

Chapter 1. finding the structure of words

Human language is a complicated thing. We use it to express our thoughts and through language, we receive information and infer its meaning.

Linguistic expressions are not unorganized, though. They show structure of different kinds and complexity and consist of more elementary components whose co-occurrence in context refines the notions they refers to in isolation and implies further meaningful relations between them.

In this chapter, we explore how to identify words of distinct types in human languages, and how the internal structure of words can be modeled in connection with the grammatical properties and lexical concepts the words should represent. The discovery of word structure is morphological. parsing words and their components!

words are defined in most languages as the smallest linguistic units that can form a complete utterance by themselves. The minimal parts of words that deliver aspects of meanings to them are called morphemes.

Tokens:

Sentence divided into segments, punctuations and words , eliminating white space

Suppose, for a moment, that words in English are delimited only by whitespace and punctuation

example: will you read the newspaper? will you read it? I won't read it.

lexemes:

The content of the concept provided by the lexical analyzer using set of alternative forms in the algorithm. The lexemes also consider words in the context called lemmas

example: Did you see him? I didn't see him. I didn't see anyone

Presents the problem of tokenization of didn't and the investigation of the internal structure of anyone.

Morphemes:

Morphemes properties of words and their components are usually called segments. The parts of speech of word are called morphemes or morphs

Morphological theories differ on whether and how to associate the properties of word forms with their structural components.

These components are usually called segments or morphs.

Example: 1. Whether will + you - read this the - newspapers?

a. Whether will + you - read + it? not - will I - read + it

Typology:

Morphological typology divides into group of languages that phenomena provides prevalent characters.

Isolating (or) analytic :

Languages include no or relatively few words that would comprise more than one morpheme.

Synthetic: synthetic languages can combine more morphemes in one word and are further divided into agglutinative and fusional languages.

Agglutinative: Agglutinative languages have morphemes in one word associated with only a single function at a time

Fusional: language are defined by their feature-per-morpheme ratio higher than one.

Concatenative: language linking morphs and morphemes one after another.

Issues and challenges.

Morphological parsing tries to eliminate or alleviate the variability of word forms to provide higher-level linguistic units whose lexical and morphological properties are explicit and well defined. It attempts to remove unnecessary irregularity and give limits to ambiguity both of which are present inherently in human language.

- Irregularity
- Ambiguity
- Productivity
- Irregularity:

Morphological parsing is motivated by the quest for generalization and abstraction in the world of words. Immediate descriptions of given linguistic data may not be the ultimate ones due to either their inadequate accuracy or inappropriate complexity, and better formulations may be needed. The design principles of the morphological model are therefore very important.

Irregular words like if they are in sentence provide indirect meaning made abstractions for irregular morphology.

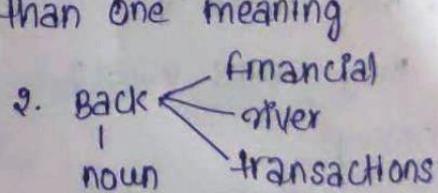
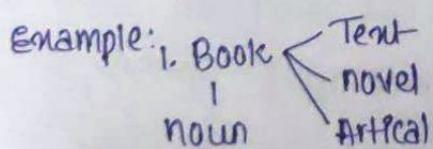
- Ambiguity:

Morphological Ambiguity is the possibility that word forms can be understood in multiple ways in the context.

- Word, sentence, phrase are ambiguous words contains multiple meanings.

Types of Ambiguity:

Lexical Ambiguity: word has more than one meaning



3. I can play cricket

give me that can

Syntactic Ambiguity:

To present in sentence structure specify possible arrangements of word in sentence.

This ambiguous words depends on grammatical roles

- Abhi saw the man with book
- Abhi saw the man carrying book
- Abhi saw the man through the book

Syntactic Ambiguity:

Each word in the sentence have more than one meaning

- If a sentence has 10 words each having three meaning associated.
- Number of possible interpretations

Example:-

$$\begin{array}{c} \text{10 words} \\ \text{in in in} \end{array} \quad 3 * 3 * 3 * 3 * 3 * 3 * 3 * 3 * 3 = 59049$$

Madu loves his cat

deepu loves too

deepu also loves cat

Anaphoric Ambiguity:

To identify meaning of the sentence provides parts of speech in the sentence and identify when, where, by whom occurrence of the sentence what was said

example: 1. Monkeys eat banana,
when they wake up

Here they means what?

- monkeys

2. Monkeys eat banana,

when they are ripe

Here they means what?

- Banana.

Pragmatics:

(3)

understanding speakers intension and also known as informative ambiguity

example: you're late (criticism)

you're good (compliment)

Challenges of Ambiguity:=

1. Elongated words:

I am sooo sorry → too long words

It was too yummy

Shortcuts: please = pls

By the way = BTW

what = wat

OK = K

Ellipsis :

example: peter worked hard and passed exam, Kelvin too

Interpretation.

Peter worked hard, Peter passed exam

Kelvin worked hard, Kelvin passed exam

Kelvin worked hard and also passed exam

Ambiguity punctuations:

women, without her man, is nothing

women! without her, men, is nothing

Productivity:

1. In India we have INNLTK (India natural language tool kit)
API (Application programming interface)
from nltk import setup

Tokenize:

```
from nltk import tokenize  
tokenize(text, 'code-of-language')
```

Out:- Native language	Code.
Hindi	hi
Punjabi	pa
Gujarathi	ga
Telugu	te
English	en
Urdu	ur
malayalam	ml
Tamil	tm

Code mixed language: combination of languages

example: English + Hindi	shpt	Code
= (Hindi + Eng) Hinglish	latm	hi-en
Tamil + English → Tanglish	latm	tm-en
Malayalam + English → Manglish	latm	ml-en

Get embedded Vectors:

If we give ur then the code converted to Urdu and gives output as Urdu

This return array of "Embedded Vectors"

```
from nltk import get_embedding_vectors
```

```
Vectors = get_embedding_vectors('text', 'code-of-language')
```

```
Vectors = get_embedding_vectors('Bharath', 'hi')
```

Output: [1128]

len(Vectors)

Out:- 4

Predict -n- words

④

```
from nltk import predict_ngram_words  
predict_ngram_words(text, n, 'code-of-language')
```

Identify the language

```
from nltk import identify_language  
reset_language_id(identify_models)  
identify_language(text)  
identify_language("भारत")
```

Output:- Hindi-hi

Bharath

In Indian language the product of the NLP is INLTK.

Command: pip install nltk (aps)

Here INLTK runs on CPU as a desired behaviour most of the learning models in product supported languages.

API

To set up INLTK in the CPU before that to run and Import INLTK Application

Command :- from nltk import setup

- To Import native languages and code of the language

```
from nltk
```

```
setup('code-of-language')
```

→ Morphological models:

Morphology: It is study of structure of word and word formation

Morpheme: The smallest meaning unit of a word that has meaning

Type of morphemes:

1. free morphemes

- 1. Lexical morphemes

- 2. Functional morphemes

2. Bound Morphemes

- 1. Derivational

- 2. Inflectional

free morphemes: The word can stand by themselves means it has own meaning

Eg:- Boy, girl, car, beauty.

Lexical morpheme: Classify the words using POS

Eg:- Nouns, Verbs and Adjectives

Functional Morphemes: Built the class of the words using POS that have grammatical functions

Eg:- Prepositions : of, to

Conjunctions : And, But

Determiners: this, that

Pronouns : I, you

Verbs : Is, can

Bound Morphemes: These words cannot stand alone By themselves cannot exist independently

Eg:- Boy + s

mouses + s = mouses.

↓
Bound morpheme

Derivational Morpheme: It changes the category of the word and there are some exceptions. changes the grammatical category contains two types:

1. Class changing

2. Class maintaining

1. Class changing

Eg:- Driver Drive Beaty Beautiful
 ↓ ↓
 Noun Verb ↓ ↓
 Noun Adjective

2. Class maintaining.

Eg:- Adjust Adjustment (3rd person)
 ↓ ↓
 Verb noun

Help Helpful
↓ ↓
Verb Adjective

Inflectional morpheme:

That are used to indicate grammatical functions of a word and provide tenses to the sentence.

Eg:- Reach + ed (past tense)

Eat + en (past participle)

Happier (-er) (comparative)

Happiest (superlative) (-est)

Dictionary lookup:

Morphological passing is a process by which word forms of a language are associated with corresponding linguistic descriptions in this context a dictionary is understood as a data structure that directly enables obtaining some precomputed results, in our case word analysis the data structures can be shared lookup operations are relatively simple as usually quick dictionaries can be implemented for instance, as lists, binary search trees, has tables etc.

Unification-Based morphology:

It have been inspired by advances in various formal linguistic framework aiming as enabling complete grammatical descriptions of human languages.

A finite-state morphological models, both - surface and lexical forms are by themselves unstructured strings of atomic symbols In higher-level approaches. Linguistic information is expressed by more approximated data

P: $\Psi \rightarrow \{\Psi\}$ p: form $\rightarrow \{content\}$

functional morphology :-

This group of morphology models includes not only the ones following the methodology of functional morphology but even those related to it defines its modify using principles

of functional programming and type theory

It implementation are intended to be refuted as programming libraries capable of handling the complete morphology of a language and to be incorporated into various kinds of applications

I: lexeme $\rightarrow \{ \text{parameter} \} \rightarrow \{ \text{form} \}$

D: lexeme $\rightarrow \{ \text{parameter} \} \rightarrow \{ \text{lexemes} \}$

L: context $\rightarrow \{ \text{lexemes} \}$

Morphology Induction:

We have focused on finding structures of words in diverse language supposing we know what we are looking for

We have not considered the problem of discovering and inducing word structure without the human insight

There are several challenging issues about detecting words structure just from the forms and their context.

They are ambiguity and irregularity in morphology.

In order to improve the changes of statical inference, the learning of morphology for multiple language is proposed by resulting in discovery of abstract morphemes.

Finding the structure of documents:

In human language, words and sentences do not appear randomly but usually have a structure. For example, combinations of words form sentences - meaningful grammatical units, such as statements, requests, and commands. Likewise, in written text, sentences form paragraphs - self-contained units of discourse about a particular point or idea. Sentences may also be related to each other by explicit discourse connectives such as therefore.

Automatic extraction of structure of documents helps subsequent natural language processing (NLP) tasks; for example, parsing, machine translation, and semantic

Sentence Boundary Detection:

Sentence boundary detection (also called sentence segmentation) deals with automatically segmenting a sequence of words tokens into sentence units. In written text in English and some other languages, the beginning of a sentence is usually marked with an uppercase letter, and the end of a sentence is explicitly marked with a period (.), a question mark (?), an exclamation mark (!) or another type of punctuation.

- It was the best of times
- It was the worst of times
- It was the age of wisdom
- It was the age of foolishness

Division of tokens

'It', 'was', 'the', 'best', 'of', 'times', 'worst', 'age', 'wisdom', 'foolishness'.

Here we removed repeated words, white spaces and punctuations

- It was the best of times [1, 1, 1, 1, 1, 1, 0, 0, 0, 0]
- It was the worst of times [1, 1, 1, 0, 1, 1, 1, 0, 0, 0]
- It was the age of wisdom [1, 1, 1, 0, 1, 0, 0, 1, 1, 0]
- It was the age of foolishness [1, 1, 1, 0, 1, 0, 0, 1, 0, 1]

Topic Boundary Detection:

Topic segmentation (sometimes called discourse or text segmentation) is the task of automatically dividing a stream of text or speech into topically homogeneous blocks. That is, given a sequence of (written or spoken) words, the aim of topic segmentation is to find the boundaries where topics change.

Topic segmentation is an important task for various language-understanding applications, such as information extraction and retrieval and text summarization.

Append the parts of speech

> Import corpus

```
nlp = Corpus("package Name");
```

```
doc = nlp(para)
```

"This is nlp document" Document contains some paragraph".

```
list = nlp(list)
```

```
out: ["'this'", "'is'", "'nlp'", "'document'", "'Document'", "'contains'", "'some'", "'paragraph'"]
```

```
out: "'This',
```

```
'is',
```

```
'nlp',
```

```
'document',
```

```
'Document',
```

```
'contains',
```

```
'some',
```

```
'Paragraph'"
```

```
In: sent = list(doc.sents)
```

```
out: ['This', 'is', 'nlp', 'document',  
      'document', 'contains', 'some', 'paragraph']
```

```
In: len(sents) # length of the sentence]
```

```
In: print(sents)
```

```
Out: 2
```

Table formation of text:

```
In: text:"Apple is looking at buying U.K starts for $1 billion"
```

```
In: pip install Table
```

```
In: from table import Table;
```

```
In: Table.columns.header = ["text", "pos", "Tag", "Explam.Tag", "shape",  
                           "is_stop"]
```

```
In: for token in doc:
```

```
    Table.rows.append([token.text, token.pos_, token.tag_,  
                      token.explam_tag, token.shape_, token.is_stop])
```

Discriminative classification model:

(5)

①

In: print (Table)

out:	text	pos	TAG1	Explain-TAG	shape	ps-stop
	Apple	Noun	xxx	smg-Noun	xxxxx	0

Methods:

Sentence segmentation and topic segmentation have mainly been considered as a boundary classification problem. Given a boundary candidate ℓ between two word tokens for sentence segmentation and between two sentences for topic segmentation, the goal is to predict whether or not the candidate ℓ is an actual boundary (sentence or topic boundary).

In this we have two methods

- Generative
- Discriminative

Generative

The most commonly used generative sequence classification method for topic and sentence segmentation is the hidden markov model (HMM).

Naive Baye's algorithm

e.g. frequency table.

S.No	Outlook	Play
0	Ramy	Yes
1	Sunny	Yes
2	overcast	Yes
3	Overcast	Yes
4	Sunny	No
5	Ramy	Yes
6	Sunny	Yes
7	Overcast	Yes
8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Ramy	No
12	Overcast	Yes
13	Overcast	Yes

(3)

Descriptive classification model:

Descriptive classifiers aim to model the distribution of attributes. The most likelihood table.

	yes	No
overcast	5	0 = $\frac{5}{14} = 0.35$
Rainy	2	2 = $\frac{4}{14} = 0.29$
Sunny	3	2 = $\frac{5}{14} = 0.35$
	<hr/> $\frac{10}{14} = 0.7$	<hr/> $\frac{4}{14} = 0.29$

frequency table

↓

likelihood table

↓

Probability

↓

decision

By applying Bayes' theorem

for Yes

$$P\left[\frac{\text{yes}}{\text{sunny}}\right] = \frac{P\left[\frac{\text{sunny}}{\text{yes}}\right] * P(\text{yes})}{P(\text{sunny})}$$

$$P\left[\frac{\text{sunny}}{\text{yes}}\right] = \frac{3}{10} = 0.3$$

$$P(\text{sunny}) = \frac{5}{14} = 0.35$$

$$P(\text{yes}) = \frac{10}{14} = 0.7$$

$$\Rightarrow \frac{0.3 \times 0.7}{0.35} = \frac{0.21}{0.35} = 0.60$$

for No

$$P\left[\frac{\text{No}}{\text{sunny}}\right] = \frac{P\left[\frac{\text{sunny}}{\text{No}}\right] * P(\text{No})}{P(\text{sunny})}$$

$$P\left[\frac{\text{sunny}}{\text{No}}\right] = \frac{2}{4} = 0.5$$

$$P(\text{No}) = \frac{4}{14} = 0.29$$

$$P(\text{sunny}) = \frac{5}{14} = 0.35$$

$$= \frac{0.5 \times 0.29}{0.35} = \frac{0.145}{0.35} = 0.41$$

$$P\left[\frac{\text{yes}}{\text{sunny}}\right] > P\left[\frac{\text{No}}{\text{sunny}}\right]$$

Discriminative classification Model:

(3)

Discriminative classifiers aim to model $P(y_i/x_i)$ directly. The most important distinction is that whereas class densities, $p(x/y)$, are model assumptions in generative approaches, such as naive Bayes, in discriminative methods, discriminant functions of the feature space define the model.

e.g.: Import numpy as np;

import pandas as pd

dataset = pd.read_csv('Social_Network_Ads.csv')

Dataset.head()

O/p:

	userId	RegNo	G	Age	FstSal	Purchased
0	156250	M	19	19000	0	
1	156251	M	35	20,000	0	
2	156252	F	26	43,000	0	
3	156253	F	27	57,000	0	
4	156254	M	19	76,000	0	

Dataset.isna().sum

O/p:

	userId	RegNo	gender	Age	fst sal	Purchased
	0	0	0	0	0	0

X = data.iloc[0:2, 2:4].values

y = data.iloc[0:0, 0:4].values

out: X = [0:0; 19, 19000]
[0:0; 35, 20000]
[0:0; 26, 25000]
[0:0; 27, 30000]
[0:0; 29, 35,000]

out y = [0, 0, 0, 0, 0]

Regression:-

import logistic Regression
(x.data, y.data, x.age, y.Purchased).values

Overall x & y data is displayed.

Discriminative sequence classification methods.

In segmentation tasks, the sentence or topic decision for a given example (word, sentence, paragraph) highly depends on the decision for the examples in its vicinity. Discriminative sequence classification methods are in general extensions of local discriminative models with additional decoding stages that find the best assignment of labels by looking at neighbouring decisions to label an example.

It is used to remove repeated data, irrelevant data, noisy data, unwanted data.

Sequence classification model

↓

Pre-processing model

$$S_1 = \sum_{\pi \in C_1} (\pi - \mu_1)(\pi - \mu_1)^T$$

$$C_1 = \{(4, 1), (2, 4), (2, 3), (3, 6), (4, 6)\}$$

$$C_2 = \{(9, 10), (6, 8), (9, 5), (8, 7), (10, 8)\}$$

$$\mu_1 = \left[\frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}, \frac{y_1 + y_2 + y_3 + y_4 + y_5}{5} \right]$$

$$= \left[\frac{4+2+2+3+4}{5}, \frac{1+4+3+6+6}{5} \right] = [3, 4] = \mu_1$$

$$\mu_2 = \left[\frac{9+10+9+8+10}{5}, \frac{10+8+5+7+8}{5} \right] = [8.4, 7.6] = \mu_2$$

Apply coherent matrices

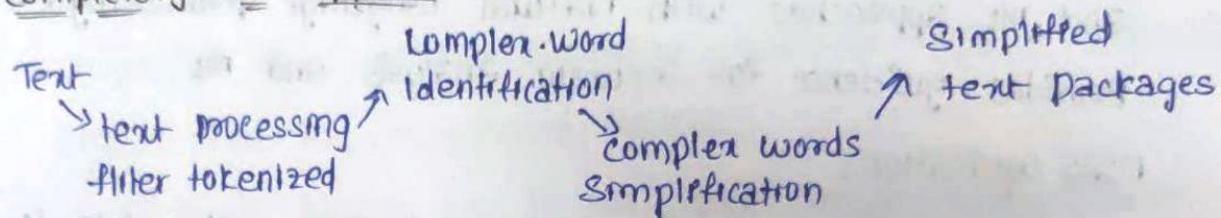
$$(\pi - \mu_1)(\pi - \mu_2)^T$$

$$[\pi - \mu_1] = [1 \ 1 \ 1 \ 0 \ 1]$$

Generate predictions:

Finally the output is the predictions that the ML Model will make on unknown data to understand information with in the context are related to each other.

Complexity of Approaches



Example: `import nltk.data
nltk.data.load('corpora')
from nltk.corpus import dataset
Vocabulary(categories)`

In linguistics, complexity is a characteristic of a text but there are multiple measures and hence multiple implied definitions in practice. In NLP, the measures are useful for descriptive statistics.

Complexity measures:

most popular ways of answering text

Complexity: how readable your text is (textual readability) and how rich it is (textual richness)

Performances of the Approaches:

For sentence segmentation in speech, performance is usually evaluated using the error rate (ratio of number of errors to the no of example). F₁-measure (The harmonic mean of recall is defined as the ratio of the no. of correctly returned sentence boundaries to the

$$\begin{aligned} [x_{-4}] &= [9 - 8.4/6 - 8.4, 9 - 8.4/10 - 8.4] \\ &= [0.6, -2.4/10.6, -0.4/1.6] \end{aligned} \quad (4)$$

Hybrid Approaches:

The Hybrid approach combine best of rule based and ML approaches with Natural language processing. It provides guidance for accurate analysis and ML

Data acquisition:

Collection is the process of getting the data to the model data can come from valid resources like data websites, data where houses, like strings of the data like audio and video cleaning and preprocessing:

Raw of data in most cases contains information that may not help in the analysis therefore after data acquisition is completely perform some cleaning and preprocessing that includes like removing out layers (Outliers)

Data analysis and Interpret:

It is important at the beginning of the analysis it reveals the information that means organizing the data to predict to analysis process and predicate the data using patterns and algorithm after this analysis process to interact the data for training set. To model for machine learning process.

Building ML Model:

Once data is ready for analysis process it start building machine learning model. Here the model is to provide plenty of data to learn and algorithm uses to learn about patterns.

no. of correctly returned sentence boundaries to the no. of
all automatically estimated sentence boundaries. (10)

filtering spam emails, language model, text-generate
for sentence segmentation in text, researchers have reported
error rate results on a subset of tree the wall street journal
corpus of about 27,000 sentences, Mikheev reports that wall
street that his rule-based system.

If we search hoc it will give results many similar
words usmg named entity recognisatlon.

Syntactic Analysis:

Parsing uncovers the hidden structure of linguistic input. In many application involving natural language, the underlying predicate argument structure of sentence can be useful. The syntactic analysis of language provides a means to explicitly discover the various predicate argument dependencies that may exist in a sentence.

Parsing Natural Language:

In a text-to-speech application, input sentence are to be converted to a spoken output that should sound like it was spoken by a native speaker of the language.

Parsing:

modeling the data in tree format

Parse → Latin word

Parsing → Tree → Named entity recognition

We have two types of parsing

1. Bottom up approach

2. top down approach

Example:

$E \rightarrow T$ (expression)

$T \rightarrow T * F$ (traverse)

$F \rightarrow T$ (field)

$F \rightarrow id$ (field)

$id * id$ (Identifier)

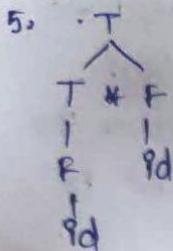
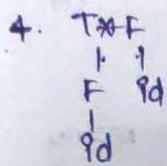
1. $id * id \rightarrow F$ (id)

2. $F * id \rightarrow T$ (id replaced with F)

$\frac{1}{id}$

3. $T * id$

$\frac{1}{F}$
 $\frac{1}{id}$



Top down approach.

$$S \rightarrow VP$$

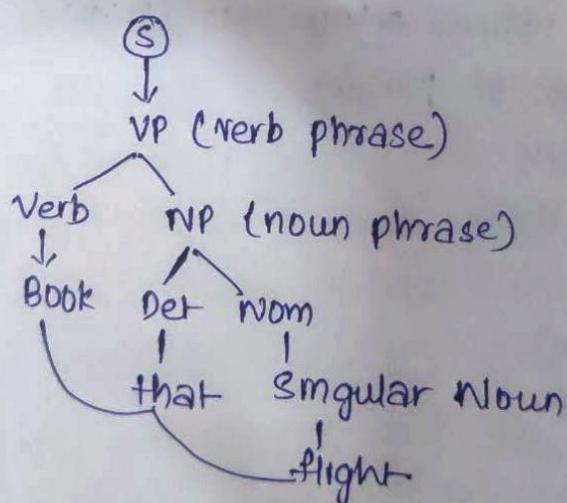
$$VP \rightarrow Verb \text{ or } NP$$

$$NP \rightarrow Det \text{ or } Nom$$

$$Det \rightarrow That$$

$$Nom \rightarrow Singular\ Noun$$

$$Singular\ noun \rightarrow flight$$

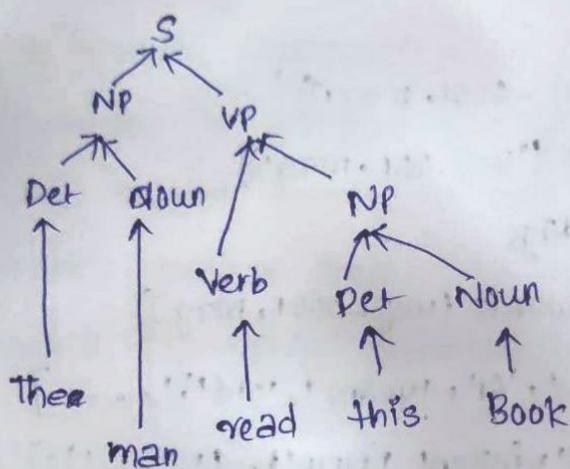


It is also known as recursive parsing starts from the "start" symbol and transforming into the input symbol mean parsing the nodes or extension of the links from top to bottom

Bottom up parsing

It is used to construct a partstree from an input string. Parsing starts with the input symbol and construct the partstree input string (bottom) up to the start symbol (top)

Bottom up approach



Tree bank: A Data Driven Approach to Syntax

To construct a tree using lexical analysis and maintaining combination of parts of speech to the data. Here parser is a compiler that used to break data into smaller elements in lexical analysis phase. A parser takes input sequence of elements produce output in the form of tree using text to speech application

construct partstree top down and bottom up to the following sentence using example note.

Treebank.nlp

```
# download required files from nltk
```

```
import nltk
```

```
nltk.download('treebank')
```

```
True
```

```
# check what's in the treebank  
from nltk.corpus import treebank  
print(treebank.fileids())  
print(len(treebank.fileids()))
```

out: ['wsj-0001.msg', 'wsj-0002.msg', ..., 'wsj-009.msg', '99']

see words

```
print(treebank.words('wsj-0001.msg'))
```

```
print(len(treebank.words('wsj-001.msg')))
```

see word with tags

```
print(treebank.tagged_words('wsj-0001.msg'))
```

out: ['pierre', 'vinken', ',', ',', '61', 'years', 'old', ',', ...]
[('pierre', 'NNP'), ('vinken', 'NNP'), ('', 'PRT'), ('', 'ADJ')]

See the parsed sentence tree

```
print(treebank.parsed_sents('wsj-0001.msg')[0])
```

```
print(* * * * * * * * *)
```

```
print(treebank.parsed_sents('wsj-0001.msg')[1])
```

out: (S(NP(pierre, NNP) (vinken NNP))

(ADJP(NP(Det 61) (NS years)) (N old))

(NP(VB will) (VB join))

(NP(Det the) (NN board))

(NP(as) (VB a) (Det Non-executive) (NN Director))

pierre vinken 61 years old will join

The board as a Non-executive Director

* * * * *

T(S(NP(NNP vinken))

(NP(98)(NN chairman))

(VB of)(Det the))

(NNP Daten)(VB publishing))

(NN group))

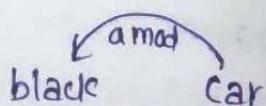
Representation of syntactic structure.

Syntax Analysis using Dependency Graphs

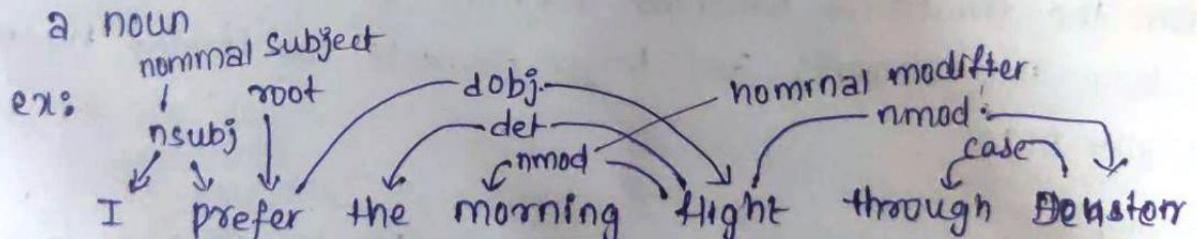
The main philosophy behind dependency graphs is to connect a word - the head of a phrase - with the dependents in that phrase. The notation connects a head with its dependent using a directed connection.

It is a process to analyze the grammatical structure and find out related words as well as the type of the relationship between them.

Example dependency relation between two words



There exists a relationship between car and black because black modifies the meaning of car. Here "car" acts as the "head" and "black" is a "dependency" of the "head". The nature of relationship here is "a mod" stands for "Adjective modifier". It is an "adjective" or an "adjective phrase" that modifies a noun.



The syntactic structure of sentence is described solely in terms of the words and a set of directed binary grammatical relations.

The relations between every word of the sentence is indicated using directed arc in a typed dependency structure.

- Here, the root the tree "prefer" form head of above sentence
- The relationship between any 2 word is marked a dependency tag. for instance, the word "flight" modifies meaning the noun "Denver"
- You can notice a dependency from 'flight' \rightarrow Denver where the flight is head Denver is child or dependent.
- It is denoted by nmod which is nominal modifier

I \Rightarrow nsubj

prefer \Rightarrow root

The \Rightarrow det \Rightarrow flight

morning \Rightarrow n mod \Rightarrow flight

Flight \Rightarrow dobj \Rightarrow prefer

through \Rightarrow case \Rightarrow denver

Denver \Rightarrow nmod \Rightarrow flight

Syntax Analysis using phrase structure trees,

A phrase structure syntax analysis of the sentence derives from the traditional sentence diagrams in that partitions of a sentence into elements and larger elements are formed by merging smaller ones

Example

- It tooks me more than two hours to translate a few pages of english .

The sentence identify pos tag and the root node provide proper noun phrase (ansubj - ansubject) and linked with using grammar tree

The another example sentence "the team is not performing well in the match"

To identify pos tag and dependency parsing and also provide phrase structure tree.

Program:

```
import spacy  
nlp = spacy.load("en")  
doc = nlp("The team is not performing well in the match")  
for token in doc:  
    print(str(token.text), str(token.pos_), str(token.dep_),  
          str(token.head_))
```

Point 1)

>>> The, Noun, nsubj => performing
team, Noun, nsubj (root) => performing
is, Verb, aux => performing
not, ADV, neg => performing
Performing, Verb, root => team
well, ADV, Advmod => performing
in, ADP, Prep => performing
The, Noun, pobj => In
Match, Noun, pobj => In

Dependency parsing:



Phrase structure tree

(S (NP (N the) (N team),
(VP (V Is) (ADV not),
(VP (V performing)
(ADV well)
(ADP Pn),
(NP (N The) (N match))))))

Parsing Algorithm:

To given an input sentence a parser produces an output analysis of that sentence. Here parser do not need an explicit grammar consider the following simple sentence can be derive (a and b or c) using and, or symbols

To start symbol with N, $N \rightarrow 'a' / 'b' / 'c'$

Here ' \rightarrow ' the symbol represents a sequence of steps called derivation

this sentence actions separated with symbols and append with N

$N \rightarrow 'N' \text{ 'and'}' N$

$N \rightarrow N \text{ 'or'}' N$

Here the sentence derived from the start symbol N to expand with sentence actions.

$\Rightarrow N \text{ 'or'}' N$

$\Rightarrow N \text{ 'or'}' c$

$\Rightarrow N \text{ 'and'}' N \text{ 'or'}' c$

$\Rightarrow N \text{ 'and'}' b \text{ 'or'}' c$

$\Rightarrow a \text{ 'and'}' b \text{ 'or'}' c$

form $N \rightarrow 'a' / 'b' / 'c'$

$N \rightarrow N \text{ and } N$

$N \rightarrow N \text{ or } N$

$N \rightarrow a \text{ and } b \text{ or } c$

The sentence to expand with right most sequence form we arrange the derivation following rule.

Right side Rule

a and b or c

$\Rightarrow 'N' \text{ and } b \text{ or } c \# N \Rightarrow a$

$\Rightarrow 'N' \text{ and } 'N' \text{ or } c \# N \Rightarrow b$

$\Rightarrow 'N' \text{ and } 'N' \text{ or } 'c' \# N \Rightarrow 'N' \text{ and } 'N'$

$\Rightarrow 'N' \text{ or } 'c' \# N \Rightarrow c$

$\Rightarrow N \text{ or } 'c'$

$\Rightarrow N \text{ or } N \# N \Rightarrow 'N' \text{ or } N$

$\Rightarrow N$

Shift-Reduce parsing:

To build a parser, we need to create an algorithm that can perform the steps in the preceding rightmost derivation for any grammar and for any input string. Every CFG turns out to have an automaton that is equivalent to it, called a pushdown automaton.

Parse Tree	Stack	Input	Action
a		a and b or c	init
a	a	and b or c	shift a
(Na)	N	and b or c	reduce $N \Rightarrow a$
(Na) and	N and	b or c	shift and
(Na) and (Nb)	N and b	or c	shift b
(N(Na)) and (Nb))	N	or c	reduce $N \Rightarrow b$
(N(Na)) and (Nb)) or	N or	c	shift or

$(N(Na) \text{ and } (Nb)) \text{ or } c$	$N \text{ or } \emptyset$	reduce $\emptyset \rightarrow N$ or shift c
$(N(Na) \text{ and } (Nb)) \text{ or } (Nc)$	$N \text{ or } N$	reduce $N \rightarrow \emptyset$
$(N(N(Na) \text{ and } (Nb)) \text{ or } (Nc))$	N	reduce $N \rightarrow N \text{ or } N$
$(N(N(Na) \text{ and } (Nb)) \text{ or } (Nc))$	N	Accept!

Shift reduce parsing algorithm for the input "a and b or c" for the grammar rules using and append with N symbol

To construct partstree following section

a and b or c

$(Na) \cdot \#$ using $N \rightarrow a$

$(Na) \text{ and }$

$(Na) \text{ and } b$

$(Na) \text{ and } (Nb) \# \text{ use } N \rightarrow b$

$(N(Na) \text{ and } (Nb)) \# \text{ use } N \rightarrow (Na) \text{ and } (Nb)$

$(N(Na) \text{ and } (Nb)) \text{ or }$

$(N(Na) \text{ and } (Nb)) \text{ or } c$

$(N(Na) \text{ and } (Nb)) \text{ or } (Nc) \# N \rightarrow c$

$(N(N(Na) \text{ and } (Nb)) \text{ or } (Nc))) \# \text{ use } N \rightarrow (N(Na) \text{ and } (Nb)) \text{ or } (Nc))$

Allocation.

To store the data in the memory to represent the size by using array

$$\cdot \quad a[5] = \begin{array}{|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 \\ \hline 0 & 1 & 2 & 3 & 4 \\ \hline \end{array}$$

Index 0 - starts

Index 4 - ends

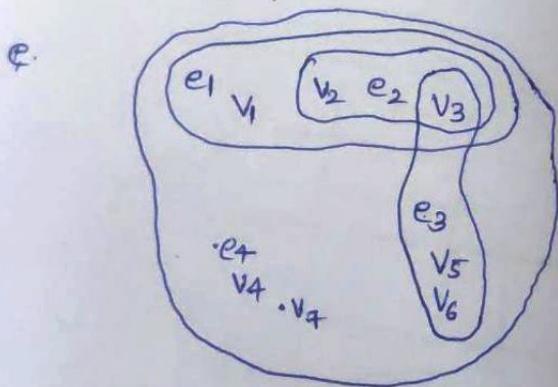
Hyper Graphs: and chart parsing:

Shift-reduce parsing allows a linear time parse but requires access to an oracle. For general CFGs in the worst case, such a parser might have to resort to backtracking, which means reparsing the input, which leads to a time that is exponential in the grammar size in the worst case.

Hyper graph is the graph of in which an edges can join any no. of vertices. Here we have two types

1. Undirected graph
2. Directed graph.

Undirected graph:



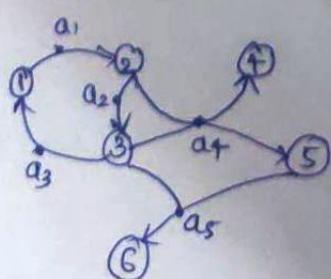
$$e_1 = \{v_1, v_2, v_3\}$$

$$e_2 = \{v_2, v_3\}$$

$$e_3 = \{v_3, v_5, v_6\}$$

$$e_4 = \{v_4\}$$

Directed graph:



$$n = \{1, 2, 3, 4, 5, 6\}$$

$$E = \{e_1, e_2, e_3, e_4, e_5\}$$

$$a_1 = \{(1, 2)\}, a_2 = \{(2, 3)\}$$

$$a_3 = \{(3, 1)\}, a_4 = \{(2, 4)\}$$

$$a_5 = \{(5, 6)\}$$

Hyper graph parsing:

Example: Input grammar rules

$N \rightarrow 'N'$ and ' N '

$N \rightarrow [a] \wedge \dots \wedge N$ or ' N '

$N \rightarrow [a, b, c]$

$N \rightarrow NN$

$NN \rightarrow 'and' \wedge N$

$N \rightarrow N NV$

$N \rightarrow 'or' \wedge N$

$N \rightarrow [a, b, c] \wedge N$

$N[0, 5] \rightarrow N[0, 1] \wedge N^*[1, 5]$

$N[0, 3] \rightarrow N[0, 1] \wedge N^*[1, 3]$

$N^*[1, 3] \rightarrow 'and' \wedge [1, 2] \wedge N[2, 3]$

$N^*[1, 5] \rightarrow 'and' \wedge [1, 2] \wedge N[2, 5]$

$N[0, 5] \rightarrow N[0, 3] \wedge NV[3, 5]$

$N[2, 5] \rightarrow N[2, 3] \wedge NV[3, 5]$

$NV[3, 5] \rightarrow 'or' \wedge [3, 4] \wedge N[4, 5]$

$N[0, 1] \rightarrow 'a' [0, 1]$

$N[2, 3] \rightarrow 'b' [2, 3]$

$N[4, 5] \rightarrow 'c' [4, 5]$

CKY Algorithm. (Coicke, Kasami, younger)

Rules

$S \rightarrow NP VP [0.8]$

$NP \rightarrow Det N [0.3]$

$VP \rightarrow V NP [0.2]$

$V \rightarrow \text{Includes} [0.05]$

$Det \rightarrow \text{The} [0.4]$

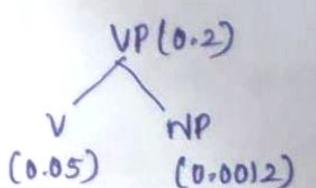
$Det \rightarrow a [0.4]$

$N \rightarrow \text{meal} [0.01]$

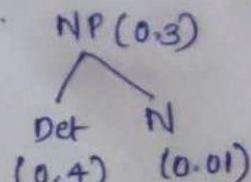
$N \rightarrow \text{flight} [0.02]$

	1	2	3	+	5
0	$Det \leftarrow$	$NP \leftarrow$			$S \downarrow$
1			$N \downarrow$		
2			$N \leftarrow$	$VP \downarrow$	
3				$Det \leftarrow$	$NP \downarrow$
4					$N \downarrow$

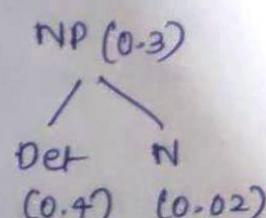
Probabilistic CKY



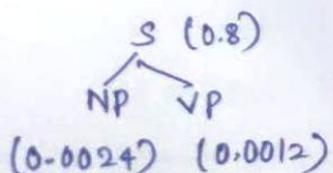
$$= 0.2 \times 0.05 \times 0.0012 \\ = 0.000012$$



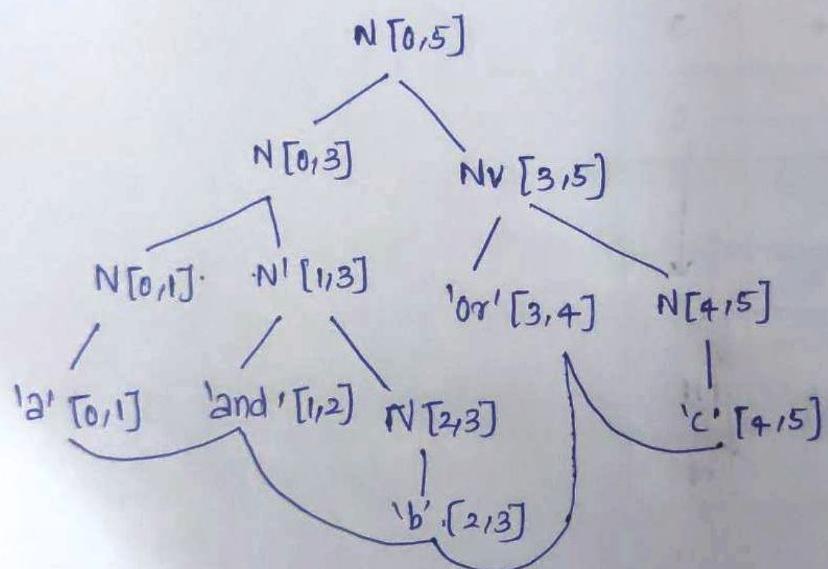
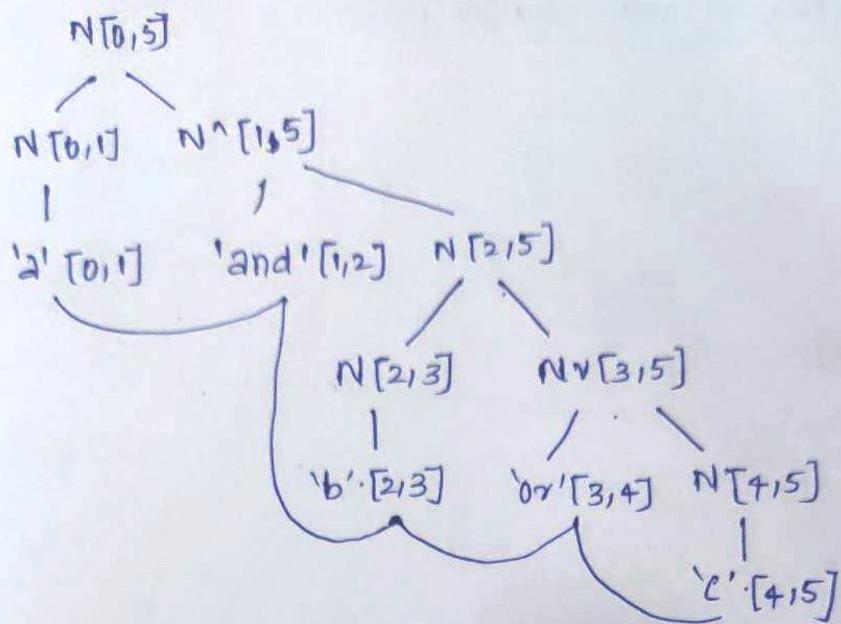
$$= 0.3 \times 0.4 \times 0.01 \\ = 0.0012$$



$$= 0.3 \times 0.4 \times 0.02 \\ = 0.0024$$

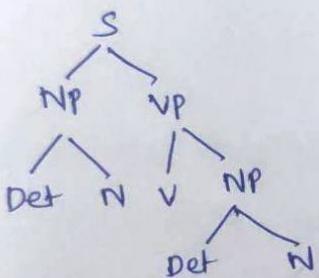


$$0.8 \times 0.0024 \times 0.000012$$



	1	2	3	4	5
0	Det ← 0.4	NP ← 0.0024 ↓ N		S	2.30
1					
2		✓ ← 0.05		VP	0.00002
3			Det ← 0.4	NP 0.0012	
4				N	0.01

Tree



The sentence parsing with set of grammar rules to construct parsing tree in this algorithm content free grammar rules converted to CNF by using CNF rules (comskey Normal form) and it is parsing in matrix from or chart table

CNF Rules:

Here using same CFGI rules but some of the rules are changes

e.g. $NP \rightarrow$ The combination of direct terminal and non-terminal changes to direct terminal to non terminal

$NP \rightarrow Det\ N$

$NP \rightarrow Det\ noun$

$pp \rightarrow changes\ to$

$NP \rightarrow Det\ nom$
 $nom \rightarrow noun\ pp.$

Spanning Tree Algorithm:

- Minimum Spanning Trees and dependency parsing.

$F \leftarrow \Gamma \rightarrow$ function

$T \leftarrow \Gamma \rightarrow$ constant tree

Source $\Gamma \rightarrow$ final score

for each $v \in V$

$\text{Min Edge} \leftarrow \arg \min_{e \in E} \text{Score}[e] \rightarrow$ select min edge weight

$F \leftarrow F \cup \text{MinEdge}$

for each $e = (u, v) \in E$ do \rightarrow stop condition

if $T = (V, F)$

else

return it spanning tree

$c \leftarrow$ a cycle in F

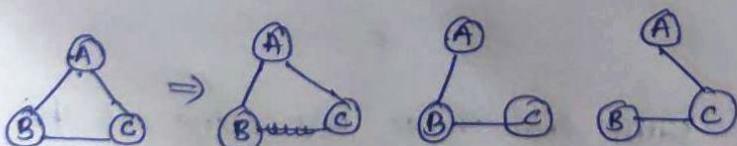
$G_1 \leftarrow \text{contract}(G_1, c) (G_1, \text{root}, \text{Score})$

$T \leftarrow \text{EXPAND}(T, c) \rightarrow$ Expanded graph

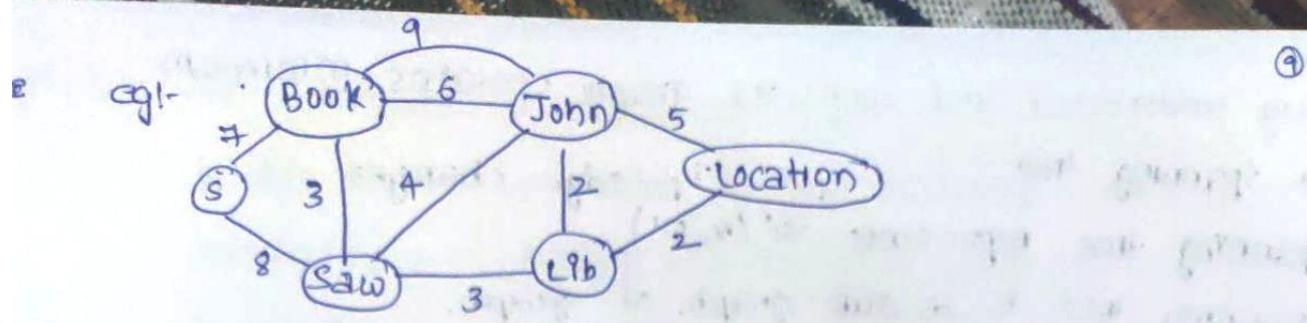
Return T

Minimal Spanning Tree

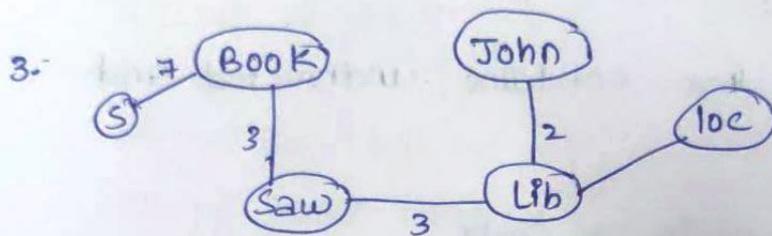
1. There is no connection or loop forming in the spanning tree



minimum spanning tree based on connected graph in directed graph to choose minimum weight of edge in the graph



1. Select the root node is!
2. Remove all the loop in the graph.



trace the tree with minimum weight if we have same weight we can consider two edges without forming loop

$$7 + 3 + 3 + 2 + 2 = 17 \text{ [minimum weight]}$$

2nd method.

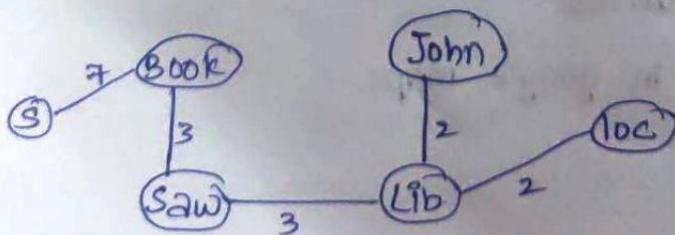
Arrange all the edges ascending order of weight

$$\{ \{ \text{John}, \text{Lib} \}, \{ \text{Lib}, \text{Loc} \}, \{ \{ \text{Book} \}, \text{Saw} \}, \{ \text{Saw}, \text{Lib} \}, \{ \text{Saw}, \text{John} \} \}$$

2 2 3 3 4

$$\{ \{ \text{John}, \text{Loc} \}, \{ \text{Book}, \text{John} \}, \{ \{ \text{S}, \text{Book} \}, \{ \text{S}, \text{Saw} \}, \{ \text{Book}, \text{John} \} \}$$

5 6 7 8 9



$$7 + 3 + 3 + 2 + 2 = 17$$

- Every undirected and connected graph contains minimum one spanning tree \rightarrow edge changes
- Spanning tree represents $G'(N, E')$
- Spanning tree is a sub graph of graph.

Properties:

- In directed graph, consider minimum edge weight of graph
- the rule of spanning tree contains undirected and connected graph
- Spanning tree never contain a cycle
- If it contains cyclic graph it can divided into undyclic possible graph
- possible to have more than minimum spanning tree the weight of same edge are same in the graph.

Applications:

- routing protocols
- Civils networks
- Clustering process
- communication b/w group
- to calculate distance b/w nodes
- Telecommunication networks
- finding shortest paths in google maps.

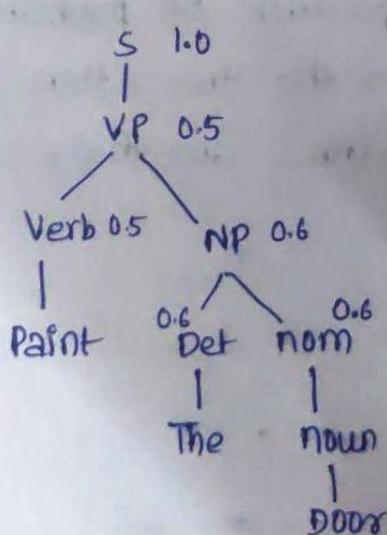
Models for Ambiguity Resolution in parsing.

In this section we focus on the modeling aspect of parsing; how to design features and ways to resolve ambiguity in parsing.

Probabilistic context-free grammar.

$$\begin{aligned}
 S &\rightarrow NP \quad VP \quad 0.8 \\
 S &\rightarrow VP \quad 0.2 \\
 \} &\qquad\qquad\qquad 1.0 \\
 NP &\rightarrow Det \quad noun \quad 0.6 \\
 NP &\rightarrow Proper \quad Noun \quad 0.2 \\
 \} &\qquad\qquad\qquad 1.0 \\
 NP &\rightarrow prep \quad noun \quad 0.1 \\
 \} &\qquad\qquad\qquad 1.0 \\
 NP &\rightarrow Det \quad NP \quad 0.1 \\
 VP &\rightarrow Verb \quad NP \quad 0.5 \\
 VP &\rightarrow Verb \quad 0.3 \\
 \} &\qquad\qquad\qquad 1.0 \\
 VP &\rightarrow NP \quad PP \quad 0.2 \\
 PP &\rightarrow prep \quad p \quad NP \quad 1.0
 \end{aligned}$$

Paint the door



$$\begin{aligned}
 P(D) &= 1.0 \times 0.5 \times 0.5 \times 0.6 \times 0.6 \times 0.2 \\
 &= 0.018
 \end{aligned}$$

Assigning probabilities to possible parser of a sentence

It is also known as probabilistic context free grammar

In this assign probabilities most likely parses and tagging the

POS provide final resultant to each phrase is "1" means

Partition each phrase with assigning value to possibilities
of the phrase and final resultant is the Output

Like find \Rightarrow Assign \Rightarrow Return

Output: PCFG1

Example : The sentence

"paint the door with the hole"

Grammar Rules:

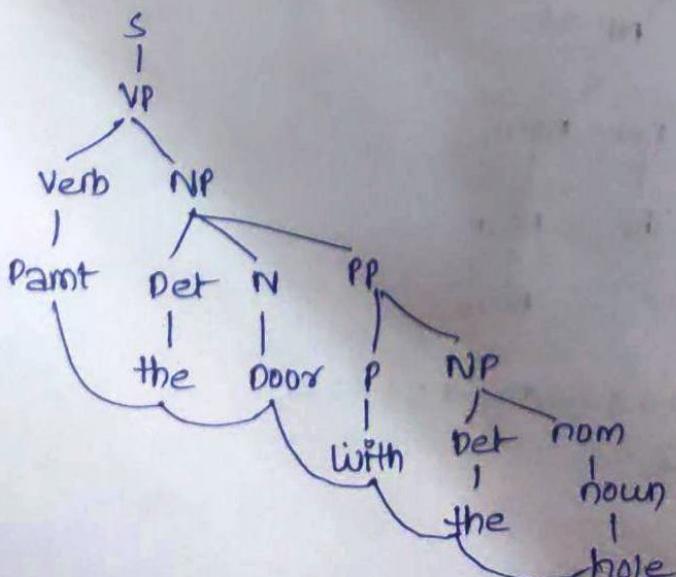
Noun \rightarrow paint / door / hole / bird

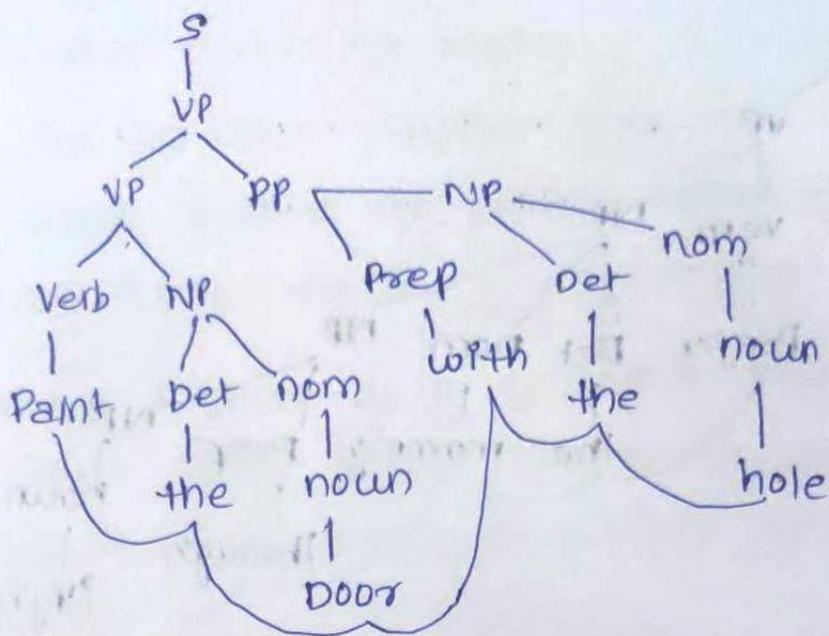
Verb \rightarrow sleeps / sings / open / saw / paint

Prep \rightarrow from / with / on / to

Pronoun \rightarrow she / he / they

To construct parsing tree with product of probability values and calculate derivation to the tree. Here to construct CFG1 then two possibilities are their.





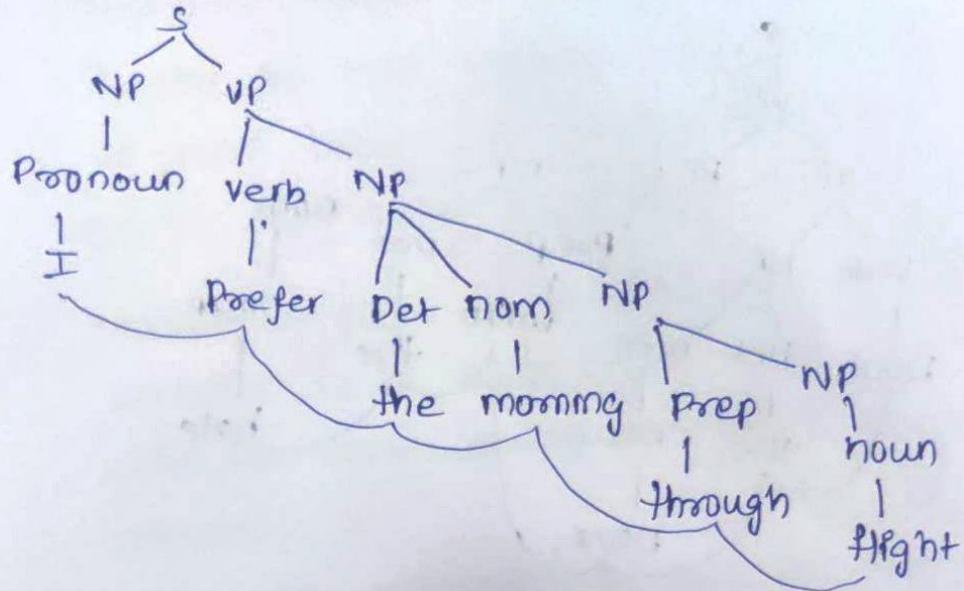
$NP \rightarrow Det\ Nom\ 0.3$
 $NP \rightarrow P\ noun\ 0.2$
 $NP \rightarrow Pronoun\ 0.2$
 $NP \rightarrow Det\ N\ P\ 0.3$

$$\begin{aligned}
 PCFG_1 &= P(D_1) + P(D_2) \\
 &= \text{Result}
 \end{aligned}$$

The probability context free grammar
and possibility

I prefer the morning through flight

$S \rightarrow NP\ VP$
 $NP \rightarrow Det\ nom$
 $NP \rightarrow Det\ nom\ NP$
 $NP \rightarrow Prep\ NP$
 $NP \rightarrow Pronoun$
 $NP \rightarrow Verb\ NP$



The giraffe dreams.

$S \rightarrow NP \ VP$

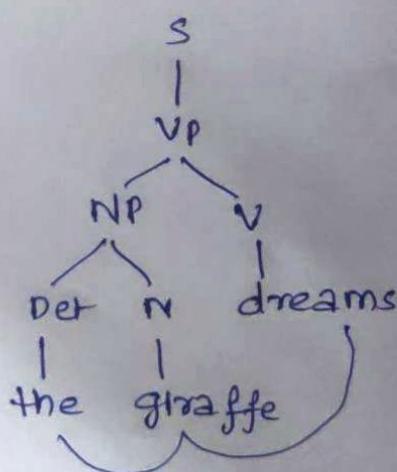
$NP \rightarrow Det \ N$

$VP \rightarrow V \ NP$

$det \rightarrow \text{the} / \text{a} / \text{an}$

$n \rightarrow \text{giraffe} / \text{apple}$

$v \rightarrow \text{dreams} / \text{cats}$



Generative Models for parsing. (12)

To find the most plausible parse tree, the parser has to choose between the possible derivations each of which can be represented as a sequence of decisions.

Modeling a parse as a generative process allows the partstree to be broken into manageable parse. for which the corresponding probabilistic can be reliable estimated.

Probability estimation is easy for this sorts of models but choosing how to breakup the parse is something of a black art.

PCFG provide a redimate solution to this statical parsing Problem.

Representation.

1. The set of parts of Speech tags
2. Whether to pass lexical heads of the tree.
3. whether to replace words with their morphological stems

Decomposition:

1. The Order in which to generate the tree.
2. In order of decision made in generating the tree
3. These decisions do not have to correspond to parsing decisions.

Assumptions:

1. Each rule in generative model

$$P(t) = \prod_{i=1}^n q_i(a_i \rightarrow b_i)$$

where $q_i(a_i \rightarrow b_i)$

The probability for rule using $A \rightarrow B$

$$q(VP \rightarrow VT\ NP) = 0.4$$

$$S \Rightarrow NP \quad NP \quad 1.0$$

$$VP \rightarrow V\{(\text{verb}) \quad 0.4$$

$$VP \rightarrow N \quad NP \quad 0.4$$

$$VP \rightarrow VP \quad PP \quad 0.2$$

$$NP \rightarrow DT \quad NN \quad 0.3$$

$$NP \rightarrow NP \quad PP \quad 0.7$$

$$PP \rightarrow P \quad NP \quad 1.0$$

Lexicon

$$V\{ \rightarrow \text{sleeps} \quad 1.0$$

$$VT \rightarrow \text{saw} \quad 1.0$$

$$NN \rightarrow \text{man} \quad 0.7$$

$$NN \rightarrow \text{woman} \quad 0.2$$

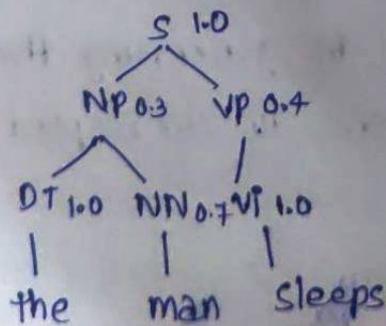
$$NN \rightarrow \text{telescope} \quad 0.1$$

$$DT \rightarrow \text{The} \quad 1.0$$

$$IN \rightarrow \text{The} \quad 0.5$$

$$IN \rightarrow \text{in} \quad 0.5$$

The man sleeps



$$1.0 \times 0.3 \times 0.4 \times 1.0 \times 0.7 \times 1.0$$

$$= 0.084$$

Generative Rule

(13)

$$P(t) = \prod_{i=1}^n q_i(\alpha_i \rightarrow \beta_i)$$

$$\text{where } q_i(\alpha_i \rightarrow \beta_i)$$

$$NP \rightarrow DT \quad 0.3$$

$$VP \rightarrow V \quad 0.4$$

$$VP \rightarrow NP \ PP \quad 0.2$$

Derivation	Rule	Prob
S	NP VP	1.0
NP NP	NP \rightarrow DT NN	0.3
DT NN NP	DT \rightarrow the	1.0
The NN VP	NN \rightarrow man	0.7
The man NP	VP \rightarrow V	0.4
The man V	V \rightarrow sleeps	1.0
The man sleeps		

$$P(t) = 1.0 \times 0.3 \times 1.0 \times 0.7 \times 0.4 \times 1.0$$

$$= 0.084$$

Discriminative model for parsing:

Discriminative model is a class of logistic model used for classification or regression

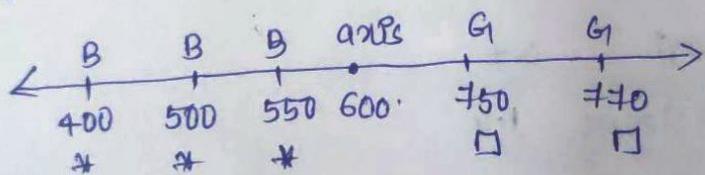
It makes predictions based on conditional probability on the data the distribution of data set to return by using example letters determined representing CIBIL score in HDFC Bank customers.

- * Here axis point is known as linear regression boundary

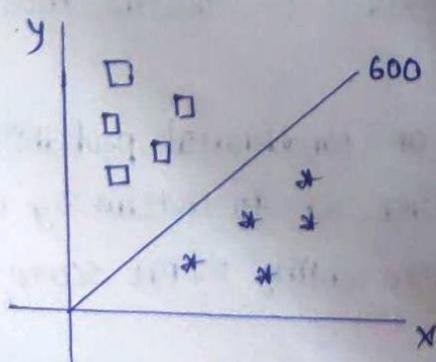
S.NO	CIBIL	Decision
1	500	B
2	550	B
3	750	G
4	770	G

Regression line estimate Average value in a given data set. It slightly differentiate sufficient and Insufficient predictions in the data set.

Here dataset contains another dimension the differentiate by using x and y axis line is known as quantile regression or necessary condition analysis. In this situation regression created boundary analysis.



Years	S.NO	CIBIL	Decision
2	1	500	B
4	2	550	B
3	3	750	G
5	4	770	G



The data set contains collection of sufficient and insufficient predicates and recreate using amis boundary in existing graph. ④

Multilingual Issues :

Tokenization, case and encoding.

Tokenization:

It is a process of converting a sequence of characters into tokens. It performs lexical analysis contains lexer or tokenizer. These generally combine with a parser together analyse the syntax of program.

Lexical Syntax:

The regular expression alphabets consists of single characters of source code

Grammatical Syntax:

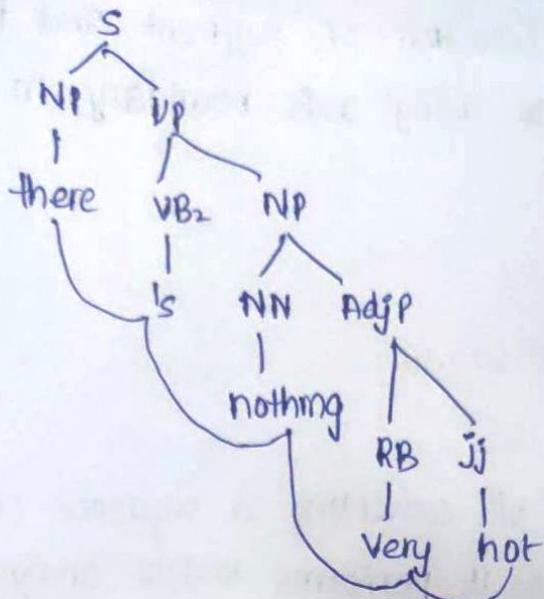
It is usually a context free language whose alphabet consist of tokens produced by lexer.

Let us determine example sentence

e.g. there is nothing very hot

tokenization :

s ((NP (There))
 (VP (VB2's))
 (NP (NN nothing))
 (Adjs (RB very))
 (JJ hot)))))



Tokenization: Classify sentence into smallest units

Tokenize produces a tree instead of grammar being return in natural language in terms of individual characters in the parsing

Word Segmentation

This is a technique to divide set of words in a sentence we generally compute a score for each word to signify importance in the sentence. This used in information retrieval and text mining.

The segmentation measures frequency of the word in a sentence this highly depends on the sentence length.

Eg: S_1 : Earth is the third planet from the sun

S_2 : Jupiter is the largest planet

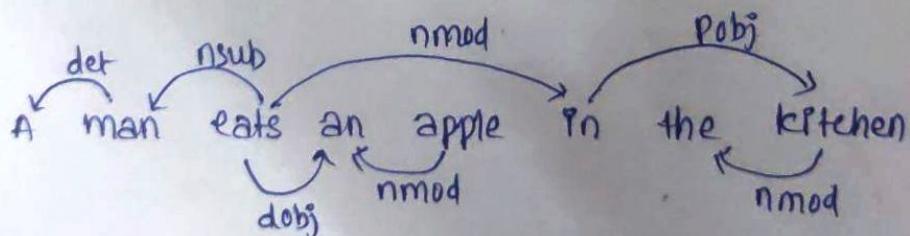
Word	S_1	S_2	frequency
earth	1/8	0/5	$\log(2/1) = 0.3010$
is	1/8	1/5	$\log(2/2) = 0$
the	2/8	1/5	$\log(2/3) = -0.1760$

third	1/8	0/5	$\log(2/1) = 0.3010$
Planet	1/8	1/5	$\log(2/2) = 0$
from	1/8	0/5	$\log(2/1) = 0.3010$
sun	1/8	0/5	$\log(2/1) = 0.3010$
Jupiter	0/8	1/5	$\log(2/1) = 0.3010$
largest	0/8	1/5	$\log(2/1) = 0.3010$

Dependency parsing analysis shows syntactic dependencies
need to be aware of the morphemes with in the words
e.g- morpheme boundaries with in a boundary are
shown using '+'

Morphology :

In many languages the notion of splitting up tokens using spaces is problematic because each word can contain several components, called morphemes, such that the meaning of a word can be thought of as composed of the combination of the meaning of the morphemes.



the + kitchen \Rightarrow NP \Rightarrow Det N

A + man \Rightarrow NP \Rightarrow Det N

The is learning piano \Rightarrow learning + piano
 $\text{VP} = \text{V} + \text{N}$

semantic Analysis

- NLP is a field which aims to give machine the ability to understand natural language.
- semantic analysis is a sub-topic of it.
- semantic analysis means giving exact meaning of text.
- Role of semantic analysis is to check text for meaningfulness.

semantic Analysis

- meaning of individual word
(lexical semantics)
- combination of words,
(sentences)
 - Relationship exist between words
 - The word order
 - context
 - semantic structure
 - Real world knowledge.

Composition Semantics

It involves how word combines to form larger meaning.

→ It says that meaning of each word matters. Syntax and way in which the sentence are constructed plays an important role as well.

Ex: I like you and you like me are same words but meaning are different

I → subject | you → sub

y → obj | me → obj

→ lexical semantic + composition semantic =

Semantic analysis

→ called principle of compositionality.

Basic Building blocks of semantic system

Entities: It represents the individual such

as particular person, location, etc.

ex: Ramu, Chennai.

Concepts: It represents general category of

Individuals: like person, vehicle, animal

Relations: relations between entities and concepts

ex: chennai is a location.

Predicates: It represent Verb structure.

e.g.: case grammar and semantic rules.

case grammar: A form of grammar in which structure of sentence

Approaches:

→ first order predicate logic FOL

→ semantic nets

→ case grammar.

→ convert sentence into logical form

e.g.) Jack loves Jill

loves (Jack, Jill)

or pharma.

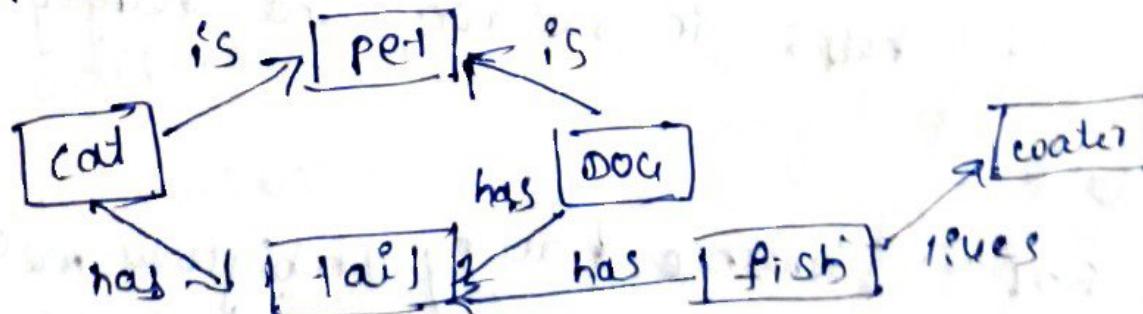
2) Suman takes Engg

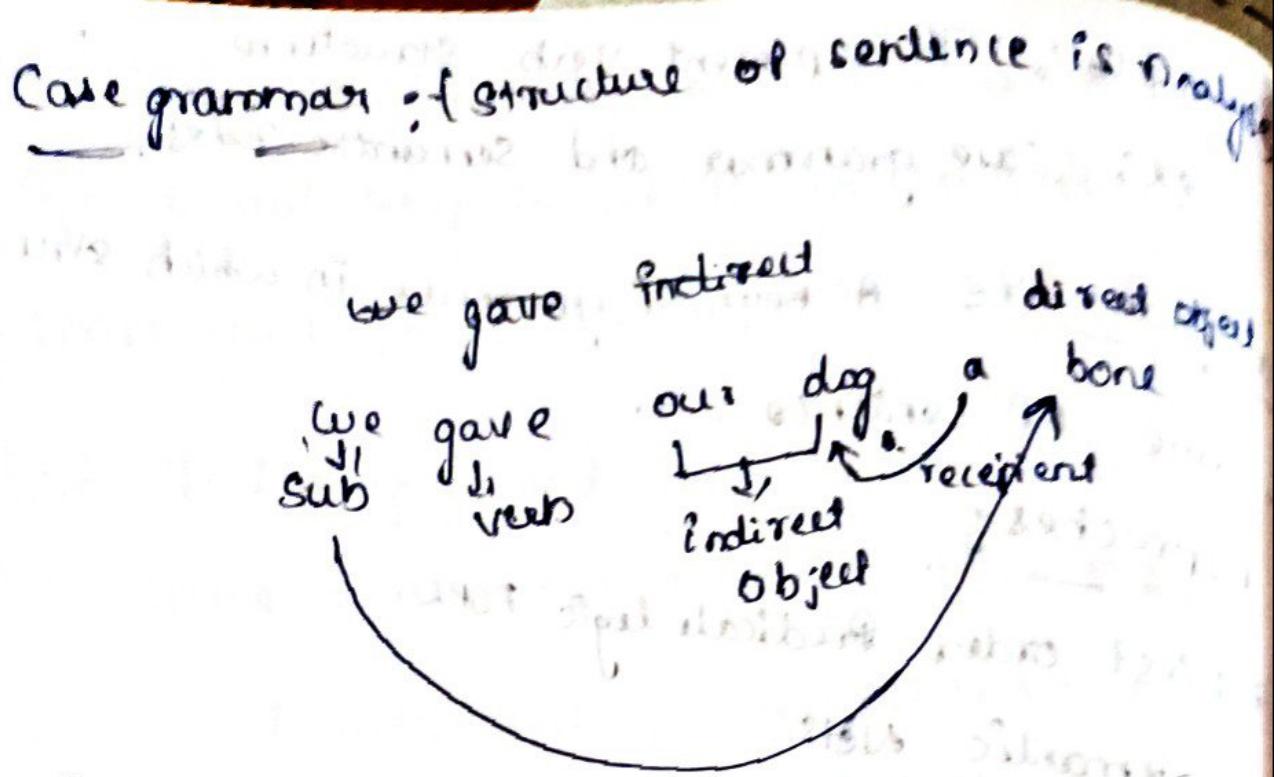
takes (suman, Engg) & takes (suman, pharm).

Semantic Nets:

→ It shows Semantic Relations b/w concept.

→ It is used for knowledge representation.





Early semantic Analysis is important means in real word posts, reviews, feedbacks, likes, dislikes purpose using this analysis.

Semantic Interpretations

- It helps in linking the linguistic elements of a sentence to the non-linguistic elements.
- It helps in representing unambiguous data at lexical level.
- It helps in reasoning and verifying data.
- Explain Sentence having ambiguous meaning and it should account

- 8) resolve the ambiguities of words in context
 - 9) identify meaningless but syntactically well-formed sentences
 - 10) identify syntactically or transformationally unrelated paraphrases of a concept having the same semantic content
- details about the Relations;

Hyponymy : It illustrates the connection b/w a generic term is known as hypernym, while the occurrences are known as hyponyms.

Homonymy : It may be described as words with the same spelling or form but diverse and unconnected meanings.

Polysemy : polysemy is a term or phrase that has a different but comparable meaning. To put it another way, polysemy has the same spelling but various and related meanings.

⇒ Antonymy: It is the relationship between two lexical items that includes semantic components that are symmetric with respect to an axis.

⇒ metonymy: It is described as a logical arrangement of letters and words indicating a comparison of or member or anything.

Structural Ambiguity:

When we talk of structure, we generally refer to the syntactic structure.

word sense disambiguation

It is an automatic process of identifying the context of any word, in which it is used in the sentence. In natural language one word can have many meanings.

e.g., The word "light" could be meant as not very dark or not this is done by word sense.

Relationship Extractions

In a sentence, there are a few entities that are co-related to each other. Relationship extraction is the semantic relationship b/w these entities. In a sentence,

"I am learning mathematics",
2 entities 'I' & 'mathematics' The relation
b/w them 'learn'.

Predicate Argument Structure

Predicate argument structure is based on the function features of lexical items (most often verbs). The function features determines the thematic roles to be played by the other words.

Ques. The police detained the suspect at w
who did what to whom at where
The police officer detained the suspect at the scene of
the crime

Understanding a sentence = "knowing" who
what (to whom, when, where, why ...).

Verbs correspond to predicates (what are
their arg (and modifiers) identify who did
to whom, where, when, why, etc).

Syntactic parsing:

The police officer detained the suspect at the scene.
~~~~~  
| | | |  
| | | |  
sub root direct obj modifier

### meaning Representation:

Verbs describe events or states &

agents - it is Tom broke the window with a rock.

The window broke.

The window was broken by Tom / by a rock.

We could translate verbs to (logical) predicates

er's Sub = first arg, Obj = second arg

does not capture that similarities in meaning.

break ( Tom, window, rock )

break ( window )

break ( window, Tom )

break ( window, rocks ).

### Semantic / thematic Roles

Verbs describe events or states

Tom broke the window with a rock

The window broke

The window was broken by Tom / by a rock.

The window was

Role(s) refer to participants of these events:

Agent ( who performed the action ): Tom

Patient ( who was the action performed on ): window

Tool/Instrument ( what was used to perform the action ): rock

Semantic Roles ( agent, patient ) are different

to Grammar Rules ( sub, obj ).

e.g. John broke the window ( or ) with a rock

Agent

Patient Instrument

Break ( Agent, patient ).

Break ( Agent, patient, instrument ).

## System Paradigms

### 1. system Architecture

(a) knowledge based: As the name suggests, these systems use a predefined set of rules of a knowledge base to obtain a solution to a new problem.

(b) unsupervised: These systems involve the methods. It is based on the fact that words that are related to each other be found in the definitions.

### (c) supervised methods

Sense-annotated corpora are used to train machine. The rules of the grammar will have an extra field to hold semantic attachment.

## Predicate Argument structure

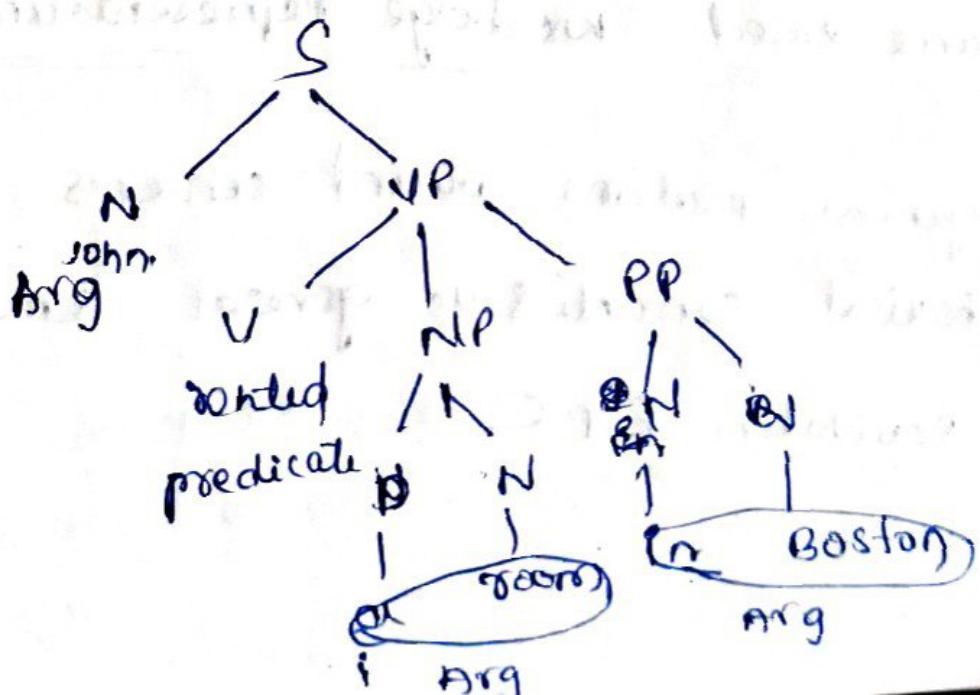
It is an Verb-centered approach

I want Turkish food NP want NP  
I want to spend less than five dollars NP want PNP<sup>10</sup>  
I want it to be close by here NP want NP In PNP

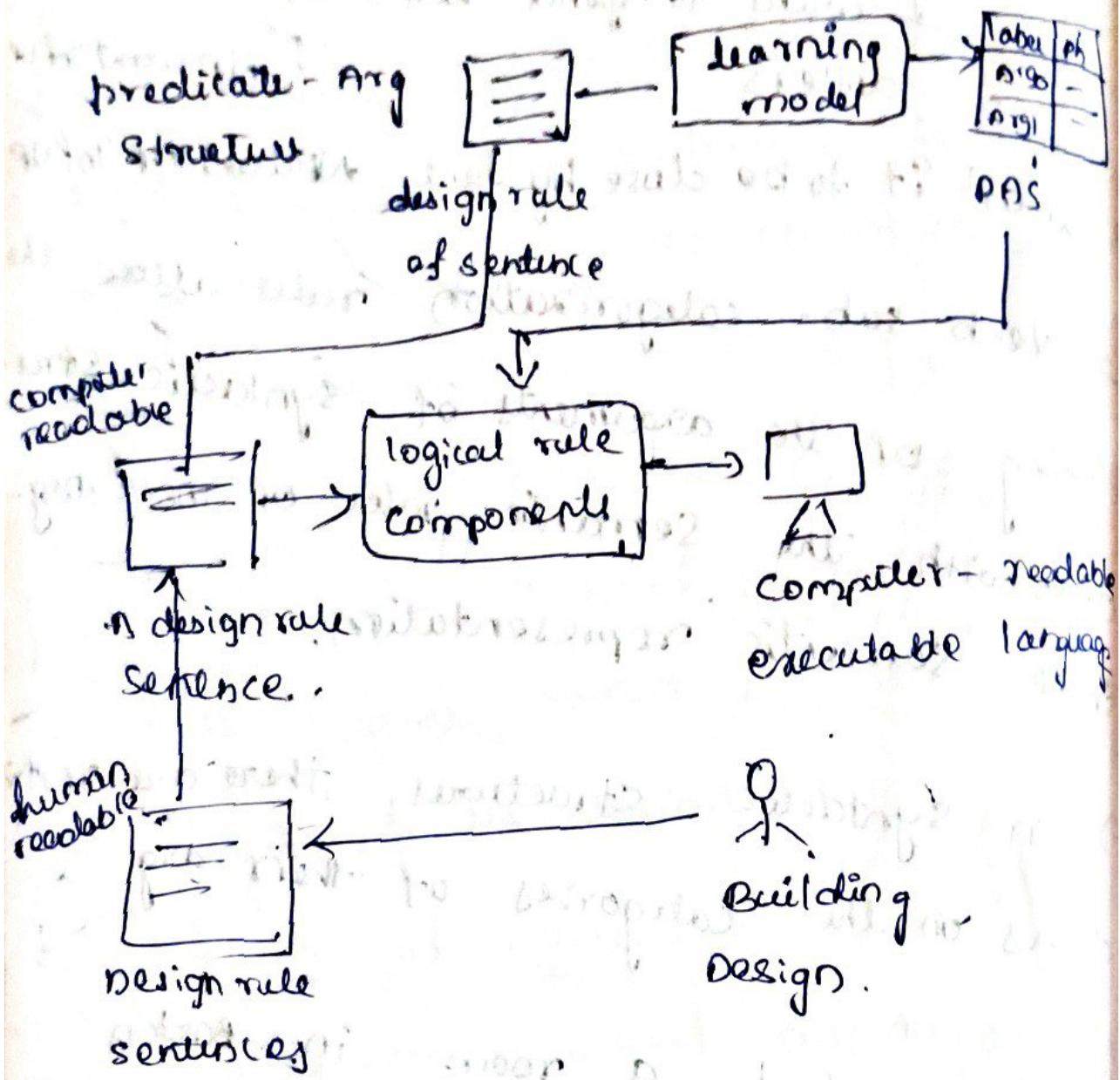
Verbs sub-categorization rules allow the linking of the arguments of syntactic structure with the semantic roles of these arguments in the semantic representation.

In syntactic structures, there are restrictions on the categories of their arguments on the categories of their arguments

John rented a room in Boston

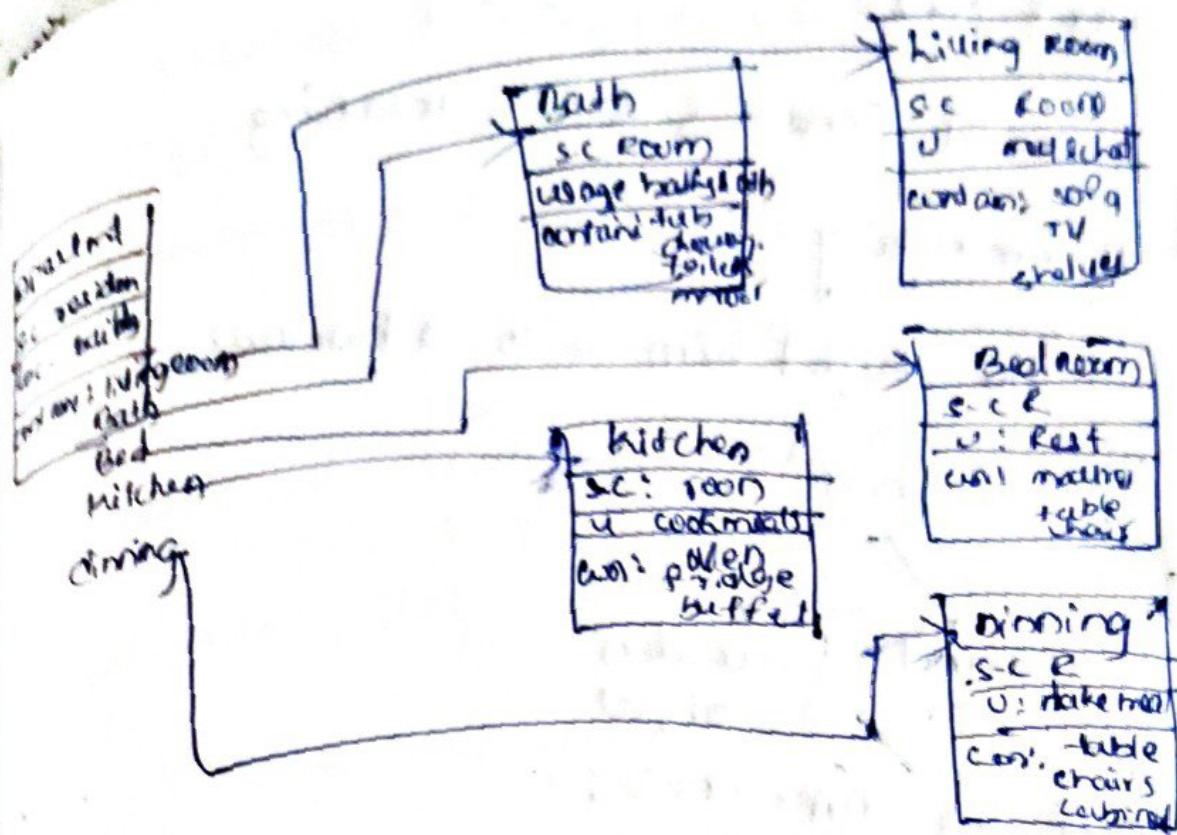


## Representation



## Frame Based knowledge Representing

- organize relations around concepts
- lexical semantics via general semantic
- equivalent POP C.



frequency, speed, time

- John broke the window      breaker, frequency
- The window broke      broken thing.
- John always breaking things      breaking frequency.

We ate dinner.      eater

We already ate.      eaten thing,

The pies were eaten up quickly.      eating speed.

→ Open up!      opened.

→ Some left the door open.      opened thing.

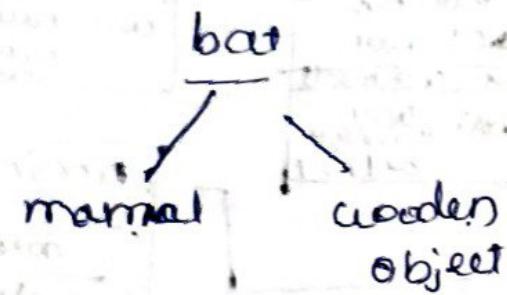
→ John opens the window at night.      opened time, at night

## word sense:

word sense & sense meaning

### # Ambiguity

She killed him with a base ball



### Different approaches:

→ knowledge Based approach:

Resources required

- chart corp. 1) Raw corpora  
vocab. 2) machine readable dictionary  
etc. Thesaurus.

how's done?

The bank can guarantee deposits will eventually cover failure.

Let's use example of the word "bark"

one meaning is → outer covering of the tree

2nd meaning → sound made by a dog.

so here same word different meaning.

1 sentence: "cinnamon comes from the bark  
of the cinnamon tree".

2 sentence: "the dog barked at the stranger".

This sentence is passed to an algorithm.

For sentiment analysis, "bark" and "barked"

might mean the same meaning.

So, we can understand that, same

words can mean differently based on

usage of the word in a particular sen-

the usage of words depends a lot

on tense. The usage of words depends a lot

about their meaning. But the problem

lies that, in NLP, while dealing with text

data, we need some way to interpret the

different words with different meaning.

### WSD Applications:

WSD has many applications in various

text processing and NLP fields.

→ WSD can be used alongside lexicography.

Much of the modern lexicography is corpus based. WSD, used in lexicography can provide

## Significant lexical Predicators.

→ WSD can also be used in Text mining and information extraction tasks. As the major purpose of WSD is accurately understand the meaning of a word.

→ As a security point of view, a system should be able to understand the difference b/w a coal "mine" and a land "mine".

→ WSD used for Information Retrieval purposes. It works through <sup>text</sup> data primarily based textual info. knowing the relevance of using a word in any sentence.

## Challenges in word sense Disambiguation

WSD faces a lot of challenges and problems.

→ The most common problem is the difference b/w various dictionaries of text corpora.

different applications need different alg  
and that is often a challenge for WSD.  
A problem also arises is that words  
cannot be divided into discrete meanings.  
words often have related meanings and  
this causes a lot of problems.

This causes rule based  
systems : rule based  
Dictionary - and knowledge - based methods

These methods rely on text data like  
dictionaries, thesaurus, etc. It is based  
on the fact that words that are related  
to each other can be found in the def.  
of each other can be found in the def.  
nitions. The popularly used test method.

Supervised methods :  
sense-annotated corpora are used.  
train machine learning models. But, a  
problem that may arise is that such  
corpora are very tough and time consuming  
to create.

## Semi-supervised methods

Due to the lack of such corpora, word sense disambiguation algorithms use semi-supervised methods. The process starts with a small amount of data, which is often manually created.

## unsupervised methods

Unsupervised methods pose the greatest challenge to researchers and NLP professionals. A key assumption of these

models is that similar meanings and senses occur in a similar context.

They are not dependent on manual effort.

Hence can overcome the knowledge acquisition

deadlock;

## unsupervised task algorithms

The task alg is a classical word sense disambiguation alg introduced.

The task alg is based on the idea that words in a given region of the text

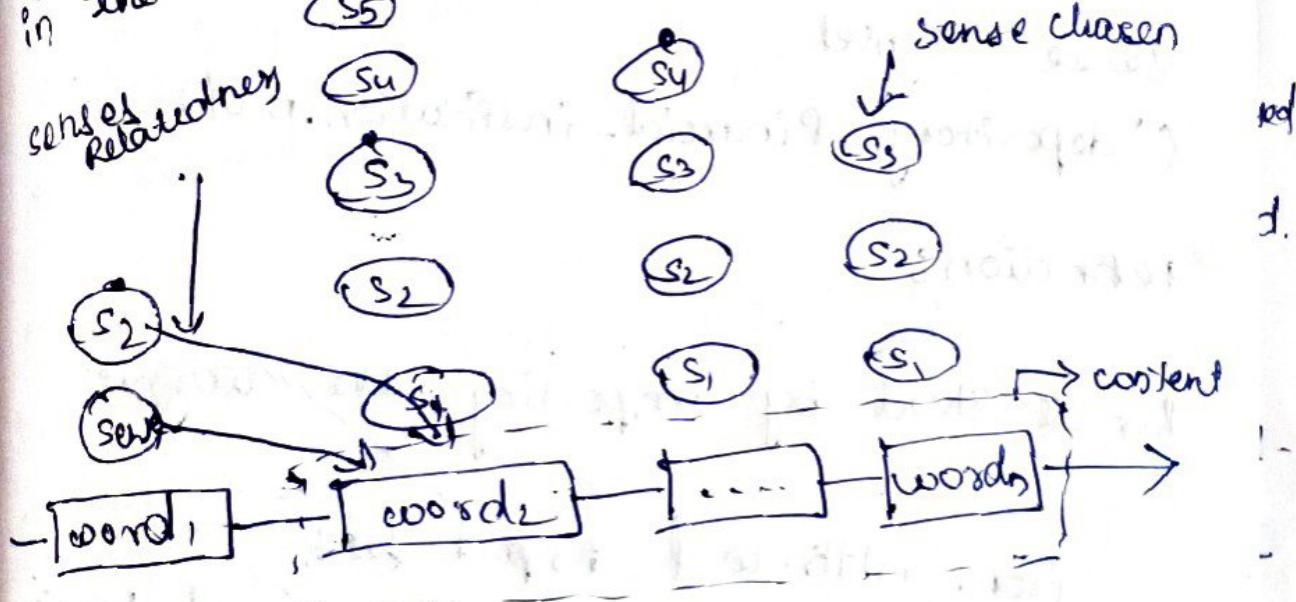
will have a similar meaning.

For ex: in the sentences below, the word "bank" has different meanings based on the content of a sentence.

Text 1 = I went to the bank to deposit my money.

Text 2 = The river bank was full of dead fishes.

It is basis on the hypothesis that words used together in text are related to each other and that the relation can be observed in the definitions of words and their sense classes.



Text

Use first install the library "pywsd".  
python implementation of word sense disambiguation (wsd)

pip install pywsd

```
from pywsd import simple_lesk
```

Sentences = { 'I went to the bank to deposit my money',

'The river bank was full of dead fish'

```
print('content-1:', sentences[0])
```

```
answer = simple_lesk(sentences[0], 'bank')
```

```
print('sense:', answer)
```

```
print('definition:', answer.definition())
```

out: content-1:  
I went to the bank to deposit my money

sense - synset

('depository - financial-institution.n.01')

definition:

Let's start by importing the libraries.

```
from nltk.wsd import lesk
```

```
from nltk.tokenize import word_tokenize
```

out  
a1 = lesk(word\_tokenize('this device is used  
to jam the signal'), 'jam')

print (a1, a1 · definition(c))

a2 = leek (word-tokenize ('I am stuck in.  
- a traffic jam'), 'jam')

print (a2, a2 · definition(c))

o/p: synset ('jamming.v.01') deliberate radiation or  
reflection of electromagnetic energy for the  
purpose of disrupting energy use of electronic  
devices or systems.

synset ('jam.v.05') get stuck and <sup>immobilized</sup>  
first; it is meant to stop the correct signal.  
and second implies a traffic jam.

ex2:  
b1 = leek (word-tokenize ('apply spices to the  
chicken to season it'), 'season')

print (b1, b1 · definition(c))

o/p: ('season.v.01') add flavour to

'season' is used for cooking sense of view.

gen: b2 = lex ( word\_tokenise ('India receives  
lot of train in the rainy season'), 'season')  
print ( b2, b2.definition())

Q/p: synset ('season.n.01') a period of the year  
marked by special events or activities in  
some field.

"geographic season."

i) c1 = lex ( word\_tokenise ('The current'), 'current')  
print ( c1, c1.definition())

Q/p: synset ('stream.n.02') suggestive of running.

- c1 = lex ( word\_tokenise ('the current time  
is 2 AM'), 'current' )

print ( c1, c1.definition())

Q/p :- synset ('current', a.o.1) occurring now beh-  
ging to the present time

```
from nltk.corpus import wordnet  
syn = wordnet.synsets('hello')[0]  
print("synset name: ", syn.name())  
print("in synset meaning: ", syn.definition());  
print("in synset example: ", syn.examples())  
  
if synset.name == 'hello.n.01':  
    print("synset meaning: an expression of greeting")  
    print("synset ex: ['every morning they exchanged n-  
    polite hellos']")
```

wordnet.synsets(word) can be used to get a list  
of synsets. This list can be empty (if no such  
word is found) or can have few elements.

### hypernyms and hyponyms

Hypernyms: more abstract terms  
Hyponyms: more specific terms  
Both come to picture as sets are orga-  
nized in a structure similar to that of an  
inheritance tree.

Supervised Pure  
Algorithm task 6

procedure: Simplified - best(word, sentence) ret  
best sense of word

best sense of word → most frequent sense of word

1) best-sense → most frequent sense of word

2) max-overlap ← 0

3) content ← set of word in sentence.

4) for all sense ∈ senses of word do

5) signature ← set of words in gloss and

examples of sense.

computeOVERLAP(sign, content)

6) overlap ← computeOVERLAP(sign, content)

7) if overlap ≥ max-overlap then

max-overlap ← overlap

best-sense ← sense

end if

end for

return best-sense.

Supervised: Ironically, the simplest form of

word sense disambiguating systems. The supervised approach which tends to transfer all

the complexity to the machine learning

classifiers probably the most common and highest performing classifiers are support vector machines (SVM) and maximum entropy (MaxEnt) classifiers. Many good-quality, freely available distributions of each are available and can be used to train word sense disambiguation models. Typically, because each lemma has a separate sense inventory, it is almost always the case that a separate model.

features: we discuss a more commonly found subset of features that have been useful in supervised learning of word sense.

lexical context: This feature comprises the words and lemmas of words occurring in the online paragraph or a smaller window.

part of speech: This feature comprises the POS information for words in the window surrounding the word that is being labeled.

Bag of words context: This featured comprising the POS information for words in the window. A threshold is typically tuned to include most informative words in the larger context.

local collocations: local collections are an ordered sequence of phrases near the target word, provide semantic context for disambiguation - ally every small window of about three tokens on each side of the target word.

Span ex: He bought a box of nails from the hardware store.

Here  $c_{1,1} \rightarrow$  bought would be word form.  
 $c_{1,2} \rightarrow$  from ~~the~~ store would be the string.  
 $c_{1,3} \rightarrow$  ~~the~~ store would be the string.

Syntactic relations: If the phrase of the sentence containing the target word available then we can use syntactic features.

Topic features: the broad topic, or domain, of the article that the word belongs to.

is also a good indicator of what sense of the word.  
voice of the sentences this ternary feature.  
indicates whether the sentence in which  
the word occurs is active, semiactive, or  
active sentence.

presence of subject object: This binary feature  
indicates whether the target word has a subject  
or object. Given a large amount of training  
data, we could also use the actual name  
sentential complement: This binary feature  
indicates whether the word has a sentential  
complement.

prepositional phrase adjunct: This feature in-  
dicates whether the target word has a senten-  
tial complement.

prepositional phrase adjunct: This feature  
indicates whether the target word has a pre-  
positional phrase.

## Algorithm: Rules for selecting Syntactic role

as features

- 1) if  $w$  is noun then
- 2)   Select parent head word ( $h$ )
- 3)   Select part of speech of  $h$
- 4)   select voice of  $h$ .
- 5)   Select position of  $h$  (left, right)
- 6) else if  $w$  is a verb then
- 7)   select nearest word  $l$  to the left of  $w$  such that  $w$  is the parent head word of  $l$ .
- 8)   select nearest word  $r$  to the right of  $w$  such that  $w$  is the parent head word of  $r$ .
- 9)   Select part of speech of  $l$ .
- 10)   Select part of speech of  $r$ .
- 11)   Select part of speech of  $w$ .
- 12)   Select voice of  $w$ .
- 13) else if  $w$  is a adjective then
- 14)   select parent head word, ( $h$ )
- 15)   Select part of speech of  $h$ .
- 16) end if.

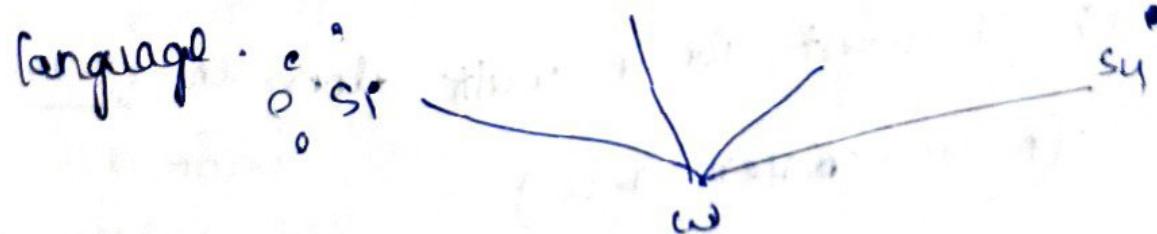
Named entity: This feature is the named entity of the proper nouns and certain types of common nouns.

wordnet: wordnet synsets of the hyponyms of head nouns of the noun phrase arguments of verbs and preposition.

path: this feature is the path from the target verb to verb's arg.

subcategorization: The subcategorization frame is essentially the string formed by joining the verb phrase type with that of its children.

unsupervised: progress in word sense disambiguation is stymied by the dearth of labeled training data to train a classifier for every sense of each word in a given language.



- training phase:
- features are extracted from the corpus.
  - an ordered decision list of the form
  - an ordered decision list of the form
  - a feature value, sense, score? is created.
  - The score of a feature  $f$  is the log-likelihood ratio of the sense given the feature  $f$  as,

$$\text{score } (s_i) = \max_f \log \left[ \frac{P(s_i | f)}{\sum_{f' \neq f} P(s_i | f')} \right]$$

### Decision lists

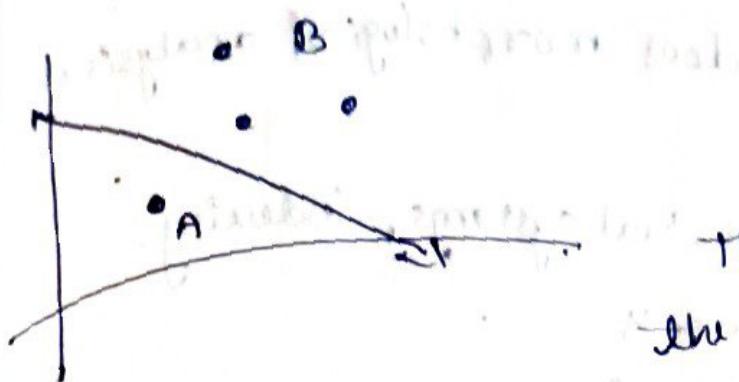
The decision list for the word "bank".

| <u>feature</u>      | <u>prediction</u> | <u>score</u> |
|---------------------|-------------------|--------------|
| of with the bank    | bank/fin          | 4.33         |
| standing in bank    | bank/fin          | 3.35         |
| bank of blood       | bank/supply       | 2.48         |
| walks in banks      | bank/fin          | 2.33         |
| The left river bank | bank/River        | 1.12         |

Test: I went for a walk along the river bank  
 (Noun - Verb - 1.12)

Support Vector Machines

e.g.: If a word has 4 features



| Sign | A              | B                               |
|------|----------------|---------------------------------|
| 1    | s <sub>1</sub> | s <sub>2</sub> , s <sub>3</sub> |
| 2    | s <sub>2</sub> | s <sub>3</sub> , s <sub>4</sub> |
| 3    | s <sub>3</sub> | s <sub>4</sub>                  |
| 4    | s <sub>4</sub> | s <sub>2</sub> , s <sub>3</sub> |

The distance gives the confidence score for each sign.

Program of WSD

Supervised:

Nominalization: It converts plural form of words to singular and present tense to past tense.

was → be

mice → mouse

meeting → meet

import Spacy.cli ('en\_core\_web\_sm')

spacy.cli.download('en\_core\_web\_sm')

import Spacy

nlp = spacy.load('en\_core\_web\_sm')

Nominalization is the process of grouping together the different inflected forms of a word so they can be analyzed as a single

item. Lemmatization is similar to stemming.  
it brings context to the words. so it links  
words with similar meaning the one way  
Lemmatization does morphological analysis of  
words.

Applications: retrieval systems, indexing.

ex: rocks : rock  
corpora : corpus  
better : good

Import nltk

nltk.download('wordnet')  
from nltk.stem import PorterStemmer  
from nltk.corpus import stopwords  
from nltk.corpus import stopwords

paragraphs "I have three visions for India

import nltk

from nltk.stem import PorterStemmer  
from nltk.stem import PorterStemmer

porter = PorterStemmer()

porter = Lemmatizer()

print ("Porter Stemmer")

print (porter.stem ("cats"))

print (porter.stem ("trouble"))

("troubling")

homographs:  
print (porter. stem ("troubled"));  
print ("hemmatizer");  
print ("hemmatizes. stem ("cats"))  
("trouble")  
("troubling")  
("troubled")

all porter stems

cat  
trouble  
troubl  
troubl

hemmatizer stems

cat  
troubl  
troubl  
troubl

porter stemmer uses suffix stripping to produce stems. It is giving the root (stem) of the word "cats" removing the "

stemming: It is process of reducing a word to its word stem that affixes to suffixes and prefixes or to the roots of words known as a lemma.

porter stemmer) & it includes the Porter Stemming technique. We can implement the Porter Stemming instance of Porter. One we convert an Porter alg. to Stem stemmer) and we the the list of words.

→ Snowball Stemmer this method utilizes in instance of more precise and is referred to as "English Stemmer". It is somewhat faster and more logical than original Porter Stemmer.

ex: ntlk-stem import snowball Stemmer  
Snowball = snowball stemmer(language='eng')  
words = ['generous', 'generate', 'generously',

'generation']

for word in words:

print(word, " -> ", Snowball.stem(word))

Op: generous -> generous

generate -> general

generously -> generous

generation -> general.

Software: several software programs are made available by the research community for word sense disambiguation, ranging from similarity-aware modules to full disambiguation systems. It is not possible to list all of them here so we list a selected few.

Step 1: Preprocessing  
for each sense of a word  $w_s$ , determine the synsets of wordnet in which  $w_s$  appears for each such synset, determine monosemous words included in that synset.

Step 2: search  
form search phrases using the following procedures in order of preference:  
1) → if they exist, extract monosemous synonyms.  
2) select each of the unambiguous passed constituents search phrase.

- 3) After parsing the gloss, replace all the words with a new operator and build a query.
- 4) we only the food phrase combined words in the synset using AND operator.

### Steps? Post processing

keep only those sentence in which the word under consideration belongs to the same pos as the selected sense, and delete the others.

## Predicate Argument Structure:

shallow semantic parsing, or what people know today as semantic role labeling is the process of identifying the various arguments or predicates in a sentence. The linguistic community has been debating for a few decades over the constituents the set of arguments and what the granularity of such arguments label should be, for various predicates.

Resources:

## Predicate - Argument Structure

Unit - IV

shallow semantic Parsing, or what is popularly known today as semantic role labeling. is the process of identifying the various arguments of predicates in a sentence.

### Resources:

The late 1990s saw the emergence of two important corpora that are semantically tagged. One is FrameNet and the other is Propbank. These resources have begun a transition from a long tradition of predominately rule-based approaches toward more data-oriented approaches. These approaches focus on transforming linguistic insights into features, rather than into rules, and letting machine learning learn from such resources. FrameNet is based on the theory of frame semantics

## framenet:

frame net corpus is a lexical database on English that is both human- and machine-readable, based on annotating example - of how words are used, in actual texts. framenet based on a theory of meaning called frame semantics, deriving from the work of charles. fillmore and colleague.

goals : It is convert semantic knowledge into something that is machine readable for use in NLP systems etc.

→ A semantic frame is a structure used to define the semantic meaning of a word.

→ For ex: there are abstract frames (such as the "replacement" frame); as well-as specific frame (such as the "applyheat" frame).

## part Replacement Frame

- part replaced [old The curtains] [new with wooden blinds].
- the framenet lexicon contains one lexical unit for the association of the verb place with the replacement frame.

## apply - heat frame

Boil [Food The rice] [cooking for 3 minute]  
(medium in water) Then drain.

Here Boil is unassociated with the Apply-heat frame to form another lexical unit.

Some additional possible associations with Apply-heat include : char, fry, grill, microwave, etc.

## Inherited frames

frames can also be inherited from other frames, a good example is the transport frame.

frame (transportation)

frame-elements (rooter (0), road, etc.)

scene (mover(s)) move along with the road

frame (driving)

inherent (trans)

frame-elements (mover (car), vehicle (car))

rider (0) (=mover (0)), cargo (=mover (0))

scene (Driver starts vehicle, driver controls, vehicle drives, stops vehicle)

frame (riding)

inherent (trans)

frame-elements (rider (=mover (0)), vehicle (car))

scene (Rider starts vehicle,

Vehicle carries rider along path,

Rider leaves vehicle).

Here "driving" and "riding" inheriting from frame.

transportation frame.

propbanks From the start, project was -

developed with the idea of generating training data for machine learning-based semantic

role labeling system in mind. It requires

that all arguments to a verb be syntactic

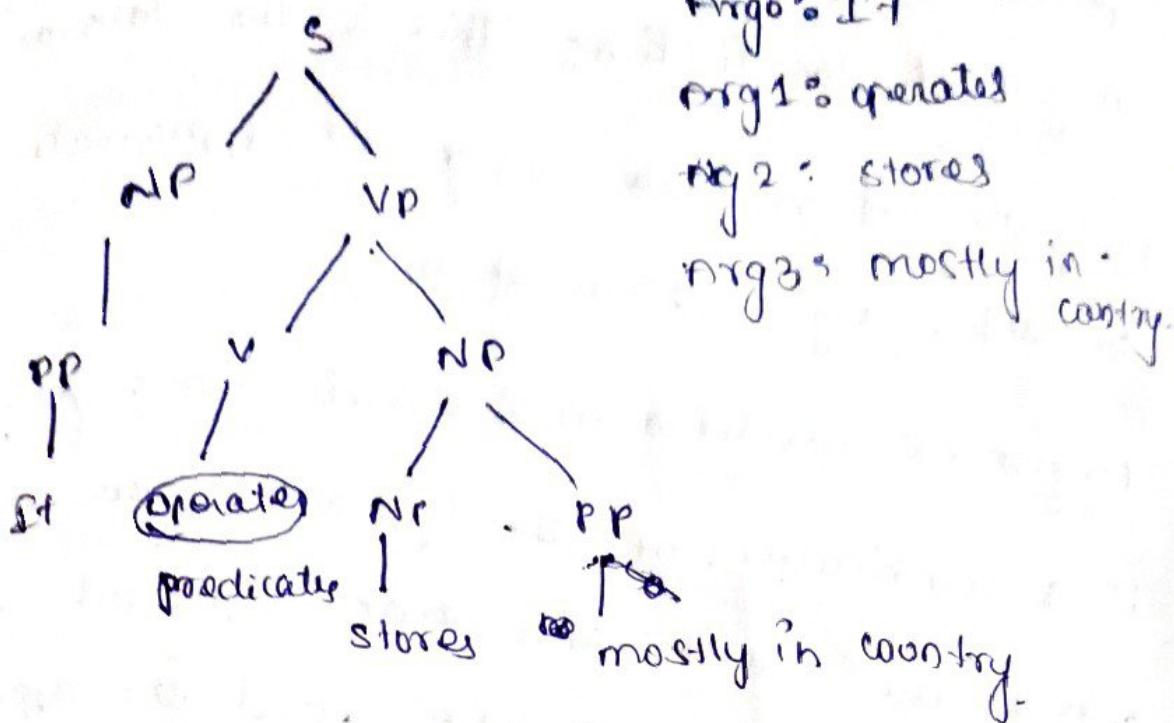
constituents do a word be syntactic constituents  
and different senses of a word are distinguished if the differences bear the  
ments.

- propbank: each verb has its own roles
- propbank more used, because it's layered  
over one treebank (and so has greater coverage  
plus parses)

Ex: ① sense: shift from one thing to another  
Arg0: cause of shift  
Arg1: thing being changed  
Arg2: old thing  
Arg3: new thing

② sense: move downward  
A1: thing falling  
A2: extent distance fallen  
A3: start point  
A4: end point

ex: It operates stores mostly in country



[Arg0 It] [medicat operates] [Arg1 stores]

[Arg2 mostly in country].

systems: WSD, little research has gone into learning predicate-argument structures from unannotated corpora, perhaps because it is closer to the actual applications and has been more or less absorbed in the area of information extractions.

Argument identification: This is the task of identifying all and only the pure constituents

That represents valid semantic arguments of predicate.

Argument classification: This is the task of identifying all arguments, assign the appropriate valid semantic arg label to them.

argument identification & classification: This task is a combination of the previous two tasks where the constituents that represent arg of a predicate are identified and the approp. arg label is assigned to them.

Syntactic Representation:

As we have seen, PropBank was created at a layer of annotation on top of Penn TreeBank style phrase structure trees and in some of the early work in recov.-ing propBank annotations, Gildea and Jurafsky [13] added argument labels to parses obtained from a parser trained on penn treeBank.

- Algorithm: The semantic role labeling (SRL)
- procedure : SRL (sentence) returns best semantic role labeling
- Input: syntactic parse of the sentence.
- 1) generate a full syntactic parse of the sentence
  - 2) identify all the predicates
  - 3) for all predicate  $\in$  sentence do .
  - 4) extract a set of features for each node in the tree relative to the predicate
  - 5) classify each feature vector using the model created in training
  - 6) select the class of highest scoring classifier
  - 7) Return best Semantic role labeling.
  - 8) end for .

phrase structure Grammar (PSG)

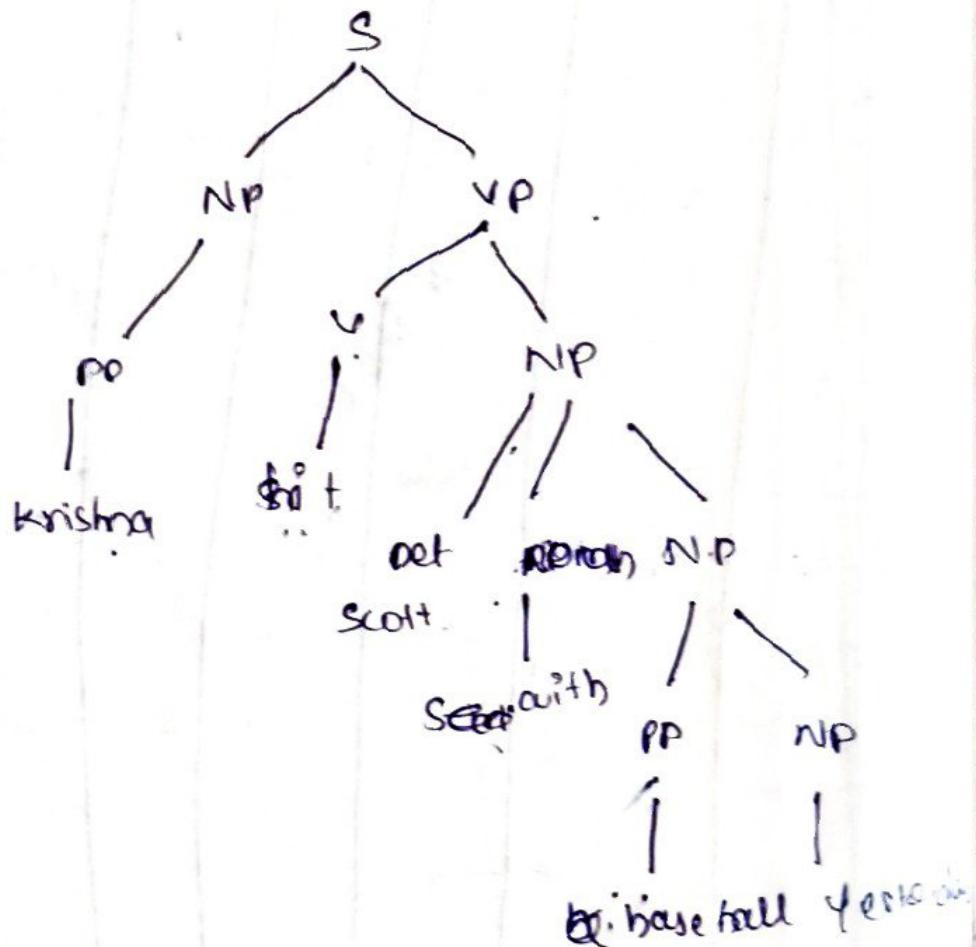
→ who hit Scott with a base ball  
→ with

→ who hit it?  
→ whom did Krishna hit with a baseball?  
hit Scott with

→ arham did hit scott with  
→ arved did krishna hit scott with

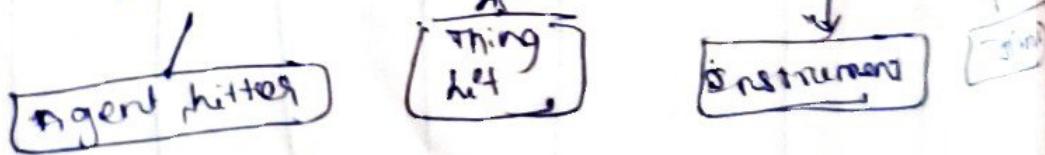
→ cuten did Krishna hit Scott with a baseball.

## Analysis



## Variations

Krishna hit Scott with a baseball.



phrase structure grammar (PSG): Frame net.  
marks word spans in sentence to represent  
arg, whereas PropBank tags nodes in a tree-  
bank tree with arguments. Because several  
high-quality statiscal parsers existed that  
could produce phrase structure tree.

path: this feature is the syntactic path thro-  
ugh the parse tree from the base consti-  
uent to the predicate being classified.

Predicate: The identity of the predicate lemma  
is used as a feature.

phrase type: This feature is a syntactic -  
category (NP, PP, S, etc).

position: This feature is a binary feature  
identifying whether the phrase is before or  
after the predicate.

voice: this feature indicates whether the  
predicate is realized as an active or passive  
construction.

head word: this feature is the syntactic head of the phrase. It is calculated using a head word table.

sub categorization: This feature is the phrase structure rule expanding the predicate's parent node.

verb clustering: the predicate is one of the most salient features in predicting the arguments class. given various syntactic semantic constraints, that a predicate can appear in, any amount of hand-tagged training data.

content word: Because the head word feature of some constituents, such as pp and

(1) restricts the development to system that operate in limited domains.

(ii) they are hard to maintain & scale up as the problem becomes more complex & more domain independent

(iv) they tend to be brittle

Software:

Not a lot of software programs are available for older, more rule based system, but the following are available for download

WASP [<http://www.cs.utexas.edu/~ml/wasp/>]

CHILL [<http://www.cs.utexas.edu/~ml/chill/>]

## UNIT- IV

### PREDICATE - ARGUMENT STRUCTURE

Shallow semantic parsing , or what is popularly known today as semantic role labelling is the process of identifying the various arguments of predicates in a sentence.

#### Resources :

The late 1990s saw the emergence of two important corpora that are semantically tagged . One is permanent & the other is propbank. These resources have begun a transition from a long tradition of predominately rule-based approaches toward more data-oriented approaches . These approaches focus on transforming linguistic insights into features , rather than into rules , and letting a machine learning framework use this in such resources . Framenet is based on the theory of frame semantics .

#### Frame Net :

Frame net corpus is a lexical database on English that is both human and machine-readable based on amulating example of how words are used in actual texts . Frame net based on a theory of meaning called Frame semantics , deriving from the work of Charles \*

full more and colleagues.

### FrameNet Example

| Frame: Awareness |            |
|------------------|------------|
| cognizer         | Manner     |
| content          | Expressor  |
| Evidence         | Role       |
| Topic            | Phenomenon |
| Degree           | Practice   |

### Goals:

It is convert semantic knowledge into something that is machine readable for use in NLP systems etc.

→ A semantic frame is a structure used to define the semantic meaning of a word.

→ For ex: There are abstract frames such as the "Replacement frame", as well as specific frames (such as the "Apply-heat" frame)

### Replacement frame

- ⇒ pat replaced [old The curtains] (new with wooden blinds).
- ⇒ The FrameNet lexicon contains one lexicon unit for the association of the verb place with the replacement frame.

Apply - heat frame.

Boil [ food the rice] [ duration for 3 minutes]  
[ medium in water] then drain.

Note Boil is associated with the apply - heat frame to form another lexical unit with some additional possible associations with apply - heat include : char , fry , grill , wave , etc

Inherited frames :

Frames is also be inherited from other frames, a good example is the transportation frame.

Ex: Frame (transportation)

frame . elements (mover(s), means, path)

Scene (mover(s) move along path by means)

frame (driving)

inherit (trans)

frame . elements (driver (= mover) . vehicle (= means), rider(s) (= mover(s)), cargo (= mover(s))

Scene (Driver starts vehicle , drive controls

Vehicle . Driver . stops . vehicle)

frame (riding)

Inherit (trans)

Frame' elements (Rider (= mover(s)), vehicle (= means))

Scene (Rider enters vehicle,  
vehicle covers Rider along path,  
Rider leaves vehicle)

PROPER BANK (PROP BANK):

From the start prop bank was developed with the idea provide training data to machine learning based on the semantic rules. It requires all arguments to a verb & different senses of a word distinguished by the argument.

→ In sentence each verb follows CFG rules it layered over the tree bank.

Ex: Sense: shift from one thing to another

①

Argo : consider shift  
Arg1 : change (replace)

Arg 2 : old thing

Arg 3 : new thing

②  
Sense : move downward

A1 : thing falling

A2 : exert distance fallen

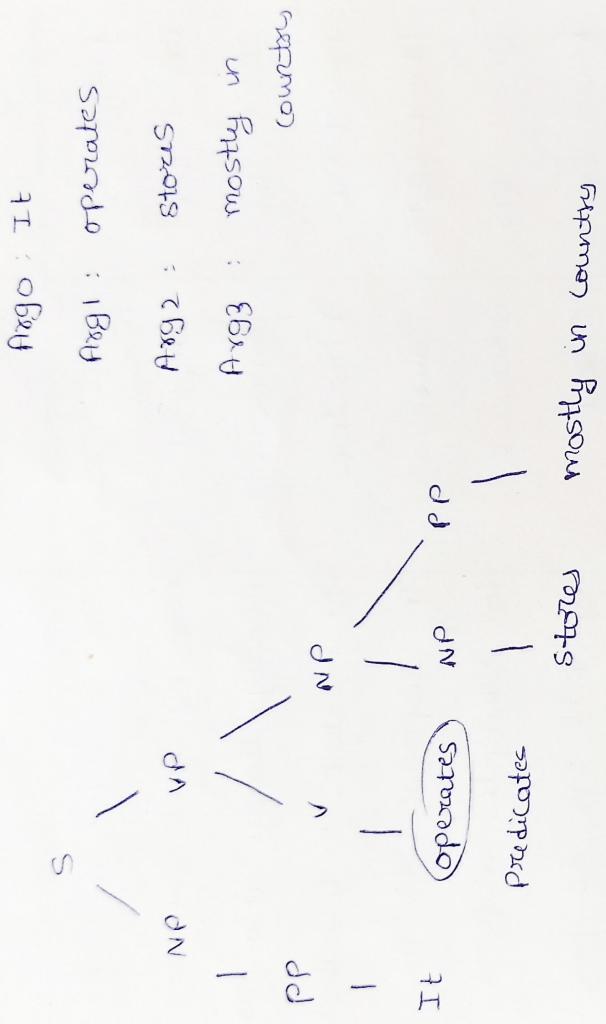
A3 : Start point

A4 : end point

List of adjunctive arguments in PropBank

| Tag      | Description | Example                          |
|----------|-------------|----------------------------------|
| ARGM-LOC | Locative    | the museum in Westborough, Mass. |
| ARGM-TMP | Temporal    | now, by next summer              |
| ARGM-MNR | Manner      | heavily, dearly, at a rapid rate |
| ARGM-DIR | Direction   | to market, to Bangkok            |
| ARGM-CAU | Cause       | In response to the ruling        |
| ARGM-DIS | Discourse   | for example, in part, similarly  |
| ARGM-EXT | Extent      | at \$38.375,50 points            |
| ARGM-PRP | Purpose     | to pay for the plant             |
| ARGM-NEG | Negation    | not, n't                         |
| ARGM-MOD | Modal       | can, might, should, will         |
| ARGM-REC | Reciprocals | each other                       |
| ARGM-PRO | Secondary   | to become a teacher              |
|          | Predication |                                  |
| ARGM     | Bare ARG    | with a police escort             |
| ARGM-ADV | Adverbials  | (none of the above)              |

Ex: It operates stores mostly in country



[Argo It] [medicate operates] [Arg2 stores]

[Arg-hoc mostly in country]

Systems:

WSD, little research has gone into learning medicate - argument structures from annotated corpora, perhaps because It is closer to the actual applications and has been more or less absorbed in the area of information extractions.

Argument Identification.

This is the task of identifying all and only the pause constituents that represents valid semantic arguments of the predicate.

### Argument classification:

This is the task of identifying all predicate, assign the appropriate value semantic arg label to them.

### Argument Identification & Classification:

This task is a combination of the previous few tasks where the constituents that represent arg of a predicate are identified and the appropriate arg label is assigned to them.

### Syntactic Representation:

As we have seen, PropBank was created as a layer of annotation on top of Penn TreeBank style phrase structure trees and in some of the early work in recovering propBank annotations, Gildea and Jurafsky [113] added argument labels to pastes obtained from a parser trained on Penn Tree Bank.

### Algorithm : The Semantic role labeling (SRL)

procedure : SRL (sentence) returns best semantic role labeling.

Input : Syntactic parse of the sentence.

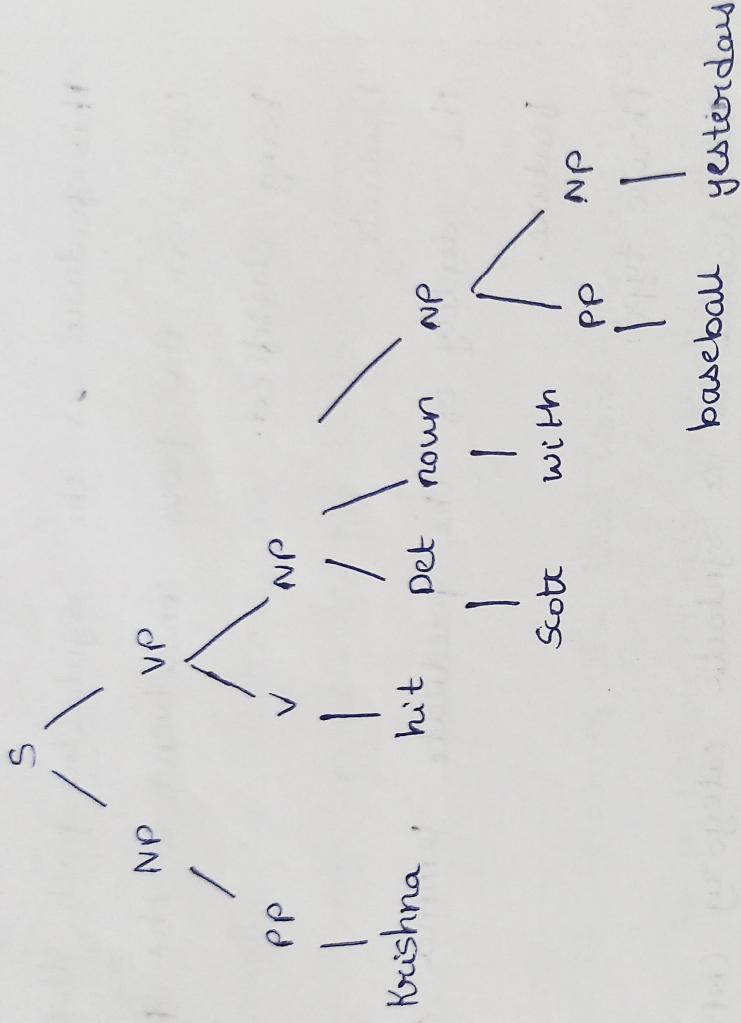
- 1) Generate a full syntactic parse of the sentence.

- 2) Identify all the predicates
- 3) for all predicate  $\in$  Sentence do
  - 4) Extract a set of features for each node in the tree relative to the predicate
  - 5) classify each feature vector using the model created in training
  - 6) select the class of highest scoring classifier.
- 7) Return best semantic node labelling
- 8) end for

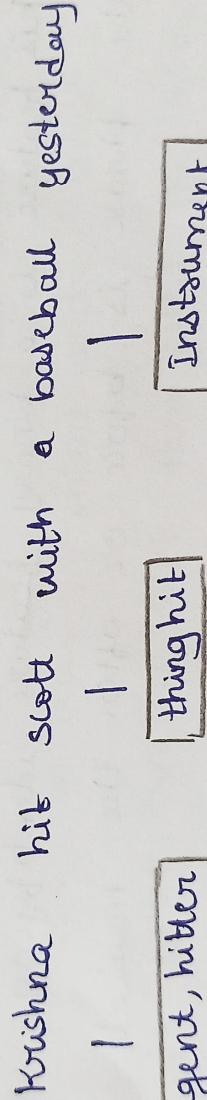
Ex:

- |         | who | which | what            | when      |
|---------|-----|-------|-----------------|-----------|
| Krishna | hit | scott | with a baseball | yesterday |
- Who hit scott with a base ball  
 → whom did Krishna hit with a baseball  
 → what did Krishna hit scott with  
 → When did Krishna hit scott with a baseball.

### Analysis:



Variations:



### PHRASE STRUCTURE GRAMMAR (PSG):

Frame net models and spans in sentence to represent arguments, propBank tags nodes in a tree. Bank tree with arguments because several high quality statistical parsers existed that could produce phrase structure tree.

tree.

Path :

This feature is the syntactic path through the parse tree from the parse constituent to the predicate being classified.

Predicate :

The identity of the predicate lemma is used as a feature.

Phrase type :

This feature is a syntactic category (NP, PP, S, etc.)

Position :

This feature is a binary feature identifying whether the phrase is before or after the mediate.

Voice :

This feature indicates whether the predicate is realised as an active or passive construction.

Head word :

This feature is the syntactic head of the phrase - it is calculated using a head word table

Sub categorization :

This feature is the phrase structure rule expanding the mediates parent node.

## Verb clustering:

The mediate is one of the most salient features in predicting the argument class given various syntactic constructions that a predicate can appear in, any amounts of hand-tagged training data.

word:

Adding part of speech of the head word & content word of a constituent as a feature to help generalize in the task of arg.

H1: If phrase type is PP then select the right most child.

Ex: phrase = "in Texas", content word = "Texas"

H2: if phrase type is SBAR then select the leftmost sentence ( $s^4$ ) clause.

Example: phrase = "that occurred yesterday"  
content word = "occurred"

H3: if phrase type is VP then  
if there is a VP child  
else  
select the head word

Ex: phrase = "had placed", content word = "placed"

H4: If phrase type is ADVP then select the right most child, not IN or TO

Ex : phrase = " more than", content word => " more"

H5 : if phrase type is ADVP then select the rightmost adjective, verb , noun or ADVP

Ex : phrase 4 = " 61 years old", content word - " 61 "

H6 : for all other phrase types the head word

Ex : phrase = " red house" content word = " red "

Named entities in constituents :

Reported a performance improvement on classifying the semantic role of the constituent by using the presence of a named entity in the constituent.

verb sense information:

The arguments that a predicate can take depend on the sense of the predicate . Each predicate on the sense in which it is used. the two sense of predicates talk in propbank corpus.

talk . 01

Description

Arg 0

talker

Arg 1

Subject

Arg 2

hearer

talk . 02

Description

Arg 0

talker

Arg 1

talked to

Arg 2

Secondary  
actr.

clause-based path variations:

position of the clause node (S, SBAR) seems to be an important feature in argument identification.

Accordingly, experimented with four clause-based path feature variations.

- Replacing all the nodes in a path other than clause (S, SBAR) seems to be an important feature in arg identification.

→ Replacing all the nodes in a path - other than clause nodes with an (+). for ex: the path NP ↑ S ↑ VP ↑

SBAR ↑ NP ↑ VP ↓ VBD becomes NP ↑ S ↑ VS ↑ V ↑ ↓ VBP

Path n-grams:

This feature category decomposes a path into a series of trigrams. Ex:

path NP ↑ S ↑ VP ↑ SBAR ↑ NP ↑ VP ↑ VBD becomes  
NP ↑ S ↑ VP , S ↑ VP ↑ SBAR , VP ↑ SBAR ↑ , NP,  
SBAR ↑ NP ↑ VP , and so on. shorter paths were padded

with

Predicate Context:

predicate sense variations. Two words before and two words after were added as features.

Punctuations:

for some adjunctive arguments punctuations plays an important

note. This set of features captures whether punctuation amends immediately.

Combinatory Categorical Grammar (CCG):

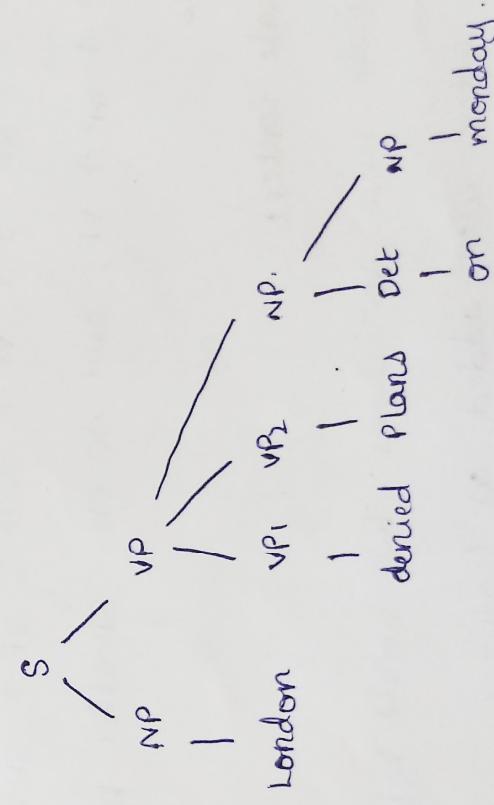
The path feature is very important from the path feature to argument identification task. It is one of the most sparse features and may be difficult to train or generate.

Phrase type:

This is category of the maximal projection between the two words, the mediate and the independent word.

London denied plans on monday

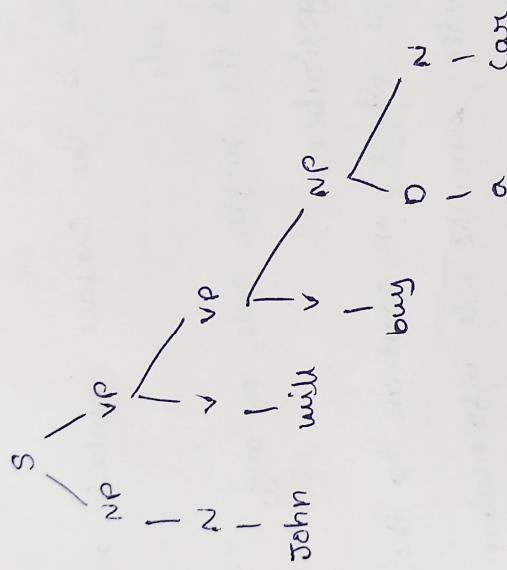
denied London ( $S \setminus NP_1$ )  $NP$   
denied plans ( $S \setminus NP_1 \setminus VP_1$ )  
on denied ( $S \setminus N \setminus VP_1$ )  
monday ( $S \setminus N \setminus NP_1$ )



Tree adjoining grammar (TAG):

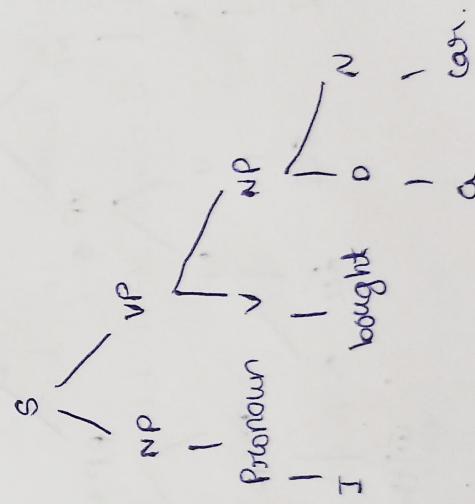
- Tree generating systems
  - what are the elementary objects of CG
  - what are the elementary objects of TAG.

"John will buy a car"



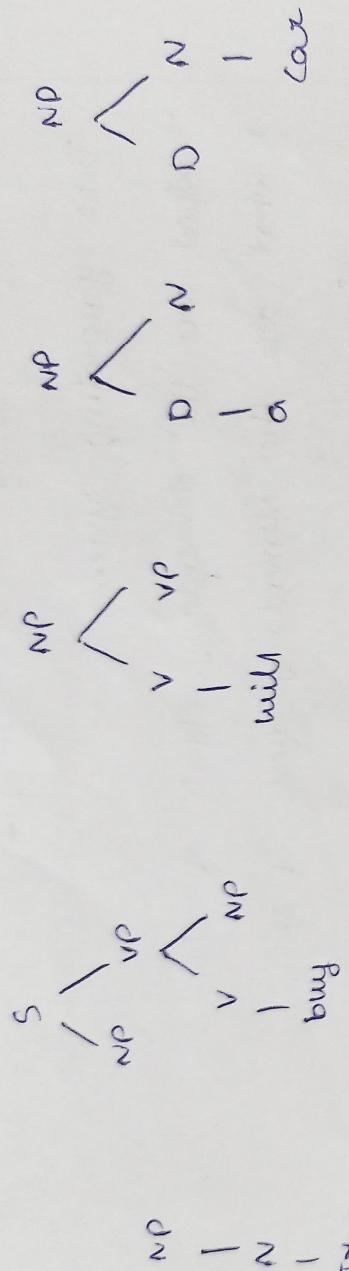
Strings: S, NP, N, VP, V, NP, O, D, N

I bought a car



S, NP, Pronoun, VP, V, NP, O, D, N

Trees:



super tag path:

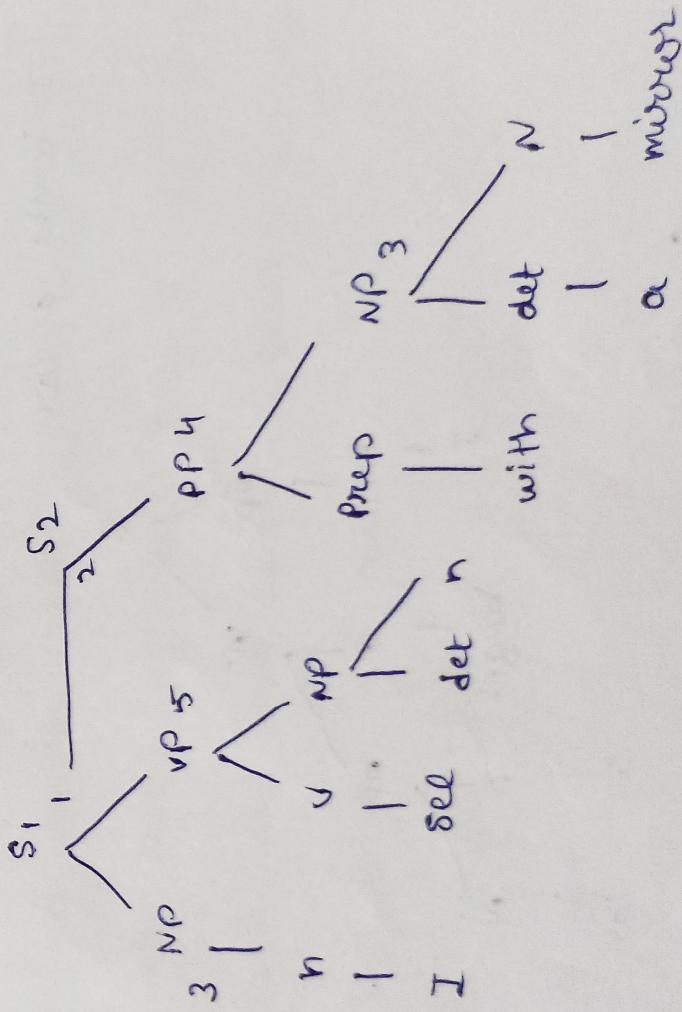
This feature checker is the surface syntactic rule of any surface syntactic rule.

This feature is the surface role of org.

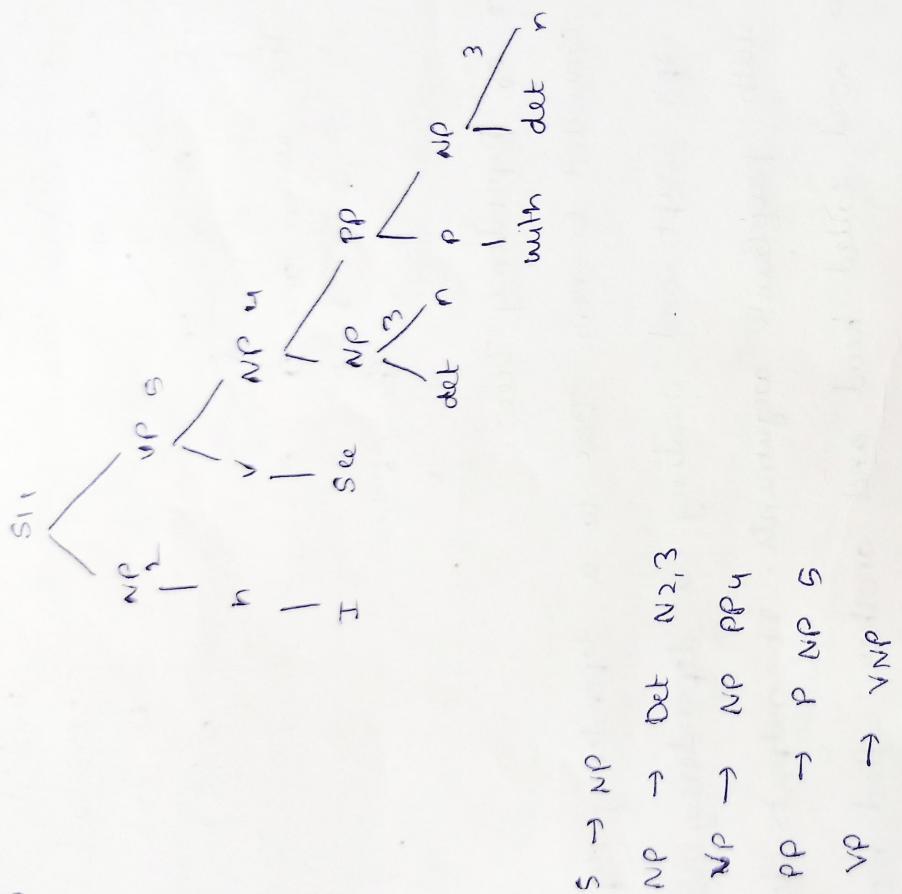
Semantic Subcategorization:

This is a frame where, in addition to the syntactic categories, the feature includes semantic role information.

I see a man with a mirror"



(ii)



### Dependency Trees:

Performance of a system depends on the exact span of arg.

Annotated according to the constituents in the Penn Tree bank. Labels are scored as correct only if they match.

Because propBank and most syntactic-parsers are developed on the Penn Tree Bank.

Ex:

$$S \rightarrow NP VP$$

$$VP \rightarrow V NP$$

$$V \rightarrow \text{loves}$$

$$D \rightarrow \text{the}$$

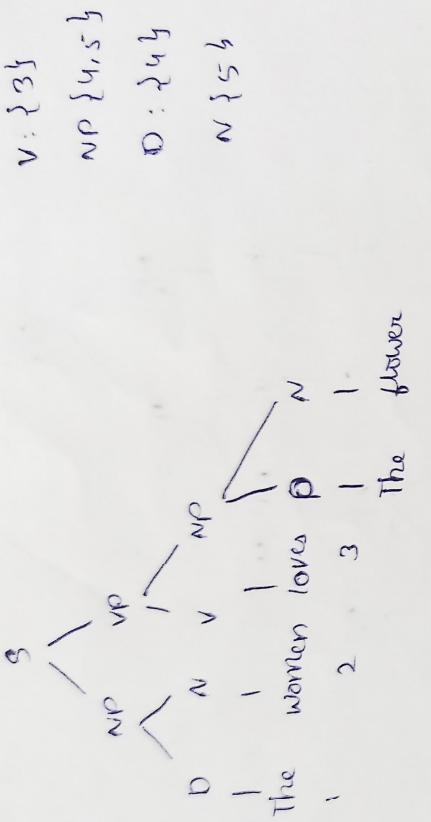
$$N \rightarrow \text{women} \mid \text{flower}$$

$$S: \{1, 2, 3, 4, 5\}$$

$$NP: \{1, 2\}$$

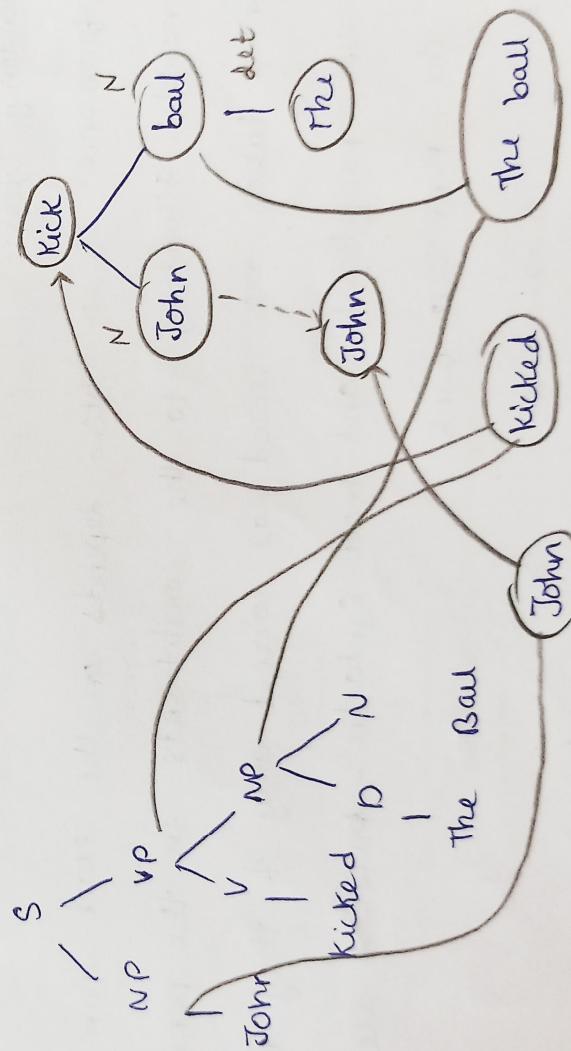
$$D: \{1\}$$

$$VP: \{3, 4, 5\}$$



Architecture of Dependency Tree:

It is converted to Penn Tree to a dependency representation using a spot script and creating a dependency structure labelled with propBank arguments. It outputs dependencies between a word called head and another called modifier.



Head word: the word representing the node in the dependency tree

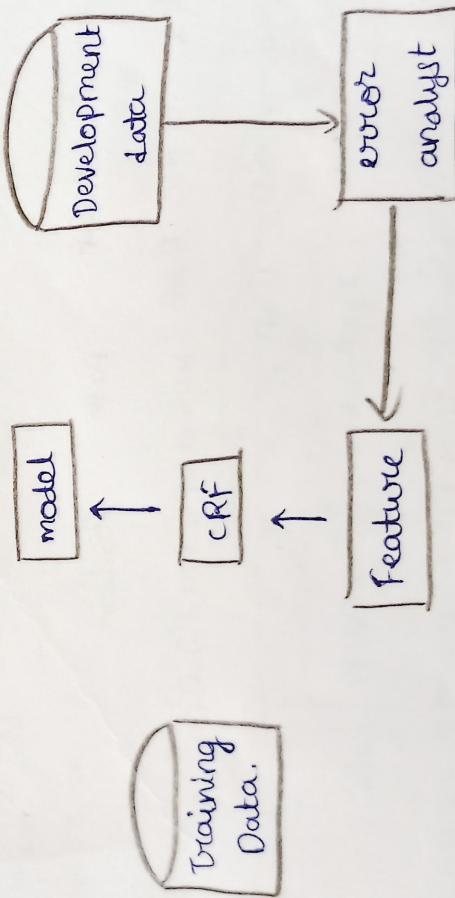
Head word nos: pos of head word

pos path: the path from the predicate to the head word through the dependency tree.

Dependency Path: each word that is connected to the head word has a dependency relationship to the word.

Voice position: voice of the predicate whether the node is before or after the predicate.

Base phrase chunks: A common question is how much does the full syntactic representation help the task of semantic role labelling. In other words, how important is it to create a full syntactic tree before classifying the arguments of a predicate? A chunk representation can be faster & more robust to phrase reordering as happens in speech-data



defines the chunks using IOB tags:

- It specifies where the chunk begins and ends along with its types.

- A part of speech tagger can be trained on these IOB tags to further train a chunker sub class.
- First using chunked sentences from the corpus, a tree obtained and is then transformed to a format usable by a pos tagger

- conll-tag-chunks() uses tree conll tags () to convert a sentence tree into a list of 3 tuples of the form (word, pos, iob)

Ex: B-NP, I-NP, O - NP

Beginning Inside Outside

\* Sales declined 3% to \$ 524.5 mill from \$ 539.4 mill.

↳ chunk into syntactic base phrases

[Sales NP] [declined VP] [to VP] [NP \$ 524.5 m] [from] [539.4 NP]

Phrase Headword posBP Path Position

NP Sales NNS B-NP NNS → NP → Pred → vBD

Pred declined VBZ B-VP -

NP / NN I-NP NN → NP → Pred → vBD

PP to TO B-PP to PP → NP → NP → Pred → vBD

NP mill CD T-NP CD → NP → PP → NP - Pred → vBD

from TN B-PP PN - PP → NP → PD → NP → Pred → vBD

PP mill CD I-NP CD → NP → DD → ND → PD → NP → Pred → vBD

NP mill CD I-NP CD → NP → DD → ND → PD → NP → Pred → vBD

## Overcoming the Independence Assumption

As mentioned earlier, various post processing stages have been proposed to overcome the limitations of treating Semantic role labelling as a series of independent argument classification.

Ex: But [nobody] [knows] [at what level] [the futures]  
P Arg1 Arg2

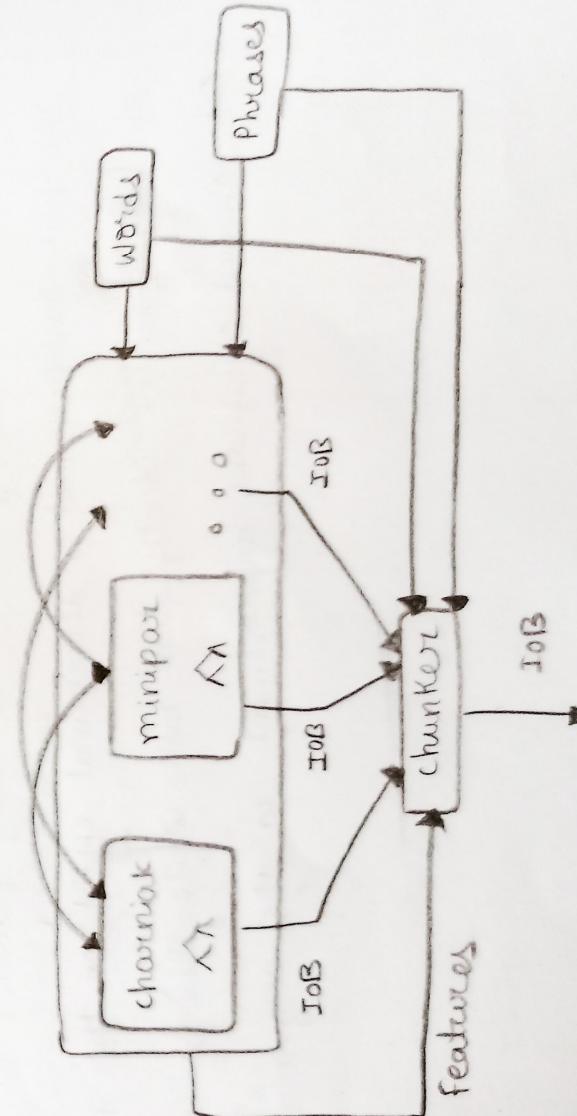
### Feature Performance:

Not all features are equally useful in each task. Some features add more noise than information in one context than in another.

### Feature Salience

In analyzing the performance of the system, it is useful to estimate the relative contribution of the various feature sets used.

### New architecture:



Semantic role labels

### Overcoming Parsing Errors!

After performing a detailed error analysis, Prodromou et al. found that the identification problem poses a significant bottleneck to improving overall system performance. The base line system's accuracy on the task of labelling nodes known to represent semantic arguments is 90.1. On the other hand, the system's performance on the identification task is quite a bit lower, achieving only 80.1 recall with 86.1 precision. The two sources of these identification errors are induced by the system to identify all and only those constituents that correspond to semantic roles, where these constituents are present in the syntactic analysis and failures by the syntactic analyzer to provide the constituents that align with correct arguments.

### Feature Selection:

The fact that adding the named entity features to the null/non-full classifier had a detrimental effect on the performance of the argument identification task, while the same feature set showed significant improvement to the argument classification task

(2)

(2)

### Argument classification

Table: Effect of each feature on the argument classification task when added to the base task & argument identification task when added to the base task. An asterisk indicates that the improvement is statistically significant.

| Features                     | Argument classification |      |      | Argument identification |      |
|------------------------------|-------------------------|------|------|-------------------------|------|
|                              | A                       | P    | R    | F1                      | F1   |
| Baseline [120]               | 87.9                    | 93.7 | 88.9 | 91.3                    |      |
| + Named entities             | 88.1                    | 93.3 | 88.9 | 91.0                    |      |
| + Head POS                   | 88.6                    | 94.4 | 90.1 | 92.2                    | ship |
| + Verb cluster               | 88.1                    | 94.1 | 89.0 | 91.5                    |      |
| + Partial Path               | 88.2                    | 93.3 | 88.9 | 91.1                    |      |
| + Verb Sense                 | 88.1                    | 93.7 | 89.5 | 91.5                    |      |
| + Noun head PP (only pos)    | 88.6                    | 94.4 | 90.0 | 92.2                    |      |
| + Noun head PP (only head)   | 89.8                    | 94.0 | 89.4 | 91.7                    |      |
| + Noun head PP (both)        | 89.9                    | 94.7 | 90.5 | 92.6                    |      |
| + First POS in constituent   | 89.0                    | 94.4 | 91.1 | 92.7                    |      |
| + Last POS in constituent    | 89.4                    | 93.8 | 89.4 | 91.6                    |      |
| + Ordinal const. pos. const. | 88.4                    | 94.4 | 90.6 | 92.5                    |      |
| + const. tree distance       | 88.3                    | 93.6 | 89.1 | 91.3                    |      |
| + Parent constituent         | 87.7                    | 93.7 | 89.2 | 91.4                    |      |
| + Parent head                | 88.0                    | 93.7 | 89.5 | 91.5                    |      |
| + Parent head POS            | 87.9                    | 94.2 | 90.2 | 92.2                    |      |
| + Right sibling constituent  | 85.8                    | 94.2 | 90.5 | 92.3                    |      |
| + Right sibling head         | 88.5                    | 94.3 | 90.3 | 91.9                    |      |
| + Right sibling head POS     | 87.9                    | 94.0 | 89.9 | 92.1                    |      |
| + Left sibling constituent   | 87.9                    | 94.4 | 89.9 | 92.0                    |      |

mention that a

|   |                       |      |      |      |      |
|---|-----------------------|------|------|------|------|
| * | right sibling head    | 88.1 | 91.1 | 86.1 | 89.9 |
| + | left sibling Head POS | 88.6 | 93.6 | 89.3 | 91.4 |
| - | Temporal cue words    | 86.9 | 93.9 | -    | -    |
| + | Dynamic class context | 88.8 | 93.5 | -    | -    |

Table : Performance of various feature combinations on the task of argument classification.

Features                          Accuracy

|                              |      |
|------------------------------|------|
| All features [120]           | 91.0 |
| All except path              | 90.8 |
| All except phrase type       | 90.8 |
| All except Hw & Hw-POS       | 90.7 |
| All except all phrases       | 83.6 |
| All except predicate         | 82.4 |
| All except Hw & fw & Lw info | 75.1 |
| only path & predicate        | 78.4 |
| only path & phrase type      | 47.2 |
| Only Head word               | 37.1 |
| Only Path                    | 28.0 |

BordenSearch: Another approach is to broaden the search by selecting constituents in a best paths of using a packed context representation which more efficiently represents variations over much larger n. Using a parse forest shows an absolute improvement of 1.2 points over single best paths and 0.5 points over most paths.

Example of nominalization:

she complained about the attack  
she made an official complain about the attack

John walked around the university  
John took a walk around the university

identify and label the arguments of nouns. The only works that come close are the rule-based system described by Hull & Gomez and the work of Lapata on interpreting the relationship b/w the head of a nominalized compound & its modifier noun.

Robustness across genre  
one possible shortcoming of this & other approaches is that all the parsers are trained on the same Penn Treebank, when evaluated on sources other than WST, seems to degrade in performance. whereas & Marquez show that the drop in performance on test data from the Brown corpus was 10%+ such points lower than the test data from WST. When trained & tested on WST prepositions the syntactic parser's performance is the main source of error, and the classification performance is quite good.

Meaning Representation:

We now turn to the third, deeper level of semantic interpre-  
lation whose objective is to take natural language input  
& transform it into an unambiguous representation that a

machine can act on. This is the form that would be more likely to be as incomprehensible to humans as it would be more likely comprehensible to machines. We can think of a parallel to programming languages that are much closer to the way humans manipulate information and the low-level machine code that the computer executes. Although compilers and interpreters impose various specific syntactic & semantic restrictions on a program written in a high level programming language, no such restrictions are imposed on the form that natural language can take.

#### Resources:

A number of projects have created representations & resources that have promoted experimentation in this area.

#### ATIS

The Air Travel Information System project is considered one of the first concerted efforts to build systems to transform natural language into a representation that could be used by an end application to make decisions. The task involved a machine to transform a user query in spontaneous speech, using a restricted vocabulary, about flight information. It then formed a representation that was compiled into a SQL query, which was used to extract answers from a flight database.

While the communicator program runs the follow-on to ATIS, while ATIS was more focused on user-initiated dialog, communicator involved a mixed-initiative dialog, whereby the human and machine had a dialog with each other with the computer helping presenting users with real-time travel information & helping them negotiate a preferred itinerary.

Rule Based:

Semantic parsing systems that performed very well for both the ATIS & communicator projects were rule-based systems in the sense that they used an interpreter whose semantic grammar was handcrafted to be robust to speech recognition errors. The underlying philosophy was that the traditional syntactic explanation of a sentence is much more complex than the underlying semantic information, so parsing the meaning units in the sentence into semantic primitives to be a better approach.

Supervised:

Although rule-based techniques are relatively easy to craft in the beginning & served a good purpose to formulate solutions to various tasks, they have several downsides:

- i) they need some effort upfront to create rules
- ii) the time & specificity required to write rules usually

## UNIT - V

language model :

A language model is a probability distribution over sequence of words. Given such a sequence of length ' $m$ ', a language model assigns a probability  $P$  to the whole sequence. Language models generate probabilities by training of text corpora in one or many languages.

Language model from the backbone of the NLP. They are a way of transforming qualitative information about text into quantitative information that machines can understand.

N-gram model :

N-gram are continuous sequences of words or symbols or tokens in a document. In technical terms, they can be defined as the neighbouring sequences of items. They come into play when we deal with text data in NLP.

What are n-grams:

N-grams are continuous sequence of words or symbols or tokens in a document. In technical terms, they can be defined as the neighbouring sequence of items.

n-grams classifier:

' $n$ ' is just a variable that can have the integer values including 1, 2, 3 and so on.

→ 'n' basically refers to multiple.

| n | Term    |
|---|---------|
| 1 | unigram |
| 2 | Bigram  |
| 3 | Trigram |
| n | n-gram  |

when  $n=1$ , it is said to be unigram

$n=2$ , it is said to be a bigram

$n=3$ , it is said to be a trigram

Ex: Language model:

$$p(x, y) = p(x/y)p(y)$$

$$p(x, y, z) = p(x)p(y/x)p(z/x,y)$$

$$p(\text{she}, \text{is}, \text{dead}) = p(\text{she})$$

$$p(\text{is}/\text{she})$$

$$p(\text{dead}/\text{she is})$$

Markov Assumption:

I wish was a unicorn

$$\rightarrow p(\text{unicorn}/\text{I wish was a}) \approx p(\text{unicorn}/a)$$

$$\text{or } p(\text{unicorn}/\text{was a})$$

Ex: N-gram such as unigram

("This", "article", "is", "on", "NLP") or

bigram: ("This article", 'article is', 'is on', 'on NLP')

Now, we will establish a relation on how to find the next word in the sentence.

$P(NLP / \text{This article is on})$

equation:  $P(w_5 / w_1 w_2 w_3 w_4)$  or  $P(w)$

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

calculation:  $P(A/B) = \frac{P(A, B)}{P(B)}$

$$P(A, B) = P(A/B)P(B)$$

Corpus:

<ss> I am a human </ss>

<ss> I am not a stone </ss>

<ss> I I live in mumber </ss>

-> check the probability of

-> <ss> I am not </ss> using

$\rightarrow (I \ I am \ not)$

$$= P(I / \text{ss}) P(I/I) P(am/I) P(not/am) P(ss/not)$$

$$\approx \frac{c(ss/I)}{c(ss)} \cdot \frac{c(I/I)}{c(I)} \cdot \frac{c(I/am)}{c(I)} \cdot \frac{c(am/not)}{c(am)}$$

$$\frac{c(not/ss)}{c(not)}$$

$$= \frac{3}{3} \cdot \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{1}{2} \cdot 0\%$$

$$= 0 \quad //$$

→ Consider the following training data

<ss> I am Jack </ss>

<ss> Jack I am </ss>

<ss> Jack I like </ss>

<ss> Jack I do like </ss>

<ss> do I like Jack </ss>

Assume that we use a bigram.

language model based on above data.

and is most probable next word predicted by the model

1. <ss> Jack 2. <ss> Jack I do ...

3. <ss> Jack I am Jack ...

4. <ss> do I like ...

$$p(I/\text{ss}) = \frac{c(\text{ss}/I)}{c(\text{ss})} = 1/5$$

$$p(\text{Jack}/\text{ss}) = \frac{c(\text{ss}/\text{Jack})}{c(\text{ss})} = 3/5$$

$$p(\text{do}/\text{ss}) = \frac{c(\text{ss}/\text{do})}{c(\text{ss})} = 1/5$$

$$p(\text{am}/I) = \frac{c(I/\text{am})}{c(I)} = 2/5$$

$$p(\text{like}/I) = \frac{c(I/\text{like})}{c(I)} = 1/5$$

$$p(\text{do}/I) = \frac{c(I/\text{do})}{c(I)} = 2/5$$

$$p(\text{ss}/\text{Jack}) = \frac{c(\text{Jack}/\text{ss})}{c(\text{Jack})} = 2/5$$

$$p(\text{ss}/\text{like}) = \frac{c(\text{like}/\text{ss})}{c(\text{like})} = 2/3$$

$$p(\text{ss}/\text{am}) = \frac{c(\text{am}/\text{ss})}{c(\text{am})} = 1/2$$

$$P(I \mid \text{Jack}) = c(\text{Jack} \mid I) / c(\text{Jack}) = 3/5$$

$$P(\text{like} \mid \text{do}) = c(\text{do} \mid \text{like}) / c(\text{do}) = 1/2$$

$$P(\text{Jack} \mid \text{like}) = c(\text{like} \mid \text{Jack}) / c(\text{like}) = 1/3$$

$$P(\text{Jack} \mid \text{am}) = c(\text{am} \mid \text{Jack}) / c(\text{am}) = 1/2$$

1) Jack I

2) Jack I & like

3) Jack I am Jack --

4) <sy do I like </sy

Trigram model:

A Trigram model approximates the probability of a word given all the previous words by using only the conditional probability of the preceding words while trigram model looks two words into the past.

Ex: <sy sun is a </sy

<sy is a kind </sy

<sy a kind soul </sy

<sy kind soul she </sy

<sy soul she will </sy

$$\text{for: } P(x, y, z) = P(x) P(y) [P(z/x, y)] = P(x, y, z)$$

$$\rightarrow P(\text{<sy, sun, is}) = P(\text{<sy, sun, is}) = 1/2 = 0.5$$

$$P(\text{sun, is, a}) = P(\text{sun, is/a}) = 1/3 = 0.333$$

$$P(\text{is, a, </sy}) = P(\text{is, a/is}) = 2/5 = 0.4$$

$$P(\text{<sy, is, a}) = P(\text{<sy, is/a}) = 1/3 = 0.333$$

$$P(rs, a, kind) = P(rs, a / kind) = \frac{2}{3} = 0.6666$$

$$P(a, kind, <rs>) = P(a, kind / <rs>) = \frac{2}{5} = 0.4$$

$$P(<rs>, a, kind) = P(<rs>, a / kind) = \frac{1}{3} = 0.3333$$

$$P(a, kind, soul) = P(a, kind / soul) = \frac{2}{3} = 0.666$$

$$P(kind, soul, <rs>) = P(kind, soul / <rs>) = \frac{2}{5} = 0.4$$

$$P(<rs>, kind, soul) = P(<rs>, kind / soul) = \frac{1}{3} = 0.333$$

$$P(kind, soul, she) = P(kind, soul / she) = \frac{2}{2} = 1$$

$$P(soul, she, <rs>) = P(soul, she / <rs>) = \frac{2}{5} = 0.4$$

$$P(<rs>, soul, she) = P(<rs>, soul / she) = \frac{1}{2} = 0.5$$

$$P(soul, she, will) = P(soul, she / will) = \frac{2}{1} = 2$$

$$P(she, will, <rs>) = P(she, will / <rs>) = \frac{1}{5} = 0.2$$

$$P = 0.5 \times 0.3333 \times 0.4 \times 0.333 \times 0.6666 \times 0.4 \times 0.333 \times 0.666 \times \\ 0.4 \times 0.333 \times 1 \times 0.4 \times 0.5 \times 2 \times 0.5$$

$$P = 0.00001403821 //$$

∴ Total probability  $P = 0.00001403821$

## Language Model Evaluation:

Before descriptive parameter estimation method further refinements of the basic n-gram modeling approach.

Let us consider the problem of judging the performance of a language model.

Unigram

|      |    |     |      |    |      |
|------|----|-----|------|----|------|
| This | is | Big | Data | AI | Book |
|------|----|-----|------|----|------|

Bi-gram

|      |    |        |          |         |         |
|------|----|--------|----------|---------|---------|
| This | is | is Big | Big Data | Data AI | AI Book |
|------|----|--------|----------|---------|---------|

Tri-gram

|      |        |             |             |              |
|------|--------|-------------|-------------|--------------|
| This | is Big | is Big Data | Big Data AI | Data AI Book |
|------|--------|-------------|-------------|--------------|

Text

N-gram

Data

1-gram

great info

2-gram

Nice meeting

3-gram

you

Nice to meet

4-gram

you

Evaluating the string:

as I want to eat chnese food <1st

No. of tokens = 8

No. of sentences = 1

No. of paragraphs = 0

No. of types = 1

|                  | N-gram | unigram | Bigram | Trigram |
|------------------|--------|---------|--------|---------|
| Text probability | 962    | 170     | 109    |         |

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

$$\text{Bigram}(P) = \frac{C(w_n)}{C(w_1)}$$

$$\text{Trigram}(T) = C\left(\frac{(w_1, w_2)}{w_n}\right)$$

### Parameter Estimation:

#### 1. Maximum-Likelihood estimation and smoothing:

so far, we have discussed estimating the mean and variance of a distribution. our methods have been somewhat and have more specifically. if it is not clear how we can estimate other parameters.

now we would like to talk about a systematic way of parameter estimation.

**Definition:** The process of redistributing probability means search to some small & non zero value is called smoothing. The most common smoothing technique is back-off. It includes splitting n-grams into those counts exceed the threshold.

- Ex:
  - <s> I am a human </s>
  - <s> I am not a stone </s>
  - <s> I I live in Lahore </s>

$$\rightarrow P(I \text{ am a human}) = P(I \mid \text{ssr}) P(\text{am} \mid I) P(\text{a} \mid \text{am}) P(\text{human} \mid \text{a})$$

$$= \frac{3}{3} \cdot \frac{2}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}$$

$$= 0.125$$

$$\rightarrow \text{split vocabulary: } P(I \text{ am human}) = P(I \mid \text{ssr}) P(\text{am} \mid I) P(\text{human} \mid \text{am})$$

$$= \frac{3}{3} \cdot \frac{2}{4} \cdot \frac{1}{2} = 0.25$$

General Bigram probability:  $P(x \mid y) = C(xy) / C(y)$

→ Bigram probability with "Laplace smoothing":

$$P(x \mid y) = \frac{C(xy) + 1}{C(y) + v}$$

where  $v = \text{distinct words}$

$$P(I \mid \text{am}) = \frac{2+1}{2+11} = 0.230$$

Trigram Training set:

denied the allegations

denied the reports

denied the claims

denied the request

→ probability using simple trigram model

Denied the - offer - ?

$$\rightarrow P(\text{offer} \mid \text{Denied the}) = \frac{\text{count(Denied the offer)}}{\text{count(Denied the)}} = 0/4 = 0$$

→ probability using "Laplace smoothing or add one smoothing"

$$P(\text{offer/denied the}) = \frac{\text{count(denied the offer)}}{\text{count(denied the)}} + v$$

$$= \frac{0+1}{4+6} = \frac{1}{10} = 0.1$$

$\Rightarrow$  steal probability mass to generalize better

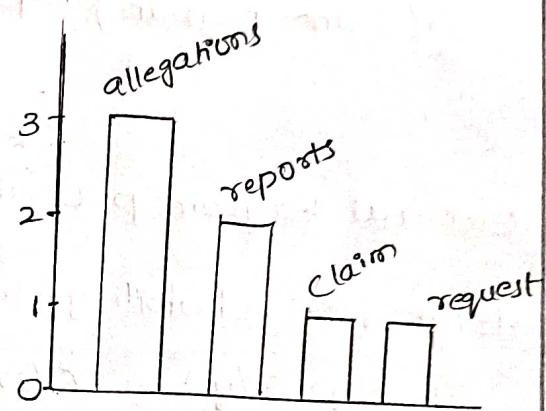
denied the allegation = 3

denied the reports = 2

denied the claim = 1

denied the request = 1

Total : 7



$\Rightarrow$  Bigram Example given table:

|         | i  | want | to  | eat | chinese | food | lunch | spend |
|---------|----|------|-----|-----|---------|------|-------|-------|
| i       | 5  | 827  | 0   | 9   | 0       | 0    | 0     | 2     |
| want    | 2  | 0    | 608 | 1   | 6       | 6    | 5     | 1     |
| to      | 2  | 0    | 4   | 686 | 2       | 0    | 6     | 211   |
| eat     | 0  | 0    | 2   | 0   | 16      | 2    | 42    | 0     |
| chinese | 1  | 0    | 0   | 0   | 0       | 82   | 1     | 0     |
| food    | 15 | 0    | 15  | 0   | 1       | 4    | 0     | 0     |
| lunch   | 2  | 0    | 0   | 0   | 0       | 1    | 0     | 0     |
| spend   | 1  | 0    | 1   | 0   | 0       | 0    | 0     | 0     |

$\rightarrow$  Bigram counts for 8 of words ( $V=1446$ )

I want to eat chinese food lunch spend.

2533 927 2417 746 158 1093 341 278

→ what is the probability of the sentence "I want to"

$$\begin{aligned}
 p(\text{I want to}) &= p(\text{want/I}) * p(\text{to/want}) * p(\text{to/want}) \\
 &= \frac{c(\text{I want})}{c(\text{I})} * \frac{c(\text{want want})}{c(\text{want})} * \frac{c(\text{want to})}{c(\text{to})} \\
 &= \frac{827}{2533} * \frac{0}{927} + \frac{608}{2417} = 0 //
 \end{aligned}$$

→ using Laplace or add 1 smoothed bigram model

$$= \frac{827+1}{2533+1446} * \frac{0+1}{927+1446} * \frac{608+1}{2417+1446} = 1$$

→ Table after adding 1 in all values:

|         | I  | want | to  | eat | chinese | food | lunch | spend |
|---------|----|------|-----|-----|---------|------|-------|-------|
| I       | 6  | 828  | 1   | 10  | 1       | 1    | 1     | 3     |
| want    | 3  | 1    | 609 | 2   | 7       | 7    | 6     | 2     |
| to      | 3  | 1    | 5   | 687 | 3       | 1    | 7     | 212   |
| eat     | 1  | 1    | 3   | 1   | 17      | 3    | 43    | 1     |
| chinese | 2  | 1    | 1   | 1   | 1       | 83   | 2     | 1     |
| food    | 16 | 1    | 16  | 1   | 2       | 5    | 1     | 1     |
| lunch   | 3  | 1    | 1   | 1   | 1       | 2    | 1     | 1     |
| spend   | 2  | 1    | 2   | 1   | 1       | 1    | 1     | 1     |

### Bayesian Parameter Estimation:

Bayesian probability estimation is an alternative parameter estimation method whereby the set of parameters of a model is itself viewed as a random variable.

governed by a prior statistical distribution. given training sample is:

training set:

"he is he is <sup>he is</sup> going abroad is going to study in the field"

$$\Rightarrow P(\text{is going abroad to study}) = ?$$

→ using bi-gram model:

$$\begin{aligned} &= P(\text{going | is}) * P(\text{abroad | going}) * P(\text{to | abroad}) * P(\text{study | to}) \\ &= \frac{\text{count(is going)}}{\text{count(is)}} * \frac{\text{count(going abroad)}}{\text{count(going)}} * \frac{\text{count(abroad to)}}{\text{count(abroad)}} * \frac{\text{count(to study)}}{\text{count(to)}} \\ &= \frac{2}{4} * \frac{1}{2} * \frac{0}{1} * \frac{1}{1} = 0 // \end{aligned}$$

→ using Laplace smoothing

$$\text{Distinct vocabulary } v = 9$$

→  $P(\text{is going abroad to study})$

$$\begin{aligned} &= \left( \frac{2+1}{4+9} \right) * \left( \frac{1+1}{2+9} \right) * \left( \frac{0+1}{1+9} \right) * \left( \frac{1+1}{1+9} \right) \\ &= \frac{3}{13} * \frac{2}{11} * \frac{1}{10} * \frac{2}{10} = 0.000839 \end{aligned}$$

→ using BayesPm:

"he is he is he is going abroad is going to study in the field"

$$\Rightarrow P(\text{is going abroad to study}) = ?$$

$$\text{total distinct words } v = 9$$

(is going, going abroad, abroad to, to study)

pairs of distinct values = 4 //

pairs of distinct values = 9

-> using Laplace smoothing add  $\nu=4$

$$= \left( \frac{2+1}{4+4} \right) * \left( \frac{1+1}{2+4} \right) * \left( \frac{0+1}{1+4} \right) * \left( \frac{1+1}{1+4} \right)$$

$$= \frac{3}{8} * \frac{2}{6} * \frac{1}{5} * \frac{2}{5} = \frac{12}{1200} = 0.01$$

=

Large-scale language model:

These language models (Lsm) estimates the relative likelihood of different phrases and are useful in many different NLP applications.

Ex: They have been used in Twitter bots for robot accounts to generate own sentences.

Defining Language models: The goal of probabilistic language modeling is to calculate the prob of a sentence or sequence of words.

$$p(w) = p(w_1 w_2 w_3 \dots w_n)$$

→ And can be to find the probability of the next word in the sequence.

$$p(w_5 | w_1 w_2 w_3 w_4)$$

Conditional probability:

Let A, B be two events with  $p(B) \neq 0$

conditional probability of A given B is:

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Chain Rule:

In general cases, the formula is as follows:

$$P(x_1, x_2, \dots, x_n) = P(x_1) P(x_2|x_1) \dots P(x_n|x_{n-1})$$

Ex:  $P(\text{"this water is so transparent"})$

$$\rightarrow P(\text{This}) * P\left(\frac{\text{water}}{\text{this}}\right) * P\left(\frac{\text{is}}{\text{water}}\right) *$$

$$P\left(\frac{\text{so}}{\text{this water is}}\right) * P\left(\frac{\text{transparent}}{\text{this water is so}}\right)$$

Markov Property:

$\Rightarrow$  stochastic process has a Markov Property if the conditional probability distribution of future states of the process depends only on the present state, not on sequence of events that proceeded if in other words probability of next word can be estimated given only the previous  $k$  no. of words.

$\rightarrow$  if  $k=1$ ;

$$P(\text{transparent} | \text{this water is so}) \approx P(\text{transparent} | \text{so})$$

or if  $k=2$ ;

$$P(\text{transparent} | \text{this water is so}) \approx P(\text{transparent} | \text{is so})$$

General eqn for all the Markov Assumption

if  $k=1$ ; then  $P(w_1 | w_1 w_2 \dots w_{1-1})$

Types of Language Models:

i) class-based language model:

class-based N-gram model represent a common meaning for improving the robustness of language modeling by providing adequate representation of words or sequences using Markov assumption not existing or under represented in training or corpus.

$$p(w_i | w_{i-1}) = p(w_i | c_i) p(w_i | c_{i-1})$$

Ex:  $p(I \text{ saw a cat on a mat})$

$$\rightarrow p(I)$$

$$p(a | I \text{ saw})$$

$$p(cat | I \text{ saw. a})$$

$$p(on | I \text{ saw a cat})$$

$$p(a | I \text{ saw a cat on})$$

$$p(mat | I \text{ saw a cat on a})$$

& after class-based model eliminate or split before word in the denominator in the bigram

$$p(I \text{ saw a cat on a mat})$$

$$p(I) \rightarrow p(I)$$

$$p(saw | I) \rightarrow p(saw | I)$$

$$p(a | I \text{ saw}) \rightarrow p(a | I \text{ saw})$$

$$p(cat | I \text{ saw a}) \rightarrow p(cat | saw a)$$

$$p(on | I \text{ saw a cat}) \rightarrow p(on | a \text{ cat})$$

$$p(a | I \text{ saw a cat on}) \rightarrow p(a | cat on)$$

$P(\text{mat} \neq \text{sqw} \text{ or } \text{cat} \text{ on a}) \rightarrow P(\text{mat on a})$

Estimating probability using sentence simplify bigram

using class-based model:

$\langle \text{S} \rangle \text{ I am sam } \langle \text{I} \rangle \text{S}$

$\langle \text{S} \rangle \text{ I am I am } \langle \text{I} \rangle \text{S}$

$\langle \text{S} \rangle \text{ I fly } \langle \text{I} \rangle \text{S}$

→ distinct words: I am, sam, I, fly

Bigram:  $\langle \text{S} \rangle$  and  $\langle \text{I} \rangle$  are also tokens.

There are  $6(u+2)$  tokens and  $6 \times 6 = 36$  bigrams.

$$\rightarrow P(I/\langle S \rangle) = \frac{2}{6}, P(\frac{\text{sam}}{\langle S \rangle}) = \frac{1}{6}, P(\frac{\text{am}}{\langle S \rangle}) = P(\frac{\text{fly}}{\langle S \rangle}) = 0, P(\frac{\langle S \rangle}{\langle S \rangle}) = 0$$
$$P(\langle I \rangle / \langle S \rangle) = 0, P(\frac{\text{sam}}{I}) = P(\frac{\text{am}}{I}) = P(\frac{\text{fly}}{I}) = P(\frac{\langle S \rangle}{I}) = P(\frac{\langle I \rangle}{I}) = 0$$

$$P(\frac{I}{I}) = 0, P(\frac{\text{sam}}{I}) = P(\frac{\text{am}}{I}) = P(\frac{\text{fly}}{I}) = P(\frac{\langle S \rangle}{I}) = P(\frac{\langle I \rangle}{I}) = 0$$
$$P(\frac{I}{\text{am}}) = 0, P(\frac{\text{sam}}{\text{am}}) = P(\frac{\text{am}}{\text{am}}) = P(\frac{\text{fly}}{\text{am}}) = P(\frac{\langle S \rangle}{\text{am}}) = P(\frac{\langle I \rangle}{\text{am}}) = 0$$

$$P(\frac{I}{\text{sam}}) = \frac{1}{2}, P(\frac{\text{sam}}{\text{sam}}) = 0, P(\frac{\text{am}}{\text{sam}}) = 0, P(\frac{\text{fly}}{\text{sam}}) = P(\frac{\langle S \rangle}{\text{sam}}) = P(\frac{\langle I \rangle}{\text{sam}}) = 0$$

$$P(\frac{I}{\text{fly}}) = 0, P(\frac{\text{sam}}{\text{fly}}) = P(\frac{\text{am}}{\text{fly}}) = P(\frac{\text{fly}}{\text{fly}}) = 0, P(\frac{\langle S \rangle}{\text{fly}}) = P(\frac{\langle I \rangle}{\text{fly}}) = 0$$

$$P(\frac{I}{\langle I \rangle}) = P(\frac{\text{sam}}{\langle I \rangle}) = P(\frac{\text{am}}{\langle I \rangle}) = P(\frac{\text{fly}}{\langle I \rangle}) = P(\frac{\langle S \rangle}{\langle I \rangle}) = P(\frac{\langle I \rangle}{\langle I \rangle}) = 0$$

trigram are  $6 \times 6 \times 6 = 216$

→ Compute the probability of sentence.

I want English food

$P(\text{I want English food})$

## ② Variable-length language model:

Since probabilities are less than or equal to 1, the more probabilities we multiply together, the smaller the product becomes. multiply enough n-grams together result in numerical underflow.

By using log probabilities instead of raw probabilities, we get numbers that are not as small.

Adding in log space is equivalent to multiplying in linear space, so we combine log.

→ adding is faster than multiplying the result of doing all computing and storage in log space is that we only need to convert back into prob if we need to report them at the end.

$$\log(P_1 P_2 P_3 P_4) = \log P_1 + \log P_2 + \log P_3 + \log P_4$$

$$P_1 * P_2 * P_3 * P_4 = \exp(\log P_1 + \log P_2 + \log P_3 + \log P_4)$$

## ③ Discriminative language model:

Standard n-gram models is a general model for assigning a probability to a given word sequence 'w'. However, in practical applications like machine translation or speech recognition, the task of a language model is often to separate good sentence hypotheses from bad sentence hypotheses such that word strings of wordly differing quality receive maximally.

Distinct probability estimates Recent attempts at such discriminative language model.

Ez:

Given data:

Jane saw will

Mary saw Jane

→ Collect data with POS tagging

Jane saw will      Mary saw Jane  
N      V      N      N      V      N

Create lookup table:

|      | N | V |
|------|---|---|
| Mary | 1 | 0 |
| saw  | 0 | 2 |
| Jane | 2 | 0 |
| will | 1 | 0 |

→ three training data set:

Mary <sup>N</sup> will <sup>M</sup> see <sup>V</sup> Jane <sup>N</sup>

will <sup>N</sup> will <sup>M</sup> see <sup>V</sup> Mary <sup>N</sup>

Jane <sup>N</sup> will <sup>M</sup> see <sup>V</sup> will <sup>N</sup>

Create lookup table:

|      | N | V | M |
|------|---|---|---|
| Mary | 2 | 0 | 0 |
| will | 2 | 0 | 3 |
| see  | 0 | 3 | 0 |
| Jane | 2 | 0 | 0 |

Bigram: N followed by M

M followed by V

V followed by N

• Bigram lookup table:

|           | N-M | M-V | V-N |
|-----------|-----|-----|-----|
| Mary-will | 1   | 0   | 0   |
| will-see  | 0   | 3   | 0   |
| see-Jane  | 0   | 0   | 1   |
| will-will | 1   | 0   | 0   |
| see-marry | 0   | 0   | 1   |
| Jane-will | 1   | 0   | 0   |
| see-will  | 0   | 0   | 1   |

→ let us an example proposed and appropriate tag sequence for a sentence.

"John can see will"

↓      ↓      ↓      ↓  
N    modal   verb   Noun

first noun word after comes modal and a modal after a verb,  
verb after noun.

→ "Tred will spot will"

↓      ↑      ↓      ↑  
N    m    v    N

tagged noun, modal, verb, noun. These probabilities are  
emission prob.

Ex(2): mary Jane can see will

spot will see mary

will Jane spot mary

mary will pat sat

- Create lookup table

|      | N | V | M |
|------|---|---|---|
| Mary | 4 | 0 | 0 |
| Jane | 2 | 0 | 0 |
| can  | 0 | 0 | 1 |
| see  | 0 | 2 | 0 |
| will | 1 | 0 | 3 |
| spot | 1 | 1 | 0 |
| pat  | 0 | 1 | 0 |
| sat  | 1 | 0 | 0 |

→ Identify Noun, Model(M), verb(V). Let's divide each word appeared by the total no. of every pos in the set of sentence.

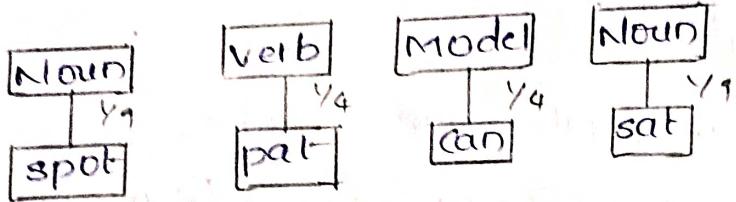
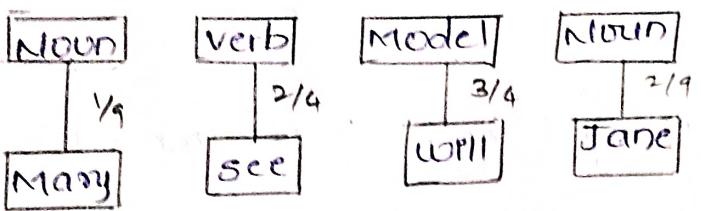
→ No. of N = 9, No. of V = 4, No. of M = 4

|      | N   | V   | M   |
|------|-----|-----|-----|
| Mary | 4/9 | 0/4 | 0   |
| Jane | 2/9 | 0   | 0   |
| can  | 0   | 0   | 1/4 |
| see  | 0   | 2/4 | 0   |
| will | 1/9 | 0   | 3/4 |
| spot | 1/9 | 1/4 | 0   |
| pat  | 0   | 1/4 | 0   |
| sat  | 1/9 | 0   | 0   |

→ Mary: N = 4/9, V = 0, M = 0

Jane: N = 2/9, V = 0, M = 0

;



$$\rightarrow \text{total probability} = \frac{4}{9} \cdot \frac{2}{4} \cdot \frac{3}{4} \cdot \frac{2}{9} \cdot Y_9 \cdot Y_4 \cdot Y_4 \cdot Y_9 \\ = 0.00694$$

#### ④ Syntax-Based Language Models:

→ Well known drawback of n-gram language model is that they take account of similar words in the history. That fall outside the limited window of the directly preceding  $(n-1)$  words. However Natural language exhibits many types of long distance dependency the choice of the current word is dependent on the words that are relatively removed in terms of sequence position.

our goal is to compute the probability of sentence or sequence of words  $w = (w_1, w_2, \dots, w_n)$

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

$$= P(w_n | w_1, w_2, w_3, \dots, w_{n-1})$$

→ And compute probability of a sequence and computing the conditional probability of a word given previous words

$$\text{Conditional probability: } P(B|A) = \frac{P(A, B)}{P(A)}$$

$$\rightarrow P(A, B) = P(A)P(B|A)$$

$$\rightarrow P(A, B, C, D) = P(A)P(B|A)P(C|A, B)P(D|A, B, C)$$

Unigram:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k)$

Ex:  $P(<\text{ss}> \text{the man from jupiter came } <\text{ss}>) \approx$   
 $P(\text{the}) P(\text{man}) P(\text{from}) P(\text{jupiter}) P(\text{came})$

Bigram:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$

Ex:  $P(<\text{ss}> \text{the man from jupiter came } <\text{ss}>) \approx$   
 $P(\text{the}/<\text{ss}>) P(\text{man}/\text{the}) P(\text{from}/\text{man}) P(\frac{\text{jupiter}}{\text{from}}) P(\frac{\text{came}}{\text{jupiter}})$   
 $\therefore P(<\text{ss}>/\text{jupiter came})$

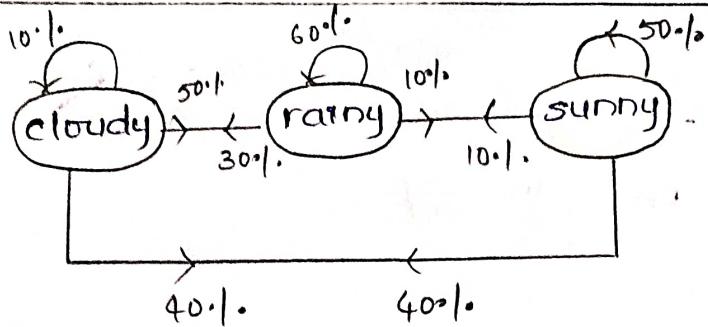
Trigram:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1}, w_{k-2})$

Ex:  $P(<\text{ss}> \text{the man from jupiter came } <\text{ss}>) \approx$   
 $= P(\frac{\text{man}}{\text{ss} \text{ the}}) P(\frac{\text{from}}{\text{the man}}) P(\frac{\text{jupiter}}{\text{man from}}) P(\frac{\text{came}}{\text{from jupiter}})$   
 $\qquad \qquad \qquad P(\frac{\text{ss} \text{ jupiter came}}{\text{jupiter came}})$

## ⑤ Max Ent language model:

This model combined with Max-likelihood-based prob estimation for language model is that the constraints imposed by estimates solely derived from the training data are too strong.

Ex: The diagram below represents three states representing weather of the day. (cloudy, rainy, and sunny).



→ There are different states such as cloudy, rainy, and sunny.  
 Following represents the transition probabilities based on above diagram.

- If sunny today, then tomorrow:

$$p(\text{sunny}) = 50\%$$

$$p(\text{rainy}) = 10\%$$

$$p(\text{cloudy}) = 40\%$$

- If rainy today, then tomorrow:

$$p(\text{rainy}) = 60\%$$

$$p(\text{sunny}) = 10\%$$

$$p(\text{cloudy}) = 30\%$$

- If cloudy today, then tomorrow:

$$p(\text{cloudy}) = 10\%$$

$$p(\text{rainy}) = 50\%$$

$$p(\text{sunny}) = 40\%$$

→ cloudy(wednesday) sunny(tuesday)

probability to cloudy(wednesday)

→ sunny - sunny(Tuesday) - cloudy(wed)

the prob to a cloudy on wednesday

$$= 0.5 * 0.4 = 0.2$$

→ sunny - rainy (tuesday) - cloudy (wed)

the prob to a cloudy on wednesday

$$= 0.1 * 0.3 = 0.03$$

→ sunny - cloudy (tuesday) - cloudy (wed)

the prob to a cloudy on wednesday

$$\Rightarrow 0.1 * 0.4 = 0.04$$

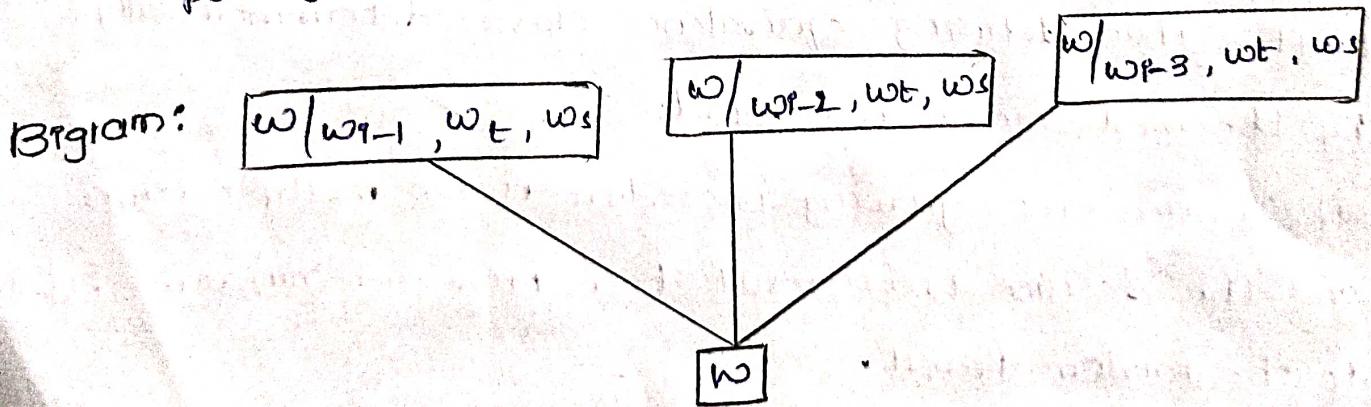
→ total prob to a cloudy on wednesday =  $0.2 + 0.03 + 0.04 = 0.27$

⑥ Factored language models:

The factored language model (FLM) approach [62, 63] builds on the observations that the prediction of words is dependent on the surface form of the preceding words and that better generalizations can be made when additional information, such as the word's POS or morphological class, is taken into account. In particular, word n-gram counts might be insufficient to robustly estimate the probability of word  $w_i$  given word  $w_{i-1}$ , but if we know that word  $w_{i-1}$  is of a particular class, say a determiner, we might be able to obtain a good probability estimate for  $P(w_i/\text{determiner})$ .

This is reminiscent of the class-based models described previously; however, in an FLM, many such class-based estimates are combined and structured hierarchically through the use of a generalized backoff strategy.

- > PFLMs assume a factored word representation, where words are considered feature vectors rather than individual surface forms; that is,  $w \in \mathbb{R}^{k \times k}$ .
  - > It is an extension of conventional language model.
  - > The probabilistic model provides where the prediction of a factor 'f' is based on 'n' parents ( $f_1 - f_n$ ).
  - > If 'w' represents a word and 't' represents POS tag expression is  $P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n})$ .
  - > PFLM used to allow stemming, the word and provides pos tagging like morphological information.
- Ex: word : stock prices are rising  
 stem : Stock price is rise  
 pos tag: N N-Plural v-plural v-part

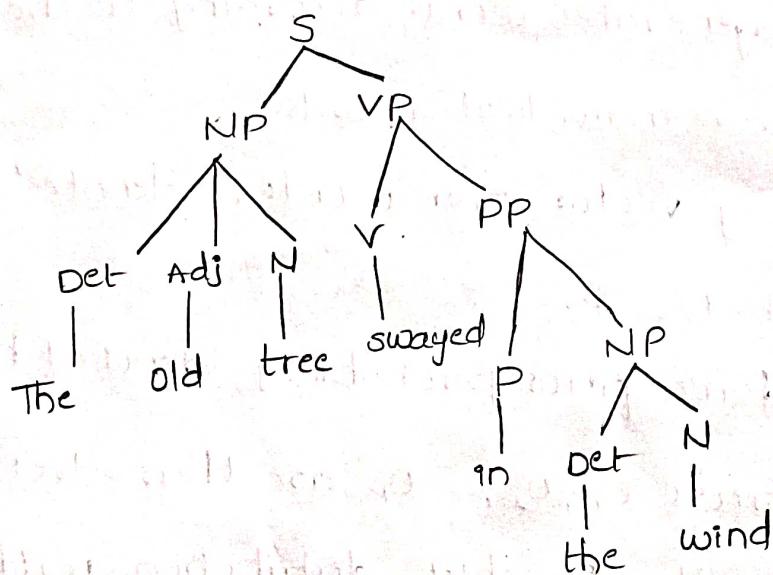


② Other Tree-Based Language Models:

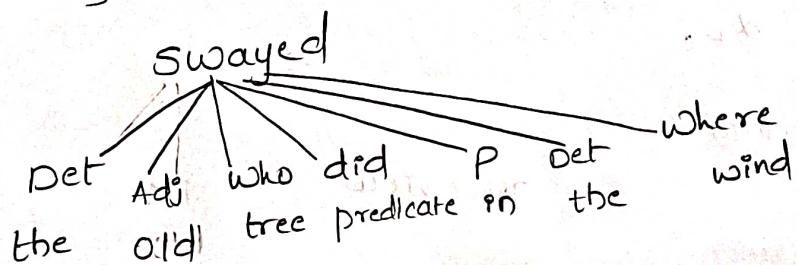
These model can be used for regular expression using CFG rules to generate prediction applying pos tagging to classify predicting categorization values.

- Random Forest language (RFLMs) Models, (produced) proposed by Xu and Jelinek, model all word histories seen in the training data as a collection of randomly grown decision trees (a random forest).
- Nodes in the decision trees are associated with sets of histories, with the root node containing all histories.
- Trees are grown by dividing the set of histories into two subsets according to the word identity at a particular position in the history.
- Out of a number of possible splits chosen to undergo the log-likelihood test is selected randomly as well.
- After the growing process has stopped, each leaf node in a decision tree can be viewed as a cluster of similar word histories forming an equivalence class.
- Rather than defining equivalence classes deterministically by the most recent  $n-1$  words.
- The decision tree-growing procedure is run multiple times, and the decision trees result from each run are added to the random forest.
- Supposing that  $M$  decision trees have been obtained, the RFLM probabilities are computed as averages of the individual decision tree probabilities.

Ex: The old tree swayed in the wind  
 who predicate where.  
 ↓ did



dependency tree:



Bayesian Topic-Based Language Models:

- Bayesian Topic-Based Language modeling is a significant recent trend in statistical language modeling.
- A significant recent trend in statistical language modeling is Bayesian modeling of latent topic structure in documents.
- It is a probabilistic algorithm based on applying Bayes theorem with Naïve assumption of conditional independence between every pair of feature. Given example training data.
- One of the first models of this type was the latent Dirichlet allocation model proposed by Blei, Ng, & Jordan.

- The LDA model assumes that a document is composed of  $k$  topics, denoted as  $z_1, \dots, z_k$ .
- Each topic generates words according to the topic-specific distribution over individual words.
- The probability vector over words is denoted by  $\theta_k$  for each topic  $k=1, \dots, K$ .
- Each topic has a prior probability, denoted by  $\phi_k$ .
- The topic priors  $\phi_1, \phi_2, \dots, \phi_K$  are themselves distributed according to the Dirichlet distribution with hyperparameters  $\alpha_1, \dots, \alpha_K$ .

Ex: Given training data

| Text                               | Review   |
|------------------------------------|----------|
| I liked the movie                  | positive |
| It's a good movie nice story       | positive |
| Nice songs but sadly boring ending | Negative |
| Hero acting is bad overall nice    | positive |
| movie heroine acting is good       | Negative |
| sad boring movie                   |          |

→ Classify the text and calculate the probability of  $p(\text{"overall liked the movie"})$ . Is.

$$P(+ve) = 3/5, P(-ve) = 2/5$$

→  $P(\text{overall liked the movie}) = P(\text{overall}) P(\text{liked}) P(\text{the}) P(\text{movie})$

→ Calculate  $P(\text{overall liked the movie} | +ve) =$

$$P\left(\frac{\text{overall}}{+ve}\right) P\left(\frac{\text{liked}}{+ve}\right) P\left(\frac{\text{the}}{+ve}\right) P\left(\frac{\text{movie}}{+ve}\right)$$

$$= Y_3 \cdot Y_3 \cdot Y_3 \cdot Y_3 = 0.04938 //$$

→ Calculate  $P(\text{overall liked the movie} | -ve) =$

$$P\left(\frac{\text{overall}}{-ve}\right) P\left(\frac{\text{liked}}{-ve}\right) P\left(\frac{\text{the}}{-ve}\right) P\left(\frac{\text{movie}}{-ve}\right)$$

$$= Y_2 \cdot Y_2 \cdot Y_2 \cdot Y_2 = 0.25 //$$

### ⑨ Neural Network Language Models:

Neural Network language models (NNLMs) take a different approach: discrete word sequences are first mapped into a continuous representation, and n-gram probabilities are then

estimated within this continuous space.

→ The assumption is that words having similar distributional properties will receive similar continuous representations, which will in turn result in smoother probability estimates.

→ The Neural Network is typically a multilayer perceptron with an input layer, projection layer, a hidden layer, and an output layer of nodes.

→ A graphical representation of the architecture of an NNLM is shown below.

- Adjacent layers are fully interconnected by weights
- For a vocabulary of  $V$  words, the input consists of a concatenation of  $n-1$ ,  $V$ -dimensional binary feature vectors representing the history of  $n-1$  words (eg: The first two words in a trigram).
- The projection layer ' $i$ ' of a given fixed dimensionality  $d'$  encodes the shared continuous representation of the words, which is learned during training.
- The hidden layer ' $h$ ' also has a fixed no. of  $J$  nodes.

$$h_j = \tanh \left( \sum_{k=1}^d w_{jk}^h i_k + b_j^h \right) \quad \forall j, i = 1, \dots, J$$

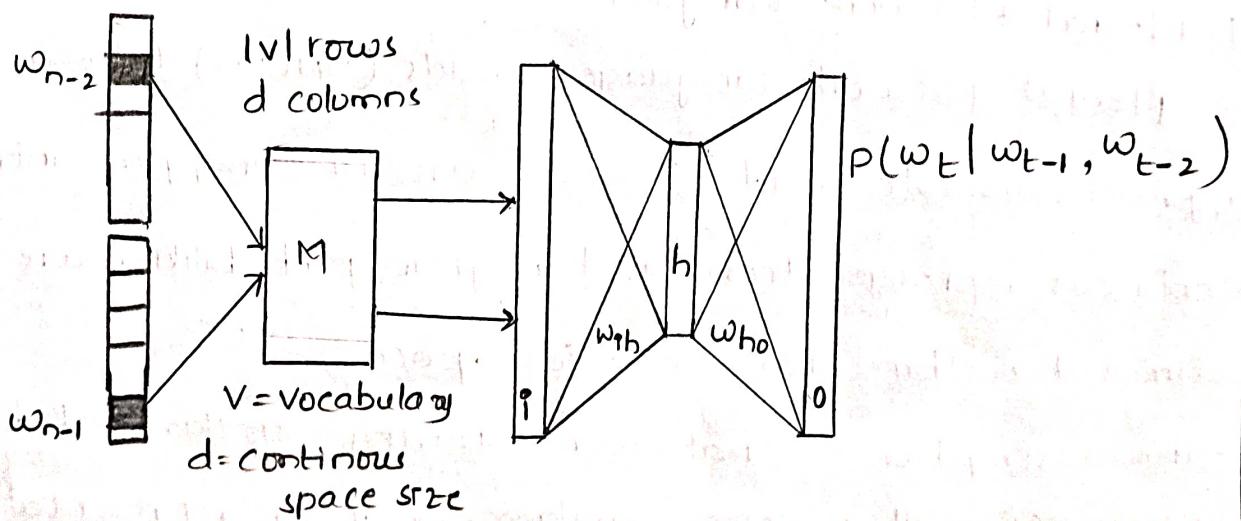


Fig: NNLM

### Language - Specific Modeling Problems:

The majority of language modeling research has focused on the English language.

- However speech and language processing technology has been ported to a range of other languages.

→ some of which have highlighted problems with the standard n-gram modeling approach and have necessitated modifications to the traditional language modeling framework. In this section, we look at three types of language-specific problems: morphological complexity, lack of word segmentation and spoken versus written languages.

### ① Language Modeling for Morphologically Rich Languages:

- A morphologically rich language is characterized by a larger number of different unique word forms in relation to the no. of word tokens in a text that is due to productive morphological processes in the language.
- A morpheme is the smallest meaning-bearing unit in a language.
- Morphemes can be either free (i.e., they can occur on their own) or they are bound.
- Morphological processes include compounding (forming a new word out of 2+ independently existing free morphemes), derivation (combination of a free and a bound morpheme to signal a particular grammatical feature).
- Many languages have rich inflectional paradigms.
- Morphological complexity causes problems for language modeling due to the high ratio of types of tokens, which

results in a lack of training data: many n-grams in the test data are not seen at all in the training data or are not observed frequently enough, leading to unreliable probability estimates.

Table : MSA inflectional paradigms for present tense verb forms and possessive pronouns (affixes are separated from the word stem by hyphens)

| Word         | Meaning           | Word        | Meaning         |
|--------------|-------------------|-------------|-----------------|
| 'a-skun(u)   | I live            | kitaab-iy   | my book         |
| ta-skun(u)   | you(MASC) live    | kitaabu-ka  | your(MASC) book |
| ta-skun-iyna | you(FEM) live     | kitaabu-ki  | your(FEM) book  |
| ya-skun(u)   | he lives          | kitaabu-hu  | hrs book        |
| ta-skun(u)   | she lives         | kitaabu-haa | her book        |
| na-skun(u)   | we live           | kitaabu-na  | our book        |
| ta-skun-uwna | you(MASC-PL) live | kitaabu-kum | your book       |
| ya-skun-uwna | they live         | kitaabu-hum | their book      |

## ⑪ Selection of subword units:

- The identification of subword units can be performed in a data-driven, unsupervised way
- It can be based on linguistic information (e.g. morphological analyzer), or it can be based on combination of both.
- Linguistically based methods typically involve a handcrafted morphological analyzer tool, such as the Buckwalter Morphological Analyzer, which analyzes each word form into its

morphological components.

- Data-driven approaches may incorporate varying degrees of information about the specific language.
- Some approaches are intended to discover units corresponding to linguistically defined morphemes.
- Automatic algorithms for identifying linguistic morphemes are as old as Zellig Harris's approach of estimating the perplexity of different letters.
- Morfessor package is one example. It attempts to derive a morpheme inventory  $M'$  from a corpus  $C$  by maximizing its posterior probability:

$$M = \underset{M'}{\operatorname{argmax}} P(M/C) = P(C/M)P(M)$$

### III Modeling with Morphological Categories:

- Most of the work on language modeling with subword units has focused on linear decomposition of words, primarily for agglutinative languages.
- As a result subword units are most frequently used in a standard n-gram model.
- One problem is that the n-gram context needs to be increased in order to be able to model interword dependencies in addition to dependencies b/w subword units.
- This in turn places demands on the amount of training data needed.

- The language model assigns probabilities to sequences of entire words, but does so by taking into account the constraints proved by morpheme-based feature functions.
- This model was shown to provide slightly better results (0.3% absolute word-error rate reduction on Broadcast News recognition task) than a discriminative language model using word-based feature.
- The language model predicts the probability of the word string and the associated tree jointly, using MaxEnt probability estimation.
- The model was evaluated on a English-to-Arabic translation task improved the points over a word trigram model.

#### IV Languages without Word Segmentation:

- Although agglutination produces a large number of long and complex word forms in many languages, other languages do not possess any explicit segmentation of character strings into words at all.
- In languages such as Chinese and Japanese, sentences are written as sequence of characters delimited by punctuation signs but without any intervening whitespaces.
- Rather than trying to match linguistically defined words, it is also possible to optimize segmentation directly for language model performance.

→ Sproat et al. showed that the dictionary used for Chinese word segmentation significantly influences the perplexity of a bigram language model trained on the segmented text and that is possible to iteratively optimize the dictionary by merging frequent word co-occurrences.

such that the perplexity is reduced at each iteration.

Note: such approaches are similar to the data-driven algorithms for deriving morpheme-like subword units.

→ Another example of a data-driven approach, in this case for Japanese, uses chunks of characters for language modeling.

→ chunks are derived from training data by selecting the highest-frequency n-grams and patterns.

## ⑥ Spoken versus Written Languages:

→ statistical language modeling crucially relies on large amounts of written text data, and significant trend within language modeling research is the development of methods for scaling current language modeling techniques to ever-larger databases. However, many of the world's 6,900 languages are spoken languages, that is, languages without a writing system.

They are either indigenous languages without a literary

tradition.

- For example, with the many dialects of Arabic which are used in everyday conversation but are almost never found in written form.
- Other languages may be spoken as well as written but may not have a standardized orthography.
- Only way of obtaining language model training data is to manually transcribe the language or dialect.
  - This is a costly and time-consuming process because it involves (i) the development of a writing standard
  - (ii) training native speakers to use in writing system accurately
  - (iii) the actual transcription effort.

## Multilingual and Crosslingual Language Modeling

### ① Multilingual Language Modeling

Up to this point we have discussed problems that arise when tailoring statistical language models to a particular language or language type, such as agglutinative languages or languages without word segmentation.

- The tacit assumption has been that the resulting language model is used in an application where only the language of interest is encountered.
- Here speakers may use several languages or dialects

side by side, often within the same utterance.

→ The phenomenon of code switching exists in a variety of bilingual and multilingual communities or where diglossia exists, such as where a formal std language is used in addition to colloquial or dialectal varieties.

## ② Crosslingual Language Modeling:

Another question that might be asked is whether data in one language can be leveraged to improve a language model for a different language.

→ The most straightforward way of applying this idea is to automatically translate the foreign-language text into the desired language, and to use the translated text as additional language training data.

→ A unigram extracted from the translated text was interpolated with a trigram baseline language model.

→ The use of machine translation technology for processing out-of-language data is likely to fail when the desired language does not have sufficient data to train a machine translation system.

→ A drawback of the previous approaches is that the quality of the resulting model is heavily dependent on the translation accuracy

$$\text{for: } P_{EN}(w) = \sum_k \phi_k^{EN}(w) \theta_k^{\text{ch}}$$

Assuming that the source language is Chinese (ch)

and target language is English (en).

$P_{EN}(w)$  is the marginal word probability prob  
distribution of english.

$\phi_k(w)$  is the probability assigned to word  $w$

$\theta_k$  is the prior for the  $k^{\text{th}}$  topic.

$\rightarrow$  target-language marginals are then incorporated into  
the target-language model as follows:

$$P_{\text{target}}(w/h) \propto \left( \frac{P_{\text{BLSA}}(w)}{P_{\text{base}}(w)} \right)^B P_{\text{base}}(w/h)$$

where  $P_{\text{BLSA}}$  is the adapted probability

$P_{\text{base}}$  is the baseline probability

$\rightarrow$  this approach enforces a one-to-one topic correspondence  
across languages.