

# **Mapping Influential Articles in Local Wikipedia Networks via Personalized PageRank and Hyperlink-Induced Topic Search**

Soham Jain and Saye Karthikeyan

School of Computer Science, Carnegie Mellon University

Matrices and Linear Transformations (21-241)

Pavel Kovalev

December 5, 2025

# 1 Mathematical Background

## 1.1 Markov Chains

A probability vector is a vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$$

whose entries are all nonnegative and sum to one. In other words,

$$x_1, x_2, \dots, x_n \geq 0 \quad \text{and} \quad x_1 + x_2 + \dots + x_n = 1.$$

Some example of probability vectors are

$$\begin{bmatrix} 0.3 \\ 0.2 \\ 0.5 \end{bmatrix}, \quad \begin{bmatrix} 0.25 \\ 0.75 \\ 0 \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}.$$

A Markov matrix is a matrix  $M \in M_{n \times n}(\mathbb{R})$  whose columns are all probability vectors. Thus,

$$M = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n] \text{ where } \mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n \text{ are probability vectors.}$$

Some examples of Markov matrices are

$$\begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 1 & 0.2 \\ 0 & 0.8 \end{bmatrix}.$$

Let  $M \in M_{n \times n}(\mathbb{R})$  be a Markov matrix and let  $\mathbf{x}_0 \in \mathbb{R}^n$  be an initial probability vector. Then, the following sequence of probability vectors is called a Markov chain:

$$\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots \text{ where each } \mathbf{x}_k = M^k \mathbf{x}_0$$

Consider the Markov matrix  $M = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix}$ . We will attempt to show that 1 is an

eigenvalue of  $M$  via showing  $\det(M - I_n) = 0$ . We can now observe:

$$= \det \begin{bmatrix} a_{1,1} - 1 & \dots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,n} - 1 \end{bmatrix} \quad (\text{Subtraction by } I_n)$$

$$= \det \begin{bmatrix} a_{1,1} + \dots + a_{m,1} - 1 & \dots & a_{1,n} + \dots + a_{m,n} - 1 \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,n} - 1 \end{bmatrix} \quad (\text{From addition of columns to the first row})$$

$$= \det \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,n} - 1 \end{bmatrix} \quad (\text{From definition of stochastic matrix: } M^T \mathbf{1} = \mathbf{1})$$

$$= 0 \quad (\text{From definition of determinant of a matrix with a zero row})$$

Therefore, 1 is an eigenvalue of any Markov matrix  $M$ . This eigenvector associated to  $\lambda = 1$  will be of interest when considering the topic of convergence to a unique steady state.

## 1.2 Random Walks

In the context of Markov chains, a random walk is a process where each state corresponds to a node in a graph, and the probabilities of moving from one state to another are given by a Markov matrix  $M \in M_{n \times n}(\mathbb{R})$ . In particular, each  $M_{ij}$  represents the probability of moving from state  $j$  to state  $i$ . Starting at an initial probability vector  $\mathbf{x}_0$ , the probability distribution after  $k$  steps is  $\mathbf{x}_k = M^k \mathbf{x}_0$ .

Random walks can be used to model the way that a user navigates through links in a website. At each step, the user can choose the next page to visit according to the probabilities in  $M$ . If a page  $j$  has  $d_j$  outgoing links, one simple example of a Markov matrix representation is

$$M_{ij} = \begin{cases} \frac{1}{d_j} & \text{if page } j \text{ links to page } i \\ 0 & \text{if page } j \text{ does not link to page } i \end{cases}$$

The goal of a random walk is often to use an algorithm like PageRank to find a steady state vector  $\mathbf{w}$  where the probabilities converge (i.e.  $\mathbf{x}_k \rightarrow \mathbf{w}$  as  $k \rightarrow \infty$ ). The steady-state distribution can be used to indicate which pages are the most “important” or frequently visited in the long run [6].

## 1.3 PageRank

PageRank refines on this idea of a random walk by capturing real user behavior in two parts.

A specific PageRank score of a page  $i$  can be computed by the sum of all the pages  $w$  that link to  $i$ , divided by the number of outbound links on our selected page  $w$ . We can represent this through  $P(v) = \sum_{i \in Q} \frac{P(i)}{L(i)}$ , where  $Q$  is the set of all pages linked to  $v$ , and  $L(i)$  is the amount of links on each page  $i$ . This gives way to PageRank being defined recursively: if page  $v$  links to an important page  $w$ , then  $v$  must be of relative importance as well.

One issue that we have yet to address is the problem surrounding dangling nodes. Dangling nodes are those that do not have pages linking to them, which would make our algorithm come to a halt. We address this through the damping factor  $d$ , which represents the probability that a user will stop following links and jump to a random page in our network. For the purposes of this paper, we will use  $d = 0.85$  for our damping factor. We can account for dangling nodes with adding a  $\frac{(1-d)}{N}$  term to our PageRank score, which represents the probability that our user chooses any of the  $n$  webpages at random after a dangling node.

Let  $M \in M_{n \times n}(\mathbb{R})$  be a Markov matrix that represents a network of  $n$  pages, where  $M_{ij} > 0$  if page  $j$  links page  $i$  and  $M_{ij} = 0$  otherwise.  $M$  is called an importance matrix. Let  $M' \in M_{n \times n}(\mathbb{R})$  be the same as the matrix  $M$ , but with any zero columns replaced by  $\frac{1}{n}\mathbf{1}$ . Still, the matrix  $M'$  may have zero entries because it may contain columns where some entries are zero and some are nonzero. Thus, we cannot necessarily say that  $M'$  is a positive Markov matrix (a Markov matrix with all positive entries).

In order to ensure a positive Markov matrix representation, PageRank uses this damping factor for dangling nodes. Let  $B = \frac{1}{n}\mathbf{1}\mathbf{1}^T \in M_{n \times n}(\mathbb{R})$ . Then, we can let the Google matrix  $G = dM' + (1 - d)B \in M_{n \times n}(\mathbb{R})$ . Because the values in  $M'$  are originally nonnegative and we are adding a positive  $(1 - d) \times \frac{1}{n}$  to each entry in the matrix,  $G$  is a positive Markov matrix. By the Perron-Frobenius Theorem, there exists a positive steady-state vector  $\mathbf{w}$  that satisfies  $\mathbf{x}_k \rightarrow \mathbf{w}$  with  $\lambda = 1$  as  $k \rightarrow \infty$ . We can use  $\mathbf{w}$  to identify the most important pages in the network (namely, the greatest values in the normalized steady-state vector). A PageRank score  $P_i$  represents the long-run probability that a random-walker will land on page  $i$  in their walk.

## 1.4 Power Iteration

We will attempt to compute our PageRank vector  $P_r$  through power iteration on our stochastic matrix  $M := (K^{-1}A)^T$ , where  $A \in M_{n \times n}(\mathbb{R})$  is our constructed adjacency matrix from graph  $G$ , where  $A_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to } j \\ 0 & \text{otherwise} \end{cases}$ .

First, let  $K = \text{diag}(k_1, \dots, k_n)$ , where  $k_i$  is the number of edges in  $G$  s.t.  $(i, w) \in E, \exists w \in V$ , where  $G$  is our subgraph, consisting of vertices in  $V$  and edges in  $E$ . We can then show that  $K = \text{diag}(\frac{1}{k_1}, \dots, \frac{1}{k_n})$ . We can now show the following:

$$\begin{aligned} K^{-1}A &= \begin{bmatrix} \frac{1}{k_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{k_n} \end{bmatrix} \begin{bmatrix} a_{1,1} & \cdots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \cdots & a_{m,n} \end{bmatrix} \\ &= \begin{bmatrix} \frac{a_{1,1}}{k_1} & \cdots & \frac{a_{1,n}}{k_1} \\ \vdots & \ddots & \vdots \\ \frac{a_{m,1}}{k_n} & \cdots & \frac{a_{m,n}}{k_n} \end{bmatrix} \end{aligned} \quad (\text{Expansion})$$

We know that for some row  $i$ ,  $\sum_{k=1}^n a_{ik} = k_i$  by the way we defined matrix  $K$ . Therefore,  $(K^{-1}A)$  is row-stochastic, and  $(K^{-1}A)^T$  is column stochastic, thus a Markov matrix.

We can now apply our iteration operation, where

$$T(v + 1) = dMT(v) + \frac{1 - d}{n}\mathbf{1}, \mathbf{1} := \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix}$$

We will check on each iteration whether  $|T(v + 1) - T(v)| < \epsilon, \epsilon = 10^{-6}$ , on which case we will consider the steady state to have been stabilized.

## 2 Project Background

The goal of this project is to investigate how PageRank behaves on a subgraph of Wikipedia constructed from a small set of seed articles. By treating Wikipedia as a directed graph, where articles

serve as vertices and hyperlinks act as directed edges, we aim to understand how importance acts in this subgraph. Our objective is to determine whether a seed-focused PageRank procedure can reveal meaningful structures and relations, which could serve useful for further application. We seek to also compare these results with Hyperlink-Induced Topic Search (HITS) to observe the differences that each algorithm gives in importance.

To accomplish this, we first extract a subset of Wikipedia centered around three to four seed articles, recursively gathering linked pages to form a well-defined induced subgraph. We then construct the corresponding Markov transition matrix and apply personalized PageRank, biasing the random-walk distribution toward the chosen seeds. This personalization allows us to emphasize specific topics in our subgraph, effectively influencing steady-state probabilities within the subgraph. Finally, we visualize the resulting directed graph, emphasizing according to the computed PageRank scores to provide an interpretable presentation of article importance.

### 3 Methodology

#### 3.1 Extension: Personalized PageRank

While the PageRank algorithm is efficient in producing consistent results, a problem arises when we want to tailor it towards personalized preferences. Consider a subset  $S \subseteq V$  of our total articles, which we consider to be the seed from where we begin our surfing. Then, some of the articles in our graph are not of much relevance to us.

We can remedy this by making a small change to our PageRank algorithm. We know that in our original algorithm, we compute the chance that our web-surfer makes a random jump to an article in our graph. However, we desire articles closely linked to  $S$  to be of more relevance. Personalized PageRank's impact comes from modifying this idea of a random jump, giving preference to articles that match our interests [2].

In the standard PageRank formula, the teleportation probability distributes uniformly across all nodes. In Personalized PageRank, we replace this uniform distribution with a personalized teleportation vector  $v$ , where  $v_i = \begin{cases} \frac{1}{|S|} & \text{if page } i \in S \\ 0 & \text{otherwise} \end{cases}$ , which concentrates probability mass on the seed pages [3]. We can apply power iteration on this matrix  $M$  until we find a steady state vector that our system converges to [1]. This vector  $p$  represents the long-run probabilities that a web-surfer will land on page  $i$ , which is in other words, the relative importance of a page  $i$ .

We will attempt to make a Column-stochastic Matrix  $M$  from our graph  $G$ , then we will iterate on

a PageRank vector  $p$  which will be initialized to  $\begin{bmatrix} \frac{1}{n} \\ \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{bmatrix}$ , which will terminate once the differences between iterations are within range of our tolerance  $\epsilon$ .

The Google Matrix was not explicitly used, but can be shown to be equivalent to applying power

iteration on our stochastic matrix  $M$ .

### 3.2 Hyperlink-Induced Topic Search

Rather than modeling a random walk, another way to measure importance in a network is by looking at how pages reinforce each other through two roles. One role is that a page can act as a hub if it points to many useful sources. Another role is that a page may act as an authority if many strong hubs point to it. These two features feed into each other through an iterative process, as strong hubs boost authority for the pages they link to, and vice versa. HITS utilizes this approach in order to identify important pages in a network.

Suppose that our directed graph with  $n$  pages has adjacency matrix  $A \in M_{n \times n}(\mathbb{R})$ , where

$$A_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to page } j \\ 0 & \text{if otherwise} \end{cases}$$

Let  $\mathbf{a} \in \mathbb{R}^n$  be the authority vector and let  $\mathbf{h} \in \mathbb{R}^n$  be the hub vector. At each iteration  $k$ , we update these vectors:

$$\mathbf{a}_{k+1} = A^T \mathbf{h}_k, \quad \mathbf{h}_{k+1} = A \mathbf{a}_{k+1}$$

Applying L2 normalization at each iteration ensures that these vectors remain bounded:

$$\mathbf{a}_{k+1} = \frac{\mathbf{a}_{k+1}}{\|\mathbf{a}_{k+1}\|_2}, \quad \mathbf{h}_{k+1} = \frac{\mathbf{h}_{k+1}}{\|\mathbf{h}_{k+1}\|_2}$$

These updates repeat until both vectors stabilize. Substituting the update for  $\mathbf{a}_{k+1}$  into the update for  $\mathbf{h}_{k+1}$ , and vice versa, we obtain:

$$\mathbf{h}_{k+1} = A A^T \mathbf{h}_k, \quad \mathbf{a}_{k+1} = A^T A \mathbf{a}_k$$

Thus, the iterative process of HITS is equivalent to computing the eigenvectors of the matrices  $A^T A$  and  $A A^T$ , which are symmetric and positive definite [4]. Therefore, the algorithm converges when the hub and authority vectors converge to the eigenvectors for  $A A^T$  and  $A^T A$ , respectively. This mathematical interpretation makes it clear how authorities and hubs complement one another. Because authority scores are computed from the hubs that point to them, a page's importance as an authority grows when it is referenced by strong hubs. Similarly, hub scores are computed from the authorities they point to, so a page becomes a stronger hub if it links to important authorities.

After normalization, the largest entries in  $\mathbf{a}$  correspond to the strongest authorities and the largest entries in  $\mathbf{h}$  correspond to the strongest hubs. Thus, we receive two sets of articles from a given seed of articles with HITS, compared to just one set with PageRank [5].

### 3.3 Code: Setup and Graph Initialization

```

1  import wikipediaapi
2  import networkx as nx
3  import matplotlib.pyplot as plt
4  import numpy as np
5  from collections import deque
6  import time
7  import requests
8  import json
9  import random
10 # Included all necessary libraries. NetworkX used for graph creation. WikiAPI
11 # used for scraping necessary data. Matplotlib and Numpy for
    calculation/visualization
12
13 # Initialize Wikipedia API
14 wiki_wiki = wikipediaapi.Wikipedia(
15     user_agent='WikipediaGraphAnalysis/1.0 (https://example.com)',
16     language='en',
17     extract_format=wikipediaapi.ExtractFormat.WIKI
18 )

```

First, we imported all necessary libraries. Wikipedia-API is the primary API we use to load the dataset of Wikipedia articles. We also imported NetworkX, a module for modeling relations between objects and working with graphs, to construct the graph for visualizing the network of Wikipedia articles. In addition, we imported libraries like NumPy and Matplotlib for mathematical calculations and visualizing our results.

```

1  def search_wikipedia_first_result(search_term):
2      """
3      Searched for the very first result for the search term.
4      Doing a search instead of direct lookup since syntactically
5      whatever we enter as a seed might not be the exact page name,
6      so this helps us circumvent that.
7      """
8      try:
9          url = "https://en.wikipedia.org/w/api.php"
10         params = {
11             "action": "query",
12             "list": "search",
13             "srsearch": search_term,
14             "srlimit": 1,
15             "format": "json",
16             "srnamespace": 0
17         }
18
19         headers = {
20             'User-Agent': 'WikipediaGraphAnalysis/1.0 (https://example.com;
                contact@example.com)'
21         }
22
23         response = requests.get(url, params=params, timeout=15, headers=headers)
24
25         # Check response status
26         if response.status_code != 200:
27             print(f"HTTP {response.status_code} for '{search_term}'")
28             return None
29

```

```

30     # Parsing JSON
31     try:
32         data = response.json()
33     except ValueError as e:
34         print(f"Response preview")
35         return None
36
37     if "query" in data and "search" in data["query"] and
38         len(data["query"]["search"]) > 0:
39         title = data["query"]["search"][0]["title"]
40         page = wiki_wiki.page(title)
41         if page.exists():
42             return title
43         return None
44     except Exception as e:
45         print(f"Error searching")
46         return None
47
48     return None
49
50 def get_wiki_links(page_title, max_links=50):
51     """
52     Get internal Wiki links from a page.
53     Based on parameters that we tuned, especially
54     the maximum number of links, which is discussed
55     more thoroughly in the report. Filtering for actual
56     links and not files or categories.
57
58     Randomized since we initially found that
59     our links were returned in alphabetical order
60     in the event of a tie between pagerank scores.
61     """
62     try:
63         page = wiki_wiki.page(page_title)
64
65         # Collect all valid links
66         all_links = []
67         for link_title in page.links.keys():
68             if link_title and not link_title.startswith('Category:') and not
69                 link_title.startswith('File:') and not
70                 link_title.startswith('Template:'):
71                 all_links.append(link_title)
72
73         # Randomize to avoid uniformity, then returning up until the maximum
74         # amount of links
75         random.shuffle(all_links)
76         return all_links[:max_links]
77     except Exception as e:
78         print(f"Error getting links")
79         return []

```

These were our methods for data collection. Instead of a direct match with a Wikipedia article, we used the search feature of the Wikipedia API to gather results. Then, we scraped the most similar result from the platform. This process was done in order to circumvent the need to match our seeds to the exact naming convention that the articles use. Furthermore, we filtered out pages starting with “Category:”, “File:”, or “Template:”, since these are metadata or auxiliary pages. We also used Wikipedia API’s “links” feature to extract all the outgoing Wiki links from the article.



```

1  def build_wikipedia_graph(seed_terms, max_depth=2, max_links_per_page=20,
2      max_total_pages=200):
3      """
4          Build a directed graph from Wikipedia pages starting with seed terms. (u,
5          w) and (w, u) valid edges in this graph.
6          - seed_terms: List of search terms
7          - max_depth: How many levels of links to follow
8          - max_links_per_page: Maximum links to extract per page
9          - max_total_pages: Maximum total pages to include in graph
10         """
11         G = nx.DiGraph()
12         # BFS Mark set to check membership easily
13         visited = set()
14         seed_pages = []
15
16         for term in seed_terms:
17             page_title = search_wikipedia_first_result(term)
18             if page_title:
19                 print(f"Found: '{term}' -> '{page_title}'")
20                 seed_pages.append(page_title)
21                 visited.add(page_title)
22                 G.add_node(page_title, is_seed=True)
23             else:
24                 print(f"No result found for '{term}'")
25
26         # Running BFS
27         queue = deque([(page, 0) for page in seed_pages]) # (page_title, depth)
28
29         while queue and len(visited) < max_total_pages:
30             current_page, depth = queue.popleft()
31
32             if depth >= max_depth:
33                 continue
34
35             # Stop if we have reached the maximum number of pages
36             if len(visited) >= max_total_pages:
37                 break
38
39             # Get links from current page - use max_links_per_page parameter
40             links = get_wiki_links(current_page, max_links=max_links_per_page)
41
42             for link in links:
43                 # Stop adding nodes if we've reached the limit
44                 if len(visited) >= max_total_pages:
45                     break
46
47                 if link != current_page:
48                     if link not in visited and len(visited) < max_total_pages:
49                         # Can add new node if the link is not already in visited
50                         visited.add(link)
51                         G.add_node(link, is_seed=False)
52                         G.add_edge(current_page, link)
53
54                         # Add to queue if we haven't reached max depth
55                         if depth + 1 < max_depth:
56                             queue.append((link, depth + 1))
57
58         time.sleep(0.1) # To avoid timeouts from excessive API calls

```

```

57
58     return G

```

We used the NetworkX module to construct our directed graphs of Wikipedia articles. Our graph is represented as  $G = (V, E)$ , where  $V$  is the set of vertices (Wikipedia articles) and  $E$  is the set of edges  $(v, w)$ , which represents article  $v$  linking to article  $w$ . Note that  $(v, w) \neq (w, v)$ , since one Wikipedia article does not necessarily have to link to its source. We run Breadth First Search, starting with our seed pages, to populate our graph for PageRank and HITS.

### 3.4 Code and Algorithm Design: Personalized PageRank

```

1  def personalized_pagerank(G, seed_pages, damping=0.85, max_iter=100, tol=1e-6):
2      """
3      Runs Personalized PageRank
4
5      Using power iteration to calculate the Pagerank vector, within a tolerance of
6      1e-6.
7
8      Only addition to Power iteration is that we restart at seed pages, making
9      algorithm focus on nodes relevant to the seed topics.
10     """
11
12     nodes = list(G.nodes())
13     n = len(nodes)
14
15     if n == 0:
16         return {}
17
18     # Create node-to-index mapping, since networkX doesn't enumerate nodes
19     node_to_idx = {node: i for i, node in enumerate(nodes)}
20     idx_to_node = {i: node for i, node in enumerate(nodes)}
21
22     # Build adjacency matrix A (for directed graph)
23     # A[i,j] = 1 if there's a directed edge FROM node i TO node j
24     A = np.zeros((n, n))
25     for i, node in enumerate(nodes):
26         # For directed graph: check outgoing edges (successors)
27         for neighbor in G.successors(node):
28             j = node_to_idx[neighbor]
29             A[i, j] = 1.0 # Edge from node i to node j
30
31     # Making markov matrix out of adjacency matrix A
32     # Compute column sums (out-degrees)
33     col_sums = np.sum(A, axis=0)
34
35     # Build transition matrix M: handle dangling nodes (columns with sum 0)
36     # For columns with sum > 0: M[i,j] = A[i,j] / col_sum[j]
37     # For columns with sum = 0 (dangling nodes): M[i,j] = 1/n (uniform
38     # distribution)
39     M = np.zeros((n, n))
40     for j in range(n):
41         if col_sums[j] > 0:
42             M[:, j] = A[:, j] / col_sums[j]
43         else:
44             # Dangling node: replace with uniform distribution 1/n
45             M[:, j] = 1.0 / n

```

```

44     # Create personalization vector v (teleportation vector)
45     v = np.zeros(n)
46     seed_weight = 1.0 / len(seed_pages)
47     for node in nodes:
48         if node in seed_pages:
49             # Giving equal weight to seed pages, nonzero only for seeds
50             v[node_to_idx[node]] = seed_weight
51             # Normalize to ensure it sums to exactly 1
52     v = v / np.sum(v) if np.sum(v) > 0 else np.ones(n) / n
53
54     # Initialize PageRank vector
55     pr = np.ones(n) / n
56
57     # Power method
58     for iteration in range(max_iter):
59         # Compute new PageRank
60         pr_new = (1 - damping) * v + damping * (M @ pr)
61
62         # Check if it is within tolerance
63         diff = np.sum(np.abs(pr_new - pr))
64         pr = pr_new
65
66
67         if diff < tol:
68             break
69
70     # Normalize entries
71     pr = pr / np.sum(pr) if np.sum(pr) > 0 else pr
72
73     # Convert back to dictionary
74     pagerank_scores = {idx_to_node[i]: float(pr[i]) for i in range(n)}
75
76     return pagerank_scores

```

For calculation of the PageRank scores for each article, we first need to find the transition matrix  $M$ . This was done by creating an adjacency matrix  $A$  from our directed graph, and then normalizing by the sum of each column such that  $M_{ij} = \frac{A_{ij}}{\sum_{j=0} A_{ij}}$ . We also deal with the problem of dangling nodes by replacing them with  $\frac{1}{n}$  to ensure a random jump from the side of the web-surfer when faced with a Wiki article with no outgoing links, which ensures that we can find a unique steady state vector  $w_q$  when applying our PageRank algorithm.

We are interested in articles related to our seed articles, therefore, we replace  $\mathbf{1}$ , which is represented as

$$\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^n, \text{ with a personalization vector } v, \text{ where } v = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ such that } x_i = 0 \text{ if } x_i$$

does not represent the long run probability for one of our selected seed-pages. We assign a non-zero probability for our seed-pages, which we encoded as  $\frac{1}{n}$ .

As for performing the Personalized PageRank on this directed graph, we used Power Iteration, which was a choice we made due to the scalability of the algorithm. Let  $p^{t+1}$  be the  $t+1^{\text{th}}$  iteration of our process to find the steady-state. Then,  $p^{t+1} = (1 - d)v + dMp^t$ , where  $d$  is the damping factor. We set our tolerance  $\epsilon = 10^{-6}$  to be the degree at which we consider our PageRank scores

to have stabilized [1]. A corresponding probability from this vector  $p$  will then be the PageRank score for its corresponding webpage.

```
1  def visualize_graph(G, pagerank_scores, seed_pages, top_n=10):
2      '''
3          Visualize the graph with node sizes and colors based on PageRank scores.
4      '''
5      if G.number_of_nodes() == 0:
6          print("Graph is empty, cannot visualize.")
7          return
8
9      plt.figure(figsize=(16, 12))
10
11     # Get top influential nodes - sort by score (descending)
12     def sort_key(item):
13         node, score = item
14         # tie breaker for same scores
15         return (-score, hash(node) % 1000000)
16
17     sorted_nodes = sorted(pagerank_scores.items(), key=sort_key)
18     top_nodes = [node for node, _ in sorted_nodes[:top_n]]
19
20     node_sizes = []
21     node_colors = []
22     node_labels = {}
23
24     for node in G.nodes():
25         # Node size based on PageRank score
26         size = pagerank_scores.get(node, 0) * 10000
27         node_sizes.append(max(size, 50)) # Minimum size
28
29         # Color: red for seeds, orange for top influential, blue for others
30         if node in seed_pages:
31             node_colors.append('red')
32         elif node in top_nodes:
33             node_colors.append('orange')
34         else:
35             node_colors.append('lightblue')
36
37         # Only label seed pages and top influential nodes
38         if node in seed_pages or node in top_nodes:
39             # Truncate long titles
40             label = node[:30] + '...' if len(node) > 30 else node
41             node_labels[node] = label
42
43     # Use spring layout for better visualization
44     pos = nx.spring_layout(G, k=1, iterations=50, seed=42)
45
46     # Draw the graph
47     nx.draw_networkx_nodes(G, pos, node_size=node_sizes, node_color=node_colors,
48                            alpha=0.7)
49     nx.draw_networkx_edges(G, pos, alpha=0.2, width=0.5, arrows=True,
50                            arrowsize=10, arrowstyle='->')
51     nx.draw_networkx_labels(G, pos, node_labels, font_size=8, font_weight='bold')
52
53     plt.title('Wikipedia Graph with Personalized PageRank\n(Red=Seed, Orange=Top
54               Influential, Blue=Others)',
55              fontsize=14, fontweight='bold')
56     plt.axis('off')
```

```

54 plt.tight_layout()
55 plt.show()
56
57 # Print top influential nodes
58 print(f"\nTop {top_n} Most Influential Pages (by Personalized PageRank):")
59 print("-" * 70)
60 prev_score = None
61 for i, (node, score) in enumerate(sorted_nodes[:top_n], 1):
62     # Show if score is same as previous (tied)
63     print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
64     prev_score = score

```

We proceeded with visualizing the graph in terms of influence. Seed pages are given red-coloring, pages with significant PageRank scores are given orange coloring, and every other page is given a blue coloring.

### 3.5 Code and Algorithm Design: Hyperlink-Induced Topic Search

```

1 def hits(G, max_iter=100, tol=1e-6):
2     """
3     Runs Hyperlink-Induced Topic Search (HITS)
4     Updates authority and hub vectors to convergence until tolerance of 1e-6 is
5     reached:
6     a_{k+1} = A^T * h_k
7     h_{k+1} = A * a_{k+1}
8     a_{k+1} = a_{k+1} / ||a_{k+1}||_2
9     h_{k+1} = h_{k+1} / ||h_{k+1}||_2
10    """
11    nodes = list(G.nodes())
12    n = len(nodes)
13
14    if n == 0:
15        return {}, {}
16
17    # Create node-to-index mapping, since networkX doesn't enumerate nodes
18    node_to_idx = {node: i for i, node in enumerate(nodes)}
19    idx_to_node = {i: node for i, node in enumerate(nodes)}
20
21    # Build adjacency matrix A (for directed graph)
22    # A[i,j] = 1 if there's a directed edge FROM node i TO node j
23    A = np.zeros((n, n))
24    for i, node in enumerate(nodes):
25        # For directed graph: check outgoing edges (successors)
26        for neighbor in G.neighbors(node):
27            j = node_to_idx[neighbor]
28            A[i, j] = 1.0 # Edge from node i to node j
29
30    # Initialize authority vector and hub vector
31    hub = np.random.rand(n) * 1e-6
32    auth = np.random.rand(n) * 1e-6
33
34    # Iterative updates of a and h
35    for iteration in range(max_iter):
36        # Compute new authority vector: a_{k+1} = A^T * h_k
37        new_auth = A.T @ hub
38        # Compute new hub vector: h_{k+1} = A * a_{k+1}
39        new_hub = A @ new_auth

```

```

39
40     # Normalize authority and hub vectors
41     a_norm = np.linalg.norm(new_auth, 2)
42     h_norm = np.linalg.norm(new_hub, 2)
43     if a_norm == 0 or h_norm == 0:
44         break
45     new_auth = new_auth / a_norm
46     new_hub = new_hub / h_norm
47
48     # Check if it is within the tolerance
49     diff_a = np.linalg.norm(new_auth - auth, 2)
50     diff_h = np.linalg.norm(new_hub - hub, 2)
51     auth = new_auth
52     hub = new_hub
53     if diff_a < tol and diff_h < tol:
54         break
55
56     # Convert back to dictionary
57     authority_scores = {idx_to_node[i]: float(auth[i]) for i in range(n)}
58     hub_scores = {idx_to_node[i]: float(hub[i]) for i in range(n)}
59
60     return authority_scores, hub_scores

```

For each article, we calculated the authority and hub scores through an iterative process. We did so by first constructing an adjacency matrix  $A$  from the directed graph using the following formula:

$$A_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to page } j \\ 0 & \text{otherwise} \end{cases}$$

Once the matrix is transformed, we initialized the authority and hub vectors. During each update, the authority vector is reassigned to  $\mathbf{a}_{k+1} = A^T \mathbf{h}_k$  and the hub vector is reassigned to  $\mathbf{h}_{k+1} = A \mathbf{a}_{k+1}$  at step  $k$ . Then, we L2 normalize  $\mathbf{a}_{k+1}$  and  $\mathbf{h}_{k+1}$ , which involves scaling each vector so that the sum of squares of all entries adds to 1. Note that this means the sum of all authority scores and the sum of all hub scores may be greater than 1, since squaring a value between 0 and 1 may be less than the value itself. The normalization inside each loop is important, since without it the values in both vectors would grow without bound. The loop continues until the difference between consecutive updates becomes smaller than the tolerance  $\epsilon = 10^{-6}$ .

After convergence, the greatest entries of the authority vector identify pages supported by reliable hubs, and the greatest entries of the hub vector identify pages that direct users to strong authorities. The output gives us two rankings tied to the seed set of articles, allowing us to compare how the Wikipedia network's structure differs from the long-run behavior captured by Personalized PageRank.

```

1     def visualize_hits(G, authority_scores, hub_scores, seed_pages, top_n=10):
2         '''
3         Visualize the graph with node sizes and colors based on HITS authority and
4         hub scores.
5         '''
6         if G.number_of_nodes() == 0:
7             return
8
9         # Get top influential nodes - sort by score (descending)
10        def sort_key(item):

```

```

10     node, score = item
11     # tie breaker for same scores
12     return (-score, hash(node) % 1000000)
13
14 sorted_auth = sorted(authority_scores.items(), key=sort_key)
15 sorted_hub = sorted(hub_scores.items(), key=sort_key)
16 top_auth = [node for node, _ in sorted_auth[:top_n]]
17 top_hub = [node for node, _ in sorted_hub[:top_n]]
18
19 # Use spring layout for better visualization
20 pos = nx.spring_layout(G, k=1, iterations=50, seed=42)
21
22 # Function to draw the graph for each type of score
23 def draw_graph(title, score_map, top_nodes):
24     plt.figure(figsize=(16, 12))
25
26     node_sizes = []
27     node_colors = []
28     node_labels = {}
29
30     for node in G.nodes():
31         # Node size based on HITS authority or hub score
32         size = score_map.get(node, 0) * 10000
33         node_sizes.append(max(size, 50)) # Minimum size
34
35         # Color: red for seeds, orange for top influential, blue for others
36         if node in seed_pages:
37             node_colors.append('red')
38         elif node in top_nodes:
39             node_colors.append('orange')
40         else:
41             node_colors.append('lightblue')
42
43         # Only label seed pages and top influential nodes
44         if node in seed_pages or node in top_nodes:
45             # Truncate long titles
46             lbl = node[:30] + '...' if len(node) > 30 else node
47             node_labels[node] = lbl
48
49     # Draw the graph
50     nx.draw_networkx_nodes(G, pos, node_size=node_sizes,
51                             node_color=node_colors, alpha=0.7)
52     nx.draw_networkx_edges(G, pos, alpha=0.2, width=0.5, arrows=True,
53                             arrowsize=10, arrowstyle='->')
54     nx.draw_networkx_labels(G, pos, node_labels, font_size=8,
55                             font_weight='bold')
56     plt.title(title, fontsize=14, fontweight='bold')
57     plt.axis('off')
58     plt.tight_layout()
59     plt.show()
60
61 # Draws graph for authority and hub scores
62 draw_graph("Wikipedia Graph with HITS Authority Scores\n(Red=Seed, Orange=Top
63             Influential, Blue=Others)", authority_scores, top_auth)
64 draw_graph("Wikipedia Graph with HITS Hub Scores\n(Red=Seed, Orange=Top
65             Influential, Blue=Others)", hub_scores, top_hub)
66
67 # Print top influential nodes for authorities
68 print(f"\nTop {top_n} Most Influential Pages (by HITS Authority Scores):")

```

```

64     print("-" * 70)
65     prev_score = None
66     for i, (node, score) in enumerate(sorted_auth[:top_n], 1):
67         # Show if score is same as previous (tied)
68         print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
69         prev_score = score
70
71     # Print top influential nodes for hubs
72     print(f"\nTop {top_n} Most Influential Pages (by HITS Hub Scores):")
73     print("-" * 70)
74     prev_score = None
75     for i, (node, score) in enumerate(sorted_hub[:top_n], 1):
76         # Show if score is same as previous (tied)
77         print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
78         prev_score = score

```

Similar to the graph visualization for PageRank, we visualized the graphs for HITS in terms of influence for both authority scores and hub scores, but separately. Like before, seed pages are given red-coloring, pages with significant scores are given orange coloring, and every other page is given a blue coloring.

## 4 Results and Discussion

We tested our implementation on few local Wikipedia networks with varying sizes and compared results across Personalized PageRank and HITS.

### 4.1 Small Wikipedia Network

```

1     # Input params
2     seed_terms = [
3         "planets",
4         "supernova",
5         "hubble"
6     ]
7
8     # Build the graph
9     G = build_wikipedia_graph(
10         seed_terms,
11         max_depth=2,
12         max_links_per_page=2,
13         max_total_pages=200
14     )
15
16     # Initialize pages for seed articles
17     seed_pages = [node for node, attrs in G.nodes(data=True) if attrs.get('is_seed',
18         False)]

```

First, we tested our Personalized PageRank and HITS algorithms with the set of seed terms  $S = \{\text{"planets", "supernova", "hubble"}\}$ , while setting `max_links` to 2 and `max_links_per_page` to 2 in order to observe algorithmic behavior closely. The `build_wikipedia_graph` function identified the top Wikipedia page search results, matching the seed terms to the pages “Star,” “Planet,” “Supernova,” and “Hubble Space Telescope,” respectively. Then, we ran the algorithms using the functions `personalized_pagerank` and `hits`.



The vectors representing Personalized PageRank scores, HITS authority scores, and HITS hub scores, respectively, are shown below:

0.16930952	0.36873486	0.52146983
0.16930952	0.36873486	0.50747636
0.13846833	0.35883998	0.48380106
0.05762713	0.35883998	0.29250832
0.05762713	0.34209901	0.28786906
0.05762713	0.34209901	0.18861364
0.05762713	0.20683462	0.17741506
0.05762713	0.20683462	0.03150495
0.02134336	0.20355417	0.00000000
0.02134336	0.20355417	0.00000000
0.02134336	0.13336998	0.00000000
0.02134336	0.13336998	0.00000000
0.02134336	0.12545139	0.00000000
0.02134336	0.12545139	0.00000000
0.02134336	0.02227737	0.00000000
0.02134336	0.02227737	0.00000000
0.02134336	0.00000000	0.00000000
0.02134336	0.00000000	0.00000000
0.02134336	0.00000000	0.00000000

Furthermore, the raw outputs below show the scores corresponding with article names in the network.

#### Personalized PageRank:

1. Hubble Space Telescope (Score: 0.16930952)
2. Planet (Score: 0.16930952)
3. Supernova (Score: 0.13846833)
4. Super Jupiter (Score: 0.05762713)
5. Kepler's laws of planetary motion (Score: 0.05762713)
6. Space Flyer Unit (Score: 0.05762713)
7. Broad Band X ray Telescope (Score: 0.05762713)
8. Guest star (astronomy) (Score: 0.05762713)
9. List of most massive stars (Score: 0.02134336)
10. Kosmos 2312 (Score: 0.02134336)
11. Vis viva equation (Score: 0.02134336)
12. Newton's laws of motion (Score: 0.02134336)
13. Solar Orbiter (Score: 0.02134336)
14. TENKO 100 (Score: 0.02134336)
15. List of proper names of exoplanets (Score: 0.02134336)
16. JAXA (Score: 0.02134336)
17. SPHEREx (Score: 0.02134336)
18. History of supernova observation (Score: 0.02134336)
19. Planetesimal (Score: 0.02134336)

#### HITS authority scores:

1. Space Flyer Unit (Score: 0.36873486)
2. Broad Band X ray Telescope (Score: 0.36873486)
3. Super Jupiter (Score: 0.35883998)
4. Kepler's laws of planetary motion (Score: 0.35883998)
5. TENKO 100 (Score: 0.34209901)
6. Guest star (astronomy) (Score: 0.34209901)
7. Vis viva equation (Score: 0.20683462)
8. Newton's laws of motion (Score: 0.20683462)
9. List of proper names of exoplanets (Score: 0.20355417)
10. Planetesimal (Score: 0.20355417)
11. Kosmos 2312 (Score: 0.13336998)
12. JAXA (Score: 0.13336998)
13. Solar Orbiter (Score: 0.12545139)
14. SPHEREx (Score: 0.12545139)
15. List of most massive stars (Score: 0.02227737)
16. History of supernova observation (Score: 0.02227737)
17. Hubble Space Telescope (Score: 0.00000000)
18. Supernova (Score: 0.00000000)
19. Planet (Score: 0.00000000)

#### HITS hub scores:

1. Hubble Space Telescope (Score: 0.52146983)
2. Planet (Score: 0.50747636)
3. Supernova (Score: 0.48380106)
4. Kepler's laws of planetary motion (Score: 0.29250832)
5. Super Jupiter (Score: 0.28786906)
6. Space Flyer Unit (Score: 0.18861364)
7. Broad Band X ray Telescope (Score: 0.17741506)
8. Guest star (astronomy) (Score: 0.03150495)
9. List of most massive stars (Score: 0.00000000)
10. Kosmos 2312 (Score: 0.00000000)
11. Vis viva equation (Score: 0.00000000)
12. Newton's laws of motion (Score: 0.00000000)
13. Solar Orbiter (Score: 0.00000000)
14. TENKO 100 (Score: 0.00000000)
15. List of proper names of exoplanets (Score: 0.00000000)
16. JAXA (Score: 0.00000000)
17. SPHEREx (Score: 0.00000000)
18. History of supernova observation (Score: 0.00000000)
19. Planetesimal (Score: 0.00000000)

We found that Personalized PageRank emphasizes pages directly related to the seed terms, with “Hubble Space Telescope” (0.16930952), “Planet” (0.16930952), and “Supernova” (0.13846833) at the top due to their centrality across multiple related pages. This aligns with the mathematical intuition that PageRank, especially Personalized PageRank, amplifies nodes with strong link connections back to the seeds.

On the other hand, HITS authority scores lacked seeds as the most influential pages since many Wikipedia articles do not link bidirectionally. In fact, all three of the seed terms had an authority score of 0.00000000, emphasizing that there were no articles that linked back to them. There were a few tied authority scores, which makes sense because they likely had the same amount of articles linking to them. For example, if the page “Supernova” linked to “History of supernova observation,” but nothing else linked to this page, it received a score of 0.02227737. This same logic could be applied to “List of most massive stars,” and several other pages too. Furthermore, some pages like “Space Flyer Unit” (0.36873486) were highly ranked because many hub pages linked to it.

Hub scores revealed a different pattern. Seed terms like “Hubble Space Telescope”, “Planet,” and “Supernova” appeared as some of the most influential, because naturally, these pages point to the rest of the articles in the network. Thus, they have high hub scores. In this way, we noticed similarities between HITS hub scores and Personalized PageRank, because both emphasized the seed articles as the most important in the network. However, Personalized PageRank and HITS authority scores differed, as authority showed articles that were linked to the most in the network. Notably, some hub scores had a value of 0.00000000, because these were dangling nodes, lacking outgoing links to authoritative nodes. Since we set our `max_depth` to 2, it makes sense that some articles did not link to others.

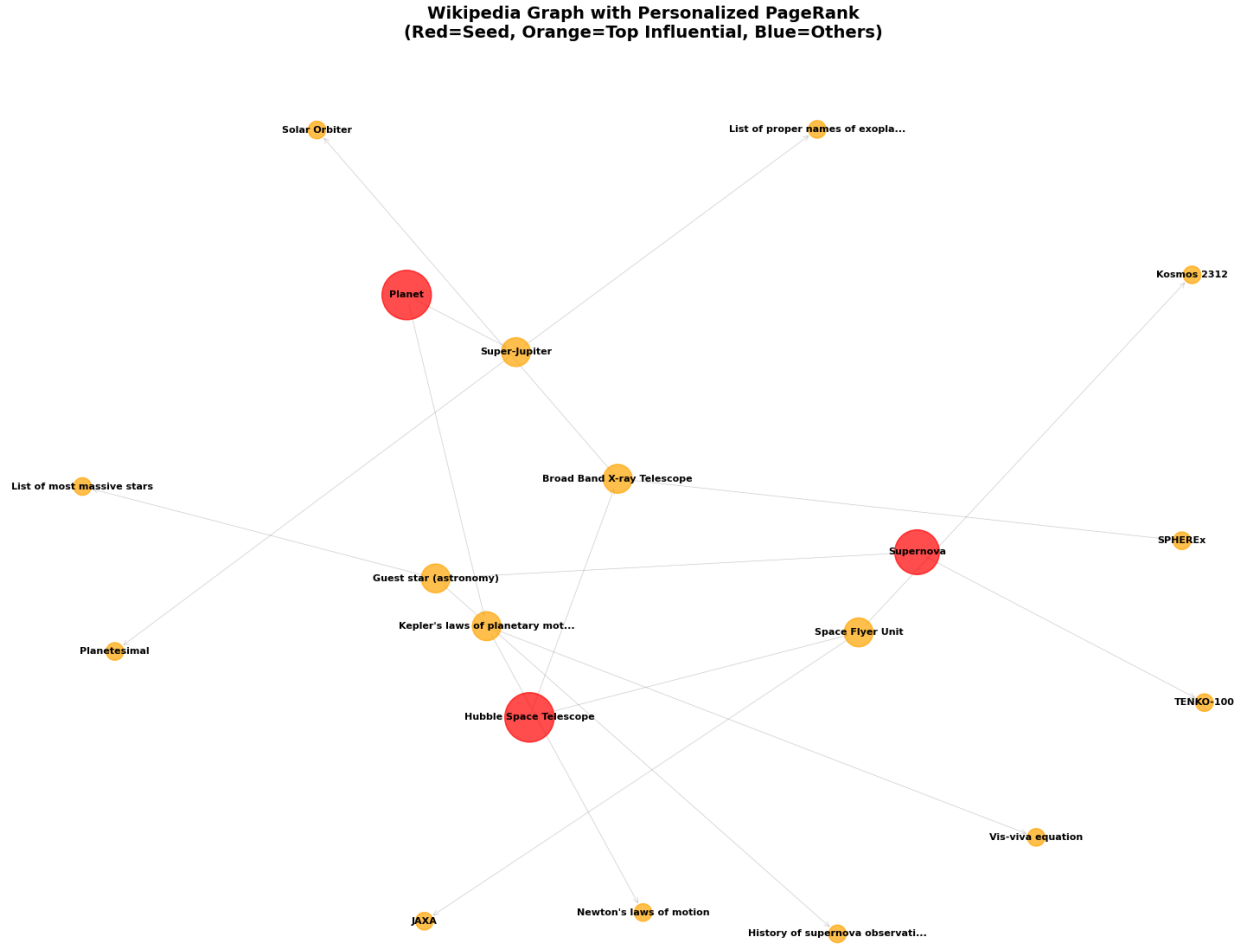


Figure 1: Personalized PageRank visualization for small graph

Figure 1 shows the Wikipedia graph network using Personalized PageRank, with red circles indicating the seed articles, orange indicating top influential, and blue indicating other articles. The proximity of top influential articles to the seeds highlights that Personalized PageRank effectively prioritizes pages that are both strongly linked to the seeds and widely referenced by other nodes. The numerous blue nodes scattered around the edges correspond to pages with low influence scores. Their peripheral placement emphasizes the steep decay in PageRank values, where a few highly central pages dominate the network, while the majority have minimal influence. Thus, the graph concurs with our mathematical results.

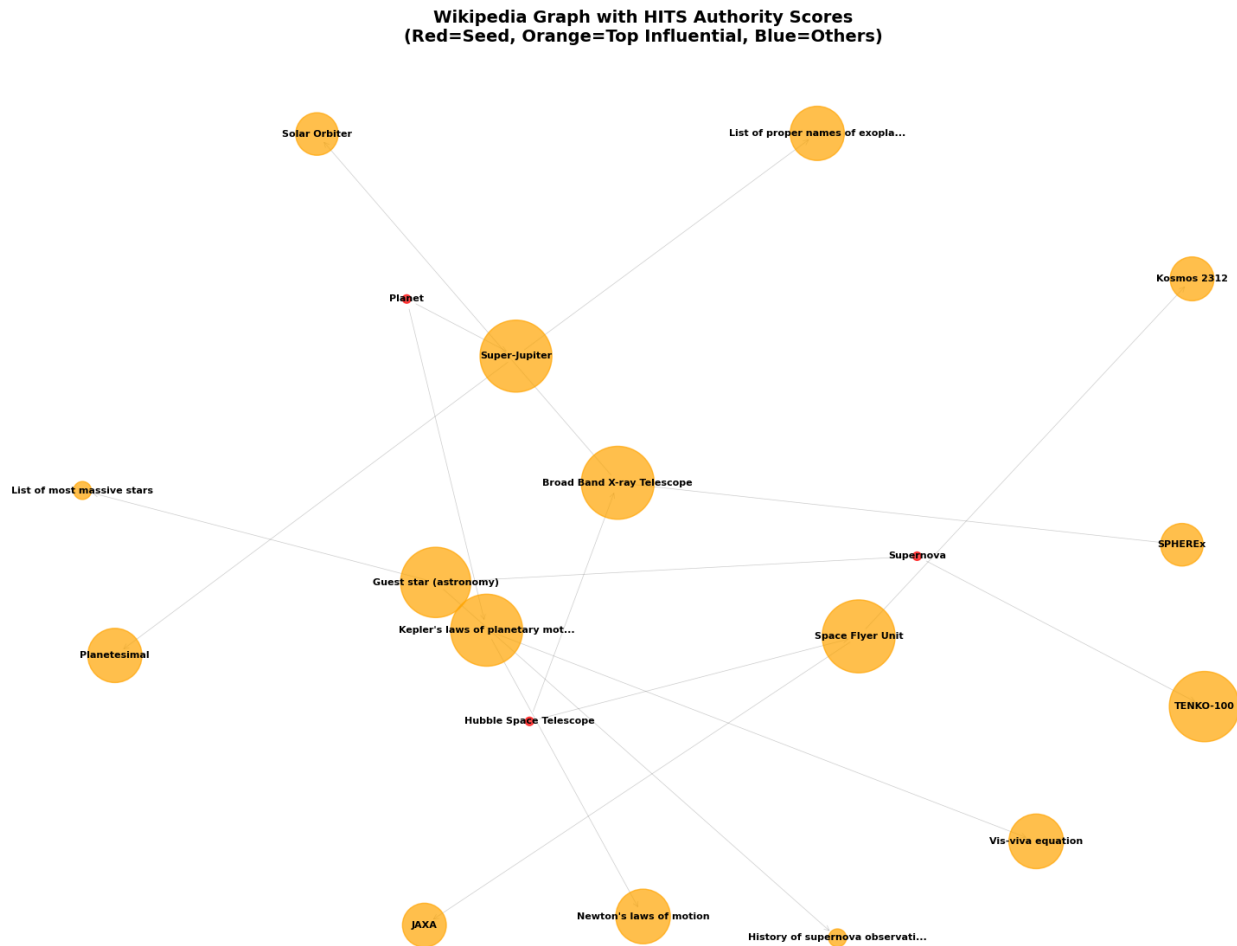


Figure 2: HITS authorities visualization for small graph

Figure 2 illustrates the HITS authority scores for the small Wikipedia network. The seed nodes, represented in red, were significantly smaller than the nodes surrounding them, indicating a lower authority score. As noted before, this makes sense because the seed nodes act as the base, linking to other articles, but since many Wikipedia pages are not bidirectional, they are not necessarily high authority. On the other hand, several peripheral nodes showed high authority since many Wikipedia pages linked to them.

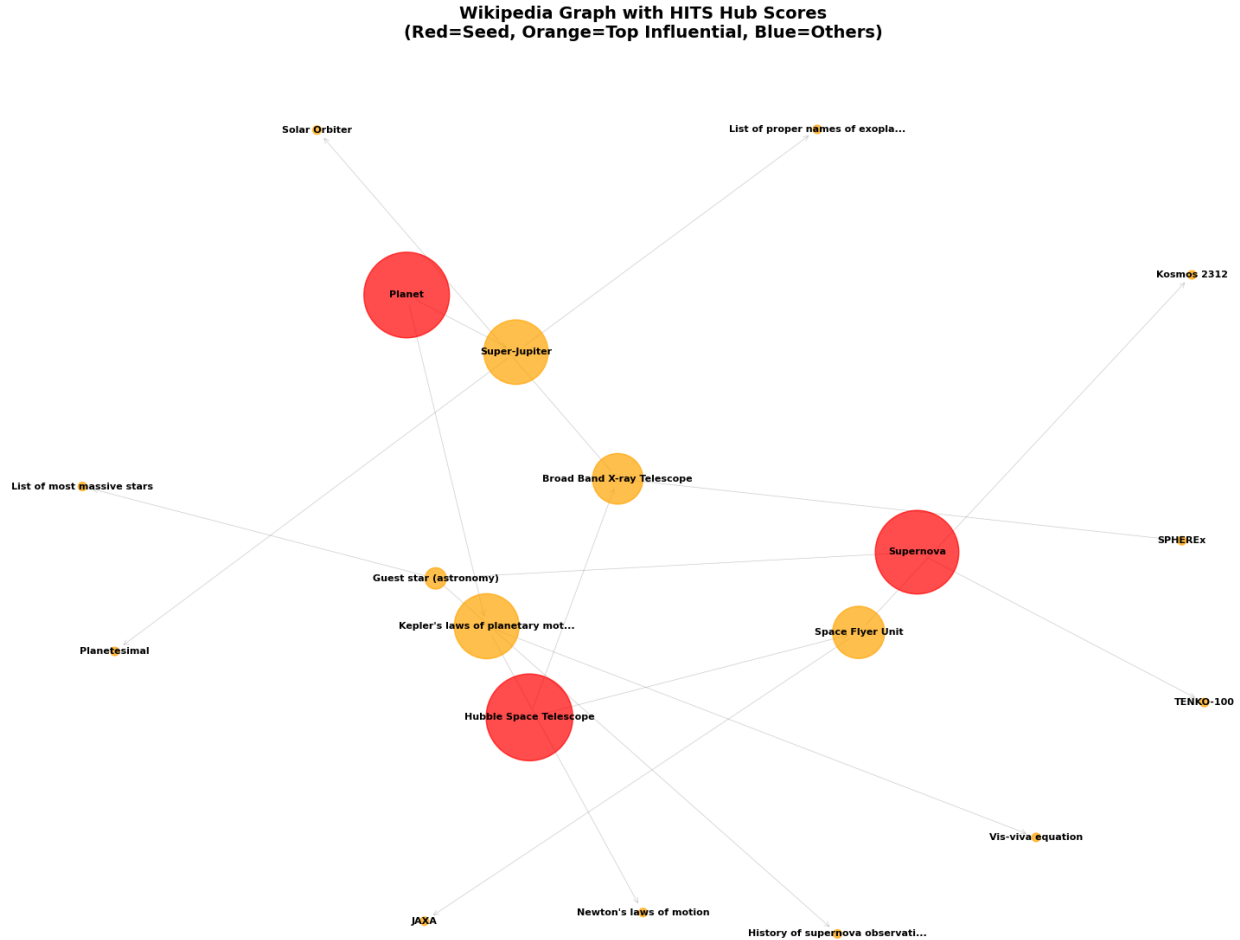


Figure 3: HITS hubs visualization for small graph

Figure 3 shows the HITS hub scores. In contrast to the authority scores, seeds had the highest hub scores because they linked to the most authority pages. This result makes sense because seeds served as the baseline for constructing the graph network. Similar to Personalized PageRank, the peripheral nodes had low scores, because they did not link to many other pages in the graph.

Overall, Personalized PageRank produced sharper concentration around the seed pages, with all three seeds dominating the ranking due to the sparse and highly interconnected structure. HITS produced more varied authority and hub lists, driven by sensitivity of the algorithm to a small network.

## 4.2 Extension: Medium-Sized Wikipedia Network

```

1  # Input params
2  seed_terms = [
3      "linear algebra",
4      "eigenvalues",
5      "dot product"

```

```

6 ]
7
8 # Build the graph
9 G = build_wikipedia_graph(
10     seed_terms,
11     max_depth=2,
12     max_links_per_page=5,
13     max_total_pages=200
14 )
15
16 # Initialize pages for seed articles
17 seed_pages = [node for node, attrs in G.nodes(data=True) if attrs.get('is_seed',
    False)]

```

Next, we ran our implementation on a medium-sized network, with a Wikipedia structure containing 5 max\_links and 5 max\_links\_per\_page as parameters. Here, we used the set of seed terms  $S = \{\text{“linear algebra”, “eigenvalues”, “dot product”}\}$ .

Personalized PageRank found the top 30 most influential pages to be:

1. Dot product (Score: 0.13858620)
2. Linear algebra (Score: 0.13858620)
3. Eigenvalues and eigenvectors (Score: 0.13858620)
4. Sparse matrix (Score: 0.01994972)
5. Van der Waerden notation (Score: 0.01994972)
- ⋮
28. Moon (Score: 0.00379994)
29. Absolutely convex set (Score: 0.00379994)
30. Finite set (Score: 0.00379994)

Also, HITS found these 30 articles to have the highest authority scores:

1. Bar (unit) (Score: 0.18171734)
2. Angular momentum (Score: 0.18171734)
3. Multilinear algebra (Score: 0.18171734)
4. Symmetry (Score: 0.18171734)
5. Vector (geometry) (Score: 0.18171734)
- ⋮
28. Algebraic variety (Score: 0.13018781)
29. Parallelogram (Score: 0.13018781)
30. Quantum mechanics (Score: 0.13018781)

Then, HITS found these 30 articles to have the highest hub scores:

1. Cauchy stress tensor (Score: 0.40633233)
2. Arthur Cayley (Score: 0.40071487)
3. Graph theory (Score: 0.34656665)
4. Eigenvalues and eigenvectors (Score: 0.31498772)
5. Hecke eigensheaf (Score: 0.30715994)

```
28. Moon (Score: 0.00000000)
29. Absolutely convex set (Score: 0.00000000)
30. Finite set (Score: 0.00000000)
```

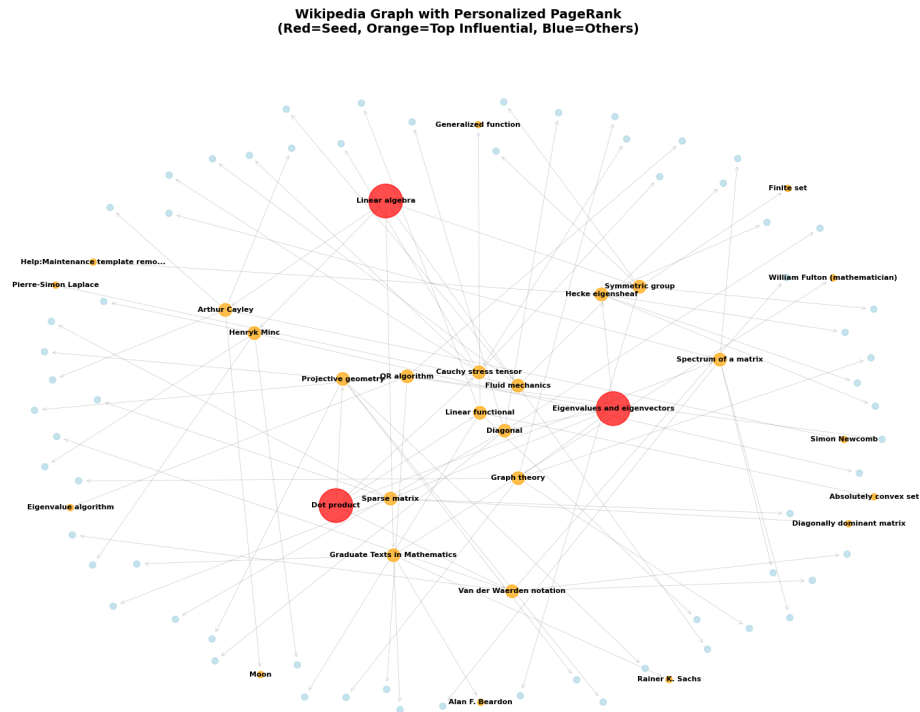


Figure 4: Personalized PageRank visualization for medium-sized graph



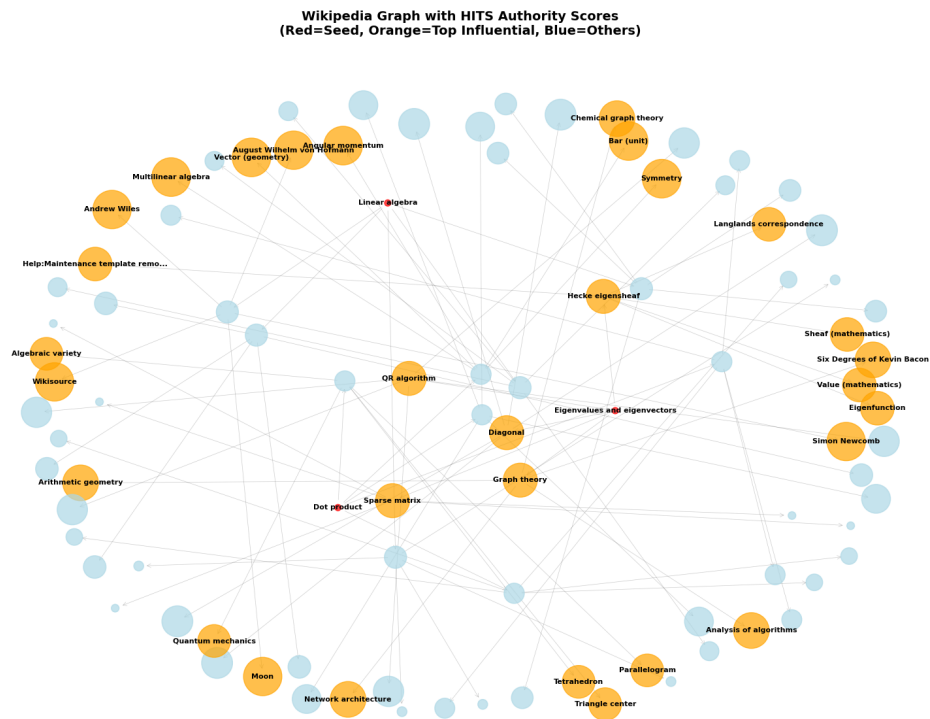


Figure 5: HITS authority visualization for medium-sized graph

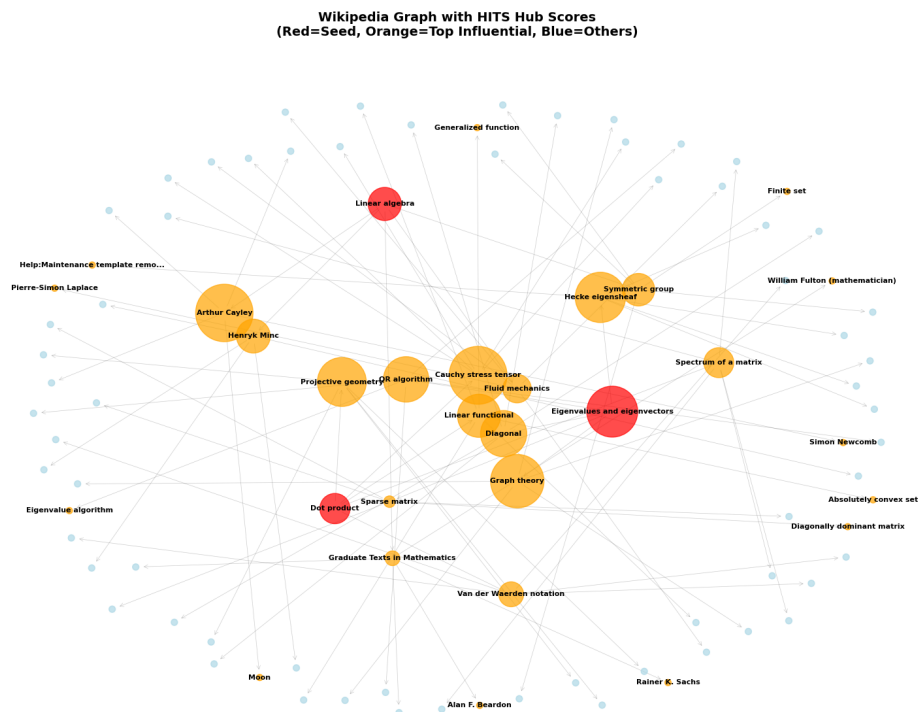


Figure 6: HITS hub visualization for medium-sized graph

For the medium-sized graph, Personalized PageRank continued to prioritize the seeds “Linear algebra,” “Eigenvalues and eigenvectors,” and “Dot product,” localizing strongly around seed-adjacent neighborhoods. Unlike the small network, the effect of larger connectivity becomes more apparent here. Specifically, the decay of PageRank values is sharper for nodes further from the seeds, reflecting the algorithm’s tendency to favor nodes within dense clusters, as shown in Figure 4. Even pages like “Sparse matrix” and “Van der Waerden notation” appear with moderate scores, showing how Personalized PageRank can propagate influence. Figure 4 confirms this pattern, with high-scoring nodes tightly clustered near the seeds, and peripheral nodes being loosely connected and less influential.

HITS surfaced with different patterns, however. Authority scores, as seen in Figure 5, prioritize pages that receive links from many other nodes, regardless of proximity to the seeds. This explains why some high-authority pages like “Bar (unit)” and “Angular momentum” appear prominently despite not being seed-adjacent. The medium-sized graph allowed these distinctions to emerge more clearly than in the small network, where limited connectivity often produced ties in scores.

Hub scores, visualized in Figure 6, highlight articles that point to influential authorities, including both seed-related pages and other hubs like “Cauchy stress tensor” or “Graph theory.” The variation in hub scores demonstrates how the algorithm captures influence across the network.

### 4.3 Extension: Large Wikipedia Network

```

1  # Input params
2  seed_terms = [
3      "computer science",
4      "artificial intelligence",
5      "robotics"
6  ]
7
8  # Build the graph
9  G = build_wikipedia_graph(
10     seed_terms,
11     max_depth=5,
12     max_links_per_page=8,
13     max_total_pages=200
14 )
15
16 # Initialize pages for seed articles
17 seed_pages = [node for node, attrs in G.nodes(data=True) if attrs.get('is_seed',
    False)]

```

Finally, we ran our implementation on a large network with `max_links` set to 8 and `max_links_per_page` set to 8 as well. Furthermore, we set `max_depth` to 5. We used the set of seed terms  $S = \{\text{“computer science”, “artificial intelligence”, “robotics”}\}$ .

For Personalized PageRank, the 30 most influential articles were:

1. Computer science (Score: 0.14327787)
2. Artificial intelligence (Score: 0.14203077)
3. Robotics (Score: 0.12083008)
4. Snake (Score: 0.01346351)

5. Wireless sensor network (Score: 0.01346351)

⋮

28. NarrowBand IOT (Score: 0.00172609)

29. Windows 1.0 (Score: 0.00172609)

30. Curt Weldon (Score: 0.00172609)

**HITS authority scores were the greatest for these 30 articles:**

1. Empirical observation (Score: 0.12191533)

2. Completely integrable (Score: 0.12191533)

3. Orbital mechanics (Score: 0.12191533)

4. Electric field (Score: 0.12191533)

5. Lunar space elevator (Score: 0.12191533)

⋮

28. Chimera (Barth novel) (Score: 0.11055698)

29. Invisible Man (Score: 0.11055698)

30. John Casey (novelist) (Score: 0.11055698)

**Also, HITS hub scores were the greatest for the following 30 articles:**

1. Glossary of aerospace engineering (Score: 0.34482863)

2. Computer science (Score: 0.32345325)

3. History of the graphical user interface (Score: 0.31970106)

4. Jonathan Franzen (Score: 0.31270236)

5. Atanasoff-Berry computer (Score: 0.28817081)

⋮

28. NarrowBand IOT (Score: 0.00000000)

29. Windows 1.0 (Score: 0.00000000)

30. Curt Weldon (Score: 0.00000000)

Wikipedia Graph with Personalized PageRank  
(Red=Seed, Orange=Top Influential, Blue=Others)

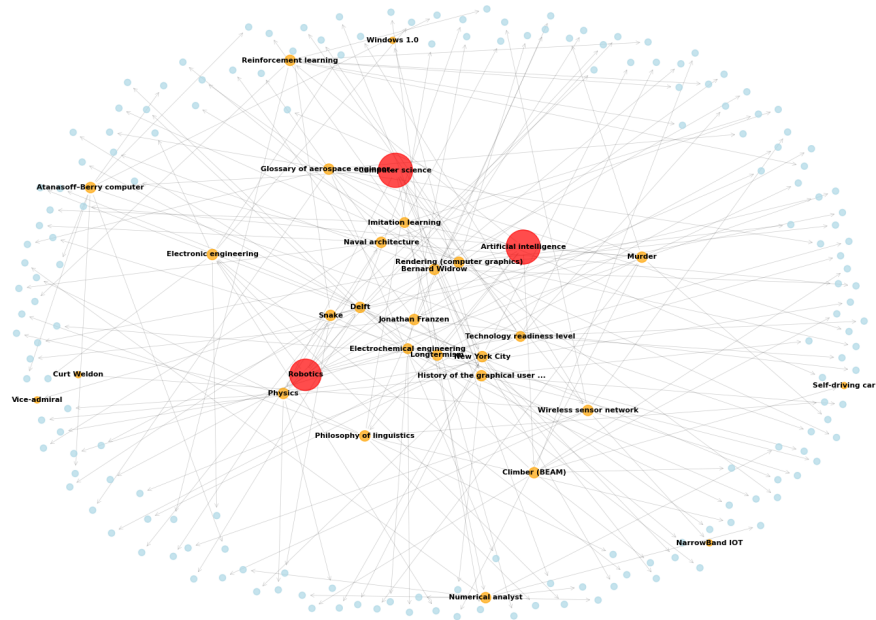


Figure 7: Personalized PageRank visualization for large graph

Wikipedia Graph with HITS Authority Scores  
(Red=Seed, Orange=Top Influential, Blue=Others)

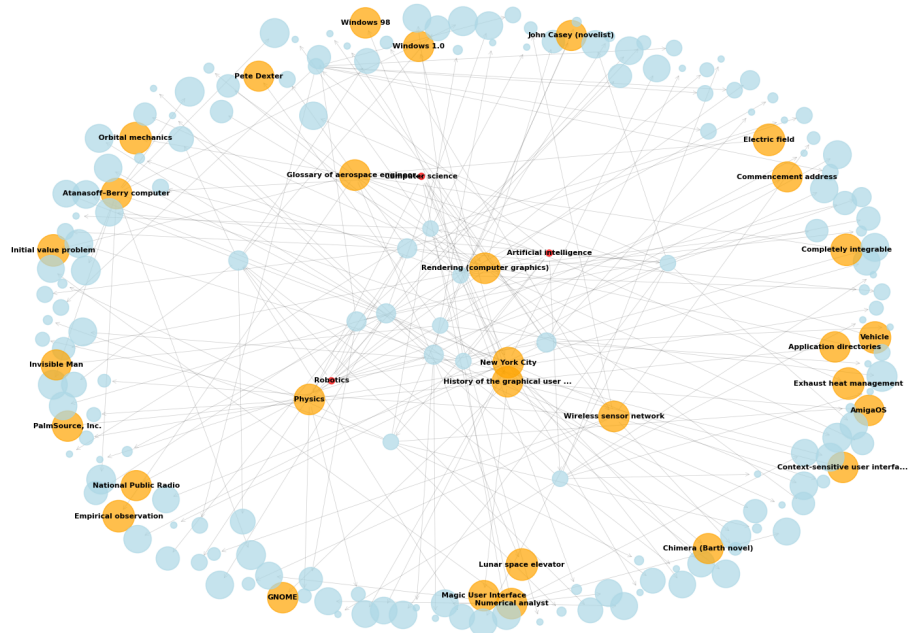


Figure 8: HITS authority visualization for large graph

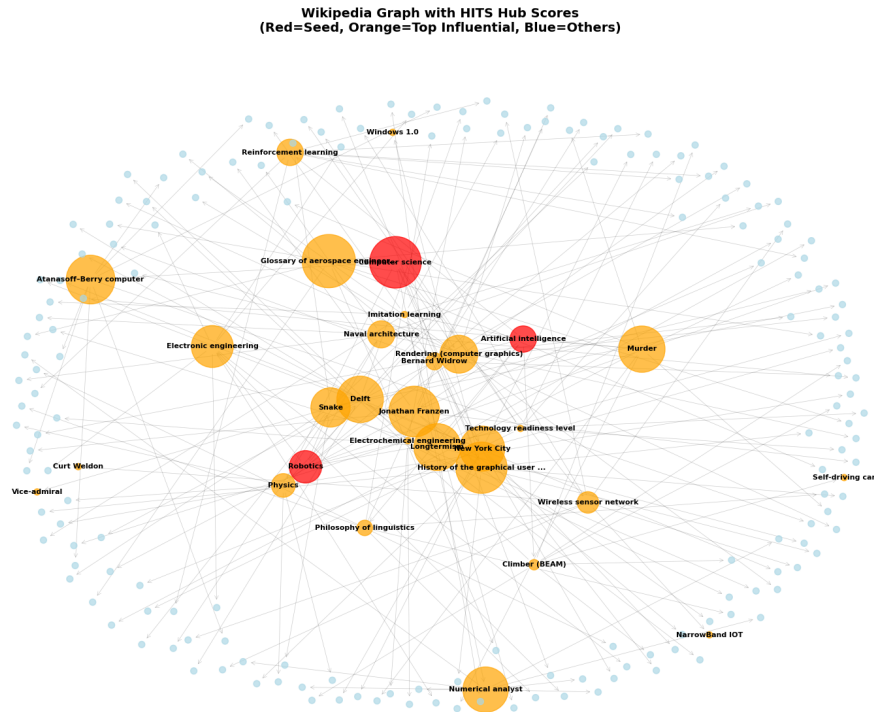


Figure 9: HITS hub visualization for large graph

In the large Wikipedia network, Personalized PageRank continued to prioritize the seed articles. While the seeds still dominated, peripheral nodes such as “Snake” and “Wireless sensor network” received moderate scores, demonstrating that PageRank can propagate influence, as we saw in the medium-sized graph. For the large graph specifically, though, increasing `max_depth` to 5 shows that not all the influential articles directly link to the seeds. Figure 7 confirms this pattern, as the red nodes are central, but there are several clusters of orange nodes that indirectly link to the seeds. However, there were several hubs in the interior of the graph that repeatedly linked to these articles.

HITS authority scores in the large network, shown in Figure 8, reveal more diversity in influential articles. Many high-authority pages like “Empirical observation” and “Orbital mechanics” are not directly connected to the seeds, emphasizing that authority is determined purely by in-links rather than proximity. As the network grows, the separation between authorities and seed-adjacent nodes becomes more pronounced.

Hub scores, depicted in Figure 9, similarly show that seed articles remain highly ranked hubs due to their extensive linking to authority pages.

Overall, we can identify a clear correlation between PageRank and its emphasis on seed proximity, while HITS highlights asymmetric behavior within the subgraph.

## 5 Conclusion

### 5.1 What We Learned

Through this project, we learned how different link-analysis algorithms (Personalized PageRank and HITS) capture importance within real-world Wikipedia networks. Personalized PageRank emphasized pages closely tied to the selected seed topics, demonstrating how biasing in teleportation affects random-walk behavior, while HITS distinguished authoritative pages from strong hubs, often identifying influential non-seed nodes from the asymmetry of Wiki links. We understood how graph and seed selection, and algorithmic design shape the ranking outcomes in large networks like this one.

### 5.2 Limitations

We found that a potential limitation could be the inconsistency of Wikipedia hyperlink structure. Some important articles might be poorly linked, having many outbound links while others have close to none, which could bias the PageRank and HITS scores away from desired influential pages. In particular, PageRank may inflate the importance of seed pages that simply sit in denser local neighborhoods, while HITS may assign low authority scores to pages that are well-known but receive few incoming links. When the seed set  $S$  is small, PageRank will overweight the seed  $i$  with the largest local neighborhood.

Our set of vertices might also underrepresent the true PageRank score of many of these articles since relevant articles might get omitted.

### 5.3 Future Work

We will consider performing optimizations on the PageRank algorithm in order to accomodate for large and dense graphs. Considering that PageRank's power iteration method requires repeated matrix multiplication of the Markov matrix  $M \in M_{n \times n}(\mathbb{R})$ , this process will be computationally expensive when considering large graphs and networks.

A possible extension to this project will be the exploration of low-rank approximations of our stochastic matrix and a smaller graph  $G'$  to iterate on, and using these results to approximate the PageRank vector for the original graph  $G$ . A potential avenue would be exploring the CUR-trans method, expressing the stochastic matrix  $M$  as  $M = CUR$  [7], pooling together the most "informative" rows and columns in order to approximate PageRank.

## References

- [1] Artificial Intelligence - All in One. Lecture 8: Pagerank power iteration — stanford university. YouTube video, 2016. Published April 12.
- [2] P. Chau. Pagerank algorithm explained with examples. Medium article, 2023. Published November 14.
- [3] P. Lofgren. Efficient algorithms for personalized pagerank, 2015.
- [4] Pi Math Cornell. Hits algorithm: Hubs and authorities on the internet, n.d.
- [5] Wikipedia Contributors. Hits algorithm, 2019. Page updated May 26.
- [6] Wikipedia Contributors. Pagerank, 2019. Page updated June 18.
- [7] S. Wu, D. Wu, J. Quan, T. N. Chan, and K. Lu. Efficient and accurate pagerank approximation on large graphs. *Proceedings of the ACM on Management of Data*, 2(4):1–26, 2024.

## Appendix I. Code

```
1 import wikipediaapi
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from collections import deque
6 import time
7 import requests
8 import json
9 import random
10 # Included all necessary libraries. NetworkX used for graph creation. WikiAPI
11 # used for scraping necessary data. Matplotlib and Numpy for
    calculation/visualization
12
13 # Initialize Wikipedia API
14 wiki_wiki = wikipediaapi.Wikipedia(
15     user_agent='WikipediaGraphAnalysis/1.0 (https://example.com)',
16     language='en',
17     extract_format=wikipediaapi.ExtractFormat.WIKI
18 )
19
20 def search_wikipedia_first_result(search_term):
21     """
22     Searched for the very first result for the search term.
23     Doing a search instead of direct lookup since syntactically
24     whatever we enter as a seed might not be the exact page name,
25     so this helps us circumvent that.
26     """
27     try:
28         url = "https://en.wikipedia.org/w/api.php"
29         params = {
30             "action": "query",
31             "list": "search",
32             "srsearch": search_term,
33             "srlimit": 1,
34             "format": "json",
35             "srnamespace": 0
36         }
37
38         headers = {
39             'User-Agent': 'WikipediaGraphAnalysis/1.0 (https://example.com;
40                 contact@example.com)'
41         }
42
43         response = requests.get(url, params=params, timeout=15, headers=headers)
44
45         # Check response status
46         if response.status_code != 200:
47             print(f"HTTP {response.status_code} for '{search_term}'")
48             return None
49
50         # Parsing JSON
51         try:
52             data = response.json()
53             print(f"Response preview")
54             return None
```



```

55         if "query" in data and "search" in data["query"] and
56             len(data["query"]["search"]) > 0:
57             title = data["query"]["search"][0]["title"]
58             page = wiki_wiki.page(title)
59             if page.exists():
60                 return title
61         return None
62     except Exception as e:
63         print(f"Error searching")
64         return None
65
66     return None
67
68 def get_wiki_links(page_title, max_links=50):
69     """
70     Get internal Wiki links from a page.
71     Based on parameters that we tuned, especially
72     the maximum number of links, which is discussed
73     more thoroughly in the report. Filtering for actual
74     links and not files or categories.
75
76     Randomized since we initially found that
77     our links were returned in alphabetical order
78     in the event of a tie between pagerank scores. Randomized to introduce
79     variation in topics chosen.
80     """
81     try:
82         page = wiki_wiki.page(page_title)
83
84         # Collect all valid links
85         all_links = []
86         for link_title in page.links.keys():
87             if link_title and not link_title.startswith('Category:') and not
88                 link_title.startswith('File:') and not
89                 link_title.startswith('Template:'):
90                 all_links.append(link_title)
91
92         # Randomize to avoid uniformity, then returning up until the maximum
93         # amount of links
94         random.shuffle(all_links)
95         return all_links[:max_links]
96     except Exception as e:
97         print(f"Error getting links")
98         return []
99
100 def build_wikipedia_graph(seed_terms, max_depth=2, max_links_per_page=20,
101     max_total_pages=200):
102     """
103     Build a directed graph from Wikipedia pages starting with seed terms. (u,
104     w) and (w, u) valid edges in this graph.
105     - seed_terms: List of search terms
106     - max_depth: How many levels of links to follow
107     - max_links_per_page: Maximum links to extract per page
108     - max_total_pages: Maximum total pages to include in graph
109     """
110     G = nx.DiGraph()
111     # BFS Mark set to check membership easily
112     visited = set()

```

```

107 seed_pages = []
108
109 for term in seed_terms:
110     page_title = search_wikipedia_first_result(term)
111     if page_title:
112         print(f"Found: '{term}' -> '{page_title}'")
113         seed_pages.append(page_title)
114         visited.add(page_title)
115         G.add_node(page_title, is_seed=True)
116     else:
117         print(f"No result found for '{term}'")
118
119 # Running BFS
120 queue = deque([(page, 0) for page in seed_pages]) # (page_title, depth)
121
122 while queue and len(visited) < max_total_pages:
123     current_page, depth = queue.popleft()
124
125     if depth >= max_depth:
126         continue
127
128     # Stop if we have reached the maximum number of pages
129     if len(visited) >= max_total_pages:
130         break
131
132     # Get links from current page - use max_links_per_page parameter
133     links = get_wiki_links(current_page, max_links=max_links_per_page)
134
135     for link in links:
136         # Stop adding nodes if we've reached the limit
137         if len(visited) >= max_total_pages:
138             break
139
140         if link != current_page:
141             if link not in visited and len(visited) < max_total_pages:
142                 # Can add new node if the link is not already in visited
143                 visited.add(link)
144                 G.add_node(link, is_seed=False)
145                 G.add_edge(current_page, link)
146
147                 # Add to queue if we haven't reached max depth
148                 if depth + 1 < max_depth:
149                     queue.append((link, depth + 1))
150
151     time.sleep(0.1) # To avoid timeouts from excessive API calls
152
153 return G
154
155 def personalized_pagerank(G, seed_pages, damping=0.85, max_iter=100, tol=1e-6):
156     """
157     Runs Personalized PageRank
158
159     Using power iteration to calculate the Pagerank vector, within a tolerance of
160     1e-6.
161
162     Only addition to Power iteration is that we restart at seed pages, making
163     algorithm focus on nodes relevant to the seed topics.
164     """

```

```

164 nodes = list(G.nodes())
165 n = len(nodes)
166
167 if n == 0:
168     return {}
169
170 # Create node-to-index mapping, since networkX doesn't enumerate nodes
171 node_to_idx = {node: i for i, node in enumerate(nodes)}
172 idx_to_node = {i: node for i, node in enumerate(nodes)}
173
174 # Build adjacency matrix A (for directed graph)
175 # A[i,j] = 1 if there's a directed edge FROM node i TO node j
176 A = np.zeros((n, n))
177 for i, node in enumerate(nodes):
178     # For directed graph: check outgoing edges (successors)
179     for neighbor in G.successors(node):
180         j = node_to_idx[neighbor]
181         A[i, j] = 1.0 # Edge from node i to node j
182
183 # Making markov matrix out of adjacency matrix A
184 # Compute column sums (out-degrees)
185 col_sums = np.sum(A, axis=0)
186
187 # Build transition matrix M: handle dangling nodes (columns with sum 0)
188 # For columns with sum > 0: M[i,j] = A[i,j] / col_sum[j]
189 # For columns with sum = 0 (dangling nodes): M[i,j] = 1/n (uniform
    distribution)
190 M = np.zeros((n, n))
191 for j in range(n):
192     if col_sums[j] > 0:
193         M[:, j] = A[:, j] / col_sums[j]
194     else:
195         # Dangling node: replace with uniform distribution 1/n
196         M[:, j] = 1.0 / n
197
198 # Create personalization vector v (teleportation vector)
199 v = np.zeros(n)
200 seed_weight = 1.0 / len(seed_pages)
201 for node in nodes:
202     if node in seed_pages:
203         # Giving equal weight to seed pages, nonzero only for seeds
204         v[node_to_idx[node]] = seed_weight
205     # Normalize to ensure it sums to exactly 1
206 v = v / np.sum(v) if np.sum(v) > 0 else np.ones(n) / n
207
208 # Initialize PageRank vector
209 pr = np.ones(n) / n
210
211 # Power method
212 for iteration in range(max_iter):
213     # Compute new PageRank
214     pr_new = (1 - damping) * v + damping * (M @ pr)
215
216     # Check if it is within tolerance
217     diff = np.sum(np.abs(pr_new - pr))
218     pr = pr_new
219
220
221     if diff < tol:

```

```

222         break
223
224     # Normalize entries
225     pr = pr / np.sum(pr) if np.sum(pr) > 0 else pr
226
227     # Convert back to dictionary
228     pagerank_scores = {idx_to_node[i]: float(pr[i]) for i in range(n)}
229
230     return pagerank_scores
231
232 def visualize_graph(G, pagerank_scores, seed_pages, top_n=10):
233     """
234     Visualize the graph with node sizes and colors based on PageRank scores.
235     """
236     if G.number_of_nodes() == 0:
237         print("Graph is empty, cannot visualize.")
238         return
239
240     plt.figure(figsize=(16, 12))
241
242     # Get top influential nodes - sort by score (descending)
243     def sort_key(item):
244         node, score = item
245         # tie breaker for same scores
246         return (-score, hash(node) % 1000000)
247
248     sorted_nodes = sorted(pagerank_scores.items(), key=sort_key)
249     top_nodes = [node for node, _ in sorted_nodes[:top_n]]
250
251     node_sizes = []
252     node_colors = []
253     node_labels = {}
254
255     for node in G.nodes():
256         # Node size based on PageRank score
257         size = pagerank_scores.get(node, 0) * 10000
258         node_sizes.append(max(size, 50)) # Minimum size
259
260         # Color: red for seeds, orange for top influential, blue for others
261         if node in seed_pages:
262             node_colors.append('red')
263         elif node in top_nodes:
264             node_colors.append('orange')
265         else:
266             node_colors.append('lightblue')
267
268         # Only label seed pages and top influential nodes
269         if node in seed_pages or node in top_nodes:
270             # Truncate long titles
271             label = node[:30] + '...' if len(node) > 30 else node
272             node_labels[node] = label
273
274     # Use spring layout for better visualization
275     pos = nx.spring_layout(G, k=1, iterations=50, seed=42)
276
277     # Draw the graph
278     nx.draw_networkx_nodes(G, pos, node_size=node_sizes, node_color=node_colors,
279                            alpha=0.7)

```

```

279 nx.draw_networkx_edges(G, pos, alpha=0.2, width=0.5, arrows=True,
280                        arrowsize=10, arrowstyle='->')
281 nx.draw_networkx_labels(G, pos, node_labels, font_size=8, font_weight='bold')
282
283 plt.title('Wikipedia Graph with Personalized PageRank\n(Red=Seed, Orange=Top
284           Influential, Blue=Others)',
285           fontsize=14, fontweight='bold')
286 plt.axis('off')
287 plt.tight_layout()
288 plt.show()
289
290 # Print top influential nodes
291 print(f"\nTop {top_n} Most Influential Pages (by Personalized PageRank):")
292 print("-" * 70)
293 prev_score = None
294 for i, (node, score) in enumerate(sorted_nodes[:top_n], 1):
295     # Show if score is same as previous (tied)
296     print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
297     prev_score = score
298
299 def hits(G, max_iter=100, tol=1e-6):
300     """
301     Runs Hyperlink-Induced Topic Search (HITS)
302     Updates authority and hub vectors to convergence until tolerance of 1e-6 is
303     reached:
304     a_{k+1} = A^T * h_k
305     h_{k+1} = A * a_{k+1}
306     a_{k+1} = a_{k+1} / ||a_{k+1}||_2
307     h_{k+1} = h_{k+1} / ||h_{k+1}||_2
308     """
309     nodes = list(G.nodes())
310     n = len(nodes)
311
312     if n == 0:
313         return {}, {}
314
315     # Create node-to-index mapping, since networkX doesn't enumerate nodes
316     node_to_idx = {node: i for i, node in enumerate(nodes)}
317     idx_to_node = {i: node for i, node in enumerate(nodes)}
318
319     # Build adjacency matrix A (for directed graph)
320     # A[i,j] = 1 if there's a directed edge FROM node i TO node j
321     A = np.zeros((n, n))
322     for i, node in enumerate(nodes):
323         # For directed graph: check outgoing edges (successors)
324         for neighbor in G.neighbors(node):
325             j = node_to_idx[neighbor]
326             A[i, j] = 1.0 # Edge from node i to node j
327
328     # Initialize authority vector and hub vector
329     hub = np.random.rand(n) * 1e-6
330     auth = np.random.rand(n) * 1e-6
331
332     # Iterative updates of a and h
333     for iteration in range(max_iter):
334         # Compute new authority vector: a_{k+1} = A^T * h_k
335         new_auth = A.T @ hub
336         # Compute new hub vector: h_{k+1} = A * a_{k+1}
337         new_hub = A @ new_auth

```

```

335     # Normalize authority and hub vectors
336     a_norm = np.linalg.norm(new_auth, 2)
337     h_norm = np.linalg.norm(new_hub, 2)
338     if a_norm == 0 or h_norm == 0:
339         break
340     new_auth = new_auth / a_norm
341     new_hub = new_hub / h_norm
342
343     # Check if it is within the tolerance
344     diff_a = np.linalg.norm(new_auth - auth, 2)
345     diff_h = np.linalg.norm(new_hub - hub, 2)
346     auth = new_auth
347     hub = new_hub
348     if diff_a < tol and diff_h < tol:
349         break
350
351     # Convert back to dictionary
352     authority_scores = {idx_to_node[i]: float(auth[i]) for i in range(n)}
353     hub_scores = {idx_to_node[i]: float(hub[i]) for i in range(n)}
354
355     return authority_scores, hub_scores
356
357 def visualize_hits(G, authority_scores, hub_scores, seed_pages, top_n=10):
358     """
359     Visualize the graph with node sizes and colors based on HITS authority and
360     hub scores.
361     """
362     if G.number_of_nodes() == 0:
363         return
364
365     # Get top influential nodes - sort by score (descending)
366     def sort_key(item):
367         node, score = item
368         # tie breaker for same scores
369         return (-score, hash(node) % 1000000)
370
371     sorted_auth = sorted(authority_scores.items(), key=sort_key)
372     sorted_hub = sorted(hub_scores.items(), key=sort_key)
373     top_auth = [node for node, _ in sorted_auth[:top_n]]
374     top_hub = [node for node, _ in sorted_hub[:top_n]]
375
376     # Use spring layout for better visualization
377     pos = nx.spring_layout(G, k=1, iterations=50, seed=42)
378
379     # Function to draw the graph for each type of score
380     def draw_graph(title, score_map, top_nodes):
381         plt.figure(figsize=(16, 12))
382
383         node_sizes = []
384         node_colors = []
385         node_labels = {}
386
387         for node in G.nodes():
388             # Node size based on HITS authority or hub score
389             size = score_map.get(node, 0) * 10000
390             node_sizes.append(max(size, 50)) # Minimum size
391
392             # Color: red for seeds, orange for top influential, blue for others

```

```

393         if node in seed_pages:
394             node_colors.append('red')
395         elif node in top_nodes:
396             node_colors.append('orange')
397         else:
398             node_colors.append('lightblue')
399
400     # Only label seed pages and top influential nodes
401     if node in seed_pages or node in top_nodes:
402         # Truncate long titles
403         lbl = node[:30] + '...' if len(node) > 30 else node
404         node_labels[node] = lbl
405
406     # Draw the graph
407     nx.draw_networkx_nodes(G, pos, node_size=node_sizes,
408                             node_color=node_colors, alpha=0.7)
409     nx.draw_networkx_edges(G, pos, alpha=0.2, width=0.5, arrows=True,
410                             arrowsize=10, arrowstyle='->')
411     nx.draw_networkx_labels(G, pos, node_labels, font_size=8,
412                             font_weight='bold')
413     plt.title(title, fontsize=14, fontweight='bold')
414     plt.axis('off')
415     plt.tight_layout()
416     plt.show()
417
418     # Draws graph for authority and hub scores
419     draw_graph("Wikipedia Graph with HITS Authority Scores\n(Red=Seed, Orange=Top
420                 Influential, Blue=Others)", authority_scores, top_auth)
421     draw_graph("Wikipedia Graph with HITS Hub Scores\n(Red=Seed, Orange=Top
422                 Influential, Blue=Others)", hub_scores, top_hub)
423
424     # Print top influential nodes for authorities
425     print(f"\nTop {top_n} Most Influential Pages (by HITS Authority Scores):")
426     print("-" * 70)
427     prev_score = None
428     for i, (node, score) in enumerate(sorted_auth[:top_n], 1):
429         # Show if score is same as previous (tied)
430         print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
431         prev_score = score
432
433     # Print top influential nodes for hubs
434     print(f"\nTop {top_n} Most Influential Pages (by HITS Hub Scores):")
435     print("-" * 70)
436     prev_score = None
437     for i, (node, score) in enumerate(sorted_hub[:top_n], 1):
438         # Show if score is same as previous (tied)
439         print(f"{i:2d}. {node[:50]:50s} (Score: {score:.8f})")
440         prev_score = score
441
442     # Input params
443     seed_terms = [
444         "planets",
445         "supernova",
446         "hubble"
447     ]
448
449     # Build the graph
450     G = build_wikipedia_graph(
451         seed_terms,

```

```

447     max_depth=2,
448     max_links_per_page=2,
449     max_total_pages=200
450 )
451
452 # Initialize pages for seed articles
453 seed_pages = [node for node, attrs in G.nodes(data=True) if attrs.get('is_seed',
454     False)]
455
456 # Calculate PageRank scores
457 pagerank_scores = personalized_pagerank(G, seed_pages, damping=0.85)
458 # Visualize the graph
459 visualize_graph(G, pagerank_scores, seed_pages, top_n=30)
460
461 # Calculate HITS authority and hub scores
462 authority_scores, hub_scores = hits(G)
463 # Visualize the graphs
464 visualize_hits(G, authority_scores, hub_scores, seed_pages, top_n=30)

```