

Article

Matlab algorithm to generate adjacency matrix from connection pairs that nodes are represented by strings

WenJun Zhang¹, Yanhong Qi²

¹School of Life Sciences, Sun Yat-sen University, Guangzhou 510275, China; International Academy of Ecology and Environmental Sciences, Hong Kong

²Libraries of Sun Yat-sen University, Sun Yat-sen University, Guangzhou, China

E-mail: zhwj@mail.sysu.edu.cn, qiyh@mail.sysu.edu.cn

Received 12 November 2019; Accepted 18 December 2019; Published 1 December 2020



Abstract

It is harder to generate an adjacency matrix from connection (i.e., link) pairs that nodes are represented by strings than by numerical ID numbers. In present article, we developed a Matlab algorithm to generate adjacency matrix from connection pairs in which nodes are represented by strings. Full codes and executable program of the algorithm were given.

Keywords Matlab algorithm; connection pairs; connection weights; adjacency matrix; transformation.

Selforganizology
ISSN 2410-0080
URL: <http://www.iaees.org/publications/journals/selforganizology/online-version.asp>
RSS: <http://www.iaees.org/publications/journals/selforganizology/rss.xml>
E-mail: selforganizology@iaees.org
Editor-in-Chief: WenJun Zhang
Publisher: International Academy of Ecology and Environmental Sciences

1 Introduction

Adjacency matrices are widely used in many algorithms of graph theory and network science (Zhang, 2015, 2016a-k, 2017a-d, 2018; Zhang and Li, 2016; Zhang and Feng, 2017). On the other hand, we have usually collected connection pairs, in which nodes are always represented by strings (names, abbreviations, ABC codes, etc.). Nevertheless, it is much harder to generate an adjacency matrix from connection (i.e., link) pairs that nodes are represented by strings than by numerical ID numbers. In this article, I present a Matlab algorithm to generate adjacency matrix from connection pairs in which nodes are represented by strings. Full codes and executable program of the algorithm are given also.

2 Method

Suppose there are *num* connection pairs, which are recorded in a space separated text file as, for example

```
v1 v9  
v1 v20  
v3 v5  
...
```

(1)

or

$$\begin{array}{l} v_1 \quad v_5 \quad w_1 \\ v_1 \quad v_2 \quad w_2 \\ v_2 \quad v_9 \quad w_3 \\ \dots \end{array}$$

(2)

where v_i are nodes represented by strings and w_i are connection weights (numerical values), $i=1,2,\dots$. We want to find the series of unique nodes in (1) or (2):

$$v_1, v_2, \dots, v_m \quad (3)$$

where m is the total number of unique nodes, and transform (1) to an adjacency matrix $d=(d_{ij})_{m \times m}$. $d_{ij}=1$, if two nodes v_i and v_j are adjacent (connected), and $d_{ij}=0$, if v_i and v_j are not adjacent; $i, j=1,2,\dots, m$; or transform (2) to an adjacency matrix $d=(d_{ij})_{m \times m}$. $d_{ij}=w_{ij}$, if two nodes v_i and v_j are adjacent (connected), and $d_{ij}=0$, if v_i and v_j are not adjacent; $i, j=1,2,\dots, m$. Or for the situation (2), if a threshold for the connection weight, $h>0$, is given, we have $d_{ij}=1$ if $w_{ij} \geq h$ and, v_i and v_j are adjacent (connected), and $d_{ij}=0$ if $w_{ij} < h$ or v_i and v_j are not adjacent; $i, j=1,2,\dots, m$.

The IDs of rows and columns of the resultant adjacency matrix correspond to the natural IDs of unique nodes in the resultant series of unique nodes.

According to the above principle, the full Matlab algorithm, used in Matlab environment, is developed as the following.

```
[newdel,OK]=listdlg('liststring',{'Node-Node Data','Node-Node-Weight Data'},'listsize',[300
100],'OkString','OK','CancelString','Cancel','promptstring','Two Columns or Three Columns Data?','selectionmode','single');
if (newdel==2)
[choice,OK]=listdlg('liststring',{'Generate Weight Matrix','Generate 0-1 Adjacency Matrix'},'listsize',[300
100],'OkString','OK','CancelString','Cancel','promptstring','Two Columns or Three Columns Data?','selectionmode','single');
if (choice==2)
threshold=input('Input the threshold of the weight to produce 0 or 1 (=1 if >=weight, and =0 if <weight): ');
end
end
if (newdel==2)
str=input('Input the text file name of node-node-weight matrix w (w=(wij)numx3): ','s');
[w(:,1) w(:,2) w(:,3)]=textread(str,'%s%s%s');
elseif (newdel==1)
str=input('Input the text file name of node-node matrix w (w=(wij)numx2): ','s');
[w(:,1) w(:,2)]=textread(str,'%s%s');
end
num=length(w(:,1));
[st,m]=uniqNodes(w);
mm=0;
for i=1:num
for k=1:m
if (isequal(w{i,1},st{k})==1)
```

```

wid(i,1)=k;
end
if (isequal(w{i,2},st{k})==1)
wid(i,2)=k;
end
end
if (newdel==2)
wid(i,3)=str2num(w{i,3});
end
end
d=zeros(m,m);
for k=1:num
if (newdel==1) d(wid(k,1),wid(k,2))=1; d(wid(k,2),wid(k,1))=1;
elseif ((choice==1) && (newdel==2)) d(wid(k,1),wid(k,2))=wid(k,3); d(wid(k,2),wid(k,1))=wid(k,3);
elseif ((choice==2) && (newdel==2)) d(wid(k,1),wid(k,2))=wid(k,3)>threshold; d(wid(k,2),wid(k,1))=wid(k,3)>threshold;
end
end
fprintf('Matrix of Interactions\n');
d
fprintf(['Names of nodes from ID 1 to ' num2str(m) ':\n']);
st

```

```

function [st,m]=uniqNodes(w)
num=length(w(:,1));
w2=[w(:,1);w(:,2)];
m=1;
st{1}=w2{1};
for i=2:2*num
s=0;
for j=1:m
if isequal(w2{i},st{j})==0
s=s+1;
end
end
if (s==m)
m=m+1;
st{m}=w2{i};
end
end

```

The executable GUI software (see supplementary material) of the algorithm above is partly indicated in Fig. 1.

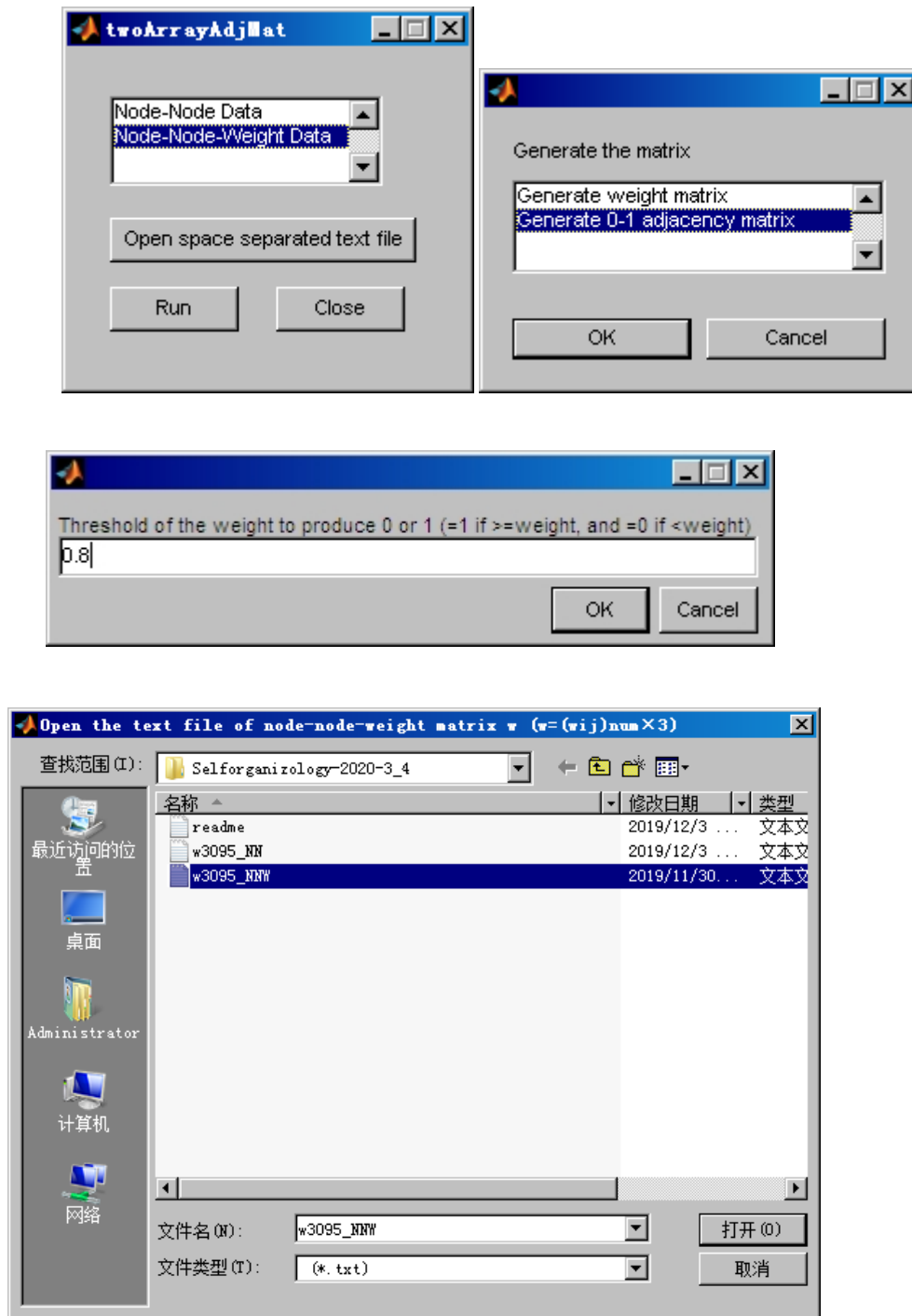
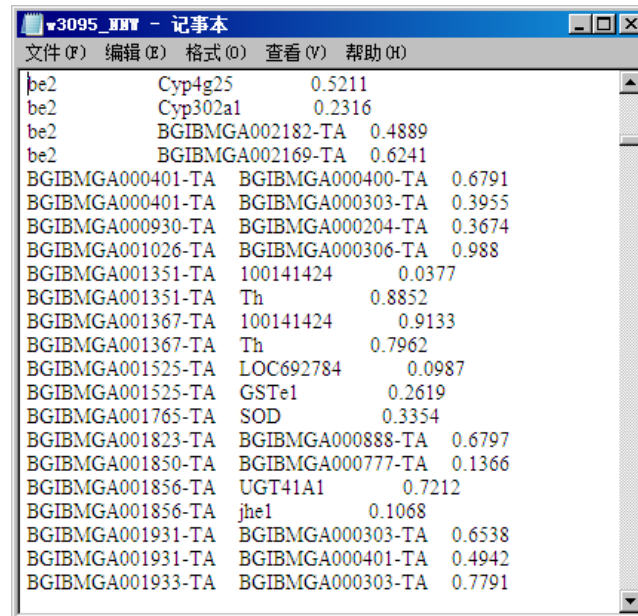


Fig. 1 The executable GUI software of the algorithm.

3 Application Example

There is a dataset for connection pairs with connection weights, as indicated in Fig. 2.

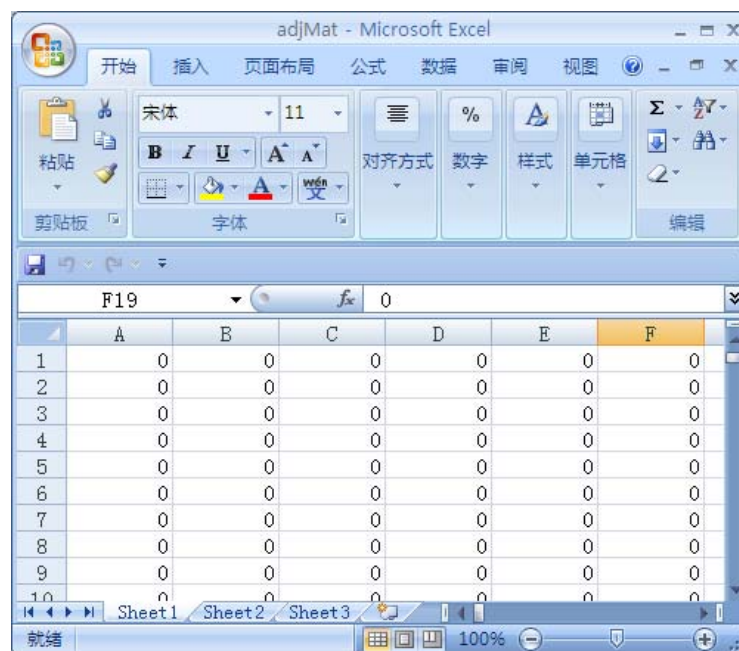


The screenshot shows a Notepad window titled 'w3095_1111 - 记事本'. The text inside is a list of connection pairs and their weights, separated by spaces. The data is as follows:

be2	Cyp4g25	0.5211
be2	Cyp302a1	0.2316
be2	BGIBMGA002182-TA	0.4889
be2	BGIBMGA002169-TA	0.6241
BGIBMGA000401-TA	BGIBMGA000400-TA	0.6791
BGIBMGA000401-TA	BGIBMGA000303-TA	0.3955
BGIBMGA000930-TA	BGIBMGA000204-TA	0.3674
BGIBMGA001026-TA	BGIBMGA000306-TA	0.988
BGIBMGA001351-TA	100141424	0.0377
BGIBMGA001351-TA	Th	0.8852
BGIBMGA001367-TA	100141424	0.9133
BGIBMGA001367-TA	Th	0.7962
BGIBMGA001525-TA	LOC692784	0.0987
BGIBMGA001525-TA	GSTe1	0.2619
BGIBMGA001765-TA	SOD	0.3354
BGIBMGA001823-TA	BGIBMGA000888-TA	0.6797
BGIBMGA001850-TA	BGIBMGA000777-TA	0.1366
BGIBMGA001856-TA	UGT41A1	0.7212
BGIBMGA001856-TA	jhe1	0.1068
BGIBMGA001931-TA	BGIBMGA000303-TA	0.6538
BGIBMGA001931-TA	BGIBMGA000401-TA	0.4942
BGIBMGA001933-TA	BGIBMGA000303-TA	0.7791

Fig. 2 An example dataset for connection pairs with connection weights. The dataset is stored in a space separated text file.

In the software interface, choose Node-Node-Weight Data and a dialog appears. In the dialog, choose generate 0-1 adjacency matrix and an input dialog appears for inputting the threshold h , for example, 0.8 is entered. The data file is loaded by open the space separated text file dialog. Finally, run the software (Fig. 1) and the two sets of results (adjacency matrix and unique node series) can be saved in to two specified excel files (Fig. 3a and 3b).



The screenshot shows a Microsoft Excel window titled 'adjMat - Microsoft Excel'. The spreadsheet displays an adjacency matrix with columns labeled A through F and rows numbered 1 through 10. All cells in the matrix contain the value 0.

	A	B	C	D	E	F
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0
6	0	0	0	0	0	0
7	0	0	0	0	0	0
8	0	0	0	0	0	0
9	0	0	0	0	0	0
10	0	0	0	0	0	0

(3a)

	A	B	C	D	E	F
37	ace1					
38	Adar					
39	Adgf					
40	ae13					
41	ae25					
42	ae3					
43	ae40					
44	ae41					
45	ae49					
46	Ago2					
47	aif					
48	Akhr					
49	akr2e					
50	Aox1					
51	Aox2					
52	app					
53	ask1					
54	Bbx-b8					

(3b)

Fig. 3 Two sets of results, adjacency matrix (3a), and unique node series (3b).

4 Discussion

It is obvious that the algorithm is applicable to datasets that nodes are represented by both strings and numerical ID numbers.

If the dataset is stored in an excel file. Users can save the file as a space separated text file (.prn), and open the “.prn” text file and re-save it as the text file (.txt) that can be used in the algorithm.

Acknowledgment

We are thankful to the support of The National Key Research and Development Program of China (2017YFD0201204), and Discovery and Crucial Node Analysis of Important Biological and Social Networks (2015.6-2020.6), from Yangling Institute of Modern Agricultural Standardization.

References

Zhang WJ. 2015. Prediction of missing connections in the network: A node-similarity based algorithm. *Selforganizology*, 2(4): 91-101

- Zhang WJ. 2016a. Detecting connectedness of network: A Matlab program and application in tumor pathways and a phylogenetic network. *Selforganizology*, 3(4): 117-120
- Zhang WJ. 2016b. A mathematical model for dynamics of occurrence probability of missing links in predicted missing link list. *Network Pharmacology*, 1(4): 86-94
- Zhang WJ. 2016c. A node-similarity based algorithm for tree generation and evolution. *Network Biology*, 6(3): 55-64
- Zhang WJ. 2016d. Screening node attributes that significantly influence node centrality in the network. *Selforganizology*, 3(3): 75-86
- Zhang WJ. 2016e. How to find cut nodes and bridges in the network? A Matlab program and application in tumor pathways. *Network Pharmacology*, 1(3): 82-85
- Zhang WJ. 2016f. A method for identifying hierarchical sub-networks / modules and weighting network links based on their similarity in sub-network / module affiliation. *Network Pharmacology*, 1(2): 54-65
- Zhang WJ. 2016g. Finding trees in the network: Some Matlab programs and application in tumor pathways. *Network Pharmacology*, 1(2): 66-73
- Zhang WJ. 2016h. A Matlab program for finding shortest paths in the network: Application in the tumor pathway. *Network Pharmacology*, 1(1): 42-53
- Zhang WJ. 2016i. A node degree dependent random perturbation method for prediction of missing links in the network. *Network Biology*, 6(1): 1-11
- Zhang WJ. 2016j. A random network based, node attraction facilitated network evolution method. *Selforganizology*, 3(1): 1-9
- Zhang WJ. 2016k. *Selforganizology: The Science of Self-Organization*. World Scientific, Singapore
- Zhang WJ, Li X. 2016. A cluster method for finding node sets / sub-networks based on between- node similarity in sets of adjacency nodes: with application in finding sub-networks in tumor pathways. *Proceedings of the International Academy of Ecology and Environmental Sciences*, 6(1): 13-23
- Zhang WJ. 2017a. Maximum matching of the network: A Matlab program and application. *Selforganizology*, 4(4): 65-68
- Zhang WJ. 2017b. Finding minimum cost flow in the network: A Matlab program and application. *Selforganizology*, 4(2): 30-34
- Zhang WJ. 2017c. Finding fundamental circuits in the network: A Matlab program and application in tumor pathway. *Selforganizology*, 4(1): 14-17
- Zhang WJ. 2017d. Finding the shortest tree in the network: A Matlab program and application in tumor pathway. *Network Pharmacology*, 2(1): 13-16
- Zhang WJ, Feng YT. 2017. Metabolic pathway of non-alcoholic fatty liver disease: Network properties and robustness. *Network Pharmacology*, 2(1): 1-12
- Zhang WJ. 2018. *Fundamentals of Network Biology*. World Scientific Europe, London, UK