

Grace Chong  
DS210

## Project Write Up

### Data:

From the Stanford Network Analysis Project (SNAP)

<https://snap.stanford.edu/data/ego-Twitter.html>

I chose a social network dataset that is a txt file of two columns of Twitter users. In a folder, you can also find datasets that include node features (profiles), circles, and ego networks from Twitter. It is said that this data set can be valuable for studying social network structures and dynamics, particularly focusing on how users organize their social circles on platforms like Twitter.

### Use:

In my graph.rs I constructed a graph data structure using Struct, which uses an adjacency list where a hashmap tracks the nodes to a hashmap of edges. And an in degree list where my hashmap tracks the in degree number of incoming edges for each node. The methods I used were new() to create new 'Graph' with the empty adjacency list and indegree list, add\_edge9src, dest) to update the adjacent list and in degree list, the 'number\_of\_edges()' to calculate the total number of edges in the graph, number\_of\_nodes() to calculate the total number of unique nodes in the graph, and a print summary to see what the first 5 nodes would be. The Function I used was the parse\_dataset(file\_path: &st) -> io::Result<Graph> in order to parse a dataset from file path to create the graph.

For my algorithms, I create 2 separate mods (algorithm.rs and algorithm2.rs). These modules focused on six degrees of separation and average distance, but it also covers the functionalities of Breadth-First Search (BFS) and six degrees calculation.

In algorithm.rs, I implemented BFS to calculate the distance from my given starting node to all the other nodes in the graph. I also implemented a function to calculate degrees of separation, which determines whether two nodes are within a specific number of degrees of separation. This was to see whether users are six or fewer social connections away from each other.

In algorithm2.rs, the function computes the average distance between pairs of nodes in the network. This calculation is significant as it offers insights into how close the overall network is.

For my main.rs, which is the initial creator of the graph and to the implementations of my algorithms, I included key operations that calculate the average distance between pairs of nodes in the graph and checked the degrees of separation between specific

nodes. This algorithm provides an understanding of the overall interconnectedness and reachability within the network.

## Tests

I included a total of 4 tests. 2 for algorithm.rs and 1 for algorithm2.rs. The tests I had for algorithm.rs was one using the twitter dataset and one from a constructed dataset. The test I had for algorithm2.rs was a constructed dataset. Not only were tests required for this project, but it helped validate the correctness of both the graph creation and the algorithms implemented.

## Results:

In order to get my shown results, I implemented cargo build, cargo run --release, and cargo test.

I used cargo build in order to compile the project and to manage my dependencies.

```
▼ TERMINAL
• (base) gracechong@crc-dot1x-nat-10-239-28-134 retry % cargo build
  Compiling retry v0.1.0 (/Users/gracechong/Downloads/DSproject/finalproject)
  Finished dev [unoptimized + debuginfo] target(s) in 1.57s
• (base) gracechong@crc-dot1x-nat-10-239-28-134 retry %
```

I used cargo --release in order to add optimization, and because it is ideal when testing the performance of my code.

```
• (base) gracechong@crc-dot1x-nat-10-239-28-134 retry % cargo run --release
  Finished release [optimized] target(s) in 0.25s
  Running target/release/retry
Number of nodes: 81306
Number of edges: 1768149
Node 14137737 has connections to: {37729461, 27902315, 31353077, 78993414, 22461427, 18715024, 64776896, 63485337, 103943240, 14631433, 38318973, 93969417, 62613968, 44136194, 21056214, 70587360, 98795
912, 23487429, 35066483, 24852559, 57082636, 15901190, 23238377, 2100521, 14224219, 17064600, 69945997, 38138597, 24422716, 155688134, 5576542, 16523953, 167992638, 10810102, 52551600, 46426258, 245854
98, 14103291, 111136741, 26281970, 14303746, 154175423, 77059348, 16303106, 534725405, 44134863, 115485051, 386697756, 12195792, 22556717, 1405071, 115216851, 19287097, 5603992, 111492860, 16799897, 15
485441, 47667972, 265650530, 73731253, 36457353, 15279452, 10671602, 5162861, 13461662}
Node 310533612 has connections to: {}
Node 373403727 has connections to: {}
Node 294325981 has connections to: {114452371}
Node 225842922 has connections to: {6037672}
Average Distance: 4.87
Nodes 22261763 and 88323281 are within 6 degrees of separation.
Average six degrees separation (test = true): 0
Average six degrees separation (test = false): 4.951219512195122
• (base) gracechong@crc-dot1x-nat-10-239-28-134 retry %
```

It returned the accurate number of nodes and edges.

```
Number of nodes: 81306
Number of edges: 1768149
```

And the image below is a return of performing bfs to 5 different RANDOM nodes, so it is different every time I run it again. This was to find the shortest path in unweighted graphs, to help visualize how the nodes are connected, and its connections.

```
Node 14137737 has connections to: {37729461, 27902315, 31353077, 78993414, 22461427, 18715024, 64776896, 63485337, 103943240, 14631433, 38318973, 93969417, 62613968, 44136194, 21056214, 70587360, 98795
912, 23487429, 35066483, 24852559, 57082636, 15901190, 23238377, 2100521, 14224219, 17064600, 69945997, 38138597, 24422716, 155688134, 5576542, 16523953, 167992638, 10810102, 52551600, 46426258, 245854
98, 14103291, 111136741, 26281970, 14303746, 154175423, 77059348, 16303106, 534725405, 44134863, 115485051, 386697756, 12195792, 22556717, 1405071, 115216851, 19287097, 5603992, 111492860, 16799897, 15
485441, 47667972, 265650530, 73731253, 36457353, 15279452, 10671602, 5162861, 13461662}
Node 310533612 has connections to: {}
Node 373403727 has connections to: {}
Node 294325981 has connections to: {114452371}
Node 225842922 has connections to: {6037672}
```

The average distance also changes every time I run it again, because it is random sampling without replacement. Though it may not be as accurate or efficient, it was a strategic choice to allow for a practical analysis of this large graph. It was to prevent expensive calculations and long waiting time.

```
Average Distance: 4.87
```

```
Average six degrees separation (test = true): 0  
Average six degrees separation (test = false): 4.640449438202247
```

The shown result of the image below indicates the average approximation of 4.64 degrees of separation which suggests a highly interconnected network. Which shows that in this graph, the twitter users are well connected.