# Assignment 5

CS834-F16: Introduction to Information Retrieval

Fall 2016

Erika Siregar

# Question 10.3

Compute five iterations of HITS (see Algorithm 3) and PageRank (see Figure 4.11) on the
graph in Figure 10.3. Discuss how the PageRank scores compare to the hub and authority
scores produced by HITS.

## Answer

Figure 1 shows the directed graph from the textbook [1] on which we will calculate the scores of
HITS and PageRank. Computing HITS (authorities and hubs) and PageRank scores are pretty easy
since we can just utilize the Link Analysis procedure that is provided by python library 'networkx'
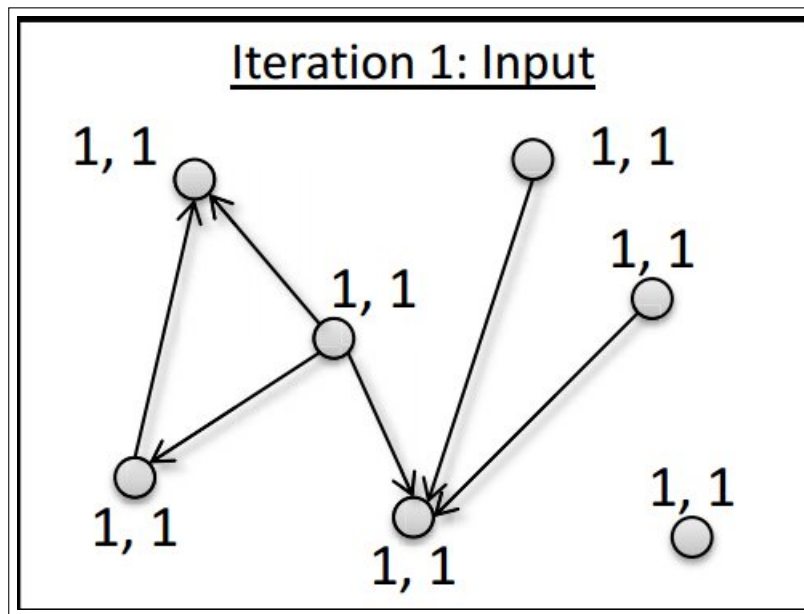[2].

Figure 1: Figure 10.3 from the textbook [1]

Figure 2 shows the scores of HITS (authorities and hubs) and PageRank, which are obtained by
running the code in listing 1. We only need to set the number of iterations.

Figure 2: HITS and Pagerank for Figure 10.3 with 5 Iterations

To make the analysis and comparison easier, I transformed the output int figure 2 into a neat table format as can be seen on table 1. From table 1, we can see that, generally, the authorities values are linearly proportional to those of PageRank. After 5 iterations, node 3 gets the highest score for 'authorities' and the second highest score for 'PageRank'. Nodes 1 and 2 get lower 'authorities' score than that of node 3, but higher 'authorities' score compare to nodes 5 and 6. The same thing can also be concluded by comparing the PageRank scores for those five nodes (1, 2, 3, 5, and 6). The strange thing happens on node 4, where its 'authorities' score is lower than node 3, but its 'PageRank' score is higher than node 3. This anomaly takes place probably because we only do 5 iterations. Maybe, if we continue iterating until the values converge into certain number, this anomaly will not happen.

| Node | Score | | |
|------|-------|------|---------|
| | **Hubs** | **Authorities** | **PageRank** |
| 1 | 0.198707510685935 | 0.201534170153416 | 0.154152105903131 |
| 2 | 0.198707510685935 | 0.201534170153416 | 0.154152105903131 |
| 3 | 0.258854060655404 | 0.252324500232450 | 0.215183655595341 |
| 4 | 0.178963973132505 | 0.188168293816829 | 0.273070054356128 |
| 5 | 0.082383472420110 | 0.078219432821943 | 0.089524353405041 |
| 6 | 0.082383472420110 | 0.078219432821943 | 0.089524353405041 |
| 7 | 0.000000000000000 | 0.000000000000000 | 0.024393371432184 |

Table 1: HITS and Pagerank for Figure 10.3 with 5 Iterations

```
1 #!/usr/bin/python
2
3 import networkx as nx
4
```

```python
def hits(G, iter=100, nstart=None, normalized=True):
    if type(G) == nx.MultiGraph or type(G) == nx.MultiDiGraph:
        raise Exception("hits() not defined for graphs with multiedges.")
    if len(G) == 0:
        return {},{}
    # choose fixed starting vector if not given
    if nstart is None:
        h=dict.fromkeys(G,1.0/G.number_of_nodes())
    else:
        h=nstart
    # normalize starting vector
    s=1.0/sum(h.values())
    for k in h:
        h[k]*=s
    i=0
    while True: # power iteration: make up to max_iter iterations
        if i >= iter: break

        hlast=h
        h=dict.fromkeys(hlast.keys(),0)
        a=dict.fromkeys(hlast.keys(),0)
        # this "matrix multiply" looks odd because it is
        # doing a left multiply a^T=hlast^T*G
        for n in h:
            for nbr in G[n]:
                a[nbr]+=hlast[n]*G[n][nbr].get('weight',1)
        # now multiply h=Ga
        for n in h:
            for nbr in G[n]:
                h[n]+=a[nbr]*G[n][nbr].get('weight',1)
        # normalize vector
        s=1.0/max(h.values())
        for n in h: h[n]*=s
        # normalize vector
        s=1.0/max(a.values())
        for n in a: a[n]*=s

        i+=1
    if normalized:
        s = 1.0/sum(a.values())
        for n in a:
            a[n] *= s
        s = 1.0/sum(h.values())
        for n in h:
            h[n] *= s
    return h,a

def pagerank(G, alpha=0.85, personalization=None,
             iter=100, nstart=None, weight='weight',
             dangling=None):
    if len(G) == 0:
        return {}

    if not G.is_directed():
        D = G.to_directed()
    else:
        D = G

    # Create a copy in (right) stochastic form
```

```python
64  W = nx.stochastic_graph(D, weight=weight)
65  N = W.number_of_nodes()
66
67  # Choose fixed starting vector if not given
68  if nstart is None:
69      x = dict.fromkeys(W, 1.0 / N)
70  else:
71      # Normalized nstart vector
72      s = float(sum(nstart.values()))
73      x = dict((k, v / s) for k, v in nstart.items())
74
75  if personalization is None:
76      # Assign uniform personalization vector if not given
77      p = dict.fromkeys(W, 1.0 / N)
78  else:
79      missing = set(G) - set(personalization)
80      if missing:
81          raise nx.NetworkXError('Personalization dictionary '
82          'must have a value for every node. '
83          'Missing nodes %s' % missing)
84      s = float(sum(personalization.values()))
85      p = dict((k, v / s) for k, v in personalization.items())
86
87  if dangling is None:
88      # Use personalization vector if dangling vector not specified
89      dangling_weights = p
90  else:
91      missing = set(G) - set(dangling)
92      if missing:
93          raise nx.NetworkXError('Dangling node dictionary '
94          'must have a value for every node. '
95          'Missing nodes %s' % missing)
96      s = float(sum(dangling.values()))
97      dangling_weights = dict((k, v/s) for k, v in dangling.items())
98  dangling_nodes = [n for n in W if W.out_degree(n, weight=weight) == 0.0]
99
100 # power iteration: make up to max_iter iterations
101 for _ in range(iter):
102     xlast = x
103     x = dict.fromkeys(xlast.keys(), 0)
104     danglesum = alpha * sum(xlast[n] for n in dangling_nodes)
105     for n in x:
106         # this matrix multiply looks odd because it is
107         # doing a left multiply x^T=xlast^T*W
108         for nbr in W[n]:
109             x[nbr] += alpha * xlast[n] * W[n][nbr][weight]
110         x[n] += danglesum * dangling_weights[n] + (1.0 - alpha) * p[n]
111
112     return x

113
114 if __name__ == '__main__':
115     iter = 5
116     G = nx.Graph()
117
118     # Add 7 nodes
119     G.add_nodes_from(range(1,8))
120
121     # Add 6 edges
122     G.add_edges_from([(1,2), (3,1), (3,2), (3,4), (5,4), (6,4)])
```

```
123
124 # Compute hubs and authorities normalized values using hits
125 h, a = hits(G, iter=iter)
126
127 print 'HITS Algorithm ({} iterations)'.format(iter)
128 print '=================='
129 print 'Hubs values = {}'.format(h)
130 print 'Authorities values = {}'.format(a)
131 print ''
132
133 # Compute pagerank of each nodes
134 pr = pagerank(G, iter=iter)
135
136 print 'Pagerank Algorithm ({} iterations)'.format(iter)
137 print '=================='
138 print 'Pagerank values = {}'.format(pr)
```

Listing 1: Computing HITS and PageRank

## Question 0.3

**Answer:**

answer 2

## Question 8.5

question 3

**Answer**

answer 3

## Question 8.7

question 4

**Answer**

answer 4

## References

[1] Bruce Croft, Donald Metzler, and Trevor Strohman. *Search Engines: Information Retrieval in Practice*. Addison-Wesley Publishing Company, USA, 1st edition, 2009.

[2] NetworkX Developers. Networkx - Link Analysis. `https://networkx.github.io/documentation/networkx-1.9/reference/algorithms.link_analysis.html`, 2016. [Online; accessed 14-December-2016].