

# NLP ETE

Sol 1:- (a) CKY Parsing Table:-

eat sushi with chopsticks with you

	1	2	3	4	5	6	7
0	NP	.	S		S	.	S
1		V <sub>nb</sub>	VP		VP	.	VP
2			NP		NP	.	NP
3				Prep	PP	.	PP
4					NP	.	NP
5						Prep	PP
6							NP

Sol 1 (b) CNF:-

$$\begin{cases} A \rightarrow \epsilon \\ A \rightarrow AA \\ A \rightarrow a \end{cases}$$

Given,  $S' \rightarrow S$  # add new production  
 $S \rightarrow ASB$   
 $A \rightarrow aASA | a | \epsilon$   
 $B \rightarrow Sbs | A | bb$

(1) Remove  $\epsilon$  productions

$$S' \rightarrow S$$

$$S \rightarrow ASB | SB | AS | S$$

$$A \rightarrow aASA | a | aS | aAS | aSA$$

$$B \rightarrow Sbs | bb | A$$

(2) Remove unit products:

$$S' \rightarrow ASB \mid SB \mid AS$$

$$S \rightarrow ASB \mid SB \mid AS$$

$$A \rightarrow aASA \mid a \mid aS \mid aAS \mid aSA$$

$$B \rightarrow sbs \mid bb \mid aASA \mid a \mid aS \mid aAS \mid aSA$$

(3) Remove useless productions.

$$S' \rightarrow XB \mid SB \mid AS$$

$$S \rightarrow XB \mid SB \mid AS$$

$$A \rightarrow ZA \mid a \mid yS \mid xX \mid UA$$

$$B \rightarrow mS \mid TT \mid ZA \mid a \mid yS \mid xX \mid UA$$

$$X \rightarrow AS$$

$$y \rightarrow a$$

$$Z \rightarrow yX$$

$$U \rightarrow yS$$

$$T \rightarrow b$$

$$m \rightarrow ST$$

Sol 3(a):- Mean Reciprocal Rank Method:

Q No:-	1	2	3	4	5	6	7
Correct Ans							
No :	4	3	6	2	8	1	2

$$\text{Mean Reciprocal Rank} = \frac{1}{101} \sum_{i=1}^{101} \frac{1}{\text{rank}_i}$$

$$\text{MRR} = \frac{1}{7} \left[ \frac{1}{4} + \frac{1}{3} + \frac{1}{6} + \frac{1}{2} + \frac{1}{8} + \frac{1}{1} + \frac{1}{2} \right] \Rightarrow \frac{1}{7} \left[ \frac{69}{24} \right] = \frac{23}{56}$$

$\Rightarrow \frac{23}{56}$  Ans

Sol 3(b) :- Rouge-2 score:

Human

Water Spinach ✓  
Spinach is ✓  
is a ✓  
a commonly  
commonly eaten ✓  
eaten leaf  
leaf vegetable ✓  
vegetable of  
of Asia ✓

System

Water spinach ✓  
Spinach is ✓  
is a ✓  
a leaf  
leaf vegetable ✓  
vegetable commonly  
commonly eaten ✓  
eaten in  
in tropical  
tropical areas  
areas of  
of Asia ✓

# (Human Bigrams) = 9

# (Common bigrams) = 6

# (System Bigrams) = 12

$$\text{Recall} = \frac{6}{9}$$

$$\text{Precision} = \frac{6}{12}$$

$$F_1 \text{ score} = \frac{2PR}{R+P} = \frac{2 \times \frac{6}{12} \times \frac{6}{9}}{\frac{6}{9} + \frac{6}{12}}$$

Rouge 1 :- Pr Re F1

10/13

10/10

$$= \frac{2 \times \frac{10}{13} \times \frac{10}{10}}{1 + \frac{10}{13}}$$

$$\Rightarrow \frac{20}{23}$$

$$= \frac{4}{7}$$



# Sol 4:- Cosine Similarity.

	Doc1	Doc2	Doc3	Doc4	Q
new	1	0	0	1	0
home	1	1	1	1	1
sales	1	1	1	1	1
top	1	0	0	0	0
forecasts	1	0	0	0	0
rise	0	1	0	1	0
in	0	1	2	0	0
july	0	1	1	1	0
increase	0	0	1	0	0
is	0	0	0	0	1
very	0	0	0	0	1
bad	0	0	0	0	1

## Similarity Matrix

	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>	D <sub>4</sub>	Q
D <sub>1</sub>	1	0.4	0.316	0.6	0.4
D <sub>2</sub>		1	0.79	0.8	0.4
D <sub>3</sub>			1	0.474	0.316
D <sub>4</sub>				1	0.4
Q					1

Using Cosine Similarity =  $\frac{D_1 \cdot D_2}{|D_1| |D_2|}$

⇒ formula, fill similarity matrix.

\*  $\lambda = 0.3$   
 $MMR = \arg \max [\lambda \text{Sim}_1(D_i, Q) - (1-\lambda) \text{Sim}_2(D_i, D_j)]$

First iteration:-

$S$  is empty set.

$\therefore$  max pairwise similarity ~~between~~ within  $S = 0$

So,

$MMR = \arg \max (\text{Sim}(d_i, q))$

$\therefore$   $d_1$  has  <sup>$d_2, d_4$</sup>  max similarity, pick any one

$S = \{d_1\}$

Second iteration:-

find max distance of an element in  $S$   
 to given  $d_i$ ,

$\text{sim}(d_1, d_i)$

for  $d_2$ ,

$\text{sim}(d_1, d_2) = 0.4$

$\text{sim}(d_2, q) = 0.4$

$MMR = 0.4 \times 0.3 - 0.6 \times 0.3$

$= 0.12 - 0.18 = \underline{-0.06}$

Similarly for  $d_3, d_4$ ,

for  $d_3$ ,

$$\text{Sim}(d_1, d_3) = 0.316$$

$$\text{Sim}(d_3, q) = 0.316$$

$$\begin{aligned} \text{MMR} &= 0.3 \times 0.316 - 0.7 \times 0.316 \\ &= -0.1264 \end{aligned}$$

for  $d_4$ ,

$$\text{Sim}(d_1, d_4) = 0.6$$

$$\text{Sim}(d_4, q) = 0.4$$

$$\begin{aligned} \text{MMR} &= 0.3 \times 0.4 - 0.7 \times 0.6 \\ &= -0.3 \end{aligned}$$

Here, for  $d_3$ ,  $\text{Max MMR}$  is there,

$$\text{So, } S = \{d_1, d_2\}.$$

Third Iteration :-

$$\text{find } \max \begin{cases} \text{Sim}(d_i, d_i) & \text{for 1st part} \\ \text{Sim}(d_i, d_2) & \text{for 2nd part} \end{cases}$$

$$\max \left( \frac{\text{Sim}(d_1, d_3)}{0.316}, \text{Sim}(d_2, d_3) \right) = 0.79$$

$$\text{Sim}(d_3, q) = 0.316$$

$$\begin{aligned} \text{MMR} &= 0.3 \times 0.316 - 0.7 \times 0.79 \\ &= -0.4582 \end{aligned}$$



for  $d_4$ ,

$$\max \left\{ \underset{0.6}{\text{sim}(d_1, d_4)} ; \underset{0.8}{\text{sim}(d_2, d_4)} \right\} = 0.8$$

$$\text{Sim}(d_4, q) = 0.4$$

$$\begin{aligned} \text{MMR} &= 0.3 \times 0.4 - 0.7 \times 0.8 \\ &= -0.44. \end{aligned}$$

mmr for  $d_4$  is max 8.,

$$S = \{d_1, d_2, d_4\}.$$

These are 3 sentences in summary set.

Sol 6:- Hidden Markov Model  $\rightarrow$  HMM.

Initial Probability:-

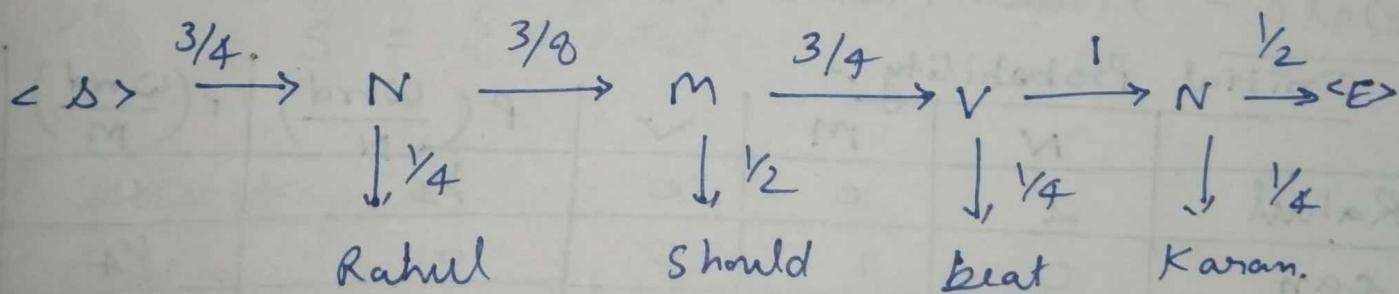
	N	M	V	$P\left(\frac{\text{word}}{N}\right)$	$P\left(\frac{\text{word}}{M}\right)$	$P\left(\frac{\text{word}}{V}\right)$
Rahul	2	0	0	$\frac{1}{4}$	0	0
can	0	1	0	0	$\frac{1}{4}$	0
play	0	0	2	0	0	$\frac{1}{2}$
chess	2	0	0	$\frac{1}{4}$	0	0
Karan	2	0	0	$\frac{1}{4}$	0	0
Should	0	2	0	0	$\frac{1}{2}$	0
Clean	0	0	1	0	0	$\frac{1}{4}$
table	1	0	0	$\frac{1}{8}$	0	0
will	0	1	0	0	$\frac{1}{4}$	0
Shayam	1	0	0	$\frac{1}{8}$	0	0
beat	0	0	1	0	0	$\frac{1}{4}$
	8	4	4			

Using HMM,

$P(NM \vee N / \text{Rahul should beat Karan})$ .

	N	M	V	$\langle E \rangle$
$\langle S \rangle$	$3/4$	$1/4$	0	0
N	0	$3/8$	$1/8$	$4/8$
M	$1/4$	0	$3/4$	0
V	$4/4$	0	0	0

"Rahul should beat Karan"



$$\begin{aligned}
 P &= \frac{3}{4} \times \frac{1}{4} \times \frac{3}{8} \times \frac{1}{2} \times \frac{3}{4} \times \frac{1}{4} \times 1 \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{2} \\
 &= \frac{27}{2^{15}} = 0.00082357.
 \end{aligned}$$



Sol 5(b):-

for,  $Sim_{path}$  :

$Path_{len}$  (Hill, Shore)

$$\Rightarrow 1 + 3 \Rightarrow \underline{4}$$

$$Sim_{path} = \frac{1}{4} = 0.25$$

for  $Sim_{sensik}$  :

$$- \log P(LCS(Hill, Shore))$$

$$- \log_2 P(\text{geo-form})$$

$$- \log_2(0.00176) = 9.15$$

for  $Sim_{Lin}$

$$\therefore \frac{2 \log P(LCS(Hill, Shore))}{\log P(Hill) + \log P(Shore)}$$

$$\Rightarrow \frac{2 \log P(\text{geo form})}{\log P(Hill) + \log P(Shore)}$$

$$\Rightarrow \frac{2 \ln(0.00176)}{\ln(0.0000189) + \ln(0.0000836)}$$

$$\Rightarrow 0.6259.$$

Sol 5 (a):-

Transition Probability :-

	jouer	croquet	la	grillion.
play	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
cricket	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
the	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
team	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Alignment Probability :-

$$\begin{array}{c} \text{Play} \\ | \\ \text{jouer} \end{array} \quad \begin{array}{c} \text{cricket} \\ | \\ \text{croquet} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\begin{array}{c} \text{the} \\ | \\ \text{la} \end{array} \quad \begin{array}{c} \text{cricket} \\ | \\ \text{grillon} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\begin{array}{c} \text{cricket} \\ | \\ \text{croquet} \end{array} \quad \begin{array}{c} \text{team} \\ | \\ \text{equipe} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\begin{array}{c} \text{play} \\ \diagdown \quad \diagup \\ \text{jouer} \quad \text{croquet} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\begin{array}{c} \text{the} \\ \diagdown \quad \diagup \\ \text{la} \quad \text{grillon} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\begin{array}{c} \text{cricket} \\ \diagdown \quad \diagup \\ \text{croquet} \quad \text{equipe} \end{array} \quad \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\text{Normalization} = \frac{\frac{1}{9}}{\frac{1}{9} + \frac{1}{9}} = \frac{1}{2}$$

weighted  
translation  
count  $\rightarrow$

	play	jouer	croquer	la	gailler	griffe
		$\frac{1}{2}$	$\frac{1}{2}$	0	0	0
the		0	0	$\frac{1}{2}$	$\frac{1}{2}$	0
Cricket		$\frac{1}{2} \times \frac{1}{3}$	$\frac{1}{3} \left( \frac{1}{2} + \frac{1}{2} \right)$	$\frac{1}{2} \times \frac{1}{3}$	$\frac{1}{2} \times \frac{1}{3}$	$\frac{1}{2} \times \frac{1}{3}$
team		0	$\frac{1}{2}$	0	0	$\frac{1}{2}$

Not realization

Alignment

~~Not realized~~ :-

Alignment Probability :-

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

$$\frac{1}{2} \times \frac{1}{6} \times \frac{1}{12}$$

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

$$\frac{1}{2} \times \frac{1}{6} \times \frac{1}{12}$$

$$\frac{1}{8} \times \frac{1}{2} = \frac{1}{16}$$

$$\frac{1}{2} \times \frac{1}{6} = \frac{1}{12}$$



Sol 7:-

To handle co-references in English Text, we can propose a novel methodology that combines traditional rule based approaches with machine learning techniques. The methodology involves the following steps:

- 1) Text Preprocessing :- Text is tokenized and parts of speech (POS) tagging and Named entity Recognition (NER) are performed to identify entities and their types.
- 2) Coreference Resolution:- Initially, a rule based approach is applied to resolve simple coreferences based on exact match of words or phrases. For ex:- in the sentence "John went to the market. He bought some fruits", the rule-based system can resolve "He" to refer to "John".
- 3) Feature Extraction:- Features such as mention headword, syntactic distance between mentions and semantic similarity between mentions are extracted to represent mentions and their contexts.

4. Machine learning models:- A machine learning model, such as neural network or a conditional random field (CRF) is trained on dataset annotated with coreferences. The model uses the extracted features to learn patterns in coreference resolution.

5) Evaluation:- The methodology is evaluated using standard coreference resolution evaluation metrics such as MUC, B-CUBED, and CEAF. These metrics measures the precision, recall, and F1 score of the resolved coreferences compared to a gold standard.

Ex:- "Tom visited Mary at her house. He brought her a gift."

Using our methodology, the system first identifies 'Tom' & 'Mary' as entities and resolves the pronoun 'He' and 'Her' to refer to 'Tom' and 'Mary' respectively. The rule-based approach handles the simple coreferences, while the machine learning model improves the resolution by considering contextual features.

Evaluation methodology performance is evaluated on a dataset of similar contexts annotated with coreferences.



Sol: PPMI

Term	Context	Auto	Comp	Money	House	Politics	$\Sigma$
1. car		10	8	1	0	3	22
2. Auto		5	1	0	0	1	7
3. Insu		1	0	4	3	0	8
4. Window		1	2	1	2	1	7
5. Comp		3	2	1	0	0	6
6. Tech		1	3	1	1	1	7
$\Sigma$		21	16	8	6	6	57

$N = 57$

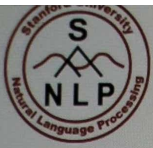
$P_{\text{Term}}$	$P_{\text{Context}}$	$\frac{21}{57}$	$\frac{16}{57}$	$\frac{8}{57}$	$\frac{6}{57}$	$\frac{6}{57}$
$\frac{22}{57}$		$\frac{10}{57}$	$\frac{8}{57}$	$\frac{1}{57}$	0	$\frac{3}{57}$
$\frac{7}{57}$		$\frac{5}{57}$	$\frac{1}{57}$	0	0	$\frac{1}{57}$
$\frac{8}{57}$		$\frac{1}{57}$	0	$\frac{4}{57}$	$\frac{3}{57}$	0
$\frac{7}{57}$		$\frac{1}{57}$	$\frac{2}{57}$	$\frac{1}{57}$	$\frac{2}{57}$	$\frac{1}{57}$
$\frac{6}{57}$		$\frac{3}{57}$	$\frac{2}{57}$	$\frac{1}{57}$	0	0
$\frac{7}{57}$		$\frac{1}{57}$	$\frac{3}{57}$	$\frac{1}{57}$	$\frac{1}{57}$	$\frac{1}{57}$



$$PPMI(w_1, w_2) = \max\left(\log_2 \frac{p(w_1, w_2)}{p(w_1) \cdot p(w_2)}, 0\right)$$

$$PPMI(w, c) = \max\left(\log_2 \frac{p(w, c)}{p(w) \cdot p(c)}, 0\right)$$

PPMI:-	Auto	Corp	Money	Hours	Relative
car					
auto					
ingr					
win					
com					
Eu					



## LESK'S ALGORITHM

**Sense Bag:** contains the words in the definition of a candidate sense of the ambiguous word.

**Context Bag:** contains the words in the definition of each sense of each context word.

E.g. "On burning **coal** we get **ash**."

### Ash

- **Sense 1**  
Trees of the olive family with pinnate leaves, thin furrowed bark and gray branches.
- **Sense 2**  
The **solid** residue left when **combustible** material is thoroughly **burned** or oxidized.
- **Sense 3**  
To convert into ash

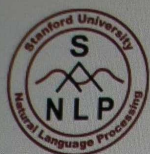
### Coal

- **Sense 1**  
A piece of glowing carbon or **burnt** wood.
- **Sense 2**  
charcoal.
- **Sense 3**  
A black **solid combustible** substance formed by the partial decomposition of vegetable matter without free access to air and under the influence of moisture and often increased pressure and temperature that is widely used as a fuel for **burning**

In this case Sense 2 of ash would be the winner sense.







## Summary: thesaurus-based similarity

$$\text{sim}_{\text{path}}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

$$\text{sim}_{\text{resnik}}(c_1, c_2) = -\log P(\text{LCS}(c_1, c_2)) \quad \text{sim}_{\text{lin}}(c_1, c_2) = \frac{2 \log P(\text{LCS}(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$$

$$\text{sim}_{\text{jiangconrath}}(c_1, c_2) = \frac{1}{\log P(c_1) + \log P(c_2) - 2 \log P(\text{LCS}(c_1, c_2))}$$

$$\text{sim}_{\text{eLesk}}(c_1, c_2) = \sum_{r, q \in \text{RELS}} \text{overlap}(\text{gloss}(r(c_1)), \text{gloss}(q(c_2)))$$