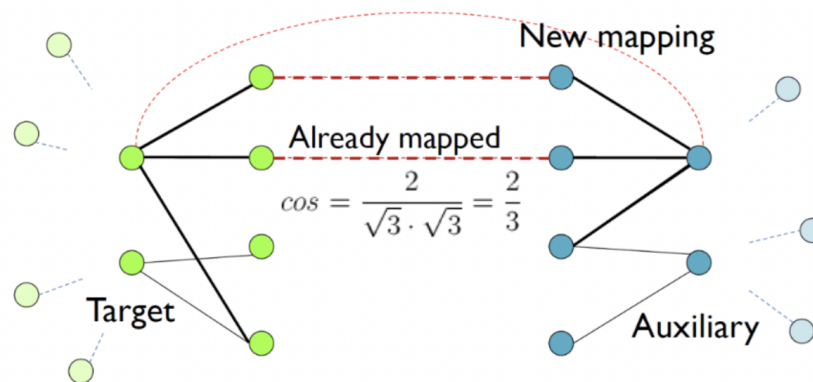# Seed-based De-anonymization

## Individual Project 1

De-anonymization is a technique used in privacy that attempts to re-identify encrypted or obscured information. De-anonymization, also referred to as data re-identification, cross-references anonymized information with other available data to identify a person, group, or transaction.

In our algorithm, we have 2 graphs. G1(V, E) and G2(V, E) .

G1(V, E) is an original graph with privacy information
G2(V, E) is an auxiliary graph with identity information.

Some of the nodes are already mapped in this. These nodes are Seed pair nodes. We need to use these already-mapped nodes to find other seed pairs.

New mapping

Already mapped

$$cos = \frac{2}{\sqrt{3} \cdot \sqrt{3}} = \frac{2}{3}$$

Target

Auxiliary

The network's topological structure and the feedback from earlier mapping iterations are used by the algorithm, which intuitively finds new mappings.

In our algorithm, we are importing graph 1 and graph 2 using networkx package by reading the edge list data.

Then we import the seed-pairs data, which is in a txt file in project data and json format in mini-dataset. When we import the json file into a variable, it is directly converted into a dictionary.

We split the seed-pair dictionary into separate nodes one has all the nodes in the seed-pair list and is graph 1 node. And the other one has all the nodes in the seed-pair list and corresponds to graph 2.

We iterate through each graph 1 node, and we find the degree of that node, then we find the degree of the node in another graph. We consider the neighbors of these nodes and check if these nodes are already mapped and if their pairs exist in the seed pair list. Here we find the unmapped nodes and measure the similarity by calculating the Score.

$$Score(node\ 1, node\ 2) = \frac{count\ of\ already\ mapped\ neighbors\ between\ (node1,\ node2)}{\sqrt{Degree\ of\ node\ 1}\ \times\ \sqrt{Degree\ of\ node\ 2}}$$

We find all the scores of these nodes. We store all of these scores into a set, and now we find eccentricity for each iteration.

$$ECCE = \frac{man1 - man2}{8}$$

Ecce = Score(max1)-Score(max2)/Standard Deviation.

We store all the scores in each iteration and find the standard deviation for this set. Then we find the first max score and second max score. We subtract these 2 scores. Then we divide it by the standard deviation. We get eccentricity value. Now we check if ECCE is greater than 0.5, we map Score(max1)'s node to the node of the first graph we are iterating. Do the

iteration until convergence. We continue this iteration until there are no other seed pair nodes to map using the above-mentioned algorithm.

After all the iterations, we append all the seed pair nodes to the initial seed-pair dictionary we imported. The final file is this dictionary loaded into a text file.

```python
[3]  1 def degree(unpaired_node_list,Graph):
     2     degree={}
     3     for i in unpaired_node_list:
     4         degree[i]=Graph.degree(i)
     5     return degree
```

```python
 1 g1 = nx.read_edgelist("seed_G1.edgelist", nodetype=int) #g1
 2 g2 = nx.read_edgelist("seed_G2.edgelist", nodetype=int) #g2
 3 sp = open('seed_node_pairs.json')
 4 ini_snp = json.load(sp) #old_g1_pair_dict/nodepairs
 5 #ini_snp = {}
 6 #with open("seed_node_pairs.txt",'r') as f:
 7 # node_pairs=f.readlines()
 8 ## pair=pair.strip('\n').split(' ')
 9    #ini_snp[int(pair[0])]=int(pair[1])
10
11 ini_g1_sp = [] #old_g1_pair
12 ini_g2_sp = [] #old_g2_pair
13 snp = {}
14 for key, value in ini_snp.items():
15    ini_g1_sp.append(int(key))
16    ini_g2_sp.append(int(value))
17    snp[int(key)] = int(value)
18 g1_sp = ini_g1_sp.copy()
19 g2_sp = ini_g2_sp.copy()
20
21 #calculating unpaired nodes
22 unpaired_g1 = []
23 unpaired_g2 = set([])
24 for i in g1.nodes:
25    if i not in ini_g1_sp:
26        unpaired_g1.append(i)
27 for j in g2.nodes:
28    if j not in ini_g2_sp:
29        unpaired_g2.add(j)
30 #calculating degree of unpaired nodes
31 unpaired_g1_degree = degree(unpaired_g1,g1)
32 unpaired_g2_degree = degree(unpaired_g2,g2)
```

```python
 1 #check in pairs
 2 for node in unpaired_g1:
 3    nodepair = set([])
 4    for neighbor in g1.neighbors(node):
 5        if neighbor in g1_sp:
 6            nodepair.add(neighbor)
 7    g2_node_score = {}
 8    for j in unpaired_g2:
 9        other_nodepair = set([])
10        for i in nodepair:
11            other_nodepair.add(snp[i])
12        count = 0
13        for x in other_nodepair:
14            if x in g2.neighbors(j):
15                count +=1
16
17        score = count/((math.sqrt(unpaired_g1_degree[node]))*(math.sqrt(unpaired_g2_degree[j])))
18        g2_node_score[j] = score
19
20    scores = set(g2_node_score.values())
21    max1 = max(scores) #max score
22    for k, v in g2_node_score.items():
23        if v == max1:
24            g2_max_score = k
25    std = np.std(list(scores))
26    scores.remove(max1)
27    max2 = max(scores)
28    if std == 0:
29        ecce = 0
30    else:
31        ecce = (max1 - max2)/std    #ecce score
32    if ecce > 0.5:
33        g1_sp.append(node)
34        g2_sp.append(g2_max_score)
35        snp[node]=g2_max_score
36        unpaired_g1.remove(node)
37        unpaired_g2.remove(g2_max_score)
38
```

Here we are writing the output into a text file.

```
38
39 with open('Indrojuproject1output.txt', 'w') as f:
40    for i in g1_sp:
41      #z = str(i)+ ' ' + str((snp[i])+'\n']
42      z = str(i)+ ' ' + str(snp[i])+'\n'
43      f.write(z)
44 f = open("Indrojuproject1output.txt", "r")
45 files.download('Indrojuproject1output.txt')
46
47
```