Sadamala Sasi Rekha
23BDS049

Theoretical Assignment - CS458 (NLP)

## Question:1

A) Tokens are the individual words obtained after splitting the text

- Tokens:

  The, children, are, playing, in, the, playgrounds, A, child, plays, happily, while, the, others, played, earlier, They, have, been, playing, every, afternoon.

B) Number of types and tokens:

- Types (unique words) : 20
  "The|the" and "A|a" are considered the same after case normalization.
- Tokens (total words) : 22

C) Lemmatize all the nouns:

children → child
Playgrounds → Playground
child → child
others → other
afternoon → afternoon

D) Lemmatization of verbs:

are → be
Playing → play
Plays → play
Played → play
have → have
been → be
Playing → play

# ① Question 2:

word pair : kitten → sitting

costs : insertion=1, Deletion=1, Substitution=2

|   |   | O | S | I | T | T | I | N | G |
|---|---|---|---|---|---|---|---|---|---|
| O |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| K | 1 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| i | 2 | 2 | 3 | 2 | 3 | 4 | 5 | 6 | 7 |
| t | 3 | 3 | 4 | 3 | 2 | 3 | 4 | 5 | 6 |
| t | 4 | 4 | 5 | 4 | 3 | 2 | 3 | 4 | 5 |
| e | 5 | 5 | 6 | 5 | 4 | 3 | 4 | 5 | 6 |
| n | 6 | 6 | 7 | 6 | 5 | 4 | 5 | 4 | 5 |

Backtracking:

Starting from cell $(n, g) = 5$

1) Insert G at end (cost = 1)

2) Match ($n \to n \Rightarrow$ cost = 0)

3) Substitution: $e \to i$ (cost = 2)

4) match : $t \to t$ (cost = 0)

5) match : $t \to t$ (cost = 0)

6) Substitution : $k \to s$ (cost = 2)

✴ Sequence of operations:

Substitution ($k \to s$ : 2)

Substitution ($e \to i$ : 2)

insertion (g : 1)

Total cost $2 + 2 + 1 = 5$