
MBRW_SPECTRAL_DIMENSION_DEMO

Table of Contents

Path to Boost library.	1
Download a protein-protein interaction network file.	1
Extract the direct protein-protein interactions.	2
Extract the largest MBRW-able subgraph.	2
Run a memory-biased random walk on the largest connected component.	4
Plot log ₂ segment mass generalized means vs segment length.	7
Use generalized means to estimate generalized spectral dimensions.	8
Plot the generalized means.	9
Take the range of finite-order generalized spectral dimensions.	9

Author: Adam Craig

Created: 2021-09-06.

Last Updated: 2021-09-06

This code demonstrates the following workflow:

1. Download a protein-protein interaction network file.
2. Extract the direct protein-protein interactions.
3. Extract the 2-core (largest subgraph with minimum degree of 2).
4. Separate the 2-core into its connected components.
5. Save the connected components to a file.
6. Run a memory-biased random walk on the largest connected component.
7. Plot log₂ segment mass generalized means vs segment length.
8. Compute generalized means of walk segments with power-of-2 lengths.
9. Use generalized means to estimate generalized spectral dimensions.
10. Plot the generalized means.
11. Take the range of finite-order generalized spectral dimensions.

Dependencies: The C++ MBRW utility uses methods from the Boost library.

You do not need to install it, just download the C++ header files.

I have tested this code with Boost versions 1.64.0 and 1.75.0.

<https://www.boost.org/>

The MATLAB scripts make use of the graph object class.

We assume behavior from MATLAB v 2018a or later.

<https://www.mathworks.com/help/matlab/graph-and-network-algorithms.html>

Path to Boost library.

Set this to where you keep the boost header files on your system.

```
boost_path = ['..' filesep 'boost_1_75_0'];
```

Download a protein-protein interaction network file.

The network on which we are working is from

Gunsalus, K. C., Ge, H., Schetter, A. J., Goldberg, D. S., Han, J. D. J., Hao, T., ... & Piano, F. (2005). Predictive models of molecular machines involved in *Caenorhabditis elegans* early embryogenesis. *Nature*, 436(7052), 861-865.

```
original_graph_file_url = ['https://static-content.springer.com/' ...  
    'esm/art%3A10.1038%2Fnature03876/MediaObjects/' ...  
    '41586_2005_BFnature03876_MOESM4_ESM.sif'];  
base_network_name = 'gunsalus_2005_c_elegans_ppi';  
print_network_name = 'C. elegans PPI from Gunsalus et al, 2005';  
networks_dir = 'networks';  
gene_names_dir = 'gene_name_lists';  
original_graph_file_name = [networks_dir filesep ...  
    base_network_name '_original.txt'];  
if ~exist(original_graph_file_name, 'file')  
    disp('Downloading original network...')  
    websave(original_graph_file_name, original_graph_file_url);  
end
```

Downloading original network...

Extract the direct protein-protein interactions.

```
direct_interactions_graph_file_name = [networks_dir filesep ...  
    base_network_name '_direct.txt'];  
direct_interactions_gene_names = [gene_names_dir filesep ...  
    base_network_name '_direct_gene_names.txt'];  
if ~exist(direct_interactions_graph_file_name, 'file') || ...  
    ~exist(direct_interactions_gene_names, 'file')  
    extract_full_gunsalus_et_al_2005( ...  
        original_graph_file_name, ...  
        direct_interactions_graph_file_name, ...  
        direct_interactions_gene_names );  
end
```

Extract the largest MBRW-able subgraph.

3. Extract the 2-core (largest subgraph with minimum degree of 2).
4. Separate the 2-core into its connected components.
5. Save the connected components to a file.

```
mbrw_able_edge_file_name = [networks_dir filesep ...  
    base_network_name '_direct_cc_1.txt'];  
mbrw_able_node_names_file_name = [gene_names_dir filesep ...  
    base_network_name '_direct_gene_names_cc_1.txt'];  
if ~exist(mbrw_able_edge_file_name, 'file') || ...  
    ~exist(mbrw_able_node_names_file_name, 'file')  
    [new_edge_file_names, new_node_name_file_names] = ...  
        make_mbrw_able(direct_interactions_graph_file_name, ...  
            direct_interactions_gene_names);  
    % Just work with the largest connected component.
```

```
mbrw_able_edge_file_name = new_edge_file_names{1};
mbrw_able_node_names_file_name = new_node_name_file_names{1};
end
G =
    read_graph(mbrw_able_edge_file_name,mbrw_able_node_names_file_name);

reading in graph...
Elapsed time is 0.023001 seconds.
has 281 nodes and 534 edges
simplifying...
Elapsed time is 0.000947 seconds.
281 nodes, 514 edges remain
taking 2-core...
iteration 1: 281 nodes, 514 edges remaining, removing 103...
iteration 2: 178 nodes, 431 edges remaining, removing 21...
iteration 3: 157 nodes, 415 edges remaining, removing 6...
iteration 4: 151 nodes, 410 edges remaining, removing 1...
iteration 5: 150 nodes, 409 edges remaining, removing 1...
done
Elapsed time is 0.004531 seconds.
149 nodes, 408 edges remain
finding connected components...
Elapsed time is 0.000145 seconds.
found 7 connected components
retrieving connected component 1...
128 nodes, 382 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_1.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_1.txt...
retrieving connected component 2...
5 nodes, 9 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_2.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_2.txt...
retrieving connected component 3...
4 nodes, 5 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_3.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_3.txt...
retrieving connected component 4...
3 nodes, 3 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_4.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_4.txt...
retrieving connected component 5...
3 nodes, 3 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_5.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_5.txt...
```

```
retrieving connected component 6...
3 nodes, 3 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_6.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_6.txt...
retrieving connected component 7...
3 nodes, 3 edges
writing to file networks
\gunsalus_2005_c_elegans_ppi_direct_cc_7.txt...
writing to file gene_name_lists
\gunsalus_2005_c_elegans_ppi_direct_gene_names_cc_7.txt...
done
Elapsed time is 0.064336 seconds.
```

Run a memory-biased random walk on the largest connected component.

bias: How much more likely are we to take a remembered edge vs an unremembered one, must be a positive integer
In general, values much less than 1000 make the effect of memory weak. Values much past 10,000 do not lead to improved sensitivity to community detection.

memory: Number of steps for which to remember an edge, must be a non-negative integer
In general, past a minimum of 4 or 5, making memory longer only weakly affects sensitivity. Most communities in networks contain shorter loops of 3 to 5 edges, and the walk agent is far more likely to find these than longer ones.

special values:
1 -> no memory except prevention of backtracking.
0 -> no memory, allow backtracking;
(Otherwise, MBRW disallows return to the previous node.)

rand_seed: seed with which to initialize the RNG
We include this for reproducibility.

log_num_steps: \log_2 (the number of steps for which to walk), must be a positive integer
Using more steps
increases the accuracy of the estimates of spectral dimension but also increases the run time.
 2^{23} steps offers a reasonable compromise for this network. Larger networks generally need more steps.

The C++ utility prints output that can help you check for convergence. As a general rule, using a value such that the number after 'min segment mass=' is the number of nodes on the last 3 lines of output is adequate.

```
% Try different parameter values.
bias = 10000;
memory = 5;
rand_seed = 0;
log_num_steps = 23;
```

```
segment_mass_log_multimean_file_name = [ ...
    'segment_mass_log_multimeans' filesep
    base_network_name '_cc_1' ...
    '_b_' num2str(bias) '_m_' num2str(memory) ...
    '_r_' num2str(rand_seed) '_c_' num2str(log_num_steps) '.csv'];
c_executable_name = 'mbrw_and_save_segment_mass_log_multimeans';
compile_command = sprintf( 'g++ -std=c++0x -I %s %s.cpp -o %s -
02', ...
    boost_path, c_executable_name, c_executable_name );
% Windows appends .exe to the file name automatically
% and does not expect ./ before the executable name when running.
% Linux does not append anything to the file name
% and expects ./ before the executable name when running.
if ispc
    dot_if_linux = '';
    exe_if_pc = '.exe';
else
    dot_if_linux = './';
    exe_if_pc = '';
end
run_command = sprintf('%s%s -i %s -o %s -b %u -m %u -r %u -c %u', ...
    dot_if_linux, c_executable_name, ...
    mbrw_able_edge_file_name, ...
    segment_mass_log_multimean_file_name, ...
    bias, memory, rand_seed, log_num_steps);
if ~exist(segment_mass_log_multimean_file_name,'file')
    if ~exist([c_executable_name exe_if_pc],'file')
        fprintf('compiling %s...\n',c_executable_name)
        tic
        compile_result = system(compile_command);
        toc
        if compile_result
            error('failed to compile %s', c_executable_name)
        end
    end
    fprintf('running %s...\n', c_executable_name)
    tic
    run_result = system(run_command);
    toc
    if run_result
        error('failed to run %s', c_executable_name);
    end
end

running mbrw_and_save_segment_mass_log_multimeans...
MBRW entropy:
input file: networks\gunsalus_2005_c_elegans_ppi_direct_cc_1.txt
output file: segment_mass_log_multimeans
\gunsalus_2005_c_elegans_ppi_cc_1_b_10000_m_5_r_0_c_23.csv
bias: 10000
memory length: 5
rand seed: 0
max doublings: 23
```

walking...

step/max segment length = 4 steps, min segment mass= 3 nodes, max mean
segment mass change: 0, max fractional change: 0, seconds: 0.000937

step/max segment length = 8 steps, min segment mass= 3 nodes, max
mean segment mass change: 1.09288e-016, max fractional change:
6.89528e-017, seconds: 0.000937

step/max segment length = 16 steps, min segment mass= 3 nodes, max
mean segment mass change: 1.27502e-016, max fractional change:
8.04449e-017, seconds: 0.000937

step/max segment length = 32 steps, min segment mass= 3 nodes, max
mean segment mass change: 2.56414e-016, max fractional change:
1.61779e-016, seconds: 0.000937

step/max segment length = 64 steps, min segment mass= 3 nodes, max
mean segment mass change: 6.66134e-016, max fractional change:
4.20284e-016, seconds: 0.000937

step/max segment length = 128 steps, min segment mass= 3 nodes,
max mean segment mass change: 6.66134e-016, max fractional change:
4.20284e-016, seconds: 0.000937

step/max segment length = 256 steps, min segment mass= 3 nodes,
max mean segment mass change: 1.77636e-015, max fractional change:
1.12076e-015, seconds: 0.001937

step/max segment length = 512 steps, min segment mass= 3 nodes,
max mean segment mass change: 1.77636e-015, max fractional change:
1.12076e-015, seconds: 0.001937

step/max segment length = 1024 steps, min segment mass= 3 nodes,
max mean segment mass change: 9.54792e-015, max fractional change:
6.02407e-015, seconds: 0.003939

step/max segment length = 2048 steps, min segment mass= 3 nodes,
max mean segment mass change: 1.24345e-014, max fractional change:
7.84529e-015, seconds: 0.005967

step/max segment length = 4096 steps, min segment mass= 20 nodes, max
mean segment mass change: 2.63697, max fractional change: 1.66374,
seconds: 0.00994

step/max segment length = 8192 steps, min segment mass= 21 nodes, max
mean segment mass change: 1.22194, max fractional change: 0.282731,
seconds: 0.017938

step/max segment length = 16384 steps, min segment mass= 21 nodes, max
mean segment mass change: 0.970476, max fractional change: 0.220948,
seconds: 0.033941

step/max segment length = 32768 steps, min segment mass= 30 nodes, max
mean segment mass change: 0.685844, max fractional change: 0.230888,
seconds: 0.06794

step/max segment length = 65536 steps, min segment mass= 35 nodes, max
mean segment mass change: 2.45698, max fractional change: 0.579211,
seconds: 0.12994

step/max segment length = 131072 steps, min segment mass= 49 nodes,
max mean segment mass change: 1.89999, max fractional change:
0.515607, seconds: 0.256963

step/max segment length = 262144 steps, min segment mass= 70 nodes,
max mean segment mass change: 2.39178, max fractional change:
0.462024, seconds: 0.507938

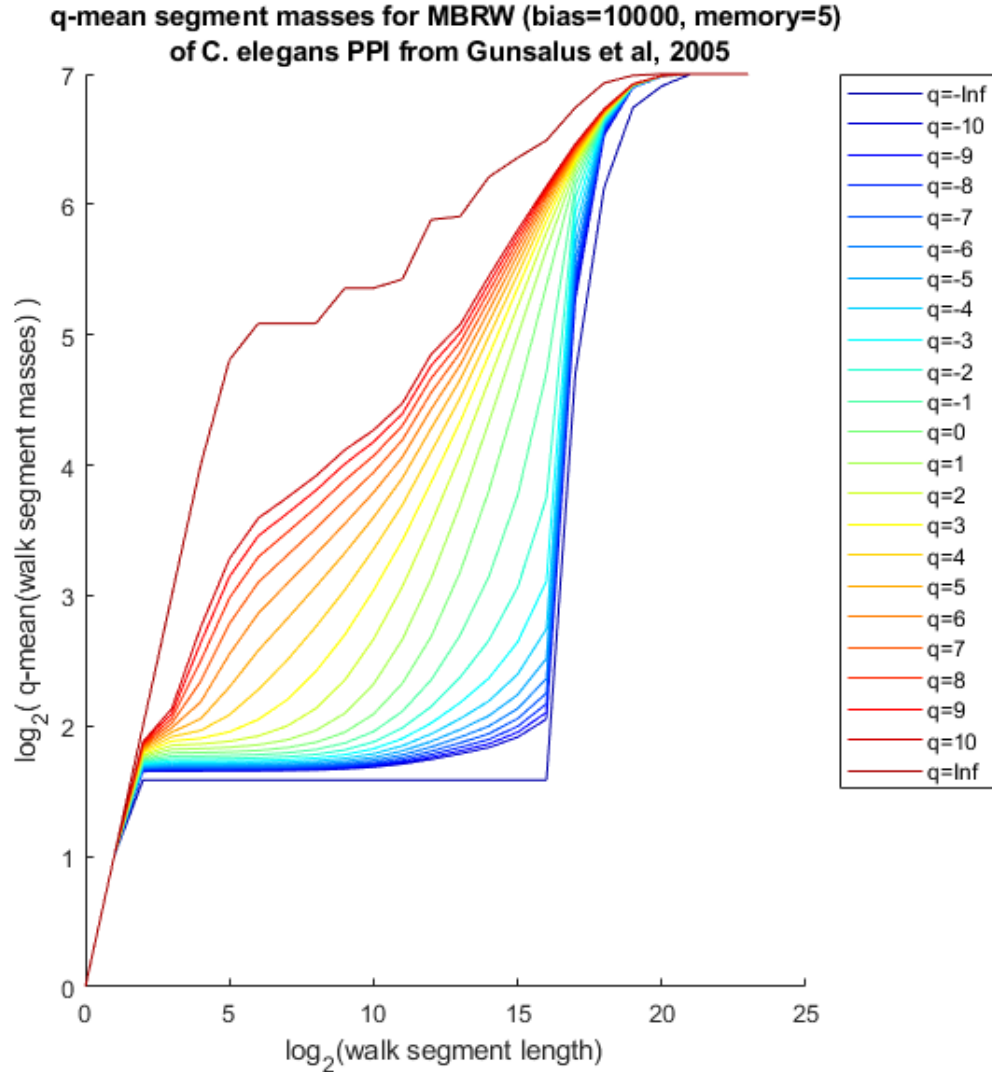
step/max segment length = 524288 steps, min segment mass= 115 nodes,
max mean segment mass change: 0.621588, max fractional change:
0.16809, seconds: 1.02094

```
step/max segment length = 1048576 steps, min segment mass= 120 nodes,
max mean segment mass change: 0.900141, max fractional change:
0.312011, seconds: 2.10216
step/max segment length = 2097152 steps, min segment mass= 128 nodes,
max mean segment mass change: 0.53382, max fractional change:
0.147795, seconds: 4.16811
step/max segment length = 4194304 steps, min segment mass= 128 nodes,
max mean segment mass change: 0.330106, max fractional change:
0.064513, seconds: 8.5367
step/max segment length = 8388608 steps, min segment mass= 128 nodes,
max mean segment mass change: 0.172001, max fractional change:
0.0357557, seconds: 16.6164
done
Elapsed time is 16.747999 seconds.
```

Plot log₂ segment mass generalized means vs segment length.

```
% The orders are hard-coded in the C++ program.
orders = [-Inf -10:10 Inf];
num_orders = numel(orders);
log2_generalized_means = readmatrix( ...
    segment_mass_log_multimean_file_name, ...
    'FileType','text','Delimiter',' ');
num_segment_lengths = size(log2_generalized_means,1);
% Segment lengths are consecutive powers of 2, starting with 2^0 = 1.
log2_segment_lengths = (0:num_segment_lengths-1)';

line_colors = jet(num_orders);
legend_items = cell(num_orders,1);
figure('Position',[0 0 600 600])
hold on
for o = 1:num_orders
    plot( log2_segment_lengths, log2_generalized_means(:,o), ...
        'Color', line_colors(o,:) )
    legend_items{o} = sprintf( 'q=%i', orders(o) );
end
hold off
legend(legend_items,'Location','northeastoutside')
xlabel('log2(walk segment length)')
ylabel('log2( q-mean(walk segment masses) )')
title( sprintf( ...
    'q-mean segment masses for MBRW (bias=%u, memory=%u)\n of %s', ...
    bias, memory, print_network_name ) )
```



Use generalized means to estimate generalized spectral dimensions.

We need to select a region of segment lengths over which to fit a line to each $\log_2(\text{q-mean segment mass})$ vs $\log_2(\text{segment length})$ curve. We start with minimum length 4, since the no-backtracking rule means all segments have mass at least 3. We stop at the last length less than the number of nodes in the network, since this is generally short enough that walk segments do not exhaust the network. For much longer segments, the growth rate of segment mass levels off due to finite network size.

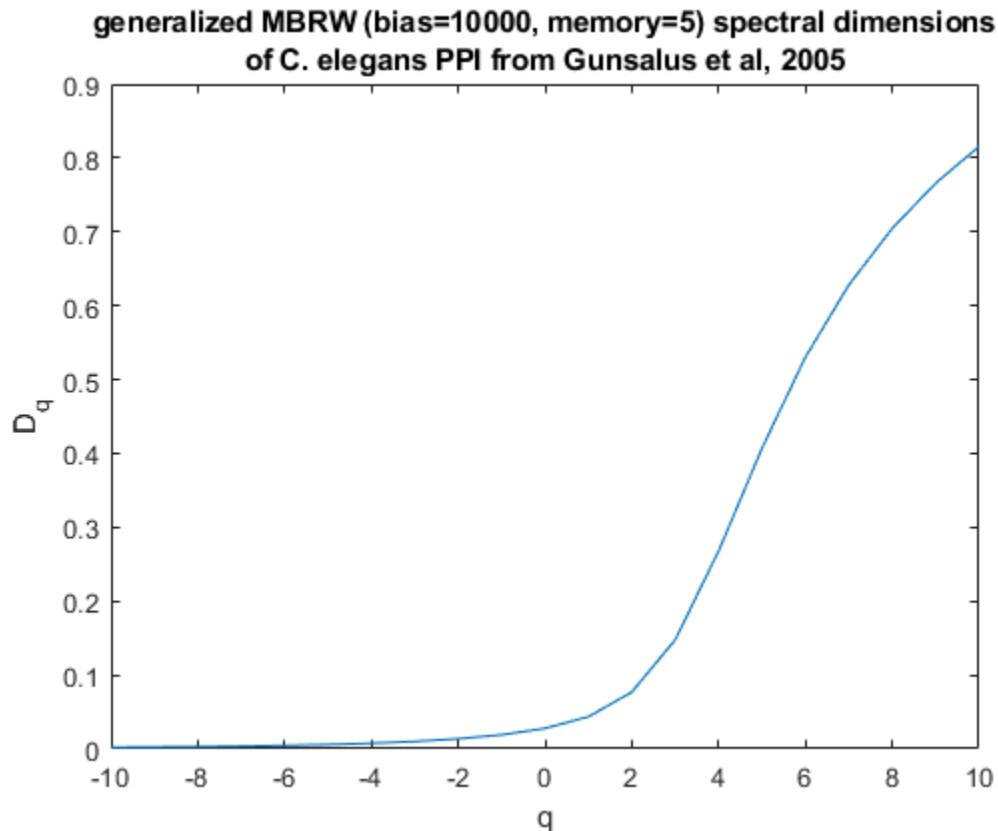
```
num_nodes = numnodes(G);
```



```
Dq = spectral_dimensions_v2(log2_generalized_means, 'length',
    num_nodes);
```

Plot the generalized means.

```
% Do not plot the values at infinity.
order_is_finite = ~isinf(orders);
finite_orders = orders(order_is_finite);
finite_order_Dqs = Dq(order_is_finite);
figure
plot(finite_orders, finite_order_Dqs)
xlabel('q')
ylabel('D_q')
title( sprintf( ...
    'generalized MBRW (bias=%u, memory=%u) spectral dimensions\n of %s
    ', ...
    bias, memory, print_network_name ) )
```



Take the range of finite-order generalized spectral dimensions.

We use this range as an order of multi-spectrality.

```
deltaDq = range(finite_order_Dqs);
```

```
fprintf('\x0394\x0044=%g\n',deltaDq)
```

```
#D=0.813504
```

Published with MATLAB® R2020b