

GRAPH ANALYSIS AND ALGORITHMS:



NETWORK ANALYSIS CASE STUDY: **ANALYZING THE RETWEET NETWORK:**

COURSE INSTRUCTOR:

DR. SOMASUNDARAM K

DONE BY:

SMRITHI. N (CB.SC.I5DAS19029)

VARUN KUMAR.M (CB.SC.I5DAS19032)

PROBLEM STATEMENT OF THE PROJECT:

- TWITTER IS A SOCIAL MEDIA PLATFORM WHICH IS SIMILAR TO THE OTHER SOCIAL MEDIA PLATFORMS AND CONTAINS MILLIONS OF USERS (APPROX.). TWITTER HAS A FEATURE WHICH ALLOWS THE USERS TO RETWEET A POST THAT A PERSON POSTS IF THEY ARE ABLE TO RELATE TO IT OR THEY LIKE IT.
- OUR PROJECT FOCUSES MAINLY ON DERIVING INSIGHTS FROM THE NETWORK THAT IS FORMED BY TAKING THE VARIOUS USERS OF TWITTER AS THE NODES AND A PERSON RETWEETING ANOTHER PERSON'S POST AS A DIRECTED EDGE.

ABOUT OUR DATASET:

- OUR DATASET CONTAINS 761 NODES AND 1030 EDGES. IN OTHER WORDS, OUR DATASET TOOK 761 TWITTER USERS INTO ACCOUNT AND FOUND 1031 RETWEET INSTANCES BETWEEN THE CONSIDERED USERS.

```
fromnode=[]
tonode=[]
txt=open("rt-twitter-copen.mtx","r")
for i in txt:
    x,y=i.split()
    fromnode.append(x)
    tonode.append(y)
```

```
df=pd.DataFrame()
df['fromnode']=fromnode
df['tonode']=tonode
```

```
df[1,columns]=tonode
```

```
a=pd.read_csv('twitter.csv',index_col=None)
a
```

fromnode tonode

	fromnode	tonode
0	761	761
1	245	1
2	512	1
3	209	2
4	255	2
...
1025	750	722
1026	735	728
1027	756	728
1028	749	737
1029	761	740

1030 rows × 2 columns

1030 rows × 2 columns

1025	761	740
1028	747	737
1029	760	758

```
len(g.nodes)
```

761

```
len(g.edges)
```

1030

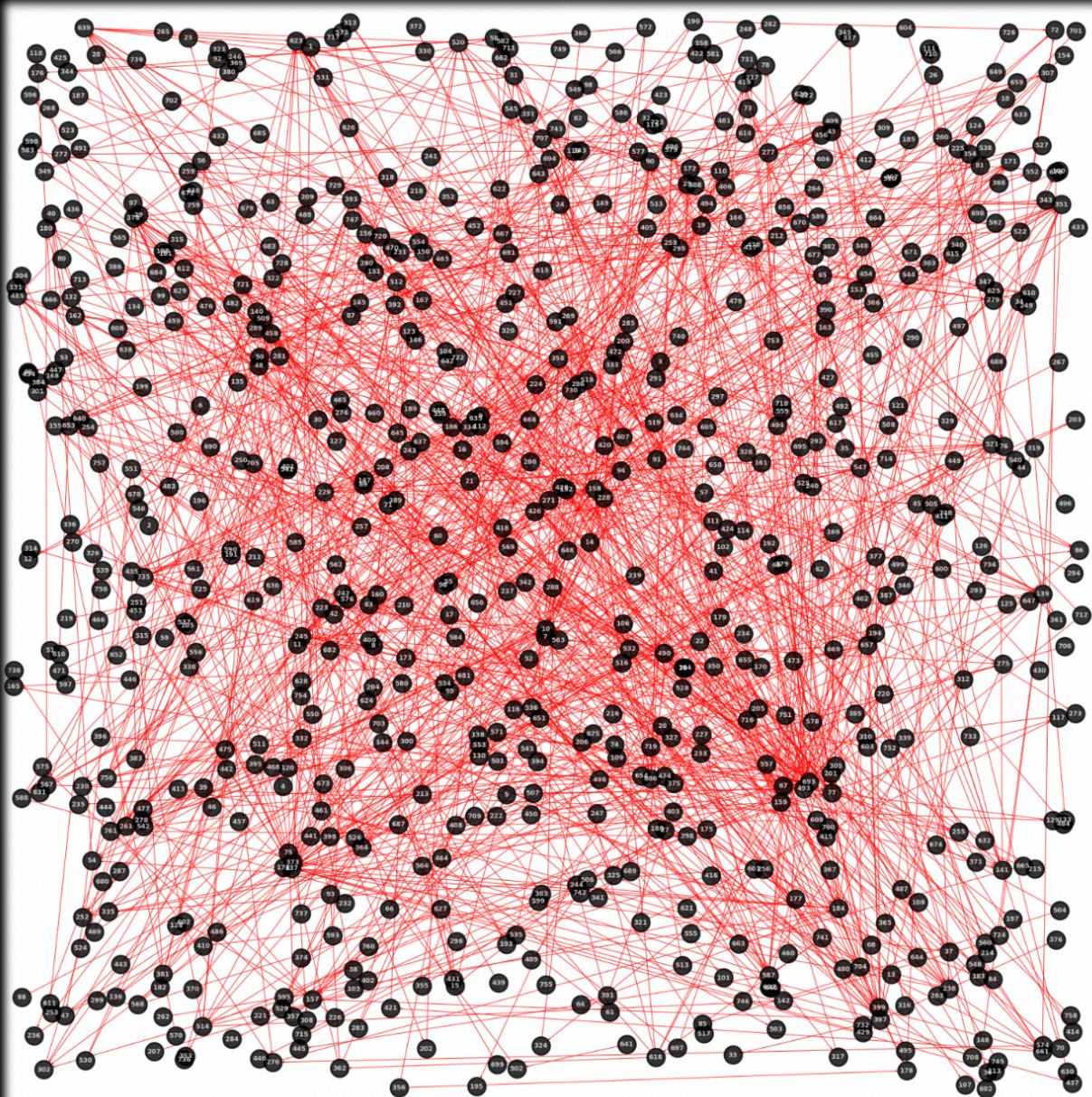
1030

VISUALIZING THE DATASET:

```
g=nx.from_pandas_edgelist(df, 'fromnode', 'tonode',create_using=nx.DiGraph())
plt.figure(figsize=(50, 50))
nx.draw_random(g,with_labels=True,node_size=2500, node_color="black", node_shape="o", alpha=0.8, linewidths=4,
                font_size=18, font_color="white", font_weight="bold", width=2, edge_color="red")
plt.show()
```

```
bff.show()
```

```
graph TD; bff --> bff
```



ABOUT GRAPH CENTRALITIES:

- GRAPH CENTRALITIES ACT AS A VERY IMPORTANT MEASURE TO IDENTIFY ALL THE IMPORTANT NODES IN A NETWORK.
- IN THIS CASE, IT IS USED TO IDENTIFY THE INFLUENTIAL PERSON IN THE RETWEET NETWORK.
- THE GRAPH CENTRALITIES SHOWN IN THIS PROJECT ARE:
 - DEGREE CENTRALITY
 - BETWEENNESS CENTRALITY
 - CLOSENESS CENTRALITY
 - EIGEN VECTOR CENTRALITY
 - CLUSTERING CENTRALITY

ABOUT DEGREE CENTRALITY:

- DEGREE CENTRALITY OF A NODE GIVES US THE TOTAL NUMBER OF EDGES THAT ARE INCIDENT ON THAT PARTICULAR NODE.
- SINCE THE NETWORK THAT WE HAVE CONSIDERED IN THIS PROJECT IS A

DIRECTED GRAPH, DEGREE CENTRALITY OF ONE NODE IN THIS NETWORK GIVES US THE TOTAL NUMBER OF INCOMING EDGES TO THAT PARTICULAR NODE.

```
v=pd.DataFrame()
deg_cent=nx.in_degree_centrality(g)
ver=list(deg_cent.keys())
deg=list(deg_cent.values())
v['vertices']=ver
v['degree_centrality']=deg
v.sort_values(by=['degree_centrality'],ascending=False)
```

	vertices	degree_centrality
105	137	0.038158
137	158	0.035526
29	10	0.031579
184	77	0.022368
264	229	0.018421
...
505	751	0.000000
504	705	0.000000
503	658	0.000000
502	578	0.000000
380	729	0.000000

761 rows × 2 columns

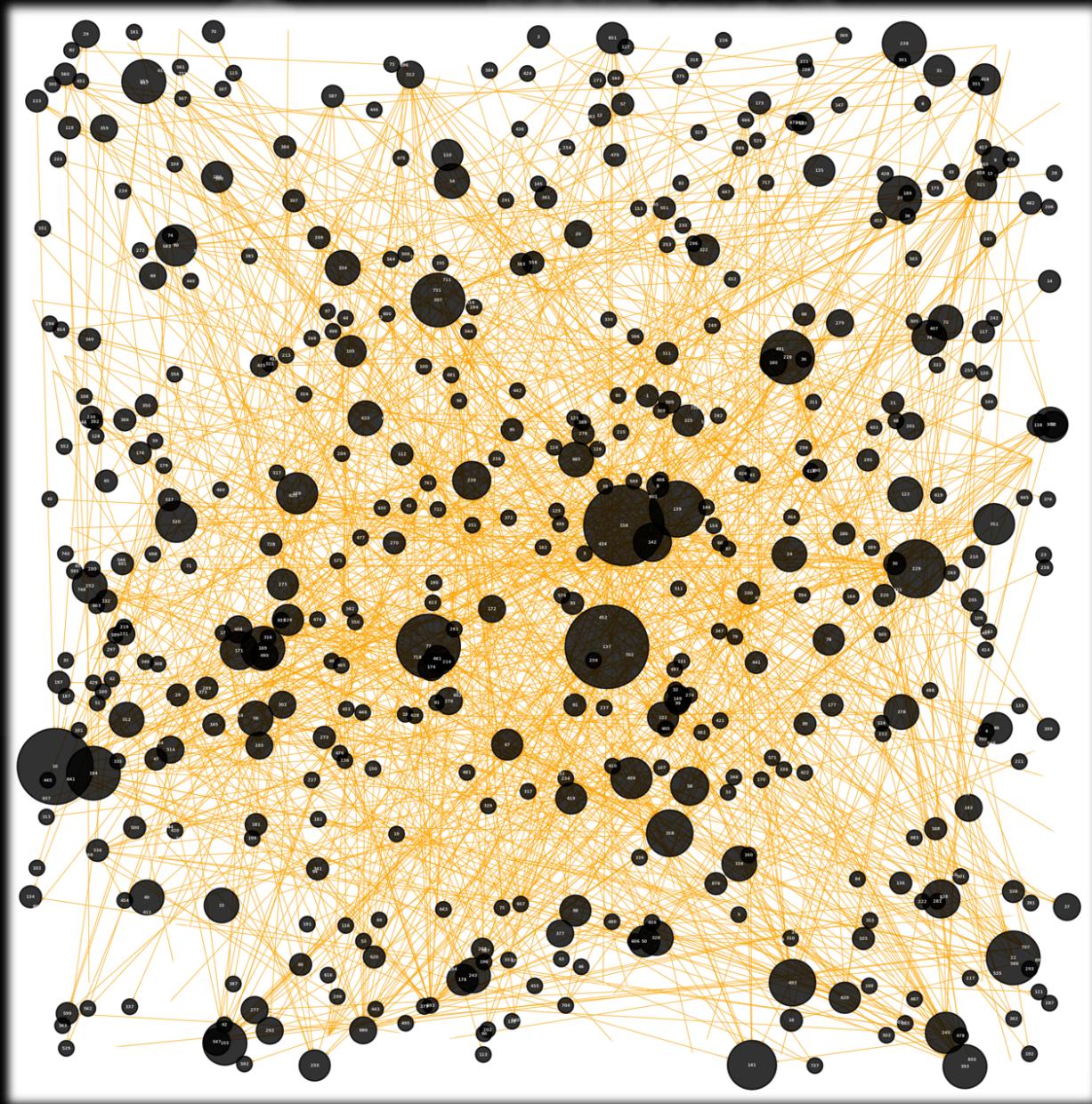
NODE 137 HAS THE HIGHEST DEGREE CENTRALITY

VISUALIZATION OF DEGREE CENTRALITY:

```
size=[d * 1500000 for d in deg]
plt.figure(figsize=(50, 50))
nx.draw_random(g, with_labels=True, node_size=size, node_color="black", node_shape="o",
               alpha=0.8, linewidths=4, font_color="white", font_weight="bold", width=2, edge_color="orange")
plt.show()
```

```
buff.show()
```

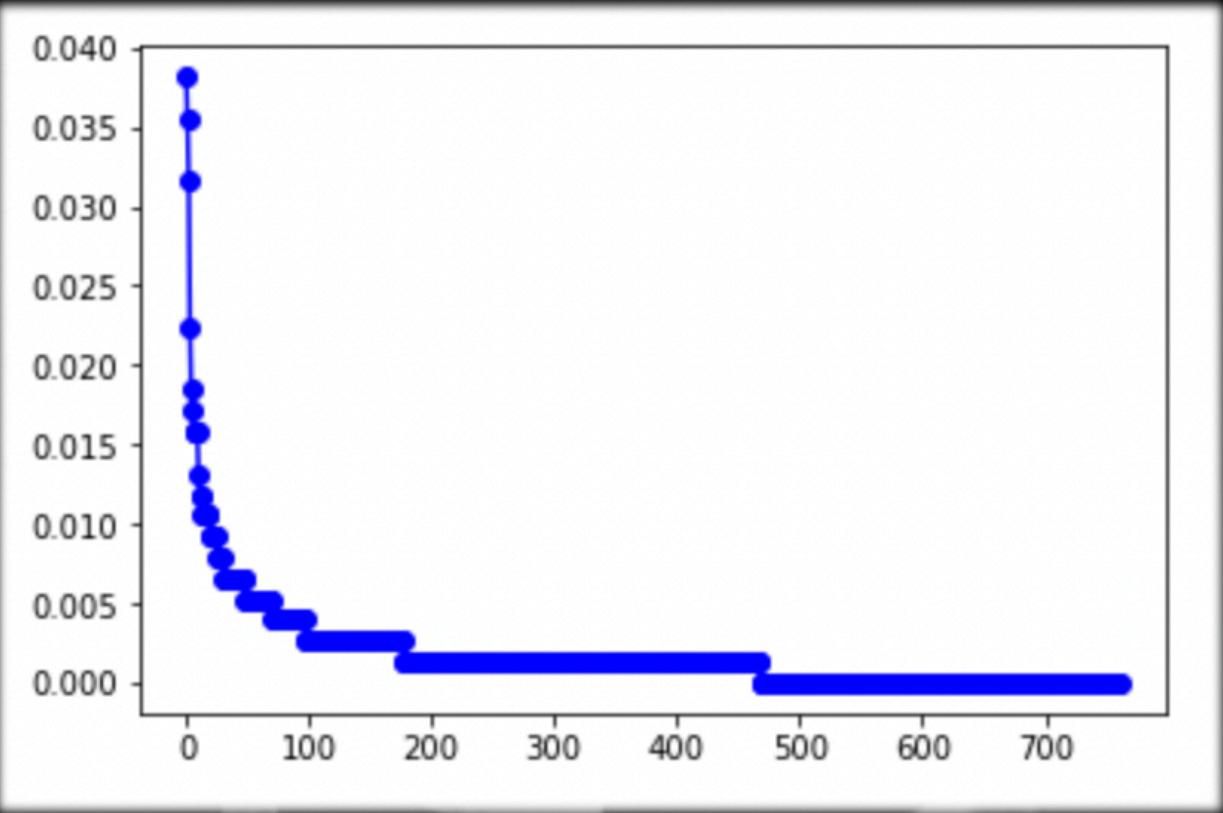
THE HIGHER THE DEGREE THE BIGGER THE NODE



INFERENCE ON DEGREE CENTRALITY:

- FROM THE ABOVE VISUALIZATION AND DEGREE CENTRALITY CALCULATION IT CAN BE SEEN THAT NODE 137 HAS THE HIGHEST DEGREE CENTRALITY, THEREFORE IT CAN BE INFERRED THAT THE POSTS IN THIS USER'S TWITTER ACCOUNT IS LIKED BY MANY OTHER USERS. SIMILARLY THE SAME INFERENCE CAN BE DRAWN FOR THE NODES 158, 10, 77 AND 229 AS THIS SET OF NODES HAS A CENTRALITY THAT IS ALMOST EQUAL TO THAT OF 105.
- **DEGREE VS RANK PLOT:**
- THE ABOVE PLOT ASSIGNS A RANK TO EACH NODE WITH RESPECT TO ITS DEGREE CENTRALITY.

```
degs=sorted(list(deg_cent.values()),reverse=True)
graph=plt.plot(degs, "b-", marker="o")
```



ABOUT BETWEENNESS CENTRALITY:

- BETWEENNESS CENTRALITY IS USED TO MEASURE THE NUMBER OF SHORTEST PATHS THAT PASSES THROUGH A NODE.
- THE FORMULA FOR BETWEENNESS CENTRALITY IS GIVEN BY:

$$C_B(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

- $\sigma_{st}(v)$ - **SHORTEST DISTANCE BETWEEN 2 VERTICES THAT PASS THROUGH THE NODE V**
- σ_{st} - **NUMBER OF SHORTEST PATHS BETWEEN THE TWO NODES**

```
bet_cen=list(nx.betweenness_centrality(g).values())
v['betweenness_centrality']=bet_cen
v.sort_values(by=['betweenness_centrality'],ascending=False)
```

	vertices	degree_centrality	betweenness_centrality
137	158	0.035526	0.002009
105	137	0.038158	0.001450
34	228	0.015789	0.001243
291	393	0.010526	0.001108
70	397	0.015789	0.000961
...
319	135	0.005263	0.000000
320	527	0.000000	0.000000
321	723	0.000000	0.000000
322	136	0.002632	0.000000
760	740	0.001316	0.000000

761 rows × 3 columns

NODE 158 HAS THE HIGHEST BETWEENNESS CENTRALITY.

VISUALIZING BETWEENNESS CENTRALITY:

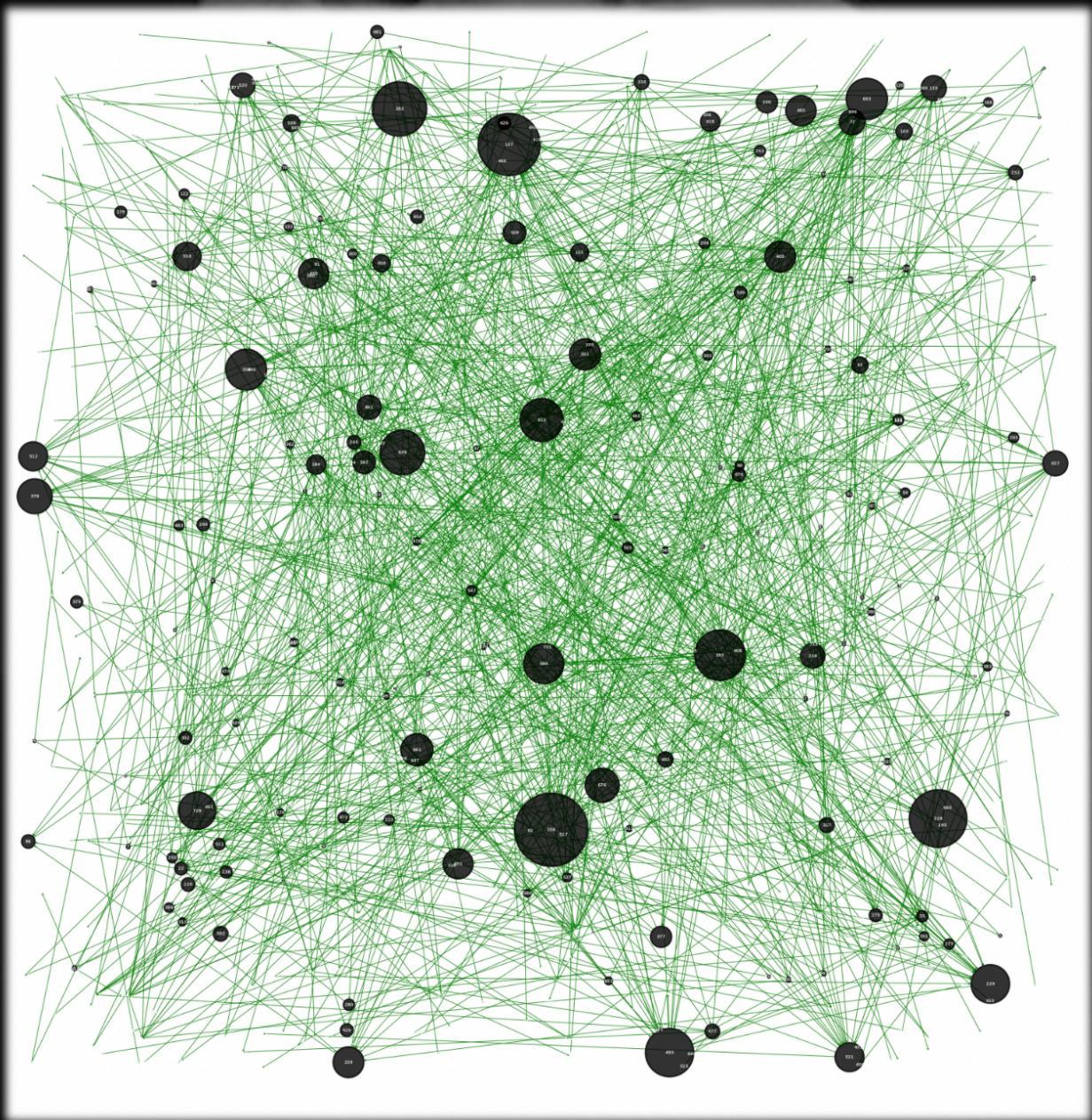
```

size1=[b * 22000000 for b in bet_cent]
plt.figure(figsize=(50, 50))
nx.draw_random(g,with_labels=True,node_size=size1, node_color="black", node_shape="o", alpha=0.8, linewidths=4,
               font_color='white',font_weight="bold", width=2, edge_color="green")
plt.show()

```

b_if::sub()

THE GREATER THE BETWEENNESS CENTRALITY THE BIGGER THE NODE



INFERENCE OF BETWEENNESS CENTRALITY:

- **BETWEENNESS CENTRALITY WILL GIVE US INFORMATION LIKE WHICH NODE OCCURS MORE WHEN WE TRAVEL FROM ONE NODE TO ANOTHER NODE.** FROM THE ABOVE GRAPH IT CAN BE SEEN THAT MOST OF THE SHORTEST DISTANCE PATHS PASS THROUGH NODE 158, WHICH TELLS US THAT THIS NODE PLAYS A REALLY IMPORTANT ROLE IN COMMUNICATION. THE NODES 137, 228, 393 AND 397 ALSO PLAY AN IMPORTANT ROLE IN COMMUNICATION AS THE BETWEENNESS CENTRALITY OF THESE NODES ARE ALSO ON PAR WITH 158.

ABOUT CLOSENESS CENTRALITY:

- **CLOSENESS CENTRALITY IS A MEASURE THAT IS USED TO MEASURE HOW FAST THE INFORMATION CAN BE PASSED FROM THAT NODE TO ALL OTHER NODES.**
- **IN OTHER WORDS, CLOSENESS CENTRALITY MEASURES HOW CLOSE ONE NODE IS TO ALL THE OTHER NODES IN A NETWORK AND IS MEASURED BY TAKING THE AVERAGE OF THE SHORTEST PATH LENGTH FROM ONE NODE TO ALL OTHER NODES.**
- **THE CLOSENESS CENTRALITY IS GIVEN BY:**

$$CC(i) = \frac{N - 1}{\sum_j d(i, j)}$$

```
close_cent=list(nx.closeness_centrality(g).values())
v['closeness_centrality']=close_cent
v.sort_values(by=['closeness_centrality'],ascending=False)
```

	vertices	degree_centrality	betweenness_centrality	closeness_centrality
184	77	0.022368	0.000244	0.069740
137	158	0.035526	0.002009	0.069524
105	137	0.038158	0.001450	0.067788
209	124	0.005263	0.000006	0.062633
282	139	0.017105	0.000239	0.060824
...
395	566	0.000000	0.000000	0.000000
393	370	0.000000	0.000000	0.000000
392	367	0.000000	0.000000	0.000000
391	340	0.000000	0.000000	0.000000
380	729	0.000000	0.000000	0.000000

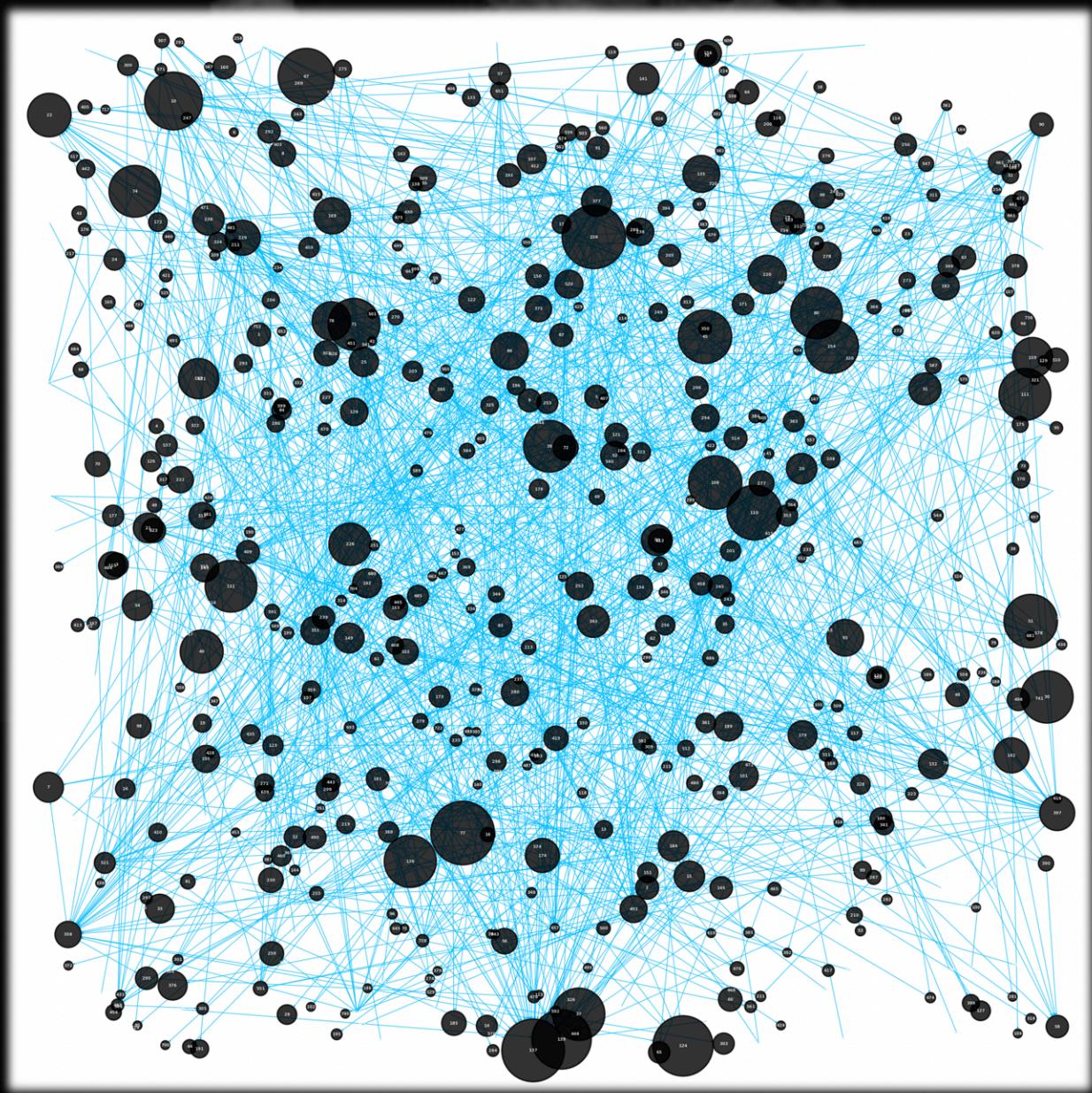
761 rows × 4 columns

NODE 77 HAS THE HIGHEST CLOSENESS CENTRALITY.

VISUALIZATION OF CLOSENESS CENTRALITY:

```
size2=[c * 500000 for c in close_cent]
plt.figure(figsize=(50, 50))
nx.draw_random(g, with_labels=True, node_size=size2, node_color="black", node_shape="o", alpha=0.8, linewidths=4,
                font_color='white', font_weight='bold', width=2, edge_color="deepskyblue")
plt.show()
```

THE GREATER THE CLOSENESS CENTRALITY THE BIGGER THE NODE



INFERENCE ON CLOSENESS CENTRALITY:

- IN TWITTER, THE SPEED OF THE INFORMATION SPREADING DEPENDS MAINLY ON HOW MANY PEOPLE RETWEET A POST AND HOW FAST IT IS DONE. THEREFORE THE HIGHER THE NUMBER OF RETWEETS THE HIGHER THE CHANCE THAT THE INFORMATION IS SPREAD. FROM THE ABOVE GRAPH IT CAN BE SEEN THAT THE NODE WITH THE LABEL 77 HAS HIGHER CLOSENESS CENTRALITY, THEREFORE THE PROBABILITY THAT THE INFORMATION IS SPREAD TO ALL THE USERS IS HIGH IF USER 77 POSTS THE INFORMATION. THE NODES 158, 137, 124 AND 139 ALSO HAVE HIGH CLOSENESS CENTRALITY.

ABOUT EIGENVECTOR CENTRALITY:

- THE EIGENVECTOR CENTRALITY IS USED TO IDENTIFY THE MOST INFLUENTIAL NODE IN A NETWORK.
- EACH NODE IS ASSIGNED AN EIGEN VALUE BASED ON THE EIGEN VALUES OF THE NODES DIRECTLY CONNECTED TO IT.

- THE EIGENVECTOR CENTRALITY OF A NODE IS CALCULATED USING:

$$\mathbf{x}_i = \frac{1}{\lambda} \sum_j^n A_{ij} x_j$$

WHERE $\lambda \mathbf{x} = \mathbf{Ax}$

```
eig_cent=list(nx.eigenvector_centrality(g).values())
v['eigenvector_centrality']=eig_cent
v.sort_values(by=['eigenvector_centrality'],ascending=False)
```

vertices	degree_centrality	betweenness_centrality	closeness_centrality	eigenvector_centrality
185	67	0.005263	0.000000	0.054892
184	77	0.022368	0.000244	0.069740
278	110	0.005263	0.000000	0.051574
283	111	0.002632	0.000000	0.046464
29	10	0.031579	0.000000	0.057693
...
396	590	0.000000	0.000000	0.000000
395	566	0.000000	0.000000	0.000000
393	370	0.000000	0.000000	0.000000
392	367	0.000000	0.000000	0.000000
380	729	0.000000	0.000000	0.000000

761 rows × 5 columns

NODE 67 HAS THE HIGHEST EIGENVECTOR CENTRALITY

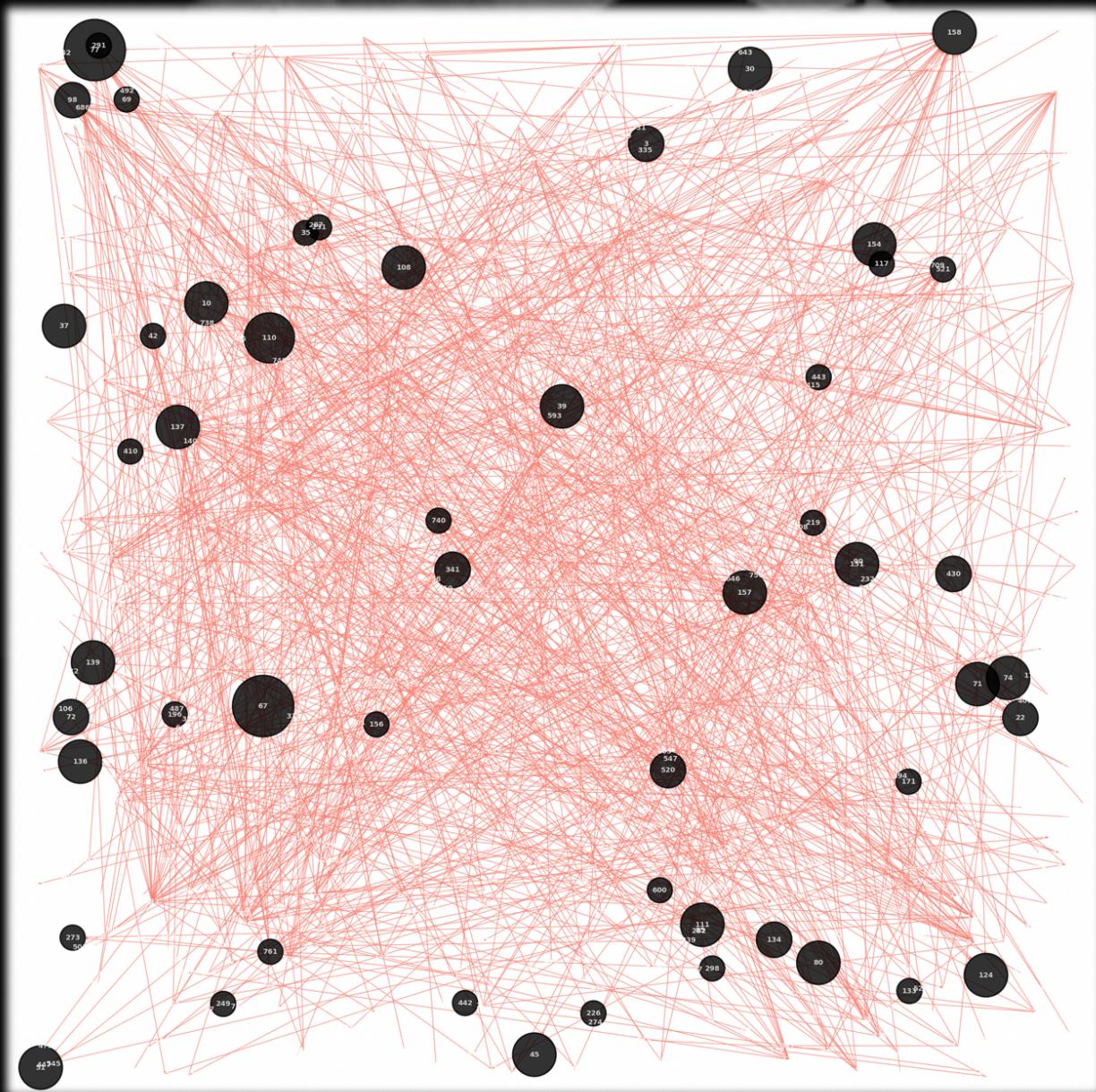
VISUALIZATION OF EIGENVECTOR CENTRALITY:

THE GREATER THE EIGENVECTOR CENTRALITY THE BIGGER THE NODE

```
size3=[e * 90000 for e in eig_cent]
plt.figure(figsize=(50, 50))
nx.draw_random(g, with_labels=True, node_size=size3, node_color="black", node_shape="o", alpha=0.8, linewidths=4,
font_color='white', font_size=20, font_weight="bold", width=2, edge_color="salmon")
plt.show()
```

```
bff.show()
```

INFERENCE ON EIGENVECTOR CENTRALITY:



- EIGENVECTOR CENTRALITY IS ALMOST SIMILAR TO DEGREE CENTRALITY EXCEPT FOR THE FACT THAT IT TAKES THE RELEVANCE OF THE NEARBY NODES INTO ACCOUNT WHEN IT IS CALCULATED. FROM THE ABOVE GRAPH IT CAN BE SEEN THAT NODE 67 HAS THE HIGHEST EIGENVECTOR CENTRALITY WHICH TELLS US THAT USER 67 IS WELL CONNECTED WITH USERS WHO HAVE CONNECTIONS WITH MANY INFLUENTIAL PEOPLE. THE NODES 77, 110, 111 AND 10 ALSO HAVE HIGH EIGENVECTOR CENTRALITIES INDICATING THAT THESE USERS ALSO HAVE INFLUENTIAL CONNECTIONS.

ABOUT CLUSTERING CENTRALITY:

- CLUSTERING CENTRALITY IS USED TO MEASURE HOW MUCH THE NODES CLUSTER TOGETHER IN OTHER WORDS IT IS USED TO MEASURE THE CONNECTIVITY OF ONE NODE TO ALL THE OTHER NODES.
- LOCAL CLUSTERING COEFFICIENT OFF A NODE SHOWS HOW CLOSE ITS NEIGHBORS ARE TO BEING A CLIQUE.

- THE LOCAL CLUSTERING COEFFICIENT IS CALCULATED BY:

$$C_l = \frac{\text{No. of neighbors}}{\frac{n(n - 1)}{2}}$$

- THE CLUSTERING COEFFICIENT OF THE WHOLE NETWORK IS CALCULATED BY:

$$C = \frac{1}{n} \sum c_i$$

- THE MAXIMUM CLUSTERING COEFFICIENT IN THE RETWEET NETWORK IS GIVEN BY **1.034693265181537E-05.**

```

n=len(g.nodes)+1
cluster=[]
adjacency_matrix=np.zeros([n,n],dtype=float)
edges=list(g.edges())
for i in edges:
    adjacency_matrix[int(i[0])][int(i[1])]=1
print('\n')
print("*****The adjacency matrix is given by*****\n")
print(adjacency_matrix)
print('\n')
print("***** CLUSTERING COEFFICIENT OF EACH NODE *****")
print('\n')
for i in range(len(adjacency_matrix)):
    neighbours=adjacency_matrix[i].sum()
    clust=((2*neighbours)/(n*(n-1)))
    cluster.append(clust)
    if(clust>0):
        print(f"The clustering coefficient of node {i} is {clust}")
print('\n')
print("***** THE CLUSTERING COEFFICIENT OF THE LEFT OUT NODES IS ZERO *****")
print('\n')
print("***** CLUSTERING COEFFICIENT OF THE WHOLE GRAPH *****\n")
print(f"The clustering coefficient of the whole graph is {sum(cluster)/n}")

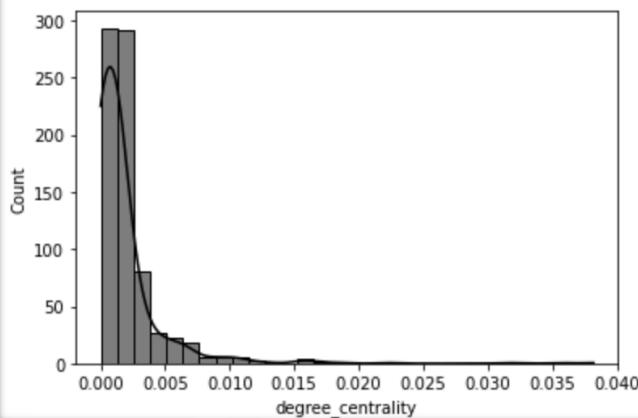
```

- THE AVERAGE CLUSTERING COEFFICIENT OF THE NETWORK IS GIVEN BY **4.6620037757523535E-06**

THE RANGE OF DIFFERENT CENTRALITIES:

- DEGREE CENTRALITY:

```
sns.histplot(x=deg,kde=True,color='k',alpha=0.5,bins=30).set(xlabel='degree_centrality')  
[Text(0.5, 0, 'degree_centrality')]
```

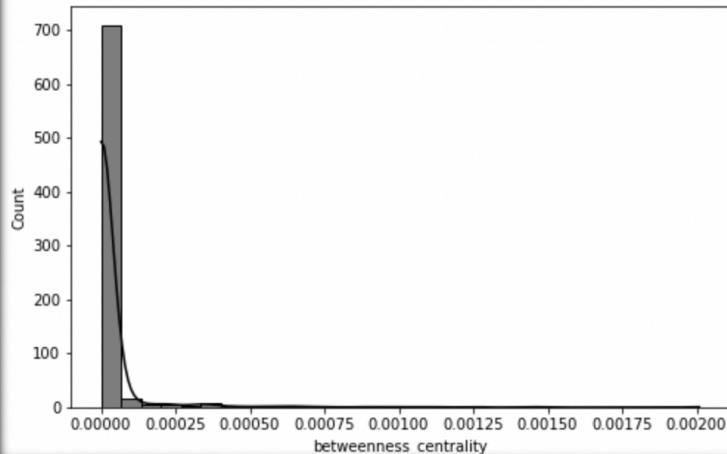


THE RANGE OF THE DEGREE CENTRALITIES OF THE 761 NODES IS IN BETWEEN 0 TO 0.038(APPROXIMATELY)

- BETWEENNESS CENTRALITY:

THE RANGE OF THE BETWEENNESS CENTRALITIES OF THE 761 NODES IS IN BETWEEN 0 TO 0.00028(APPROXIMATELY)

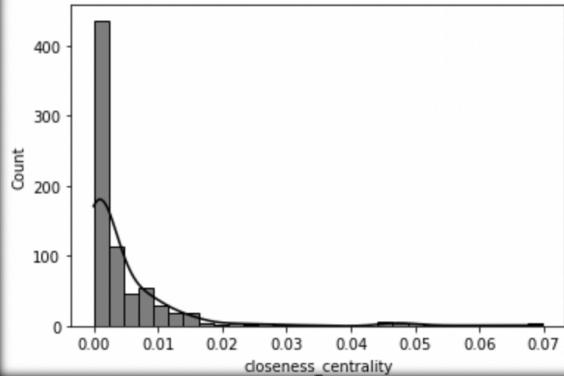
```
plt.figure(figsize=(8,5))
sns.histplot(x=bet_cent,kde=True,color='k',alpha=0.5,bins=30).set(xlabel='betweenness_centrality')
plt.show()
```



- CLOSNESS CENTRALITY:

- THE RANGE OF THE CLOSNESS CENTRALITIES OF THE 761 NODES IS IN BETWEEN 0 TO 0.02 AND 0.045 TO 0.07(APPROXIMATELY)

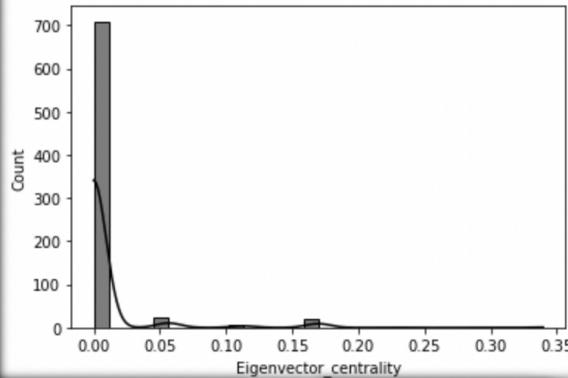
```
sns.histplot(x=close_cent,kde=True,color='k',alpha=0.5,bins=30).set(xlabel='closeness_centrality')
[Text(0.5, 0, 'closeness_centrality')]
```



- EIGENVECTOR CENTRALITY:

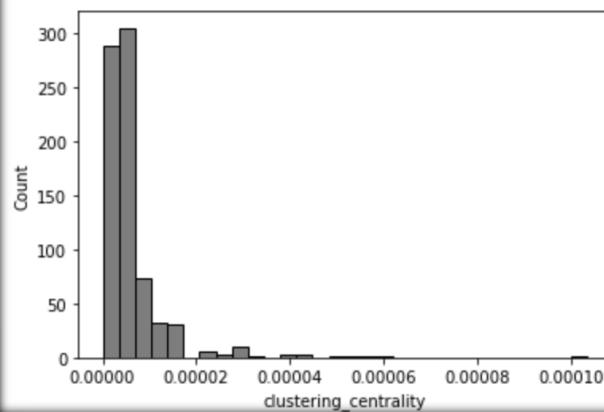
THE EIGEN VECTOR CENTRALITY VALUES OF VARIOUS NODES MOSTLY LIE ON 0 OR 0.05 OR 0.17

```
sns.histplot(x=eig_cent,kde=True,color='k',alpha=0.5,bins=30).set(xlabel='Eigenvector_centrality')  
[Text(0.5, 0, 'Eigenvector_centrality')]
```



- CLUSTERING CENTRALITY:

```
sns.histplot(x=cluster,color='k',alpha=0.5,bins=30).set(xlabel='clustering_centrality')  
[Text(0.5, 0, 'clustering_centrality')]
```



THE RANGE OF THE CLUSTERING CENTRALITIES OF THE 761 NODES IS IN BETWEEN 0 TO 0.000061 (APPROXIMATELY)

CONCLUSION:

- TAKING ALL THE GRAPH CENTRALITIES INTO ACCOUNT IT CAN BE SEEN THAT THE NODES 137, 158, 10, 77, 229, 67, 110 ARE VERY IMPORTANT AS THIS SET HAS HIGH GRAPH CENTRALITY MEASURES. THEREFORE IT CAN BE CONCLUDED THAT GRAPH CENTRALITY CAN BE USED TO FIND THE MOST INFLUENTIAL AND IMPORTANT NODES IN A NETWORK.