```
Jupyter 42_Code5 Last Checkpoint: 8 minutes ago (autosaved)
                                                                                                                                           Logout
                                             Widgets
                                                       Help
                                                                                                                                       Python 3 O
                      Insert
 File
        Edit
              View
                              Cell
                                     Kernel
                                                                                                                         Not Trusted
                             ► Run ■ C → Code
↑ ↓
                                                           ~
       In [2]:
                 1 # Netflix
                 2 file1 = open(r'C:\Users\DELL\Downloads\Compressed\archive (1)\combined_data_1.txt', 'r', encoding="utf8")
                 3 Lines = file1.readlines()
                 4 print(len(Lines))
                24058263
                Preprocessing
       In [3]:
                 1 from tqdm import tqdm
                 2 movie_list= []
                 3 user_list = []
                 4 | time_list = []
                 5 count = 0
                 6 movie_id = ""
                 7 for i in tqdm(range(0, len(Lines))):
                        count += 1
                 9
                        if(Lines[i][-2:-1])==':':
                10
                            movie_id = int(Lines[i][:-2])
                11
                        else:
                12
                            movie_list.append(movie_id)
                            user_list.append(int(Lines[i].split(",")[0]))
                13
                            time_list.append(Lines[i].split(",")[2])
                14
                15
                16 print(len(user_list))
                17 print(len(movie_list))
                18 print(len(time_list))
                            24058263/24058263 [00:45<00:00, 533813.51it/s]
                24053764
                24053764
                24053764
                 1 print(type(movie_list[0]), movie_list[0])
       In [4]:
                 2 print(type(user_list[0]), user_list[0])
                 3 print(type(time_list[0]), time_list[0])
                <class 'int'> 1
                <class 'int'> 1488844
                <class 'str'> 2005-09-06
                 1 import networkx as nx
       In [ ]: |
                 2 G = nx.Graph()
      In [12]:
                 1 movie_degree_dict = {}
                 2 user_degree_dict = {}
                 3 probability_dict = {}
                 4 total_degree = 0
                 5 count = 0
                 6 for i in tqdm(range(0, len(movie_list))):
                        count += 1
                        node1 = movie_list[i]
                 8
                       node2 = user_list[i]
                 9
                       if node1 not in movie_degree_dict:
                10
                            movie_degree_dict[node1] = 0
                11
                        movie_degree_dict[node1] += 1
                12
                13
                        if node2 not in user_degree_dict:
                14
                15
                            user_degree_dict[node2] = 0
                        user_degree_dict[node2] += 1
                16
                17
                18
                        total_degree += 2
                        probability_dict[node2] = user_degree_dict[node2]/total_degree
                19
                20
                21 #
                          if count == 10:
                              print(movie_degree_dict)
                22 #
                              print(user_degree_dict)
                23 #
                              print(probability_dict)
                24 #
                25 #
                              break
                26
                                                             24053764/24053764 [00:32<00:00, 743010.40it/s]
       In [ ]:
      In [13]:
                 1 print(len(user_degree_dict.values()))
                 print(len(probability_dict))
                470758
                470758
                1 # for key, value in probability_dict.items():
      In [62]:
                          print(key, value)
                Network Evolution : Probability vs Destination Degree
      In [14]:
                 1 import matplotlib.pyplot as plt
                 2 import matplotlib
                 3 # plt.scatter(x_axis, y_axis, label = "fraction", color='blue')
                 4 plt.loglog(list(probability_dict.values()), list(user_degree_dict.values()), "o", label = "scaled frequencies")
                 5 # plt.scatter(list(probability_dict.values()), list(movie_degree_dict.values()), label = "scaled frequencies")
                 6 # plt.xscale("log")
                 7 # plt.yscale("log")
                 8 plt.xlabel('probability')
                10 plt.ylabel('degree')
                11
                12 plt.title('Evolution')
                13 plt.legend()
                14 # plt.show()
                15 fig = matplotlib.pyplot.gcf()
                16 fig.set_size_inches(10, 7)
                                                      Evolution

    scaled frequencies

                   10^{3}
                   10<sup>1</sup>
                   10°
                                10^{-7}
                                             10<sup>-6</sup>
                                                          10-5
                                                                      10^{-4}
                                                                                   10^{-3}
                                                      probability
       In [5]:
                 1 import networkx as nx
                 2 G = nx.Graph()
       In [6]:
                 1 for edge in tqdm(range(len(user_list))):
                        node1=user_list[edge]
                        node2=movie_list[edge]
                 3
                        G.add_edge(node1,node2)
                 5 Degree_count_from={}
                 6 Degree_count_to={}
                 7 Degree_from={}
                 8 Degree_to={}
                 9 user_list_set=set(user_list)
                10 movie_list_set=set(movie_list)
                11
                12 for x in user_list_set:
                # if(x not in Degree_count_from):
                          Degree_from[x]=G.degree(x)
                14
                     # else:
                15
                     # print("hi")
                16
                      # break
                17
                18
                19 | for y in movie_list_set:
                     # if(x not in Degree_count_to):
                          Degree_to[y]=G.degree(y)
                21
                22
                                                               24053764/24053764 [01:56<00:00, 207249.97it/s]
                Degree Distribution : User
       In [9]:
                 1 list_of_degrees=list(Degree_from.values())
                 2 from_degree_counts={}
                 3 for deg in list_of_degrees:
                        if(deg in from_degree_counts):
                 5
                            from_degree_counts[deg]+=1
                 6
                        else:
                            from_degree_counts[deg]=1
                 8 max_from_degree_freq=max(from_degree_counts.values())
                10
                11 from_degree_counts1=from_degree_counts.copy()
                13 from_distis=dict((sorted(from_degree_counts1.items())))
                14 from_max_degree_count=max(from_distis.values())
                15 # total_nodes=sum(distis.values())
                16 # from_distis.update((x, round(y/from_max_degree_count,2)) for x, y in from_distis.items())
                17 # x_axis=list(from_distis.keys())[1:len(from_distis.keys())]
                18 # y_axis=list(from_distis.values())[1:len(from_distis.values())]
                20 x_axis=list(from_distis.keys())
                21 y_axis=list(from_distis.values())
                23 import matplotlib.pyplot as plt
                24 import matplotlib
                25 plt.loglog(x_axis, y_axis,"o")
                26 # plt.xscale("log")
                27 # plt.yscale("log")
                28 plt.xlabel('degree')
                30 plt.ylabel('number of nodes having degree')
                32 plt.title('Degree Distribution of User')
                33 plt.legend()
                34 # plt.show()
                35 fig = matplotlib.pyplot.gcf()
                No handles with labels found to put in legend.
                                               Degree Distribution of User
                   10^{4}
              number of nodes having degree
                   10<sup>1</sup>
                   10°
                                    10<sup>1</sup>
                                                 10^{2}
                                                                          10<sup>4</sup>
                                                                                       105
                        10°
                                                              10<sup>3</sup>
                                                        degree
                Degree Distribution: Movies
      In [10]:
                 1 list_of_degrees=list(Degree_to.values())
                 2 to_degree_counts={}
                 3 for deg in list_of_degrees:
                        if(deg in to_degree_counts):
                            to_degree_counts[deg]+=1
                 6
                        else:
                            to_degree_counts[deg]=1
                    max_to_degree_freq=max(to_degree_counts.values())
                10
                11 to_degree_counts1=to_degree_counts.copy()
                13 to_distis=dict((sorted(to_degree_counts1.items())))
                14 | to_max_degree_count=max(to_distis.values())
                15 # total nodes=sum(distis.values())
                # to_distis.update((x, round(y/to_max_degree_count,2)) for x, y in to_distis.items())
                17 # x_axis=list(to_distis.keys())[1:len(to_distis.keys())]
                18 # y_axis=list(to_distis.values())[1:len(to_distis.values())]
                20 x_axis1=list(to_distis.keys())
                21 y_axis1=list(to_distis.values())
                23 import matplotlib.pyplot as plt
                24 import matplotlib
                25 # plt.scatter(x_axis1, y_axis1, label = "fraction")
                26 plt.loglog(x_axis1, y_axis1,"o")
                27 # plt.xscale("log")
                28 # plt.yscale("log")
                29 plt.xlabel('degree')
                31 plt.ylabel('number of nodes having degree')
                32
                33 plt.title('Degree Distribution of Movies')
                34 plt.legend()
                35 # plt.show()
                36 | fig = matplotlib.pyplot.gcf()
                No handles with labels found to put in legend.
                                              Degree Distribution of Movies
                   10<sup>1</sup>
```

10°

10<sup>2</sup>

 $10^{3}$ 

degree

 $10^{4}$ 

10<sup>5</sup>