

Algorithm for Extracting a Graph from Wikipedia Links

1 Abstract

This document describes an algorithm for extracting a graph from Wikipedia links, starting from a root article. The algorithm uses a breadth-first search approach to crawl the Wikipedia links and build the graph, where nodes represent Wikipedia pages and edges represent links between them. The algorithm also includes error handling for cases such as non-existent pages and disambiguation errors.

2 Introduction

The extraction of knowledge graphs from Wikipedia is a crucial task in various domains, including natural language processing, information retrieval, and data mining. This document details a specific algorithm designed for this purpose. It leverages the `wikipedia` and `networkx` Python libraries to retrieve page content and manage the graph structure, respectively. The algorithm performs a breadth-first search, starting from a specified root article, to explore linked Wikipedia pages and construct a graph representation of their relationships.

3 Algorithm Description

The algorithm, named `extract_graph`, takes the title of a root Wikipedia article and an optional maximum number of nodes as input. It returns a NetworkX graph object representing the extracted knowledge graph. The algorithm utilizes a helper function `__crawl` which implements the breadth-first search.

3.1 Pseudo-code

```
1: procedure EXTRACT_GRAPH(root_article_title, max_nodes)
2:   __check_title(root_article_title)                                ▷ Check if root article exists
3:   graph ← nx.Graph()
4:   graph.add_node(root_article_title)
5:   visited ← {root_article_title}
6:   queue ← [root_article_title]
7:   __crawl(graph, queue, visited, max_nodes)
8:   if not nx.is_connected(graph) then
9:     raise Exception("Graph is not connected")
10:  end if
11:  return graph
12: end procedure
13: procedure __CRAWL(graph, queue, visited, max_nodes)
14:   if not queue then
15:     return
16:   end if
17:   if len(graph.nodes) ≥ max_nodes then
18:     return
19:   end if
20:   page_title ← queue.pop(0)
21:   Try
22:     page ← __check_title(page_title)
23:   Catch(Exception)
24:   return                                                         ▷ Skip this branch
```

```

25:   EndTry
26:   for link_title in page.links do
27:     if len(graph.nodes) ≥ max_nodes then
28:       return
29:     end if
30:     if link_title in visited then
31:       __add_edge(graph, page_title, link_title)
32:     else
33:       visited.add(link_title)
34:       queue.append(link_title)
35:       graph.add_node(link_title)
36:       __add_edge(graph, page_title, link_title)
37:     end if
38:   end for
39:   __crawl(graph, queue, visited, max_nodes)
40: end procedure
41: procedure __ADD_EDGE(graph, page_src, page_dst)
42:   if page_src == page_dst then
43:     return ▷ Avoid self-links
44:   end if
45:   if graph.has_edge(page_src, page_dst) or graph.has_edge(page_dst, page_src) then
46:     return ▷ Edge already present
47:   end if
48:   graph.add_edge(page_src, page_dst)
49: end procedure
50: procedure __CHECK_TITLE(page_title)
51:   Try
52:     page ← wikipedia.page(page_title)
53:   Catch(wikipedia.exceptions.PageError)
54:     raise Exception("Page not found")
55:   Catch(wikipedia.exceptions.DisambiguationError)
56:     raise Exception("Disambiguation Error")
57:   EndTry
58:   return page
59: end procedure

```

3.2 Implementation Details

The implementation uses the `wikipedia` library for fetching page content and the `networkx` library for graph manipulation. The `__check_title` function handles potential exceptions during page retrieval. The `__add_edge` function ensures no duplicate or self-referential edges are added. The `json_dump` function provides a way to serialize the graph into JSON format. The example usage demonstrates how to extract and visualize the graph using `matplotlib`.

4 Conclusion

This document presented a detailed algorithm for extracting knowledge graphs from Wikipedia. The algorithm uses a breadth-first search strategy and incorporates error handling to ensure robustness. The generated graph can be further analyzed and visualized for various applications.