# Impact of Social Media
# Information Spread Analysis on X

Progress Report

Colin Song (261049870)

# I. Abstract

This project investigates the flow of information on X (formerly Twitter) by analyzing the platform's network structure and user interactions. The primary objective of this project is to identify subgraphs and their properties to

gain deeper insight into how information flows on X. Characterizing subgraphs and their properties reveals patterns of influence, the roles of gatekeepers, and the mechanisms by which information reaches distant parts of the network. These findings allow for a better understanding of the structural and behavioral factors of subgraphs that ultimately reveal how information flows on X.

# II. Project Description

## 2.1 Purpose

The purpose of this project is to gain a better understanding of how information flows on social media platforms primarily on X. Information flow is crucial for understanding how ideas, news, and opinions spread across networks, influence public perspectives, and affect people's decision-making. By studying these processes, we can reveal patterns of influence, understand who is influenced, and calculate how quickly information reaches different parts of the network. This knowledge is key to addressing challenges such as preventing an increase in news of misinformation, optimizing the spread of accurate information, and improving communication strategies.

Analyzing communities within social media networks plays a critical role in achieving said goals. Communities, which represent smaller subsets of the network, allow us to focus on specific areas of the network. This localized analysis is valuable for understanding the behavior and dynamics of clustered communities and their impact on the network. Community analysis helps identify significant nodes—such as influencers and gatekeepers—that control the flow of information. Studying communities also makes network analysis more manageable. By focusing on smaller areas of a network, we can gain understanding that apply to the entire network without being overwhelmed by its complexity.
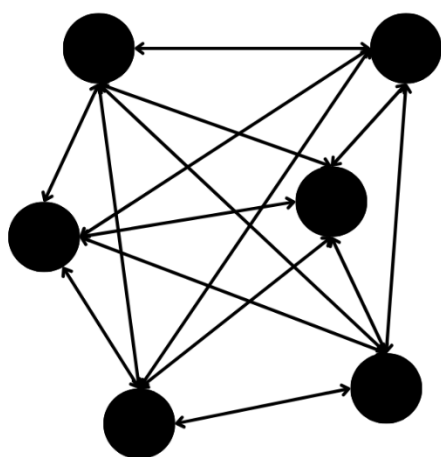
## 2.2 Description

In this project, the primary focus was to analyze the flow of information on X by exploring how different graph structures influence the network. Three distinct graph types were studied: cloistered graph, tree graph, and line branch graph. Each structure represents a different analysis and understanding of how information flows through a network. I will provide a summary of each type of subgraph
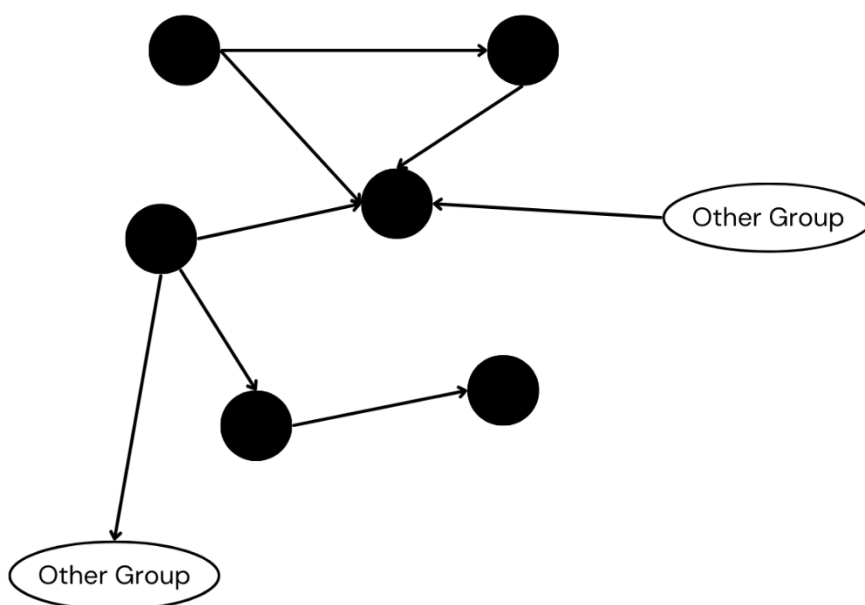
## Cloistered Graph:

The cloistered graph simulates tightly knit communities, also referred to as groups, within a network, where each user has many connections with others in the same community. An ideal cloistered graph would be a disconnected complete graph. However, since such groups are difficult to see in a real-world setting, a more practical definition for a cloistered graph would be a group that has few to none edges connected to other groups within a network. This structure would be analogous to closed groups on X, such as private discussions or niche communities where information tends to flow among active users.
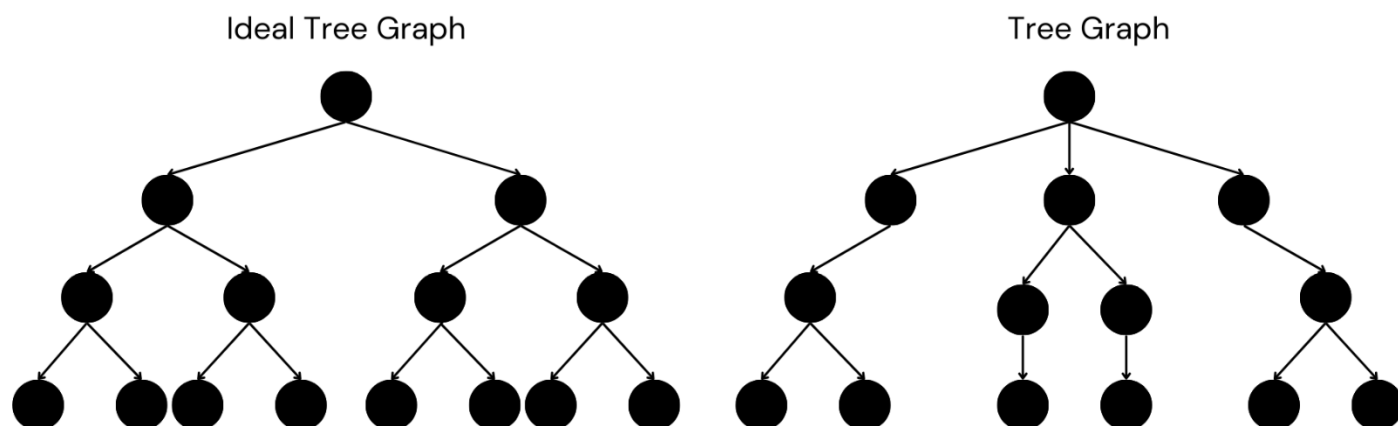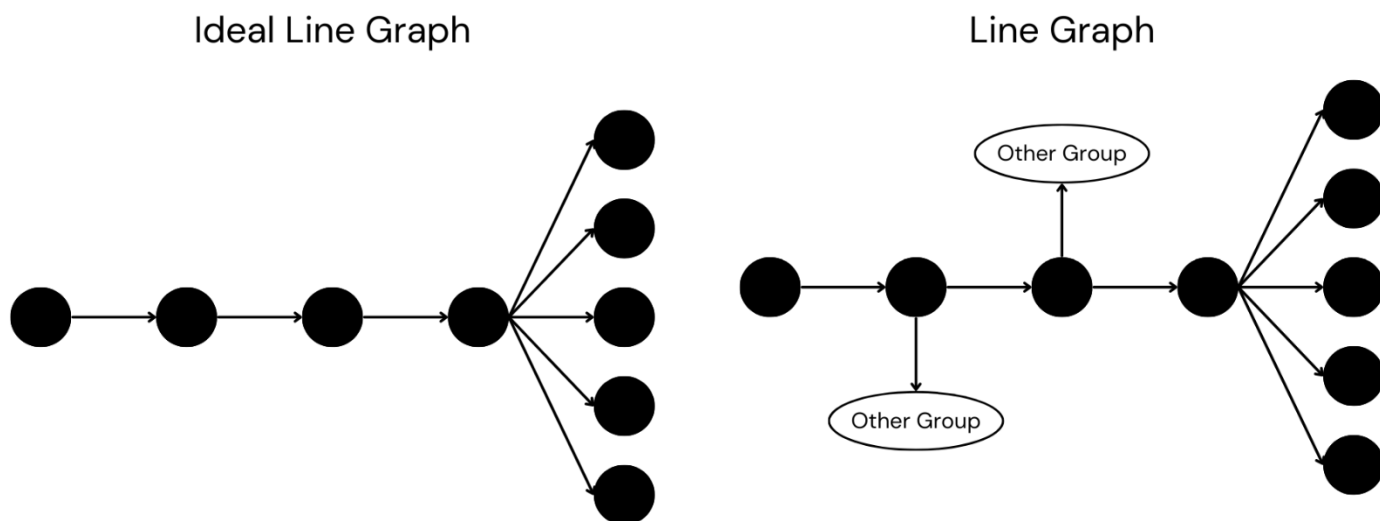


## Tree Graph:

The tree graph is a hierarchical data structure where an ideal tree graph is a noncyclic graph that develops from a single node (root node) and each node splits into 2 child nodes. This model reflects scenarios on X where a central figure (an influencer or organization) posts information to their followers, which then spreads through subsequent layers of users.

Ideal Tree Graph                                        Tree Graph



## Line Branch Graph:

The line branch graph is a line of nodes where the final node in the line is connected to several nodes. An ideal line branch graph would be a group where the line of nodes only has connections to the node before and after itself, and the final node in the line has several connections to other nodes. This sequence of nodes would be a group where information passes sequentially from one user to the next, representing the retweet chain seen in niche topics or during initial phases of information propagation.

Analyzing these three graph types revealed how certain community configurations affect the speed, reach, and flow of information in networks on X. These findings contribute to understanding the properties of real-world social networks and provide a baseline for further exploration of how information flows.

# 2.3 What was built

To implement the analysis of information flow on X, the project focused on graphing the network and analyzing the network through the lens of three distinct graph structures: a cloistered graph, tree graph, and line branch graph. These graph types were chosen to represent different community configurations commonly observed in networks.

The implementation was carried out using R, leveraging its robust libraries for graph representation and analysis. The program written in R is structured to organize the data, graph the network, and analyze the communities in networks by categorizing them into subgraphs.

The program starts with loading the necessary packages for the program.

```r
1 ▾ # startup
2 ▾ ```{r}
3   library(tidyverse)
4   library(igraph)
5   library(visNetwork)
6   library(MASS)
7
```

The program then reads from a CSV file that contains the data representing the network. The following lines in the program then organize data by splitting the data into two lists: nodes that point to other nodes and nodes that are pointed to.

```r
7
8   # read file
9   nd <- as.data.frame(read_csv("./[INSERT_FILE_NAME].csv"))
10  # users sending
11  from_id <- nd[,1]
12  # users receiving
13  to_id <- nd[,2]
14 ▴ ```
```

After organizing the data, the program then adds one to every node id in the nodes pointing to other nodes and nodes that are pointed to lists. This is due to R not being able to graph nodes that start with 0 or search through lists using an index of 0.

```r
15
16 ▾ # add 1 to from_id list and to_id list
17 ▾ ```{r}
18    # create new from_id list where we add 1 to each node id (this is because you cant have a node 0, must start at node 1)
19    # initialize from_id_plus list
20    from_id_plus <- c()
21    # iterate through nodes in from_id list
22 ▾ for (i in from_id) {
23       # add node id plus 1 to from_id_plus list
24       from_id_plus <- c(from_id_plus, i+1)
25 ▴ }
26    # create new to_id list where we add 1 to each node id (this is because you cant have a node 0, must start at node 1)
27    # initialize to_id_plus list
28    to_id_plus <- c()
29    # iterate through nodes in to_id list
30 ▾ for (i in to_id) {
31       # add node id plus 1 to to_id_plus list
32       to_id_plus <- c(to_id_plus, i+1)
33 ▴ }
34 ▴ ```
```

The program then creates a directed list that connects the nodes using the two lists created above.

```r
35
36 ▾ # create directed list connecting from and to nodes
37 ▾ ```{r}
38    # list with nodes
39    from_to <- c()
40    # index
41    j <- 1
42    # for loop to generate nodes with directed edges
43 ▾ for (i in 1:dim(nd)) {
44       # get users sending (need to add 1 as can't have nodes starting with 0 in R)
45       from_to[j] <- nd[i,1]+1
46       j <- j + 1
47       # get users receiving
48       from_to[j] <- nd[i,2]+1
49       j <- j + 1
50 ▴ }
51 ▴ ```
```

A graph object is created (from the *igraph* package) using the directed list. The program then applies the Louvain algorithm to identify communities in the network. This algorithm is a built-in function (from the *igraph* package) that detects communities in a network by optimizing a parameter called modularity, which quantifies how well a

network can be divided into distinct groups. It separates the network into groups of nodes that are more densely connected to each other than to other nodes. Since my objective was to sort communities into subgraphs, the Louvain algorithm was an excellent choice. Its efficiency on large datasets made it perfect for analyzing my data, which included networks with over 1 million connections. The following lines apply the Louvain algorithm to the graph object and manipulate the data to prepare a visual representation of the network.

```r
55 ▾ ### GRAPHING
56
57 ▾ # create graph objects
58 ▾ ```{r}
59   # create graph object
60   g <- graph(from_to)
61   # apply Louvain algorithm on graph object
62   new_g <- cluster_louvain(as.undirected(g, mode = "collapse"), weights = NULL, resolution = 1)
63 ▴ ```
64
65 ▾ # create variables used to create node/edge data frame
66 ▾ ```{r}
67   # initialize list to get number of elements in each group
68   group_lengths <- c()
69   # get all unique node ids
70   unique_node_ids <- c()
71   # iterate through nodes in new_g
72 ▾ for(i in 1:length(new_g)){
73     # add group of nodes to unique_node_ids list
74     unique_node_ids <- c(unique_node_ids, c(new_g[[i]]))
75     # add group length to group_lengths list
76     group_lengths <- c(group_lengths, length(new_g[[i]]))
77 ▴ }
78 ▴ ```
```

The program creates two data frames that contain information about the nodes and edges in the network. These two data frames are used to assist with graphing the network.

```r
80 - # create node data frame variables (for visNetwork graph)
81 -  ```{r}
82    # create node labels
83    # initialize node labels list
84    node_labels <- c()
85    # iterate through unique node ids
86 - for (i in unique_node_ids) {
87      # add node to node_labels
88      node_labels <- c(node_labels, as.character(i))
89 - }
90    # create node groups
91    # initialize node groups list
92    node_groups <- c()
93    # temp counter
94    j <- 1
95    # iterate through group_lengths list
96 - for (i in group_lengths) {
97      # repeat group id by value of i and add to node_groups list
98      node_groups <- c(node_groups, rep(j, i))
99      # increase j
100     j <- j + 1
101 - }
102   # create node titles
103   # initialize node titles list
104   node_titles <- c()
105   # iterate through unique node ids
106 - for (i in unique_node_ids) {
107     # create node title and add to list
108     node_titles <- c(node_titles, paste("Node", as.character(i), "info"))
109 - }
110   # create nodes data frame for Louvain algorithm applied graph
111   new_g_nodes <- data.frame(
112     id = unique_node_ids,
113     label = node_labels,
114     group = node_groups,
115     title = node_titles
116   )
117 - ```
```

```r
119 - # create edge data frame variables (for visNetwork graph)
120 -  ```{r}
121   # create edge labels
122   # initialize edge labels list
123   edge_labels <- c()
124   # iterate with length of from_id_plus list
125 - for (i in 1:length(from_id_plus)) {
126     # create edge label and add to list
127     edge_labels <- c(edge_labels, paste("Edge", from_id_plus[i], "-", to_id_plus[i]))
128 - }
129   # create edge titles
130   # initialize edge titles list
131   edge_titles <- c()
132   # iterate with length of from_id_plus list
133 - for (i in 1:length(from_id_plus)) {
134     # create edge title and add to list
135     edge_titles <- c(edge_titles, paste("Edge from", from_id_plus[i], "to", to_id_plus[i]))
136 - }
137   # create edges data frame for Louvain algorithm applied graph
138   new_g_edges <- data.frame(
139     from = from_id_plus,
140     to = to_id_plus,
141     label = edge_labels,
142     arrows = "to",
143     title = edge_titles
144   )
```
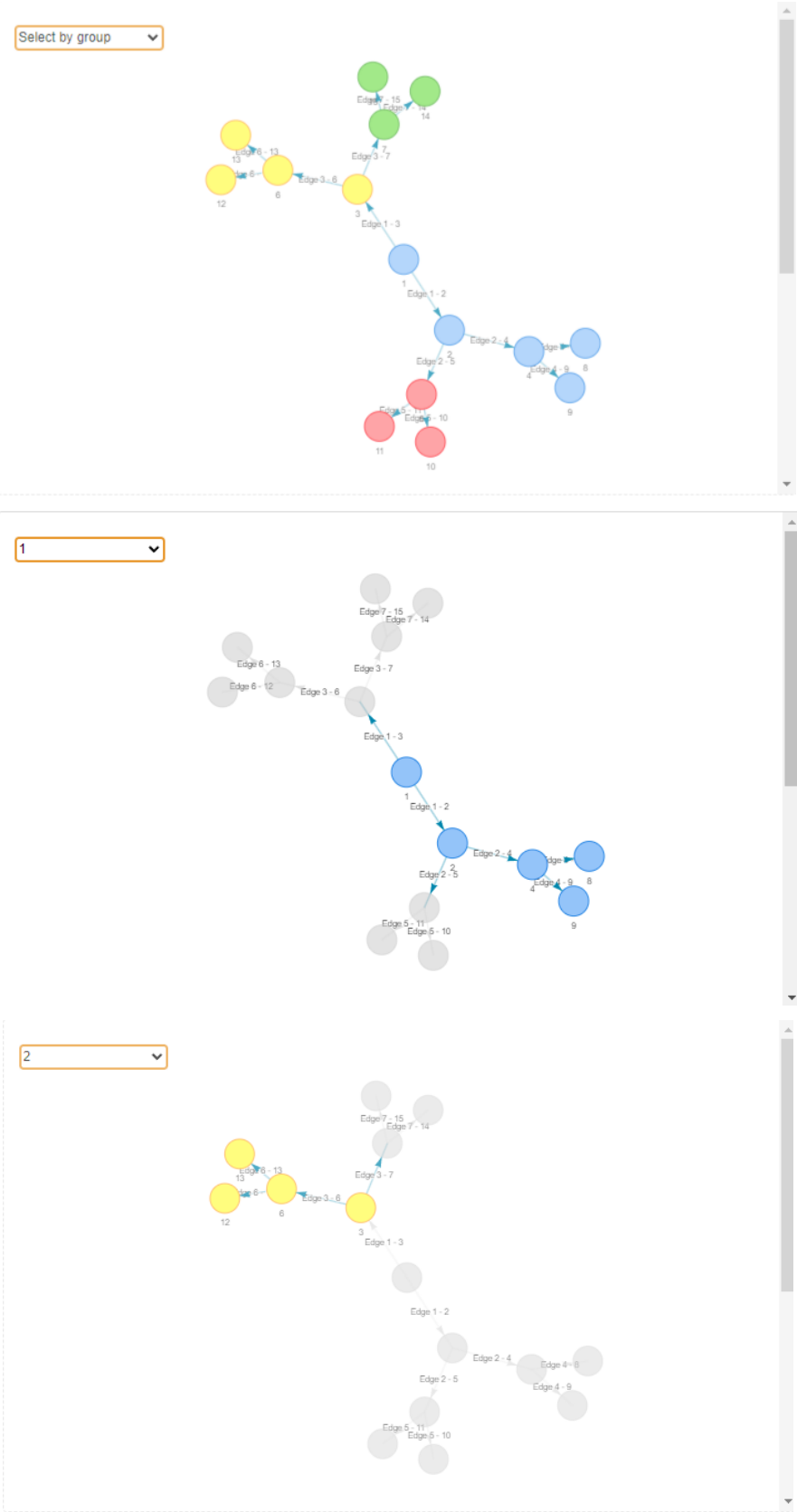
After preparing the data, the program graphs the network using *visNetwork* and *igraph*. While both graphs represent the same network, they offer slightly different perspectives and insights. *visNetwork* provides an

interactive representation of the network as well as an option to export the representation as an HTML file. *igraph* provides a static image of the network and focuses on network analysis.

```r
147 ▾ # graph using visNetwork
148 ▾ ```{r}
149   #set.seed(11) (seed for similar looking graph)
150   set.seed(24)
151   visNetwork(new_g_nodes, new_g_edges, height = "850px", width = "100%") %>%
152     visIgraphLayout() %>%
153     visNodes(
154       shape = "dot",
155       color = list(
156         background = "#0085AF",
157         border = "#013848",
158         highlight = "#FF8000"
159       ),
160       shadow = list(enabled = FALSE, size = 35)
161     ) %>%
162     visEdges(
163       shadow = FALSE,
164       color = list(color = "#0085AF", highlight = "#C62F4B")
165     ) %>%
166     visOptions(highlightNearest = list(enabled = T, degree = 1, hover = T),
167                selectedBy = "group")
168 ▴ ```
169
170 ▾ # graph using igraph
171 ▾ ```{r}
172   set.seed(2)
173   plot(g)
174   set.seed(2)
175   plot(new_g, g)
176 ▴ ```
```
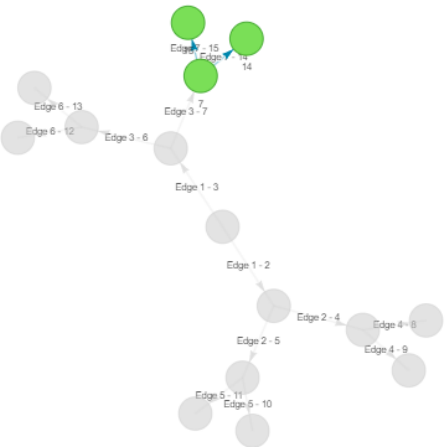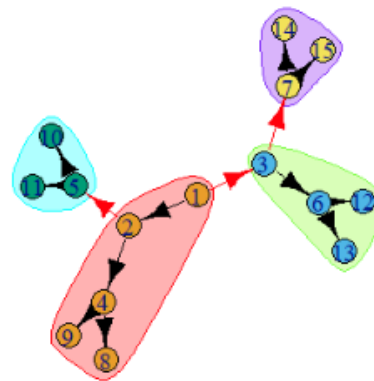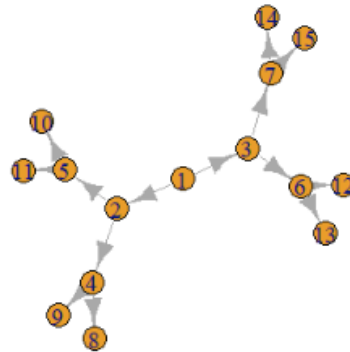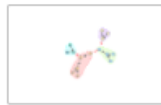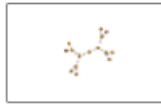
Sample graph using *visNetwork*:

After dividing the network into communities and graphing it, the program then analyzes the communities and applies the subgraph analysis.

To identify all communities that resemble cloistered graphs, the program compares the number of edges in the group to the maximum number of possible edges in the group.

I created a few helper functions to assist with the analysis. The first helper function finds every connection between nodes in the data. The second helper function finds every edge that connects to a different community and every node that belong to the same community. The third helper function finds the MPNE (maximum possible number of edges) of the group.

```r
180 ▾ ### ANALYSIS
181
182 ▾ # CLOISTERED GRAPH
183
184 ▾ # get all connections between nodes
185 ▾ ```{r}
186   # node connections list
187   node_connections <- c()
188   # iterate through all nodes in from_to list
189 ▾ for (i in seq(1, length(from_to), 2)){
190     # append each to node to respective from node
191     node_connections[[as.character(from_to[i])]] = c(node_connections[[as.character(from_to[i])]], from_to[i+1])
192 ▴ }
193 ▴ ```
194
195 ▾ # find edges that belong to same group or different group and find all node connections in same group
196 ▾ ```{r}
197   # nodes in same group list
198   group_node_connections <- list()
199   # edges that belong to a group
200   in_edge_groups <- list()
201   # edges that leave a group
202   out_edge_groups_f <- list()
203   # edges that enter a group
204   out_edge_groups_s <- list()
205
206   # iterate through all edges
207 ▾ for (i in 1:nrow(new_g_edges)){
208     # get from node
209     fnum <- new_g_edges$from[i]
210     # get to node
211     snum <- new_g_edges$to[i]
212     # find which group fnum belongs to
213     fnum_group <- new_g_nodes$group[which(new_g_nodes$id == fnum)]
214     # find which group snum belongs to
215     snum_group <- new_g_nodes$group[which(new_g_nodes$id == snum)]
216
217     # if fnum and snum belong to same group
218 ▾   if(fnum_group == snum_group){
219       # add edge to in_edge_group
220       in_edge_groups[[as.character(fnum_group)]] <- c(in_edge_groups[[as.character(fnum_group)]], sprintf('(%i, %i)', fnum, snum))
221       # add nodes to group_node_connections
222       group_node_connections[[as.character(fnum)]] <- c(group_node_connections[[as.character(fnum)]], snum)
223     # fnum and snum don't belong to same group
224 ▾   } else{
225       # add edge to out_edge_group_f
226       out_edge_groups_f[[as.character(fnum_group)]] <- c(out_edge_groups_f[[as.character(fnum_group)]], sprintf('(%i, %i)', fnum, snum))
227       # add edges to out_edge_group_s
228       out_edge_groups_s[[as.character(snum_group)]] <- c(out_edge_groups_s[[as.character(snum_group)]], sprintf('(%i, %i)', fnum, snum))
229 ▴   }
230 ▴ }
231 ▴ ```
```

```r
233 ▾  # find MPNE for each group
234 ▾  ```{r}
235    # max possible number of edges in a group
236    MPNE <- c()
237    # iterate through all groups
238 ▾  for (i in 1:length(new_g)) {
239      # max possible number of edges calculated as n * (n-1) / 2 where n = number of nodes
240      # append max possible number of edges to MPNE
241      MNPE <- c(MPNE, (length(new_g[[i]])*(length(new_g[[i]])-1))/2)
242 ▴  }
243 ▴  ```
```

Using these three helper functions, the program then performs the cloistered graph identification analysis on every group within the network.

```r
245 ▾ # Cloistered Graph Analysis
246 ▾ ```{r}
247   print("-------- Cloistered Graph Analysis --------")
248
249   # iterate over different groups
250 ▾ for (i in 1:length(MPNE)) {
251     # Group i analysis
252     print(sprintf("Group %i:", i))
253     # find total number of edges entering and leaving
254     edges_l_e <- length(out_edge_groups_f[[as.character(i)]]) + length(out_edge_groups_s[[as.character(i)]])
255     # how many edges are in each group
256     len_egroup <- length(in_edge_groups[[as.character(i)]])
257     # how close group is to that of a complete graph as a percentage
258     percent <-  len_egroup / (MPNE[[i]]) * 100
259     # how many edges are missing from that of a complete graph
260     missing_e <- MPNE[[i]] - len_egroup
261     # group is cloistered graph or not
262     cloistered <- FALSE
263
264     # if number of edges leaving/entering group is 2 percent of max possible number of edges
265     # then consider it a cloistered group
266 ▾   if (edges_l_e < MPNE[[i]] * 0.02) {
267       # print statement
268       print(sprintf("   Group %i is a cloistered graph with %i edges leaving/entering the group", i, edges_l_e))
269       # set cloistered to TRUE
270       cloistered <- TRUE
271     # if more than 2 percent then don't consider a cloistered graph
272 ▾   } else {
273       # print statement
274       print(sprintf("   Group %i is not a cloistered graph with %i edges leaving/entering the group", i, edges_l_e))
275 ▴   }
276
277     # if percentage = 100% and graph is a cloistered graph then group is an ideal cloistered graph
278 ▾   if (percent == 100 & cloistered & length(edges_l_e) == 0) {
279       # print statement
280       print(sprintf("   Group %i is an ideal cloistered graph", i))
281     # if percentage = 100% but not cloistered graph then only complete graph
282 ▾   } else if (percent == 100 & !cloistered){
283       # print statement
284       printf(sprintf("   Group %i is a complete graph", i))
285     # if percentage != 100% and not a cloistered graph
286 ▾   } else {
287       # print statement
288       print(sprintf("   Group %i is %f%% (%i/%i) of a complete graph with %i missing edges compared to that of a complete graph", i, percent, len_egroup, MPNE[[i]], missing_e))
289 ▴   }
290 ▴ }
291 ▴ ```
```

Next, the program determines which communities are tree graphs by comparing every community against the properties of an ideal tree graph.

Once again, I created several helper functions to assist with the analysis. The first helper function constructs a tree structure using the nodes within a group. The second helper function checks for cycles within a group. The third helper function examines the group in its tree structure and calculates specific statistics, which will be discussed in Section 3.1. The fourth helper function evaluates these statistics to determine how similar the tree structure is to an ideal tree graph.

```r
295 ▾ # TREE GRAPH
296
297 ▾ # function to create tree
298 ▾ ```{r}
299 ▾ create_tree <- function(tree = list(), level = 1, node) {
300     # if no nodes in current level of tree
301 ▾   if (is.null(tree[[as.character(level)]])) {
302       # add node to tree at level
303       tree[[as.character(level)]] <- node
304     # nodes exist in current level of tree
305 ▾   } else {
306       # add node to tree at level
307       tree[[as.character(level)]] <- c(tree[[as.character(level)]], node)
308 ▴   }
309
310     # get children of  node
311     children <- group_node_connections[[as.character(node)]]
312
313     # if node has children
314 ▾   if (!is.null(children)) {
315       # iterate through children
316 ▾     for (child in children) {
317         # recursive call|
318         tree <- create_tree(tree, level + 1, child)
319 ▴     }
320 ▴   }
321
322     # return tree
323     return(tree)
324 ▴ }
325 ▴ ```
```

```r
327 ▾ # check for cycle function
328 ▾ ```{r}
329 ▾ has_cycle <- function(node_list, group_node_connections) {
330      # helper function for DFS
331 ▾   dfs <- function(node) {
332
333        # if a node has not been visited
334 ▾     if (!visited[[node]]) {
335          # mark node as visited
336          visited[node] <- TRUE
337          # add node to recursion stack
338          rec_stack[node] <- TRUE
339
340          # if the node is connected to other nodes in group
341 ▾       if (!is.null(group_node_connections[[as.character(node)]])) {
342            # iterate through other nodes in group
343 ▾         for (near_node in group_node_connections[[as.character(node)]]) {
344              # if the node has not been visited and recursion call on near_node
345 ▾           if (!visited[[near_node]] && dfs(near_node)) {
346                # return TRUE
347                return (TRUE)
348              # otherwise if the node near is on the recursion stack
349 ▾           } else if (rec_stack[[near_node]]) {
350                # return TRUE
351                return (TRUE)
352 ▴           }
353
354 ▴         }
355 ▴       }
356 ▴     }
357        # remove the node from recursion stack
358        rec_stack[node] <- FALSE
359        # return FALSE
360        return (FALSE)
361 ▴   }
362
363      # list for visited nodes
364      visited <- list()
365      # list for recursion stack
366      rec_stack <- list()
367      # initialize visited and recursion stack
368 ▾   for (node in node_list){
369        visited[node] <- FALSE
370        rec_stack[[node]] <- FALSE
371 ▴   }
372
373      # check each node in the group
374 ▾   for (node in node_list) {
375        # if dfs(node is true then there is a cycle)
376 ▾     if (dfs(node)){
377          # return TRUE
378          return (TRUE)
379 ▴     }
380 ▴   }
381      # there is no cycle so return FALSE
382      return (FALSE)
383 ▴ }
384 ▴ ```
```

```r
386 ▾ # function to find statistics
387 ▾ ```{r}
388 ▾ tree_statistics <- function(tree, group_node_connections){
389       # set level error to 0
390       level_error <- 0
391       # set extra nodes to 0
392       extra_nodes <- 0
393       # set missing nodes to 0
394       missing_nodes <- 0
395       # set height balance to number of levels
396       height_balance <- length(names(tree))
397       # set nodes relative to ideal nodes to 0
398       nrin <- 0
399       # set total nodes to 1 (root node)
400       total_nodes <- 1
401       # set largest number of nodes in a level to 1 (root node)
402       lnol <- 1
403       # set number of nodes in a level to 1 (root node)
404       nol <- 1
405
406       # iterate through all levels in tree (skip level 0)
407 ▾     for (level in names(tree)[-1]){
408         # increase total nodes by number of nodes on level
409         total_nodes <- total_nodes + length(tree[[as.character(level)]])
410         # get level number of previous level
411         prev_level <- as.character(as.numeric(level)-1)
412
413         # if current number of nodes is greater than 2 ^ number of nodes from previous level
414 ▾       if (length(tree[[as.character(level)]]) > 2^length(tree[[prev_level]])){
415           # increase extra nodes by the difference
416           extra_nodes <- extra_nodes + length(tree[[as.character(level)]]) - 2^length(tree[[prev_level]])
417           # increase level error by 1
418           level_error <- level_error + 1
419         # if current number of nodes is less than 2 ^ number of nodes from previous level
420 ▾       } else if (length(tree[[as.character(level)]]) < 2^length(tree[[prev_level]])){
421           # increase missing nodes by difference
422           missing_nodes <- missing_nodes + 2^length(tree[[prev_level]]) - length(tree[[as.character(level)]])
423           # increase level error by 1
424           level_error <- level_error + 1
425 ▴       }
426
427         # nol equal to number of nodes on current level
428         nol <- length(tree[[as.character(level)]])
429         # if nol less than or equal to lnol
430 ▾       if (nol <= lnol){
431           # decrease height_balance by 1
432           height_balance <- height_balance - 1
433         # if nol greater than lnol
434 ▾       } else {
435           # set lnol to nol
436           lnol <- nol
437 ▴       }
438
439 ▴   }
```

```r
440
441       # set left_nodes to left half of children of root node
442       left_nodes <- head(tree[['1']], floor(length(tree[['1']])/2))
443       # set total left nodes to length of left_nodes
444       total_left_nodes <- length(left_nodes)
445       # set right_nodes to right half of children of root node
446       right_nodes <- tail(tree[['1']], ceiling(length(tree[['1']])/2))
447       # set total right nodes to length of right_nodes
448       total_right_nodes <- length(right_nodes)
449
450       # iterate through all nodes of left_nodes
451 ▾     while (length(left_nodes) != 0){
452         # if first node in left_nodes has children
453 ▾       if (!is.null(group_node_connections[[as.character(left_nodes[1])]])){
454           # increase total_left_nodes by number of children
455           total_left_nodes <- total_left_nodes + length(group_node_connections[[as.character(left_nodes[1])]])
456           # append children to left_nodes
457           left_nodes <- c(left_nodes, group_node_connections[[as.character(left_nodes[1])]])
458 ▴       }
459         # pop off first node in left_nodes
460         left_nodes <- left_nodes[-1]
461 ▴   }
462       # iterate through all nodes of right_nodes
463 ▾     while (length(right_nodes) != 0){
464         # if first node of right_nodes has children
465 ▾       if (!is.null(group_node_connections[[as.character(right_nodes[1])]])){
466           # increase total_right_nodes by number of children
467           total_right_nodes <- total_right_nodes + length(group_node_connections[[as.character(right_nodes[1])]])
468           # append children to right_nodes
469           right_nodes <- c(right_nodes, group_node_connections[[as.character(right_nodes[1])]])
470 ▴       }
471         # pop off first node in right_nodes
472         right_nodes <- right_nodes[-1]
473 ▴   }
474
475       # find nrin
476       nrin <- sprintf("%i/%i", total_nodes, 2^(length(names(tree)))-1)
477
478       # return statistics
479       return (c(extra_nodes, missing_nodes, level_error, total_left_nodes, total_right_nodes, nrin, height_balance))
480 ▴ }
481 ▴ ```
```

```r
483 # function to find value of tree based on analysis
484 ```{r}
485 tree_analysis <- function(temp_tree_statistic){
486   left_nodes <- as.numeric(temp_tree_statistic[4])
487   right_nodes <- as.numeric(temp_tree_statistic[5])
488   # if left_nodes < right_nodes
489   if (left_nodes < right_nodes) {
490     # side_balance is ratio of left_nodes to right_nodes
491     side_balance <- left_nodes / right_nodes
492   # left_nodes > right_nodes
493   } else{
494     # side_balance is ratio of right_nodes to left_nodes
495     side_balance <- right_nodes / left_nodes
496   }
497
498   # split nrin into numerator and denominator by /
499   fraction <- strsplit(temp_tree_statistic[6], "/")[[1]]
500   # get numerator of fraction
501   numerator <- as.numeric(fraction[1])
502   # get denominator of fraction
503   denominator <- as.numeric(fraction[2])
504   # calculate nrin
505   nrin <- numerator / denominator
506
507   node_error <- as.numeric(temp_tree_statistic[1]) + as.numeric(temp_tree_statistic[2])
508
509   # node_error has weight of 0.35
510   node_error_weight <- (1 - (node_error / denominator)) * 0.35
511   # level_error has weight of 0.25
512   level_error_weight <- (1 - (as.numeric(temp_tree_statistic[3]) / log((denominator+1), 2))) * 0.25
513   # nrin has weight of 0.15
514   nrin_weight <- nrin * 0.15
515   # side_balance has weight of 0.15
516   side_balance_weight <- side_balance * 0.15
517   # height_balance has weight of 0.1
518   height_balance_weight <- (1 - (as.numeric(temp_tree_statistic[7]) / log((denominator+1), 2))) * 0.1
519
520   # sum up weights to calculate temp_tree_statistic and return max(0, weights)
521   return (max(0, (node_error_weight + level_error_weight + side_balance_weight + height_balance_weight + nrin_weight)))
522 }
523 ```
```

Using these four helper functions, the program then performs the tree graph identification analysis on every group within the network.

```r
525  # Tree Graph Analysis
526  ```{r}
527  print("-------- Tree Graph Analysis --------")
528
529  # create new from_to list for finding cycle
530  tree_cycle_graph_list <- list()
531  # iterate through nodes in same group
532  for (group in unique(new_g_nodes$group)){
533    # add nodes in same group to tree_cycle_graph_list
534    tree_cycle_graph_list[[as.character(group)]] <- new_g[[as.character(group)]]
535  }
536
537  # create potential root nodes list
538  potential_root_nodes <- list()
539
540  # get each group in tree_cycle_graph_list
541  for (group in names(tree_cycle_graph_list)){
542    # group must have more than 1 node
543    if (length(tree_cycle_graph_list[[group]]) > 1) {
544      # if there is a cycle
545      if (has_cycle(tree_cycle_graph_list[[group]], group_node_connections)) {
546        # print statement
547        print(sprintf("  Group %s has a cycle so not a tree", group))
548      # group does not have cycle
549      } else {
550        # look at edges in edges leaving groups
551        for (edge in out_edge_groups_f[[group]]){
552          # extract number from edge
553          number <- gregexpr("[0-9]+", edge)
554          result <- regmatches(edge, number)
555          #get first number (the from node)
556          node <- as.numeric((result[[1]][1]))
557          # add node to potential root nodes list
558          potential_root_nodes[[group]] <- unique(c(potential_root_nodes[[group]], node))
559        }
560        # look at edges in edges entering groups
561        for (edge in out_edge_groups_s[[group]]){
562          # extract number from edge
563          number <- gregexpr("[0-9]+", edge)
564          result <- regmatches(edge, number)
565          #get first number (the from node)
566          node <- as.numeric((result[[1]][2]))
567          # add node to potential root nodes list
568          potential_root_nodes[[group]] <- unique(c(potential_root_nodes[[group]], node))
569        }
570
571        # if there are no nodes for the group in potential_root_nodes list
572        if (length(potential_root_nodes[[group]]) == 0){
573          # add every node in the group to potential_root_nodes list
574          for (node in tree_cycle_graph_list[[group]]){
575            # only keep unique nodes
576            potential_root_nodes[[group]] <- unique(c(potential_root_nodes[[group]], node))
577          }
578        }
579
580      }
581    }
582  }
```

```r
583
584  # create best root node analysis list (keeps track of best root node, best tree value, and best tree statistics)
585  broot_node_analysis <- list()
586
587  # group index
588  group_i <- 1
589  # iterate through groups in potential_root_nodes list
590  for (group in names(potential_root_nodes)){
591    # statistic of best tree in group
592    best_tree_statistic <- c()
593    # analysis of best tree in group
594    best_tree_analysis <- 0
595    # root node in group that led to best tree
596    best_node <- 0
597    # iterate through nodes in groups in potential_root_nodes
598    for (node in potential_root_nodes[[group]]){
599      # create a tree for root node
600      tree <- create_tree(list(), 1, node)
601
602      # get statistics of temp_tree
603      temp_tree_statistic <-  tree_statistics(tree, group_node_connections)
604      # get analysis of temp_tree
605      temp_tree_analysis <- tree_analysis(temp_tree_statistic)
606
607      # if temp_tree analysis >= best_tree_analysis
608      if (temp_tree_analysis >= best_tree_analysis){
609        # set best_tree_statistic to temp_tree_statistic
610        best_tree_statistic <- temp_tree_statistic
611        # set best_tree_analysis to temp_tree_analysis
612        best_tree_analysis <- temp_tree_analysis
613        # set best_node to node
614        best_node <- node
615      }
616    }
617    # save best_tree_statistic in broot_node_analysis
618    broot_node_analysis[[group_i]] <- c(best_node, best_tree_analysis, best_tree_statistic)
619
620    # Group i analysis
621    print(sprintf("Group %s:", group))
622    # if the tree statistic for the group is greater than or equal to 0.8 then the group is a tree
623    if (best_tree_analysis >= 0.7){
624      # print statements
625      print(sprintf("  Group %s is a tree graph with node %i being the root node", group, as.numeric(broot_node_analysis[[group_i]][1])))
626      print(sprintf("  %f%% of the group shares the same properties as those of an ideal tree graph", 100*as.numeric(broot_node_analysis[[group_i]][2])))
627      print(sprintf("  %i Node Errors (%i Extra Nodes, %i Missing Nodes), %i Level Errors, %i Left Nodes and %i Right Nodes, %s (%f%%) NRIN, %i Height Balance",
628  as.numeric(broot_node_analysis[[group_i]][3])+as.numeric(broot_node_analysis[[group_i]][4]), as.numeric(broot_node_analysis[[group_i]][3]), as.numeric(broot_node_analysis[[group_i]][4]),
629  as.numeric(broot_node_analysis[[group_i]][5]), as.numeric(broot_node_analysis[[group_i]][6]), as.numeric(broot_node_analysis[[group_i]][7]), broot_node_analysis[[group_i]][8],
630  100*as.numeric(eval(parse(text = broot_node_analysis[[group_i]][8]))),  as.numeric(broot_node_analysis[[group_i]][9])))
631    # if the tree statistic for the group is less than 0.8 then the group is not a tree
632    }else{
633      # print statements
634      print(sprintf("  Group %s is not a tree and has a structure closest to a tree graph when node %i is the root node", group, as.numeric(broot_node_analysis[[group_i]][1])))
635      print(sprintf("  %f%% of the group sharing the same properties as those of an ideal tree graph", 100*as.numeric(broot_node_analysis[[group_i]][2])))
636      print(sprintf("  %i Node Errors (%i Extra Nodes, %i Missing Nodes), %i Level Errors, %i Left Nodes and %i Right Nodes, %s (%f%%) NRIN, %i Height Balance",
637  as.numeric(broot_node_analysis[[group_i]][3])+as.numeric(broot_node_analysis[[group_i]][4]), as.numeric(broot_node_analysis[[group_i]][3]), as.numeric(broot_node_analysis[[group_i]][4]),
638  as.numeric(broot_node_analysis[[group_i]][5]), as.numeric(broot_node_analysis[[group_i]][6]), as.numeric(broot_node_analysis[[group_i]][7]), broot_node_analysis[[group_i]][8],
639  100*as.numeric(eval(parse(text = broot_node_analysis[[group_i]][8]))),  as.numeric(broot_node_analysis[[group_i]][9])))
640    }
641    # increase group_i by 1
642    group_i <- group_i + 1
643  }
644  ```
```

Next, the program determines which communities are line branch graphs by comparing every community against the properties of an ideal line branch graph.

Like the analyzation of the cloistered graph and tree graph, I created several helper functions to assist with the analysis. The first helper function finds the longest line that ends with a given node in a group. The second function finds every reversed node connection within the same group. The third helper function examines the group in its line branch structure and calculates specific statistics, which will be discussed in Section 3.1. The fourth helper function evaluates these statistics to determine how similar the line branch structure is to an ideal line branch graph.

```r
642  #LINE BRANCH GRAPH
643
644  # function to find longest line in group
645  ```{r}
646  longest_line <- function(node, group_node_connections, visited = c()) {
647     # if the initial node has been visited
648     if (node %in% visited){
649        # return length of 0 and an empty path
650        return(list(length = 0, path = c()))
651     }
652
653     # initial node has not been visited so add node to visited
654     visited <- c(visited, node)
655
656     # if node does not have connections to other nodes
657     if (is.null(group_node_connections[[as.character(node)]]))
658        # return length of 1 (just initial node) and the node as the path
659        return(list(length = 1, path = c(node)))
660
661     # longest line length
662     max_length <- 0
663     # path of longest line
664     max_path <- c()
665
666     # explore all nodes connected to node
667     for (near_node in group_node_connections[[as.character(node)]]) {
668        # recursive call
669        line <- longest_line(near_node, group_node_connections, visited)
670        # if the length of the line is greater than max_length
671        if (line$length > max_length) {
672           # set max_length to the line's length
673           max_length <- line$length
674           # set max_path to the line's path
675           max_path <- line$path
676        }
677     }
678
679     # return max_length + 1 to include current node and max_path
680     return(list(length = max_length + 1, path = c(node, max_path)))
681  }
682  ```
```

```r
684   # create reverse group node connections list
685   ```{r}
686   # function to create reverse group node connections
687   reverse_connections <- function(connections) {
688     # initialize reverse group node connections list
689     reversed <- list()
690     # iterate through all from nodes in connections list (from nodes: nodes that point to other nodes)
691     for (from in names(connections)) {
692       # iterate through all to nodes (to nodes: nodes that are pointed to)
693       for (to in connections[[from]]) {
694         # if to node is not in the names of reversed list
695         if (!as.character(to) %in% names(reversed)) {
696           # initialize nodes pointed at to node in reversed list
697           reversed[[as.character(to)]] <- c()
698         }
699         # add from node to reversed list at index to node
700         reversed[[as.character(to)]] <- c(reversed[[as.character(to)]], as.numeric(from))
701       }
702     }
703     # return reversed list
704     return(reversed)
705   }
706
707   # create reversed connections
708   reversed_gn_connections <- reverse_connections(group_node_connections)
709   ```
```

```r
711   # function to find statistics
712   ```{r}
713   lb_statistics <- function(line, group_nodes, group_node_connections){
714     # set branch_nodes to number of nodes connected to last node in line
715     branch_nodes <- group_node_connections[[as.character((tail(line, n=1)))]]
716
717     # set line_error to (number of nodes in group) - ((number of nodes in line) + (branching nodes))
718     line_error <- length(group_nodes) - (length(line) +  length(branch_nodes))
719     # set llrill to (number of nodes in line) / (2*log(length(branch_nodes) + 1) + 5)
720     llrill <- sprintf("%i/%i", length(line), as.integer((2*log(length(branch_nodes) + 1) + 5)))
721
722     # initialize branch_continuation to 0
723     branch_continuation <- 0
724     # iterate through all branch_nodes
725     for (node in branch_nodes){
726       # add number of nodes connected to node to branch_continuation
727       branch_continuation <- branch_continuation + length(group_node_connections[[as.character(node)]])
728     }
729
730     # initialize branch_positions
731     branch_positions <- c()
732     # set line_index to 0
733     line_index <- 0
734
735     # iterate through nodes in line
736     for (node in line[-length(line)]){
737       # increase line_index by 1
738       line_index <- line_index + 1
739       # if the node has more than 1 connection in group
740       if (length(group_node_connections[[as.character(node)]]) > 1){
741         # add position relative to line to branch_positions
742         branch_positions <- c(branch_positions, sprintf("%i/%i", line_index, length(line)))
743       }
744     }
745   }
746   if (length(branch_positions) == 0){
747     branch_positions = 0
748   }
749
750   # return statistics
751   return (c(line_error, llrill, branch_continuation, branch_positions))
752 }
753 ```
```

```r
755 ▾ # function to find value of line branch based on analysis
756 ▾ ```{r}
757 ▾ lb_analysis <- function(lb_statistic){
758       # split llrill into numerator and denominator by /
759       fraction <- strsplit(lb_statistic[2], "/")[[1]]
760       # get numerator of fraction
761       line_length <- as.numeric(fraction[1])
762       # get denominator of fraction
763       i_line_length <- as.numeric(fraction[2])
764       # calculate llrill
765       llrill <- line_length / i_line_length
766
767       # calculate how many branch nodes there are in group (isolate branch_nodes in ideal line length = 2*log(length(branch_nodes) + 1) + 5)
768       branch_nodes <- exp((i_line_length - 5)/2)-1
769       # calculate line_error such that line_error is relative to number of nodes in ideal line length + branching nodes
770       line_error <- as.integer(lb_statistic[1]) / (i_line_length + as.integer(branch_nodes))
771
772       # calculate branch_continuation relative to number of nodes in line
773       branch_continuation <- as.integer(lb_statistic[3]) / line_length
774
775       # initialize penalty_sum
776       penalty_sum <- 0
777       # iterate through positions of branching nodes in line relative to line
778 ▾     for (position_r in lb_statistic[4]){
779         # split position relative into numerator and denominator by /
780         fraction <- strsplit(position_r, "/")[[1]]
781         # get position (numerator)
782         position <- as.numeric(fraction[1])
783         # add ((number of nodes in line) - (position of branching node in line)) / (number of nodes in line) to penalty_sum
784         penalty_sum <- penalty_sum + ((line_length - position) / line_length)
785 ▴     }
786
787       # calculate branch_position (1 - penalty_sum / (number of branching nodes in line))
788       branch_positions <- 1 - penalty_sum/length(lb_statistic[4])
789
790       # line_error has weight of 0.4
791       line_error_weight <- (1 - line_error) * 0.4
792       # llrill has weight of 0.3
793       llrill_weight <- llrill * 0.3
794       # branch_continuation has weight of 0.2
795       branch_continuation_weight <- (1 - branch_continuation) * 0.2
796       # branch_positions has weight of 0.1
797       branch_positions_weight <- branch_positions * 0.1
798
799       # sum up weights to calculate lb_statistic and return max(0, weights)
800       return (max(0, (line_error_weight + llrill_weight + branch_continuation_weight + branch_positions_weight)))
801 ▴ }
802 ▴ ```
```

Using these four helper functions, the program then performs the line branch graph identification analysis on every group within the network.

```r
803
804 # LINE BRANCH GRAPH ANALYSIS
805 ```{r}
806 print("-------- Line Branch Graph Analysis --------")
807
808 # create new from_to list for finding cycle
809 lb_cycle_graph_list <- list()
810 # iterate through nodes in same group
811 for (group in unique(new_g_nodes$group)){
812   # add nodes in same group to lb_cycle_graph_list
813   lb_cycle_graph_list[[as.character(group)]] <- new_g[[as.character(group)]]
814 }
815
816 # create potential end nodes list
817 potential_end_nodes <- list()
818
819 # get each group in lb_cycle_graph_list
820 for (group in names(lb_cycle_graph_list)){
821   # group must have more than 1 node
822   if (length(lb_cycle_graph_list[[group]]) > 1) {
823     # if there is a cycle
824     if (has_cycle(lb_cycle_graph_list[[group]], group_node_connections)) {
825       # print statement
826       print(sprintf("   Group %s has a cycle so not a line branch", group))
827     # group does not have cycle
828     } else {
829       potential_end_nodes[[group]] <- lb_cycle_graph_list[[group]]
830     }
831   }
832 }
```

```r
833
834 # create best start node analysis list (keeps track of best start node, best line branch value, and best line branch statistics)
835 bstart_path <- list()
836
837 # iterate through groups in potential_end_nodes list
838 for (group in names(potential_end_nodes)){
839   # path of longest line
840   temp_longest_path <- c()
841   # most branching nodes
842   largest_branch <- 0
843   # longest path with most branching
844   longest_path <- list()
845   # iterate through nodes in groups in potential_end_nodes
846   for (node in potential_end_nodes[[group]]){
847     # if the node has more branches than largest_branch
848     if (length(group_node_connections[[as.character(node)]]) > largest_branch){
849       # set largest_branch to the number of branches of the node
850       largest_branch <- length(group_node_connections[[as.character(node)]])
851       # set the longest_path to the longest line of nodes in the group that ends with node
852       longest_path <- longest_line(node, reversed_gn_connections)$path
853     # if the node has equal number of branches to largest_branch
854     } else if (length(group_node_connections[[as.character(node)]]) == largest_branch){
855       # find the longest line of nodes in the group that ends with node
856       temp_longest_path <- longest_line(node, reversed_gn_connections)$path
857       # if the length of this line is greater than the length of longest_path
858       if (length(temp_longest_path) > length(longest_path)){
859         # set longest_path to the line
860         longest_path <- temp_longest_path
861       }
862     }
863   }
864   # save path using best start node to bstart_path
865   bstart_path[[group]] <- c(rev(longest_path))
866 }
```

```r
867
868 bstart_node_analysis <- list()
869 # iterate through groups in potential_end_nodes list
870 for (group in names(potential_end_nodes)){
871   # get line branch statistic using path in bstart_path at the group as the start node
872   lb_statistic <- lb_statistics(bstart_path[[group]], lb_cycle_graph_list[[group]], group_node_connections)
873   # get line branch value from lb_analysis
874   lb_value <- lb_analysis(lb_statistic)
875   # add line branch value and lb_statistic to bstart_node_analysis
876   bstart_node_analysis[[group]] <- c(lb_value, lb_statistic, bstart_path[[group]])
877   # Group i analysis
878   print(sprintf("Group %s:", group))
879   # if the line branch analysis for the group is greater than or equal to 0.8 then the group is a line branch
880   if (lb_value >= 0.7){
881     # print statements
882     print(sprintf("   Group %s is a line branch with node %i being the start node", group, as.numeric(bstart_node_analysis[[group]][6])))
883     last_i <- length(bstart_node_analysis[[group]])
884     print(sprintf("   The line is from node %s", paste(as.character(as.integer(bstart_node_analysis[[group]][6:last_i])), collapse=" to ")))
885     print(sprintf("   Node %s then branches into nodes %s", as.character(as.integer(bstart_node_analysis[[group]][last_i])),
     paste(as.character(group_node_connections[[as.character(bstart_node_analysis[[group]][last_i])]][1:length(group_node_connections[[as.character(bstart_node_analysis[[group]][last_i])]])]),
     collapse=", ")))
886     print(sprintf("   %f%% of the group shares the same properties as those of an ideal line branch", 100*as.numeric(bstart_node_analysis[[group]][1])))
887     print(sprintf("   %i Line Errors, %s (%f%%) LLRBN, %i Branch Continuation, %s Branch Positions", as.numeric(bstart_node_analysis[[group]][2]), bstart_node_analysis[[group]][3],
     100*as.numeric(eval(parse(text = bstart_node_analysis[[group]][3]))), as.numeric(bstart_node_analysis[[group]][4]), bstart_node_analysis[[group]][5]))
888   # if the line branch analysis for the group is less than 0.8 then the group is not a line branch
889   }else{
890     # print statements
891     print(sprintf("   Group %s is not a line branch and has a structure closest to a line branch when node %i is the start node", group, as.numeric(bstart_node_analysis[[group]][6])))
892     last_i <- length(bstart_node_analysis[[group]])
893     print(sprintf("   The line is from node %s", paste(as.character(as.integer(bstart_node_analysis[[group]][6:last_i])), collapse=" to ")))
894     print(sprintf("   Node %s then branches into nodes %s", as.character(as.integer(bstart_node_analysis[[group]][last_i])),
     paste(as.character(group_node_connections[[as.character(bstart_node_analysis[[group]][last_i])]][1:length(group_node_connections[[as.character(bstart_node_analysis[[group]][last_i])]])]),
     collapse=", ")))
895     print(sprintf("   %f%% of the group shares the same properties as those of an ideal line branch", 100*as.numeric(bstart_node_analysis[[group]][1])))
896     print(sprintf("   %i Line Errors, %s (%f%%) LLRBN, %i Branch Continuation, %s Branch Positions", as.numeric(bstart_node_analysis[[group]][2]), bstart_node_analysis[[group]][3],
     100*as.numeric(eval(parse(text = bstart_node_analysis[[group]][3]))), as.numeric(bstart_node_analysis[[group]][4]), bstart_node_analysis[[group]][5]))
897   }
898 }
899 ```
```

## 2.4 What does it do

The program allows for both visualization and analysis of networks by graphing the entire network and identifying subgraphs among communities within the network. The program also incorporates a community detection algorithm to split the network into distinct communities. The network is then visually represented as a directed graph, highlighting the connections between nodes (users) and the direction of the connection (flow of the information).

Once the network is split into communities, the program studies the communities by analyzing properties to identify subgraphs. By analyzing structural properties and the flow of information, the program calculates if a community is a cloistered graph and assesses the similarities of communities to the ideal form of tree and line branch graphs. The program is a dynamic tool designed to visualize the network, analyze its structure, and evaluate the flow of information.

## 2.5 How does it work

To provide a clear understanding of how the program operates and identifies subgraphs within the network, the following lines of pseudocode break down and explain each segment of the code and its underlying ideas

```
1 load packages
2 read csv file
3 split data in csv file into nodes that point to other nodes (from id list) and nodes that
are pointed to (to id list)

4 iterate through every row in csv file to create directed list (from_to list) using from
and to id list

### GRAPHING
5 create a graph object (from igraph package)
6 apply Louvain algorithm on graph object to create new graph object

7 iterate through nodes in graph object to find unique nodes id's and number of nodes in
each group

8 iterate through unique node id's to find node labels (node id but as a char)
9 iterate through number of nodes in each group to find node groups (group that node
belongs to)
10 iterate through unique node id's to create node titles (title of the node)
11 create node data frame that holds the following info
    unique node id's
    node labels
    node groups
    node titles
```

12 iterate through from id list and add 1 to every node id to create from id plus one list
(nodes/indexing in R must start at 1)
13 iterate through to id list and add 1 to every node id to create to id plus one list

14 iterate through from id plus one list to create edge labels (every edge between two
nodes)
15 iterate through from id plus one list to create edge titles (title of every edge)
16 create edge data frame that holds the following info
    from id plus one list
    to id plus one list
    edge labels
    edge titles

17 using node data frame and edge data frame, graph network using visNetwork
18 using graph object, graph network using igraph


### ANALYSIS
# CLOISTERED GRAPH
19 iterate through from_to list and find all node connections for every node

20 iterate through every edge
21     if nodes in edges belong to the same group
22         add the edge to same group edge list
23         add nodes to group node connections list
24     else
25         add edge to different group edge list

26 iterate through graph object to find maximum possible number of edges (MPNE) for every
group

27 iterate through every group
28      find number of edges leaving/entering group
29      find number of edges in group
30      calculate how close group is to a complete graph using number of edges in group and
MPNE (percentage)
31     find how many edges are missing for group to be a complete graph

32     if number of edges leaving/entering a group is less than 2% of MPNE
33         then group is a cloistered graph
34         print statistics and analysis
35     else
36         group is not a cloistered graph
37         print statistics and analysis

38     if percentage is 100% and group is a cloistered graph and group is disconnected
39         then group is an ideal cloistered graph

```
40            print analysis
41        else if percentage is 100% and group is not a cloistered graph
42                then group is only a complete graph
43                print analysis
44        else
45                group is value of percentage of a complete graph
46                print statistics and analysis



# TREE GRAPH
47 function create_tree(tree = list(), level = 1, node)
     """
48      create_tree makes a tree structure given a node
49      :tree: tree structure is contained in a list
50      :level: current level of tree as an int
51      :node: node as an int to either add to tree or add children of node to tree
52      :return: returns tree structure with node as the root node
     """
53      if tree has no nodes
54            add node to tree at level 1
55      else
55            add the node to tree at current level

56      get children of node (its connections in group node connections list)

57      if there are children
58            iterate through every child
59                    recursively call create_tree(tree, level+1, child)

60      return tree



61 function has_cycle(node_list, group_node_connections)
     """
62      has_cycle checks if there a cycle given the nodes in a group and the group node
connections list
63      :node_list: vectors of ints for nodes in a group
64      :group_node_connections: group node connections list
65      :return: boolean value indicating whether a group has a cycle
     """
66            function dfs(node)
          """
67            dfs uses dfs to check if a node has been seen before
68            :node: int of node id that needs to be checked if node has been seen before
69            :return: boolean value that indicates whether node has been seen before
          """
```

70          if the node has not been visited before
71                set value in visited and rec_stack list at index node to TRUE
72                if node has other connections in group
73                    iterate through all connections
74                        if other nodes have not been visited and dfs(other nodes)
75                            return TRUE
76                        else if other nodes are in rec_stack
77                            return TRUE

78          set value in rec_stack at index node to FAlSE
79          return FALSE


80     create visited and recursion stack list

81     iterate through nodes in node_list
82          set every value of visited and rec_stack at index node to FALSE as we have not
seen these nodes yet

83     iterate through nodes in node_list
84          if dfs(node)
85                return TRUE

86     return FALSE


87 function tree_statistics(tree, group_node_connections)
    """
88     tree_statistics finds the statistics of a tree
89     :tree: tree structure as a list
90     :group_node_connections: group node connections list
91     :return: statistics of a tree as a vector
    """
92     initialize line error, extra nodes, missing nodes, height balance, nrin, total
nodes, lnol, nol variables

93     iterate through every level in tree
94          add number of nodes on level to total nodes

95          if number of nodes on level is greater than 2 ^ (number of nodes on previous
level)
96                add difference to extra nodes
97                increase level error by 1
98          else if less
99                add difference to missing nodes
100               increase level error by 1

```
101          set nol to number of nodes on level
102          if nol is less than or equal to lnol
103              increase height balance by 1
104          else
105              set lnol to nol

106    set left nodes list to left half of children of root node
107    set total left nodes to number of left nodes
108    set right nodes list to right half of children of root node
109    set total right nodes to number of right nodes

110    while left nodes list is not empty
111        if first node in left nodes list has children
112              append children to left nodes list
113              add number of children to total left nodes
114        pop off first node in left nodes list

115    while right nodes list is not empty
116        if first node in right nodes list has children
117              append children to right nodes list
118              add number of children to total right nodes
119        pop off first node in right nodes list

120    set nrin = (total number of nodes in group) / (2 ^ (number of levels in tree) - 1)

121    return (extra nodes, missing nodes, level error, total left nodes, total right
nodes, nrin, height balance)


122 function tree_analysis(tree statistics)
    """
123    tree_analysis analyzes statistics of a tree and gives a rating of how well the tree
is
124    :tree statistics: statistics of a tree as a vector
125    :return: float
    """
126    get left nodes and right nodes from tree statistics (total left nodes and total
right nodes)
127    if left nodes is less than right nodes
128        side balance = left nodes / right nodes
129    else
130        side balance = right nodes / left nodes

131    calculate nrin as a float

132    get node error = extra nodes + missing nodes
```

```
133    node error has weight of 0.35
134    level error has weight of 0.25
135    nrin has weight of 0.15
136    side balance has weight of 0.15
137    height balance has weight of 0.1


138    find value of node error given its weight and its relative value
139    find value of level error given its weight and its relative value
140    find value of nrin given its weight
141    find value of side balance given its weight
142    find value of height balance given its weight and its relative value
142    sum up values


144    return max(0, sum of values)



145 create potential root nodes list
146 iterate through every group in network
147    if group has more than one node
148        check if group has a cycle by calling has_cycle(group, group_node_connections)
149            if group has a cycle then group is not a tree
150        else
151            iterate through edges leaving/entering group
152                add nodes in group connected to other groups to potential root nodes
list

153            if there are no nodes in potential root nodes list
154                add every node in group to potential root nodes list

155 create best root node analysis list
156 iterate through every group
157    set best tree statistics, best tree analysis, and best node to 0
158    iterate through nodes in potential root nodes list that belong to same group
159        create tree by calling create_tree(list(), 1, node)
160        get statistics of tree by calling tree_statistics(tree,
group_node_connections)
161        get analysis of tree by calling tree_analysis(tree statistics)

162        if tree analysis > best tree analysis
163            set best tree statistics to tree statistics
164            set best tree analysis to tree analysis
165            set best node to node


166    save best node, best tree analysis, and best tree statistics to best root node
analysis list at index group
```

```
167    if best tree analysis is greater than or equal to 0.7
168        then group is a tree with best node being the root node
169        print statistics and analysis
170    else
171        group is not a tree but best node being the root node leads to a tree structure
closest to an ideal tree compared to other nodes in group as the root node
172        print statistics and analysis


# LINE BRANCH GRAPH
173 function longest_line(node, group_node_connections, visited = c())
    """
174    longest_line finds the longest line of nodes in a group that ends with node
175    :node: node id as an int of last node in longest line of nodes for a group
176    :group_node_connections: group node connections list
177    :visited: visited nodes in group as a list
178    :return: list with length of longest line and nodes in longest line
    """
179    if node is in visited list
180        return (length = 0, path = c())


181    add node to visited list


182    if node does not have connections in group
183        return (length = 1, path = c(node))


184    set max length and max path to 0
185    iterate through nodes that node is connected to in group
186        recursively call longest_line(near_node, group_node_connections, visited)
187        if the length of line is greater than max length
188            set max length to length of line
189            set max path to path of line


190    return (length = max length, path = c(node, max path))


191 function reverse_connections(connections)
    """
192    reverse_connections reverses a connections list (reverses 'from nodes' -> 'to nodes'
to 'to nodes' -> 'from nodes')
193    :connections: connections list
194    :return: connections list but reversed
    """
195    initialize reversed list
196    iterate through all from nodes in connections list
```

197        iterate through all to nodes that are connected to from node
198             add to node and its from node to reversed connections list


199     return reversed connections list



200 create reversed group node connections by calling reverse_connections(group node connections)




201 function lb_statistics(line, group_nodes, group_node_connections)
    """
202     lb_statistics finds the statistics of a line branch
203     :line: line of nodes as a vector of ints
204     :group_nodes: all nodes in the group as a vector of ints
205     :group_node_connections: group node connections list
206     :return: statistics of a line branch as a vector
    """
207     set branching nodes to the connections of last node in line


208     set line error = group_nodes - ((number of nodes in line) + (number of branching nodes))


209     set llrill = (number of nodes in line) / (2*log(number of branching nodes + 1) + 5)


210     initialize branch continuation to 0
211         iterate through all nodes in branching nodes
212         add number of connections each node has to branch continuation


213     initialize branch positions to 0
214         iterate through all nodes in line
215             if node has more than one connection in group
216                 add the position of node to branch positions


217     return (line error, llrill, branch continuation, branch positions)



218 function lb_analysis(lb_statistic)
    """
219     lb_analysis analyzes statistics of a line branch and gives a rating of well the line branch is
220     :lb_statistic: statistics of a line branch as a vector
221     :return: float
    """
222     set line length to numerator of llrill

223     set ideal line length to denominator of llrill
224     calculate llrill as a float


225     calculate how many branch nodes there are in group


226     calculate line error = line error / ((ideal line length) + (number of branch nodes))


227     calculate branch continuation = branch continuation / line length



228     set penalty sum to 0
229     iterate through branch positions
230         get node position
231         add (line length - node position) / (line length) to penalty sum
232     calculate branch positions = 1 - (penalty sum) / (number of positions in branch
positions)


233     line error has weight of 0.4
234     llrill has weight of 0.3
235     branch continuation has weight of 0.2
236     branch position has weight of 0.1


237     find value of line error given its weight and its relative value
238     find value of lrrill given its weight
239     find value of branch continuation given its weight and its relative value
240     find value of branch position given its weight
241     sum up values


242     return max(0, sum of values)



243 create potential start nodes list
244 iterate through every group in network
245     if group has more than one node
246         if has_cycle(group, group_node_connections)
247             group has a cycle so not a line branch
248         else
249             add every node in group to potential start nodes list

250 create best start path list
251 iterate through every group
252     set temporary longest path, largest branch, longest path to 0
253     iterate through nodes in potential start nodes list that belong to same group
254         if node has more connections than largest branch
255             set largest branch to number of connections node has

256            set longest path to path ending with node by calling longest_line(node,
257 reversed group node connections) and getting its path
258         else if node has number of connections same as largest branch
259             set temporary longest path to path from longest_line(node, reversed group
node connections)
260             if number of nodes in temporary longest path is greater than number of
nodes in longest path
261                 set longest path to temporary longest path

262    save longest path to best start path list at index group

263 create best start node analysis list
264 iterate through every group
265    get statistics of line branch by calling lb_statistics(path at index group in best
start path list, nodes in group, group node connections)
266    get analysis of line branch by calling lb_analysis(lb_statistic)
267    add line branch analysis and statistic and path to best start node analysis list at
index group

268    if line branch analysis is greater than or equal to 0.7
269        then group is a line branch
270        print statistics and analysis
271    else
272        group is not a line branch
273        print statistics and analysis

# III. Detailed Project Description

## 3.1 Domain Model

1. Cloistered Graph

Definitions:

**Cloistered Graph:** Group that has few to none edges leaving or entering the group

**Ideal Cloistered Graph:** Complete graph that is disconnected

**MPNE:** Maximum possible number of edges in a group. Calculated by taking every node in group and finding all possible connections between every node in same group.

Method to identify cloistered graphs: Compare Number of Edges to Max Possible Number of Edges (MPNE)

1. Find how many edges are in the group
2. Find the number of edges that leave/enter the group
3. Calculate MPNE
4. If number of edges leaving/entering the group is less than 2% of MPNE of that group, then group is a cloistered group. If greater than 2% of MPNE, then group is not a cloistered group
5. Divide number of edges by MPNE. This percentage represents how close the group is to that of a complete graph.

Example:

Group has 12 nodes with 11 inner group edges (1 node connects to the other 11 nodes in the group). The group also has 2 outer group edges.

11 edges inside group
2 edges leaving/entering the group
MPNE: 12 nodes * (12 nodes – 1 node) / 2 = 66 edges
*Because 2 / 66 * 100% = 3.03% (> 2%) of nodes leave/enter the group, the group is not a cloistered group*

Accuracy Percentage: 11 edges / MPNE * 100% = 16.67%
*Group is 16.67% of a complete graph*

2. Tree

Definitions:

**Ideal Tree:** A noncyclic graph where the graph develops from a single node (root node) and each node splits into 2 child nodes

**Tree Finder Algorithm (TFA):** Properties to check/find to identify whether a group is a tree

What the TFA checks for and analyzes:
- Node Error
- Level Error
- Side Balance
- Height Balance
- Nodes Relative to Number of Ideal Nodes

**Node Error:** Sum of missing nodes and extra nodes that should not exist in a tree (if there is an extra node that should not exist, then assume for the remaining levels that the deeper nodes should have two additional child nodes to ensure that the error is not counted more than once). This statistic helps identify deviations from an ideal tree by analyzing node counts at ever level.

**Level Error:** Number of levels with an incorrect number of nodes. For each level, calculate the expected number of nodes based on the previous level (e.g. if level 2 has 3 nodes, level 3 should have 6 nodes). This statistic is assessed on a per-level basis and evaluates how closely the tree follows the property of each parent node having two children nodes.

**Side Balance:** Determines how evenly distributed the nodes are between the left and right sides of the tree. Split the root node into two sides: left and right. Then count the number of nodes on each side. This statistic finds how balanced the tree is in terms of side distribution.
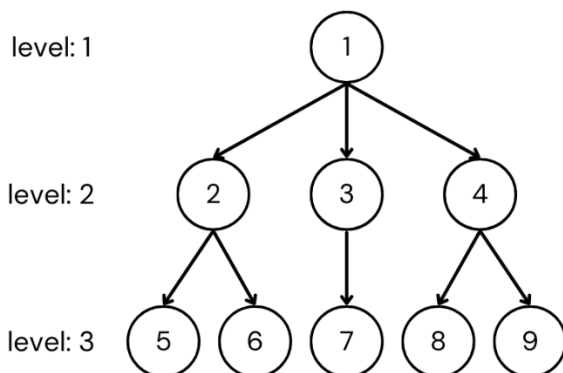
**Height Balance:** Evaluates the distribution of nodes across the levels of the tree. Iterate through all levels and find how many levels have fewer nodes than the largest number of nodes observed so far on a level (lnol). This statistic determines whether the tree continues to grow in terms of number of nodes per level.

**Nodes Relative to Number of Ideal Nodes (NRIN):** Total number of nodes in tree relative to number of nodes in an ideal tree with the same depth. This statistic determines how close the tree is to an ideal tree by comparing the actual number of nodes to the theoretical ideal number of nodes given the depth level.

Method to identify trees:

a. Check for cycles: if the group contains a cycle, it cannot be a tree
b. Identify potential root nodes: create a list of potential root nodes for each group if no cycle exists in that group
c. Evaluate each root node: for each node in the group's potential root nodes list
    i. Consider current node to be root node
    ii. Create tree structure using root node
    iii. Apply TFA on tree structure to get statistics
    iv. Compare statistics to best statistics observed so far. If statistics are better, then save statistics as best statistics and save root node as best root node
d. Repeat until group potential root nodes exhausted: repeat step c) until list of nodes from same group in potential root nodes list is exhausted
e. Output results: from the list of best statistics, identify the group and output the best node and statistics that form a tree closest to an ideal tree when the best node is used as the root node

Example:



Node Error: 2
Node 1 should only have 2 children (has 1 extra child).
Node 3 should have 2 children (has 1 missing child)
1 extra child + 1 missing child = 2 node errors

Level Error: 2
Level 2 should have 2 nodes (has 3 nodes)
Level 3 should have 6 nodes (has 5 nodes)
Level 2 and Level 3 have errors, so Level Error = 2

Side Balance = 1
Both Left Side and Right Side have same number of nodes

Height Balance: 3
Number of nodes increases every level, because there are 3 levels then Height Balance = 3

NRIN: 9/7
9 total nodes and ideal tree of height 3 should have 7 nodes

Value of Tree:
$(1 - (2/7)) * (0.35) + (1 - (2/3)) * 0.25 + 9/7 * 0.15 + 1 * 0.15 + (1 - (3/3)) * 0.1 = 0.67619047619$

## 3. Line Branch

Definitions:

**Ideal Line Branch:** a sequence of nodes connected in a path sequence where the final node in the path branches into multiple children nodes

**Line Branch Finder Algorithm (LBFA)**: Properties to check/find to identify whether a group is a line branch

What the LBFA checks for and the analyzes:

- Line Error
- Line Length Relative to Ideal Line Length
- Branch Continuation
- Branch Positions

**Line Error:** Finds how many errors there are in the line (nodes in the line should not branch before the final node in line). Calculate the number of nodes that branch prematurely, as well as the total number of nodes in branches that should not exist. This statistic determines how much the line branch deviates from an ideal line branch.

**Line Length Relative to Ideal Line Length (LLRILL):** Compares actual number of nodes in a line to the ideal line length. The ideal line length L is calculated by

$$L = a \cdot log(B + 1) + c$$

where $a$ is the scaling factor, $B$ is the number of branching nodes after the final node in the line, and $c$ is the baseline line length. This statistic evaluates how accurately the group aligns with the ideal line length.
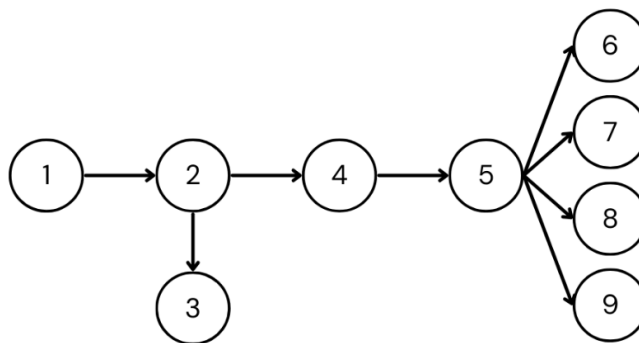
**Branch Continuation:** Finds whether additional branching occurs after the initial branching from the final node in the line. This statistic determines how closely the group follows the structure of an ideal line branch.

**Branch Positions:** Finds the positions of nodes in the line that branch early. This statistic will tell the user whether branching occurs, and if so, whether it happens earlier or later in the line branch.

Method to identify line branch:

a. Check for cycles: if a cycle exists in the group, it cannot be a line branch

b. Identify potential ending nodes: create a list of potential ending nodes for each group if no cycle exists for that group. Then initialize most branching nodes and best path (longest path in the group that ends the line with the node that led to most branching nodes) and iterate over each node in the group's potential ending nodes list

    i. Consider current node as the ending node of the line

    ii. Find number of branching nodes ending node has

    iii. If number of branching nodes is greater than most branching nodes, get longest path in the group that ends with the ending node and set best path to the current path and set most branching nodes to number of branching nodes from ending node. If number of branching nodes is same as most branching nodes, get longest path in the group that ends with the ending node. Compare the length of longest path to length of best path. If longest path is longer than best path, then set best path to longest path and most branching nodes to number of branching nodes from ending node

c. Evaluate each best path for each group: apply LBFA on each best path for each group to get statistics and analysis.

d. Output results: identify each group and output the path and statistics that form a line branch closest to an ideal line branch

# Example:



Line Error: 1
Node 2 branches to Node 3

LLRILL: $4/(2*\log(4+1) + 5) = 4/(2*\log(5)+5)$
4 Nodes in line
Ideal length of line = $a*\log(B+1)+c$
where $a = 2$, $B = 4$, $c = 5$

Branch Continuation = 0

Branch Positions: 2/4
Node 2 branches and is second node in line

Value of Line Branch:
$(1 - (1/(2*\log(5)+5+4))) * 0.4 + (4/(2*\log(5)+5)) * 0.3 + (1 - (0/4)) * 0.2 + (1 - ((4-2)/4)/4) * 0.1 = 0.83659121245$
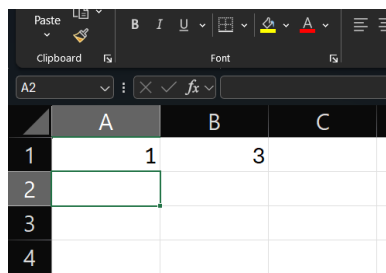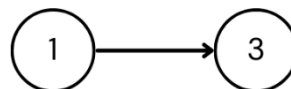
## 3.2 Deployment Description

This section outlines how to run the program, including organizing the CSV file, understanding the libraries used, and executing the program.

The program requires a CSV file as input to represent the network. The file must follow a specific structure for the data to be properly processed and analyzed. The first column should contain the source nodes, representing the nodes that initiate connections (nodes that point towards other nodes). The second column should contain the target nodes, representing the nodes that receive connections (nodes that are pointed towards). Therefore, each row in the CSV file corresponds to a directed edge in the network, where connections point from the source nodes to target nodes. For example, a row with 1 in the first column and 3 in the second column represents node 1 pointing to node 3. This format for the CSV file must be maintained to ensure no errors occur during the program's execution.



The program utilizes several R libraries to carry out data manipulation, graph analysis, and visualization
Tidyverse: R packages that assist with data wrangling, parsing, reading, and writing.
igraph: Allows for graph visualisation and graph analysis
visNetwork: Allows for interactive graph visualization as well as dynamic network diagrams
MASS: R packages that include statistical tools and functions used to analyze data

Steps to run the program (R Markdown File)
1. Prepare the CSV file and specify file path in the program
2. Install required libraries
3. Run each code chunk to view analysis and results

## 3.3 Information for Next Development Team

This project was done individually and, while it explores a similar topic to that of another team (the project done on Information Spread Models on X), it branches away significantly in its approach and focus. The emphasis here was on visually representing a network and identifying subgraphs within a network.

For more information about this project or any inquiries, contact Colin Song

Email: colin.song@mail.mcgill.ca

# IV. Experiment or Test Cases

## Describe Experiment and Results

The primary focus of this project was to visually represent the network and identify subgraphs within a network. As a result, extensive testing was not a central component of this project. However, deeper analysis was conducted on cloistered graphs, tree graphs, and line branch graphs. After defining these subgraphs, I evaluated their characteristics to understand the differences in their structures and how these subgraphs influence information flow.

This analysis involved studying the structural properties of each subgraph, using metrics such as connectivity and node distribution. For cloistered graphs, I analyzed the density of internal connections and the role of gatekeeper nodes. Analysis revealed that cloistered graphs often create self-contained cycles of information with limited spread beyond the group. For tree graphs, I examined properties such as height and branching patterns. The analysis showed that the branching factor and depth of the tree significantly influences the spread of information, with bottlenecking at critical nodes slowing the spread of further information. For line branch graphs, I focused on the length of the line and the branching behavior. I found that the linear structure of a line branch allowed for predictable spread of information, however, the line often ended shortly and soon branched into many nodes. The line branch is a structure where the source of the information is less influential than the final user, who interacts with and impacts many of their followers. Although these analyses provide valuable information into the nature of the subgraphs, more testing must be done to further understand the impact of these subgraphs and improve the program's capabilities.

# VI. Reflection

## 6.1 What would you do differently (what have you learned negative)

My primary goals for this project were to develop a program to visually represent a network and implement three algorithms to identify three subgraphs: cloistered graphs, tree graphs, and line branch graphs. While I was successful in achieving these goals, I ended up coding a program of 900 lines. I often found myself getting lost in

the complexity of the code as well as the sheer number of variables used in the program. This made the code harder to follow and required me to rewrite certain sections of the program multiple times to fix issues or improve the clarity.

Looking back, I would focus on better organization for not only the coding process but also for the planning process of implementing the three algorithms. Although I had clear definitions and steps for each algorithm, I did not effectively plan out how to translate those steps into code. This created numerous challenges during the development of the program, as I often found myself unsure of how to proceed with the next step. I would have also dedicated more time to designing the structure of the code and planning how to implement the algorithms programmatically. This approach would have relieved much of the stress I experienced from being overwhelmed by the complexity of the code and ensured a smoother implementation process.

## 6.2 What would you do again (what have you learned positive)

Looking back, one decision I would make again is choosing to code this project in R. R proved to be an excellent choice due to many tools for statistical analysis, data visualization, and graph modeling the software offers. Libraries such as *igraph* and *tidyverse* allowed for constructing and analyzing the network, while *visNetwork* allowed for effective and interactive visualizations. These tools streamlined many aspects of the project and let me focus more on my primary goals.

Another decision I would repeat again is the decision to create my own definitions for the subgraphs—cloistered graphs, tree graphs, and line branch graphs—and to develop properties for analyzing these structures. The definitions and properties I came up with proved to be the foundation of this project and provided a unique perspective for evaluating groups of nodes within a network. By breaking down the steps to analyze each subgraph, I was able to have a better understanding of not only the influences subgraphs have on the network but how information flows. The process of defining these subgraphs provided clarity and a consistent method of analyzing communities.

Finally, the algorithms I designed to identify and classify subgraphs within a network is a decision I would make again. By breaking down the steps systematically, I developed a procedure to identify any community as a specific subgraph, given it matches the properties of the ideal subgraph it is being identified as. Although the algorithms can be improved, they currently succeed in achieving the goals I set out and demonstrated to be very effective when analyzing large complex network structures. These decisions not only contributed to the success of the project but also provides a strong foundation for future development.

## Future Work

## Things to Fix

One thing to fix involves tuning the parameters used in the program. One example of this is when calculating the ideal line length for a line branch graph. The values of the scaling factor $a$ and the baseline length $c$ in the formula $L = a \cdot log(B + 1) + c$ require further tuning to align with real-world data. Testing numerous datasets would help determine optimal values for these parameters, ensuring the formula better represents realistic line structures. This adjustment would improve the accuracy of the line branch graph analysis.

One major challenge I faced during my research was graphing an entire network. Early on, I attempted to graph a dataset with one million connections, which took the program an hour to complete. However, the resulting graph was unreadable as such a dense network appeared as a blob of colors. Networks with many nodes quickly became visually overwhelming, making it hard to distinguish meaningful patterns or relationships. The only fix I had around this issue was to focus on smaller areas of the network to make graphing and analyzing more manageable.

Lastly, the need for new data remains an unresolved issue. The current project heavily relies on existing datasets, which may limit the scope of analysis and the ability to simulate diverse scenarios. Developing a method to generate or acquire new data—such as creating datasets, simulating user interactions, or collaborating with data providers—would expand the project's capabilities and allow for more robust testing of the algorithms. Access to new and varied data would also enable the program to adapt to evolving network dynamics and remain relevant for future applications.

## Features to Add/Future Work

**1. Testing and Tuning Parameters**

Future development should involve more extensive testing to fine-tune the formula for the ideal line length. The current formula, $L = a \cdot log(B + 1) + c$, uses parameters $a$ (scaling factor) and $c$ (baseline length). I currently have $a$ set to be 2 and $c$ set to be 5, however these values should be adjusted to better reflect real-world network properties.

Weights used when analyzing tree and line branch graphs should also be adjusted. These weights determine whether a group is fit to be a tree or a line branch graph. Due to the significant role these weights play in identifying subgraphs within a network, they must be tested rigorously against large datasets to ensure accurate identification

Experimenting with these parameters can help refine the identification process to ensure it aligns more closely with observed patterns in social media networks. Testing various datasets and edge cases will provide valuable insights into how well the program adapts to different types of networks.

**2. Using the Program to Evaluate Information Spread**

The program has the potential to simulate and evaluate how information flows in networks. Future teams can build on my program to analyze the flow of information through specific subgraphs, focusing on factors such as:

- Speed of Propagation: How quickly information reaches distant nodes.
- Influence of Key Nodes: Identifying nodes that play pivotal roles in spreading information.
- Impact of Communities: How tightly-knit groups influence or impede the flow of information.

These factors can be extended to simulate real-world scenarios, such as tracking viral tweets or examining how misinformation spreads. By using the program to analyze real data, the algorithms could be tested in more practical scenarios and reveal how information flows when subgraphs are identified. This could include studying how tightly connected cloistered graphs gatekeep information, how hierarchical tree structures spread information, and how line branch graphs influence the reach and longevity of information. Analyzing real data would uncover patterns in information spread that are unique to real-world networks.

**3. Adding More Properties to Subgraphs**

To further solidify the definitions of each subgraph, additional properties and metrics can be incorporated into the analysis. Examples include:

- Clustering Coefficient: Measure how tightly nodes are connected within a subgraph.
- Edge Weight Analysis: Include weights on edges to represent the strength of relationships or interaction frequencies.
- Centrality Metrics: Analyze measures such as betweenness or eigenvector centrality to better understand the roles of nodes within subgraphs.

Improving the subgraph properties will provide a more comprehensive understanding of their structure and behavior, enabling more accurate identification and deeper analysis. While the current metrics provide a strong foundation, adding more properties would help solidify the definitions of cloistered graphs, tree graphs, and line branch graphs. Incorporating metrics could provide deeper insights into the behavior of each subgraph. This would not only improve the identification process but also offer a better understanding of how these subgraphs function within larger networks.

## Suggestions

1. Plan ahead: Before starting, organize the steps required to complete the project and outline each task
2. Research: Take the time to read articles to become familiar with the topics involved in the project
3. Focus on your goals: Always relate your work back to the project's main objective

# VII. Supplemental Materials

During the course of this project, several materials and resources were studied that were ultimately not used for the final implementation of my research. However, this information could be valuable for future work. These include unused statistical metrics and hydrating tweets using Twarc.

## Unused Statistical Metrics

In the process of analyzing subgraphs, I gathered information about several additional statistical metrics that were not incorporated into the project. These metrics could provide deeper insights into network behavior and help refine subgraph analysis.

1. Modularity
   - Compares edges inside communities to edges between communities. Quantifies how well a graph can be divided into separate communities

2. Number of Communities
   - Total number of distinct communities in network. Large number of communities may indicate highly fragmented network. Small may indicate more centralized structure

3. Distribution of Community Size
   - Distribution of sizes of detected communities. Shows presence of large communities vs. small

4. Average # of Edges in Community
   - Average number of edges in communities. Shows how well-connected per node in community

5. Number of In-Degrees and Out-Degrees
   - Number of edges coming to a vertex and number of edges coming out of vertex. Can find users who are at central areas within community

6. Conductance
   - Ratio of in-degrees to total number of edges. Lower conductance indicates more tightly-knit, well-defined community, higher indicates community has more connections outside

7. Conductance of Whole Network
   - Sum of conductance values of all communities. Measures how well network is divided into communities. Lower value suggests well-defined community structures

8. Edge Density

  – Ratio of number of edges to max number of possible edges. An idea of how dense a graph is in terms of edge connectivity

9. Clustering Coefficient

  – The degree to which nodes in a community tend to cluster together. High coefficient suggests that nodes form tight-knit groups

10. Eigenvector Centrality

  - Finds centrality for a node, proportional to other nodes. Nodes that are pointed to by "important" nodes will also be considered "important"

11. Betweenness Centrality

  – Measure of centrality based on shortest path. Node with vlaue will have more information passed through

12. Intra-Community/Inter-Community Interactions

  – Ratio of edges within communities to edges between communities. High ratios suggest strong community boundaries, low ratios suggest overlap/interaction between communities

13. Overlapping Communities

  – Identifies nodes that belong to multiple communities. Reveals users that participate in various groups

# Hydrate Tweets

Another supplemental material involves the use of *Twarc*, a Python library for working with Twitter data. *Twarc* allows for hydrating tweet IDs, which retrieves the full metadata of tweets from their IDs (this requires an X Developer Account and access to the X API). While this approach was not utilized in this project due to limited API call rates, I invested considerable time exploring it and believed it had the potential to provide valuable datasets.

# References and Background Material (accumulated)

*Communities.* Communities - NetworkX 3.4.2 documentation. (n.d.).
    https://networkx.org/documentation/stable/reference/algorithms/community.html

Datta, S. S. (2024, March 18). *Graph density*. Baeldung on Computer Science.
    https://www.baeldung.com/cs/graph-
    density#:~:text=Graph%20density%20represents%20the%20ratio,new%20edges%20to%20the%20network

Du, D. (n.d.). *Social Network Analysis: Lecture 3-network characteristics*. Social Network Analysis.
    https://ddu.ext.unb.ca/6634/Lecture_notes/Lec3_network_statistics_handout.pdf

*Eigenvector_centrality*. eigenvector_centrality - NetworkX 3.4.2 documentation. (n.d.).
    https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.centrality.eig
    envector_centrality.html

*Finding community structure by multi-level optimization of modularity*. R igraph manual pages. (n.d.-a).
    https://igraph.org/r/doc/cluster_louvain.html

GeeksforGeeks. (2022, July 21). *Betweenness centrality (centrality measure)*.
    https://www.geeksforgeeks.org/betweenness-centrality-centrality-measure/

GeeksforGeeks. (2024, March 20). *Difference between graph and tree*. https://www.geeksforgeeks.org/difference-between-graph-and-tree/

Junker, N. (2020, March 4). *Community detection with Louvain and Infomap: R-bloggers*. R-bloggers. https://www.r-bloggers.com/2020/03/community-detection-with-louvain-and-infomap/

Junker, N. (2022, August 17). *Interactive network visualization with R*. statworx. https://www.statworx.com/en/content-hub/blog/interactive-network-visualization-with-r/

Lizardo, O., & Jilbert, I. (2020, January 6). *2.7 average degree: Social Networks: An introduction*. Social Networks: An Introduction. https://bookdown.org/omarlizardo/_main/2-7-average-degree.html

Louf, B., McDiarmid, C., & Skerman, F. (n.d.). *Modularity and graph expansion - drops*. Modularity and Graph Expansion. https://drops.dagstuhl.de/storage/00lipics/lipics-vol287-itcs2024/LIPIcs.ITCS.2024.78/LIPIcs.ITCS.2024.78.pdf

*Modularity*. modularity - NetworkX 3.4.2 documentation. (n.d.). https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.community.quality.modularity.html#r6937dc4d2017-4

Newman, M. (n.d.). *Modularity and community structure in networks*. Proceedings of the National Academy of Sciences. https://www.pnas.org/doi/full/10.1073/pnas.0601602103

Pisharody, A. (2022, March 8). *Learn TWARC!*. Twarc Command Basics · Learn Twarc! https://scholarslab.lib.virginia.edu/learn-twarc/06-twarc-command-basics

Popova, K. (2024, June 5). *Network analysis and community detection on political tweets*. Medium. https://medium.com/@pop.kristina1/network-analysis-and-community-detection-on-political-tweets-9e0f21294e31

R igraph manual pages. (n.d.-b). https://igraph.org/r/doc/communities.html

Terzi, E. (n.d.). *Graph clustering*. CAS CS 565, Data Mining. https://cs-people.bu.edu/evimaria/cs565-16/lect10.pdf

Wolitzky, A. (n.d.). *Lecture 3: Eigenvector centrality measures*. MIT 14.15 Networks. https://ocw.mit.edu/courses/14-15-networks-spring-2022/mit14_15s22_lec3.pdf

## Databases:

*Higgs twitter dataset*. SNAP. (n.d.-a). https://snap.stanford.edu/data/higgs-twitter.html

*Social circles: Twitter*. SNAP. (n.d.-b). https://snap.stanford.edu/data/ego-Twitter.html

*Stanford Large Network Dataset Collection*. SNAP. (n.d.-c). https://snap.stanford.edu/data/index.html#twitter

*Twitter Interaction Network for the US congress*. SNAP. (n.d.-d). https://snap.stanford.edu/data/congress-twitter.html

## Databases Not Used (not APA):

https://github.com/shaypal5/awesome-twitter-data?tab=readme-ov-file#twitter-datasets

https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi%3A10.7910%2FDVN%2FPDI7IN

https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LW0BTB

https://data.world/socialmediadata/j-k-rowling-tweets-and-retweets/workspace/project-summary?agentid=socialmediadata&datasetid=j-k-rowling-tweets-and-retweets

https://github.com/echen102/COVID-19-TweetIDs