

Project 3: Graph Algorithms

Overview:

Project 3 involves testing different graph algorithms experimentally to determine the properties of models of real-world networks. In this project we are going to test 3 different algorithms:

- Diameter Algorithm – Determine the length of longest path between any two vertices in the graph.
- Clustering-Coefficient – Determine the ratio of three times the number of triangles over the number of 2-edge paths.
- Degree-Distribution – For each possible degree, determine the number of vertices with that degree.

And we are going to test these algorithms on different of random graphs based on two different models, one is the Erdős–Rényi model and the other is the Babarási-Albert model.

Erdős – Rényi Model:

Description:

- $G(n,p)$: - Graph G with n nodes, with every pair connected individually with probability p and average degree $d = (n-1)*p \sim np$

Implementation:

Pseudocode:

ALG. 1: $\mathcal{G}(n,p)$

Input: number of vertices n , edge probability $0 < p < 1$

Output: $G = (\{0, \dots, n-1\}, E) \in \mathcal{G}(n,p)$

$E \leftarrow \emptyset$

$v \leftarrow 1; w \leftarrow -1$

while $v < n$ **do**

 draw $r \in [0,1)$ uniformly at random

$w \leftarrow w + 1 + \lceil \log(1-r) / \log(1-p) \rceil$

while $w \geq v$ **and** $v < n$ **do**

$w \leftarrow w - v; v \leftarrow v + 1$

if $v < n$ **then** $E \leftarrow E \cup \{v, w\}$

C++ implementation:

```
1 Graph create_ER_graph(int n) {  
2  
3     Graph erdos_renyi = make_graph(n, vector<int> {}, vector<int> {});  
4     double p = (2 * (log(n))) / n;  
5     int v = 1;  
6     int w = -1;  
7     double r;  
8  
9     while (v < n) {  
10        mt19937 seed2 = get_seed2();  
11        r = uniform_int_distribution<int>(0, 1)(seed2);  
12        int x = (floor((log(1 - r) / log(1 - p))));  
13        w = w + 1 + x;  
14        while (w >= v and v < n) {  
15            w = w - v;  
16            v = v + 1;  
17        }  
18        if (v <= n) {  
19            erdos_renyi.add_edge(v, w);  
20        }  
21    }  
22  
23    return erdos_renyi;  
24 }
```

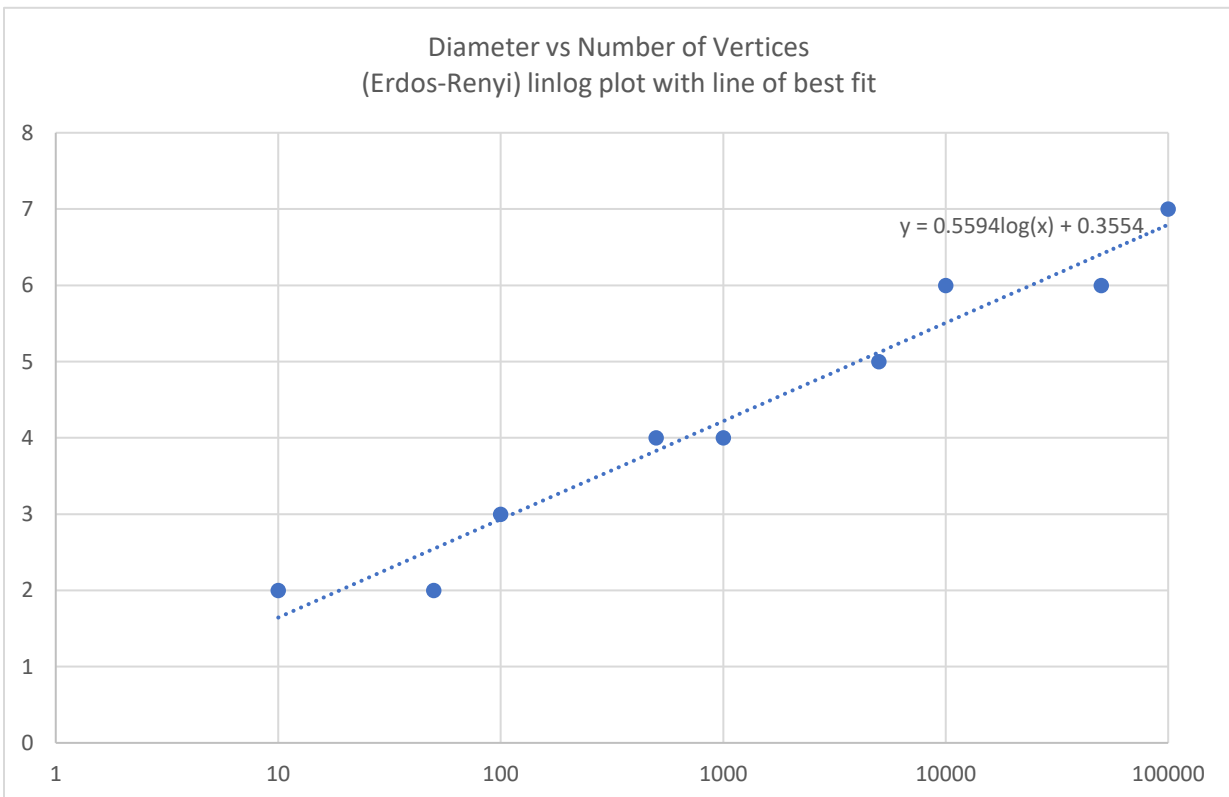
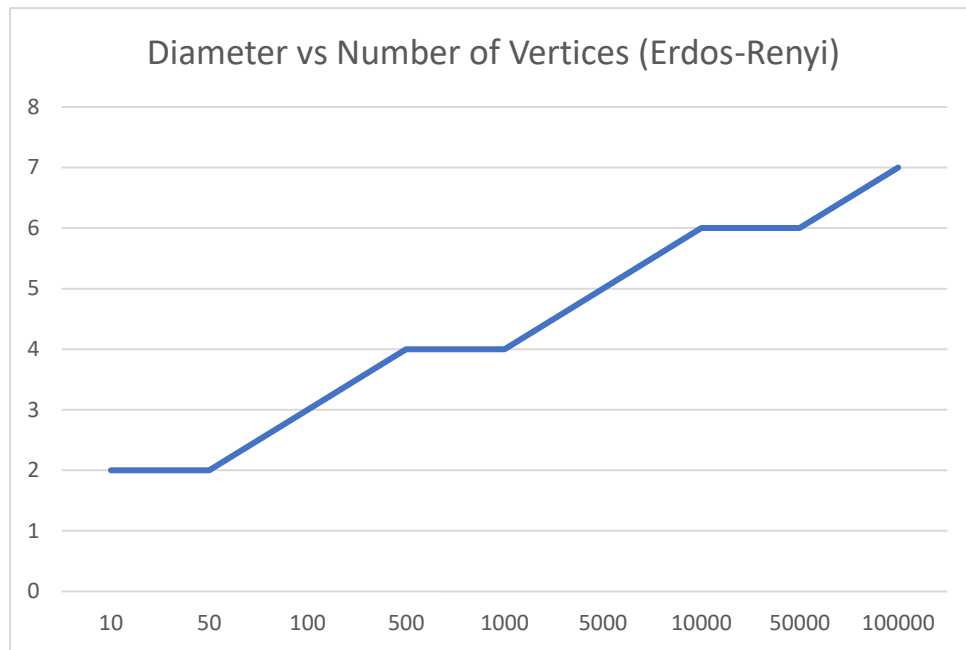
Diameter:

Implementation (using Heuristic Idea 2):

- Use a random generator to get the start node
- Using BFS from the start, traverse the graph
- BFS has been modified to return the node in the search for the maximum distance between that node and every other node in the graph
- Repeat the BFS till the node with the maximum distance has been found

C++ Code:

```
1 int get_diameter(Graph graph) {  
2     int size = graph.get_num_nodes();  
3     int diameter = 0;  
4  
5     if (size > 1) {  
6         int max = 0;  
7         mt19937 seed = get_seed();  
8         int ind_start = get_generator(1, size)(seed);  
9         Node start(ind_start);  
10        tuple < Node, int > furthest = make_tuple(start, max);  
11  
12        do {  
13            diameter = get < 1 > (furthest);  
14            furthest = breadth_first_search(graph, get < 0 > (furthest));  
15        } while (get < 1 > (furthest) > diameter);  
16    }  
17  
18    return diameter;  
19 }
```



From the above plots we can see that the diameter increases linearly as the number of vertices increased from 10 to 10^5 . And from the lin-log plot we can say that the diameter: $f(d) = O(n)$ which is proportional to $\log n$.

Clustering Coefficient:

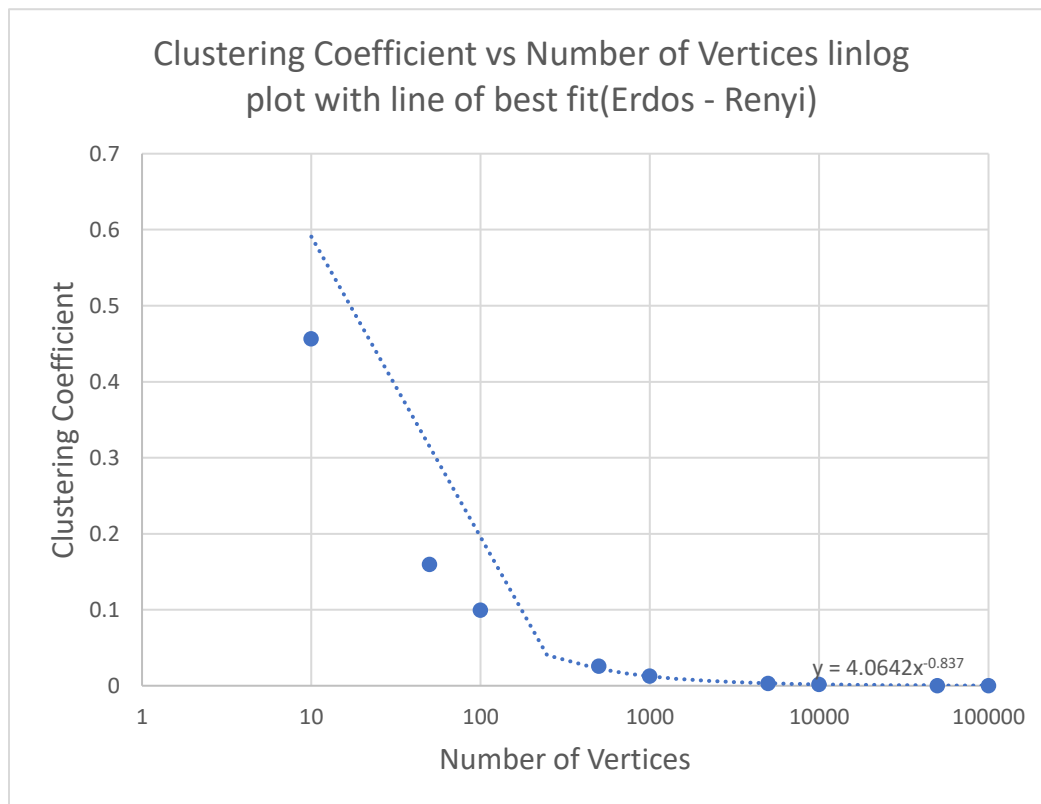
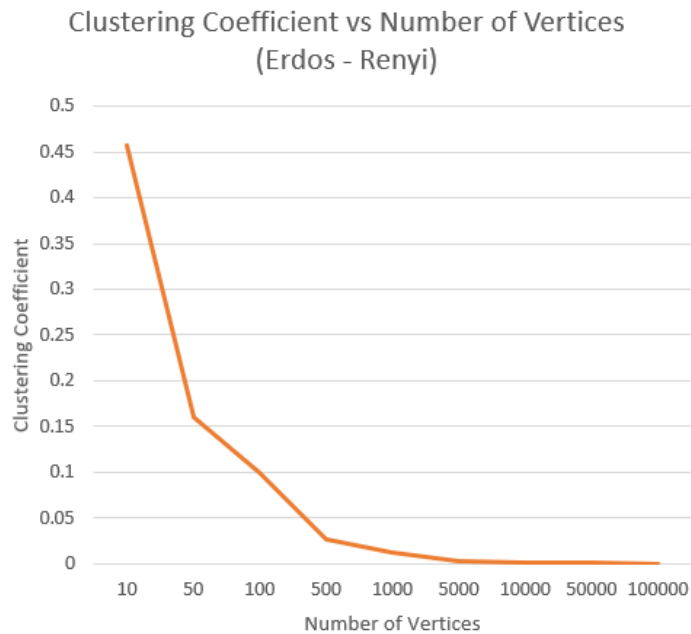
Implementation: - Return $3 \cdot (\text{Number of triangles}) / \text{Number of 2 edge paths}$

Counting Triangles implementation:

1. Initialize an output list, L , to be empty.
 2. Compute a number, d_v , for each vertex v in G , which is the number of neighbors of v that are not already in L . Initially, d_v is just the degrees of
 3. Initialize an array D such that $D[i]$ contains a list of the vertices v that are already in L for which $d_v = i$.
 4. Let N_v be a list of the neighbors of v that come before v in L . Initially, N_v empty for every vertex v .
 5. Initialize k to 0.
 6. Repeat n times:
 - Let i be the smallest index such that $D[i]$ is nonempty.
 - Set k to $\max(k, i)$.
 - Select a vertex v from $D[i]$. Add v to the beginning of L and remove it from $D[i]$. Mark v as being in L (e.g., using a hash table, H_L).
 - For each neighbor w of v not already in L (you can check this using H_L):
 - Subtract one from d_w
 - Move w to the cell of D corresponding to the new value of d_w , i.e., $D[d_w]$
 - Add w to N_v
- Compute a d -degeneracy ordering of the vertices, e.g., using the algorithm of the previous slide.
 - Process the vertices according to this ordering, L :
 - For each vertex, v :
 - For each pair of vertices, u and w , adjacent to v and earlier in the ordering, i.e., u and w are in the list N_v from the degeneracy algorithm:
 - If (u, w) is an edge in the graph, then add one to the triangle count.

C++ Code:

```
1= int count_triangles(Graph g) {
2   int triangles = 0;
3   int degeneracy = 0;
4   int num_nodes = g.get_num_nodes();
5
6   int deg[num_nodes];
7   vector< int > N[num_nodes];
8
9   stack< int > L; //Initialize an output list, L, to be empty.
10  map< int, unordered_set< int >> D;
11
12  //Initialize an array D such that D[i] contains a list of the vertices v that are not already in L for which dv
13  for (int i = 1; i <= num_nodes; ++i) {
14      Node n(i);
15      int degree = g.get_neighbors(n).size();
16      D[i].insert(i);
17      deg[i - 1] = degree;
18  }
19  //Select a vertex v from D[i]. Add v to the beginning of L and remove it from D[i]. Mark v as being in L (e.g.,
20  unordered_set< int > H;
21  for (int ind = 0; ind < num_nodes; ind++) {
22      int min = -1;
23      for (auto j = D.begin(); j != D.end(); j++) {
24          if (!j->second.empty()) {
25              min = j->first;
26              break;
27          }
28      }
29      unordered_set< int > nlist = D[min];
30      degeneracy = max(degeneracy, min);
31
32      int nid = *nlist.begin();
33      L.push(nid);
34      H.insert(nid);
35
36      D[min].erase(nid);
37      vector< Node > nlist = g.get_neighbors(Node(nid));
38      /*
39      - For each neighbor w of v not already in L (you can check this using HL):
40      • Subtract one from dw
41      • Move w to the cell of D corresponding to the new value of dw, i.e., D[dw]
42      • Add w to Nv
43      */
44      for (int x = 0; x < nlist.size(); x++) {
45          Node nwt = nlist[x];
46          int nodeID = nwt.id;
47          if (H.find(nodeID) == H.end()) {
48              int cur = deg[nodeID - 1]--;
49              D[cur].erase(nodeID);
50              D[cur - 1].insert(nodeID);
51              N[nid - 1].push_back(nodeID);
52          }
53      }
54  }
55
56  //Count the triangles
57  while (!L.empty()) {
58      int top = L.top();
59      vector< int > nodes = N[top - 1];
60      int num_nodes2 = nodes.size();
61      for (int a = 0; a < num_nodes2 - 1; a++) {
62          for (int b = a + 1; b < num_nodes2; b++) {
63              if (g.is_adjacent(Node(nodes[a]), Node(nodes[b]))) {
64                  triangles++;
65              }
66          }
67      }
68      L.pop();
69  }
70
71  return triangles;
}
```



From the above plots we can see with the equation $y = 4.0642x^{-0.837}$ that the clustering coefficient decreases proportionally as the number of vertices increased from 10 to 10^5 so we can say that the clustering coefficient: $f(C) = O(1/n)$ which is proportional to $\log n$.

Degree Distribution:

Implementation:

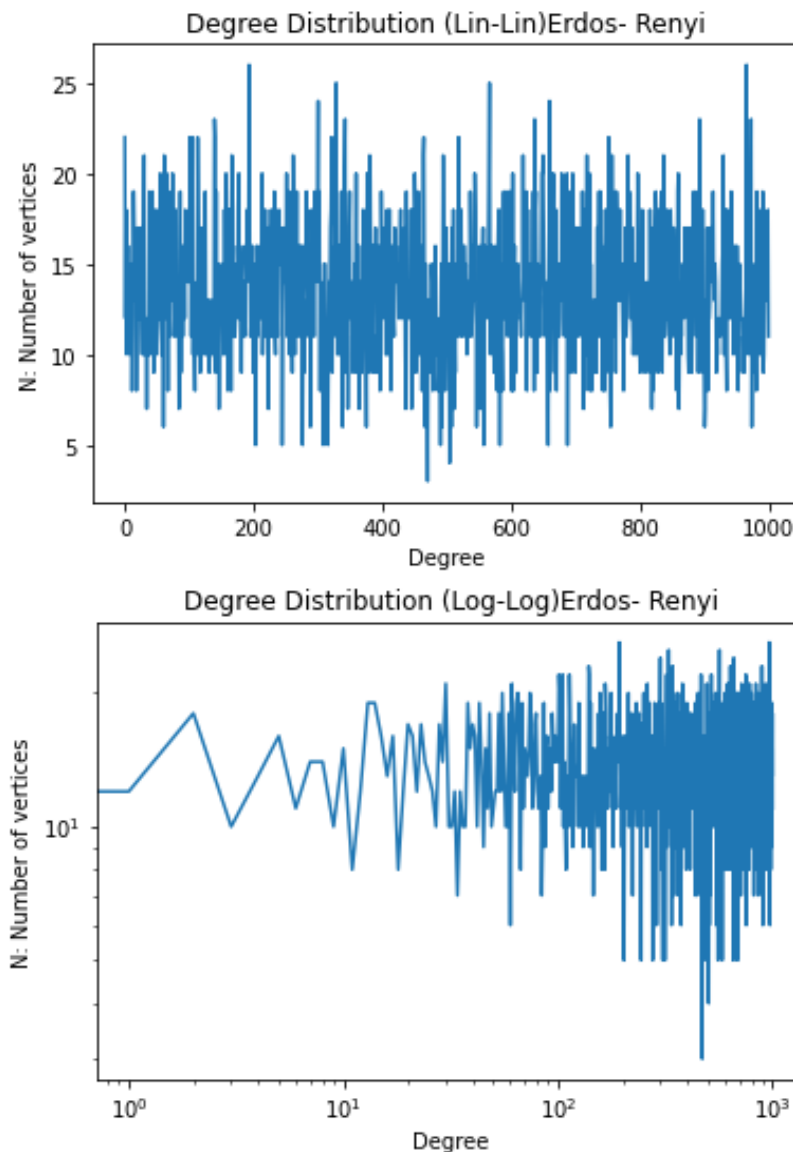
```
def degree_distribution(G):  
    map = new HashTable  
    for each vertex v in G:  
        map[v.degree]++  
    return map
```

C++ code:

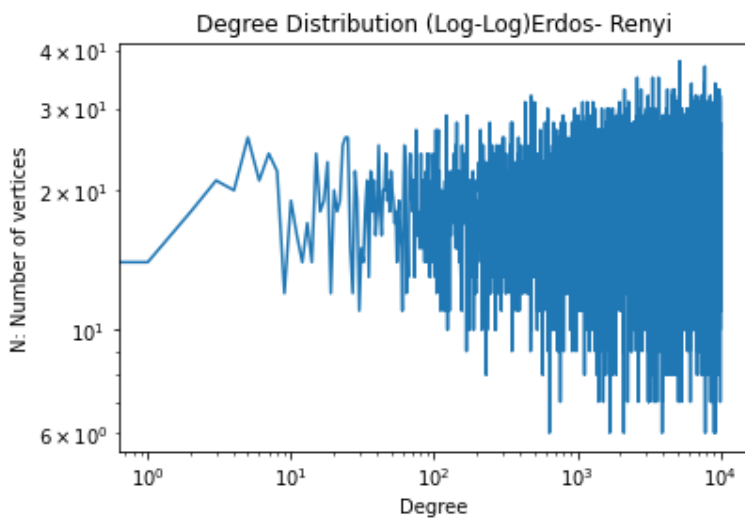
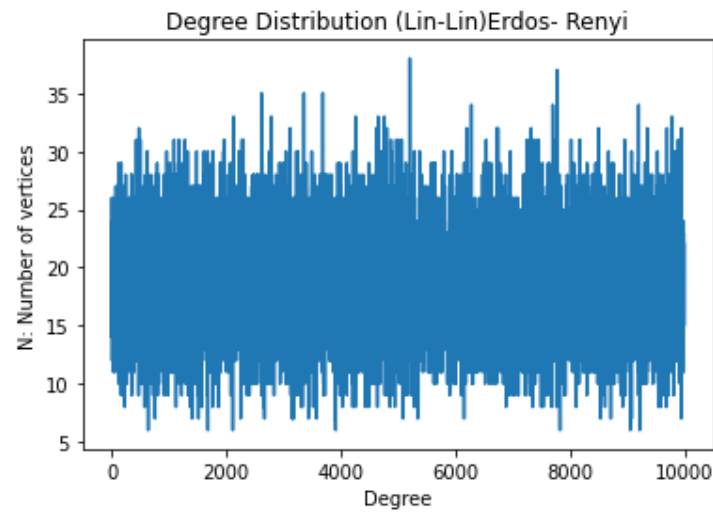
```
1 std::map < int, int > get_degree_distribution(Graph graph) {  
2     map < int, int > deg_dist;  
3     int size = graph.get_num_nodes();  
4     int degree = 0;  
5  
6     for (int i = 1; i <= size; ++i) {  
7         Node n(i);  
8         degree = graph.get_neighbors(n).size();  
9         deg_dist[degree] = deg_dist[degree] + 1;  
10    }  
11    return deg_dist;  
12 }
```

Plotting distribution for different sizes (N : number of vertices):

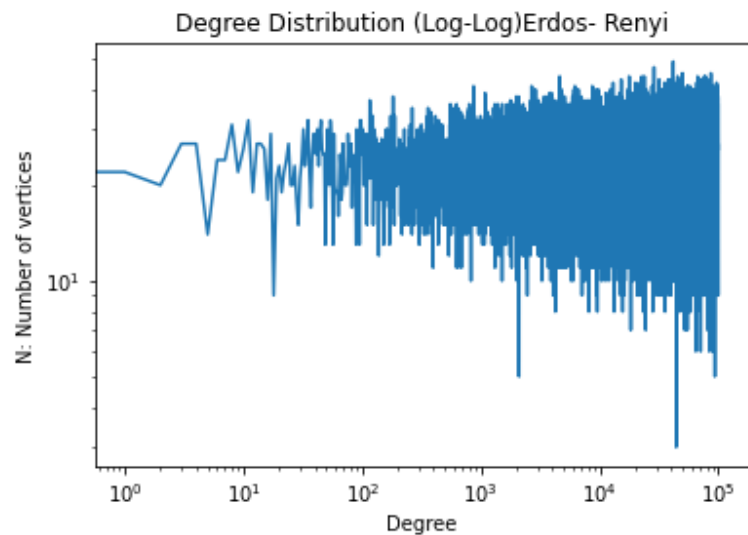
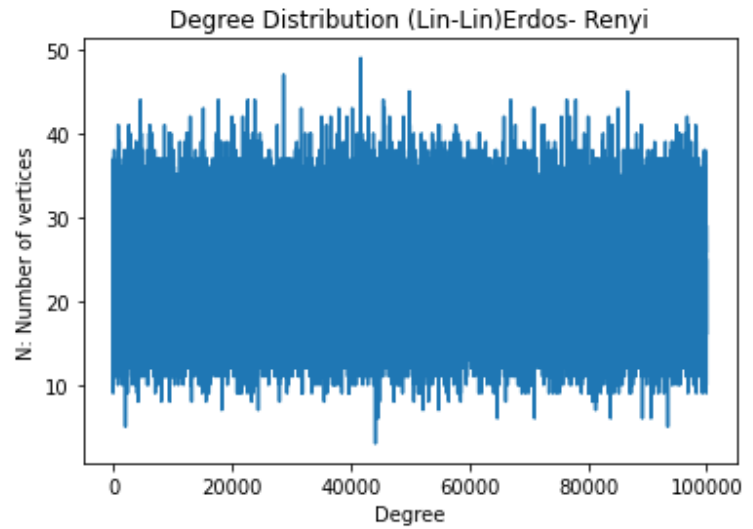
- $N = 1000$



- $N = 10,000$



- $N = 100,000$



Degree Distribution – Analysis:

As we can see the above graphs do not have a heavy/long tail, hence we can conclude that these distributions do not exhibit the power law.

Barabási – Albert Model:

Description:

- $G(n,d)$: - Graph G with n nodes, with minimum degree $d \geq 1$ with each new node being connected with a probability $= \frac{d_i}{\sum_i d_i}$ (here $\rightarrow d = 5$ (const))

Implementation:

Pseudocode:

ALG. 5: preferential attachment

Input: number of vertices n

minimum degree $d \geq 1$

Output: scale-free multigraph

$G = (\{0, \dots, n-1\}, E)$

M : array of length $2nd$

for $v=0, \dots, n-1$ **do**

for $i=0, \dots, d-1$ **do**

$M[2(vd+i)] \leftarrow v$

 draw $r \in \{0, \dots, 2(vd+i)\}$ uniformly at random

$M[2(vd+i)+1] \leftarrow M[r]$

$E \leftarrow \emptyset$

for $i=0, \dots, nd-1$ **do**

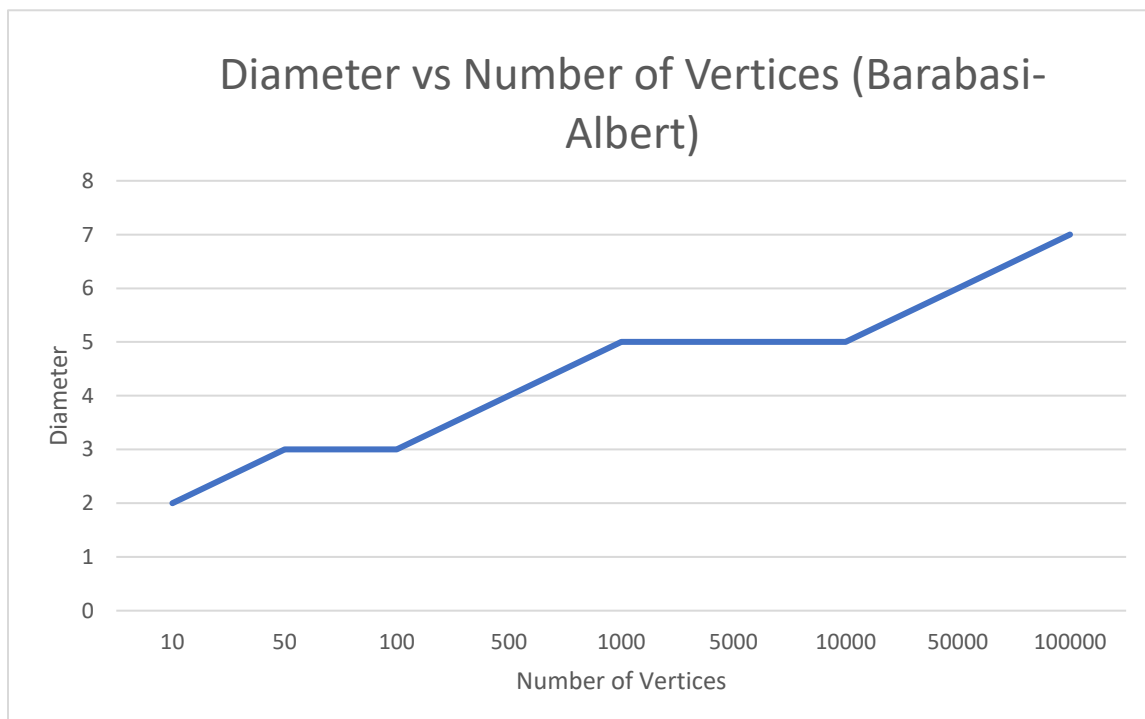
$E \leftarrow E \cup \{M[2i], M[2i+1]\}$

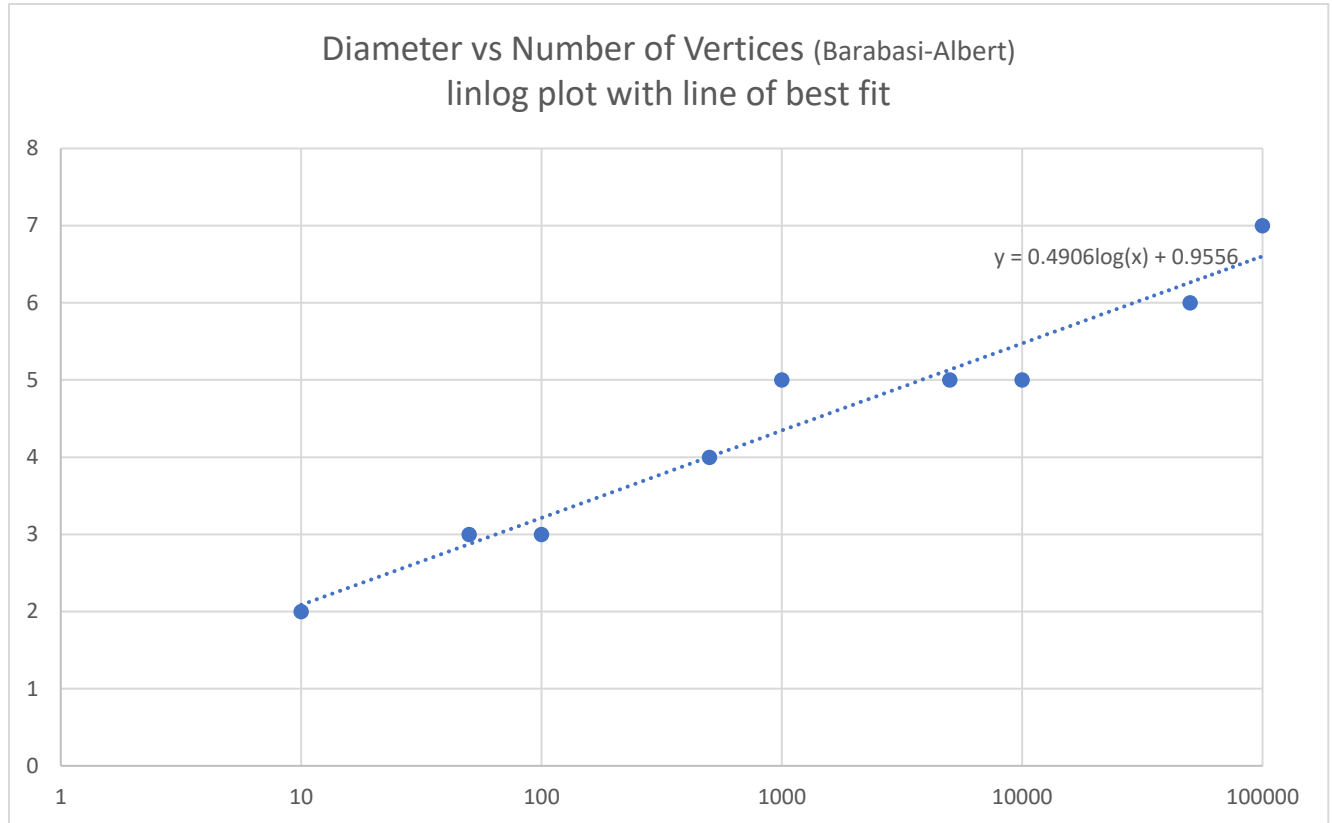
C++ Code:

```
Graph create_BA_graph(int nodes, int degree) {  
  
    Graph barabasi_albert = make_graph(nodes, vector<int> {}, vector<int> {});  
    int size = 2 * nodes * degree;  
    int L[size];  
    set<tuple<int, int>> edges;  
  
    for (int i = 0; i < nodes; i++) {  
        for (int j = 0; j < degree; j++) {  
            mt19937 seed = mt19937(chrono::system_clock::now().time_since_epoch().count());  
            L[2 * (i * degree + j)] = i;  
            int p = uniform_int_distribution<int>(0, 2 * (i * degree + j))(seed);  
            L[2 * (i * degree + j) + 1] = L[p];  
        }  
    }  
  
    for (int x = 0; x < nodes * degree; x++) {  
        int nid1 = L[2 * x] + 1;  
        int nid2 = L[2 * x + 1] + 1;  
        tuple<int, int> edge = make_tuple(nid1, nid2);  
        tuple<int, int> edge2 = make_tuple(nid2, nid1);  
        if (nid1 != nid2 && edges.find(edge) == edges.end()) {  
            barabasi_albert.add_edge(nid1, nid2);  
            edges.insert(edge);  
            edges.insert(edge2);  
        }  
    }  
  
    return barabasi_albert;  
}
```

Diameter:

(Implementation and Code on Page 2)

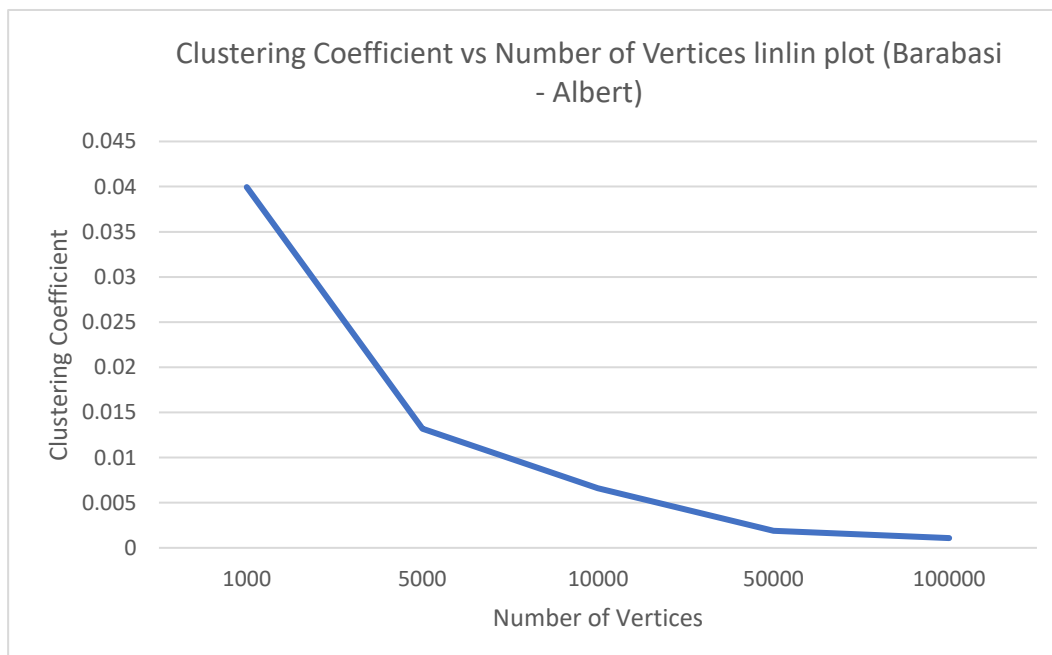


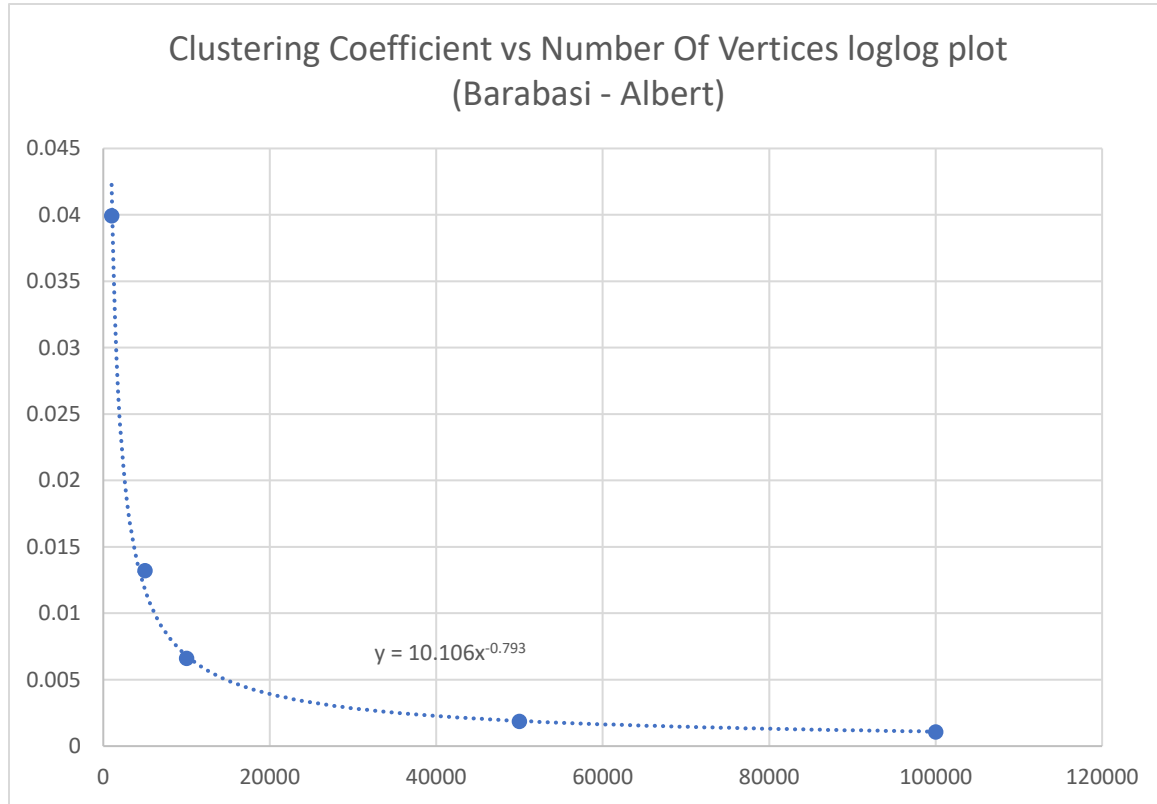


From the above plots we can see that the diameter increases linearly as the number of vertices increased from 10 to 10^5 . And from the lin-log plot we can say that the diameter: $f(d) = O(n)$ which is proportional to $\log n$.

Clustering Coefficient:

(Implementation and Code on Page 4)





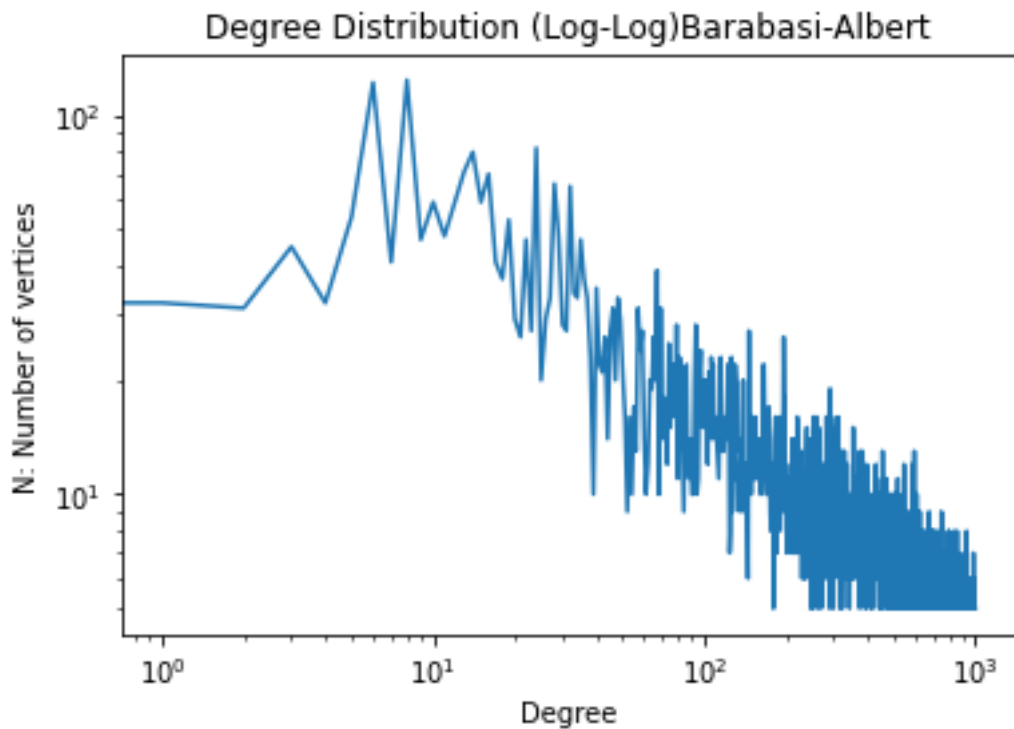
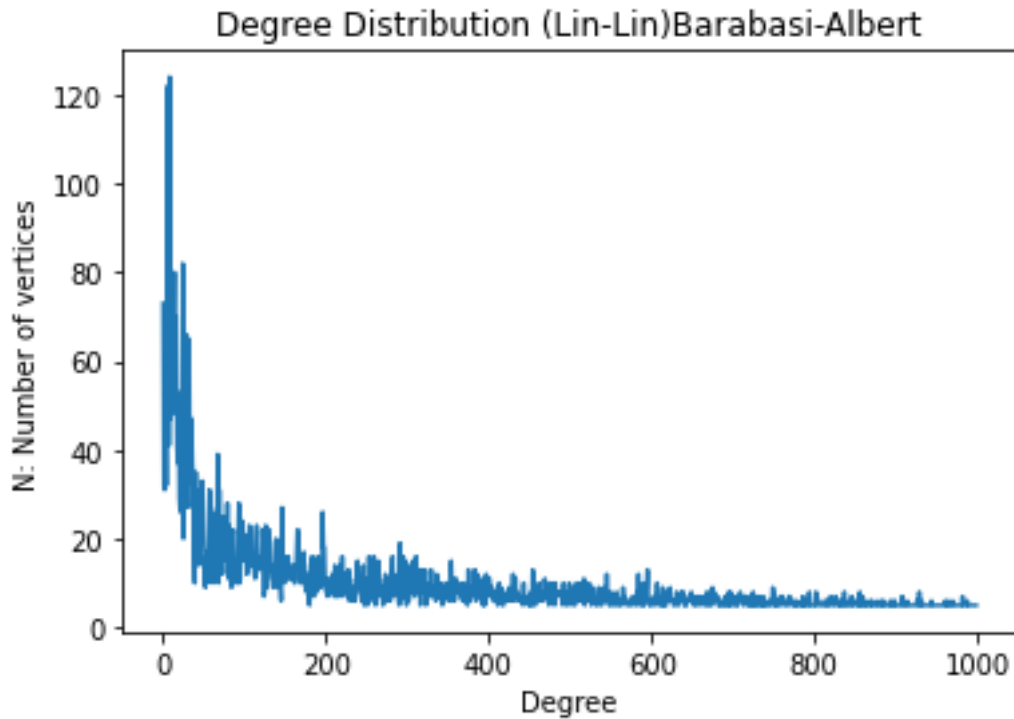
From the above plots we can see with the equation $y = 10.1062x^{-0.793}$ that the clustering coefficient decreases proportionally as the number of vertices increased from 10 to 10^5 so we can say that the clustering coefficient: $f(C) = O(1/n)$ which is proportional to $\log 1/n$.

Degree Distribution:

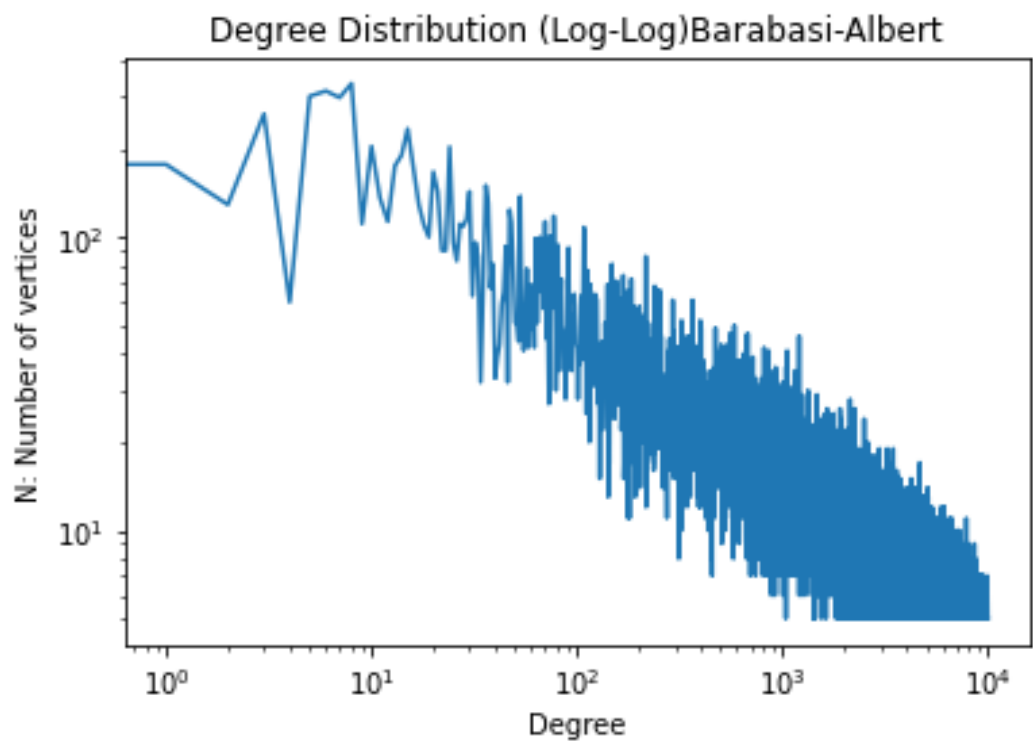
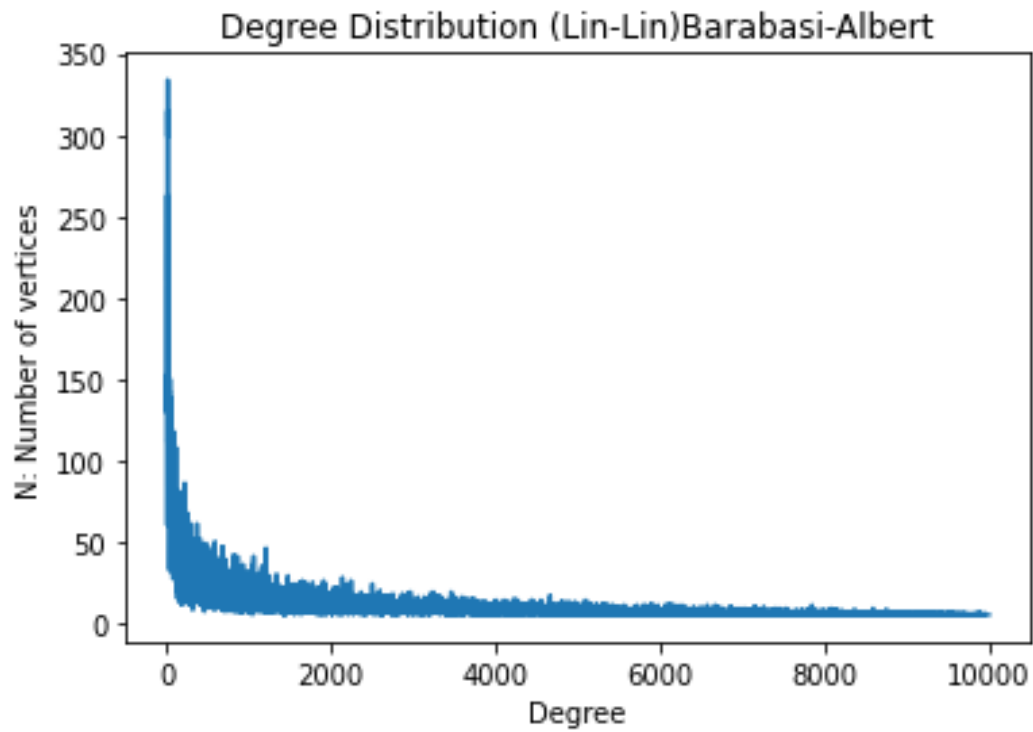
(Implementation and Code on Page 6)

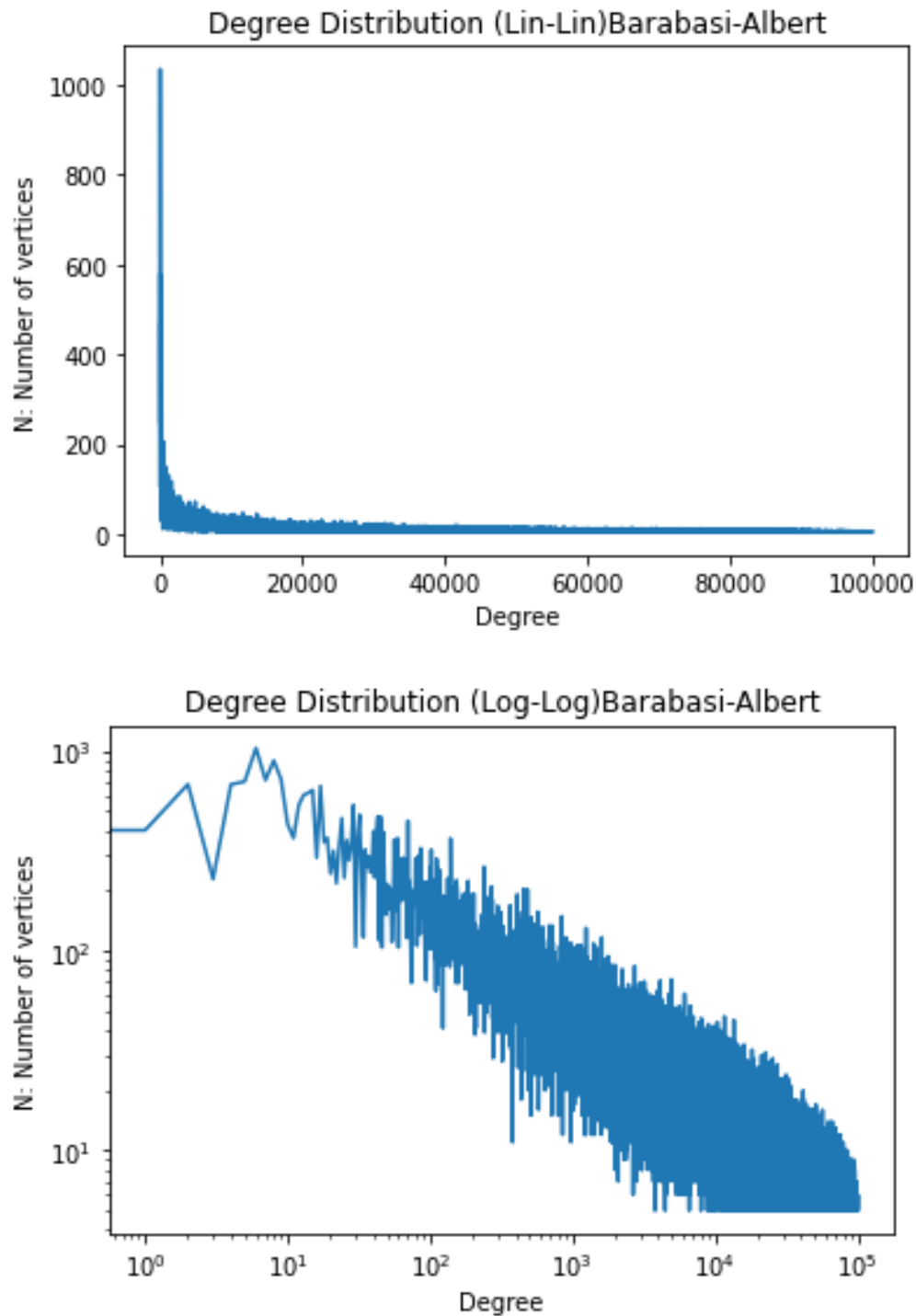
Plotting distribution for different sizes (N : number of vertices):

- $N = 1000$



- $N = 10,000$

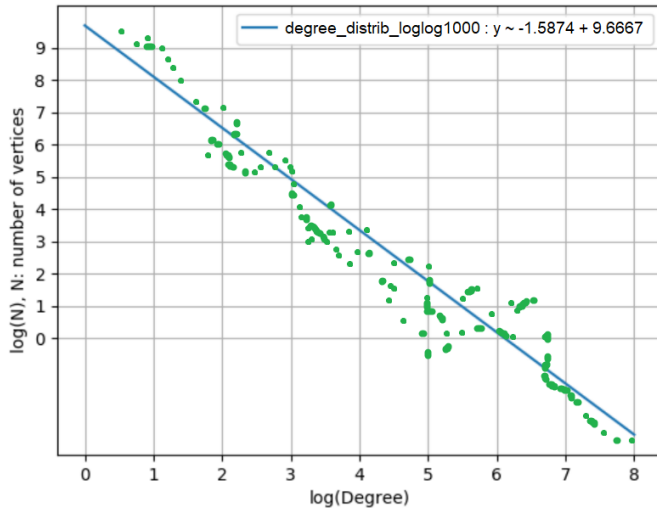




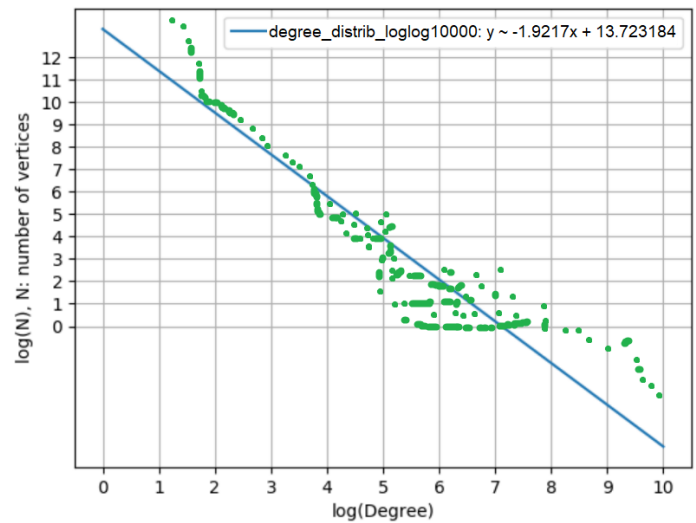
The above graphs show that the degree distribution for Barabasi-Albert models do exhibit a power law unlike the degree distributions for Erdos-Renyi models so we can perform a regression analysis on this model and find the line of best fit.

Regression Analysis for Barabasi Albert models only:

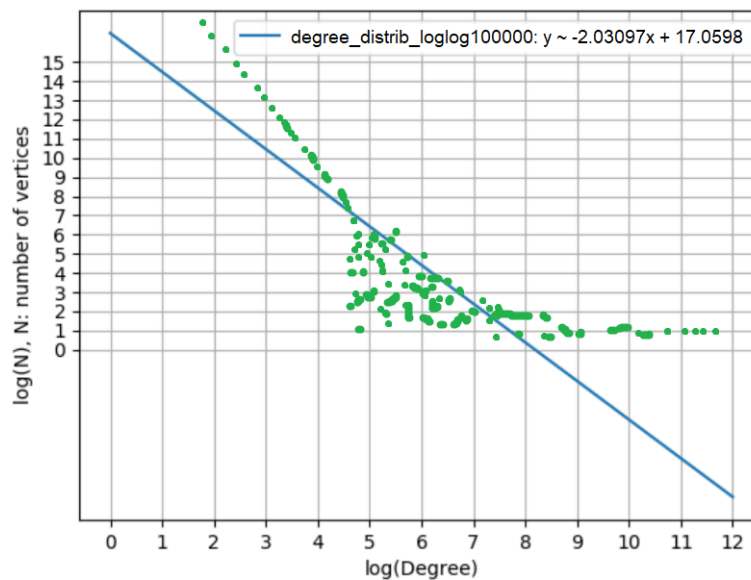
- $N = 1000$



- $N = 10,000$



- $N = 100,000$



Results:

The above graphs (for the number of vertices = 1000, 10000, 100000) clearly exhibit a power law since the log-log plots exhibit a heavy/long tail. And after doing a regression analysis we find that the exponents of the power law (slopes) of the best fitting lines are: $-1.59, -1.92, -2.03$ respectively.

Conclusion:

After experimenting the three algorithms on both models, we can conclude the following:

- Diameter $\propto \log(n)$
- Clustering Coefficient $\propto \log(1/n)$
- Degree Distribution $\propto n$

Therefore, we can say that all three properties of the graph are in some way proportional to the number of nodes in the graph when using either graph model.

Erdős – Rényi Model:

When experimenting the degree distribution on Erdős – Rényi models, we see that none of the graphs exhibit a heavy/long tail, so they do not have a power law as they follow what looks like Poisson distribution. This is due to the probability for each node being added to the edge. This probability does not model to any real-world network. Thus, if we were to do a study on real-world networks the Erdős – Rényi model would not be preferred for the experiments. But because the model is so simple and efficient, it is a good model to test the functionality of graph algorithms when implementing them.

Barabási – Albert Model:

The Barabási – Albert models, on the other hand, do exhibit a long/heavy tail in the degree distribution graphs so we can conclude that they do have a power law. This is due to the preferential attachment algorithm, vertices with higher degrees are likely to get more edges connected. This models closer to real-world networks such as social or communication networks where the more connections/friends you have (degree) the likelier you are to make more which will generally lead to more heavy-tailed distributions. Thus, if we were to do a study on real-world networks the Barabási – Albert model would be preferred for the experiments. Although both model algorithms run in $O(n + m)$ time, this model is more complex than the Erdős – Rényi model because we're calculating the probability by summing over the degrees of the remaining nodes for every vertex in the graph.