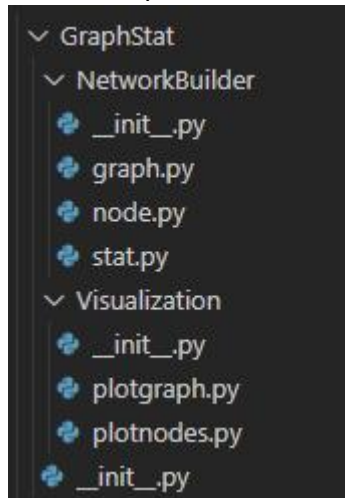


参照作业 pdf 建立包 GraphStat 的结构



本报告按照完成代码的时间逻辑排序

一、读取文件并完成序列化

为了避免每次都要重新读取两个文件并生成图这一繁琐的过程,我们先将两个文件读入并生成图,对其进行序列化:

Step1.读取 features 文件 (node.py)

```
print("-----正在导入节点属性数据-----")
f=open(filename_features,encoding='UTF-8')
title = f.readline().strip().split(',') #读取第一行
txt_features=[]
line=[1]
while line: # 直到读取完文件
    line = f.readline().strip() # 读取一行文件, 包括换行符
    txt_features.append(line)
f.close() # 关闭文件
if txt_features[-1]=='':
    txt_features.pop()
list_features=[]
for i in txt_features:
    list_features.append(i.split(','))
#print(list_features)
print("-----数据导入完成-----")
#print(title)
```

Step2.将读取的信息转换为列表 (node.py)

```
23 dict_id={} #用户属性字典
24 for item in list_features:
25     dict_index = {'views': "NULL", 'mature': "NULL", 'life_time': "NULL",
26                 'created_at': "NULL", 'updated_at': "NULL", 'dead_account': "NULL", 'language': "NULL", 'affiliate': "NULL"}
27     for i in range(9):
28         dict_index[title[i]] = item[i]
29     dict_id[item[5]] = dict_index
30 #print(len(dict_id))
31 #print(dict_id)
32 return dict_id
```

Step3.读取 edges 文件 (graph.py)

```

2 import networkx as nx
3 import pickle
4
5 def init_graph(filename_edges,dict_id):# (可以考虑用networkx中的 Graph 等。)
6     """
7     构建网络
8     """
9     print("-----正在导入边关系数据-----")
10    f=open(filename_edges,encoding='UTF-8')
11    title = f.readline().strip().split(',') #读取第一行
12    txt_edges=[]
13    line=[1]
14    while line: #直到读取完文件
15        line = f.readline().strip() #读取一行文件, 包括换行符
16        txt_edges.append(line)
17    f.close() # 关闭文件
18    if txt_edges[-1]=='':
19        txt_edges.pop()
20    list_edges=[]
21    for i in txt_edges:
22        list_edges.append(i.split(','))
23    #print(list_edges)
24    print("-----边关系数据导入完成-----")

```

Step4.使用 networkx 库生成图

```

25 G = nx.Graph()#生成一个空白图
26 for id in dict_id:
27     for key in dict_id[id]:
28         G.add_node(id) #插入点
29         G.nodes[id][key]=dict_id[id][key] #添加属性
30 G.add_edges_from(list_edges) #连接节点
31 return G

```

该图包含了所有节点的信息, 只需要将该图序列化以后可直接利用。

Step5.将生成的图序列化

```

33 def save_graph(G,filename_save):
34     """
35     序列化图信息
36     图的信息暂时均以字符串形式存储
37     """
38     f=open(filename_save,"wb")
39     pickle.dump(G,f)
40     f.close()

```

二、检验序列化并完成 NetworkBuilder 的其他功能

Step1.

反序列化

```

41
42 def load_graph(filename_save):
43     """
44     将网络加载至内存
45     """
46     f=open(filename_save,"rb")
47     result =pickle.load(f)
48     f.close()
49     return result

```

Step2.

实现其他 NetworkBuilder 的功能

① (node.py)

```

34 def get_features(G,id,features):
35     """
36     获取节点node的指定属性。
37     属性包括views,mature,life_time,created_at,updated_at,dead_account,language和affiliate
38     """
39     lis_features=['views','mature','life_time','created_at','updated_at','dead_account','language','affiliate']
40     if features in lis_features:
41         return G.nodes[id][features]
42     else:
43         return -1
44
45 def print_node(G,id):
46     """
47     显示节点全部信息（利用 format）
48     """
49     print("id为{id}的全部信息为{1}".format(id,G.nodes[id]))

```

②一些简单的图的参数函数（过于简单的函数直接使用 networkx 自带的函数较为方便）（stat.py）

```

1  #-*- coding: gbk -*-
2  import networkx as nx
3
4  def get_node_number(G):
5      """
6      计算节点数
7      """
8      number=nx.number_of_nodes(G)
9      return number
10
11 def get_edge_number(G):
12     """
13     计算边数
14     """
15     number=nx.number_of_edges(G)
16     return number
17
18 def cal_average_degree(G):
19     """
20     计算网络中的平均度
21     """
22     sum=0
23     for i in G.nodes():
24         sum+=G.degree(i)
25     return sum/len(G.node)

```

```

36 def cal_views_distribution(G):
37     """
38     计算 views 属性的分布
39     """
40     view_dic={}
41     for id in G.nodes():#生成字典
42         view_dic[G.nodes[id]["views"]]=view_dic.get(G.nodes[id]["views"],0)+1
43     sum=0
44     for key in view_dic.keys():
45         sum+=view_dic[key]
46     for key in view_dic.keys():
47         view_dic[key]=view_dic[key]/sum
48     return view_dic

```

度的分布与 views 的分布类似，可以用相同的字典计数法求得，但 nx 库存在已有函数，故直接调用

```

28 def cal_degree_distribution(G):
29     """
30     计算网络的度分布
31     """
32     degree_dis=nx.degree_histogram(G)#统计从0到最大度的频次
33     lis = [z / float(sum(degree_dis)) for z in degree_dis]#生成密度列表
34     return lis
35

```

就此 NetworkBuilder 库的内容已经完成，下面尝试在主程序中调用。

三、在主程序中尝试调用已完成的 NetworkBuilder 包进行序列化

```

1  #-*- coding: gbk -*-
2  import GraphStat.NetworkBuilder.graph as gng
3  import GraphStat.NetworkBuilder.node as gnn
4  import GraphStat.NetworkBuilder.stat as gns
5
6  filename_edges=r'C:\Users\LF\Desktop\twitch_gamers\large_twitch_edges.csv'
7  filename_features=r'C:\Users\LF\Desktop\twitch_gamers\large_twitch_features.csv'
8  filename_save=r'C:\Users\LF\Desktop\twitch_gamers\save.txt'
9
10 dict_id=gnn.init_node(filename_features)
11 G=gng.init_graph(filename_edges,dict_id)
12 gng.save_graph(G,filename_save)
13

```

进行反序列化并输出信息

```

14 G=gng.load_graph(filename_save)
15 gnn.print_node(G,'5') #随便输出一个id的信息
16 print("总结点数为%d"%gns.get_node_number(G))
17 print("总边数为%d"%gns.get_edge_number(G))
18 print("网络中的平均度为%d"%gns.cal_average_degree(G))
19

```



```

PS E:\code\py_code> python -u "e:\code\py_code\week4\main.py"
id为5的全部信息为{'views': '4987', 'mature': '1', 'life_time': '1288', 'created_at': '2015-04-03', 'updated_at': '2018-10-12', 'dead_account': '0', 'language': 'EN', 'affiliate': '1', 'numeric_id': '5'}
总边数为168114.
总节点数为6797557.
网络中的平均度为80

```

四、图的可视化

Step1.绘制节点的局部网络(plotgraph.py)

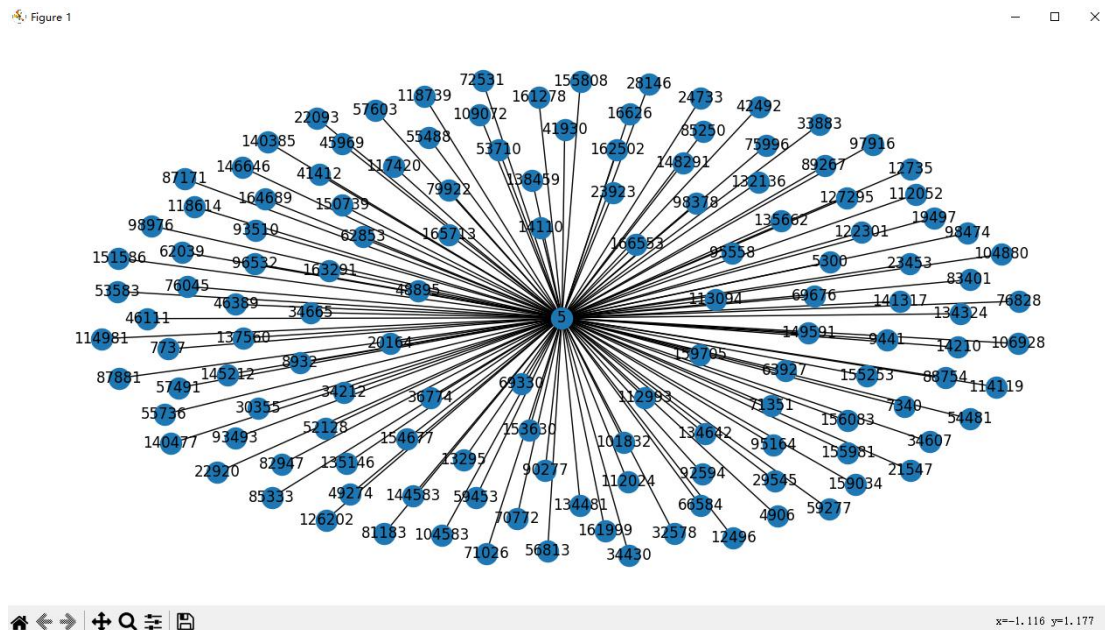
选择一个节点构建一个新的图

```

6 def plot_ego(G,node):
7     """
8     绘制节点的局部网络
9     """
10    Q=nx.Graph()
11    Q.add_node(node)
12    lis=G.nodes()
13    for i in lis:
14        if G.has_edge(node,i):
15            Q.add_node(i)
16            Q.add_edge(i,node)
17    nx.draw(Q,with_labels=True)
18    plt.show()

```

下面以 id 为 5 的节点为例输出结果



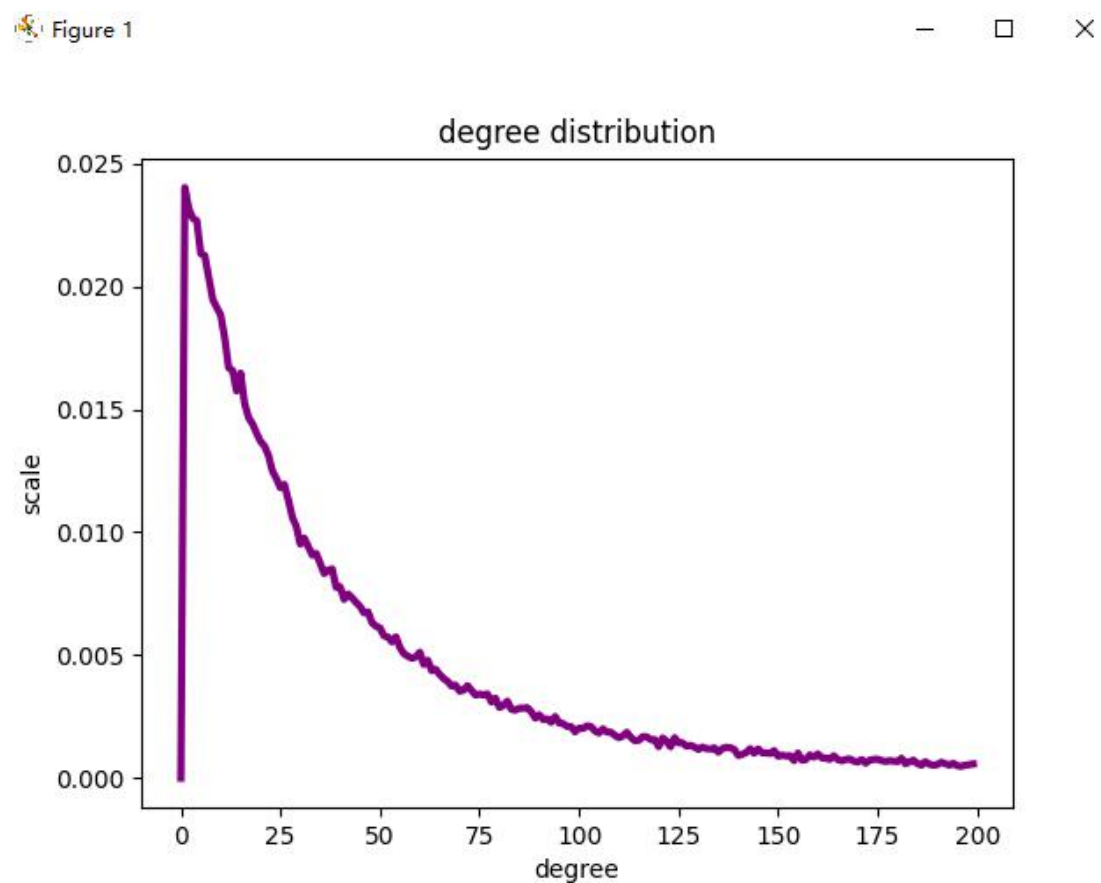
Step2.画出度的分布图

```

21 def plot_degree_distribution(degree_dis): # (观察度分布的形态)
22     """
23     度的分布图
24     由于节点分布过于分散, 仅仅画出深度为0~200的分布
25     """
26     x=[i for i in range(200)]
27     y=degree_dis[0:200]
28     plt.plot(x,y,color="purple",linewidth=3.0)
29     plt.title("degree distribution")
30     plt.xlabel("degree")
31     plt.ylabel("scale")
32     plt.show()

```

由于节点的度分布过于分散，故只画出深度为 0~200 的分布结果如下：



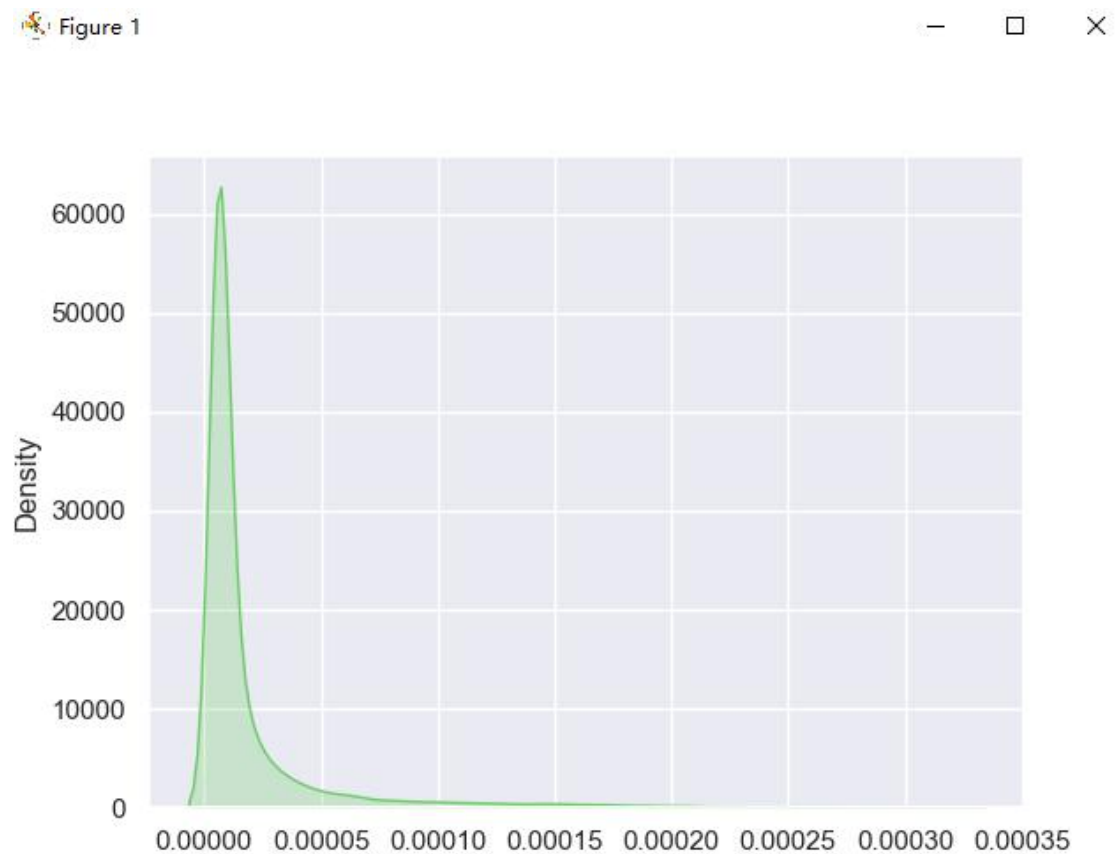
可以发现是一个近似反比例函数的图像（将点（0，0）排除在外）

Step3.观察属性的分布形态

```
1  #-*- coding: gbk -*-
2  ~ import seaborn as sns
3  import networkx as nx
4  import GraphStat.NetworkBuilder.stat as gns
5  import matplotlib.pyplot as plt
6
7  ~ def plot_degree_distribution(G,features):
8      ~ """
9      ~ 观察属性的分布形态
10     ~ """
11     node=G.nodes()
12     lis=['views', 'mature', 'life_time', 'created_at', 'updated_at', 'numeric_id', 'dead_account', 'language', 'affiliate']
13     ~ if features in lis:
14         lis1=[]
15         lis2=[0,50000,100000,150000]
16         dic=gns.cal_distribution(G,features)
17         ~ for i in dic:
18             lis1.append(dic[i])
19         sns.set(palette='muted',color_codes=True)
20         sns.distplot(lis1,lis2,hist=False,color='g',kde_kws={'shade':True})
21         plt.show()
22
```

利用特征分布函数画出对应属性的分布图

下面以‘views’属性为例



下面是在主程序的代码

```
1  -*- coding: gbk -*-
2  import GraphStat.NetworkBuilder.graph as gng
3  import GraphStat.NetworkBuilder.node as gnn
4  import GraphStat.NetworkBuilder.stat as gns
5  import GraphStat.Visualization.plotgraph as gvpkg
6  import GraphStat.Visualization.plotnodes as gvpn
7  import random
8
9  filename_edges=r'C:\Users\LF\Desktop\twitch_gamers\large_twitch_edges.csv'
10 filename_features=r'C:\Users\LF\Desktop\twitch_gamers\large_twitch_features.csv'
11 filename_save=r'C:\Users\LF\Desktop\twitch_gamers\save.txt'
12
13 #第一部分：输入数据与序列化
14 dict_id=gnn.init_node(filename_features)
15 G=gng.init_graph(filename_edges,dict_id)
16 gng.save_graph(G,filename_save)
17
18 #第二部分：反序列化与属性输出
19 G=gng.load_graph(filename_save)
20 gnn.print_node(G,'5') #随便输出一个id的信息
21 print("总结点数为%d"%gns.get_node_number(G))
22 print("总边数为%d"%gns.get_edge_number(G))
23 print("网络中的平均度为%d"%gns.cal_average_degree(G))
24
25 #第三部分：图的可视化
26 gvpkg.plot_ego(G,str(random.randint(0,10)))
27 degree_dis = gns.cal_degree_distribution(G)
28 gvpkg.plot_degree_distribution(degree_dis)
29 gvpn.plot_degree_distribution(G,'views')
```

序列化和反序列化 <https://blog.csdn.net/reuxfhc/article/details/80043320>

networkx 中文学习手册 <https://blog.csdn.net/pylittlebrat/article/details/121906479>