

Discrete Mathematics Project 2 Report

Shubham Aggarwal
2023CSB1162

1 PageRank Algorithm

1.1 Introduction

The PageRank algorithm is a widely-used algorithm for ranking web pages in search engine results. Developed by Larry Page and Sergey Brin, the founders of Google, PageRank assigns a numerical value to each element of a hyperlinked set of documents, with the purpose of measuring its relative importance within the set.

1.2 Random Walk Algorithm

The Random Walk algorithm is a stochastic process where an entity moves randomly from one state to another within a defined set of states. In the context of networks, it's often used to model the behavior of a random surfer navigating through web pages. In each step, the surfer either moves to a neighboring page with a certain probability or jumps to a random page with another probability.

1.3 My Code

Listing 1: Python code implementing the random walk algorithm and PageRank

```
1 import csv
2 import networkx as nx
3 import random
4
5 # Create a directed graph
6 G = nx.DiGraph()
7
8 # Define a function for random walk
9 def random_walk(steps):
10     # Initialize a dictionary to store ratings for each node
11     rating = dict.fromkeys(list(G.nodes), 0)
12     # Choose a random starting node
13     current_node = random.choice(list(G.nodes))
14     # Perform random walk for the specified number of steps
15     for i in range(steps):
```

```

16         # With probability 0.15, randomly jump to any node
17         if random.random() < 0.15:
18             current_node = random.choice(list(G.nodes))
19         else:
20             # Select a neighbor of the current node to move
                to
21             neighbors = list(G.neighbors(current_node))
22             if neighbors:
23                 current_node = random.choice(neighbors)
24             else:
25                 # If the current node has no neighbors,
                    randomly jump to any node
26                 current_node = random.choice(list(G.nodes))
27         # Increment the rating of the current node
28         rating[current_node] += 1
29     return rating
30
31 # Open the dataset file
32 with open("./dataset.csv") as f:
33     reader = csv.reader(f)
34     # Iterate over each row in the dataset
35     for row in reader:
36         # Skip the header row
37         if reader.line_num == 1:
38             continue
39         # Get the parent node (source) from the first column
40         parent = row[1].lower().split('@')[0]
41         # Iterate over the child nodes (targets) in the row
42         for entry in row[2:]:
43             # Skip empty entries
44             if entry == '':
45                 continue
46             # Get the child node (target) from the last word
                in the entry
47             child = entry.lower().split()[-1]
48             # Add an edge from the parent to the child in
                the graph
49             G.add_edge(parent, child)
50
51 # Define the total number of steps for the random walk
52 total = 10000000
53 # Perform random walk and get the ratings for each node
54 rating = random_walk(total)
55 # Sort nodes by their ratings in descending order
56 ranking = sorted(rating, key=rating.get, reverse=True)
57 # Print the top 10 leaders based on the random walk
    algorithm
58 print("The top 10 leaders are:")
59 for i in range(10):
60     name = ranking[i]

```

```

61     rate = rating[name] / total # Normalize rating by total
        steps
62     print(i+1, name, rate)
63
64 print("The top leader is", ranking[0], "with", rating[ranking
    [0]]/total, "points.")

```

1.4 Explanation of Code

The provided Python code implements a random walk algorithm and computes PageRank scores for a given network represented as a directed graph. Here's a breakdown of the code:

- The random walk algorithm simulates the behavior of a random surfer moving through the network. It randomly selects nodes to visit based on their neighbors and performs a certain number of steps to calculate ratings for each node.

1.5 Conclusion

In conclusion, both the Random Walk algorithm and the PageRank algorithm are powerful tools for analyzing and ranking nodes within networks. While the Random Walk algorithm provides a probabilistic model for navigating networks, PageRank offers a more sophisticated approach based on the concept of importance and connectivity. By understanding and implementing these algorithms, we can gain valuable insights into the structure and dynamics of various types of networks..

2 Missing Links Prediction Using Matrix Method

2.1 Introduction

Missing links prediction using the matrix method is a technique in network analysis where the goal is to predict the existence of edges (links) between nodes in a network by representing the network as an adjacency matrix and solving a system of linear equations. This approach leverages linear algebra concepts to infer the likelihood of connections between nodes based on the network's structure.

2.2 Matrix Method for Missing Links Prediction

The matrix method for missing links prediction involves the following steps:

1. **Adjacency Matrix Representation:** Represent the given network as an adjacency matrix, where each entry A_{ij} indicates the presence or absence of an edge between nodes i and j .

2. **Linear Combination:** Treat each row of the adjacency matrix as a vector and express one row as a linear combination of the other rows. This can be done by solving a system of linear equations.
3. **Prediction:** Once the coefficients of the linear combination are determined, use these coefficients with the corresponding column vectors to predict the missing entries in the adjacency matrix.

2.3 My Code

Listing 2: Python code implementing missing links prediction using the matrix method

```

1 import csv
2 import networkx as nx
3 import numpy as np
4
5 # Create a directed graph
6 G = nx.DiGraph()
7
8 # Open the dataset file and create the graph
9 with open("./dataset.csv") as f:
10     reader = csv.reader(f)
11     # Iterate through each row in the dataset
12     for row in reader:
13         # Skip the header row
14         if reader.line_num == 1:
15             continue
16         # Extract the parent node (source) from the second
17         # column
18         parent = row[1].lower().split('@')[0]
19         # Iterate through each entry in the row (child nodes
20         # )
21         for entry in row[2:]:
22             # Skip empty entries
23             if entry == '':
24                 continue
25             # Extract the child node (target) from the last
26             # word in the entry
27             child = entry.lower().split()[-1]
28             # Add an edge from the parent to the child in
29             # the graph
30             G.add_edge(parent, child)
31
32 # Convert the graph to adjacency matrix
33 nodes = list(G.nodes())
34 n = len(nodes)
35 adj = np.zeros(shape=(n,n))
36 for i in range(n):

```

```

33     for j in range(n):
34         # Check if there is an edge between nodes i and j
35         if G.has_edge(nodes[i], nodes[j]):
36             # If there is an edge, set the corresponding
37             # entry in the adjacency matrix to 1
38             adj[i][j] = 1
39
40 # Function to predict missing links using matrix method
41 def predict(i, j):
42     # Remove the j-th column from the adjacency matrix
43     without_col = np.delete(adj, j, axis=1)
44     # Extract the i-th row from the modified adjacency
45     # matrix
46     B = without_col[i]
47     # Remove the i-th row from the modified adjacency matrix
48     A = np.delete(without_col, i, axis=0)
49     # Solve the system of linear equations Ax = B to find
50     # coefficients x
51     X = np.linalg.lstsq(A.T, B.T, rcond=None)[0]
52     # Remove the i-th row from the original adjacency matrix
53     without_row = np.delete(adj, i, axis=0)
54     # Extract the j-th column from the modified adjacency
55     # matrix
56     C = without_row[:, j]
57     # Calculate the expected value for the missing link
58     expected = np.matmul(C, X)
59     return expected
60
61 # Find missing links and add them to the graph
62 missing_links = []
63 for i in range(n):
64     for j in range(n):
65         # Check if there is no edge between nodes i and j
66         if adj[i][j] == 0:
67             # Predict the value of the missing link using
68             # the matrix method
69             val = predict(i, j)
70             # If the predicted value is greater than or
71             # equal to 1, consider it as a missing link
72             if val > 0.5:
73                 missing_links.append((i, j))
74
75 for i, j in missing_links:
76     # Print the missing link (node i to node j)
77     print("Missing Link:", (nodes[i], nodes[j]))
78     # Add the missing link to the graph
79     G.add_edge(nodes[i], nodes[j])
80
81 # Calculate PageRank scores for the updated graph
82 rating = nx.pagerank(G)

```

```

77 # Sort nodes by their PageRank scores in descending order
78 ranking = sorted(rating, key=rating.get, reverse=True)
79 # Print the top 10 nodes with the highest PageRank scores
80 for i in range(10):
81     name = ranking[i]
82     rate = rating[name]
83     print(name, rate)

```

2.4 Code Explanation

The Python code provided below performs missing links prediction using the matrix method. Here's a detailed explanation of each section:

- The code starts by importing necessary libraries including `csv`, `networkx`, `numpy`, and `pandas`.
- It creates an empty directed graph `G` using `NetworkX`.
- The code reads a dataset from a CSV file to populate the directed graph `G`.
- The graph is then converted to an adjacency matrix representation using `NumPy`.
- A function `predict` is defined to predict missing links using the matrix method. This function takes two indices `i` and `j` corresponding to the rows and columns of the adjacency matrix.
- Missing links are predicted based on the matrix method. If the predicted value is greater than or equal to 1, the missing links are added to the graph `G`.
- PageRank scores are calculated for the updated graph using `NetworkX`, and the top 10 nodes with the highest PageRank scores are printed.

2.5 Conclusion

In conclusion, missing links prediction using the matrix method offers a powerful approach to infer connections between nodes in a network based on its structural properties. By leveraging linear algebra techniques, we can accurately predict missing links and gain insights into the underlying relationships within the network. The code provided serves as a practical implementation of the matrix method for missing links prediction, demonstrating the application of discrete mathematics concepts in network analysis.

3 Dijkstra's Algorithm

3.1 Introduction

Dijkstra's algorithm is a popular algorithm in graph theory for finding the shortest paths between nodes in a graph.

Dijkstra's algorithm, named after Dutch computer scientist Edsger W. Dijkstra, is a graph search algorithm that finds the shortest path between nodes in a weighted graph.

3.2 Algorithm Steps

The algorithm works by repeatedly selecting the node with the shortest distance from the source node and updating the distances of its neighbors accordingly.

3.3 Code

Listing 3: Python code for Dijkstra's algorithm

```
1 import csv
2 import networkx as nx
3 import numpy as np
4 import sys
5
6 # Create a directed graph
7 G = nx.DiGraph()
8
9 # Open the dataset file and create the graph
10 with open("./dataset.csv") as f:
11     reader = csv.reader(f)
12     for row in reader:
13         # Skip the header row
14         if reader.line_num == 1:
15             continue
16         # Extract the parent node (source) from the second
17         # column
18         parent = row[1].lower().split('@')[0]
19         # Iterate through each entry in the row (child nodes)
20         for entry in row[2:]:
21             # Skip empty entries
22             if entry == '':
23                 continue
24             # Extract the child node (target) from the last
25             # word in the entry
26             child = entry.lower().split()[-1]
27             # Add an edge from the parent to the child in
28             # the graph
29             G.add_edge(parent, child)
```

```

27
28 # Convert the graph to an adjacency matrix
29 nodes = list(G.nodes())
30 n = len(nodes)
31 adj = np.zeros(shape=(n,n))
32 for i in range(n):
33     for j in range(n):
34         # If there is an edge between nodes i and j, set the
           corresponding entry in the adjacency matrix to 1
35         # If there is no edge, set the entry to -1
36         if G.has_edge(nodes[i], nodes[j]):
37             adj[i][j] = 1
38         else:
39             adj[i][j] = -1
40
41 # Implementation of Dijkstra's algorithm
42 def dijkstra(adjacency_matrix, start_vertex):
43     # Get the number of nodes in the graph
44     n = len(adjacency_matrix[0])
45
46     # Initialize an array to store the shortest distances
           from the start_vertex to all other vertices
47     shortest_distances = [sys.maxsize] * n
48     # Initialize a boolean array to track whether each
           vertex has been added to the shortest path tree
49     added = [False] * n
50
51     # Set the distance from start_vertex to itself as 0
52     shortest_distances[start_vertex] = 0
53     # Initialize an array to store the parent vertices of
           each vertex in the shortest path tree
54     parents = [-1] * n
55     parents[start_vertex] = -1
56
57     # Iterate over all vertices
58     for i in range(1, n):
59         # Find the vertex with the minimum distance from the
           start_vertex that has not been added to the
           shortest path tree
60         nearest_vertex = -1
61         shortest_distance = sys.maxsize
62         for vertex_index in range(n):
63             if not added[vertex_index] and
               shortest_distances[vertex_index] <
               shortest_distance:
64                 nearest_vertex = vertex_index
65                 shortest_distance = shortest_distances[
                   vertex_index]
66
67         # Add the nearest_vertex to the shortest path tree

```



```

68         added[nearest_vertex] = True
69
70         # Update the shortest distances to all vertices
71         # adjacent to nearest_vertex
72         for vertex_index in range(n):
73             edge_distance = adjacency_matrix[nearest_vertex][vertex_index]
74             # If there is an edge from nearest_vertex to
75             # vertex_index and the distance through
76             # nearest_vertex is shorter than the current
77             # shortest distance, update the shortest
78             # distance
79             if edge_distance > 0 and shortest_distance +
80                 edge_distance < shortest_distances[
81                     vertex_index]:
82                 parents[vertex_index] = nearest_vertex
83                 shortest_distances[vertex_index] =
84                     shortest_distance + edge_distance
85
86         return shortest_distances
87
88 # Calculate the shortest distances between all pairs of
89 # nodes using Dijkstra's algorithm
90 distance = np.zeros(shape=(n,n))
91 for i in range(n):
92     arr = dijkstra(adj, i)
93     for j in range(n):
94         # If the shortest distance is infinite (no path
95         # exists), set the distance to -1
96         if arr[j] == sys.maxsize:
97             distance[i][j] = -1
98         else:
99             distance[i][j] = int(arr[j])
100
101 # Print the shortest distance matrix
102 print("The shortest distance matrix:")
103 print(distance)
104
105 #Observations based on the algorithm
106
107 max_distance = int(distance.max())
108
109 print("Maximum distance between any 2 nodes:", max_distance)
110 #This verifies that we can reach from one person to another
111 #through less than logN steps.
112
113 unreachable = 0
114 for i in range(n):
115     for j in range(n):
116         if distance[i][j] == -1:

```

```

106         unreachable += 1
107     print("Number of pairs of unreachable nodes:", unreachable)
108     #We can find out how many people were absent during the
109     exercise through this
110     print("Number of students absent:", unreachable/142)
111
112     counts = [0]*(max_distance+1)
113     for i in range(n):
114         for j in range(n):
115             x = int(distance[i][j])
116             if 0<=x<=max_distance:
117                 counts[x] += 1
118     for i,cnt in enumerate(counts):
119         print("Count of", i, ":", cnt)

```

3.4 Code Explanation

The Python code provided below performs the calculation of the shortest distance matrix between all pairs of nodes in a directed graph using Dijkstra's algorithm. Here's a detailed explanation of each section:

- The code starts by importing necessary libraries including `csv` and `networkx`.
- It creates an empty directed graph `G` using `NetworkX`.
- The code reads a dataset from a CSV file to populate the directed graph `G`. Each row in the CSV file represents a directed edge from a parent node (source) to one or more child nodes (targets).
- The graph is then converted to an adjacency matrix representation using `NumPy`. The adjacency matrix represents the connections between nodes in the graph, where a value of 1 indicates the presence of an edge, and -1 indicates no edge.
- Dijkstra's algorithm is implemented to calculate the shortest distances between all pairs of nodes in the graph. The `dijkstra` function takes the adjacency matrix and a starting vertex as input and returns an array containing the shortest distances from the starting vertex to all other vertices.
- The shortest distance matrix is initialized as a 2D `NumPy` array of zeros with shape `(n, n)`, where `n` is the number of nodes in the graph.
- The shortest distances between all pairs of nodes are calculated using Dijkstra's algorithm and stored in the matrix `distance`.
- If there is no path between two nodes, the distance is set to -1.
- Finally, the shortest distance matrix is printed to display the shortest distances between all pairs of nodes in the graph.

3.5 Observations

- We can reach from one person to another in less than $\log(N)$ steps in a network, where N are number of persons.
- All the diagonal entries in distance matrix are 0, as it represents distance from i to i .
- The count of 0 in distance entries represents number of persons.
- The count of 1 in distance entries represents the total number of entries filled in the survey.
- Number of students absent = $\frac{\text{Number of nonreachable node pairs}}{\text{Number of students} - 1}$
- Count of 2 is maximum.
- Count vs N follows bell curve with maxima at 2.

3.6 Conclusion

In conclusion, Dijkstra's algorithm is a powerful tool for finding the shortest paths between nodes in a graph. By efficiently traversing the graph and updating distances, it provides an optimal solution for various routing and navigation problems.