

The Word Ladder Problem

Artificial Intelligence

Suhas N (EC21B1113), Reddappa Reddy (EC21B1114)

April 27, 2024

1 Introduction

The Word Ladder Problem is a puzzle in which the objective is to transform one word into another by changing a single letter at a time, with the constraint that each intermediate step must also be a valid word. The challenge is to find the shortest sequence of such transformations, known as the word ladder, linking the initial word to the target word. The Word Ladder Problem has applications in various fields, including natural language processing, search engines, and puzzle-solving algorithms.

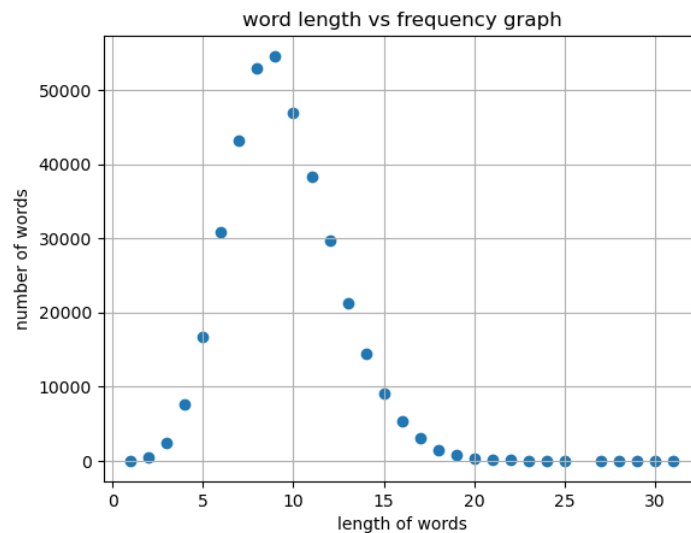
2 Solving the Problem

2.1 Gathering the data

To solve the word ladder problem, we need to form a corpus of words. We have taken the words from the following sources:

- Books listed on Project Gutenberg (<https://www.gutenberg.org/>)
- The Wordnet database (from python-nltk)

From the above sources, we have collected exactly 3,79,557 words. The following plot shows the number of words per length of words.



We may notice that a large portion of available words are 5 to 10 characters long.

2.2 Constructing a graph

Using the words and based on the problem statement, we will attempt to construct a graph using the below algorithm.

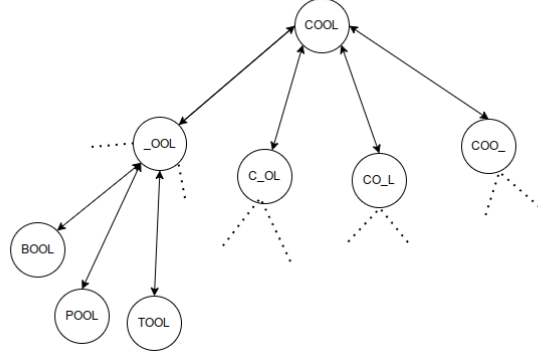
```
1 def generate(word):
2     yield from [
3         word[:i] + '_' + word[i+1:]
4         for i in range(len(word))
5     ]
6 def construct_graph():
7     graph = defaultdict(lambda : defaultdict(set))
```

```

8         for length, words in self.word_map.items():
9             for word in words:
10                 for i in range(len(word)):
11                     pattern = word[:i] + '_' + word[i+1:]
12                     graph[length][pattern].add(word)
13     return graph

```

We attempt to create separate graphs for different word lengths, as the length of a word does not change throughout the problem. The obtained graph would look like the following (for a word 'COOL'):



Now, this graph can be traversed by using algorithms, and we can obtain a chain connecting two words on the graph. Construction of the graph can happen on the fly (while traversing), but due to the size of our corpus, we are constructing the graph before traversing to facilitate analysis. The constructed graph has over 50,00,000 nodes, a much higher number than the number of words in the corpus, as a large number of words get repeated.

2.3 Traversing the graph

Due to the number of nodes in the constructed graph, we may not use the depth-first search algorithm as the it would not be able to backtrack long paths. But a breadth-first search is applicable.

```

1  from collections import deque
2
3  def bfs(begin, end, words, graph_):
4      assert not {begin, end} - words
5      assert len(begin) == len(end)
6      graph = graph_[len(begin)]
7
8      queue = deque([(begin, [begin])])
9      visited = set([begin])
10
11     iterations = 0
12     while queue:
13         w, path = queue.popleft()
14         if w == end: # goal
15             return {"length": len(path), "iterations": iterations}
16         for pattern in generate(w):
17             delta = graph[pattern] - visited
18             queue.extend([(d, path+[d]) for d in delta])
19             visited.update(delta)
20         iterations += 1
21
22     return {"length": float('inf'), "iterations": iterations} # fail

```

Upon analysis, the above algorithm takes at least 1,000-10,000 iterations for traversing the graph (for word lengths more than 5) and gives an optimal solution. To minimize the number of iterations, we can use the bidirectional breadth-first search algorithm, where we maintain two queues and check for intersection between the visited sets.

```

1  def bfs_step(queue, visited, graph):
2      q2 = deque()
3      iterations = 0
4      while queue:
5          w, path = queue.popleft()
6          for pattern in generate(w):
7              delta = graph[pattern] - visited
8              visited.update(delta)
9              q2.extend([(nw, path + [nw]) for nw in delta])

```

```

10         iterations += 1
11     return q2, iterations
12
13 def bi_bfs(begin, end, words, graph_):
14     assert not {begin, end} - words
15     assert len(begin) == len(end)
16     graph = graph_[len(begin)]
17
18     q1 = deque([[begin, [begin]]])
19     q2 = deque([[end, [end]]])
20
21     v1 = set([begin])
22     v2 = set([end])
23
24     iterations = 0
25     while q1 and q2:
26         if v1.intersection(v2): # goal
27             _, path1 = q1.popleft()
28             _, path2 = q2.popleft()
29             path = path1 + path2[::-1]
30             return {"length": len(path),
31                   "iterations": iterations}
32
33         if len(q1) <= len(q2):
34             q1, it = bfs_step(q1, v1, graph)
35         else:
36             q2, it = bfs_step(q2, v2, graph)
37         iterations += it
38
39     return {"length": float('inf'), "iterations": iterations} # fail

```

The above algorithm does not exceed 500 iterations for finding a solution and returns an optimal solution comparable to the traditional breadth-first search algorithm. But, in order to guide the search, we can use the variants of A★.

2.4 Guided Search

For using the variants of A★, we need an admissible heuristic function. From analysis, the best heuristic function is found to be:

```

1 def heuristic(w, w2):
2     return sum([a != b for a, b in zip(w, w2)])

```

The heuristic function is clearly the distance of word 'w' from the word 'w2'.

Using the above heuristic function, the following algorithm can be used to find a word ladder:

```

1 def astar(start, end, words, graph_, heuristic):
2     assert not {start, end} - words
3     assert len(start) == len(end)
4
5     pq = [[heuristic(start, end), 0, start, [start]]]
6     visited = set()
7     graph = graph_[len(start)]
8
9     iterations = 0
10    while pq:
11        _, cost, w, path = heapq.heappop(pq)
12        if w == end:
13            return {"length": len(path), "iterations": iterations}
14        for p in generate(w):
15            delta = graph[p] - visited
16            visited.update(delta)
17            priorities = [cost + heuristic(w, end) for w in delta]
18            pq.extend([(p, 1+cost, w, path + [w]) for w, p in zip(delta, priorities)])
19            heapq.heapify(pq)
20            iterations += 1
21
22    return {"length": float('inf'), "iterations": iterations}

```

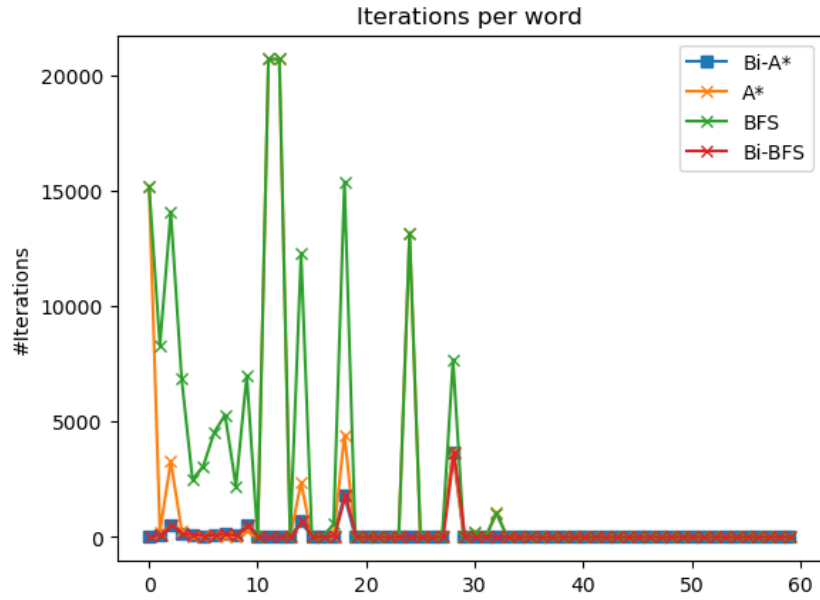
This algorithm performs worse than the Bi-directional breadth-first search algorithm, as it traverses through many nodes which are away from the end word. To solve this problem, we can use the bi-directional A★ algorithm. The following algorithm can be used to guide the search efficiently.

```

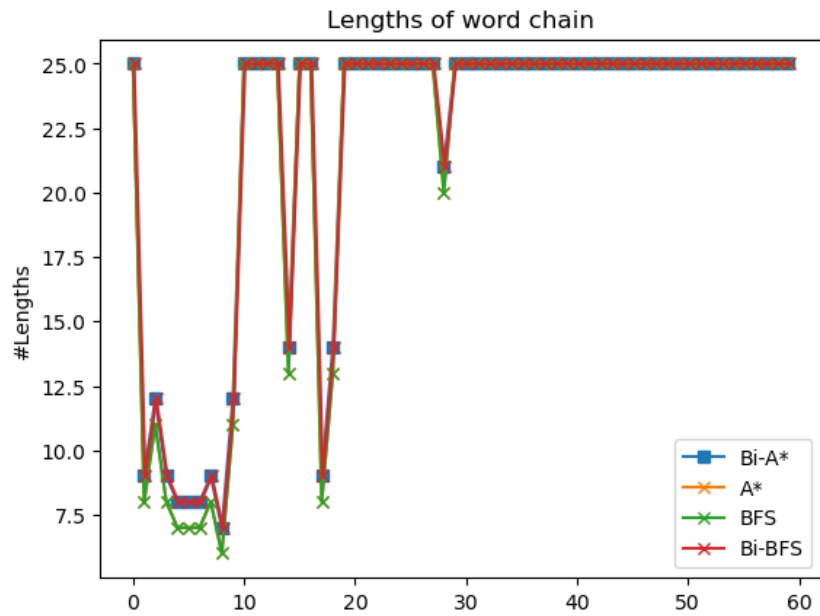
1  def bi_astar(start, end, words, graph_, heuristic):
2      assert not {start, end} - words
3      assert len(start) == len(end)
4
5      graph = graph_[len(start)]
6
7      pq1 = [(heuristic(start, end), 0, start, [start])]
8      pq2 = [(heuristic(end, end), 0, end, [end])]
9
10     v1, v2 = set(), set()
11     iterations = 0
12
13     while pq1 and pq2:
14         if v1.intersection(v2):
15             *_ , path1 = pq1[0]
16             *_ , path2 = pq2[0]
17             path = path1 + path2[::-1]
18             return {"length": len(path), "iterations": iterations}
19
20         it = 0
21
22         if len(pq1) <= len(pq2):
23             pq1_ = []
24             while pq1:
25                 it += 1
26                 _, c1, w1, path1 = heapq.heappop(pq1)
27
28                 for p in generate(w1):
29                     delta = graph[p] - v1
30                     v1.update(delta)
31                     priorities = [c1 + heuristic(w, end) for w in delta]
32                     pq1_.extend([(p, 1+c1, w, path1 + [w]) for w, p in zip(delta,
33                                     priorities)])
34                     pq1 = pq1_
35                     heapq.heapify(pq1)
36
37             else:
38                 pq2_ = []
39                 while pq2:
40                     it += 1
41                     _, c2, w2, path2 = heapq.heappop(pq2)
42
43                     for p in generate(w2):
44                         delta = graph[p] - v2
45                         v2.update(delta)
46                         priorities = [c2 + heuristic(w, end) for w in delta]
47                         pq2_.extend([(p, 1+c2, w, path2 + [w]) for w, p in zip(delta,
48                                     priorities)])
49                         pq2 = pq2_
50                         heapq.heapify(pq2)
51
52             iterations += it
53     return {"length": float('inf'), "iterations": iterations}

```

The above algorithm performs on-par with the bi-directional breadth first algorithm. The following graph compares all the above algorithms in terms of number of iterations taken.



The following plot shows the length of the word chain found through various algorithms:



The above plot shows that the BFS algorithm performs better than other algorithms in terms of the length of the word chain, but takes a large number of iterations.

3 Conclusion

In conclusion, the word ladder problem is explored using 4 different algorithms, and it can be seen that BFS algorithm produces shorter chains than the guided-search algorithms.