

Natural Language Processing and Applications

Unit 2

Word Level and Syntactic Analysis

Word Level Analysis: word classes - part of speech tagging: HMM
 POS tagging; Named Entities (NE) - Conditional Random Field
 NE recognizers; syntactic analysis: constituency - context - free
 grammar - grammar rules - tree banks; parsing: top - down - bottom -
 up - ambiguity - CYK parsing - shallow parsing - dependency
 Parsing

* Word Classes

- Word classes are used in POS tagging. This gives a significant amount of information about the word and its neighbors.
- There are 2 broad classes - (i) Closed class
 (ii) Open class

A. Closed Class - have relatively fixed membership

- mostly grammatical words to deal with the formation of sentences
- show the structural relationships words have with one another.

B. Open Class

- lexical words that deal with content and vocabulary
- They have a concrete meaning that go beyond their function in a sentence.

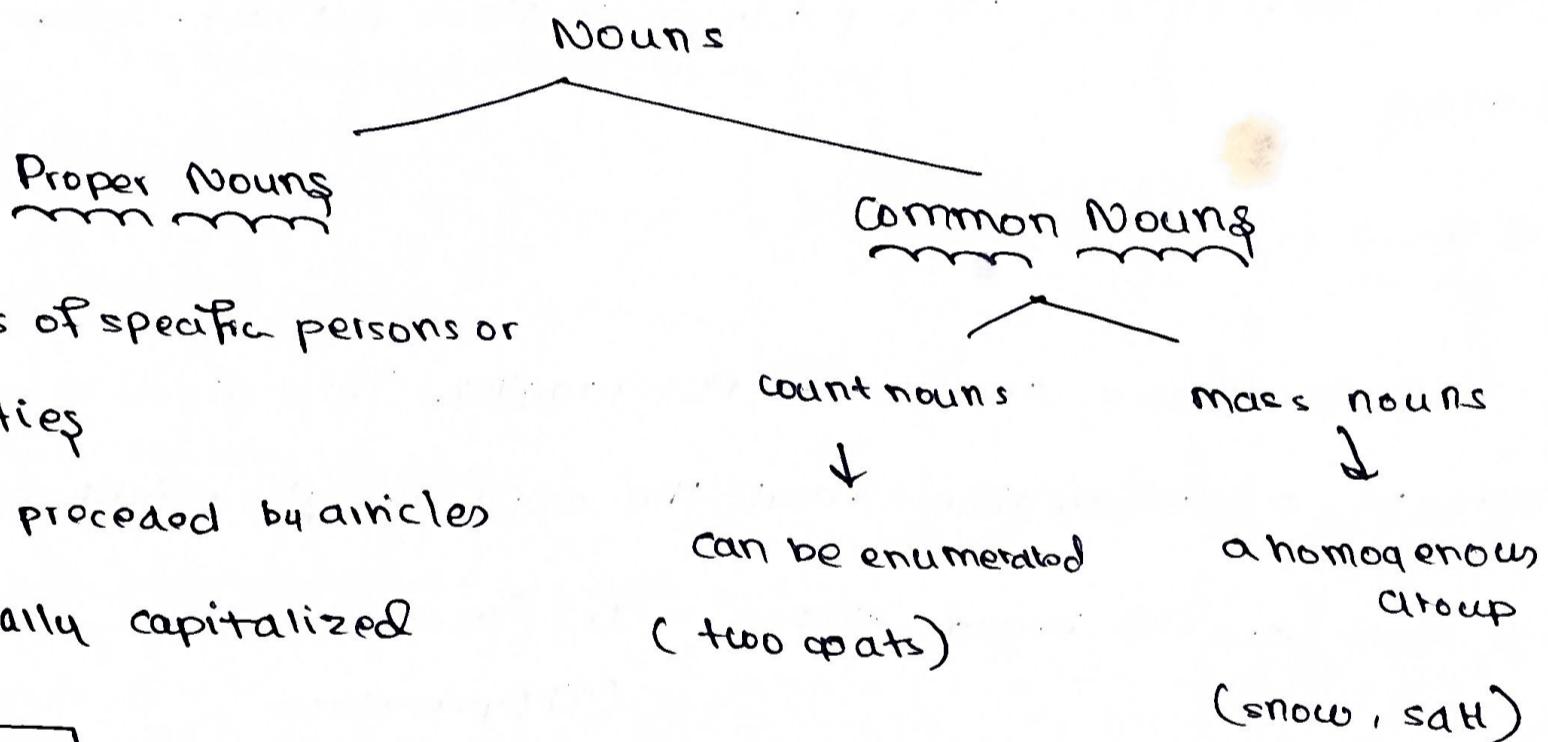
- includes nouns, verbs, adjectives and adverbs

* Open Class

A. Noun

- words for people, places and things
- nouns can be (i) concrete terms
(ii) abstractions
- can occur with determiners (a boat)
- can take possessiveness (Pooja's cat)
- can take the plural form

Q2



B. Verbs

→ refers to actions → processes

→ can take a number of morphological forms

(i) non - 3rd - person singular = eat

(ii) 3rd Person singular = eats

(iii) Progressive = eating

(iv) Past participle = eaten

c. Adjectives

- terms describing properties or qualities
- can modify nouns (large elephant, yellow cab, pale face)

D. Adverbs

→ words that modify verbs. There are 4 kinds of adverbs:

(i) Directional : specify direction / location of some action
 eq. down hill

(ii) Degree - specify the extent of some action
 eq. extremely, very, somewhat

(iii) Manner - describe the manner of some action / process
 eq. slowly, delicately

(iv) Temporal adverbs - describe the time that some action took place
 eq. yesterday, now, daily

* Closed Classes

A. Prepositions → indicate spatial or temporal relations
 eq. on, in, to, from, by

B. Particles → A particle is a word that resembles a preposition or an adverb.

→ It is often combined with a verb to form a larger unit called a phrasal verb. eq. " I went on ... ".
 " To throw off sleep "

C. Articles → used to begin a noun phrase

a, an ⇒ indefinite noun phrase

the ⇒ definite

D. Conjunctions → used to join two phrases, clauses or sentences

(i) Coordinating conjunctions: include and, or, but to join two elements of equal status.

(ii) Subordinating conjunctions: when one of the elements is of an embedded status. "I thought that you might like some milk."

E. Pronouns → act as a shorthand for referring to some noun phrase or entity or event

(i) Personal Pronouns: persons or entities (you, she, I, it, MR)

(ii) Possessive Pronouns: a form of personal pronouns indicating actual posses or just an abstract between the person and some objects (my, your, his, her, its, one's, our)

(iii) Non-pronouns - used in question forms (what, who, whom, whoever)

F. Auxiliary Verbs → marks certain semantic features of a main verb, including, (i) whether the action is in the past, present or future
(ii) whether it is completed

(iii) whether it is completed

e.g. be, can, may, must, will, should

Other closed classes

1. Interjections - oh, ah, hell, phew, alas
2. Negatives - no, not
3. Politeness markers
4. greetings

See. ~~class~~

* POS Tagging - The process of assigning a part-of-speech or other lexical class marker to each word in a corpus.

- Some tagsets in English for POS tagging are : (i) Brown Corpus
 (ii) Penn Treebank
 (iii) C5 Tagset
 (iv) C7 Tagset

* Tagsets for English

NN → preposition

EX → existential there

JJ → adjective

JJR → comparative adjective

JJS → superlative

NN → noun

NN S → noun, plural

NNP → proper noun

NNPS → proper noun, plural

PP → personal pronoun

PP\\$ → possessive pronoun

DT → determiner → that, the

RB → adverb

articles (a, an)

RBR → adverb - comparative

RBS → adverb - superlative

RP → particle (up, off)

MD → can, would, should, modal

VB → base verb

VBD → verb, past tense

VBG → gerund

VBN → past participle

VBP → Verb - non 3rd sg np

VBZ → Verb - 3 sg

Example - Tag the following sentence using the Penn Treebank

1. Book that flight

Book | VB that | DT flight | NN

2. The grand jury commented on a number of other topics.

The | DT grand | JJ jury | NN commented | VBD on | IN

a | DT number | NN of | IN other | JJ topics | NNS

3. There are 70 children there

• There | EX are | VBP 70 | CD children | NNS there | RB
 $\underbrace{\text{cardinal no}}$ $\underbrace{\text{adverb}}$

* Tagging

→ refers to the process of assigning a part-of-speech or other lexical marker to each word in a corpus

→ It is also applied to punctuation marks.

→ Tagging for NL is the same process as tokenization for computer language, though the tags for NL are much more ambiguous

→ Taggers are used in speech recognition, NL Parsing and information retrieval.

7

→ Automatically assigning a tag to a word is not trivial.

e.g. book = verb/noun

that = determiner / complementizer

- The problem of POS tagging is to resolve these ambiguities, choosing the proper tag for the context.
- Many of the ambiguous tokens are easy to disambiguate since the various tags associated w/ words are not equally likely.

* POS Tagging Algorithms

→ Tagging algorithms fall into 2 classes:

A. Rule-based Taggers

- involves a large database of hand-written disambiguation rule specifying . e.g. Enriched Tagger

B. Stochastic Tagger - resolve tagging ambiguities by using a training corpus to count the probability of a given word having a given tag in a given context.

Brill Tagger - has features of both tagging architectures

Example : Find one tagging error in each of the following sentences that are tagged with the Penn Treebank tagset

1. I / PRP need / VBP a / DT flight / noun from / IN Atlanta / NN

2. Does / VBZ this / DT flight / noun serve / VB dinner / NNS

3. I / PRP have / VB a / DT friend / noun living / VBGi in / IN

4. Can / VBP you / PRP list / VB the / DT nonstop / JJ afternoon / NN flights / NNS ?

error: can, could, should = modal

* Hidden Markov Model (HMM) Tagging

- HMM is a special case of Bayesian inference.
- Given a sequence of words w_1, \dots, w_n , in a sentence, the objective is to find the best sequence of tags which corresponds to that sequence.
- With the probabilistic/Bayesian view: consider all possible sequences of tags. Out of this sequence choose the most probable one.

→ Let $T = t_1, t_2, \dots, t_n$

$$w = w_1, w_2, \dots, w_n$$

Out of all the sequence of tags t_1, \dots, t_n

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

By Baye's rule:

$$P(x|y) = \frac{P(y|x) P(x)}{P(y)}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$$

The dr. is not needed, as it does not change across tag sequences

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \underbrace{P(w_1^n | t_1^n)}_{\text{likelihood}} \underbrace{P(t_1^n)}_{\text{prior}}$$

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} \cdot P(w_1^n | t_1^n) \underbrace{P(t_1^n)}_{\text{likelihood}} \underbrace{P(t_1^n)}_{\text{prior}}$$

Bigram assumption

The probability of a word appearing depends only on its own POS tag.

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$$

The probability of a tag appearing depends only on the previous tag

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

* Two Kinds of Probabilities

① Tag Transition Probabilities - $P(t_i | t_{i-1})$

$$P(t_i | t_{i-1}) = \frac{c(t_{i-1}, t_i)}{c(t_{i-1})}$$

② Word Likelihood Probability = $P(w_i | t_i)$

$$P(w_i | t_i) = \frac{c(t_i, w_i)}{c(t_i)}$$

* Viterbi Algorithm

~~(X)~~ Using the HMM model for POS tagging on the following corpus.

1. Mary Jane can see will.

2. Spot will see Mary

3. Will Jane spot Mary?

4. Mary will pat Spot.

Step 1 : Assign POS to each word in all the sentences

1. Mary Jane can see will

n n Modal v n

2. Spot will see Mary .

n Modal v n

3. Will Jane spot Mary

~~W~~ M n v n

4. Mary will pat Spot

n Modal v n

Step 2 : Compute the emission probabilities.

	N	M	V
Mary	4	0	0
Jane	2	0	0
will	21	3	0
spot	2	0	1
can	0	1	0
see	0	0	2
Pat	0	0	1

↓.

④ see vertically

	N	M	V
Mary	4/9	0	0
Jane	2/9	0	0
will	1/9	. 3/4	0 1/4
spot	0	0	0 1/4
can	0	1/4	2/4 0
see	0	0	1/2
Pat	0	0	1/4

Step 3: Compute the transition probability

(i) $\langle S \rangle \quad n \quad n \quad m \quad v \quad n \quad \langle E \rangle$
 Mary Jane can see will

(ii) $\langle S \rangle \quad n \quad m \quad v \quad n \quad \langle E \rangle$
 Spot will see Mary

(iii) $\langle S \rangle \quad m \quad n \quad v \quad \cancel{n} \quad \langle E \rangle$
 Will Jane spot Mary

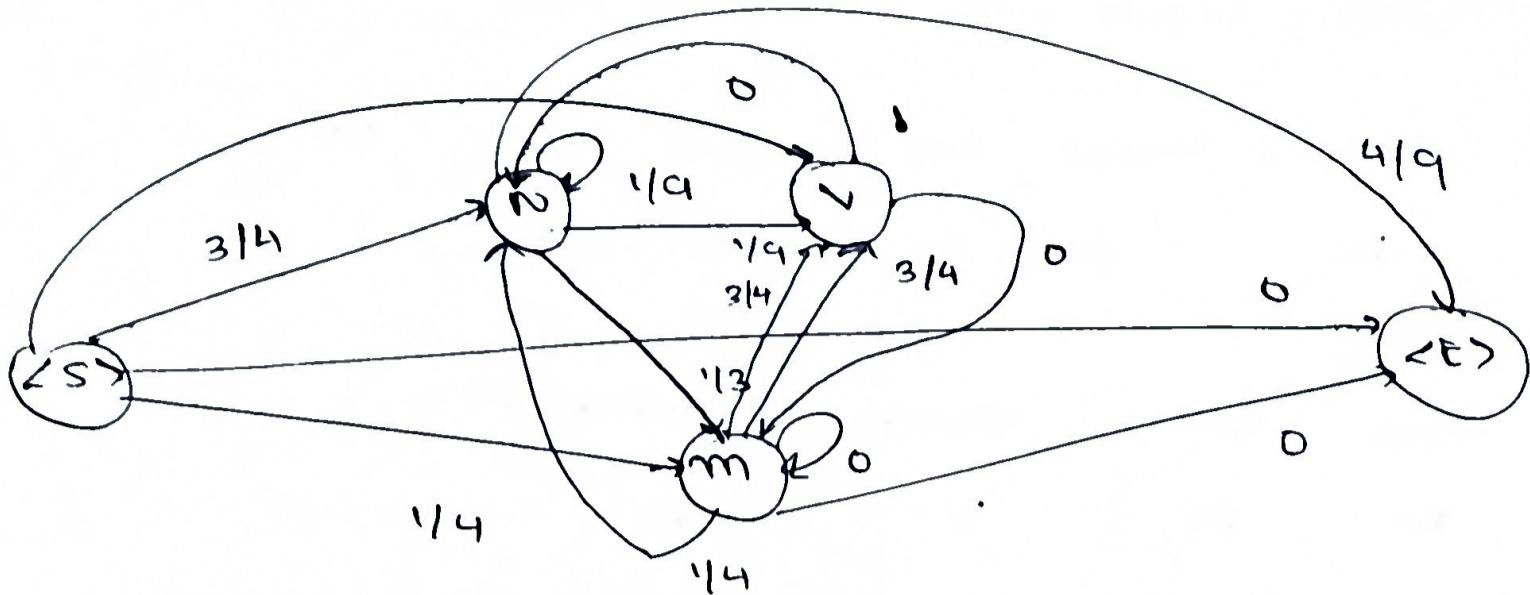
(iv) $\langle S \rangle \quad \cancel{n} \quad m \quad v \quad n \quad \langle E \rangle$
 Mary will pat spot

	n	m	v	$\langle E \rangle$
$\langle S \rangle$	2	1	0	0
n	1	3	1	4
m	1	0	3	0
v	8/4	0	0	0

④ see horizontally

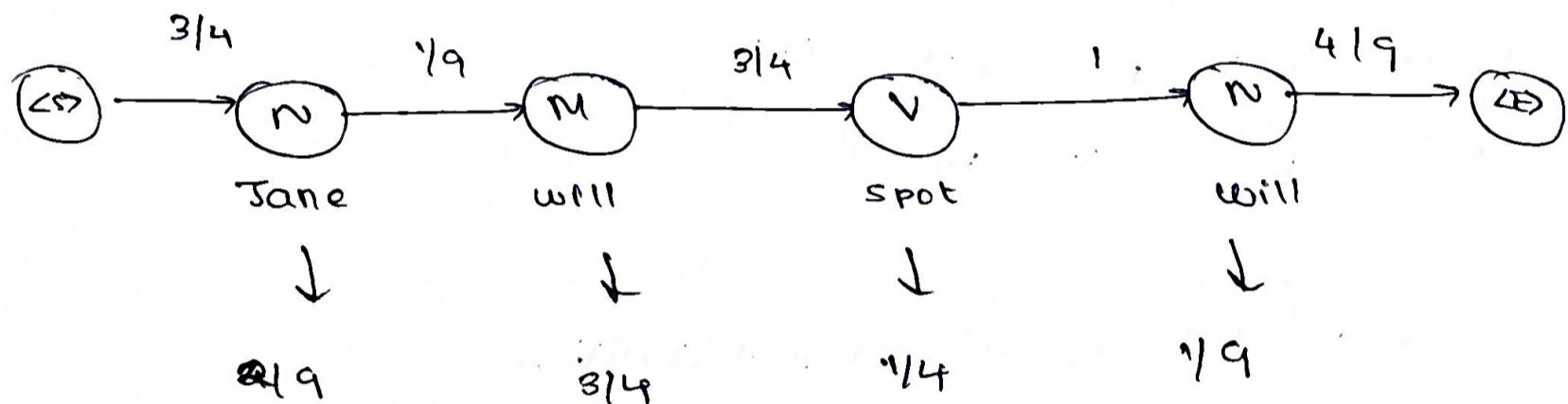
$\langle S \rangle$	n	m	v	$\langle E \rangle$
	3/4	1/4	0	0
n	1/4	1/3	1/9	4/9
m	1/4	0	3/4	0
v	1	0	0	0

Step 4 : Draw the transition diagram



Step 5 : Consider a testing sentence

Jane will spot will



$$\text{multiply} = \underline{0.0003858}$$

Example 2 : Tag the following sentence using the Viterbi algorithm.
with the following training corpus.

1. Martin Justin can watch will
2. spot will watch Martin
3. will Justin spot Martin
4. Martin will pat spot

Test: Justin will spot will .

Step 1 Training corpus with correct POS Tag.

<S> Martin Justin can'. watch will <E>
N N M V N
.

<S> Spot will watch Martin <E>
N M V N

<S> will Justin spot Martin <E>
M N V N

<S> Martin will pat spot <E>
N M V N

Step 2 : Create table of emission probabilities

	N	M	V
Martin	4/9	0	0
Justin	2/9	0	0
will	4/9	3/4	0
spot	2/9	0	1/4
can	0	1/4	0
watch	0	0	2/4
Pat	0	0	1/4

Step 3 Create table of state transition probability

(17)

$\angle S \rangle$	N	M	V	$\angle E \rangle$
	3/4	1/4	0	0
N	1/9	1/3	1/9	4/9
M	1/4	0	3/4	0
V	1	0	0	0

Testing Justin will spot will

N M V N
N N M

$$(i) \underline{\text{Justin}} \quad P(N | \text{Justin}, \angle S \rangle)$$

$$= P(\text{Justin} | N) \times P(N | \angle S \rangle)$$

$$= \frac{2}{9} \times \frac{3}{4} = \frac{1}{6}$$

see vertically

(N)

$$(ii) \underline{\text{Will}} \quad P(N | \text{will}, \angle S \rangle)$$

$$= P(\text{will} | N) \times P(N | \angle S \rangle)$$

$$= \frac{1}{9} \times \frac{1}{9} \times \frac{1}{6} = \frac{1}{486}$$

$$P(M | \text{will}, \angle S \rangle)$$

$$= P(\text{will} | M) \times P(M | \angle S \rangle)$$

$$= \frac{3}{4} \times \frac{1}{3} \times \frac{1}{6} = \frac{1}{24}$$

$$\frac{1}{24} > \frac{1}{486} \Rightarrow (M)$$

$$\begin{aligned}
 \underline{\text{spot}} \quad & P(N | \text{spot}, M) \\
 & = P(\text{spot}|N) \times P(N|M) \\
 & = 2/9 \times 1/4 \times 1/24 = 1/432
 \end{aligned}$$

$$P(V | \text{spot}, M)$$

$$\begin{aligned}
 & = P(\text{spot}|V) \times P(V|M) \\
 & = 2/9 \times 3/4 \times 1/24 = 1/432
 \end{aligned}$$

$$\frac{1}{108} > \frac{1}{432} \Rightarrow \textcircled{v}$$

$$\begin{aligned}
 \underline{\text{will}} \quad & P(N | \text{will}, V) \\
 & = P(\text{will}|N) \times P(N|V) \\
 & = \frac{1}{9} \times 1 \times \frac{1}{108} = \frac{1}{108}
 \end{aligned}$$

$$P(M | \text{will}, V)$$

$$\begin{aligned}
 & = P(\text{will}|M) \times P(M|V) \\
 & = \frac{3}{4} \times
 \end{aligned}$$

* Named Entity Recognition with CRF

Named Entity - anything that can be referred to with a proper name.

The most common tags are:

PER

LOC

ORG

GPE (Geo-political entity)

- NER aims to find spans of text that constitute proper names, and to tag the type of entity.
- NER can be used for sentiment analysis, question answering & information extraction.

Difficulties in NER

- Segmentation - entities, not words have to be segmented
- Type ambiguity - words with multiple possible entity types

Bio Tagging

- aim to assign one label to each word.

B: token that begins a span

I: tokens inside a span

O: tokens outside of any span

e.g. Jane Villanueva of United Airlines Holding discussed the

↓	↓	↓	↓	↓	↓
B-PER	I-PER	O	B-ORG	I-ORG	J-ORG

Chicago route

↓	↓
B-LOC	O

* Standard Algorithms for NER

1. Hidden Markov Models
2. Conditional Random Field (CRF) / Maximum Entropy Markov Models (MEMM)
3. Neural Sequence Models (RNNs or Transformers)
4. LLMs like BERT

* Linear Chain Conditional Random Fields

→ In a linear chain CRF, tag assignment depends only on the tag of one previous word.

→ The features for CRF are:

1. w_i → ith word of a sentence
2. embeddings → vector representation of word
3. gazetteer → list of places' names with their geographical & political information
4. word shape → notation in which letters of a word are denoted as

(i) small letter = x

(ii) capital letter = X

(iii) digits = d

e.g. Delhi १. १०३ ९. ०० = $X \times x x x^{\circ} / . d d d^{\circ} / . X X$

5. short word shape → similar to word shape, but with a slight change - remove consecutive similar letter types. The eq. would be: $\times x^{\circ} / . d^{\circ} / . X$

→ CRF for a word sequence is:

$$P(y|x) = \frac{\exp(\sum_j w_j F_j(x,y))}{\sum y' \exp(\sum_j w_j F_j(x,y'))}$$

* Syntactic Analysis

Syntax - the way words are arranged together

The main ideas of syntax are:

- (i) constituency
- (ii) grammatical relations
- (iii) subcategorization and dependencies

A. Constituency - group of words may behave as a single unit or a phrase called a constituent.

eg. Rama Krishnan

B. Grammatical Relations - A formalization of ideas from traditional grammar

C. Subcategorization and Dependencies - referring to certain kind of relations between words and phrases

eg. the verb want can be followed by infinitive / noun phrase

* Context-Free Grammar

→ a mathematical system for modelling constituent structure in English

→ A CFG has: (i) a set of rules / productions
(ii) a lexicon of words and symbols

→ The symbols used in a CFG are divided into two classes:

(i) Terminal Symbols - correspond to words ..

(ii) Non-Terminals - symbols that represent clusters / generalizations

→ A set of strings $\{d\}$ is derivable from the start symbol S

→ The formal definition of a CFG is a 4-tuple

- a set of non-terminal symbols N

- a set of terminal symbols Σ

- a set of production rules P , of the form $A \rightarrow d$ where A is a non-terminal

- a start symbol S

Example - Construct a parse tree for the following sentence using the given grammar & Lexicon.

Sentence: I prefer a morning flight

Lexicon: noun → flight | breeze | trip | morning

verb → is | prefer | like | need

adjective → cheapest | non-stop | first | latest

Pronoun → me | I | you | it

proper-noun → Alaska | Baltimore

Determiner → the | a | am | this | these

Preposition → from | to | on | near

Conjunction → and | or | but

Grammar - $S \rightarrow NP VP$

$NP \rightarrow \text{pronoun} \mid \text{proper-noun} \mid \text{det} \mid \text{nominal}$

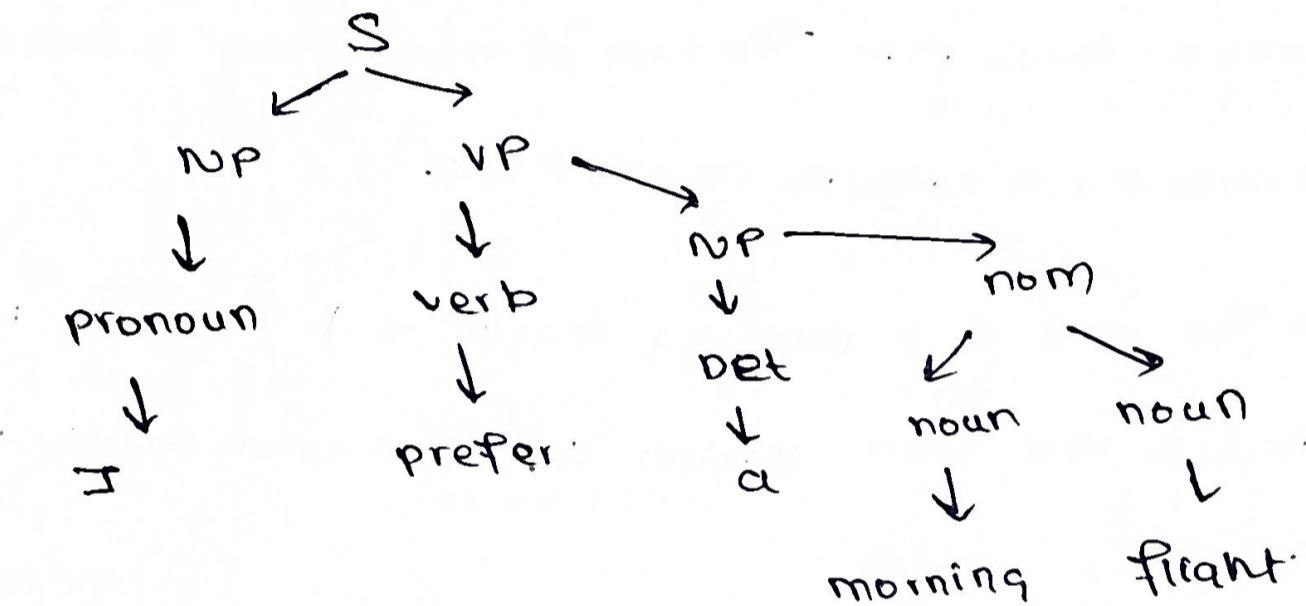
$\text{nominal} \rightarrow \text{noun} \mid \text{nominal } \backslash \text{noun}$

$VP \rightarrow \text{verb} \mid \text{verb } NP \mid \text{verb } NP \text{ PP} \mid \text{verb } PP$

$PP \rightarrow \text{preposition } NP$

Ans

I prefer a morning flight



Example - Draw the parse tree for the sentence "The man
read this book" using the grammar given below.

$S \rightarrow NP VP$

$S \rightarrow \text{Aux } NP VP$

$S \rightarrow VP$

$NP \rightarrow \text{Det Nom}$

$Nom \rightarrow \text{Noun}$

$Nom \rightarrow \text{Noun Nom}$

$VP \rightarrow \text{Verb}$

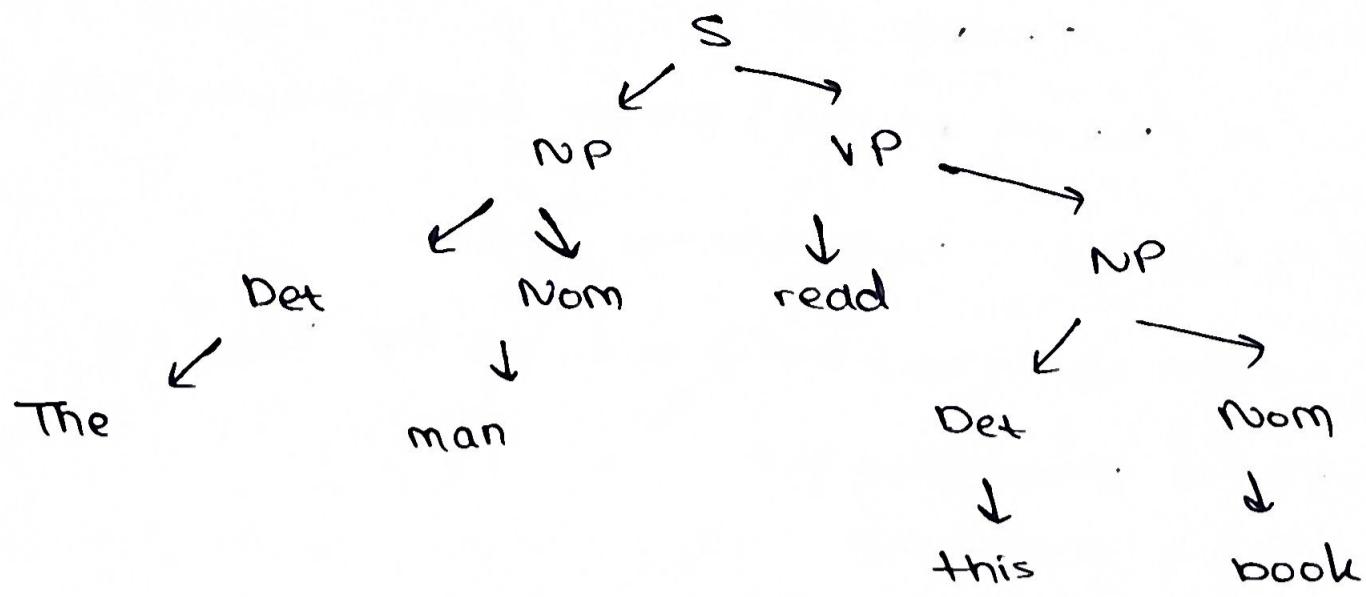
$VP \rightarrow \text{Verb } NP$

$\text{Det} \rightarrow \text{that} \mid \text{this} \mid \text{a} \mid \text{the}$

$\text{Noun} \rightarrow \text{book} \mid \text{flight} \mid \text{meal} \mid \text{man}$

$\text{Verb} \rightarrow \text{book} \mid \text{include} \mid \text{read}$

$\text{Aux} \rightarrow \text{does}$



Syntactic Parsing - The task of recognizing a sentence and assigning a syntactic structure to it.

- The goal of a parsing search is to find all trees whose root is the start symbol S, which cover exactly the words in input.
- There are 2 strategies :
 - (i) Top-down or goal-directed search
 - (ii) bottom-up or data-directed search

* Top-Down Parsing

- Uses a top-down, depth-first, left-to-right approach - expand the left-most unexpanded node in the tree

Example : Create a parse-tree using top-down parsing for "~~Book this flight~~" using the following lexicon and grammar.

Does this flight include a meal?

$S \rightarrow NP VP$

✓

$S \rightarrow Aux NP VP$

$S \rightarrow VP$

$NP \rightarrow Det Nominal$

$Nominal \rightarrow Noun$

$Nominal \rightarrow Noun Nominal$

$NP \rightarrow Proper-Noun$

$VP \rightarrow Verb$

$VP \rightarrow Verb NP$

$Det \rightarrow that / this / a$

25

$Noun \rightarrow book / flight / meal / money$

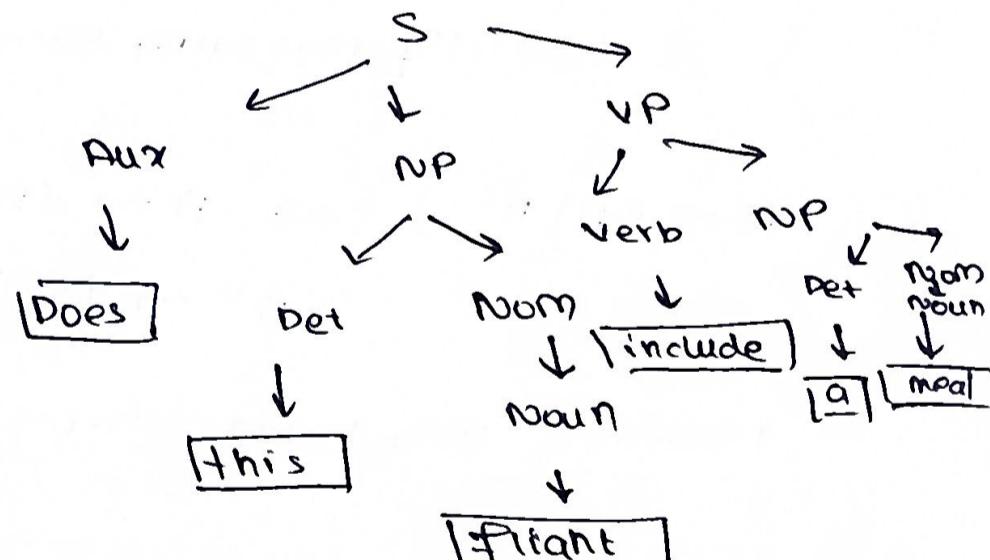
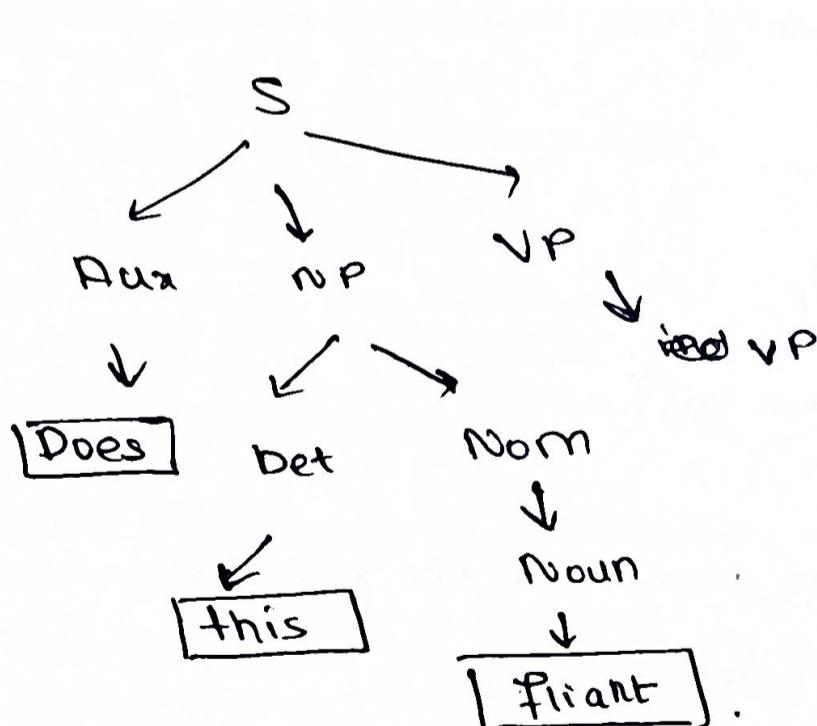
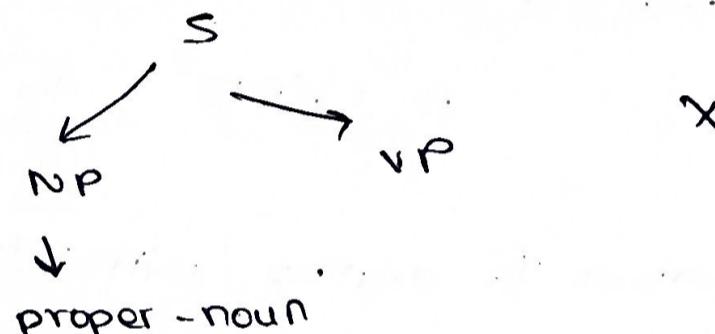
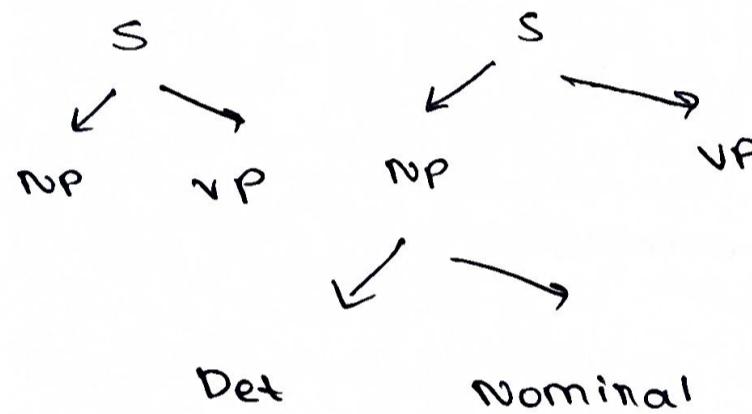
$Verb \rightarrow book / include / prefer$

$Aux \rightarrow does$

$Prep \rightarrow from / to / on$

$Proper-noun \rightarrow Houston / USA$

$Nominal \rightarrow Nominal PP$



Problems with Top-Down Parsing

1. Left Recursion
2. Ambiguity
3. Inefficient reparsing of subtree

A. Left Recursion

→ leads to infinite exploration of search spaces, when left-recursive grammars are used.

→ A grammar is left-recursive if it is of the form:

$$A \Rightarrow^\omega \alpha A \beta \quad \text{for some } \alpha \text{ and } \beta, \text{ and } \alpha \Rightarrow^* \epsilon$$

→ Two methods for dealing with left recursion include:

1. rewriting the grammar

2. explicitly managing the depth of the search during parsing

$$\begin{aligned} A \rightarrow A\beta \mid A' &\Rightarrow A \rightarrow \alpha A' \\ &A' \rightarrow \beta A' \mid \epsilon \end{aligned}$$

rewriting may make interpretation difficult

B. Ambiguity

→ A word may have more than one POS, requires disambiguation to choose the correct ^{pos} word for a word.

→ Structural ambiguity refers to how a grammar assigns more than one possible parse to a sentence.

→ There are 3 kinds of structural ambiguity:
(i) attachment ambiguity
(ii) coordination ambiguity
(iii) non-phrase bracketing ambiguity

(i) Attachment ambiguity - a particular constituent can be attached to the parse tree at more than one place.

(ii) Coordination ambiguity - different sets can be conjoined by conjunctions

old men & women

old [men & women]

[old men] & women

→ choosing the correct parse requires disambiguation, and it requires both statistical and semantic knowledge.

c. Inefficient Reparsing of Subtrees

→ The parser often builds valid trees for portions of the input, then discards them during backtracking, only to find that it has to be rebuilt again.

* Bottom-up Parsing

→ The initial goal list is the string to be parsed.

→ If the sequence in the goal list matches the RHS of a rule, then this sequence may be replaced by the LHS of the ruling.

→ Parsing is finished when the goal list has just the start symbol.

→ The standard method is shift-reduce parsing.

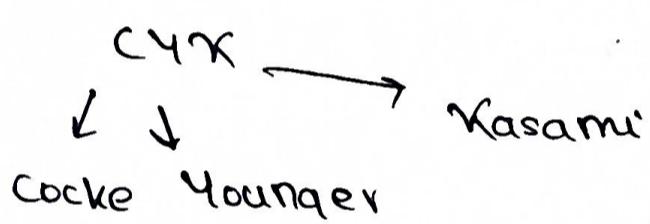
Shift → push the next input symbol from the buffer onto the stack

Reduce pop the Rule's RHS off the stack and replace it w/ the terminal on the LHS

When either shift or reduce is possible, choose arbitrarily

Example 3 Use bottom-up parsing to parse "Book that flight" using the following lexicon & grammar.

* CYK Parsing



→ The structure of rules is in a Chomsky Normal Form grammar

→ A dynamic programming approach can be used.

→ A CFG is in CNF if each rule is of the following form:

- (i) $A \rightarrow BC$
- (ii) $A \rightarrow a$ (terminal)
- (iii) $S \rightarrow \lambda$ (null string)

Ex! Apply the CKX algorithm for the string baaba w/ the following grammar.

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

$$B \rightarrow cc \mid b$$

$$C \rightarrow AB \mid a$$

S, A, C x ₁₅					
-	S, C, A x ₂₅				
-	B x ₃₄	B x ₃₅			
S, A x ₁₂	B x ₂₃	S, C x ₃₄	S, A x ₄₅		
B x ₁₁	A, C x ₂₂	A, C x ₃₃	B x ₄₄	A, C x ₅₅	

① b a a b a

$$x_{11} : b \quad B \rightarrow cc \mid b$$

$$x_{22}, \underset{33}{:} a \quad A \rightarrow BA \mid a, \quad C \rightarrow AB \mid a$$

$$x_{34} : b \quad B \rightarrow cc \mid b$$

$$x_{55} : a \quad A \rightarrow BA \mid a \quad C \rightarrow AB \mid a$$

$$② \quad \textcircled{b} \quad x_{12} : ba = B \quad \text{and} \quad A, C \\ \Downarrow$$

$$BA \quad BC$$

$$A \rightarrow BA \mid a \quad S \rightarrow AB \mid BC$$

$$\Rightarrow S, A$$

(31)

$$x_{23} = aa = A, C \text{ and } AC$$

↓

$$\begin{array}{cccc} AA & AC & CA & CC \\ \emptyset & Q & Q & B \rightarrow CC | b \end{array}$$

$$\Rightarrow \underline{\underline{B}}$$

$$x_{34} = ab = A, C \text{ and } B$$

↓

$$AB \text{ and } CB$$

$$S \rightarrow AB | BC$$

∅

$$C \rightarrow AB | a$$

$$\Rightarrow \underline{\underline{S, C}}$$

$$x_{45} = ba = B \text{ and } A, C$$

$$= BA \text{ and } BC.$$

$$A \rightarrow BA | a$$

$$S \rightarrow AB | BC$$

↓

$$\underline{\underline{S, A}}$$

(3)

$$x_{13} = bao$$

possible combinations are:

$$ba$$

$$a$$

↓

$$S, A$$

and

$$b$$

$$aa$$

↓

$$B$$

$$B$$

↓

$$SA, SC, AA, AC$$

$$\emptyset \emptyset \downarrow \downarrow$$

$$q \emptyset$$

↓

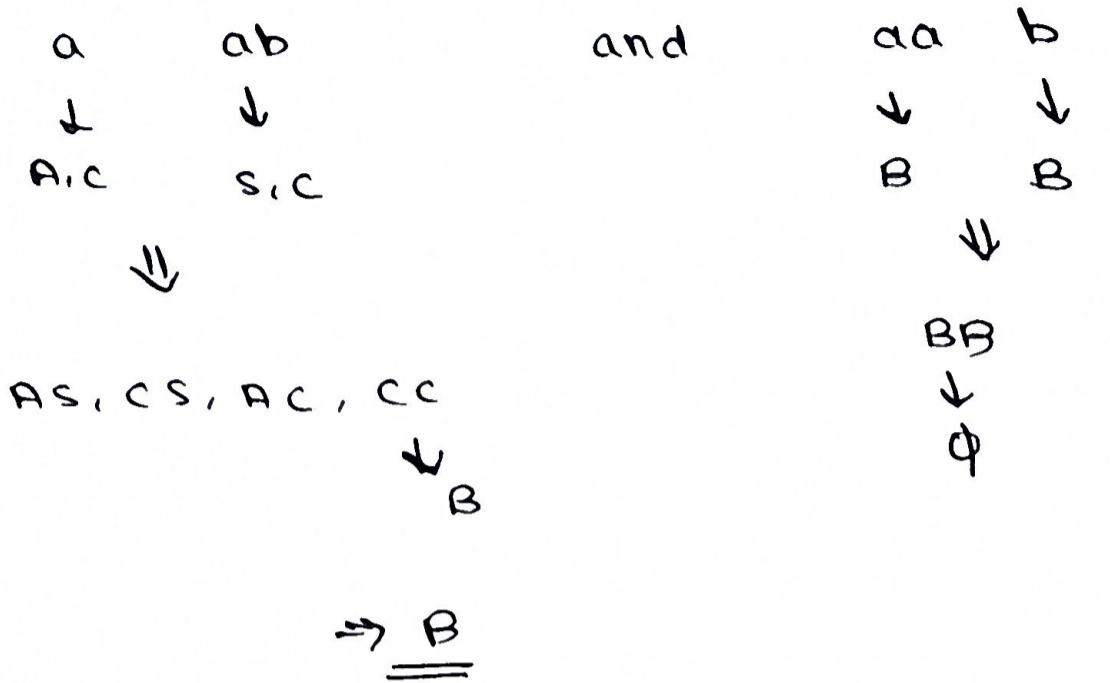
$$BB$$

↓

$$\emptyset$$

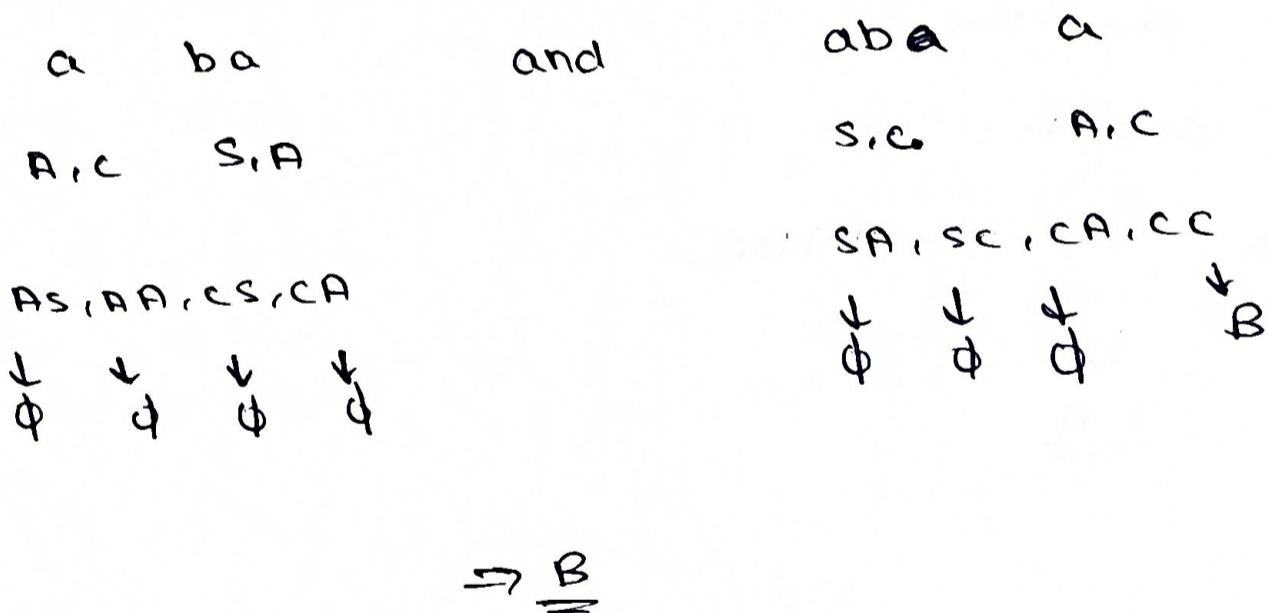
$$X_{24} = aab$$

possible combinations are:



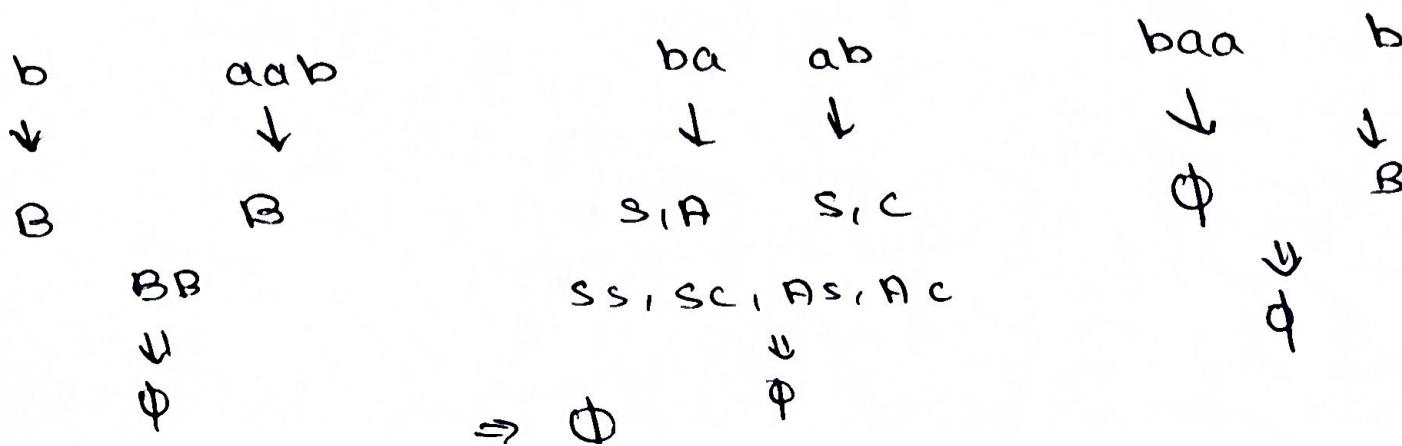
$$X_{35} = aba$$

possible combinations are



④ $X_{14} = baa$

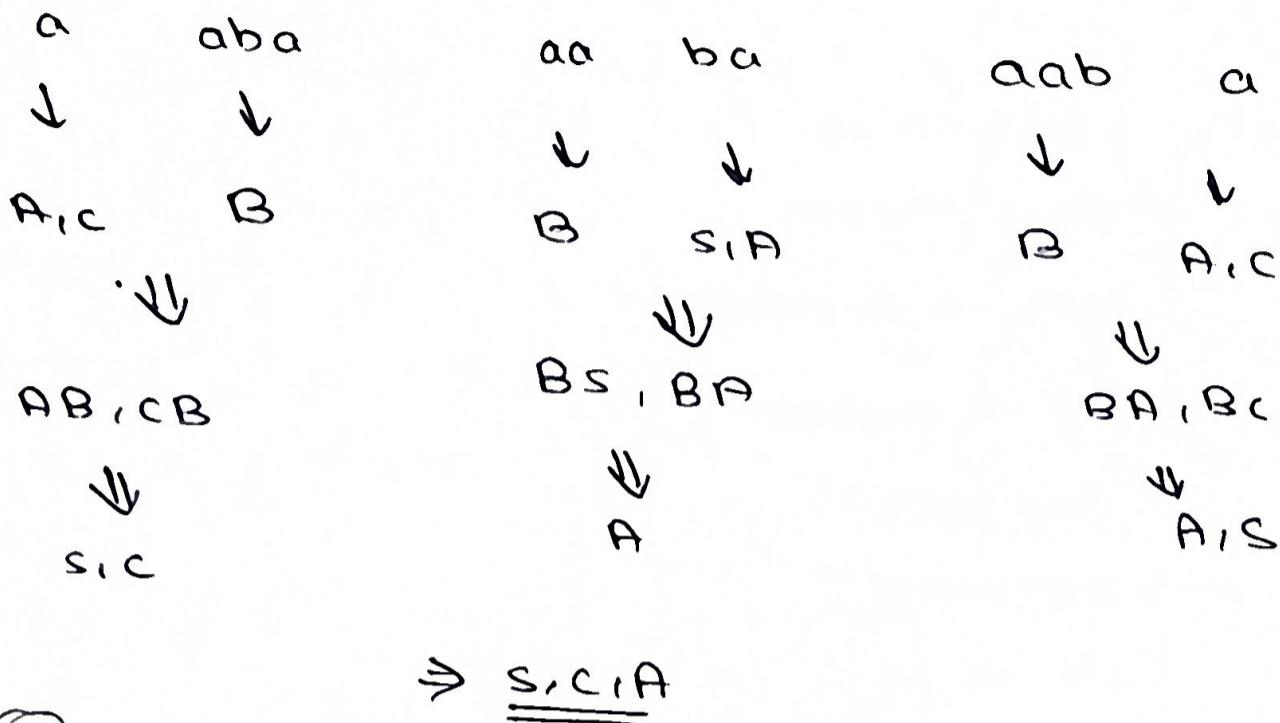
possible combinations are:



$X_{25} = aaba$

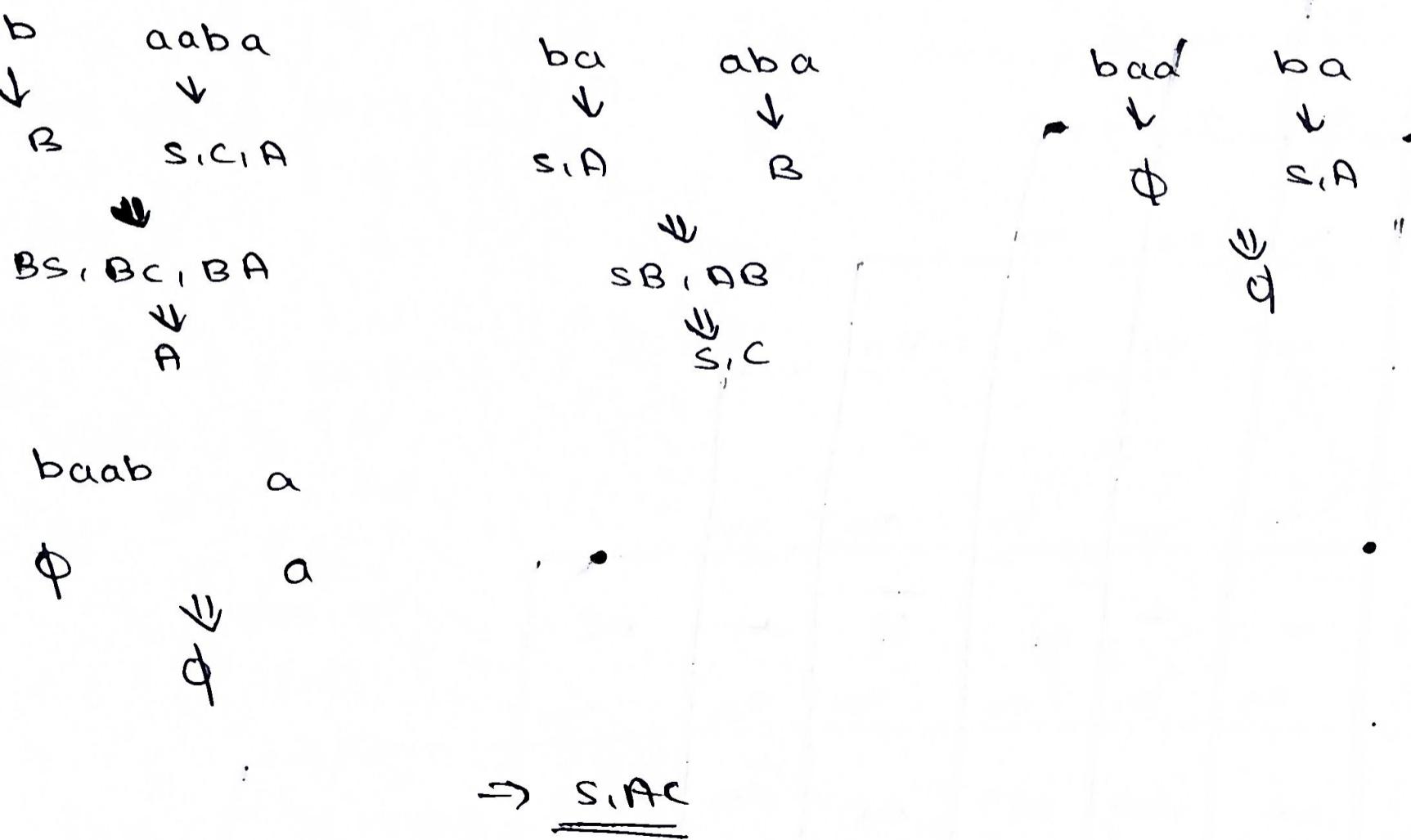
(33)

Possible combinations are:



$X_{15} = baabaa$

Possible combinations are:

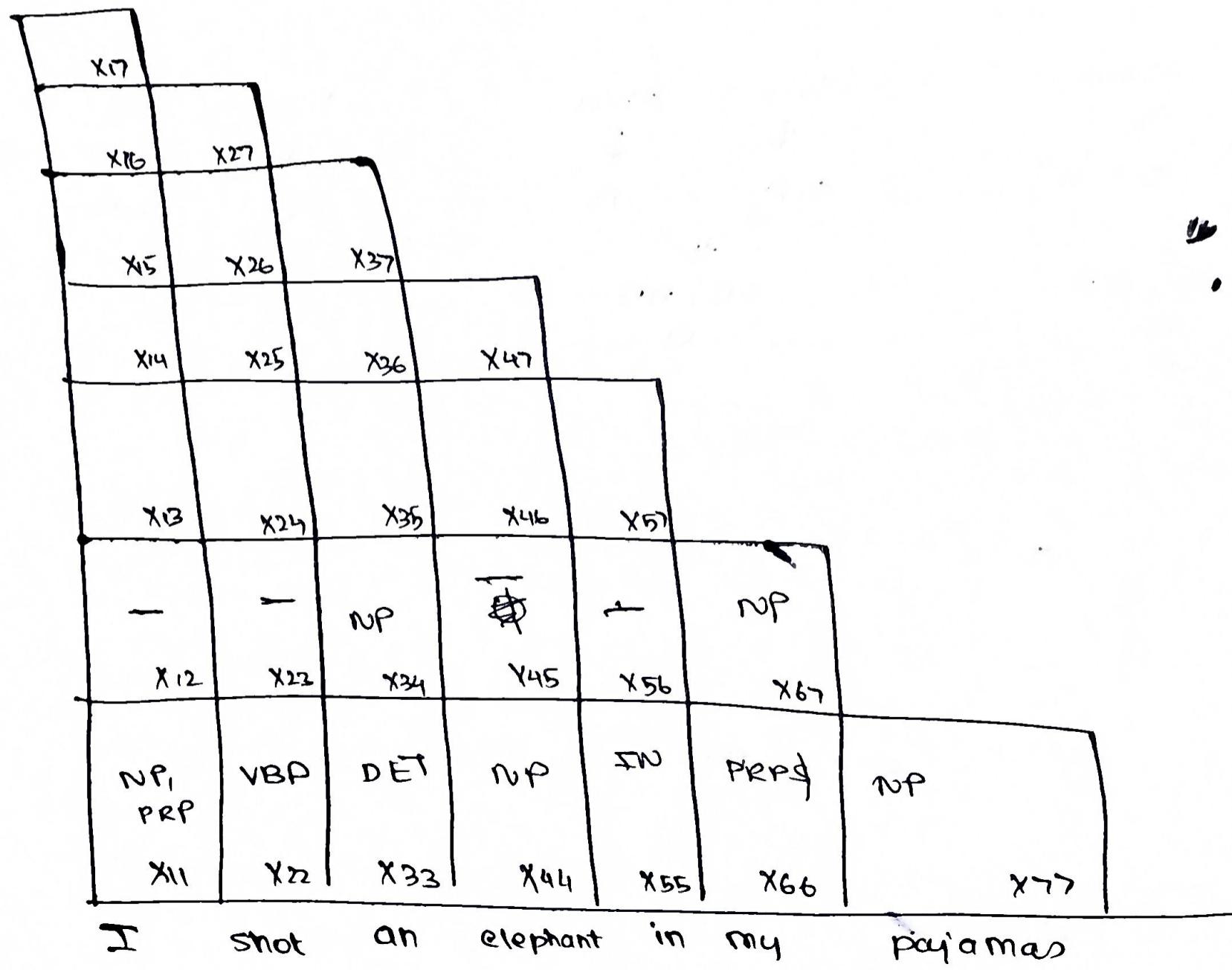


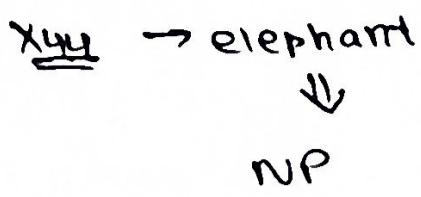
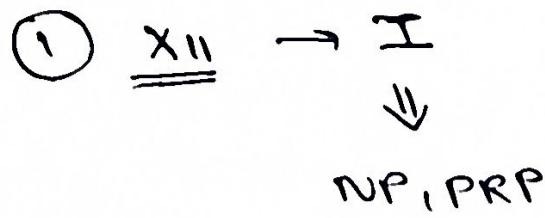
Exa Apply CYK grammar for the sentence:

"I shot an elephant in my pajamas"

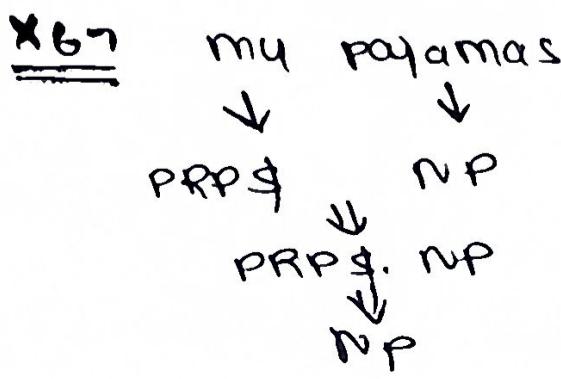
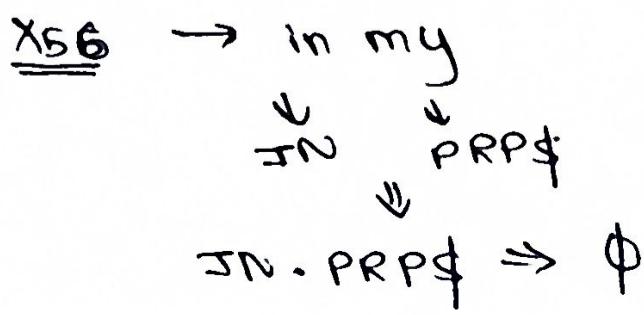
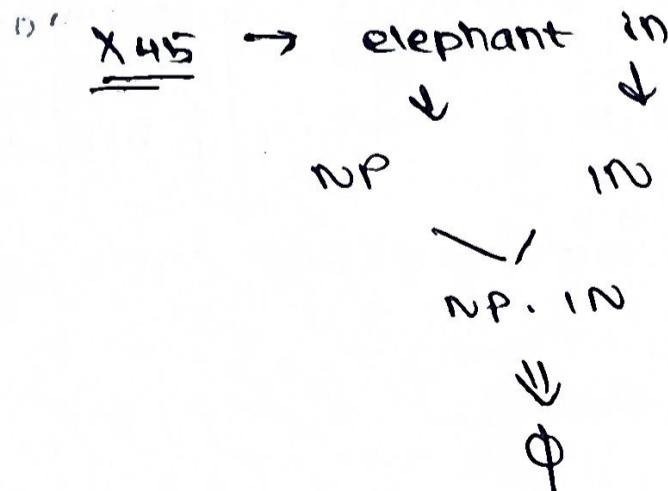
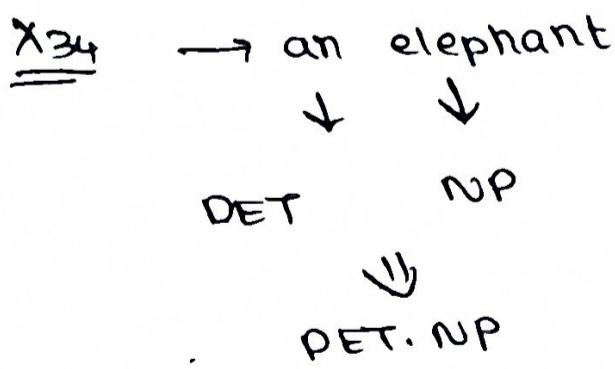
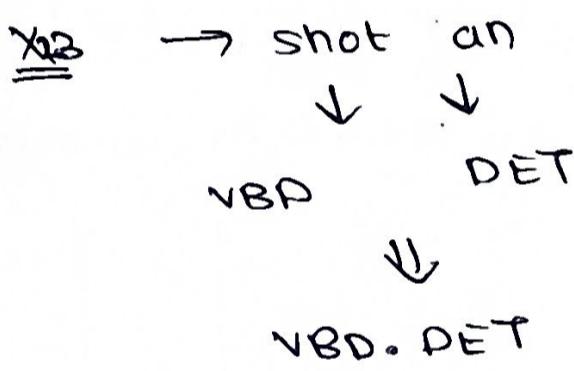
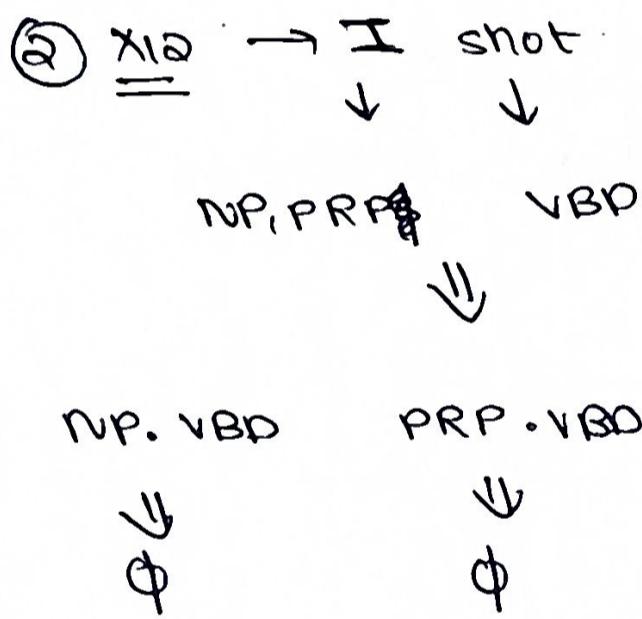
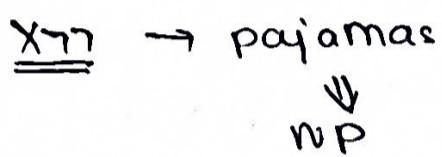
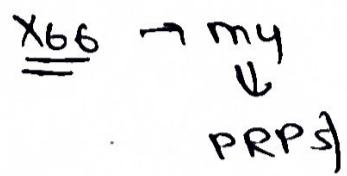
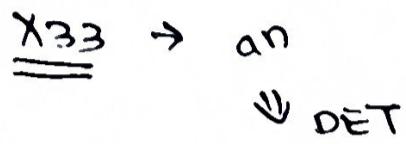
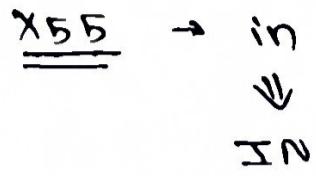
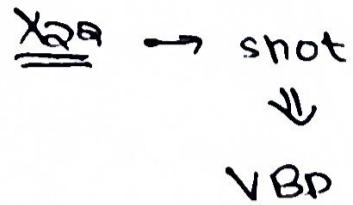
with the grammar rules:

$S \rightarrow NP VP$	DET → an
$PP \rightarrow IN NP$	VBD → shot
$NP \rightarrow DET NP$	NP → pajamas
$NP \rightarrow NP PP$	NP → elephant
$VP \rightarrow VBD NP$	NP → I
$VP \rightarrow VP PP$	PRP → I
$NP \rightarrow PRP\$ NP$	IN → in
	PRP\\$ → my





5



* Shallow Parsing

- a technique that aims to identify and extract meaningful phrases or chunks from a sentence.
- Unlike full parsing, which involves analyzing the grammatical structure of a sentence, shallow parsing focuses on identifying individual phrases like noun phrases, verb phrases etc.

Benefits

- efficient, full parsing is computationally intensive
- beneficial for processing large volumes of text like web crawling, document indexing, and machine translation

Steps

1. divide sentence into words or tokens

2. POS tagging

3. identify and extract relevant phrases / constituents. This is done using pattern matching / ML algorithms that have been trained to recognize particular types of phrases

* Dependency Parsing