

SOCIAL NETWORK ANALYSIS

Github Social Network



INTRODUCTION

A large social network of GitHub developers which was collected from the public API in June 2019. Nodes are developers who have starred at least 10 repositories and edges are mutual follower relationships between them. The vertex features are extracted based on the location, repositories starred, employer and e-mail address. The task related to the graph is binary node classification - one has to predict whether the GitHub user is a web or a machine learning developer.

Dataset [Link](#)

Gephi software was used for the visualization and analysis of this dataset.

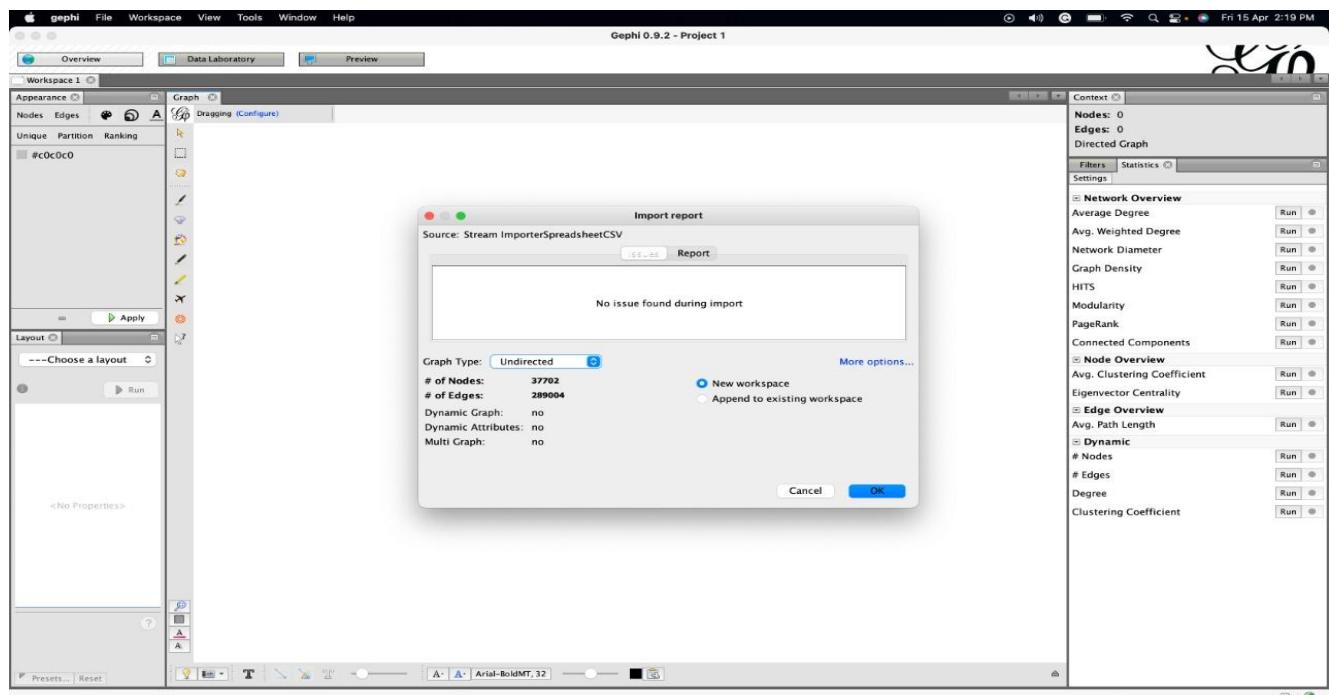
The following steps show the procedure followed for the analysis and interpretation.

Step 1:

The .gml file for the initial data cluster was imported into Gephi. The number of nodes and edges were mentioned in the following window.

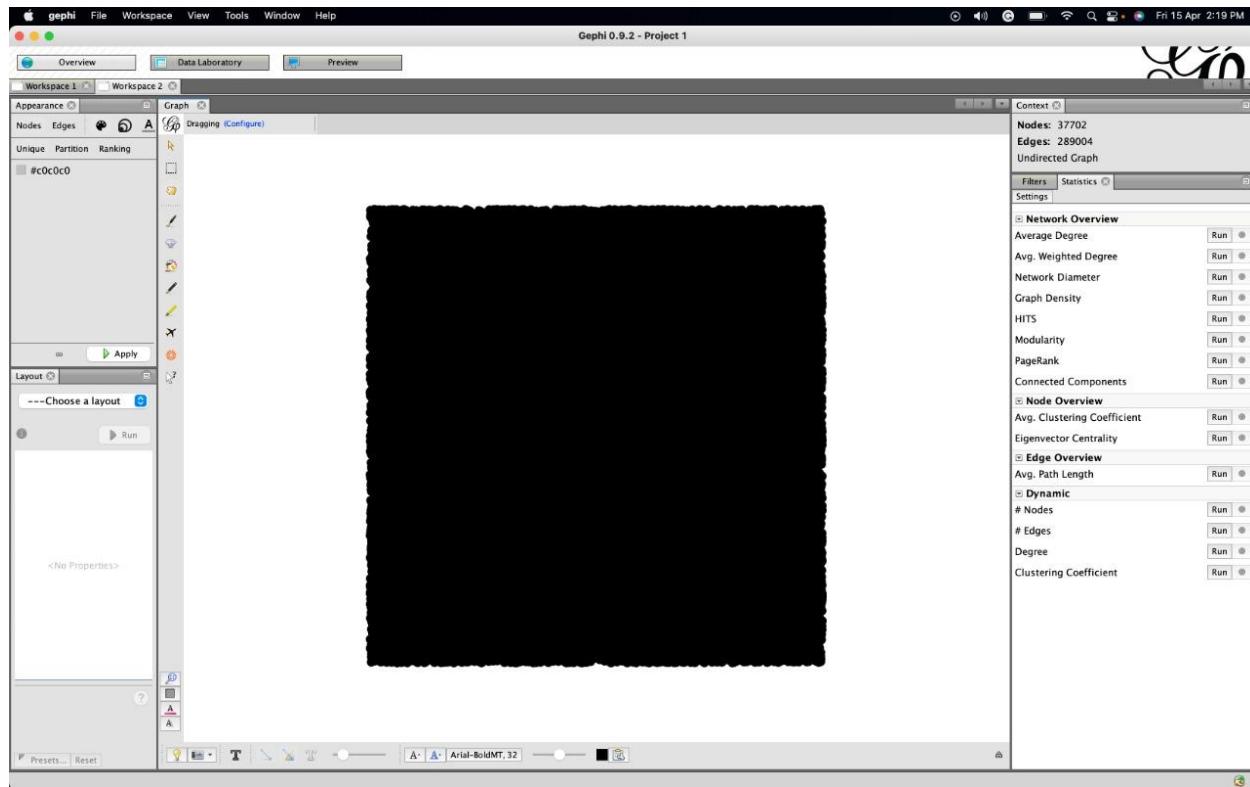
Number of Nodes- 37702

Number of Edges - 289004



Step 2:

Initially, the data cluster looked like a square with multiple nodes and edges through which no interpretation and analysis could be done.



Step 3:

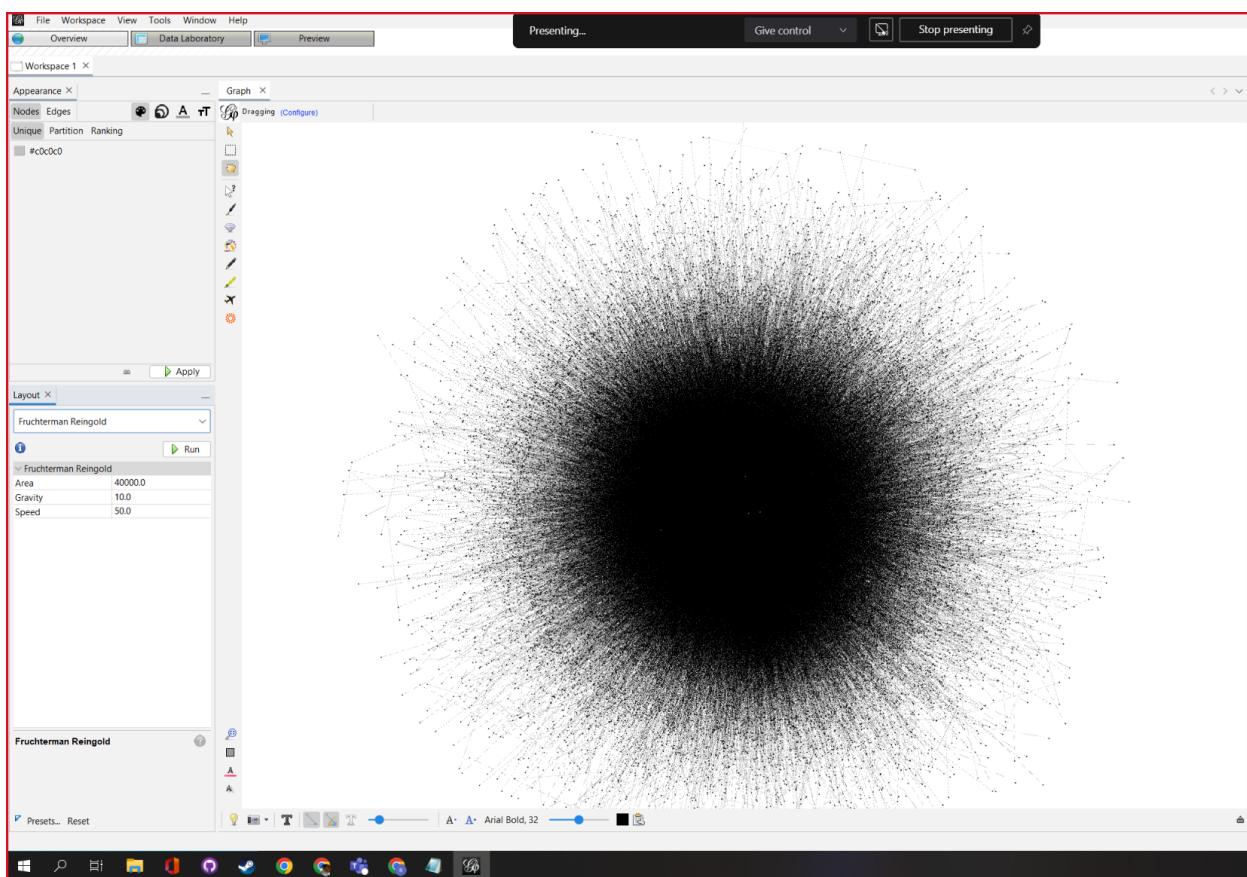
In this step, Fruchterman Reingold was selected in the Layout option. The parameter were set as-

Area- 40000

Speed- 50

Gravity- 10

The clusters rearranged themselves in a manner as shown in the picture below-



Step 4:

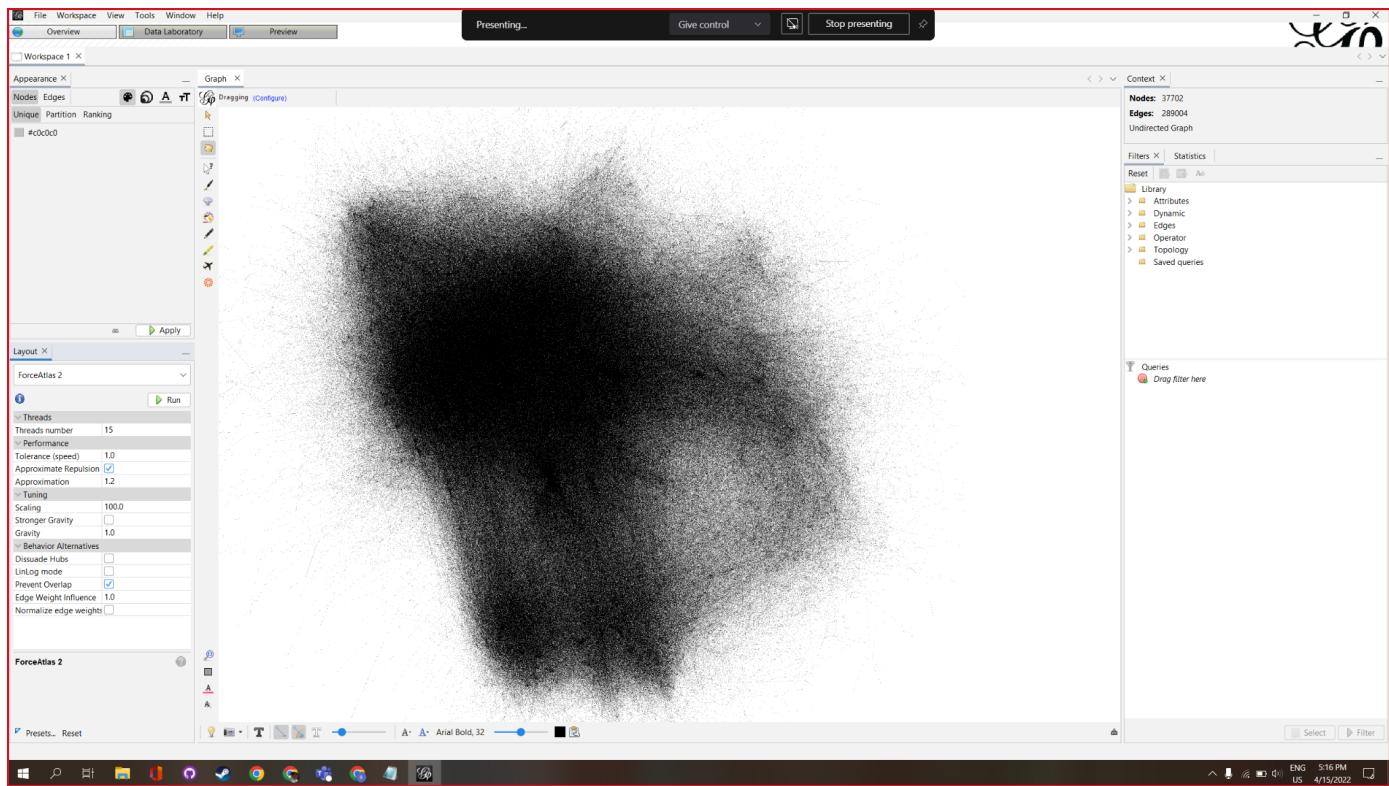
In this step, Force Atlas 2 was selected in the Layout option. The parameter were set as-

Tolerance Speed- 1.0

Scaling- 100

Prevent Overlap-

The clusters rearranged themselves in a manner as shown in the picture below- (run for 5 mins)



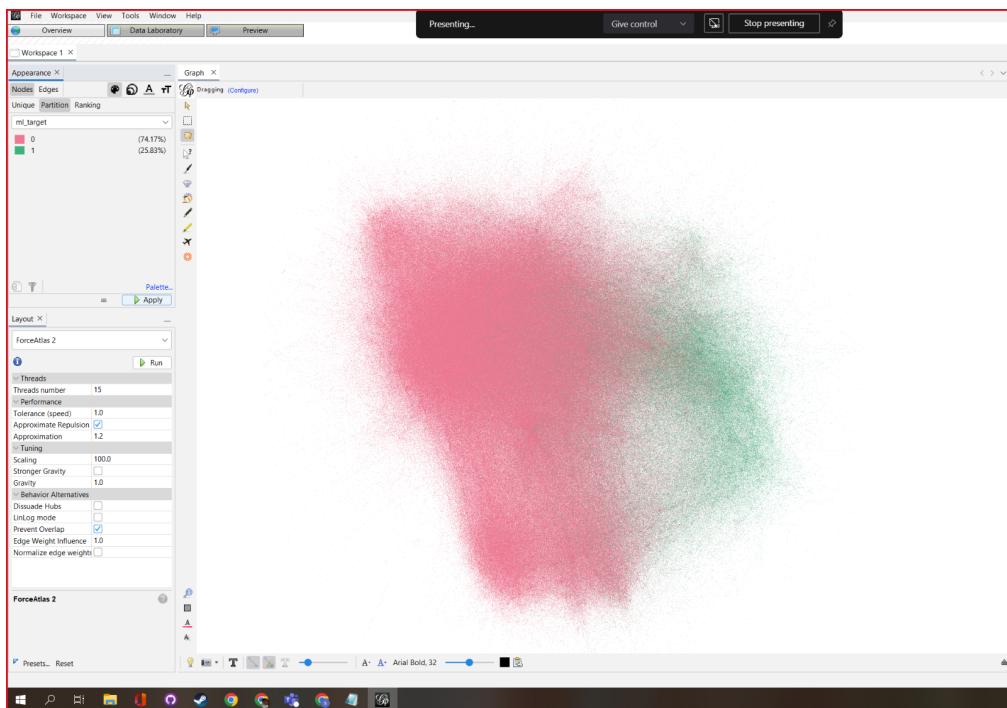
Step 5:

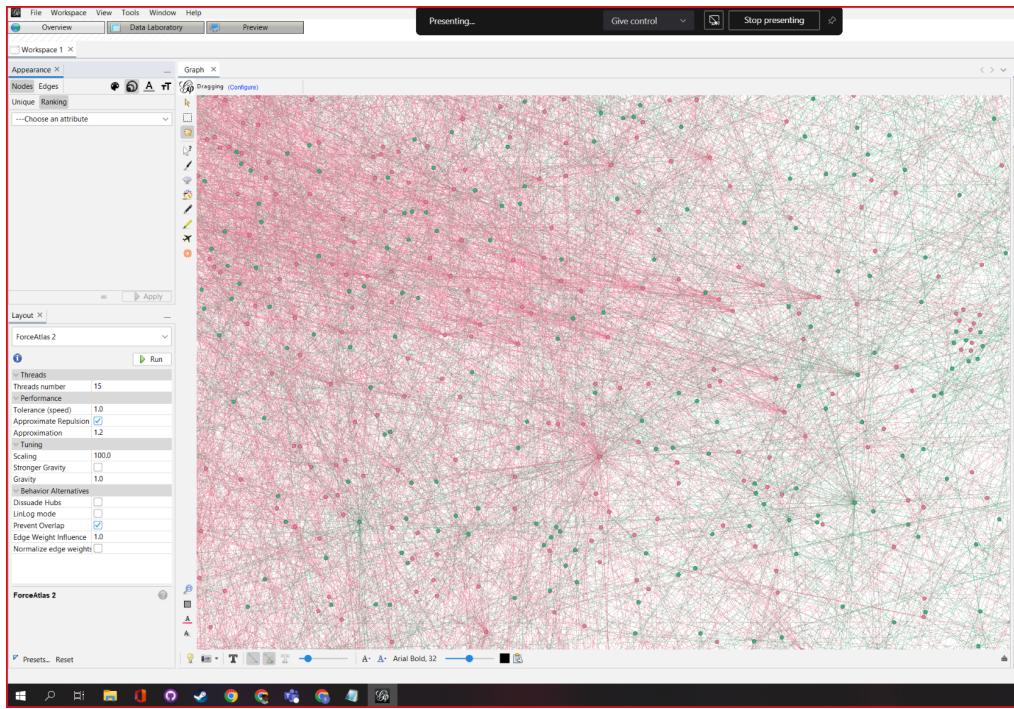
In this step, we tried to differentiate between the github users who are machine learning developers or not.

Nodes – Color Palette – Partition – “ml target” – Apply

Observation- **Pink Cluster**- Non- Machine Learning Developer

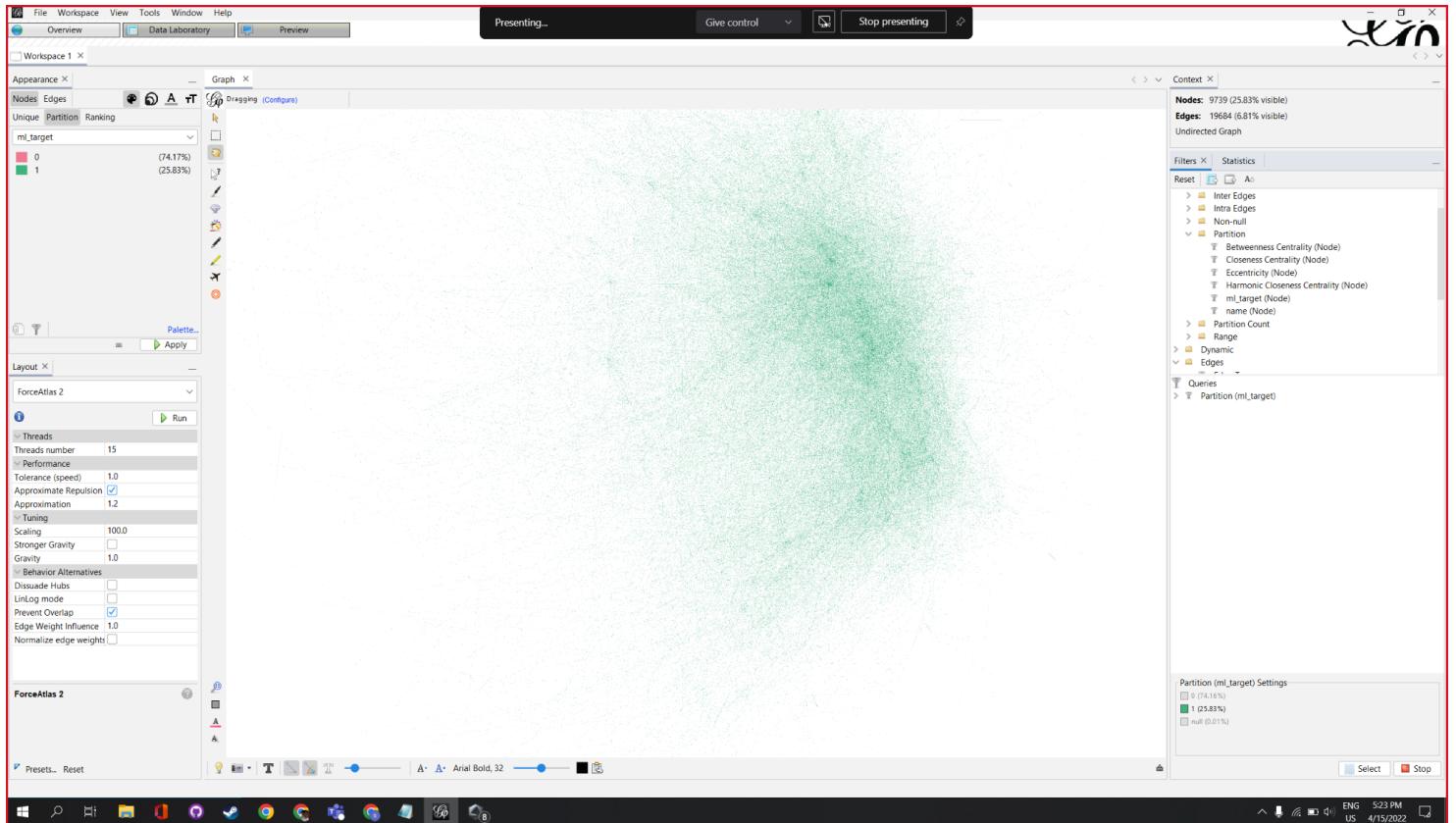
Green Cluster- Machine Learning Developer





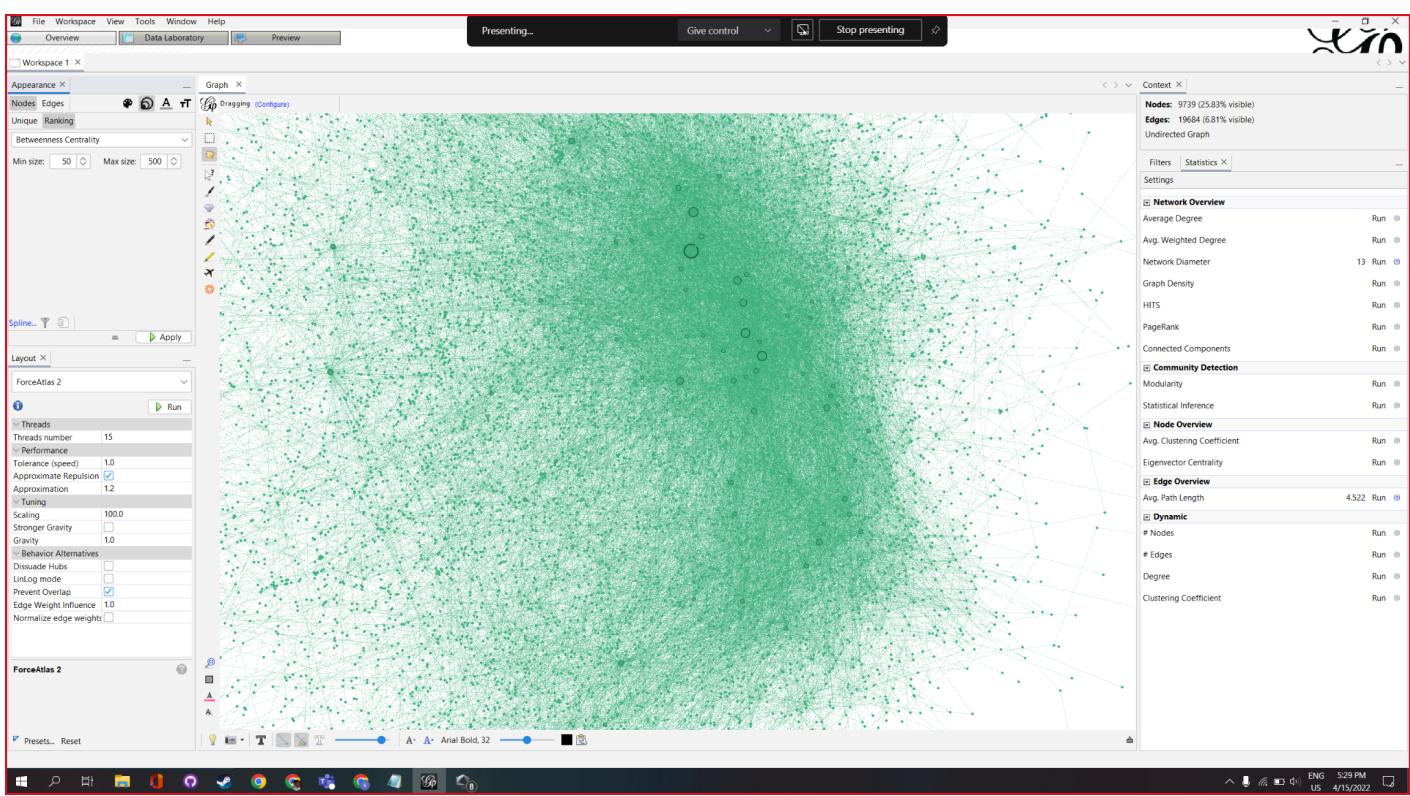
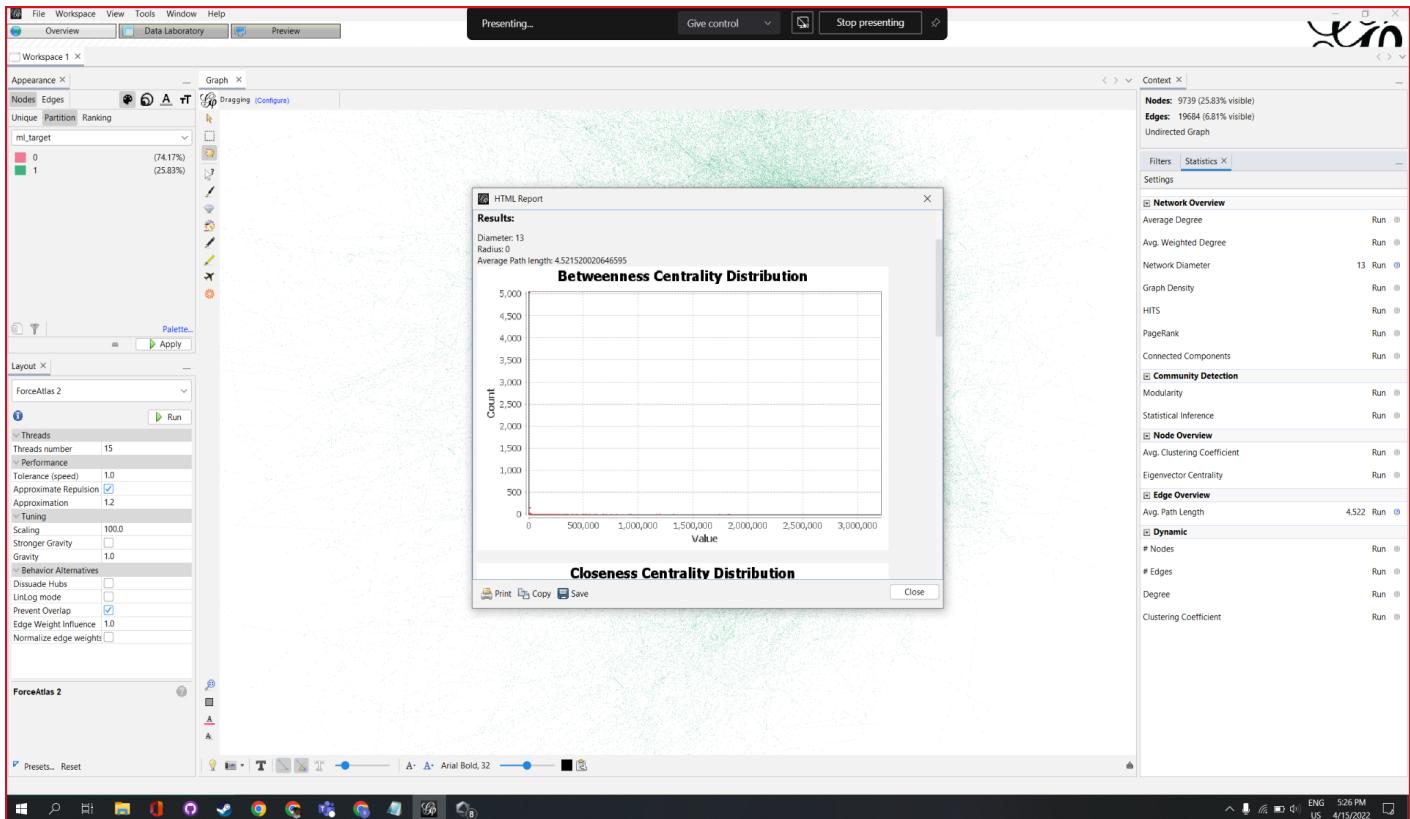
Step 6:

In this step, we applied a filter to only view the green cluster, i.e., the Machine Learning developers.



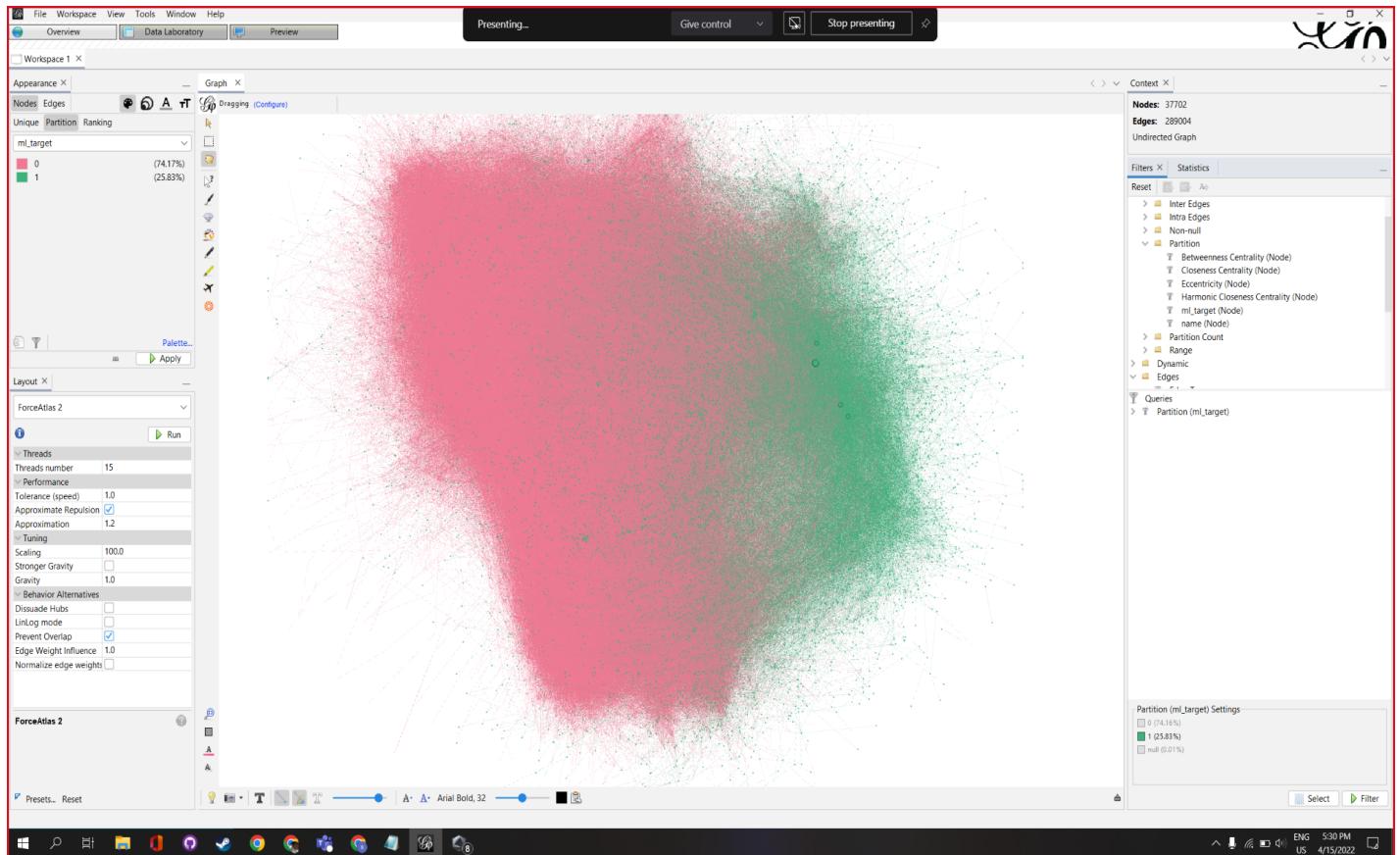
Step 7:

In this step, we calculated the centrality measure, and the centrality measure in size. Also, made the filter and then calculated



Step 8:

Removed filter on “ml_target” to view both, Machine Learning developers and Non-Machine Learning developers.



Step 9:

Modularity Report

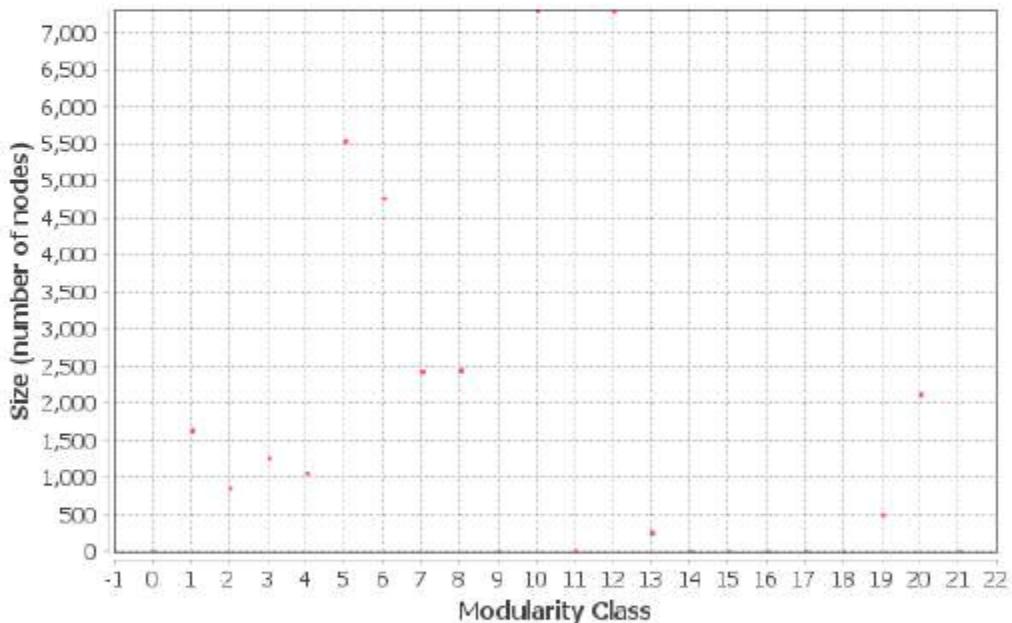
Parameters:

Randomize: On
Use edge weights: On
Resolution: 1.0

Results:

Modularity: 0.458
Modularity with resolution: 0.458
Number of Communities: 22

Size Distribution



Algorithm:

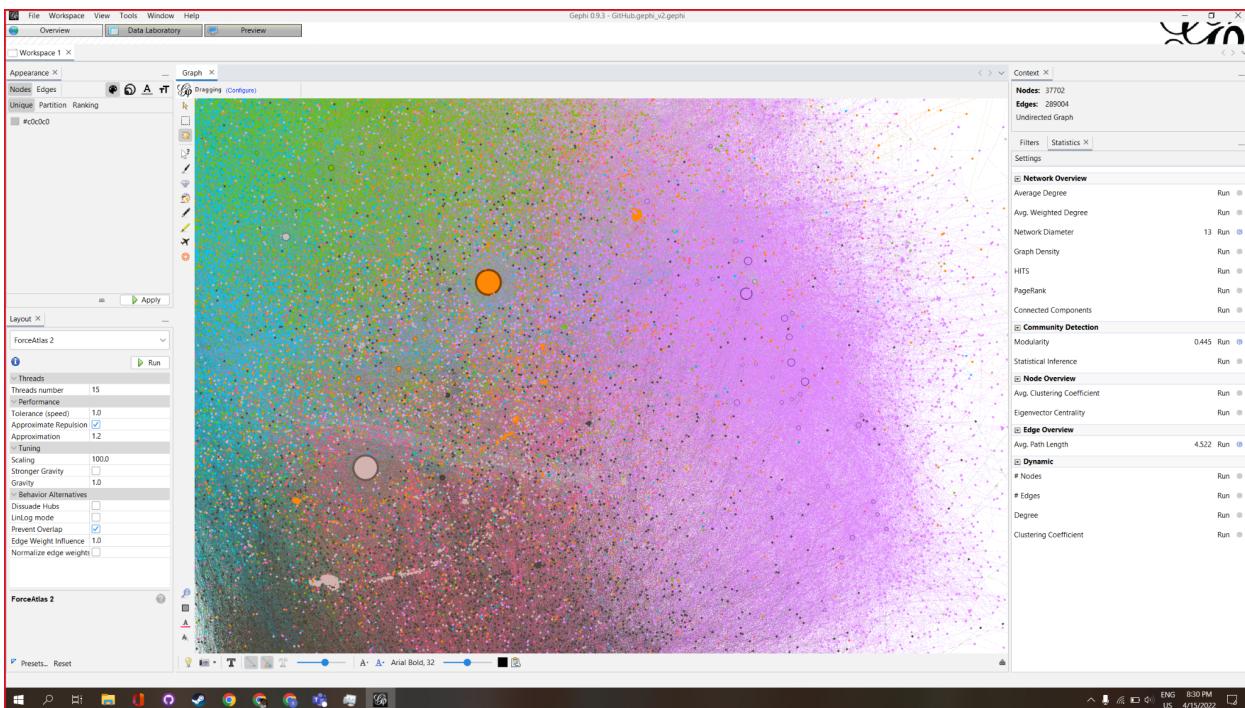
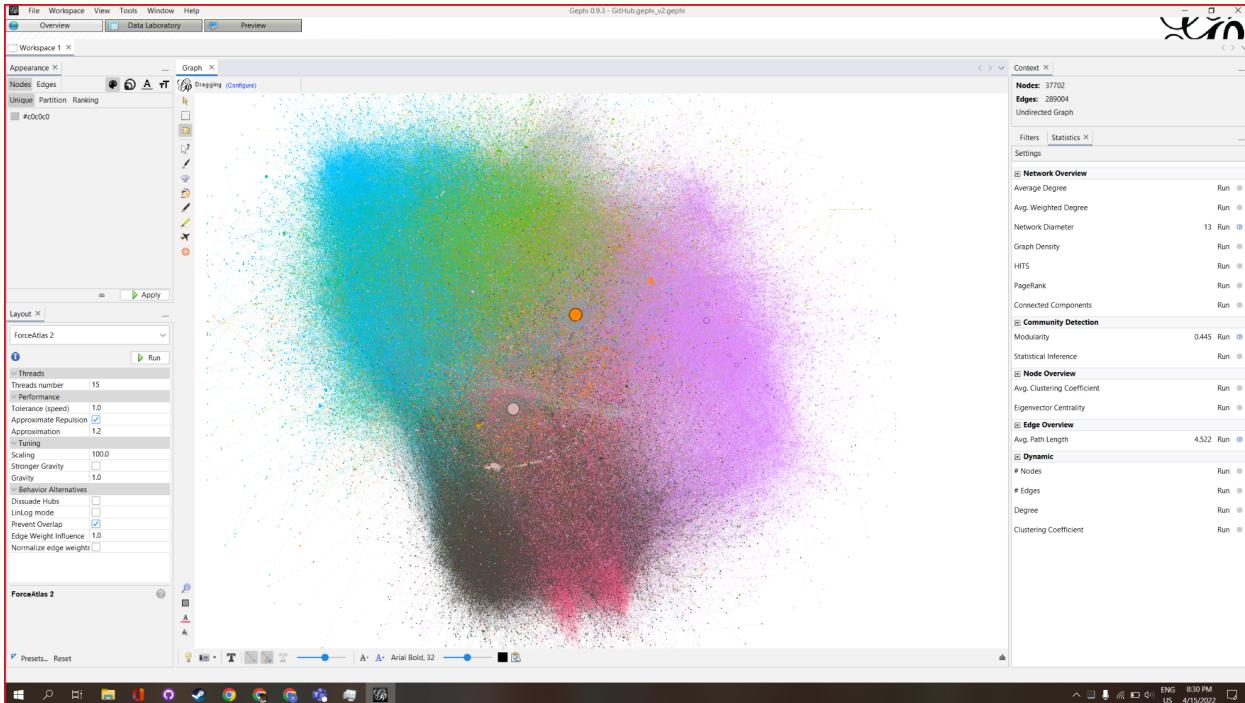
Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, *Fast unfolding of communities in large networks*, in Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P1000

Resolution:

R. Lambiotte, J.-C. Delvenne, M. Barahona *Laplacian Dynamics and Multiscale Modular Structure in Networks* 2009

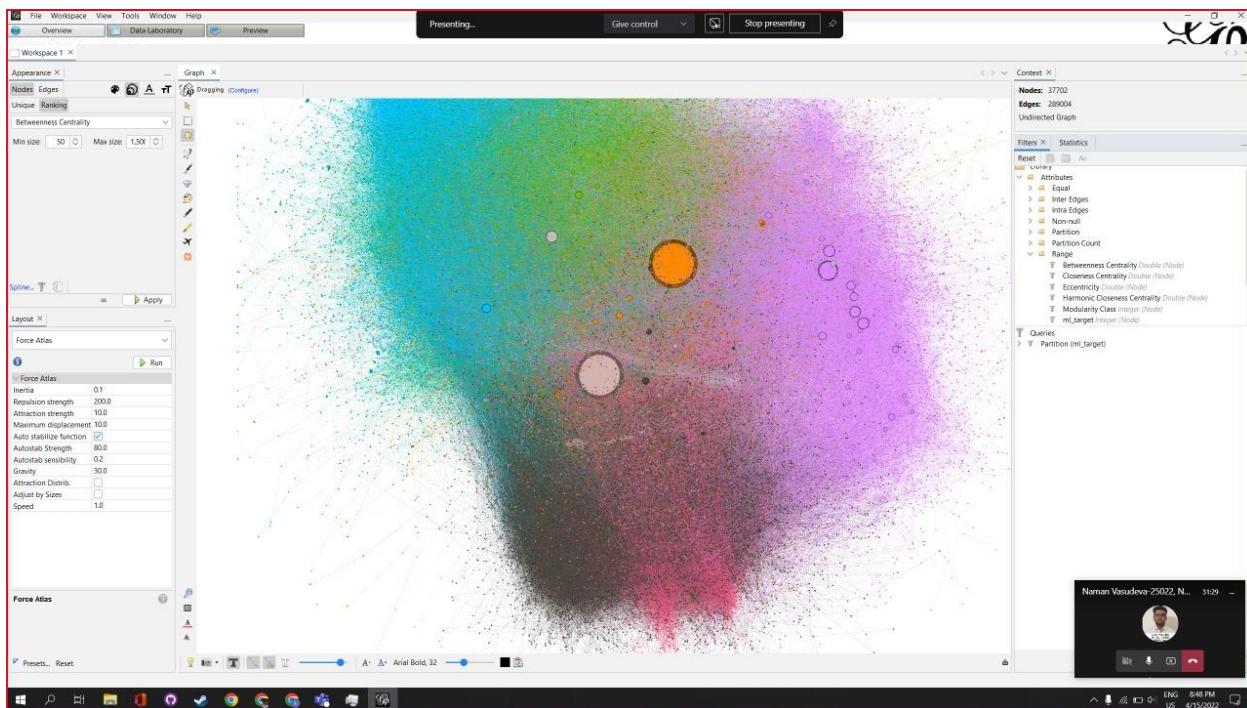
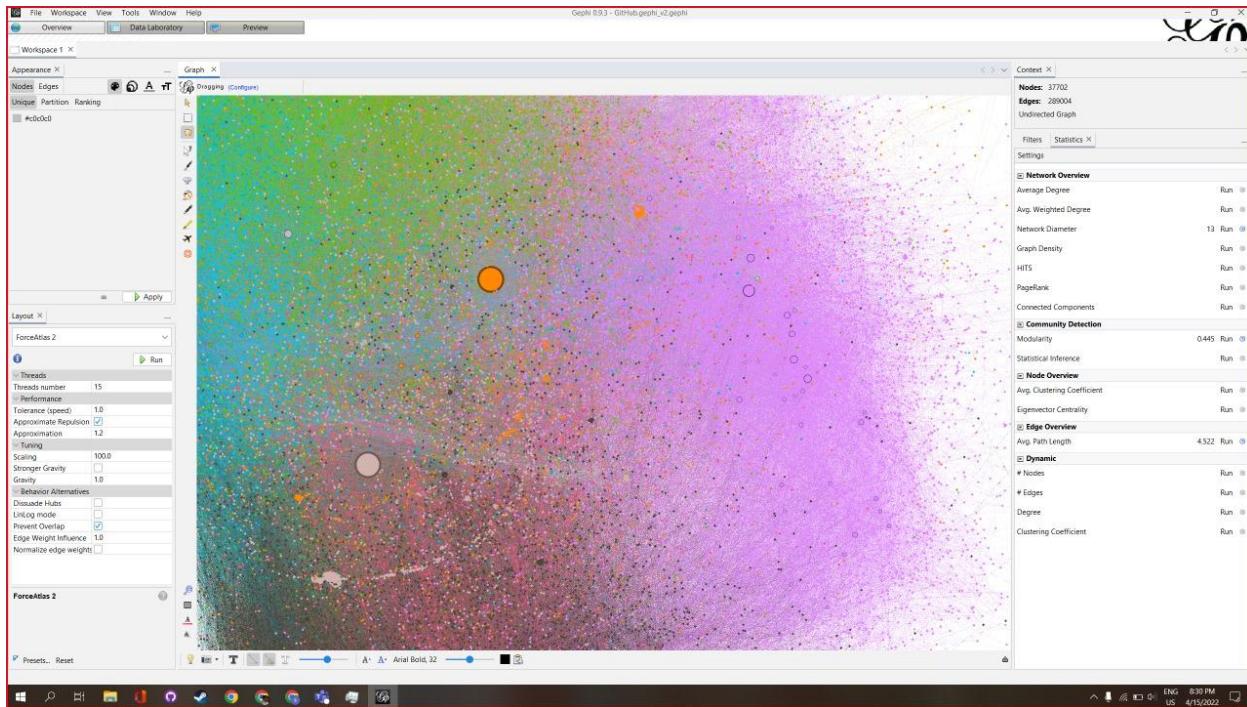
Step 10:

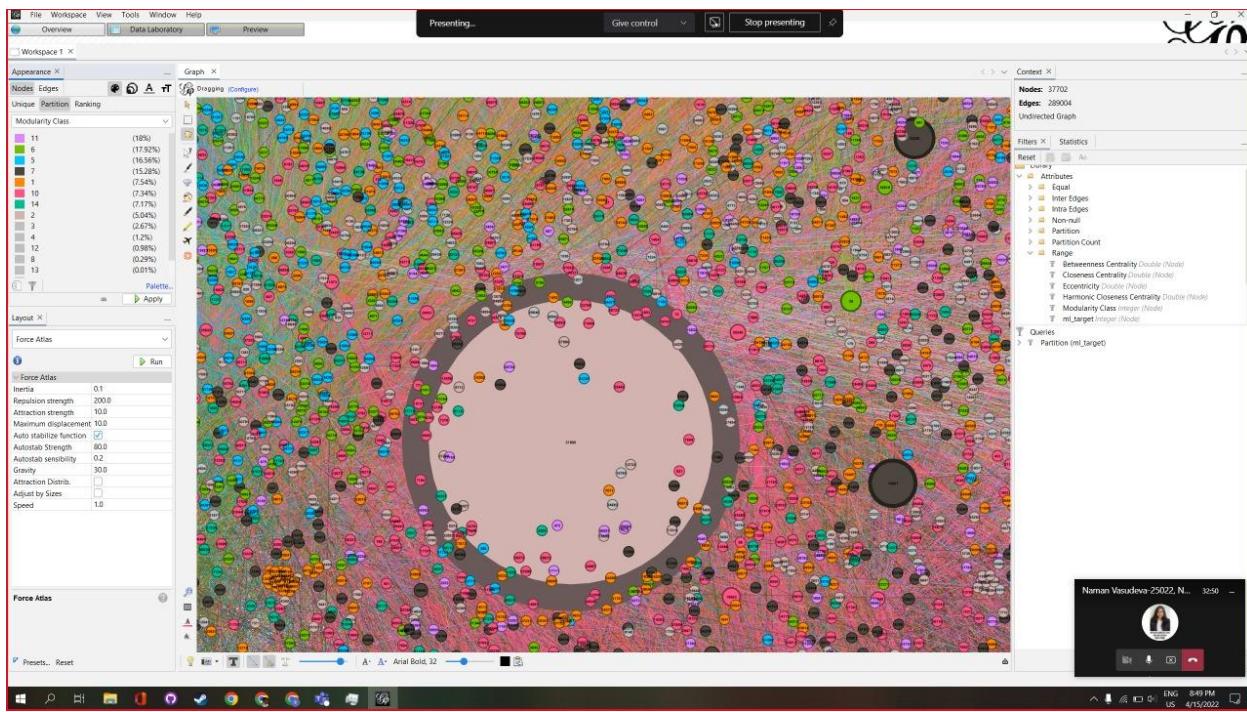
Colour as Modularity classes, size as betweenness of ML nodes



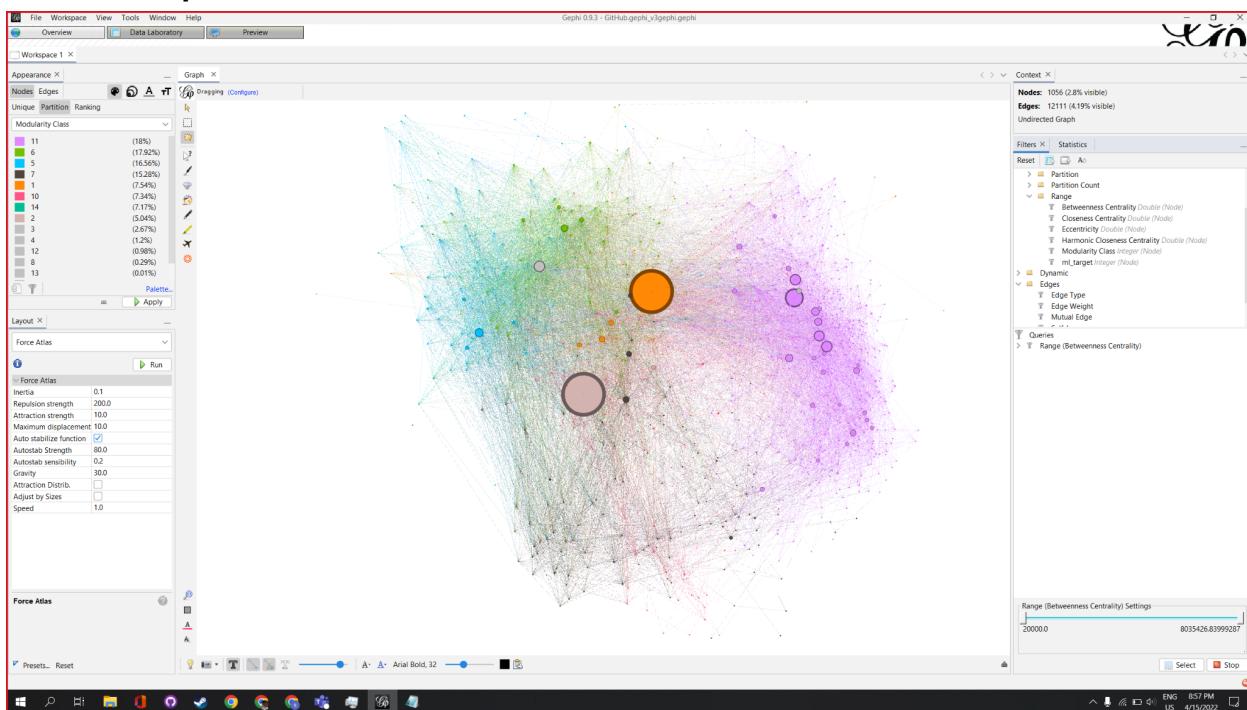
Step 11:

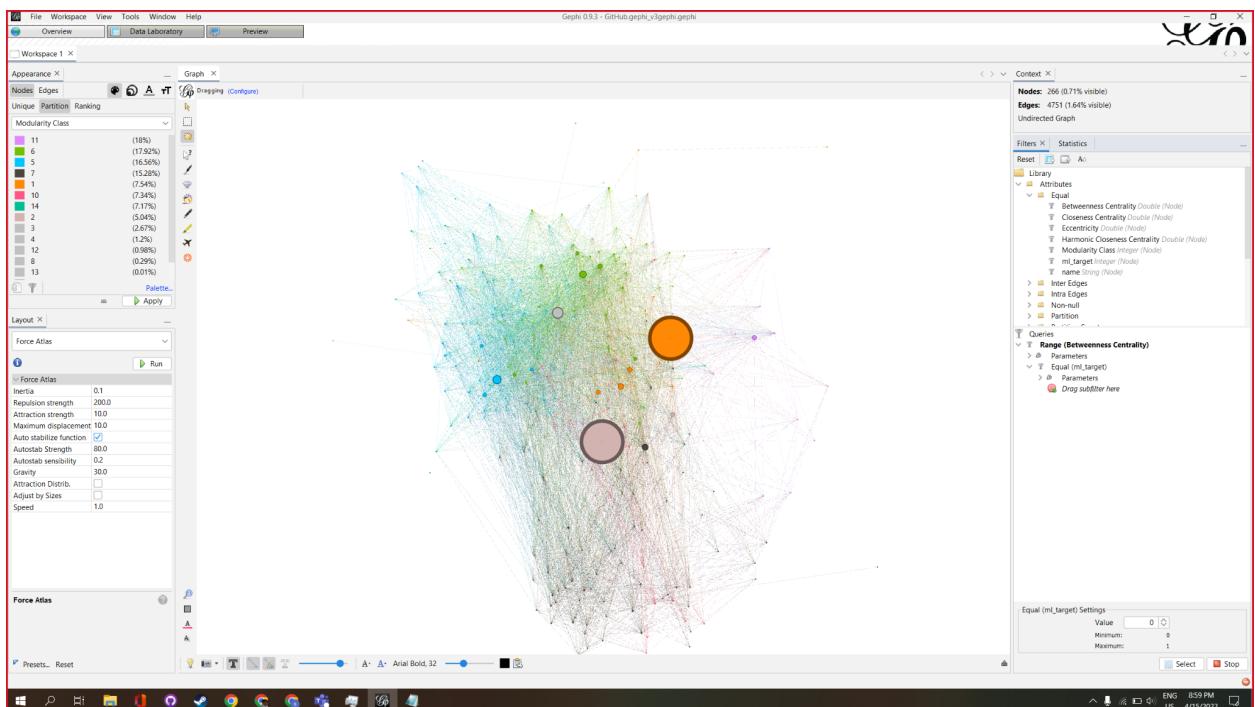
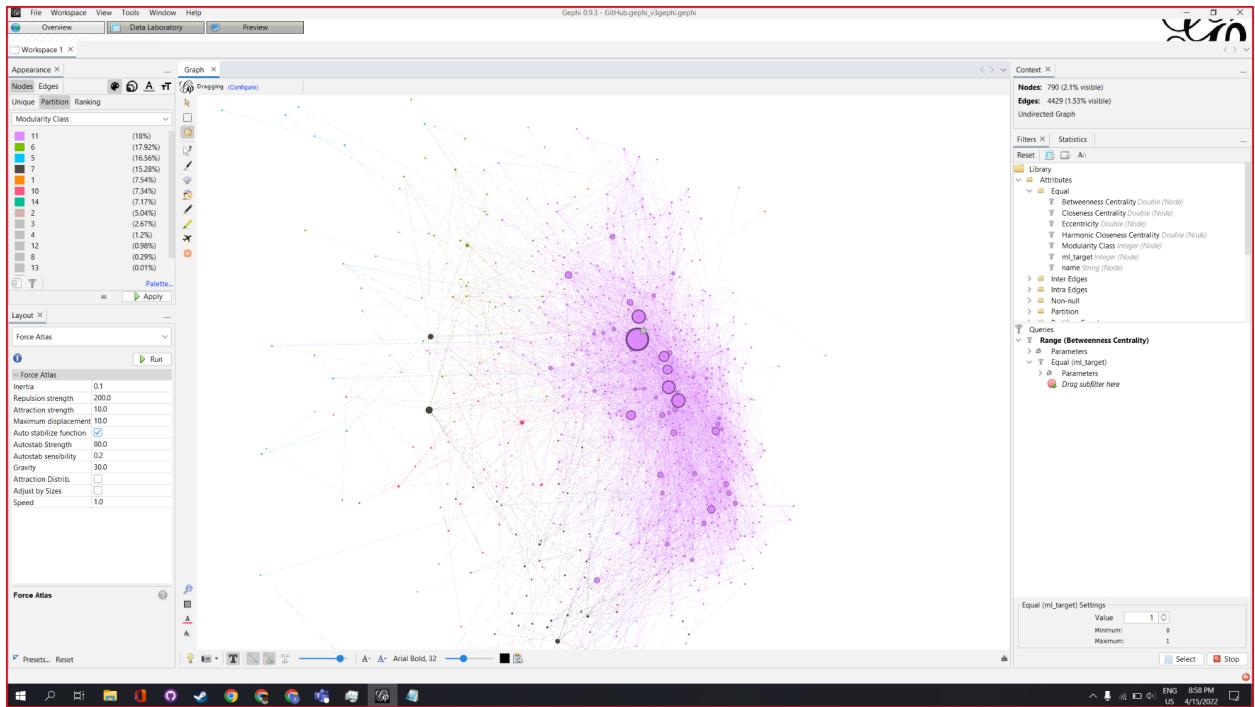
Calculated betweenness of the whole data set in this step.

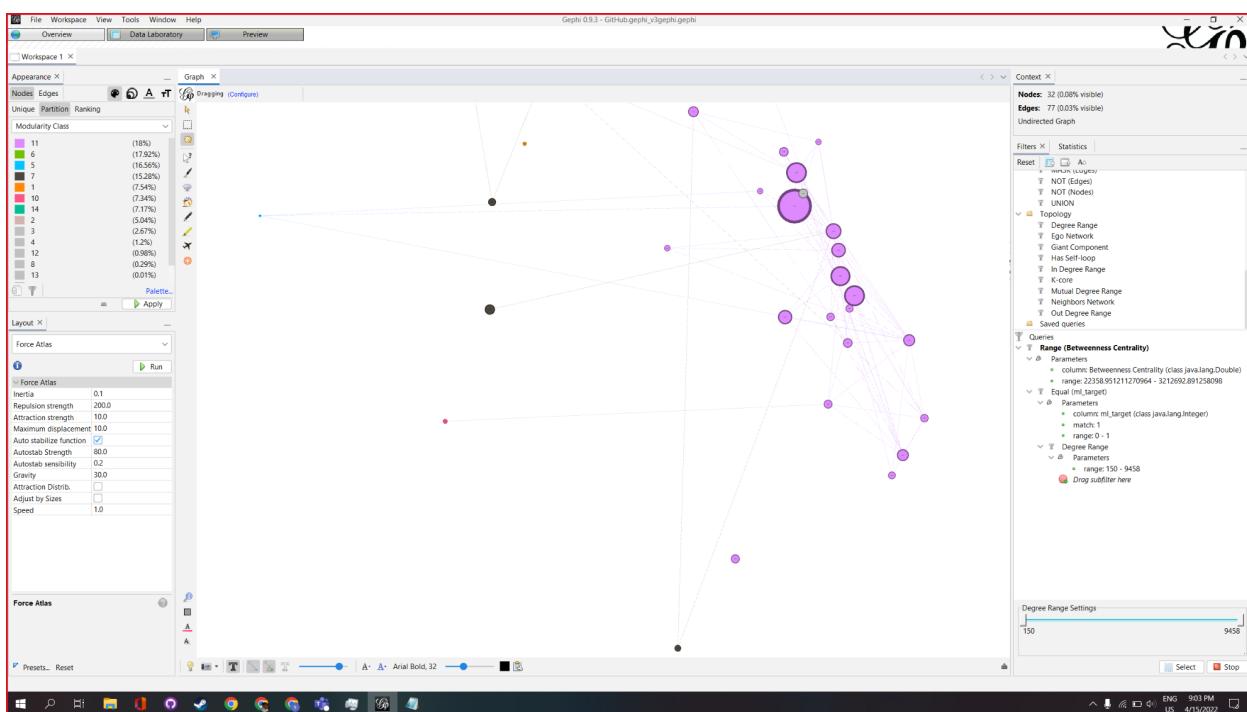
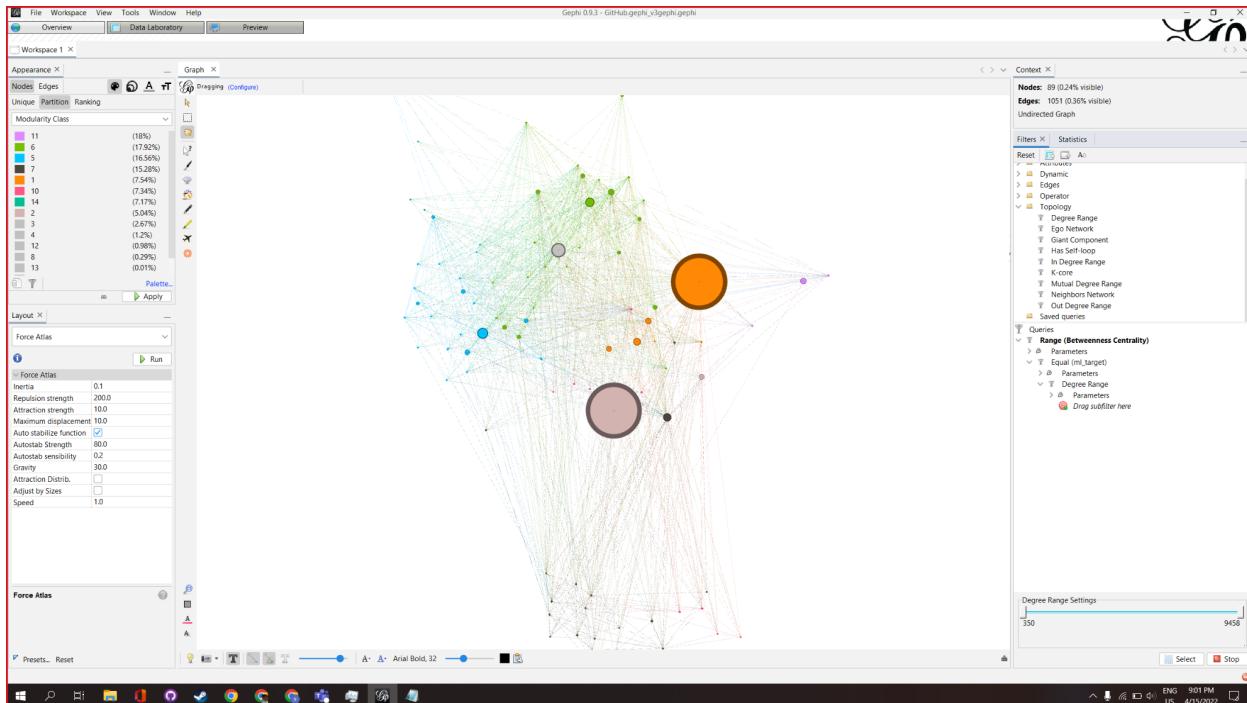




Filters to explore more







Questions

	QUESTIONS	ANSWERS
a.	Who are important entities from different points of view.	31890 , 27803 , 19222 , 14954, 20193 , 20528 , 29023 , 16631 Based on the highest degrees of nodes, the above mentioned entities can be regarded as important. Nodes 31890 and 27803 are of the most importance
b.	How many communities exist within the network? Examine the characteristics of each community. Why is a community different from other communities?	A total of 15 communities exist within this network. Communities 12,8,13 and 9 are least related to other communities. The community of machine learning developers who only do machine learning is community 11 Community 6 is possibly a community which consists of developers focused on machine learning and another topic which is common among them.
c.	Examine the relationship of nodes within and outside communities.	The dataset contains a list of all of the links, where a link represents mutual follow relationship between developers. Possibly the people using ML need help or ideas in ML coding for which they are following other developers. People who generally post ML codes are also being visited by other developers. This shows that the relevance of ML is increasing in other fields as well
d.	Any further insights that you may draw by analyzing the network	As the need for machine learning and data science has increased, it is consistent with the dataset that people from non coding background and non ML background are venturing into ML.