shift

$\$ \tau - ( . \tau )$

$\$$

reduce

$\$ \tau - \tau$

$\$$

reduce

$\$ \tau$

$\$$

accept

key :-

- uses stack & buffer

- Good for bottom-up parsing

- used in compiler design

- can handle ambiguous grammars using operator precedence.

## A chart parser

The chart parser is a type of top-down parsing technique. It is used to analyse the syntactic structure of the given sentences based on a given grammar.

The chart parser is incrementally builds parse trees by combining smaller constituents into larger ones according to the rules of the grammar.

## Grammar rules

- $S \rightarrow NP, VP$
- $NP \rightarrow \text{det}, N$
- $VP \rightarrow NP, PP$
- $VP \rightarrow V, NP$
- $VP \rightarrow S$
- $PP \rightarrow P, NP$

ex:- The cat chased the mouse  
                ↑      ↓      ↑      ↓  
                Det noun verb Det noun

chart parser

4. Regular parser  
It is also called as regular expression parser.  
It is used to identify patterns or sequences of characters. intent data based on predefined regular expressions.

The main use of regex parser, when a particular pattern is given and to find out some words in it.  
In some cases the regex parser is validating the data.

## Assignment :-

1. Explain the problems by using context free grammar for syntactic analysis of NL.
2. Write a note on tree bank approach of syntax analysis.
3. Explain robust syntax analysis using dependency graph
4. Explain the working of shift reduce parser.
5. Explain different types of parsing techniques briefly.

## Unit - III

### Qyntax II :-

- models for ambiguity resolution in parsing
- Multi lingual issues

### Semantic Parsing :-

- Introduction
- semantic Interpretation.
- System paradigms
- word sense

What is Ambiguity in parsing

Ambiguity in parsing occurs when a sentence can be analyzed (parsed) in more than one way due to unclear grammar or structure.

Ex:- "I saw the man with the telescope"

"Who have the telescope? You or the man?

These sentence has multiple valid parse trees.

Models for ambiguity resolution

3 models for ambiguity resolution

1. Probability or Probabilistic context free Grammar

2. Generative model

3. Discriminative model

1. Probabilistic context free Grammar

- PCFG is a statistical model for parsing NL

- It is an extension of CFG by assigning the probabilities to production rules

- commonly used for ambiguity resolution.

Why use it?

To resolve ambiguity by choosing the parse tree with the highest probability.

steps:-

1. Define a set of production rules which specify how sentence are formed from constituent part.

2. Train the PCFG model with the text. During training, assign production rules with probability.

probability =  $\frac{\text{frequency of the rule}}{\text{total no. of times of non-terminal.}}$

3. Try different parse trees to combine words

4. choose the best parse tree from it.

ex:-

$$\begin{aligned} S &\rightarrow NP \cup VP \\ VP &\rightarrow V \ NP \mid V \ NP \ PP \\ PP &\rightarrow P \ NP \\ NP &\rightarrow Det \ N \mid NP \ PP \end{aligned}$$

Now add probabilities (PCFG):

$$\begin{aligned} S &\rightarrow NP \ NP [1.0] \\ VP &\rightarrow V \ NP [0.3] \mid V \ NP \ PP [0.7] \\ PP &\rightarrow P \ NP [1.0] \\ NP &\rightarrow Det \ N [0.7] \mid NP \ PP [0.3]. \end{aligned}$$

sentence :-

"I saw the man with the telescope"

This has two non-terminals:

"with the telescope" modifies the man ( $NP \rightarrow NP \ PP$ )

2. "with the telescope" modifies the verb ( $VP \rightarrow V \ NP \ PP$ )  
and calculates the total probability of each tree and chooses the most probable one.

2. Generative model

It generates all possible parse trees for a given sentence and selects the most likely interpretation based on a probabilistic framework.

They model how sentences are generated - assign a joint probability to both the sentence and its parse tree.

- work by assuming a hidden structure that generates observed data.

steps:-

1. Define a set of production rules for the given sentence.

2. Generate the possible parse trees for each parse tree of sentence

3. Assign the probabilities score for each parse tree.

4. Choose one parse tree which has more probability

ex:- "The cat sat on the mat"

- Generates a parse tree :  $S \rightarrow NP \ VP$   
- generates words from terminals  
- generates words from terminals  
 $NP \rightarrow Det \ N$  :  $Det \rightarrow "The"$ ,  $N \rightarrow "cat"$ , etc  
- Assign probabilities to each step

The goal is to compute  $P(\text{sentence} | \text{parse tree}) = P(\text{parse tree}) \times P(\text{sentence} | \text{parse tree})$

Then choose the tree with the highest probability.

### 3. Discriminative models.

This model is used to predict the most likely parse tree for the given sentence without explicitly generating all possible trees. They directly learn to distinguish b/w correct and incorrect parses.

- Focus on  $P(\text{parse} | \text{sentence})$  instead of  $P(\text{sentence}, \text{parse})$

#### Steps:-

1. Define a set of features to capture relevant information about the sentence. These features include lexical information, syntactic features and contextual features.
2. Train the classifier with example sentences paired with correct parse trees
3. It predicts the most likely parse tree for the given new unseen sentence.

### Multilingual Issues

What does multilingual parsing mean?  
Multilingual parsing involves developing systems that can analyze and understand the syntax of multiple languages, not just English. But this is challenging because different languages follow different grammatical rules and have different word orders.

- extract features like POS tags, dependency paths, context
- score each parse tree
- pick the tree with highest score

### Summary :-

- PCFG
- Adds probabilities to CFG rules; resolves ambiguity using the most likely parse.
- Generative models Model how a sentence and its parse are generated together.
- Discriminative models Learn to directly select the best parse using features and data.

## 1. Tokenization

Tokenization is the process of breaking a sentence into smaller units, typically words or subwords

### Example:-

"I love playing cricket".  
→ ["I", "love", "playing", "cricket", "."]

Why it's challenging in multilingual NLP

- English: uses spaces clearly b/w words - relatively easy.
- Languages like Chinese, Japanese, Thai: no space between words.
- Ambiguity in compound words:

Tokenizer must know where to split.

- Case sensitivity can be an issue in tokenization.

### How to handle:-

- Language-specific tokenizers
- Subword tokenization
- Neural tokenizers.

## 2. Word Segmentation

Closely related to tokenization; word segmentation refers specifically to detecting word boundaries in languages that don't use space.

### Issues:-

- No spaces in many Asian languages → hard to segment.
- Ambiguous segmentation: multiple valid segmentations possible.
  - In languages like Turkish & Finnish a single word can encode complex meaning through affix.
  - Languages like German forms compound word by concatenating multiple small words.
- Solutions:
  - Statistical models
  - Dictionary-based approaches
  - Neural models.

## 3. Morphology

Morphology studies the structure of words how they are formed by root, prefixes, suffixes etc.

→ Breaking a word into morphemes.

### Issues:-

- When we try to understand sentences in some languages, because some words can change the meaning based on sentence

- By using we can identify the context of the grammar.

- Rich morphology in many languages

- A single root can have dozens of forms

- Data sparsity

- Many word forms = low frequency =

hard for models to learn.

- Ambiguity in affixes

- suffixes can mean different things in different contexts.

so [n]

- morphological analyzers

- stemming and lemmatization

- subword tokenization

- character-level models

(or)

Semantic Parsing is the process of converting natural language into machine-understandable meaning like logical form or code

- Using semantic parsing, the computer can understand natural language the way humans do.

- It is the toughest phase and not fully solved.

- Semantic → study of meaning

- Parsing → identify and relate pieces of information

- Purpose:

To make computers understand what a sentence means so they can act on it.

e.g. - "Sanction serves vegetarian food"

→ output: - serve (sanction, VegetarianFood)

This is called a meaning representation

## Semantic parsing - 1

What is semantic Parsing?

The process of understanding the meaning and interpretation of words, signs and sentence structure is called semantic parsing.

Two types:-

Type meaning

Use Case

**Deep Parsing**      detailed, domain-specific meanings  
e.g.-Travel apps, chatbots.

shallow : General meaning e.g.: word sense,  
Paring. (less detailed) sentiment app.

## Semantic Interpretation

Semantic Interpretation is the process of extracting the meaning of a sentence from its syntactic structure. Grammar tree or parsed sentence).

It converts "how words are arranged" into what they actually mean, why it is important?

- Computers need to understand not just the words, but what is being said
- without semantic interpretation, even can't

grammatically correct sentence may remain meaningless to machines.

It involves identifying predicates and arguments, resolving ambiguities and pronouns and producing a formal representation like

logic

"John gave Mary a gift."

→ predict: gava

→ arguments: John (subject), Mary (indirect object), gift (direct object)

Representation: Give John, Mary, gift

Katz and Fodor, in their 1983 paper "The Structure of a 'semantic' Theory" point forward what properties a semantic theory

A semantic theory should be able to  
should process.

### 1. Explain sentences having meanings.

2. Resolve the ambiguities of word in context

3. Identify meaningless but syntactically well-formed sentences.

## Building blocks of semantic interpretation

## 1. structural ambiguity

## 2. word sense

3. entity eq Event (entity)
4. predicate - argument structure
5. Meaning representation.

### 1. Structural Ambiguity

Structural ambiguity happens when a sentence can be parsed in multiple ways due to its structure - even though the words are the same.

Same sentence, different meanings

Ex:- "I saw the man with a telescope"

This can mean:-

- I used a telescope to see the man.
- The man I saw had a telescope.

Ambiguity is in how the phrase "with a telescope"

is attached (to verb or noun)

### 2. Word sense (Disambiguation)

Word sense refers to the correct meaning of a word in context, especially when the word has multiple meanings.

Word sense Disambiguation (WSD)

The task of selecting the correct meaning of a word based on the surrounding words.

Ex:-  
• "He went to the bank to deposit money"  
→ bank = financial institution

• "He sat by the bank of the river"  
→ bank = riverbank

### 3. Entity and Event Resolution

Entity Resolution:  
Identifying and linking different references to the same real-world entity.

Ex:- "Barack Obama was born in Hawaii. He became

President in 2008"

"He" = Barack Obama (resolved by coreference)

Event resolution:-

Identifying and tracking events and their (who did, what, when, and where)

Ex:-

"The company launched a new product. The launch was a success."

→ Both sentences refer to the same event: the product launch.

### 4. Predicate-Argument Structure

This structure defines the main action (predicate) in a sentence and the entities involved (arguments)

"Who did what to whom?"

format:- predicate (arg1, arg2, ...)

Ex:- "site gave a book to Ram"

→ Arguments:

• Arg1 (subj) = site

• Arg2 (obj) = book

• Arg3 (indirect obj) = Ram

→ Logical form: give (site, book, Ram)

## 5. Meaning representation

It is a formal way to express the meaning of a sentence - often used by machines to reason, query or respond.

Ex:- "the dog chased the cat"

Meaning representation: chase(Dog, Cat)

Type

first order logic

frames / scripts

Lambda calculus  
action: "give", giver : "sita"  
receiver: "ram"

Lambda calculus

used for functional  
compositional meanings.

## System Paradigms

System paradigms refer to the different ways semantic parsing systems are designed and functions depending on:

1. how they are built (architecture)

2. where they were used (scope)

3. How deeply they understand (coverage)

think of paradigms as different "styles" or "models" used to create semantic parses.

The historic approaches which are more prevalent and successfully generally fall into 3 categories:

a. Based on architecture

This classification is based on how the semantic parsing system is built.

1. Major architecture-based paradigms:

knowledge-based systems.

- It rely on manually written rules  
and Ontologies.

- They do not learn from data; everything is pre-programmed.

Ex:- ELIZA (early chatbot)

If sentence contains "book a flight from x to y"

Then BookFlight(x, y)

2. Supervised systems:

These are trained using labeled datasets where each sentence is paired with a correct meaning representation.

Ex:- Neural semantic parser

"Book a flight from Delhi to mumbai"

BookFlight(Delhi, Mumbai)

3. Unsupervised learning system:

It learns from unstructured data without labels using clustering or distribution similarity

#### 4. Semi-supervised system :-

combine small labeled data with large unlabeled data to improve performance. These are useful when labeled data is limited".

#### b. Scope :-

Scope refers to how general or specific the system is in terms of the domains or topics it can handle.

Two types of scope "paradigms"

scope type

Description

Domain-in-Dependent works on specific domain

Domain-Independent works across many domains (general purpose)

Ex:-

"Book a flight from Delhi to mumbai"

Book flight (from = "Delhi", to = mumbai)

Independent

• "Book a flight to delhi"

• "set an alarm for 6am"

#### c. Coverage

Coverage refers how deeply the system understands or analyzes the sentences.

#### Two types of coverage

Types description

Shallow parsing : extracts basic info like subjects, objects, verbs)

Deep Parsing : Builds full semantic meaning representations

Ex :-

shallow :-

Input : "John ate an apple"

-Subject : John

-Verb : ate

Deep :-

Eat ( John, Apple )

#### Lexical sense :-

A word sense is a specific meaning of a word when used in a particular context

One word can have multiple senses (meanings)

depending on how & where it is used

#### Types of word ambiguity

Type

Description

Example

Homonymy  
same word unrelated meaning

unrelated meaning

Bat (animal) vs Bat (used in cricket)

Polysemy  
categorical ambiguity

same word, related meaning

Bank (finance) vs Bank (store money)

word used in different parts of speech

Book (noun) vs Book (verb)

## Word sense Disambiguation (WSD) ?

- WSD is the process of identifying the correct sense of a word based on the context in which it appears.
- It's a subtask of semantic interpretation.

In NLP:

Approaches to word sense Disambiguation

There are 4 main approaches to performing

WSD:

1. Rule-based approach
  - Uses manually written rules, dictionaries and thesauri.
  - Matches words with definitions or contexts

Ex:-

If "bank" appears with "money", assign the

financial institution

- task algo

- structural semantic Inter connections (SSI)

2. Supervised learning - Approach

- Trained on labeled data where the correct sense is already known.
- Use: algo like decision trees, neural network, or transformers.

- Advantages:-

# It achieves high accuracy.

## 3. Unsupervised learning

- unlabeled data, relies on context similarity
  - groups words into clusters of similar usage.
  - No human-annotated data required.

Ex:-

"eat", "consume", "devour" may appear in similar contexts → grouped together.

Adv:-

- No need for extensive manual annotation, making it scalable & adaptable to various languages.

Limitations:-

- May require sophisticated algorithm to achieve high accuracy.
- Performance can be lower.

4. semi-supervised

- Starts with a small set of labeled data + large unlabeled data.

- uses bootstrapping or self-training to improve performance.

## Difference b/w System Paradigms

word sense

System Paradigms refer to the different architectural approaches used to build semantic parsing systems, such as rule-based or domain-specific designs.

In contrast, word sense disambiguation is the process of resolving ambiguity in word sense meanings based on individual word meanings based on context, such as distinguishing "bank" as a financial institution or a riverbank.

• word sense  
• word paradigm  
• word semantics

## Unit - IV

### Semantic Parsing - II

- predicate - argument structure
- Meaning Representation systems.

VINAY KUMAR

## Predicate - Argument Structure

-Predicate-argument structure, also known as semantic role labelling, is a method used to identify the roles of different parts of sentence.

(\*)  
— Predicate - Argument structure describe how a verb (predicate) relates to other elements (arguments) in a sentence. It captures the roles of entities involved in an event or action.

It is a framework to describe the meaning of a sentence by breaking it down into:

• **Predicate:** The main verb or action in a sentence.

• **Arguments:** The entities involved in that action.

(Subject, object etc)

e.g.: It answers the question:

Who, did what to whom, where, when and

Simple example

"Ravi gave a book to Arjun"

Predicate : gave

Arguments : Ravi (giver → subject)

a book (thing → object)

Arjun (receiver → indirect object)

Logical : Give (Ravi, Book, Arjun)

Why is it important

- Helps machine understand sentence meaning

- Crucial for:

- Semantic parsing

- Question answering

- Information extraction

- This is important for things like

translating languages, answering questions and even helping virtual assistants

understand commands better.

- Semantic Roles (Argument Roles)

Each argument has a role in the event.

Role Meaning Example  
Agent does of the action Ravi ("in Ravi kicks the ball")

Theme thing being acted upon (in above ex)

Goal end location or receiver Arjun ("Ravi gave book to Arjun")

Source starting point Delhi ("she flew from Delhi")

Instrument tool used

Knife ("he cut with a knife")

Location where something happened park ("he played in the park")

predicate - argument structure is part of semantic interpretation

- It is useful for creating meaning representation that machines can understand and reason about.

- Helps break complex sentence into action + participants format.

Tools used in NLP:-

FrameNet

User frames (e.g. Buying, Giving) to group roles

PropBank

Labels arguments as Arg0, Arg1, etc

PropBank and FrameNet are two important

tools used in semantic role labeling and predicate - argument structure analysis.

1. FrameNet:

FrameNet is a lexical database that groups verbs and nouns into semantic frames

conceptual structures representing a type of event or scenario.

PL is based on the theory of frame semantics, which suggests that the meaning of a word can be understood in terms of the typical situations it describes.

- \* A frame is a situation (like buying, giving, travelling) with expected participants (called frame elements).

Key elements :-

Frames:- A frame is a type of situation or scenario. Each frame involves certain participants, which are called frame elements.

Frame elements:- These are the roles played by the different participants in a frame.

Working:-

1. Identify frames:-  
Researchers identify common situations (frames).
2. Assign frame elements:-
  - Each frame has specific roles.
3. Label sentences:
  - Sentences are tagged with these frames and frame elements to show how words are used in context.

- \* A frame is a situation (like buying, giving, travelling) with expected participants (called frame elements).

o/p :-  
commerce - buy (Buyer = Ravi, Goods = car, Seller = Arjun)

## Q. PropBank

- PropBank is a corpus of text where each verb is annotated with its arguments, giving us a clear idea of who is doing what to whom in sentence.
- This helps in understanding the roles of different entities in relation to the verb.

PropBank stands for proposition bank.

- It helps NLP systems understand the roles of words in sentence.

## Key elements :-

Predicate: Usually a verb, it represents an action or state.

Arguments: The participants involved in the action or state described by the predicate. Arguments are categorized as core or adjunctive.

Ex:- "Ravi brought a car from Arjun"

frame : commerce - buy  
frame elements :

Buyer : Ravi

Goods : a car

Seller : Arjun

## Goal of Propbank

To provide:

- A consistent, corpus-based labeling of predicate-argument structures
- Numbered roles (ARG0 to ARG5) for different verbs.
- A foundation for training semantic role labeling (SRL) systems.

## How PropBank works

- Each verb is annotated with:

- The sense (meaning) of its arguments (entities participating in the event)
- Its arguments (entities participating in the event)
- Roles are numbered as ARG0, ARG1, ARG2, etc., based on the verb's meaning.
- Modifiers (like time, location) are labeled as ARGM-TMP or ARGM-LOC etc.
- Common PropBank argument roles

### General meaning

Role	Agent / subject / Dyer	Theme / patient / thing acted upon
ARG0		
ARG1		
ARG2		

ARGM-TMP      Temporal Modifiers (time)

ARGM-LOC      Location

ARGM-MNR      Manner

ARGM-CAU      cause

## ex :- "Sita opened the door with a key".

Predicate : opened

Arguments (based on PropBank frame for "open")

- ARG0 = sita → Agent (person doing the action)
- ARG1 = the door → Theme (thing being opened)
- ARG2 = a key → Instrument (tool used)
- ARGM-MNR = with a key (can be modifier too)

### Representation:-

OPEN.0!

ARG0 : sita

ARG1 : the door

ARG2 : the key

Both are used to help NLP systems understand sentence meaning for applications like question answering, summarization and machine translation.

## Meaning Representation Systems (MRS)

- Meaning representation is a deeper level of semantic interpretation aimed at converting natural language into format that machines can understand and act on.

- This process is similar to how programming languages are compiled into machine code that computers execute.

- Unlike artificial language's, natural language is

flexible and relies on context and general

world knowledge for understanding, which poses a challenge for machines.

- Researchers have been working for decades to

develop methods to interpret and encode this

context and knowledge for machines.

- However, current techniques are limited

to specific domains and problems and do

not scale well to arbitrary domains

A meaning representation system is a way

to formally represent the meaning of a natural language sentence so that a computer can:

- Understand it
- Store it

- Reason with it

- think of it as converting language into

a mathematical or logical format.

Goals:-

1. unambiguity: No confusion in interpretation

2. Compositional: Meaning is built from parts

3. Executability: Machines can use it to make decisions

4. Expressiveness:-

can express complex ideas & relationships.

key components:-

entities

people, objects, place mentioned  
in the sentence.

events/predicates

actions or states involving entities  
(e.g. "eat")

arguments

participants in the event  
(subject, object etc.)

relations

connections b/w entities  
(e.g. ownership, location)

Popular Meaning Representation Systems:

1. first order predicate logic (FOL)

- Based on FOL

- Represents facts and relationships using predicates, constants, variables, quantifiers.

Ex: "every student studies"

$\rightarrow \forall(x)(student(x) \rightarrow studies(x))$

sentence: "Ravi loves Priya"

$\rightarrow loves(Ravi, Priya)$

2. Frame-based

- uses a slot-filler structure to represent meaning

- Good for representing events and roles

in context.

3. Semantic Networks:

- meaning is stored as a graph of nodes and links

• Nodes = concept / entities

• Edges = relations.

## Resources:-

- ARTS (Air Travel Information System)
- Domain : flight booking / airline travel queries.

Task :-

Convert natural language into a database query.

Ex:-

"Show me flights from Boston to Denver

tomorrow"

Q/P :-

flight (from=Boston, to=Denver, date=tomorrow)

2. Communicator

• Domain :- Multi-modal travel planning & ticketing

• Use case :- Dialogue system research.

- Developed by DARPA for evaluating

Spoken dialogue systems.

Focus:-

• Robustness of dialogue understanding

• Content tracking, dialogue flow,

meaning representation.

3. Geo Query :-

- US geography database.

- Task :- convert user questions into Prolog-style logical

forms or SQL.

- Example Query

- Impact :- "What is the capital of Texas?"

Q/P :- Q/P :- answer (capital (Texas))

4. Robocup :-

It is an international competition where teams of robots play soccer and its organized by the artificial intelligence community. Their goal is to advance AI and robotics research through this

challenging and fun domain.

Ex:-

Smartie "Move the defenders back and focus on the wings"

Q/P :- Action logic or planning instructions.

These datasets help evaluate the accuracy, generalization, and compositionality of semantic parsers.

## FrameNet vs PropBank

**FrameNet** vs **PropBank**

feature	FrameNet	PropBank
purpose	Models situations (frames)	Labels verb arguments and their participants.
roles	with general semantic roles.	

**Annotation type** frame-based annotation

predicate-argument annotation

**Focus**: Meaning from the frame or scenario

meaning from the verb and its argument

**Role label** Descriptive labels (e.g. buyer, seller, item)

generic numbered roles (ARG0, ARG1, ARG2)

**Verb coverage** More thematic & lexical focused on broad sense richness, often deeper coverage

**Ex:-** "Ravi brought a car from Arjun"

frame: commerce-buy

verb : buy.01

Buyer : Ravi

ARG0 : Ravi (buyer)

Goods : car

ARG1 : car(item)

Seller : Arjun

ARG2 :

**Domain adaptability** rich for linguistic theory

easy to use in machine learning models

**Used in** FrameNet-based, theoretical linguistics

semantic role theory

## Used for FrameNet-based parsers

<b>Data size</b>	smaller (few thousand frames), high manual effort	(PropBank + Penn Treebank)
------------------	---	----------------------------

### key differences

- FrameNet :- scenario-based, uses descriptive roles

- PropBank = verb-based, uses numbered roles (ARG0, ARG1)

## Unit - V

### Language Modeling :

- Introduction
- N-Gram Models
- Language model evaluation
- Bayesian parameter estimation
- ~~Language Model Adaptation~~
- Language model Adaptation
- Language models
  - class based
  - variable length
  - Bayesian topic based
- multilingual & cross lingual language modeling.

## Language modeling

### Introduction:-

Language models are crucial for various applications in human language technology. This help predict the likelihood of word sequences based on training data.

A language model is a statistical tool used to assign a probability to a sequence of words. It helps in predicting the next word, evaluating sentence fluency, and guiding decisions in NLP such as machine translation, speech recognition and chatbots.

### Why do we need:-

language is ambiguous and unstructured, so NLP systems need a way to:

- predict the next word in a system.
- scores sentences based on how likely they are
- Help in choosing correct words in translation, correction or suggestion.

Eg:- sentence:- "The cat sat on the —"

- mat  $\rightarrow 0.6$
  - roof  $\rightarrow 0.2$
  - car  $\rightarrow 0.1$
- most likely next word: mat.

### Applications

- speech recognition
- spelling correction
- machine translation
- auto complete.
- Chatbots

Common language model:

- N-gram model
- Neural language models
- class based models
- variable-length models
- Bayesian topic model
- Multilingual model.

Before discussing this model, let's know what is

the traditional or old method used to predict the word. That is chain rule.

Chain rule :-

The chain rule in probability is used to compute the joint probability of a sequence of words by breaking it into conditional probabilities like:

$$P(w_1, w_2, w_3) = P(w_1) \times P(w_2 | w_1) \times P(w_3 | w_1, w_2)$$

(or)

$$P(w_1, w_2, \dots, w_n) = P(w_1) \times P(w_2 | w_1) \times \dots \times P(w_n | w_1, \dots, w_{n-1})$$

But computing this full history is impractical.

Bcoz it allows us to compute the probability of an entire sentence word by word.

Limitations :- computational complexity, data sparsity,

need for large storage.

→ To overcome this some model like markov assumption are used. and N-gram.

Markov assumption

The probability of a word depends only on a limited no. of previous words, not the entire history.

example:-

Diskard of computation

$$P(w_4 | w_1, w_2, w_3)$$

we assume

$P(w_4 | w_3) \rightarrow$  first order markov assumption.

-Bigram model

$$P(w_1, w_2, \dots, w_n) = P(w_n | w_{n-1})$$

-trigram model

$$P(w_1, w_2, \dots, w_n) \approx P(w_n | w_{n-2}, w_{n-1})$$

It limits the word's dependency to the last  $k$  words simplifying computation in language models

N-Gram model

An N-Gram model is a statistical language model that predicts the next word in a sentence using the previous  $(N-1)$  words.

It applies the markov assumption  
A word depends only on a pre-defined no. of previous words.

An N-Gram is a sequence of  $N$  words  
An N-Gram model uses these sequences to estimate the probability of word sequence.

It estimates

$$P(w_1, w_2, \dots, w_n) = P(w_n | w_1, w_2, \dots, w_{n-1})$$

$$\text{using } P(w_1, w_2, \dots, w_n) = P(w_n | w_{n-1}, \dots, w_{n-k+1}) = \text{N-gram}$$

Depending on the length of  $n$ , we distinguish two unigrams.

Unigrams ( $n=1$ ): The probability of each word is independent of any previous word.

$$P(w_1, w_2, \dots, w_n) = P(w_n)$$

Bigram ( $n=2$ ): The probability of each word depends on previous word

$$P(w_1, w_2, \dots, w_n) = P(w_n | w_{n-1})$$

Trigram ( $n=3$ ): uses two previous words

$$P(w_1, w_2, \dots, w_n) = P(w_n | w_{n-2}, w_{n-1})$$

Ex:- "I want to eat pizza"

$$\cdot \text{Bigram } (n=2) \quad P(I) \times P(\text{want} | I) \times P(\text{eat} | \text{want}). \quad P(\text{pizza} | \text{eat})$$

$$\cdot \text{Trigram } (n=3) \quad P(I). \quad P(\text{want} | I). \quad P(\text{to} | \text{want}). \quad P(\text{eat} | \text{want}, \text{to})$$

$$P(\text{pizza} | \text{to}, \text{want})$$

The model uses frequency counts from a corpus to estimate each probability

Advantages:- Simple & fast, efficient for short context

Good baseline, Real-time friendly

Limitations:- Many word combinations never appear in training data.

- Data sparsity

- Large storage for higher  $N$

- Needs something smoothing

Requires techniques like Laplace, backoff etc.

An  $n$ -gram model is also known as  $(n-1)^{\text{th}}$  order Markov model bcoz it only considers the  $n-1$  previous words to predict the next one.

The Markov Assumption is the idea that a word depends only on a few previous words. The  $N$ -gram model is practical application of this assumption using the last  $(N-1)$  words to predict the next word.

### LANGUAGE MODEL EVALUATION

It is the process of measuring how well a language model performs, especially in predicting word sequence correctly.

To judge the performance of the language model, how successful it is at estimating the word probabilities we use language model evaluation.

#### 1. PERPLEXITY

Perplexity tells how many word choices the model

is "confused between" at each step.

It is a measure of how well a language model predicts a sequence of words.

$\rightarrow$  Lower perplexity = better performance

formula:-

$$\text{perplexity} = 2^{-\frac{1}{M} \sum_{i=1}^M \log_2 P(w_i | \text{context}_i)}$$

- If the model has no predictive power at all, perplexity is equal to vocabulary size.

A model achieving perfect prediction, by contrast, has a perplexity of one.

## Q. Coverage rate

- coverage rate refers to the percentage of words in the test set that were seen during training.

high coverage = better vocabulary match between training and test data

If the model knows all the words, it has high coverage.

- It tells the percentage of n-grams present in the test set.

- A special case is the out-of-vocabulary rate (OOV rate), which is the percentage of unique words in the test set that the model does not know:

$$\text{formula: coverage rate} = \frac{\text{No. of known (seen) words}}{\text{Total words in test set}} \times 100\%$$

e.g. If the test set has 1000 words, and 920 of them were seen during training

$$\text{coverage rate} = \frac{920}{1000} \times 100 = 92\%$$

## 3. Cross validation:-

A method where data split into k-parts to test the model multiple times for more reliable evaluation.

## 4. Split validation (train - Test split)

split validation means dividing data into 2 fixed parts

- Training set :- to train the model
- Testing : to evaluate the model

In the ratio: 80/20 → 80% training, 20% testing

$$80/20 \rightarrow 80/20$$

## Parameter estimation

Q. What are parameters?

In NLP models, parameters are values that the model uses to make decisions.

- For example, in a language model, parameters might

include probabilities that represent how likely one word is to follow another word.

parameter estimation:-

It is process of determining the best values for the parameters of a model so that it can accurately perform tasks like

- understanding, generating, translating text  
parameters estimation is the process of calculating the probability values (parameters) that a language model uses to predict words.

In simple terms

It means learning how likely each word or word sequence is, based on the training data.

Importance :-

- The accuracy and effectiveness of an NLP model depend on how well its parameter are set.

- Good parameter estimation ensures that the model can understand & generate human-like text accurately.

there are 2 techniques

1. Maximum likelihood estimation
2. Bayesian estimation.

### Maximum likelihood estimation

MLE is about finding the parametric values that make the observed data most likely.

In the context of NLP, this means finding the probabilities of words or sequences of words that best explain the text data you have.

How it works:-

1. Gather a large amount of text data
2. Count how often different words or sequences of words appear in the data.
3. Use these counts to estimate the probabilities of words or sequences. For example, the probability of a word  $w$  given a previous word  $w_{\text{prev}}$  is estimated as:

$$P(w|w_{\text{prev}}) = \frac{\text{Count}(w_{\text{prev}}, w)}{\text{Count}(w_{\text{prev}})}$$

Here,  $\text{Count}(w_{\text{prev}}, w)$  is the no. of times  $w$  follows  $w_{\text{prev}}$ , and  $\text{Count}(w_{\text{prev}})$  is the no. of times  $w_{\text{prev}}$  appears.

ex- "I like pizza. I like coding"

$$\text{Count}("I \text{ like}") = 2$$

$$\text{Count}("like \text{ pizza}") = 1$$

$$\text{Count}("like \text{ coding}") = 1$$

Now

$$P(\text{like} | T) = \frac{2}{2} = 1.0 \quad P(\text{pizza} | \text{like}) = \frac{1}{2} = 0.5$$

$$P(\text{coding} | \text{like}) = \frac{1}{2} = 0.5$$

These probabilities are estimated from word counters - this is parameter estimation.

### Bayesian estimation

It combines prior belief (before seeing data), likelihood (from the data) to produce a posterior probability (updated belief after seeing data), based on both prior knowledge and observed data.

It combines

- prior belief (before seeing data)
- likelihood (from the data)
- To produce a posterior probability (updated belief after seeing data).

How it works:-

1. Start with a prior distribution that represents your initial beliefs about the parameters
2. Gather your text data.

3. Use Bayes Theorem to update the prior distribution with the observed data to get the posterior distribution. This posterior distribution represents your updated beliefs after seeing the data.

Bayes theorem :-

$$P(\theta | D) = \frac{P(D | \theta) \cdot P(\theta)}{P(D)}$$

where,  
 0 - the parameter to estimate (like a probability)

D - observed data

$p(\theta)$  - prior probability

$p(D|\theta)$  - likelihood

$p(\theta|D)$  - posterior (final estimated value).

Ex :- we are estimating  $p(\text{pizza} | \text{like})$

this is the conditional probability:

- what is the probability of the word "pizza"

appearing after the word "like"?

Step 1 :- word counts from training data.

lets assume your corpus had only these 2 bigrams

- like pizza → occurred 2 times
- like coding → occurred 1 time

that gives us:-

bigram count

like pizza

like coding

like dogs

like rain

so, total times "like" word followed by something

Step 2 :- Vocabulary size ( $V$ )

assume the vocabulary contains 4 possible words  
 that could "like"  
 pizza, coding', dogs, rain

Step 3 :- apply laplace smoothing (Bayesian estimation)  
 we use the Dirichlet prior ( $\alpha=1$ ), which adds  
 1 pseudo-count to every word.

Step 4 :- use the formula  
 Bayesian smoothed probability

$$p(w | \text{like}) = \frac{\text{count}(\text{like } w) + \alpha}{\text{count}(\text{like}) + \alpha \cdot V}$$

plug in values:

$$\text{for pizza: } p(\text{pizza} | \text{like}) = \frac{1+1}{2+1+4} = \frac{2}{6} \approx 0.33$$

$$\text{for coding: } p(\text{coding} | \text{like}) = \frac{1+1}{2+4} = \frac{2}{6} \approx 0.33$$

$$\text{for dogs: } p(\text{dogs} | \text{like}) = \frac{0+1}{2+4} = \frac{1}{6} \approx 0.166$$

$$\text{for rain: } p(\text{rain} | \text{like}) = \frac{0+1}{2+4} = \frac{1}{6} \approx 0.166$$

each word is assigned a probability using both  
 real data (actual count) and prior belief ( $\alpha=1$ ). This  
 prevents assigning zero probability to unseen  
 words & makes the model more robust.

### Language Model Adaptation

It is the process of modifying a trained language  
 model so that it works better for a new domain,  
 topic, language or user.

In simple terms:-

In simple terms:-  
 we "adapt" a general model (trained on  
 large data) to make it work well on specific  
 or smaller datasets.

This is important when the data used to train the model is different from the data when the model will actually be used.

#### Why adaption is needed:-

- General models are trained on broad (or para)

(e.g. news, weather)

- But in real world use, we need to perform well on:-

- Specific domains (medical, legal, tech)
- Specific users (custom vocabulary or style)
- Specific languages (dialects, regional use)

Techniques:-

#### 1. Fine Tuning :-

You take a general language model and train it finely using data from a specific domain.

→ Continue training the model on small, target-domain data

a. Common in neural language model

e.g:- a general GPT-like model fine-tuned on medical records.

#### 2. Model Interpolation

- Also called mixture models

- Combines two models

- A general model
- A domain-specific model.

Helps use both broad and focused knowledge.

#### 3. Topic-Dependent model

- These models are trained or adapted for specific topics or domains (e.g. sports, medicine etc)

to improve accuracy & relevance.  
e.g:- In a news app, you might have

- A model adapted to sports language

- If user is reading sports news, the system activates

the sports-specific language model.

#### 4. Self-adaptive models.

- Also called as trigger models

- This technique adapts the model dynamically based on the words that have already been seen in the text.

e.g:- In a financial document, seeing the word "stock" increases the likelihood of subsequently seeing "market", "uprise" or "exchange".

#### 5. Unsupervised Adaptation

- Adaptation done without labeled data, using raw unannotated text from the target domain

e.g:- Suppose you want to adapt your model for legal documents, but you don't have labeled examples.

#### Advantages :-

1. Improves accuracy in specific domain
2. Reduce errors on specialised terms
3. Better User Experience
4. Makes use of target data
5. Dynamic learning

## Language Models

It is a system that assigns probability to sequences of words and predicts the next word based on the previous words.

### Class Based

A class-based language model groups into semantic or syntactic classes (like nouns, cities, verbs) and predicts words based on their class, instead of predicting individual words directly.

### Advantages :-

- Reduces data sparsity
- Generalizes better across similar words
- Improves accuracy in low-resource settings

### Limitations:-

- Requires a good word-to-class mapping
- Might lose individual word-level nuances.

Discuss with example class based variable length model.

class based :-

Class based language models are way of simplifying language processing by grouping words that have similar meanings or functions into categories, called "classes".

- Instead of treating each word separately the model works with separately, the model works with these groups of words to make predictions about language.

Ex:- let's say we have a sentence:

"the favourite pet is a —

class prediction: The model might predict the next word belongs to the "animals" class because "pet" commonly appear in this context.

word selection: If the word "pet" belongs to the "animals" class, it might choose the word "dog", resulting in sentence

"The favourite pet is dog".

- How it works :-
- words are grouped into classes based on their roles or meaning
  - calculates the class probability
  - calculates the word probability within the class.

### Applications :-

- Low resource language modeling
- class based machine translation
- speech recognition

Variable length language model

It uses context of different lengths dynamically, instead of fixed-size n-grams.

It is also called

- variable-order n-gram model
- cache-based or Trie based model

How it works :-

- Instead of fixed history, it uses a much useful history as available
- Adapts to longer or shorter context dynamically

Ex:- "she ate a cake with a fork"

A fixed trigram :-  
a cake / a ate), p (with | cake, a)

variable-length model :-

- EP (fork | a cake with a)  $\rightarrow$  longer context gives better prediction.
- \* improves context sensitivity & prediction accuracy.

Applications :-

- + Neural language models
- speech recognition with large vocabularies
- context-sensitive predictive text

## variable length model

A variable length language model uses different amounts of context (not a fixed  $N-1$  window like  $N$ -grams) to predict the next word.

### Advantages:-

- Captures long-term dependencies better than fixed  $N$ -grams.
- More flexible in real-world data
- Reduces zero probabilities from unseen long sequences

### Limitations:-

- can be computationally heavy
- requires more memory to store variable-length histories.

## Bayesian Topic Based model

A Bayesian topic based model is a probabilistic language model uses 'Bayes Theorem' to discover hidden topics in a collection of documents, and then uses these topics to model and generate language more accurately.

most common model: Latent Dirichlet Allocation (LDA)

LDA is popular unsupervised Bayesian model for

- Topic modeling
- Document classification
- Language modeling.

key concepts:-

- Topic: A distribution over words  
e.g. Topic 1 : game, team, score | Topic 2: doctor, patient |
- Document: A mixture of topics  
e.g. A sport article 80% Topic 1, 20% Topic 2 |

- Bayesian Inference:-

Used to compute the posterior probability of topics given the words in document.

e.g. we have 3 documents

1. The doctor examined the patient

2. The team won the game

3. The hospital hired a new doctor

LDA might find:-

- . Topic A (medical), doctor, patient, hospital, medicine
- . Topic B (sports), team, game, win, score.

Then,

- Doc 1 = 90% Topic A, 10% Topic B

- Doc 2 = 45% Topic B, 55% Topic A

- Doc 3 = 70% Topic A, 30% Topic B

Advantages :-

- Captures hidden topics

- works with unlabeled data

- Generalizes across different documents

Limitations:-

- Needs long document
- requires hyperparameter tuning
- computationally expensive.

## Multilingual and cross-lingual language modeling

Multilingual language modeling involves training models on multiple languages so they can understand, generate or translate across different languages.

Cross-lingual language modeling focuses on transferring knowledge from one language (usually a high-resource language like English) to help tasks in another (often low-resource)

language.

Why its important :-

- Many languages lack enough data

- Cross and multilingual LM allow models to share knowledge across languages

- o perform translation, sentiment analysis or NLP in multiple languages.

- o understand context even in code-mixed or bilingual data.

## Feature

### multilingual

Def

Trained on multiple languages at once

goal

Understand / generate text in many languages

e.g. task

One model doing translation or Q&A in English, Hindi, Tamil

Train on English Q&A apply to Hindi Q&A

Training data

Requires multi-language datasets

May use only one language's data during training.

use case

Multilingual chatbots  
translation systems

low-resource language support using high resource training.

Model examples

BERT, mBERT, XLM-R  
XLM, LLaMA, cross-singual BERT.

example :-

multilingual :-

mBERT trained on English+Hindi+Tamil → can understand all

Cross lingual :-

Train a model on English news → use it to classify news in Hindi without training to Hindi.

## Cross Lingual

Trained in one language used to perform task in others.

Transfer knowledge from one language to another

A cross-lingual model can still answer "कौन है आप?" by transferring knowledge from

"Who are you?", by translating knowledge from English training.

### Advantages:-

- Low resource support
- Transfer learning
- Multilingual NLP tools
- Handles code-mixing.

### Limitations:-

- Bias toward high resource languages
- Vocabulary mismatch
- Expensive training.

You train a QA system on English data.

Then you ask it a question in Hindi:

"आरती गाज थारी क्या है?