

Natural Language Processing

Regular Expression And Automata

Regex	Match	Match
/[Ww]oodchuck/	Woodchuck	Woodchuck or woodchuck
/[0-9]/		a single digit
/[a-z]/		a lower case letter
/[A-Z]/		a upper case letter
/[^A-Z]/		not an upper case letter
/[^Ss]/		neither 'S' nor 's'
/[^.]/		not a period
/[e^]/		either 'e' or '^'
/a^b/		the pattern 'a^b'

Reger
colour?)

oo*h!

o*h!

~~beg.~~
beg.: n SA

^ [A-Z]

max. off seen

\. \$

\b the \b

guppy | ies

guppy (y | ies)

/Column[0-9]+ * /

/(Column[0-9]+ *)*/

Matches

Optional previous char

0 or more of previous
char.

1. or more of previous
char

any character in place
of '.'
E.g. begun, begin

Check whether

[A-Z] is at
beginning of line

Check whether '.'
is present at end
of line

Check whether

'the' is a word

guppy and ies

guppy and guppies

Here '*' denotes
multiple spaces

E.g. column0 column1

Substitution regex

S / colour / color
find replace

S / [0-9]+ / <1> → Replace pattern with pattern surrounded by '`<`' and '`>`'.
E.g 40 is replaced by `<40>`.

Find Regex of

[S-A]^*

1) the Xer they were, the Xer
they will be

E.g the Bigger they were, the bigger
they will be

Find FSA of

1) Eg One cent / forty two cent

2) two, dollar, ten cent

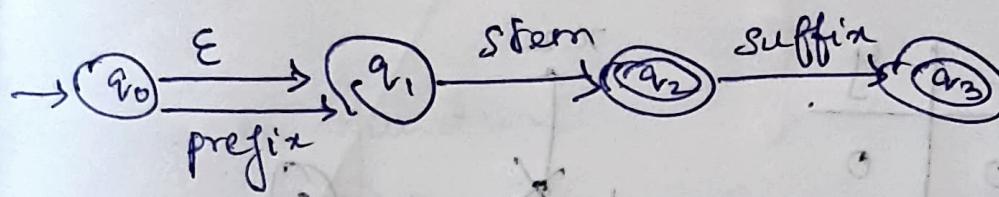
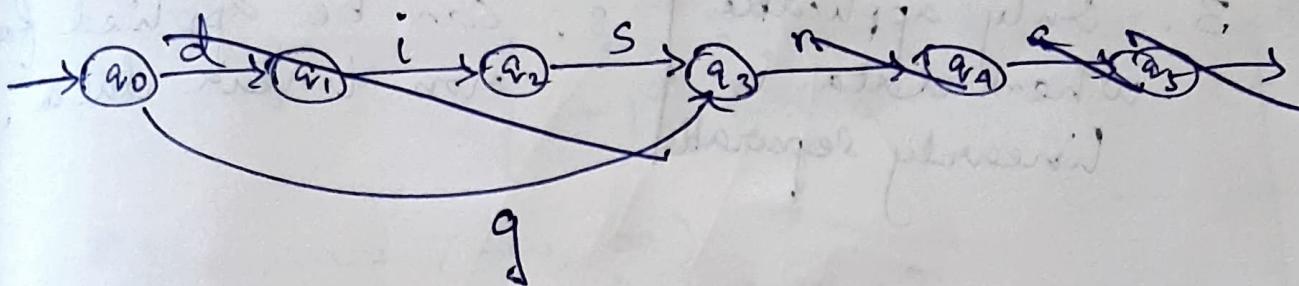
Word classes

Content words

- Nouns
- Verbs
- Adjectives
- Adverbs

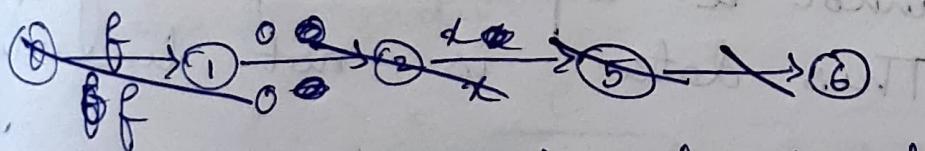
Function words

- Prepositions
- Conjunctions
- Determinants



The above automata accepts grace, graceful, disgrace, disgraceful.

Transducer



- It produces internal structure of a word

Natural Language Processing

Intrinsic evaluation of
language models
- perplexity

If a model assigns high probability to a word, it is not perplexed (~~surprised~~) to see it

$$P.P(W) = P(w_1, w_2, \dots, w_N)^{-1/N}$$

$$P(w_1, w_2, \dots, w_N) = \prod_{i=1}^N P(w_i | w_{i-1}, w_{i-2}, \dots)$$

For bigrams

$$P(w_1, w_2, \dots, w_N) = \prod_{i=1}^N P(w_i | w_{i-1})$$

Test dataset

w = (<s>, This, is, the, first, sentence, ., <\\$>, <s>, This,
is, the, second, one, ., <\\$>)

$$N = 16$$

Unigram

This = 0.04 second = 0.001

is = 0.06 one = 0.0012

the = 0.07

fast = 0.0009

sentence = 0.0008

Re Problem

Roll a die 10 times

$$PP(T) = \frac{1}{\left(\frac{1}{6}\right)^{10}} = 6^{10}$$

- 6 If we have an unfair die that rolls
6 with probability $\frac{1}{12}$ and others
with $\frac{1}{12}$ each

$$PP(T) = \frac{1}{\left(\frac{1}{12}\right)^7 \cdot \left(\frac{1}{12}\right)^3} = 3.9 \approx 4$$

- 1 Laplace smoothing
(also called Add-one estimation)

MLE estimate:

$$P_{MLE}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

Add-1 estimate:

$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1}) + V}$$

Note:

After Laplace smoothing, bigrams with respect to a certain word will have the same probability.

Good Turing Discounting

Calculate Frequency of frequencies for bigram

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

Backoff Interpolation

- use trigram if probability of all trigrams is available, otherwise move to bigrams and unigrams

Relation between perplexity and Cross Entropy

$$PP(w) = 2^{-H(w)}$$

Absolute discount

$$c^* = (c+1) \cdot 0.75$$

Kneser-Ney Smoothing

$$P_{\text{continuation}}(w) \propto |\{w_{i-1} : c(w_{i-1}, w) > 0\}|$$

$$P_{\text{continuation}}(w) = \frac{|\{w_{i-1} : c(w_{i-1}, w) > 0\}|}{\sum_{w'} |\{w_{i-1} : c(w_{i-1}, w') > 0\}|}$$

POS Tagging

Distributional hypothesis

Words that appear in similar contexts have similar representation.

Due to ambiguity (and unknown words), we cannot rely on a dictionary to look up the correct POS tags.

Hidden Markov Model

$$P(y) = P(y_1 \dots y_n)$$

$$P(y_1, \dots, y_n) \approx \prod_{i=1}^{n+1} P(y_i | y_{i-1})$$

$$\begin{aligned} P(x|y) &= P(x_1 \dots x_n | y_1 \dots y_n) \\ &\approx \prod_{i=1}^n P(x_i | y_i). \end{aligned}$$

$$P(x_t | y_t) = \frac{c(x, y)}{c(y)}$$

$$P(y_t | y_{t-1}) = \frac{c(y_{t-1}, y_t)}{c(y_{t-1})}$$

Natural Language Processing

- 3 mostly used tagsets
 - 1) Penn Tree Bank Tagset (small size - 45 Tagsets)
 - 2) CT tagset (medium size - 67 tagsets)
to tag British Nation Corporations
 - 3) CT tagset (Large size - 146 tagsets)

Penn Tree Bank POS Tags

Tag	Description	Example
CC	Conjunction	and, but
DT	Determiner	a, the, that
IN	Preposition	of, in, by
JJ	Adjective	
JJR	Adjective comparative	
JJS	Adjective superlative	
NN	Noun - Singular	
NNS	Noun - Plural	
NNP	Proper noun - singular	
NNPS	Proper noun - Plural	

Tag	Description
PP	Personal pronoun
PP\$	Possessive pronoun
RB	Adverb
RBR	Adverb - comparative
RBS	Adverb - superlative
VB	Verb - base form
VBD	Verb - past tense
VBG	Verb gerund
VBN	Verb past participle
VBZ	Verb 3 sg Pres
WDT	wh-determiner
WP	wh-pronoun
WRB	Wh-adverb
	Left Parenthesis [, C, {
E.g.	Book that flight ↑ ↑ ↑ VB DT NN

Assign POS Tags for the following sentence based on Penn Tree Bank Tagset.

- 1) Does that flight serve dinner?
- 2) This crap game is over a garage in fifty second street
- 3) I need a flight from India
- 4) I have a friend driving to Australia

5) What flight do you have from India to America?

6) Can you ~~go~~ list the non stop afternoon flights to England?

B) Answers

Does that flight serve dinner

POS tagging algorithms

- 2 classes

1) Rule based POS tagging

2) Stochastic tagging

Rule based

- suitable for large databases

- handwritten unambiguous rules for each word

E.g. ~~not~~ specify a word as a noun rather than an adjective if it follows a determiner

Stochastic

- Probability is assigned to each word with a given tag defined in the given context

- Hidden Markov Model (HMM) tagger / Maximum likelihood tagger

Natural Language Processing

Hidden Markov Model

$$\hat{t}_i^n = \operatorname{argmax}_{t_i^n} P(t_i^n | w_i^n)$$

↓

Word sequence

tag sequence

$$= \operatorname{argmax}_{t_i^n} P(w_i^n | t_i^n) \cdot P(t_i^n)$$

$$P(t_i^n) = \prod_{i=2}^n P(t_i | t_{i-1})$$

→ Prior probability
transition probability

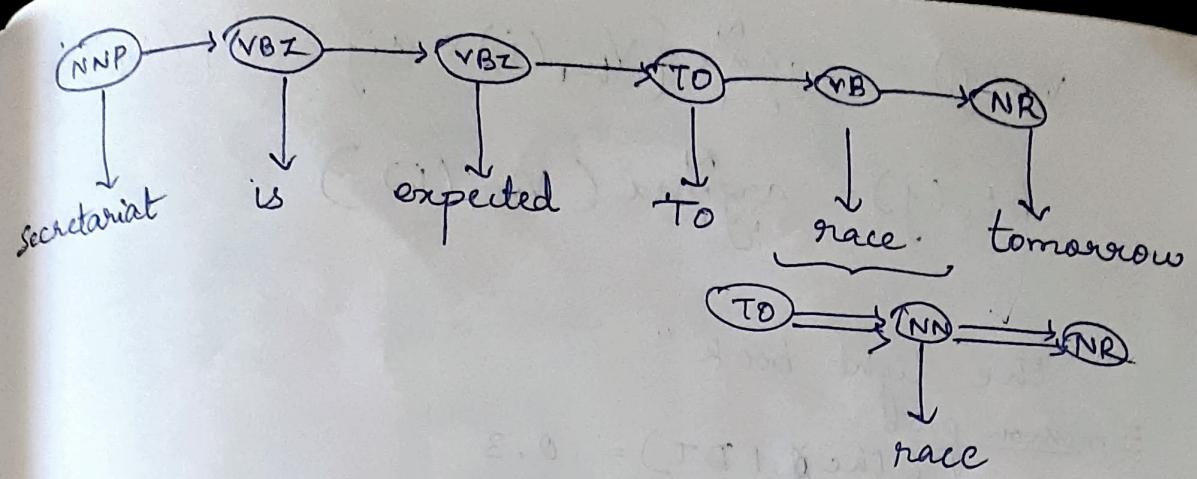
$$P(w_i^n | t_i^n) = \prod_{i=1}^n P(w_i | t_i)$$

also called

Observation likelihood / Emission probability

Problem

Secretariat / NNP is / BEZ expected / VBN
 to / TO race / VB tomorrow / NR



Given ~~'Today is hot'~~ * ~~'tomorrow will be io'~~

Given \rightarrow today = hot

tomorrow = cold

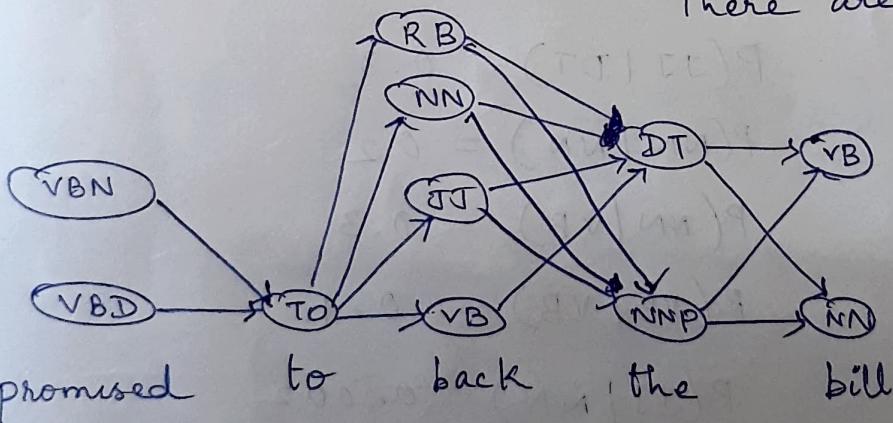
day after = warm

$$P(\text{warm}, \text{cold} | \text{hot}) = P(q_{n+1} | q_n) * P(q_{n+2} | q_{n+1}, q_n)$$

$$P(q_{n+2} | q_{n+1}, q_n)$$

Viterbi Algo

There are 32 paths



Step/
Time 1 2 3 4 5

$$V_{t+1}(K) = \max_j \{ V_t(j) * P(\cancel{w_t} | \cancel{s_j}) * P(w_{t+1} | \cancel{s_j}) \}$$

$$b_t(j) = \underset{\text{argmax}}{\{ V_t(j) \}}$$

$$v_t(j) = \max_{i=1}^N \{ v_{t-1}(i) * a_{ij} * b_j(o_t) \}$$

$$b_{t+}(j) = \arg \max_{i,j} (v_{t-1}(i))$$

"the light book"

Emission prob: $P(\text{the} | DT) = 0.3$

$$P(\text{light} | JJ) = 0.002$$

$$P(\text{book} | verb) = 0.01$$

$$P(\text{the} | NN) = 0.1$$

$$P(\text{light} | VB) = 0.06$$

$$P(\text{book} | NN) = 0.003$$

$$P(\text{light} | NN) = 0.003$$

Transition prob

$$P(VB | DT) = 0.00001$$

$$P(NN | DT) = 0.5$$

$$P(JJ | DT) = 0.3$$

$$P(NN | NN) = 0.2$$

$$P(NN | VB) = 0.3$$

$$P(VB | VB) = 0.1$$

$$P(JJ | NN) = 0.002$$

$$P(VB | NN) = 0.3$$

$$P(NN | JJ) = 0.2$$

$$P(VB | JJ) = 0.001$$

~~VB~~

Note:
 a_{ij} - transition probability
 $b_j(O_t)$ emission probability

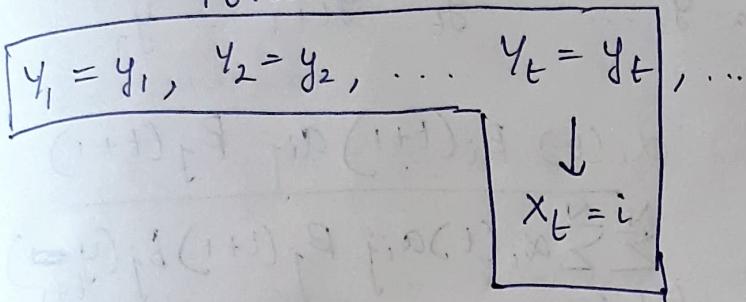
Baum-Welch Algo

$x_t \rightarrow$ hidden state $X = (x_1 = x_1, x_2 = x_2, \dots)$

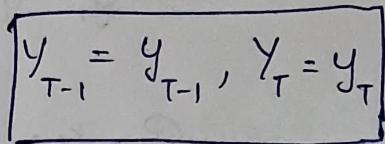
$y_t \rightarrow$ observation $Y = (y_1 = y_1, y_2 = y_2, \dots)$

$$\theta = (\pi, A, B)$$

Forward



Backward



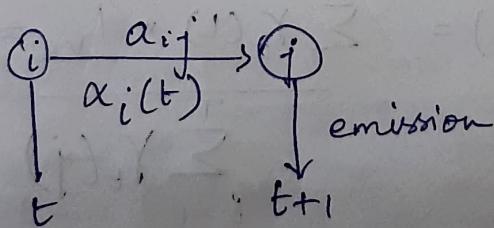
$$\alpha_i(t) = P(y_1, y_2, \dots, y_t, x_t = i | \theta)$$

$$B_i(t) = P(y_{t+1}, y_{t+2}, \dots, y_T | x_t = i, \theta)$$

Forward probabilities

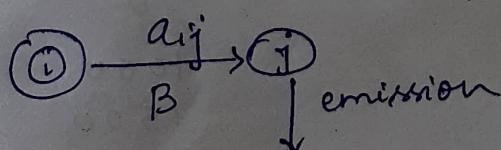
$$\alpha_i(1) = \pi_i * b_i(y_1)$$

$$\alpha_j(t+1) = \sum_{i=1}^N \alpha_i(t) a_{ij} b_j(y_{t+1})$$



Backward probabilities

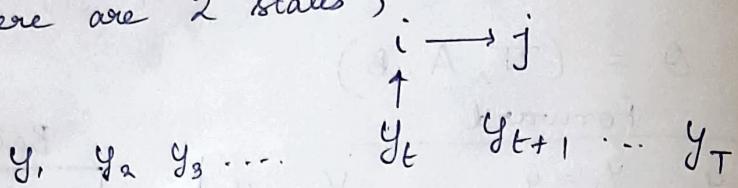
$$B_i(t) = \sum_{j=1}^N B_j(t+1) a_{ij} b_j(y_{t+1})$$



Path probabilities

$$\gamma_i(t) = \frac{\alpha_i(t) \cdot \beta_i(t)}{\sum_{j=1}^N \alpha_j(t) \beta_j(t)}$$

If there are 2 states,



$$S_{ij}(t) = \frac{\alpha_i(t) \beta_i(t+1) \alpha_{ij} \beta_j(t+1)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_i(t) \alpha_{ij} \beta_j(t+1) \beta_j(y_{t+1})}$$

$$\pi = \gamma_i(1)$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T S_{ij}(t)}{\sum_{i=1}^N \left(\sum_{j=1}^N S_{ij}(t) \right)}$$

(Updated transition probability)

$$b_j(v_k) = \frac{\sum_{t=1}^T \gamma_t(1) * I_{v_k=y_t}}{\sum_{t=1}^T \gamma_t(j)} = \begin{cases} 1 & \text{if } v_k=y_t \\ 0 & \text{otherwise} \end{cases}$$

Linear Regression

- * Predicting house prices using description
 # Vague adjectives

Amount
4

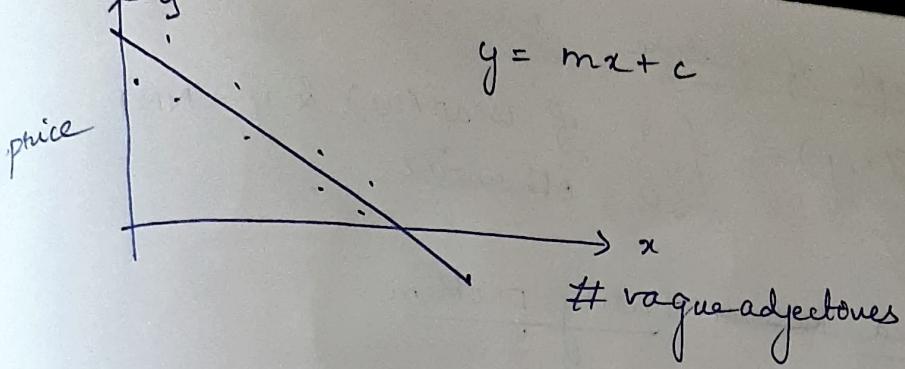
3	1000
---	------

2	1500
---	------

2	6000
---	------

1	14000
---	-------

0	18000
---	-------



Instead of m and c , we use
weights,

$$\text{price} = w_0 + w_1 (\# \text{vague adjectives})$$

In general

$$\hat{y} = \sum_{i=1}^n w_i f_i \quad [f_i - i^{\text{th}} \text{ feature}]$$

We have to minimize cost function.

$$\text{cost}(\omega) = \sum_{j=0}^m (\hat{y}_{\text{pred}} - Y_{\text{actual}})^2$$

We can obtain weights using gradient descent or directly using the below formula

$$w = (X^T X)^{-1} X^T y$$

Issues with HMM

(i) Unknown

Solution: Morphological cues

(ii) Limited context

Solution: increase the window size up to which probability is calculated

$$P_\lambda(y|x) = \frac{1}{Z_\lambda(x)} \exp \left(\sum_i \lambda_i f_i(x, y) \right)$$

Annotations for the equation:

- Normalizing: points to the denominator $Z_\lambda(x)$.
- Observation label: points to the variable y .
- ith feature: points to the term $f_i(x, y)$.
- weight for ith feature: points to the term λ_i .

Example of feature

$$f(x, y) = \begin{cases} 1 & \text{if } \text{uscap}(w) \text{ & } y = \text{NNP} \\ 0 & \text{otherwise} \end{cases}$$

Named ~~Intent~~ Entity problem

	Location	Drug	Person
	f_1	f_2	f_3
1) in Arcadia	1	0	0
2) in Quebec	1	1	0
3) taking Lantac			
4) saw Sue			

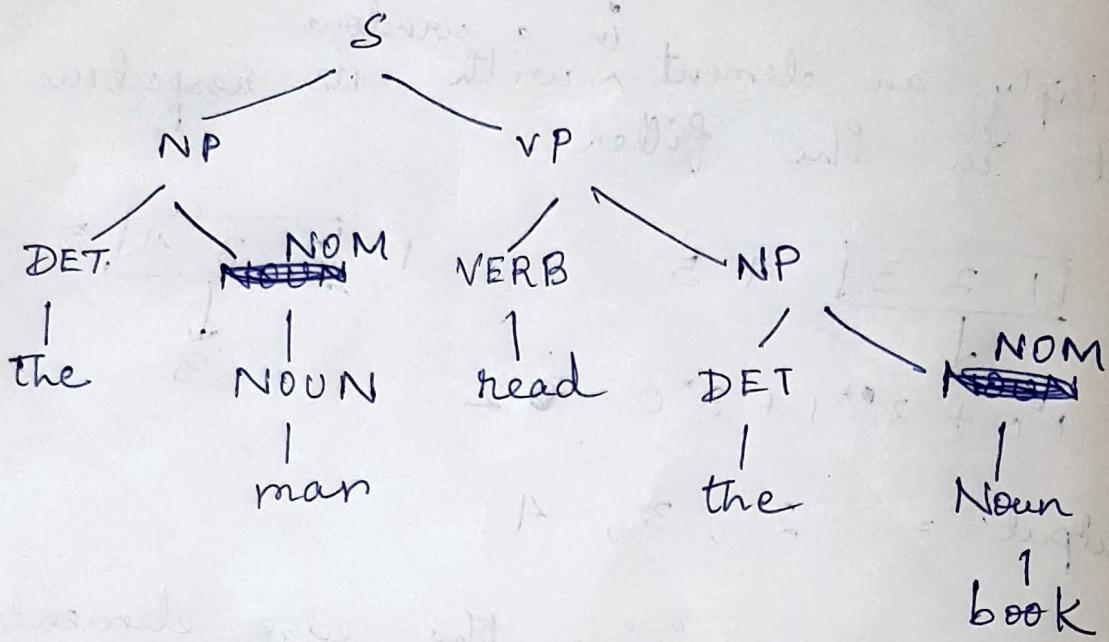
We construct 3 features

$$f_1(x, y) = [y = \text{location} \wedge w_{-1} = \text{"in"} \wedge \text{uscap}(w)]$$

$$f_2(x, y) = [y = \text{location} \wedge \text{has_central_latin_char}(w)]$$

$$f_3(x, y) = [y = \text{Drug} \wedge \text{end}(w, c)]$$

Natural Language Processing



NP (Noun Phrases)

{ Kermit the frog } comes on stage

{ They } come to Massachusetts every summer

{ December twenty sixth } comes after Christmas

PP = On dec. twenty - sixth

{ On dec 26th } I'd like to fly to Florida

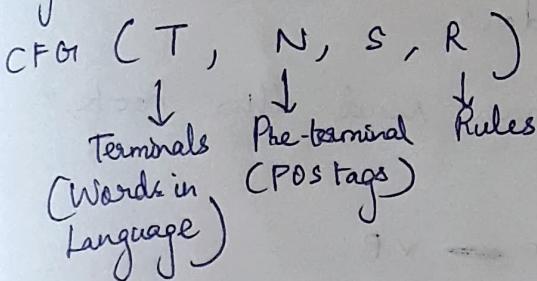
I'd like to ~~to~~ fly to Florida { on dec 26th }

The group of words ~~which~~ ^{that} ~~occur~~ ^{occur} commonly enough ~~that~~ ^{they} can be interchanged with one another is called Constituent. They always form meaningful sentences.

NP → proper noun | ! Det. NOM.

NOM → noun | noun NOM

Using Context Free Grammar we can construct a syntax tree from a given sentence.



Each rule will be of the form

$$X \rightarrow Y \quad \textcircled{X}$$

$$Y \in N^*$$

E.g

NP → Det NOM

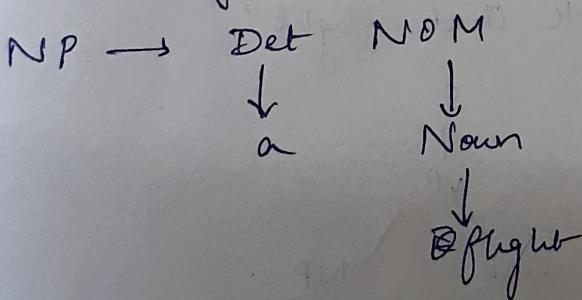
NP → Proper Noun

NOM → Noun | Noun Nominal

Det → a | the

Noun → flight

Derive "a flight"



Example 2

S → NP VP

S → AUX NP VP

S → VP

NP → Det NOM

NOM → Noun

NOM → Noun NOM

NP → verb

VP → Verb NP

Det → that | this | a | the

Noun → book | flight | meal | man

Verb → book | read

AUX → does

Derive "The man read this book"

~~S~~ → NP VP

Det Nom Verb NP
the to Noun
↓ read Det N OM
man . read this N OM
book

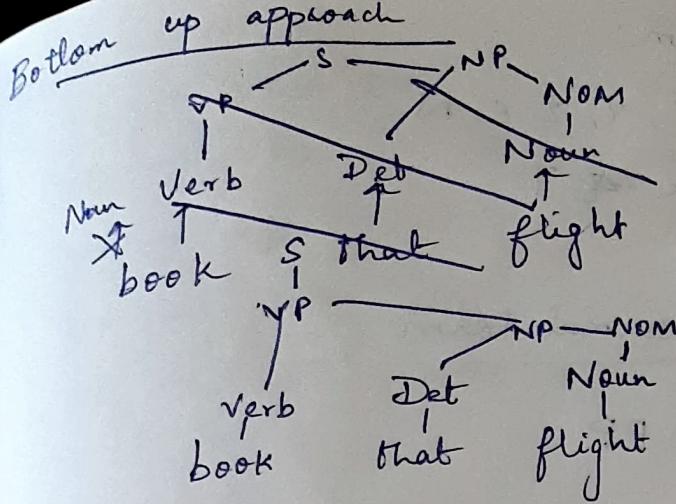
~~Adding new rules~~

VP → VP PP

PP → prep NP

Derive "Book that flight"

S.
VP
Verb / NP
Book Det / N OM
that N OM
flight

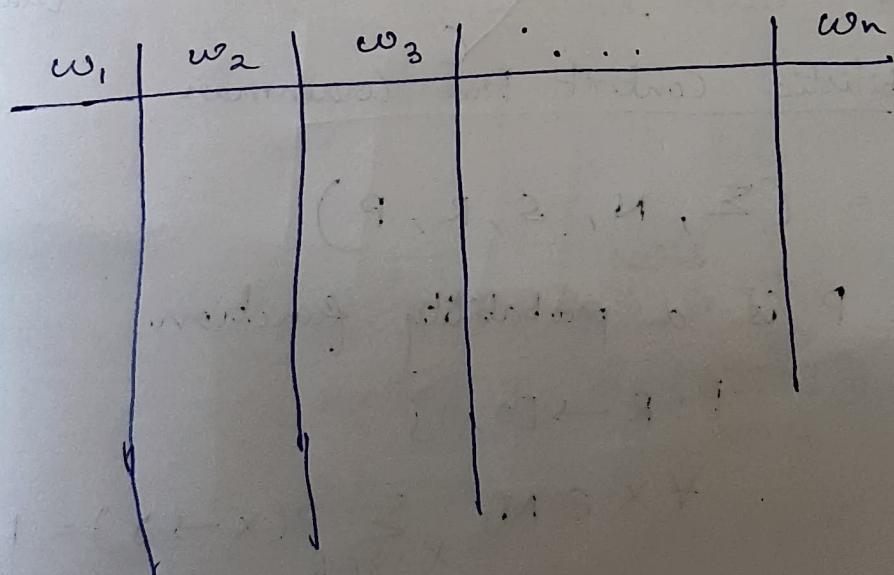


Since lot of searches are made to find which rule is correct, we make use of Dynamic Programming.

CKY algorithm

- uses dynamic programming
- requires grammar to be present in Chomsky Normal form

$$n = \# \text{ words}$$



Note:

Group of words which behaves as a single unit.

Chomsky Normal Form

$x \rightarrow yz, x, y, z \in N.$

$x \rightarrow x, x \notin T$

Convert to Chomsky Normal Form

~~or~~ $VP \rightarrow VNP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$VP \rightarrow VPP | \cancel{VNPP}$

~~or~~ $VP \rightarrow VNPP$

or $VP \rightarrow VNP NP$

or $VP \rightarrow \cancel{(people | fish | tanks)} VNPP | VNPN$

$VP \rightarrow VNP$

or $VP \rightarrow (people | fish | tanks) (people | fish | tanks | rods)$

Probabilistic Context Free Grammar

$$G_1 = (\Sigma, N, S, R, P)$$

P is a probability function

$$P: R \rightarrow [0, 1]$$

$$\forall x \in N, \sum_{\gamma \in T^*} P(x \rightarrow \gamma) = 1$$

A grammar G_1 generates a language model L .

$$\prod_{\gamma \in L^*} P(\gamma) = 1$$

$S \rightarrow NP VP$
 $VP \rightarrow VBG NNS$
 $VP \rightarrow VBZ VP$
 $VP \rightarrow VBZ NP$
 $NP \rightarrow DT NN$
 $NP \rightarrow JJ NNS$

$DT \rightarrow a$
 $NN \rightarrow pilot$
 $VBZ \rightarrow \text{likes}$
 $VBG \rightarrow \text{flying}$
 $JJ \rightarrow \text{flying}$
 $NNS \rightarrow \text{planes}$

a	pilot	likes	flying	planes
DT	NP	-	-	$S_1 \rightarrow NP VP \rightarrow NP VP_1$ $S_2 \rightarrow NP VP_2$
NN	-	-	-	-
VBZ	-	-	-	VP_1 VP_2
		VBG JJ	VP NP	
				NNS

Note :

$dp[i:j]$ - Non terminals that derive words from i^{th} to j^{th} place

$$dp[0:j] = (dp[0:j], dp[1:j])$$

$$(dp[0:j], dp[1:j])$$

DT NNL

NP VBZ

$$dp[2:j] = (dp[2:j], dp[3:j])$$

$$(dp[2:j], dp[3:j])$$

VBZ

dp[2:j], dp[3:j]

Problem

$$S \rightarrow NP VP \quad 0.9$$

$$S \rightarrow VP \quad 0.1$$

$$VP \rightarrow V NP \quad 0.5$$

$$VP \rightarrow V \quad 0.1$$

$$VP \rightarrow V @VP-N' \quad 0.3$$

$$VP \rightarrow V PP \quad 0.1$$

$$VP-N' \rightarrow NP PP \quad 1.0$$

$$NP \rightarrow NP NP \quad 0.1$$

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

$$NP \rightarrow N \quad 0.7$$

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

$$N \rightarrow people \quad 0.5$$

$$N \rightarrow fish \quad 0.2$$

$$N \rightarrow tanks \quad 0.2$$

$$N \rightarrow rods \quad 0.1$$

$$V \rightarrow people \quad 0.1$$

$$V \rightarrow fish \quad 0.6$$

$$V \rightarrow tanks \quad 0.3$$

$$P \rightarrow with \quad 1.0$$