

22B-Networks-I-Intro-Demo

November 30, 2017

1 Graph Analysis - I

The new import we are doing in this class is networkx:

<http://networkx.github.io/documentation/latest/tutorial/>

1.1 Basic graph concepts in NetworkX

Undirected Graphs

```
In [2]: g = nx.Graph()
```

Adding to the graph one node at a time

```
In [3]: g.add_node(1)
```

Adding multiple nodes at a time

```
In [4]: g.add_nodes_from([2,3])
```

Nodes are objects themselves

```
In [5]: g.add_node('ET')
g.nodes()
```

```
Out[5]: [1, 2, 3, 'ET']
```

Nodes can also be removed

```
In [6]: g.remove_node(1)
g.nodes()
```

```
Out[6]: [2, 3, 'ET']
```

Adding edges to the graph

```
In [7]: g.add_edge(1,2)
g.add_edge(3, 'ET')
g.add_edges_from([(2,3), (1,3)])
g.edges()
```

```
Out[7]: [(1, 2), (1, 3), (2, 3), (3, 'ET')]
```

```
In [8]: g.nodes()
```

```
Out[8]: [1, 2, 3, 'ET']
```

Removing edges

```
In [9]: g.remove_edge(1,2)
        g.edges()
```

```
Out[9]: [(1, 3), (2, 3), (3, 'ET')]
```

```
In [10]: g.nodes()
```

```
Out[10]: [1, 2, 3, 'ET']
```

Neighbors, degrees etc.

```
In [11]: g.neighbors(1)
```

```
Out[11]: [3]
```

```
In [12]: g.degree(1)
```

```
Out[12]: 1
```

Any networkx graph behaves like a Python dictionary with nodes as primary keys.

```
In [13]: g.add_node(1, time='5pm')
```

```
In [14]: g.node[1]['time']
```

```
Out[14]: '5pm'
```

```
In [15]: g.node[1] # Python dictionary
```

```
Out[15]: {'time': '5pm'}
```

The special edge attribute "weight" should always be numeric and holds values used by algorithms requiring weighted edges.

```
In [16]: g.add_edge(1, 2, weight=4.0 )
```

```
In [17]: g[1][2]['weight'] = 5.0 # edge already added
```

```
In [18]: g[1][2]
```

```
Out[18]: {'weight': 5.0}
```

Node and edge iterators.

```
In [19]: for node in g:
          print('nodeid: ', node, '\t degree:', g.degree(node))
```

```
nodeid: 1          degree: 2
nodeid: 2          degree: 2
nodeid: 3          degree: 3
nodeid: ET         degree: 1
```

```
In [20]: g.edges(data=True)
```

```
Out[20]: [(1, 2, {'weight': 5.0}), (1, 3, {}), (2, 3, {}), (3, 'ET', {})]
```

Directed Graphs.

Add the nodes from any container (a list, dict, set or even the lines from a file or the nodes from another graph).

```
In [21]: G = nx.DiGraph()
          G.add_node(1)
          G.add_nodes_from([2,3])
          G.add_nodes_from(range(100,110))
          H=nx.Graph()
          H.add_path([0,1,2,3,4,5,6,7,8,9])
          G.add_nodes_from(H)
```

```
In [22]: G.nodes()
```

```
Out[22]: [0,
          1,
          2,
          3,
          100,
          101,
          102,
          103,
          104,
          105,
          106,
          107,
          108,
          109,
          8,
          9,
          7,
          4,
          6,
          5]
```

G can also grow by adding edges

```
In [23]: G.add_edge(1, 2)
         G.add_edges_from([(1,2),(1,3)])
         G.add_edges_from(H.edges())
```

```
In [24]: G.edges()
```

```
Out[24]: [(0, 1),
          (1, 2),
          (1, 3),
          (2, 3),
          (3, 4),
          (8, 9),
          (7, 8),
          (4, 5),
          (6, 7),
          (5, 6)]
```

Attributes.

Each graph, node, and edge can hold key/value attribute pairs in an associated attribute dictionary (the keys must be hashable). By default these are empty, but can be added or changed using `add_edge()`, `add_node()` or direct manipulation of the attribute dictionaries named `graph`, `node` and `edge` respectively.

```
In [25]: G = nx.DiGraph(day="Friday")
         G.graph
```

```
Out[25]: {'day': 'Friday'}
```

Add node attributes using `add_node()`, `add_nodes_from()` or `G.node`

```
In [26]: G.add_node(1, time='5pm')
         G.add_nodes_from([3], time='2pm')
         print(G.node[1])
         G.node[1]['room'] = 714
         del G.node[1]['room'] # remove attribute
         G.nodes(data=True)
```

```
{'time': '5pm'}
```

```
Out[26]: [(1, {'time': '5pm'}), (3, {'time': '2pm'})]
```

Add edge attributes using `add_edge()`, `add_edges_from()`, subscript notation, or `G.edge`.

```
In [27]: G.add_edge(1, 2, weight=4.7 )
         G.add_edges_from([(3,4),(4,5)], color='red')
         G.add_edges_from([(1,2,{'color':'blue'})], (2,3,{'weight':8})))
         G[1][2]['weight'] = 4.7
         G.edge[1][2]['weight'] = 4
         G.edges(data=True)
```

```
Out[27]: [(1, 2, {'color': 'blue', 'weight': 4}),
          (2, 3, {'weight': 8}),
          (3, 4, {'color': 'red'}),
          (4, 5, {'color': 'red'})]
```

Many common graph features allow python syntax to speed reporting.

```
In [28]: 1 in G      # check if node in graph
```

```
Out[28]: True
```

```
In [29]: [n for n in G if n<3]    # iterate through nodes
```

```
Out[29]: [1, 2]
```

```
In [30]: len(G)    # number of nodes in graph
```

```
Out[30]: 5
```

```
In [31]: print(G[1]) # adjacency dict keyed by neighbor to edge attributes
          # Note: you should not change this dict manually!
```

```
{2: {'weight': 4, 'color': 'blue'}}
```

Iterating over the edges of a graph

```
In [32]: for n,nbrsdict in G.adjacency_iter():
          for nbr,eattr in nbrsdict.items():
              if 'weight' in eattr:
                  print (n,nbr,eattr['weight'])
```

```
1 2 4
2 3 8
```

or

```
In [33]: [ (u,v,edata['weight']) for u,v,edata in G.edges(data=True) if 'weight' in edata ]
```

```
Out[33]: [(1, 2, 4), (2, 3, 8)]
```

1.2 Visualizing Graphs

Visualizing a network can be quite difficult. There are many strategies that are used to draw networks in ways that communicate as much insight as possible.

```
In [34]: Ggml = nx.read_gml('data/polblogs.gml')
```

This is a directed network of hyperlinks between weblogs on US politics, recorded in 2005. Accessible [here](#).

```
In [35]: print(len(Ggml.nodes()))
         print(len(Ggml.edges()))
```

```
1490
19015
```

This is a fairly large network to try to visualize.

The standard `networkx` routine uses what is called a 'spring' layout.

Each edge has a weight parameter. The layout routine fixes a spring of that length between the nodes, and a repulsive force between each pair of nodes, and then lets the set of all forces reach its minimum energy state.

This is a kind of minimal distortion in a least-squares sense.

```
In [75]: with sns.axes_style('white'):
         fig = plt.subplots(1, figsize=(12,8))
         nx.draw_networkx(Ggml, edge_color='#a4a4a4',
                        node_size=50, with_labels=False, arrows=False)
         plt.axis('off')
```



The other kinds of possible layouts are: * `circular_layout` - position nodes on a circle * `random_layout` - position nodes randomly in the unit square * `shell_layout` - position nodes in concentric circles * `spectral_layout` - uses the eigenvectors of the graph Laplacian

Note that `networkx` is not intended as a sophisticated graph visualization package. There are more sophisticated packages available that do much more. Some examples include * `graphviz` * `gephi` * `cytoscape`

Looking for Clusters.

This graph models American football games between NCAA Div IA colleges in Fall 2000 (available [here](#)).

Each vertex represents a football team, which belongs to a specific conference (Big Ten, Conference USA, Pac-10, etc.).

An edge between two vertices v_1 and v_2 means that the two teams played each other; the weight of the edge (v_1, v_2) is equal to the number of times they played each other.

Data from M. Girvan and M. E. J. Newman, *Community structure in social and biological networks*, Proc. Natl. Acad. Sci. USA 99, 7821-7826 (2002).

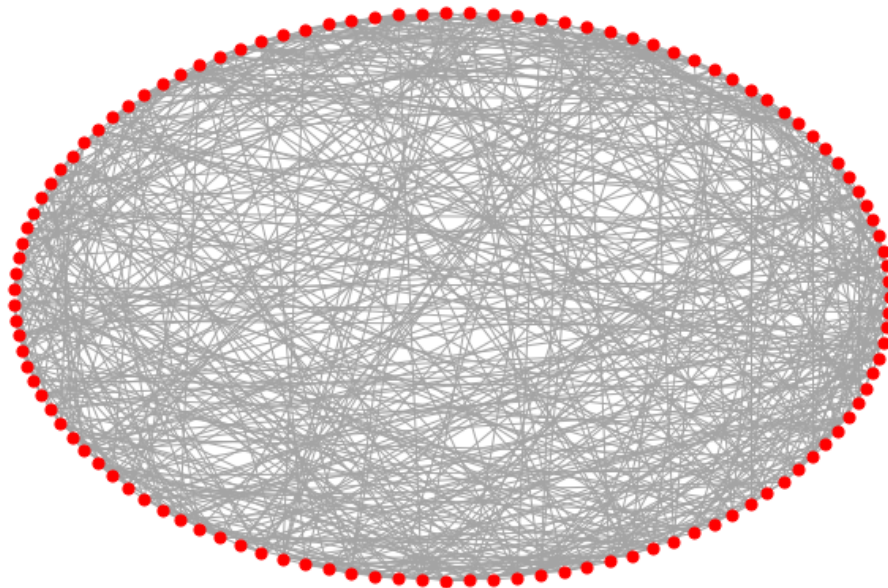
```
In [37]: with open('data/football.txt', 'r') as f:
         football = nx.parse_edgelist(f, comments='#', nodetype=int, data=False)
```

```
In [38]: print('This network has {} nodes and {} edges'.format(len(football.nodes()), len(football.edges())))
```

This network has 115 nodes and 613 edges

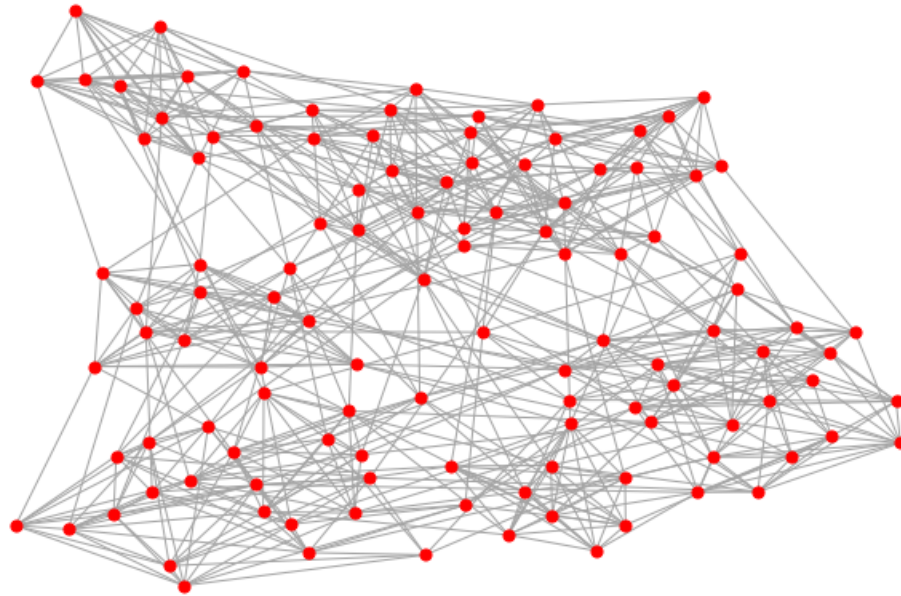
Let's start with a circular layout:

```
In [39]: with sns.axes_style('white'):
         fig = plt.subplots(1, figsize=(12,8))
         nx.draw_networkx(football, pos=nx.circular_layout(football), edge_color='#a4a4a4',
         plt.axis('off'))
```



Now, let's compare the standard spring model layout:

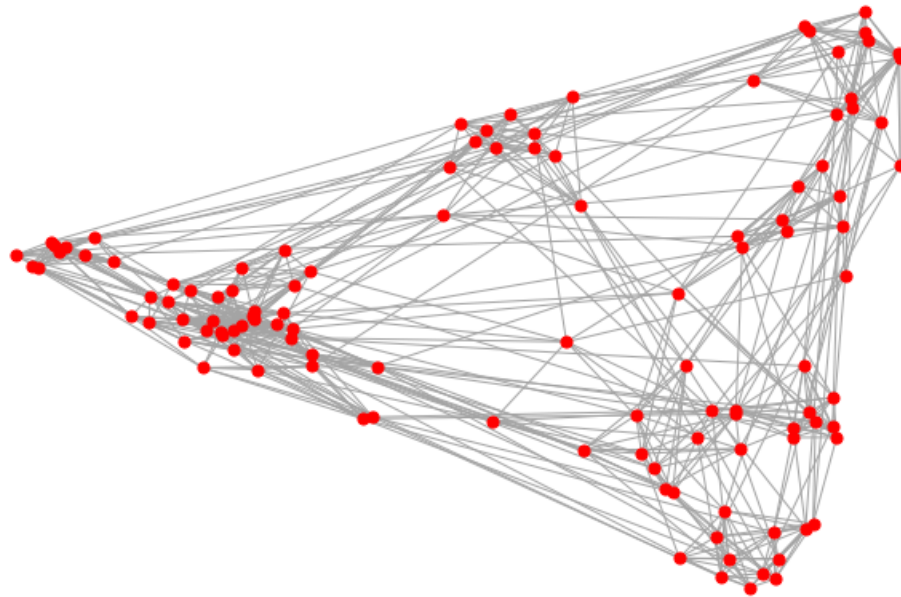
```
In [77]: with sns.axes_style('white'):
fig = plt.subplots(1, figsize=(12,8))
nx.draw_networkx(football, edge_color='#a4a4a4', node_size=50, with_labels=False)
plt.axis('off')
```



Notice how the spring layout tends to bring clusters of densely connected nodes close to each other.

Finally, we can try the spectral layout:

```
In [41]: with sns.axes_style('white'):
fig = plt.subplots(1, figsize=(12,8))
nx.draw_networkx(football, pos=nx.spectral_layout(football), edge_color='#a4a4a4',
plt.axis('off')
```

The spectral layout enhances the clustering of densely connected groups.

Generating graphs using the routines already available in python (for small data) Networkx has a wealth of data-generation routines that can be found here:

<https://networkx.github.io/documentation/latest/reference/generators.html>

This is the function that generates the Zachary's Karate club network data

```
In [42]: kn=nx.karate_club_graph()
```

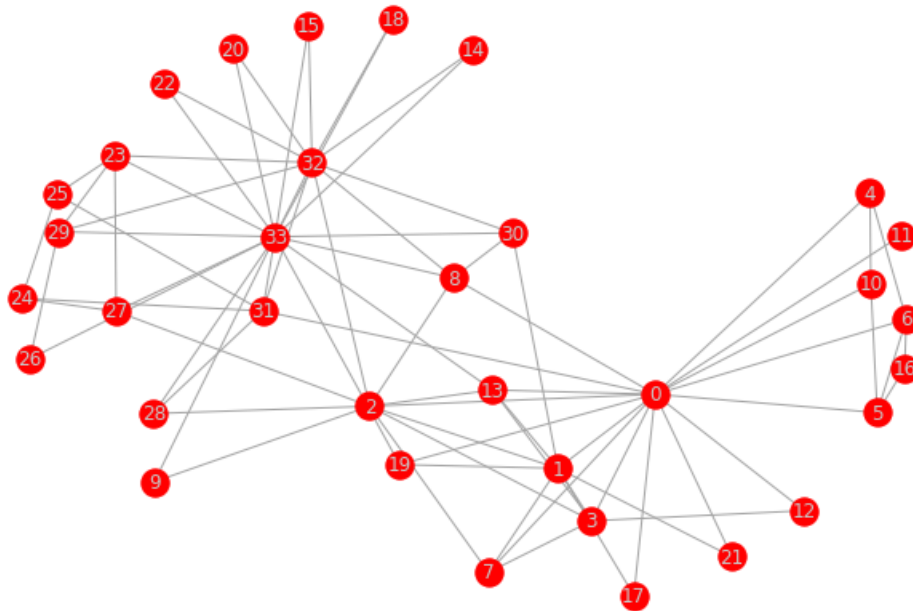
```
In [43]: num_nodes = kn.number_of_nodes()
          print('number of nodes: ' + str(num_nodes))
          num_edges = kn.number_of_edges()
          print('number of edges: ' + str(num_edges))
```

```
number of nodes: 34
```

```
number of edges: 78
```

Visualizing the network:

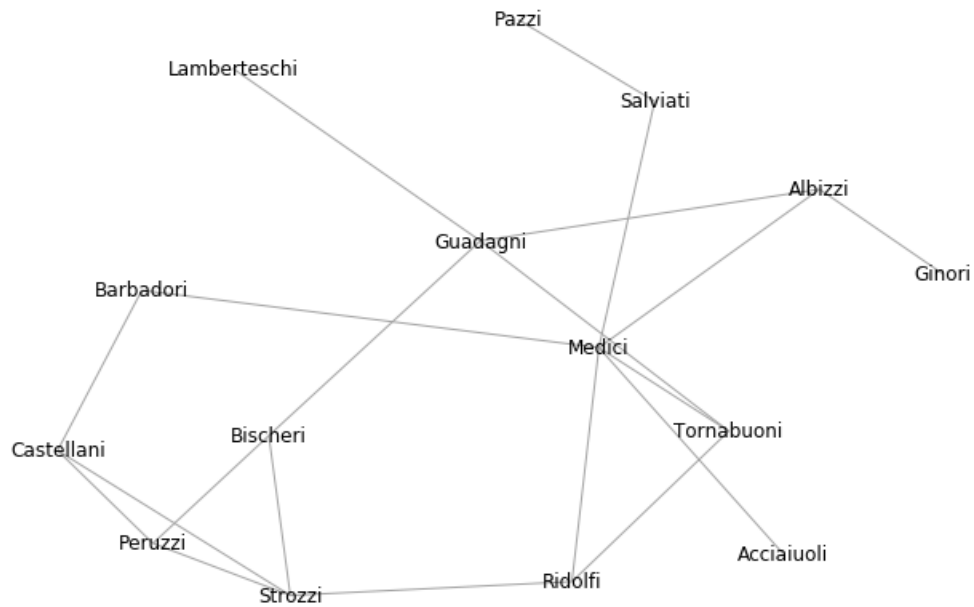
```
In [44]: with sns.axes_style('white'):
          fig = plt.subplots(1, figsize=(12,8))
          nx.draw_networkx(kn, edge_color='#a4a4a4', with_labels=True, font_color='#cacaca')
          plt.axis('off')
```



```
In [45]: fl = nx.florentine_families_graph()
num_nodes = fl.number_of_nodes()
print('number of nodes: ' + str(num_nodes))
num_edges = fl.number_of_edges()
print('number of edges: ' + str(num_edges))
with sns.axes_style('white'):
    fig = plt.subplots(1, figsize=(12,8))
    nx.draw_networkx(fl, edge_color='#a4a4a4', node_size=0, with_labels=True)
    plt.axis('off')
```

number of nodes: 15

number of edges: 20



Erdos-Renyi random graphs.

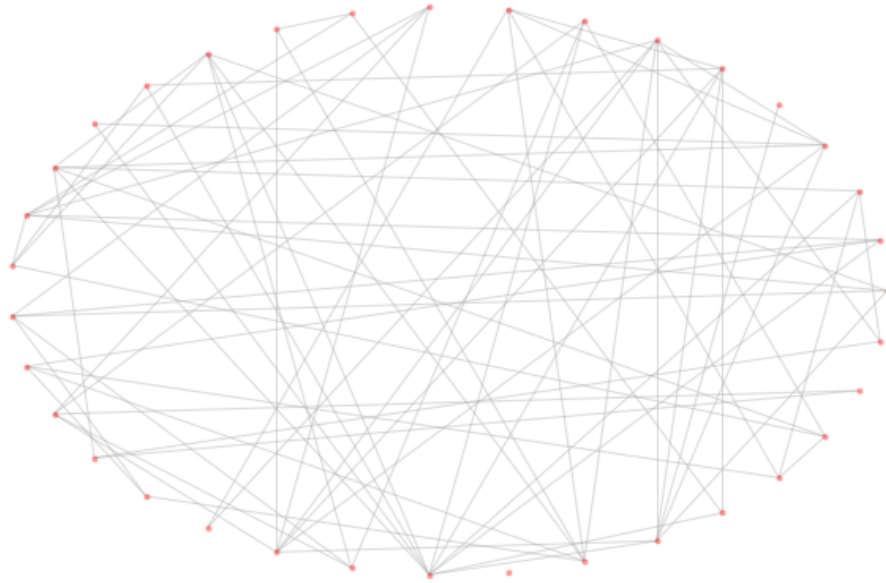
Recall that an Erdos-Renyi random graph has two parameters: * n , the number of nodes in the graph, and * p , the probability that any given pair of nodes is connected by an edge.

These graphs are sometimes called $G(n, p)$ graphs.

```

In [46]: er2=nx.erdos_renyi_graph(35,0.15)
         with sns.axes_style('white'):
             fig = plt.subplots(1, figsize=(12,8))
             nx.draw_networkx(er2, node_size=15, edge_color='#a4a4a4', pos=nx.circular_layout(er2))
             plt.axis('off')

```



Let's look at a bigger E-R graph:

```
In [47]: er=nx.erdos_renyi_graph(1000,0.15)
```

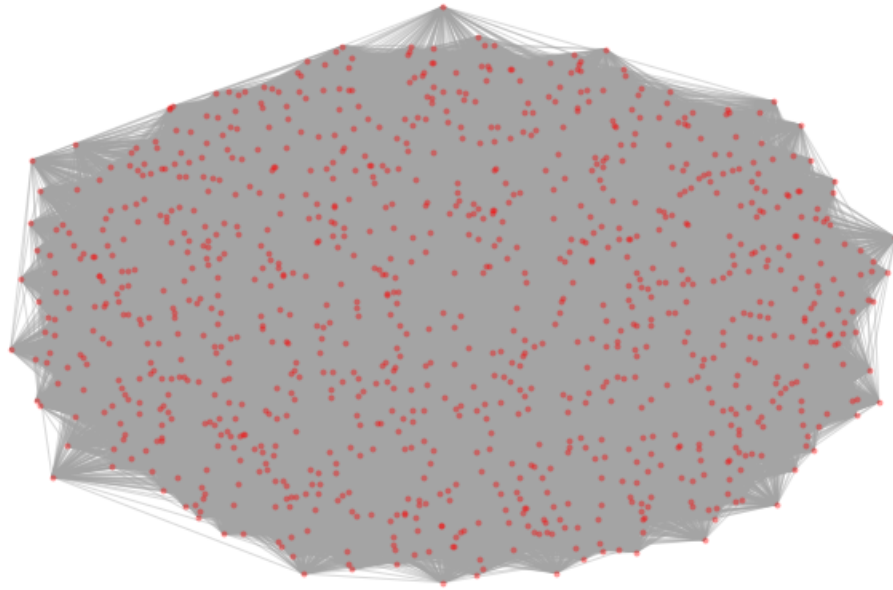
```
In [48]: er=nx.erdos_renyi_graph(1000,0.15)
print("Number of nodes in the random graph: ", er.number_of_nodes())
print("Number of edges in the random graph: ", er.number_of_edges())
```

Number of nodes in the random graph: 1000

Number of edges in the random graph: 74474

Visualizing with spring model:

```
In [49]: with sns.axes_style('white'):
fig = plt.subplots(1, figsize=(12,8))
nx.draw_networkx(er, node_size=15, edge_color='#a4a4a4', with_labels=False, alpha=.
plt.axis('off')
```



Are there clusters in this graph?

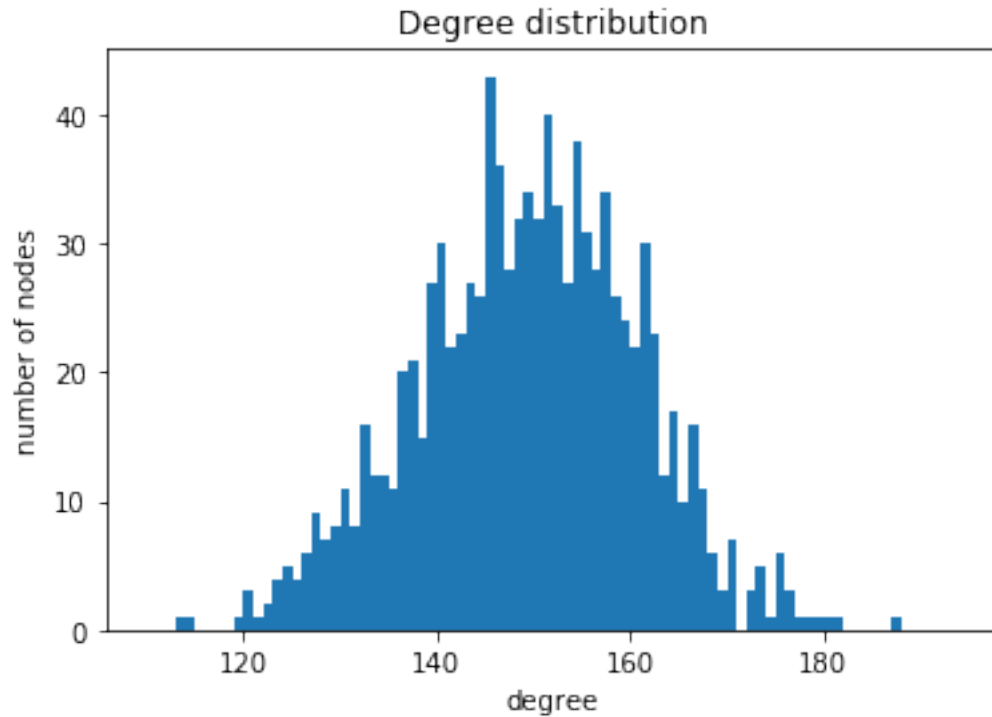
Degree distribution

```
In [50]: degree_sequence=sorted(nx.degree(er).values(),reverse=True)
         dmax=max(degree_sequence)
         dmax
```

```
Out[50]: 187
```

```
In [51]: h,bins,patches = plt.hist(degree_sequence,bins=range(110,195))
         plt.title("Degree distribution")
         plt.xlabel("degree")
         plt.ylabel("number of nodes")
```

```
Out[51]: <matplotlib.text.Text at 0x11adbef98>
```



Connected Components Two nodes of a graph belong in the same connected component if there is a path of edges of the graph that connects these two nodes.

```
In [52]: cc= nx.connected_components(er)
         print(type(cc))
         print([len(s) for s in cc])
```

```
<class 'generator'>
[1000]
```

Clustering coefficient As we discussed previously, the **clustering coefficient of a node** is defined as the number of possible triangles centered in this node, divided by the total number of possible triangles in which this node can participate in.

Formally, the clustering coefficient of a node u is defined as

$$c_u = \frac{2T(u)}{d(u)(d(u) - 1)},$$

where $T(u)$ is the number of triangles through node u and $d(u)$ is the degree of node u .

For more details for weighted graphs etc see:

<http://networkx.lanl.gov/reference/generated/networkx.algorithms.cluster.clustering.html#networkx.algorithms.cluster.clustering>

The **average clustering coefficient** is the average clustering coefficient of all the nodes in the graph.

http://networkx.lanl.gov/reference/generated/networkx.algorithms.cluster.average_clustering.html#networkx.algorithms.cluster.average_clustering

```
In [79]: ccall = nx.clustering(er)
```

```
print(ccall)
```

```
clustering_coefficient = nx.average_clustering(er)
```

```
clustering_coefficient
```

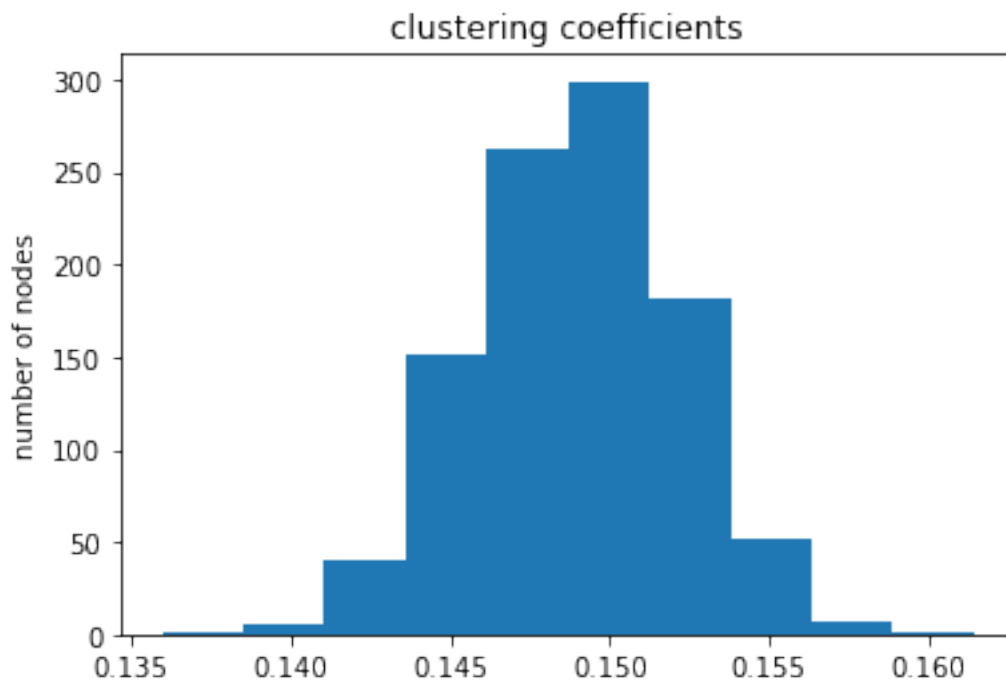
```
{0: 0.1503267973856209, 1: 0.14306712395884372, 2: 0.149908446779715, 3: 0.14006150405367626, 4:
```

```
Out[79]: 0.14894999159618452
```

```
In [80]: h,bins,patches = plt.hist(list(nx.clustering(er).values()))
```

```
plt.title('clustering coefficients')
```

```
plt.ylabel("number of nodes");
```



Triangles

```
In [81]: print(nx.triangles(er,0))
```

```
#print(nx.triangles(er))
```

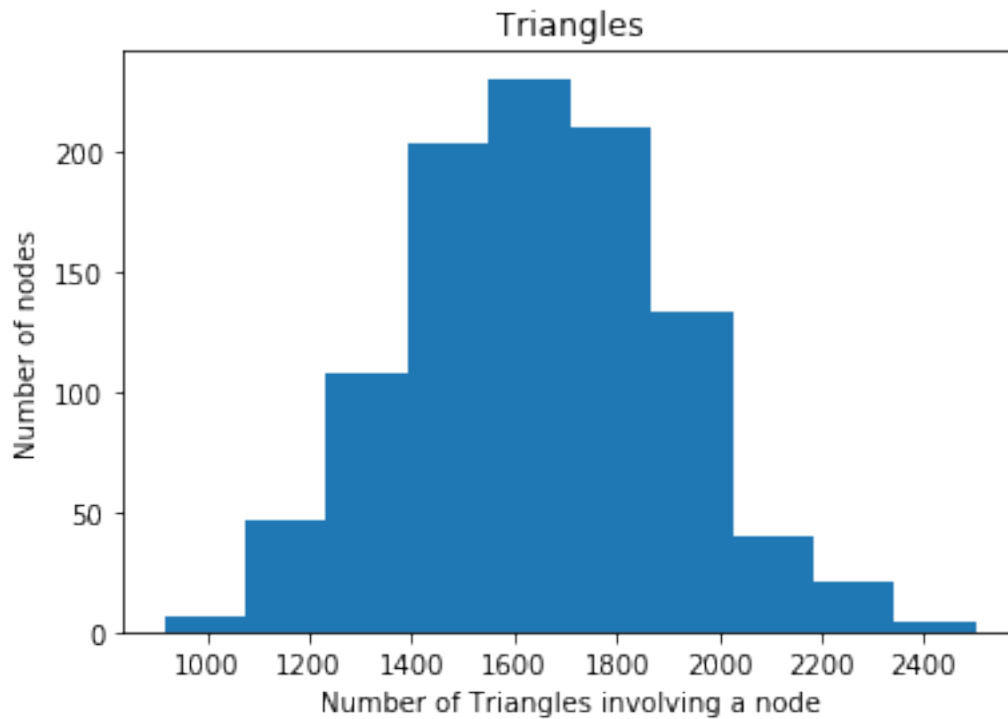
```
h,bins, patches = plt.hist(list(nx.triangles(er).values()))
```

```
plt.title('Triangles')
```

```
plt.xlabel("Number of Triangles involving a node")
```

```
plt.ylabel("Number of nodes");
```

1380



Diameter and average path length The **diameter** of a graph is defined as the largest shortest path between any two nodes in the graph

```
In [57]: print(nx.diameter(er))
```

2

The **average shortest path length** of a graph is defined as the average of all shortest path lengths in the graph

http://networkx.lanl.gov/reference/generated/networkx.algorithms.shortest_paths.generic.average_shortest_path_length.html

```
In [58]: print(nx.average_shortest_path_length(er))
```

1.850902902902903

Watts-Strogatz graphs.

```
In [59]: ws=nx.watts_strogatz_graph(500,5,0.1)
         print_cc_sizes(ws)
```

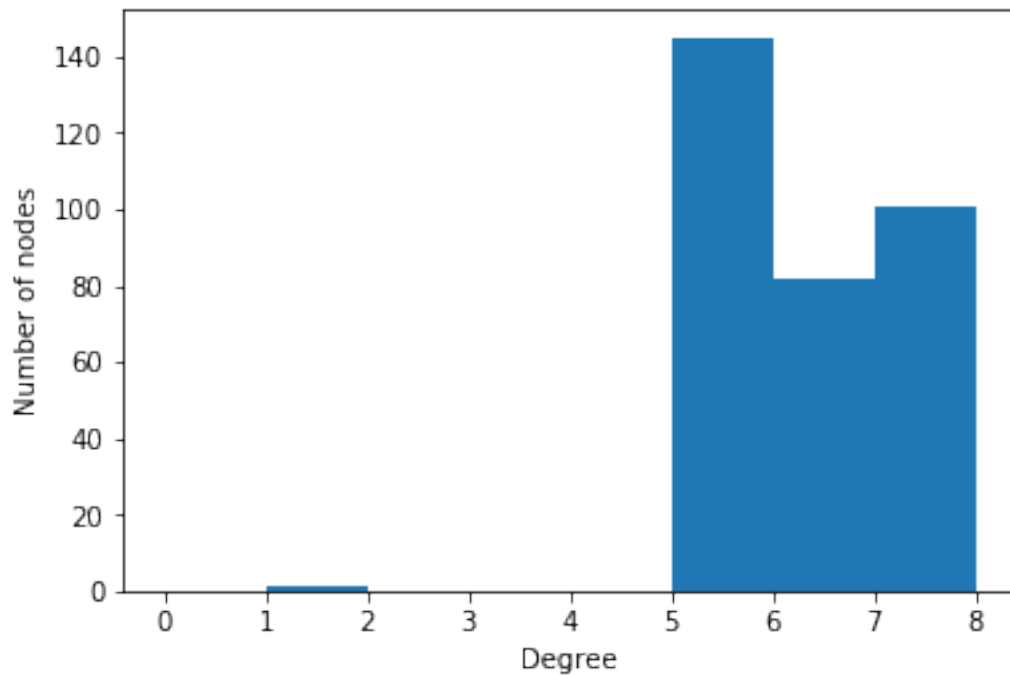
[500]

Degree distribution

```
In [60]: degree_sequence=sorted(nx.degree(ws).values(),reverse=True)
         dmax=max(degree_sequence)
         dmax
```

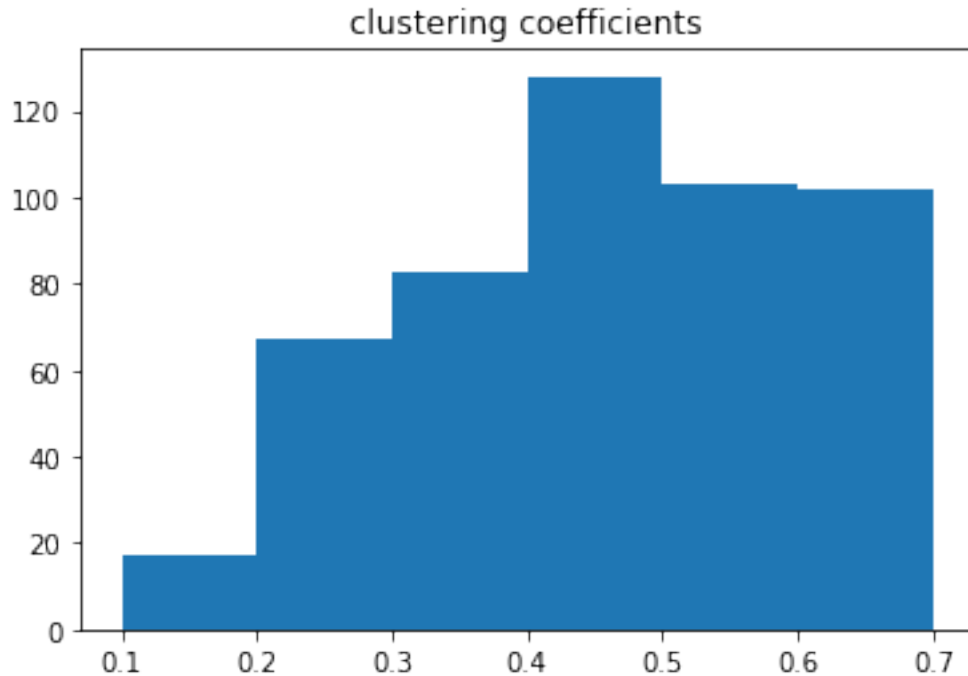
```
Out[60]: 6
```

```
In [82]: h,bins,patches = plt.hist(degree_sequence,bins=range(9))
         plt.xlabel("Degree")
         plt.ylabel("Number of nodes");
```



Clustering coefficient

```
In [83]: h,bins,patches = plt.hist(list(nx.clustering(ws).values()),bins=6)
         plt.title('clustering coefficients');
```

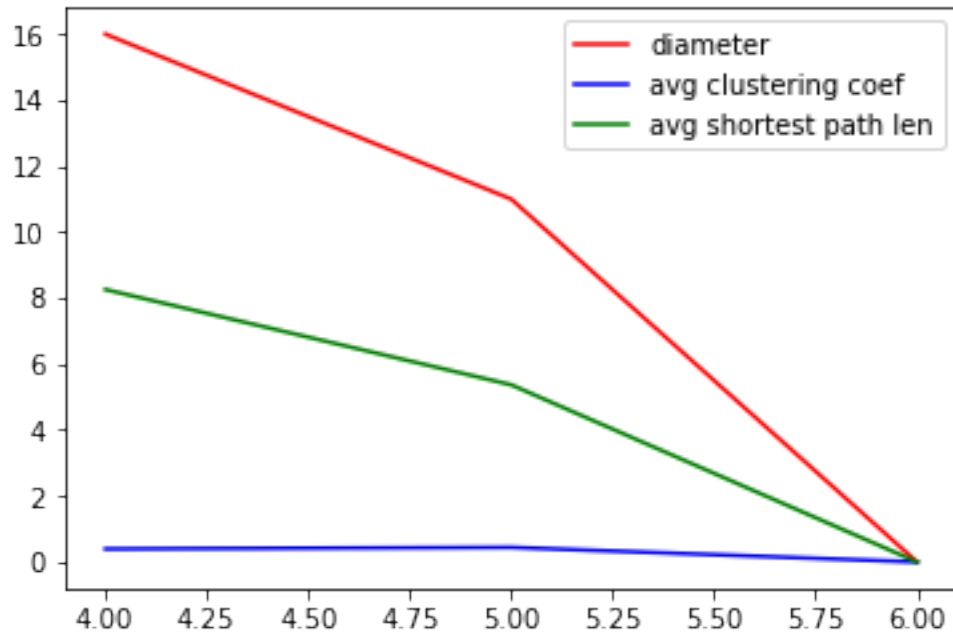


Average path length and diameter

```
In [63]: print('Diameter:', (nx.diameter(ws)))
         print('Average shortest path length:', (nx.average_shortest_path_length(ws)))
         print('Average clustering coefficient:', (nx.average_clustering(ws)))
```

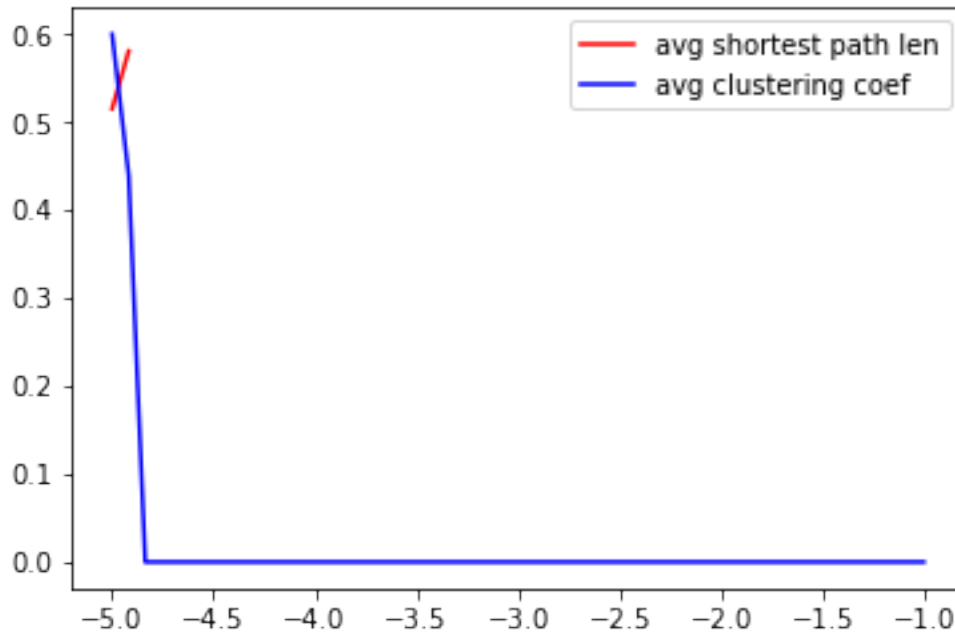
```
Diameter: 15
Average shortest path length: 7.657779559118237
Average clustering coefficient: 0.3722
```

```
In [88]: r = range(4,7)
         d = np.zeros(len(r))
         cc = np.zeros(len(r))
         pl = np.zeros(len(r))
         index = 0
         for i in r:
             ws=nx.watts_strogatz_graph(500,i,0.1)
             d[index] = nx.diameter(ws)
             cc[index] = nx.average_clustering(ws)
             pl[index] = nx.average_shortest_path_length(ws)
             index+=1
         plt.plot(r,d,'r',label='diameter')
         plt.plot(r,cc,'b',label='avg clustering coef')
         plt.plot(r,pl,'g',label='avg shortest path len')
         plt.legend(loc='best');
```



```
In [95]: r = np.logspace(-5.,-1.,50)
d = np.zeros(len(r))
cc = np.zeros(len(r))
pl = np.zeros(len(r))
index = 0
for i in r:
    ws=nx.watts_strogatz_graph(100,6,i)
    d[index] = nx.diameter(ws)
    cc[index] = nx.average_clustering(ws)
    pl[index] = nx.average_shortest_path_length(ws)
    index+=1
plt.plot(np.log10(r),pl/d,'r',label='avg shortest path len')
plt.plot(np.log10(r),cc,'b',label='avg clustering coef')
plt.legend(loc='best');
```

/Users/crovella/anaconda3/lib/python3.5/site-packages/ipykernel/__main__.py:12: RuntimeWarning:



In [94]: r

```
Out[94]: array([ 0.0001      ,  0.00011514,  0.00013257,  0.00015264,  0.00017575,
                 0.00020236,  0.000233   ,  0.00026827,  0.00030888,  0.00035565,
                 0.00040949,  0.00047149,  0.00054287,  0.00062506,  0.00071969,
                 0.00082864,  0.0009541 ,  0.00109854,  0.00126486,  0.00145635,
                 0.00167683,  0.0019307 ,  0.002223   ,  0.00255955,  0.00294705,
                 0.00339322,  0.00390694,  0.00449843,  0.00517947,  0.00596362,
                 0.00686649,  0.00790604,  0.00910298,  0.01048113,  0.01206793,
                 0.01389495,  0.01599859,  0.0184207 ,  0.02120951,  0.02442053,
                 0.02811769,  0.03237458,  0.03727594,  0.04291934,  0.04941713,
                 0.05689866,  0.06551286,  0.0754312 ,  0.08685114,  0.1          ])
```

Experimenting with Barabasi-Albert graphs <http://networkx.lanl.gov/reference/generated/networkx.generator>

```
In [65]: ba=nx.barabasi_albert_graph(500,5)
         print_cc_sizes(ba)
```

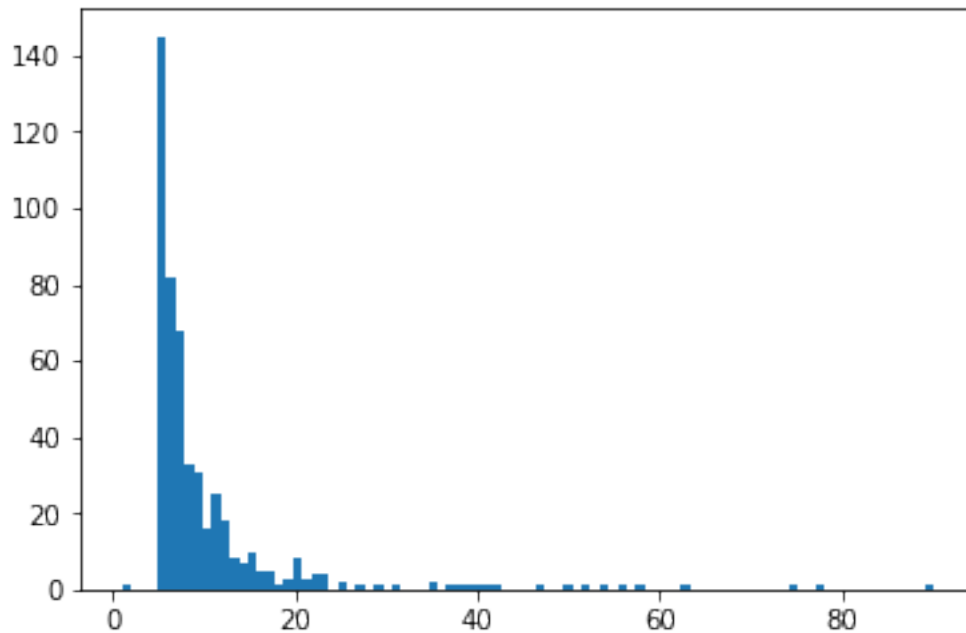
[500]

Degree distribution

```
In [66]: degree_sequence=sorted(nx.degree(ba).values(),reverse=True)
         dmax=max(degree_sequence)
         dmax
```

Out [66]: 90

```
In [67]: h,bins,patches = plt.hist(degree_sequence,bins=dmax)
```



```
In [68]: hmax=max(h)
plt.axis([1,dmax,1,hmax]) # set ranges
#x=compress(h,bins)      # remove bins with zero entries
#y=compress(h,h)         # remove corresponding entries
x=bins.compress(h)
y=h.compress(h)
plt.loglog(x,y,'bo')
plt.title("Degree distribution")
plt.xlabel("degree")
plt.ylabel("number of nodes")
plt.show()
```

