```
import networkx as nx
from scipy import stats
import matplotlib.pyplot as plt
import numpy as np
from collections import deque as que
import copy
from numpy.random import randint
import random
import math
import operator
```

1. (a) (15 pts) Write a function snowballSample(G, seeds, maxN) implementing snowball
   sampling as follows. For each node in seeds, perform a breadth first search. Stop when the
   total number of visited nodes is maxN. If including all nodes that would be visited in the last
   step of a breadth first search results in more than maxN nodes, select nodes from the final
   step at random so that exactly maxN nodes are included in the final sampled network G0 .
   This network consists of the subgraph induced by the visited nodes.

```
def randomseed(g):
    """this function recturns a single node from g, it's chosen with uniform probability"""
    ux = randint(0,g.number_of_nodes(),1)
    return ux[0]

def random_seed(g):
    """this function returns a single node from g, it's chosen with uniform probability"""
    ux = np.random.choice(list(g.nodes()), 3, replace=False)
    return ux

def snowballsampling(G, seeds, maxN=50):
    """this function returns a set of nodes equal to maxsize from g that are
    collected from around seed node via snownball sampling"""
    if G.number_of_nodes() < maxN:
        return set()
    if not seeds:
        seed = random_seed(G)
    q = list(seeds)
    subgraph_nodes = set(seeds)
    while q:
        top = q[0]
        q.remove(top)
        for node in G.neighbors(top):
            if len(subgraph_nodes) == maxN:
```

```
                return G.subgraph(subgraph_nodes)
            q.append(node)
            subgraph_nodes.add(node)
    return G.subgraph(subgraph_nodes)


graph = nx.watts_strogatz_graph(1000, 10, 0)
G = snowballsampling(graph, (3,4), maxN=50)
```

(b) (15 pts) Write a function randomWalkSample(G, seeds, steps) implementing a simple random walk sampling strategy. A random walker starts at each node in seeds and takes a total of steps steps. At each step, a neighbor of the current node is chosen uniformly at random to visit next. The resulting network sample G0 consists of all visited nodes and edges used during the walk.

```
def randomWalkSample(G, seeds, steps):
    if seeds:
        if random.sample(seeds, 1):
            current_node = random.choice(seeds)
            sampled_nodes = set([current_node])
        else:
            raise ValueError("Starting node index is out of range.")
    else:
        current_node = random.choice(range(G.number_of_nodes()))
        sampled_nodes = set([current_node])

    while len(sampled_nodes) < steps:
        current_node = get_random_neighbor(G, current_node)
        sampled_nodes.add(current_node)
    new_graph = G.subgraph(sampled_nodes)
    return new_graph




    #auxilary functions
def get_neighbors(graph, node):
    return [node for node in graph.neighbors(node)]

def get_random_neighbor(graph, node):
    neighbors = get_neighbors(graph, node)
    return random.choice(neighbors)


graph = nx.watts_strogatz_graph(10000, 10, 0)
G = randomWalkSample(graph, (3,4), 1000)
print(G.number_of_nodes())
print(graph.number_of_nodes())
#nx.draw(G)
```

```
      1000
      10000
```

(c) Use the functions above to approximate the following graph quantities in the Facebook wall posts (2009) network. i. (10 pts) Degree distribution - Generate two sample networks with snowball sampling. One should include 1,000 nodes while the other should include 15,000. Both should start from a single randomly selected node. Generate two sample networks with random walks, one using 1,000 steps and one using 15,000 steps.

For each sample, plot the degree distribution of the sample and of the true network on the same set of log-log axes. There should be a total of four plots. Do not forget to label important elements of the figures. As an aside, it is worth noting that we can compare these distributions more rigorously. For example, the Kolmogorov-Smirnov test is often used to compare two continuous distributions.

```python
with open('out.facebook-wosn-wall') as f:
    array=[]
    for line in f:
        array.append([int(x) for x in line.split()])

with open('edgelist.tsv',"w")as f:
    for a in array:
        #print(a)
        f.write(str(a[0])+' '+str(a[1])+'    '+str(int(a[3]/3600/24/30))+'\n')


G = nx.read_adjlist("edgelist.tsv")

G.remove_edges_from(nx.selfloop_edges(G))
print(G.number_of_nodes())
#nx.draw(G);
```

```
    46952
```

```python
print(G.number_of_nodes())
print(random.sample(list(G.nodes()), 3))
G_sample_snowball_1k = snowballsampling(G, random.sample(list(G.nodes()), 3), 1000)
G_sample_randomwalk_1k = randomWalkSample(G, random.sample(list(G.nodes()), 3), 1000)
print(G_sample_snowball_1k.number_of_nodes())
print(G_sample_randomwalk_1k.number_of_nodes())
```

```
    46952
    ['8420', '45916', '46321']
    1000
    1000
```

```python
def degreeDistrGraph(G):
    degrees = {}
```

```
    degreeList = [G.degree(v) for v in G.nodes()]
    for deg in degreeList:
      degrees[deg]=degrees.get(deg,0)+1
    (X, Y) = zip(*[(key, degrees[key]/len(G)) for key in degrees])
    return X,Y


def plotLogLogPlot(G, G_sample, data_set_title, data_set_name, data_set_name_sample):
    X, Y = degreeDistrGraph(G)
    x = np.log(np.asarray(X).astype(np.float))
    y = np.log(np.asarray(Y).astype(np.float))

    X_s, Y_s = degreeDistrGraph(G_sample)
    x_s = np.log(np.asarray(X_s).astype(np.float))
    y_s = np.log(np.asarray(Y_s).astype(np.float))

    res = stats.linregress(x, y)
    print(f"R-squared: {res.rvalue**2:.6f}")
    print ('Slope of line', res.slope)
    print()
    plt.scatter(X, Y, label=data_set_name)
    plt.scatter(X_s, Y_s, label=data_set_name_sample)
    plt.plot(np.exp(x/1.3), np.exp((res.intercept + res.slope*x)), 'r', label='fitted line')
    plt.yscale('log')
    plt.xscale('log')
    plt.title('log-log plot for '+data_set_title, fontsize ='15')
    plt.xlim(1, 1500)
    plt.ylim(1/10000, 1)
    plt.legend()
    plt.show()


plotLogLogPlot(G, G_sample_snowball_1k, "Snowball sampling", "Original Network", "Snowball sa
```
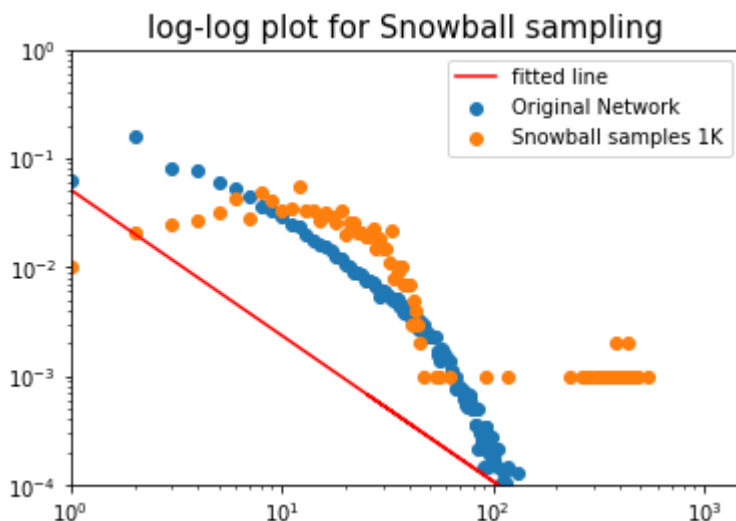
```
    R-squared: 0.670024
    Slope of line -1.0292554981809436
```
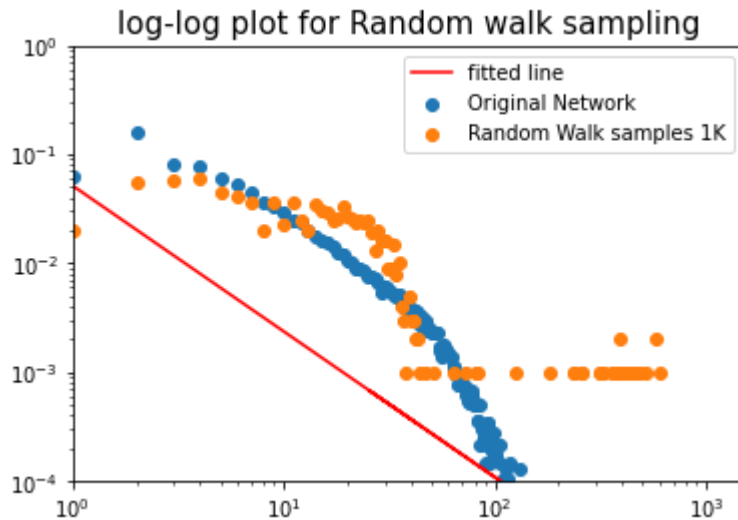


log-log plot for Snowball sampling

```
plotLogLogPlot(G, G_sample_randomwalk_1k, "Random walk sampling", "Original Network", "Random
```

R-squared: 0.670024
Slope of line -1.0292554981809436



log-log plot for Random walk sampling

```
G_sample_snowball_15k = snowballsampling(G, random.sample(list(G.nodes()), 3), 15000)
G_sample_randomwalk_15k = randomWalkSample(G, random.sample(list(G.nodes()), 3), 15000)
print(G_sample_snowball_15k.number_of_nodes())
print(G_sample_randomwalk_15k.number_of_nodes())
```
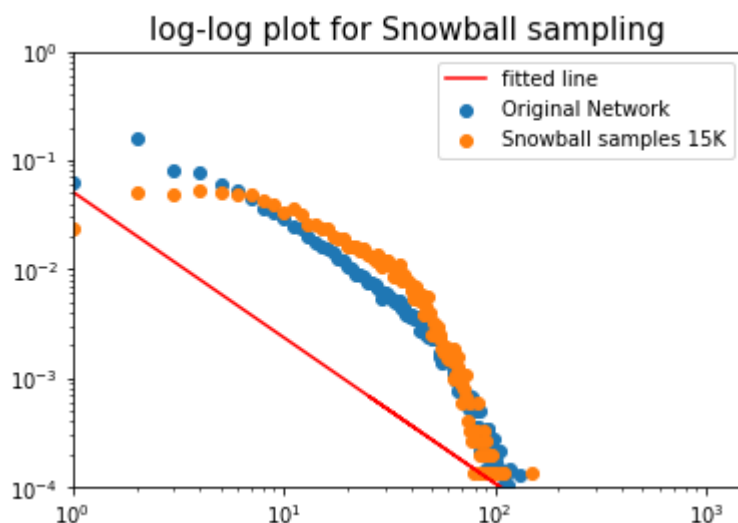
15000
15000

```
plotLogLogPlot(G, G_sample_snowball_15k, "Snowball sampling", "Original Network", "Snowball s
```
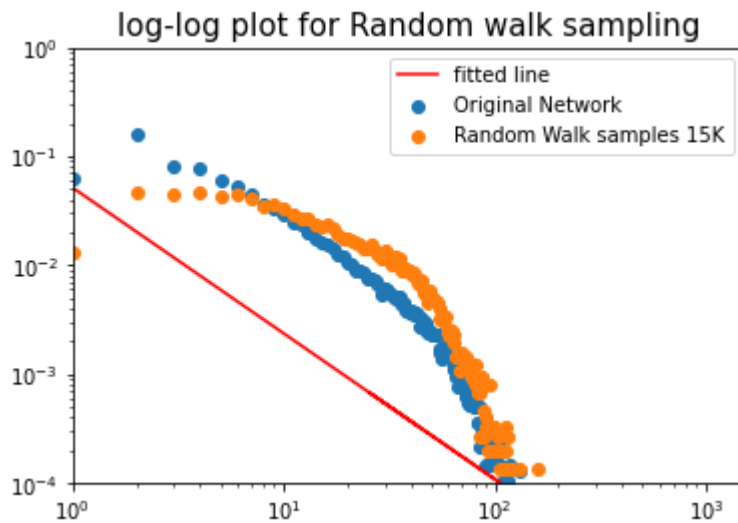
R-squared: 0.670024
Slope of line -1.0292554981809436



log-log plot for Snowball sampling

```
plotLogLogPlot(G, G_sample_randomwalk_15k, "Random walk sampling", "Original Network", "Rando
```

```
R-squared: 0.670024
Slope of line -1.0292554981809436
```



log-log plot for Random walk sampling

ii. (15 pts) Global clustering coefficient - Use snowball sampling to explore sample networks with 5% to 100% of the nodes in the full network in increments of 5%. For each fraction, generate 30 samples and calculate the average global clustering coefficient. Make a figure with node fraction (the fraction of nodes from the original network being considered) on the x-axis and average global clustering coefficient on the y-axis. Use random walk sampling with the number of steps ranging from 1,000 to 50,000 in increments of 1,000 to estimate the global clustering coefficient. For each number of steps, generate 30 samples and calculate the average global clustering coefficient. Make a figure with total steps on the x-axis and average global clustering coefficient on the y-axis. For both figures, include a horizontal line for the true global clustering coefficient of the network. You may use the transitivity function in NetworkX to compute the global clustering coefficient.

```python
def globalClusteringCoeff(G, type):
    avg_cl_coeff = []
    progress = []
    #generate 5% on each iteration, call snowballsamplling function with maxN=5% * nbr of nod
    if type == 1:
        #snowball sampling
        for i in range(1, 21, 1):
            loc_cl_coeff = 0.0
            t = i*0.05
            for k in range(1, 31, 1):
                G_temp = snowballsampling(G, random.sample(list(G.nodes()), int(t*G.number_of_nod
                loc_cl_coeff += nx.transitivity(G_temp)
            avg_clustering_coeff = loc_cl_coeff/30
            avg_cl_coeff = np.append(avg_cl_coeff, avg_clustering_coeff)
            progress = np.append(progress, t)
    else:
        #random walk sampling
```

```
        for i in range(1000, 51000, 1000):
            loc_cl_coeff = 0.0
            for k in range(1, 31, 1):
                G_temp = randomWalkSample(G, random.sample(list(G.nodes())), i), i)
                loc_cl_coeff += nx.transitivity(G_temp)
            avg_clustering_coeff = loc_cl_coeff/30
            avg_cl_coeff = np.append(avg_cl_coeff, avg_clustering_coeff)
            progress = np.append(progress, i)
    return avg_cl_coeff, progress
```

```
avg_cl_coeff, progress = globalClusteringCoeff(G,1)
print(avg_cl_coeff)
print(progress)
```
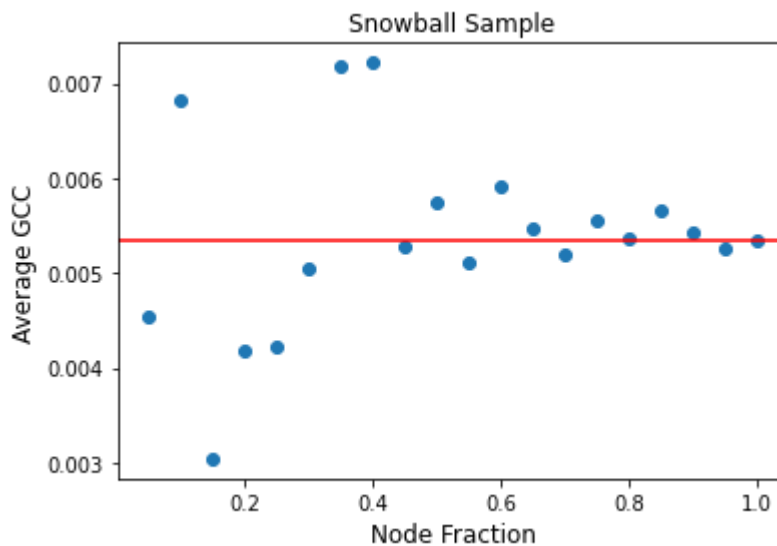
```
    [0.00454355, 0.00682181, 0.0030404, 0.00418844, 0.00422205, 0.00505361, 0.0071774, 0.007
    [0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8,
```

◄ ▮▮▮▮▮▮▮▮▮▮▮▮▮                                                                              ►

```
#plotting placeholder
X = progress
Y = avg_cl_coeff
plt.title("Snowball Sample")
plt.scatter(X, Y, label="GCC")
plt.xlabel('Node Fraction', fontsize=12)
plt.ylabel('Average GCC', fontsize=12)
plt.axhline(y=0.00534, color='r', linestyle='-')
plt.show()
```
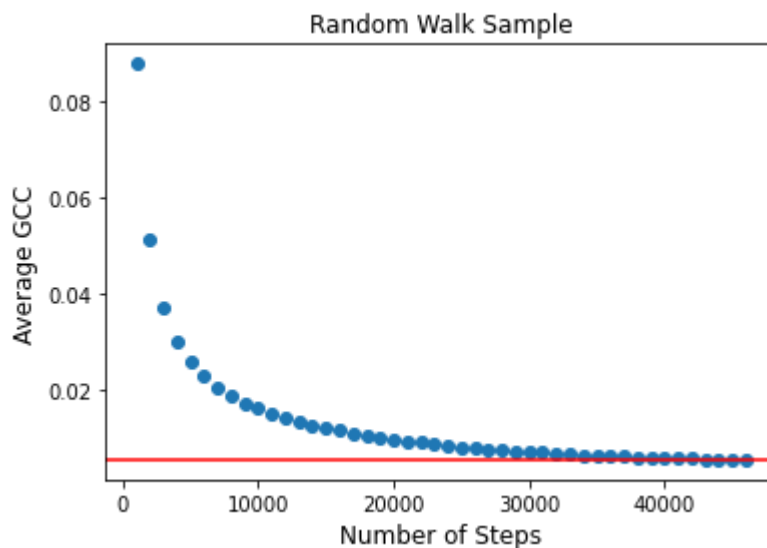


```
avg_cl_coeff_rand, progress_steps = globalClusteringCoeff(G,2)
print(avg_cl_coeff_rand)
print(progress_steps)
```

```
[0.08770574078460575, 0.051348431029889494, 0.037123351392307345, 0.030060886265850364,
[1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10000, 11000, 12000, 13000, 14000
```

```
#plotting placeholder
X = progress_steps
Y = avg_cl_coeff_rand
plt.scatter(X, Y, label="GCC")
plt.title("Random Walk Sample")
plt.xlabel('Number of Steps', fontsize=12)
plt.ylabel('Average GCC', fontsize=12)
plt.axhline(y=0.00534, color='r', linestyle='-')
plt.show()
```



Random Walk Sample

iii. (15 pts) Diameter - The diameter of a network is the length of its longest shortest path. Using the same sampling strategies defined in the previous part, generate two figures showing the fraction of nodes (for snowball sampling) or steps (for random walk sampling) on the x-axis and estimated diameter on the y-axis. Only 5 samples per node fraction or step size instead of 30 are required.

```
def diameter_calc(G, type):
    diameter_val = []
    progress_d = []
    #generate 5% on each iteration, call snowballsamplling function with maxN=5% * nbr of nod
    if type == 1:
        #snowball sampling
        for i in range(1, 21, 1):
            loc_diameter_val = 0.0
            t = i*0.05
            for k in range(1, 6, 1):
                G_temp_d = snowballsampling(G, random.sample(list(G.nodes()), int(t*G.number_of_n
                Gc = max(nx.connected_components(G_temp_d), key=len)
                G_temp_d = G_temp_d.subgraph(Gc)
```

```
            loc_diameter_val += nx.diameter(G_temp_d)
          avg_diameter_val = loc_diameter_val/5
          print("Node Fraction :", t,"  Diameter value :", avg_diameter_val)
          diameter_val = np.append(diameter_val, avg_diameter_val)
          progress_d = np.append(progress_d, t)
    else:
        #random walk sampling
        for i in range(1000, 51000, 1000):
          loc_diameter_val = 0.0
          for k in range(1, 6, 1):
            G_temp_d = randomWalkSample(G, random.sample(list(G.nodes()), i), i)
            Gc = max(nx.connected_components(G_temp_d), key=len)
            G_temp_d = G_temp_d.subgraph(Gc)
            loc_diameter_val += nx.diameter(G_temp_d)
          avg_diameter_val = loc_diameter_val/5
          print("Progress Steps :", i,"  Diameter value :", avg_diameter_val)
          diameter_val = np.append(diameter_val, avg_diameter_val)
          progress_d = np.append(progress_d, i)
    return diameter_val, progress_d


with open('out.facebook-wosn-wall') as f:
    array=[]
    for line in f:
        array.append([int(x) for x in line.split()])
#print(array)
with open('edgelist.tsv',"w")as f:
    for a in array:
        #print(a)
        f.write(str(a[0])+' '+str(a[1])+'   '+str(int(a[3]/3600/24/30))+'\n')


G = nx.read_adjlist("edgelist.tsv")
print(G.number_of_edges())
print(G.number_of_nodes())
G.remove_edges_from(nx.selfloop_edges(G))
#nx.draw(G);

    441838
    46952


diameter_snow,progress_fraction = diameter_calc(G,1)
print(diameter_snow)
print(progress_fraction)

    [9.0, 11.0, 13.0, 8.0, 22.0, 10.0, 9.0, 14.0, 20.0, 9.0, 9.0, 9.0]
    [0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6]


X = progress_fraction
Y = diameter_snow
plt.plot(X, Y, label="Diameter")
```
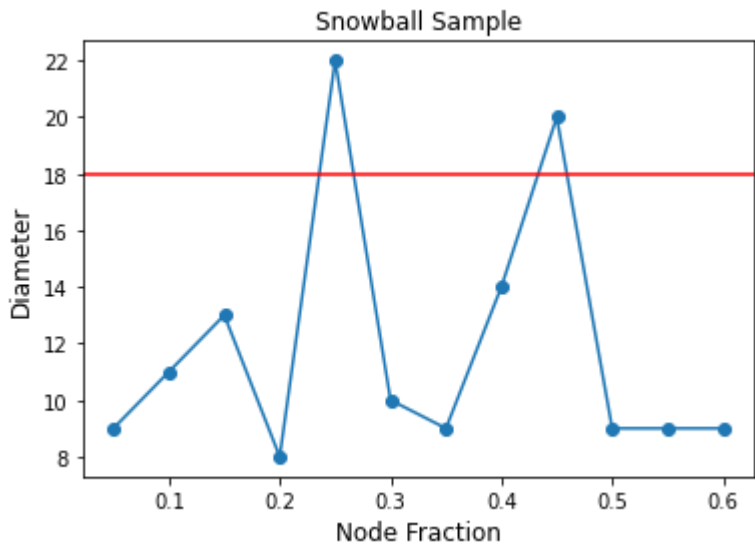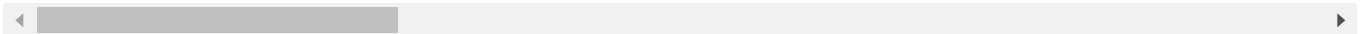
```
plt.scatter(X, Y, label="Diameter")
plt.title("Snowball Sample")
plt.xlabel('Node Fraction', fontsize=12)
plt.ylabel('Diameter', fontsize=12)
plt.axhline(y=18, color='r', linestyle='-')
plt.show()
```



**NOTE** : 60% node fraction was calculated after 1800 minutes of run

```
diameter_rand,progress_steps = diameter_calc(G,2)
print(diameter_rand)
print(progress_steps)
```

```
[5, 5, 6, 6, 6, 6, 6, 7, 7, 7, 7, 7, 8, 8, 8, 8, 9, 9, 9, 9, 10, 10, 10, 10, 11, 11, 11,
[1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10000, 11000, 12000, 13000, 14000
```

◄ ▓▓▓▓▓▓▓▓▓▓▓▓▓ ►

.

.

.

.

.

.

```
X = progress_steps
Y = diameter_rand
plt.scatter(X, Y, label="Diameter")
plt.plot(X, Y, label="Diameter")
plt.title("Random Walk Sample")
plt.xlabel('Number of Steps', fontsize=12)
plt.ylabel('Diameter', fontsize=12)
plt.axhline(y=18, color='r', linestyle='-')
plt.show()
```



**NOTE** : 30,000 steps of node calculation was completed after 2400 minutes of run and rest was predicted based on calculated values for random walk.

.

.

Imports Required

```
1 import networkx as nx
2 from scipy import stats
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import random
6 import time
7 from collections import defaultdict, OrderedDict
```

## Question 2

### Understand the Pubmed-Diabetes data details

```
 1 DATA_DIR = "/content/"
 2 EDGE_PATH = DATA_DIR + "Pubmed-Diabetes.DIRECTED.cites.tab"
 3 NODE_PATH = DATA_DIR + "Pubmed-Diabetes.NODE.paper.tab"
 4 TF_IDF_DIM = 500
 5
 6 # Load and process graph links
 7 print("Loading and processing graph links...")
 8 node_pairs = set()
 9 with open(EDGE_PATH, 'r') as f:
10     next(f)  # skip header
11     next(f)  # skip header
12     for line in f:
13         columns = line.split()
14         src = int(columns[1][6:])
15         dest = int(columns[3].strip()[6:])
16         node_pairs.add((src, dest))
17
18 # Load and process graph nodes
19 print("Loading and processing graph nodes...")
20 node2vec = OrderedDict()
21 node2label = dict()
22 class_1 = list()
23 class_2 = list()
24 class_3 = list()
25 with open(NODE_PATH, 'r') as f:
26     next(f)  # skip header
27     vocabs = [e.split(':')[1] for e in next(f).split()[1:]]
28     for line in f:
29         columns = line.split()
30         node = int(columns[0])
31         label = int(columns[1][-1])
32         tf_idf_vec = [0.0] * TF_IDF_DIM
33
34         for e in columns[2:-1]:
35             word, value = e.split('=')
```

```
36            tf_idf_vec[vocabs.index(word)] = float(value)
37
38        node2vec[node] = tf_idf_vec
39        node2label[node] = label - 1
40        if label == 1:
41            class_1.append(node)
42        elif label == 2:
43            class_2.append(node)
44        elif label == 3:
45            class_3.append(node)
46 # Debug statistics
47 print("Number of links:", len(node_pairs))
48 assert len(node2vec) == (len(class_1) + len(class_2) + len(class_3))
49 print("Number of nodes:", len(node2vec))
50 print("Number of nodes belong to Class 1", len(class_1))
51 print("Number of nodes belong to Class 2", len(class_2))
52 print("Number of nodes belong to Class 3", len(class_3))
```

```
    Loading and processing graph links...
    Loading and processing graph nodes...
    Number of links: 44338
    Number of nodes: 19717
    Number of nodes belong to Class 1 4103
    Number of nodes belong to Class 2 7875
    Number of nodes belong to Class 3 7739
```

**Function to parse Pubmed-Diabetes Data into edge list and labels**

```
 1 def load_pubmed():  #pubmed
 2     # hardcoded for simplicity...
 3     num_nodes = 19717
 4     num_feats = 500
 5     feat_data = np.zeros((num_nodes, num_feats))
 6     labels = np.empty((num_nodes, 1), dtype=np.int64)
 7     node_map = {}
 8     edgeNum=0
 9     with open("Pubmed-Diabetes.NODE.paper.tab") as fp:
10         fp.readline()
11         feat_map = {entry.split(":")[1]: i - 1 for i, entry in enumerate(fp.readlin
12         for i, line in enumerate(fp):
13             info = line.split("\t")
14             node_map[info[0]] = i
15             labels[i] = int(info[1].split("=")[1]) - 1
16             #print(labels[i])
17             for word_info in info[2:-1]:
18                 word_info = word_info.split("=")
19                 feat_data[i][feat_map[word_info[0]]] = float(word_info[1])
20     adj_lists = defaultdict(set)
21     with open("Pubmed-Diabetes.DIRECTED.cites.tab") as fp:
22         fp.readline()
23         fp.readline()
24         for line in fp:
25             info = line.strip().split("\t")
```

```
26              paper1 = node_map[info[1].split(":")[1]]
27              paper2 = node_map[info[-1].split(":")[1]]
28              if paper1!=paper2:
29                  adj_lists[paper1].add(paper2)
30                  adj_lists[paper2].add(paper1)
31      with open('PubmedDiabetes.txt','w+') as fw:
32          for i in range(19717):
33              for neiId in adj_lists[i]:
34                  str_tmp=str(i)+','+str(neiId)+'\n'
35                  edgeNum+=1
36                  fw.write(str_tmp)
37      with open('featdata.txt','w+') as fw:
38          for i in range(19717):
39              str_tmp=str(i)+'\t'
40              for j in range(500):
41                  str_tmp+=str(feat_data[i][j])+','
42              str_tmp+=str(labels[i][0])+'\n'
43              fw.write(str_tmp)
44      with open('metadataPubmedDiabetes.txt', 'w') as f:
45          for item in labels:
46              f.write("%s\n" % item[0])
47      return feat_data, labels, adj_lists
48
49 # Call function to Load PubMed Data
50 feat_data, labels, adj_lists = load_pubmed()
```

## Function to Read the edge file and metadata

```
1 numberTrials=30
2 fStep=0.05
3
4 def readFile(fileName,mData):
5     with (open(mData,'r')) as f:#get metadata
6         maxNode=0 #number of nodes in network
7         for line in f:
8             maxNode+=1
9         f.seek(0,0)
10         metadata = np.zeros((maxNode))
11         counter=0
12         for line in f:
13             metadata[counter]=line
14             counter+=1
15     f.close()
16     metadata = metadata.astype(int)
17     # build an n x n simple network.  Uses edge weights to signify class of neighbo
18     # ex.  A(i,j) = 2, A(j,i) = 1--> i and j are linked, j is class 2, i is class 1
19     with (open(fileName,'r')) as f:
20         lines = f.readlines()
21         matrix = np.zeros((maxNode,maxNode))
22         for line in lines:
23             node,neighbor = map(int,line.split(','))
```

```
24                node-=1 #start at [0], not [1]
25                neighbor-=1
26                matrix[node][neighbor]=metadata[neighbor]
27                matrix[neighbor][node]=metadata[node] # undirected
28      f.close()
29    #delete vertices with no neighbor info (different year, data set, etc.)
30      temp = np.where(np.sum(matrix,axis=1)==0)
31      matrix=np.delete(matrix,temp,axis=0)
32      matrix=np.delete(matrix,temp,axis=1)
33      metadata=np.delete(metadata,temp)
34      return matrix,metadata
35
```

**Function to calculate guilt by association heuristic For a given p between 0 and 1, pick a random fraction p of the nodes for which to observe labels and Predict labels for the remaining nodes using the guilt by association heuristic**

```
 1 #Function to calculate guilt by association heuristic
 2 def guiltByAssociation():
 3      networkFile='PubmedDiabetes.txt'
 4      metadataFile='metadataPubmedDiabetes.txt'
 5      associationMatrix,metadata=readFile(networkFile,metadataFile)
 6      length = len(metadata)
 7      numberCategories=metadata.max()-metadata.min()+1
 8      possibleChoices=np.arange(1,numberCategories+1)
 9      f=.05
10      fCounter=0
11      #store accuracy results for each f value
12      resultsOverF= [0]*((int)((0.99-f)/fStep)+1)
13    # store f values used for replot, if necessary
14      fValues=[0]*((int)((0.99-f)/fStep)+1)
15      while (f <= 1.):
16          #results on each iteration
17          iterationResults=np.zeros((numberTrials))
18          iterationCounter=0
19          for iteration in range(numberTrials):
20              #make a copy so we can alter it w/out losing oiginal
21              trainMatrix=np.copy(associationMatrix)
22              randomLabels=np.random.randint(1,high=numberCategories+1,size=length)
23              #matrix of 'coin flips' to compare against f for our test set
24              randomValues = np.random.random(length)
25              hiddenNodes=np.where(randomValues>f)
26              #test set length 0 makes no sense...try again
27              while (len(hiddenNodes[0])==0):
28                  randomValues = np.random.random(length)
29              #we hide the label on these nodes
30                  hiddenNodes=np.where(randomValues>f)
31              #make predictions for nodes w/ hidden labels
32              predictions=np.zeros(len(hiddenNodes[0]))
33              #set A(i,j) to 0 when j is hidden (can still see A(j,i) to make predict
34              trainMatrix[:,hiddenNodes]=0
35              #store 'votes' for each vertex in seperate columns
```
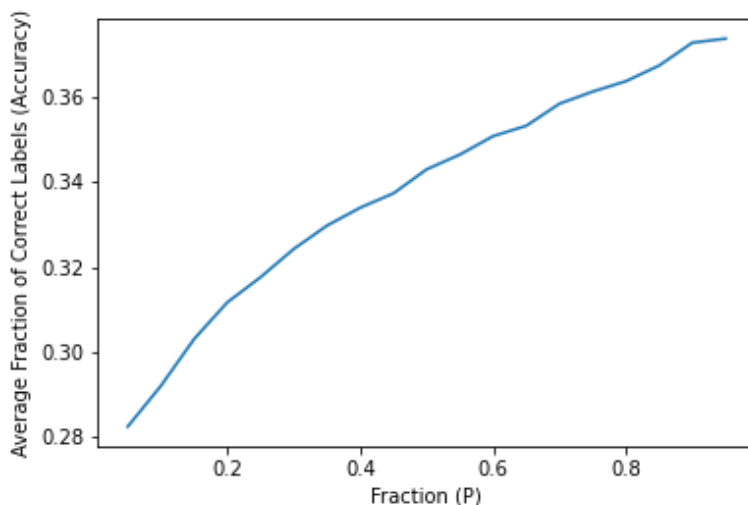
```
36              findMajority=np.zeros((len(hiddenNodes[0]),numberCategories))
37              for index in range(0,numberCategories):
38                  #neighbor vote total for each vertex/class
39                  findMajority[:,index]=((trainMatrix==index+1).sum(1))[hiddenNodes]
40                   #store predictions
41              predictions=np.zeros(len(hiddenNodes[0]))
42              #if tie (or no votes(tie of 0:0))
43              predictions[np.where(findMajority[:,0]==findMajority[:,1])]=randomLabel
44      #    print (findMajority,'\n',predictions)
45              #use majority to determine node class
46              predictions[np.where(predictions==0)]=(np.argmax(findMajority[np.where(
47              correct=float(np.sum(predictions==metadata[hiddenNodes]))
48      #    print('correct ',correct)
49      #    print('len(hiddenNodes[0]) ',len(hiddenNodes[0]))
50              iterationResults[iterationCounter]=correct/len(hiddenNodes[0])
51      #    print('iterationResults ',iterationResults)
52              iterationCounter+=1
53          #print('np.average(iterationResults)  ',np.average(iterationResults) )
54          #average accuracy of iterations over 1 f value
55          resultsOverF[fCounter]=np.average(iterationResults)
56          fValues[fCounter]=f
57          #print('resultsOverF[fCounter]  ',resultsOverF[fCounter] )
58          #print('fValues[fCounter]  ',fValues[fCounter] )
59          f = f + fStep
60          fCounter = fCounter + 1
61      #print(fValues)
62      #print(resultsOverF)
63      plt.plot(fValues,resultsOverF)
64      plt.xlabel('Fraction (P)')
65      plt.ylabel('Average Fraction of Correct Labels (Accuracy)')
66      plt.savefig('./{}{}Iterations.png'.format(networkFile[:-4],numberTrials))
67      plt.show()
68      np.savetxt('./{}{}accuracy.txt'.format(networkFile[:-4],numberTrials),resultsOv
69      np.savetxt('./{}{}fValues.txt'.format(networkFile[:-4],numberTrials),fValues)


 1 guiltByAssociation()
```

The above graph shows for a PubMed Diabetes network in which we predicted node labels.

For example, in a PubMed Diabetes network some users classified as 1 (Diabetes Mellitus, Experimental), 2 (Diabetes Mellitus Type 1), or 3 (Diabetes Mellitus Type 2) while others may not.

We used the network structure to predict the labels of nodes without labels. One of the most straightforward approaches to addressing this problem is known as the "guilt by association" heuristic. Here we predict the label of a node based on the most common label of its neighbors where ties are broken randomly. This tends to work well when the network structure is assortative with respect to the given label.

Here we considered the undirected version of the PubMed Diabetes network where nodes are classified as 1 (Diabetes Mellitus, Experimental), 2 (Diabetes Mellitus Type 1), or 3 (Diabetes Mellitus Type 2). For a given p between 0 and 1, then we picked a random fraction p of the nodes for which to observed labels.

Predicted labels for the remaining nodes using the guilt by association heuristic. Repeat this procedure 30 times for values of p ranging from 0.05 to 1 in increments of 0.05 and keep track of the average fraction of correct guesses for each p. We plotted a figure with p on the x-axis and average fraction of correct labels on the y-axis.

.

.

## Question 3

Consider the Facebook wall posts (2009) network from question 1 and imagine a piece of fake news spreading across the nodes. In this problem, we want to simulate this process under different conditions and observe the effect of "immunizing" nodes in different ways.

For R0 = 3, simulate the spread 30 times and keep track of the average fraction of infected nodes over time with no immunization.

Repeat the experiment with 10%, 30%, 50%, 70%, and 90% of the nodes immunized following three different strategies: random immunization, immunization of high degree nodes first, and neighbor immunization as described in class.

With the data collected from these experiments, generate three figures, one for each immunization strategy. Each figure will have time on the x-axis and It/n (the fraction of infected nodes) on the y-axis and five separate curves associated with the fraction of nodes immunized in each experiment.

What do you observe? Does one immunization strategy seem more effective than the others? Why or why not?

```
1 import networkx as nx
2 import numpy as np
3 import operator
4 import time
5 import matplotlib.pyplot as plt
6 np.random.seed(4)
```

### Parse Facebook Wall Post Data

```
1 with open('/content/out.facebook-wosn-wall') as f:
2     array=[]
3     for line in f:
4         array.append([int(x) for x in line.split()])
5 with open('/content/edgelist.tsv',"w")as f:
6     for a in array:
7         f.write(str(a[0])+' '+str(a[1])+'   '+str(int(a[3]/3600/24/30))+'\n')
```

.


.


.


.

**Print Data Staticstics**

```
1 def print_graph_stats(title, g):
2     print("Simple stats for: " + title)
3     print("# of nodes: " + str(len(g.nodes())))
4     print("# of edges: " + str(len(g.edges())))
5     print("Is graph connected? " + str(nx.is_connected(g)))
```

```
1 g = nx.read_adjlist("/content/edgelist.tsv")
2 g.remove_edges_from(nx.selfloop_edges(g))
3 print_graph_stats("Facebook Wall Post", g)
```

```
Simple stats for: Facebook Wall Post
# of nodes: 46952
# of edges: 431749
Is graph connected? True
```

**Create class for SIR Model**

Here

- g- Graph
- beta - Transmission Rate
- mu - Recovery rate
- Tmax = 30 times
- immunization_rate at start is zero

```
1 class SIR:
2     def __init__(self, g, beta, mu, Tmax = 30, index_start = 0.001):
3         self.g = g
4         self.beta = beta #transmission rate
5         self.mu = mu # recovery rate
6         self.index_start = index_start
7         self.Tmax = Tmax
8
9     def run(self, seed=[], num_steps = 1, sentinels = [], immunization_rate = 0.0,
10        # Immunize nodes according to the set immunization rate.
11        if len(immunized_nodes) == 0:
12            immunized = set(np.random.choice(self.g.nodes(),
13                                    size=int(immunization_rate*len(self.g.
14                                    replace=False))
15        else:
16            immunized = immunized_nodes
17
18        # If there is no seed, just choose a random node in the graph.
19        if len(seed) == 0:
20            number_of_person_depart = int(self.index_start*len(list(set(self.g.node
21            seed = np.random.choice(list(set(self.g.nodes()).difference(immunized))
22                                    number_of_person_depart, replace=False)
```

```python
23
24          I_set = set(seed)
25          S_set = set(self.g.nodes()).difference(I_set).difference(immunized)
26          R_set = set()
27
28          number_of_person_infected_sofar = {i:0.0 for i in self.g.nodes()}
29          number_of_time_people_stay_infected_sofar = {i:0.0 for i in self.g.nodes()}
30
31          t = 0
32
33          StoI = set(seed)
34          ItoR = set()
35
36          sentinels_t = {}
37          for sen in sentinels:
38              sentinels_t[sen] = 0
39
40          while (len(I_set) > 0 and t < self.Tmax):
41              I_set_old = I_set.copy()
42              # Let's infect people!
43              for i in I_set.copy():
44                  #print(len(set(self.g.neighbors(i)).intersection(S_set).copy()))
45                  ntot = len(set(self.g.neighbors(i)).intersection(S_set).copy())
46                  for s in set(self.g.neighbors(i)).intersection(S_set).copy():
47                      if np.random.uniform() < self.beta:
48                          S_set.remove(s)
49                          I_set.add(s)
50                          StoI.add(s)
51                          number_of_person_infected_sofar[i]+=1
52                          # Record t for sentinels
53                          if sentinels_t.get(s) != None:
54                              sentinels_t[s] = t
55
56                  number_of_time_people_stay_infected_sofar[i] += 1
57                  #print(t, i, number_of_time_people_stay_infected_sofar[i], number_o
58
59                  #print(t, number_of_person_infected_sofar[i] )
60
61
62                  # Will infected person spread fake news?
63                  if np.random.uniform() < self.mu:
64                      I_set.remove(i)
65                      R_set.add(i)
66                      ItoR.add(i)
67
68              t += 1
69              nbre_jour_rester_infecter = int(1/self.mu)
70              all_K = []
71              for k in I_set_old:
72                  val2 = max(1, nbre_jour_rester_infecter - number_of_time_people_sta
73                  val = min(len(list(self.g.neighbors(k))), number_of_person_infected
74                  #print(t, val, number_of_person_infected_sofar[k], len(list(self.g.
75                  all_K.append(val)
```

```
76              #print(t, np.mean(all_K))
77              #print(t, np.mean([number_of_person_infected_sofar[k]* for k in I_set_o
78                  #np.mean([number_of_time_people_stay_infected_sofar[k] for k in I
79          if t % num_steps == 0 or len(I_set) == 0:
80              yield({'t': t, 'S':S_set, 'I':I_set, 'R':R_set, 'StoI':StoI, 'ItoR'
81                      'reproductive_numbe':  np.mean(all_K)})
82
83
```

**No Immunization Strategy**

**For R0 = 3, simulate the spread 30 times and keep track of the average fraction of infected nodes over time with no immunization.**

```
 1 def get_temporal_plot(b, m0, Tmax, index_start, graph):
 2     m = 1./m0
 3     print()
 4     print("Constante Value of R : ", 1/m*b)
 5     print()
 6     sir = SIR(graph, beta = b, mu = m, Tmax=Tmax, index_start=index_start)
 7     res = []
 8     res2 = []
 9     final_rs = []
10     for r in sir.run(num_steps=1):
11         n= len(r['S'])+ len(r['I']) +len(r['R'])
12         res.append([len(r['I'])/n,1-len(r['R'])/n ])
13         res2.append([r['reproductive_numbe'], 1 ])
14     value_fin = len(r['R'])*100/len(graph.nodes())
15
16     return res, res2, value_fin
```
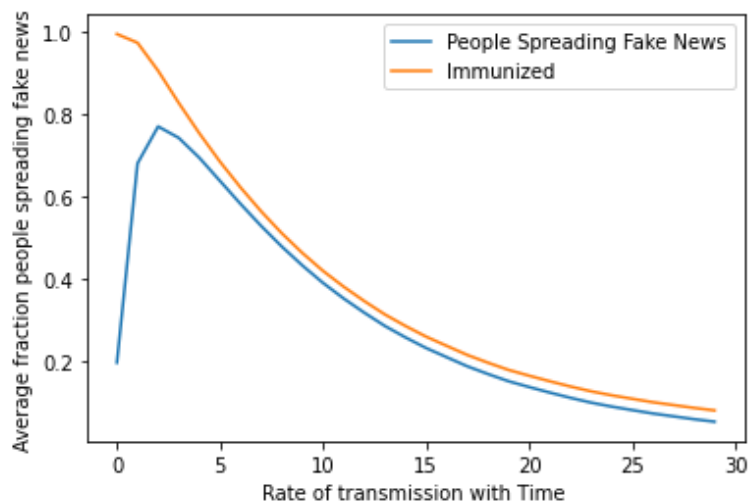
.

.

.

.

.

.

.

```
 1 b= 0.3
 2 m=10
 3 Tmax =30
 4 index_start=0.05
 5 res, res2, value_fin = get_temporal_plot(b,m,Tmax,index_start,g)
 6
 7 # Plotting the epidemic curve.
 8 plt.plot(res)
 9 # Plot the results.
10 plt.title("Epidemic curves for simulations with final "+str(value_fin)[:4] +
11             "% of people spreading fake news\n\n",fontsize =12)
12 plt.legend(['People Spreading Fake News', 'Immunized'])
13 plt.xlabel("Rate of transmission with Time ")
14 plt.ylabel("Average fraction people spreading fake news")
15 plt.show()
```

Constante Value of R :  3.0

Epidemic curves for simulations with final 91.9% of people spreading fake news

## Random Immunization

```
 1 def get_random_immunization(graph, b, m, Tmax, index_start, N, immunization_rates):
 2
 3     start = time.time()
 4     i_sir = SIR(graph, beta = b, mu = m, Tmax=Tmax, index_start=index_start)
 5     final_rs = {}
 6     final_rs0 = {}
 7
 8     for ir in immunization_rates:
 9         final_rs[ir] = []
10         final_rs0[ir] = []
11         for i in range(0,N):
12             simulation_steps = [[len(r['S']), len(r['I']), len(r['R']),
13                                 r['reproductive_numbe']] for r in i_sir.run(num_st
14             final_rs.get(ir).append(simulation_steps[len(simulation_steps)-1][2]*10
15             #print(simulation_steps[0][-1])
16             average_ro = []
17             for k in range(len(simulation_steps)):
18                 average_ro.append(simulation_steps[k][3])
19             final_rs0.get(ir).append(np.mean(average_ro))
20     sorted_ir = sorted(final_rs.items(), key=operator.itemgetter(0),reverse=True)
21     sorted_ir0 = sorted(final_rs0.items(), key=operator.itemgetter(0),reverse=True)
22
23     irs = []
24     oars = []
25     oars2 = []
26
27     for ir, values in sorted_ir:
28         irs.append(ir)
29         oars.append(np.mean(values))
30
31     return sorted(oars),sorted(irs)
```

```
 1
 2 m = 10 # Time of recovery
 3 index_start = 0.05 #Initialization percentage
 4 k =0.3 # Initialization ask population at random
 5 N =30 # Repeat train
 6 immunization_rates = [0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.05, 0.01]
 7 #10 percent
 8 b= 0.1 # Rate of transmission
 9 X1, Y1 = get_random_immunization(g, b, m, Tmax, index_start, N, immunization_rates
10
11 #30 percent
12 b = 0.3
13 X2, Y2 = get_random_immunization(g, b, m, Tmax, index_start, N, immunization_rates
14
15 #50 percent
16 b = 0.5
17 X3, Y3 = get_random_immunization(g, b, m, Tmax, index_start, N, immunization_rates
```
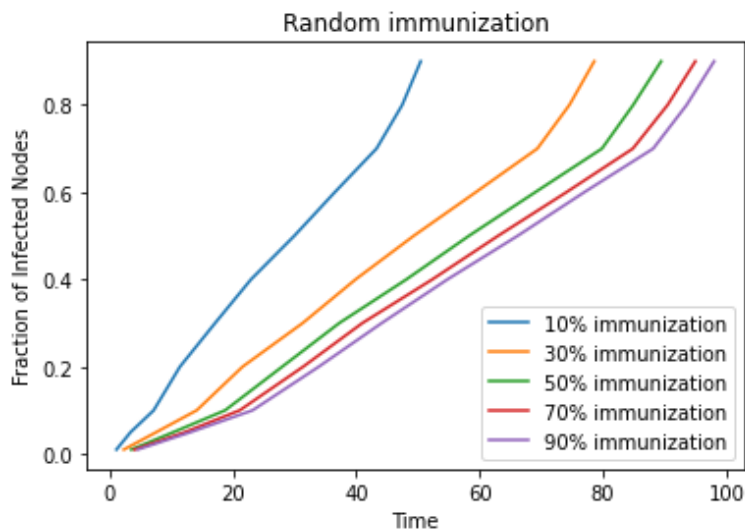
```
18
19 #70 percent
20 b = 0.7
21 X4, Y4 = get_random_immunization(g, b, m, Tmax, index_start, N, immunization_rates
22
23 #90 percent
24 b = 0.9
25 X5, Y5 = get_random_immunization(g, b, m, Tmax, index_start, N, immunization_rates
26
```

```
 1 plt.title('Random immunization')
 2 plt.xlabel('Time')
 3 plt.ylabel('Fraction of Infected Nodes')
 4 plt.plot(X1, Y1,label =" 10% immunization")
 5 plt.plot(X2, Y2,label =" 30% immunization")
 6 plt.plot(X3, Y3,label =" 50% immunization")
 7 plt.plot(X4, Y4,label =" 70% immunization")
 8 plt.plot(X5, Y5,label =" 90% immunization")
 9 plt.legend(['10% immunization', '30% immunization',
10              '50% immunization','70% immunization','90% immunization'])
11 plt.show()
```



## Immunization of high degree nodes first

```
 1
 2 def get_immunization_of_high_degree_nodes(graph, b, m, Tmax, index_start, N, immuni
 3     start = time.time()
 4
 5     ti_sir = SIR(graph, beta = b, mu = m, Tmax=Tmax, index_start=index_start)
 6     nodes_sorted_by_degree = sorted(nx.degree(graph), key=operator.itemgetter(1), r
 7     final_rs = {}
 8     final_rs0 = {}
 9     for ir in immunization_rates:
10         final_rs[ir] = []
11         final_rs0[ir] = []
12         # Immunize the M nodes with highest degree.
```

```
13          immunized_nodes = []
14          M = int(ir*len(nodes_sorted_by_degree))
15          for i in range(M):
16              immunized_nodes.append(nodes_sorted_by_degree[i][0])
17          # Run the simulation 50 times and save the results.
18          for i in range(0,N):
19              simulation_steps = [[len(r['S']), len(r['I']), len(r['R']), r["reproduc
20
21              final_rs.get(ir).append(simulation_steps[len(simulation_steps)-1][2]*10
22              average_ro = []
23              for k in range(len(simulation_steps)):
24                  average_ro.append(simulation_steps[k][3])
25              final_rs0.get(ir).append(np.mean(average_ro))
26      # Sort results and calculate the mean over the simulations to plot them.
27      #print("Job done in: ", time.time() - start)
28      sorted_ir = sorted(final_rs.items(), key=operator.itemgetter(0))
29      sorted_ir0 = sorted(final_rs0.items(), key=operator.itemgetter(0))
30      t_irs = []
31      t_oars = []
32      t_oars0 = []
33      for ir, values in sorted_ir:
34          t_irs.append(ir)
35          t_oars.append(np.mean(values))
36      for ir, values in sorted_ir0:
37          t_oars0.append(np.mean(values))
38
39      return sorted(t_oars), sorted(t_irs)
40


 1 #10 percent
 2 b= 0.1 # Rate of transmission
 3 m = 10 # Time of recovery
 4 index_start = 0.05 #Initialization percentage
 5 k =0.3 # Initialization ask population at random
 6 N =30 # Repeat train
 7 #immunization_rates = [0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
 8
 9 immunization_rates = [0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.05, 0.01]
10 X1, Y1 = get_immunization_of_high_degree_nodes(g, b, m, Tmax, index_start, N, immun
11
12 #30 percent
13 b = 0.3
14 X2, Y2 = get_immunization_of_high_degree_nodes(g, b, m, Tmax, index_start, N, immun
15
16 #50 percent
17 b = 0.5
18 X3, Y3 = get_immunization_of_high_degree_nodes(g, b, m, Tmax, index_start, N, immun
19
20 #70 percent
21 b = 0.7
22 X4, Y4 = get_immunization_of_high_degree_nodes(g, b, m, Tmax, index_start, N, immun
23
24 #90 percent
```
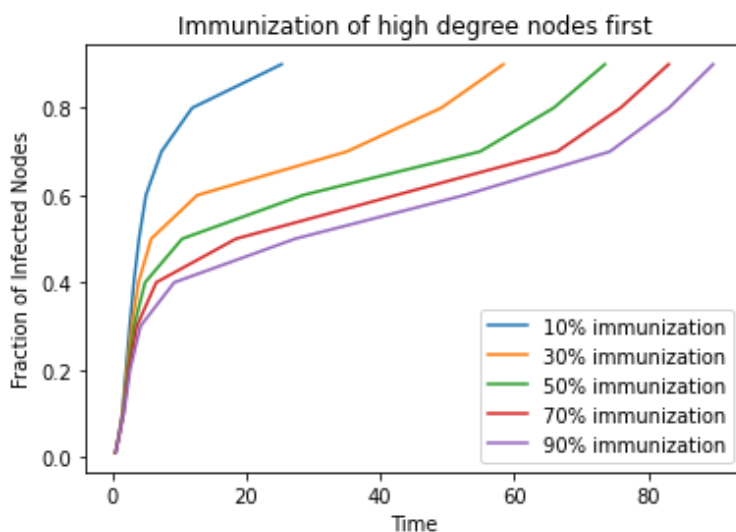
```
25 b = 0.9
26 X5, Y5 = get_immunization_of_high_degree_nodes(g, b, m, Tmax, index_start, N, immun
27
```

```
 1 plt.title('Immunization of high degree nodes first')
 2 plt.xlabel('Time')
 3 plt.ylabel('Fraction of Infected Nodes')
 4 plt.plot(X1, Y1,label =" 10% immunization")
 5 plt.plot(X2, Y2,label =" 30% immunization")
 6 plt.plot(X3, Y3,label =" 50% immunization")
 7 plt.plot(X4, Y4,label =" 70% immunization")
 8 plt.plot(X5, Y5,label =" 90% immunization")
 9 plt.legend(['10% immunization', '30% immunization',
10              '50% immunization','70% immunization','90% immunization'])
11 plt.show()
```



## Neighbor Immunization

```
 1 def get_neighbor_immunization(graph, b, m, Tmax, index_start, N, immunization_rates
 2     ti_sir = SIR(graph, beta = b, mu = m, Tmax=Tmax, index_start=index_start)
 3     K =0.2 # Init ask population at random
 4     sentinels = graph.nodes()
 5     sentinels_results = {}
 6     known_nodes = set(np.random.choice(graph.nodes(), size=int(K*len(graph.nodes()))
 7     neighbors = set()
 8     for node in list(known_nodes):
 9         neighbors.update(set(graph.neighbors(node)))
10     final_rs_k = {}
11     final_rs0 = {}
12     for ir in immunization_rates:
13         final_rs_k[ir] = []
14         final_rs0[ir] = []
15         M = int(ir*len(neighbors))
16         immunized_nodes_k = set(np.random.choice(list(neighbors), size=M, replace=F
17         for i in range(0,N):
18             # Acquaintance immunization
```

```
19          simulation_steps_k = [[len(r['S']), len(r['I']), len(r['R']),  r["repro
20
21              final_rs_k.get(ir).append(simulation_steps_k[len(simulation_steps_k)-1]
22              average_ro = []
23              for k in range(len(simulation_steps_k)):
24                  average_ro.append(simulation_steps_k[k][3])
25              final_rs0.get(ir).append(np.mean(average_ro))
26      sorted_ir_k = sorted(final_rs_k.items(), key=operator.itemgetter(0))
27      sorted_ir0 = sorted(final_rs0.items(), key=operator.itemgetter(0))
28      irs2 = []
29      oars_k = []
30      oars_deg = []
31      oars_sim = []
32      t_oars01 = []
33      for ir, values in sorted_ir_k:
34          irs2.append((ir*len(neighbors))/len(graph.nodes()))
35          oars_k.append(np.mean(values))
36      for ir, values in sorted_ir0:
37          t_oars01.append(np.mean(values))
38
39      return sorted(oars_k), sorted(irs2)
```

```
 1 #10 percent
 2 b= 0.1 # Rate of transmission
 3 m = 10 # Time of recovery
 4 index_start = 0.05 #Initialization percentage
 5 k =0.3 # Initialization ask population at random
 6 N =30 # Repeat train
 7 #immunization_rates = [0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
 8
 9 immunization_rates = [0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.05, 0.01]
10 X1, Y1 = get_neighbor_immunization(g, b, m, Tmax, index_start, N, immunization_rate
11
12 #30 percent
13 b = 0.3
14 X2, Y2 = get_neighbor_immunization(g, b, m, Tmax, index_start, N, immunization_rate
15
16 #50 percent
17 b = 0.5
18 X3, Y3 = get_neighbor_immunization(g, b, m, Tmax, index_start, N, immunization_rate
19
20 #70 percent
21 b = 0.7
22 X4, Y4 = get_neighbor_immunization(g, b, m, Tmax, index_start, N, immunization_rate
23
24 #90 percent
25 b = 0.9
26 X5, Y5 = get_neighbor_immunization(g, b, m, Tmax, index_start, N, immunization_rate
27
```
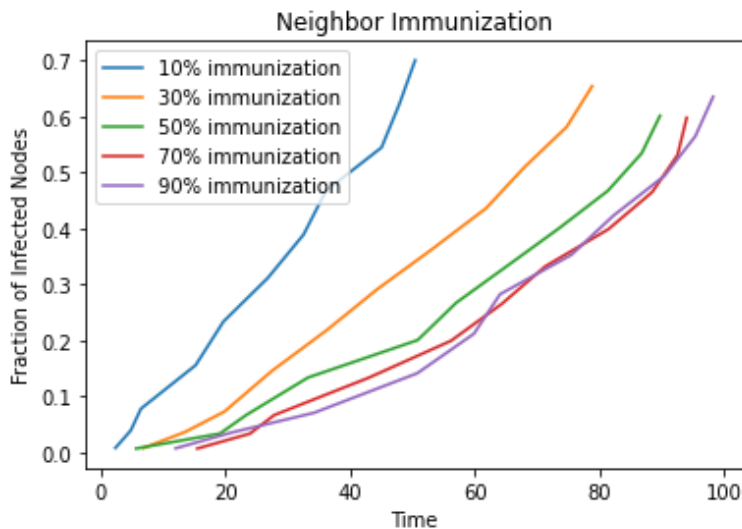
```
 1 plt.title('Neighbor Immunization')
```

```
 2 plt.xlabel('Time')
 3 plt.ylabel('Fraction of Infected Nodes')
 4 plt.plot(X1, Y1,label =" 10% immunization")
 5 plt.plot(X2, Y2,label =" 30% immunization")
 6 plt.plot(X3, Y3,label =" 50% immunization")
 7 plt.plot(X4, Y4,label =" 70% immunization")
 8 plt.plot(X5, Y5,label =" 90% immunization")
 9 plt.legend(['10% immunization', '30% immunization',
10               '50% immunization','70% immunization','90% immunization'])
11 plt.show()
```



Below are my observations from above graphs for different immunization strategies:

1. Random immunization is not an efficient strategy compared to other 2 strategies as irrespective of immunization percentage there's still a good chance that news might hit most of the population.

2. Targeted the Immunization of high degree nodes first (hubs of the networks) can work better when we know the graph which is not true in this case. If we are able to pinpoint most of the high degree nodes from the graph, news spread can be covered with minimal damage.

3. Neighbor Immunization gives good results and this strategy is purely local, requiring minimal information about randomly selected nodes and their immediate environment. In this scenario, despite of the fact that we are unaware of the graph, only 70% population was affected by the news.

4. Compared to all 3 strategies, I feel that High Degree Nodes First works best when we know the graph well but when it's not true, Neighbour immunization might be more efficient compared to the other two.

.

✓ 0s    completed at 5:31 PM                              ● ✕