山东大学<u>计算机科学与技术</u>学院 大数据分析与实践课程实验报告

学号: 202300130063 姓名: 邱珺 班级: 23 数据

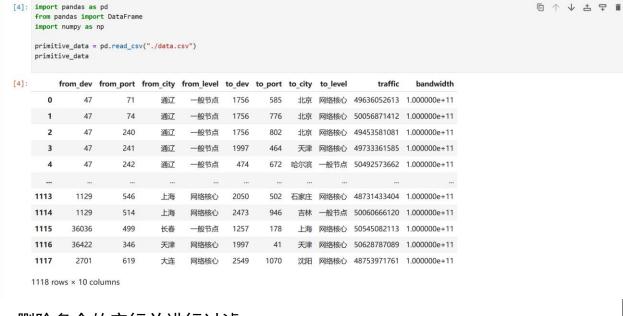
实验题目:实验一 数据采样方法实践

实验学时: 2 实验日期: 2025/9/8

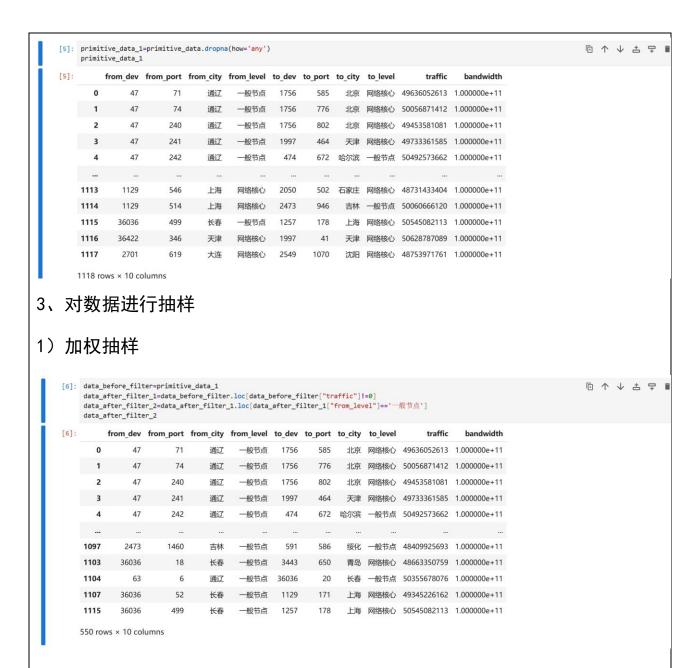
实验目标:利用 Pandas 库实现多种数据采样和过滤的方法

实验过程:

1、库的导入与数据的读入



2、删除多余的空行并进行过滤



2) 随机抽样

```
[11]: random_sample=data_before_sample
                                                                                                                     □ ↑ ↓ 占 무 🗎
         random sample finish=random sample.sample(n=50)
         random_sample_finish=random_sample_finish[columns]
         random_sample_finish
   [11]:
              from dev from port from city from level to dev to port
                                                                                    traffic
                                                                 to city to level
                                                                                            bandwidth
          111
                   787
                             54
                                    玉潔
                                           一般节点
                                                     171
                                                            422
                                                                  哈尔滨 一般节点 50571503467 1.000000e+11
          313
                            111
                                 呼和浩特
                                           一般节点
                                                    2360
                                                            197
                                                                   太原 网络核心 49309667295 1.000000e+11
         1057
                    47
                            243
                                    诵订
                                           一般节点
                                                    2473
                                                            769
                                                                   吉林
                                                                       一般节点 49117847542 1.000000e+11
          125
                   474
                           1374
                                   哈尔滨
                                           一般节点
                                                    2050
                                                            336
                                                                 石家庄 网络核心 50242784823 1.000000e+11
          706
                  2473
                            799
                                    吉林
                                           一般节点
                                                                鄂尔多斯
                                                                       网络核心 49550894885
                                                                                          1.000000e+11
          953
                                 呼和浩特
                                           一般节点
                                                     47
                                                            249
                                                                   通辽 一般节点 50233070000 1.000000e+11
                   180
                            192
           70
                   180
                             36
                                 呼和浩特
                                           一般节点
                                                    2194
                                                            406
                                                                   唐山 网络核心 50973267302 1.000000e+11
         1035
                  36036
                             54
                                           一般节点
                                                    591
                                                            23
                                                                   绥化 一般节点 50638071722 1.000000e+11
                                    长春
          416
                   591
                             60
                                    绥化
                                           一般节点
                                                     180
                                                             52 呼和浩特 一般节点 50126205393 1.000000e+11
                   591
                            586
                                           一般节点
                                                            86 劉尔多斯 网络核心 47929885030 1.000000e+11
         1093
                                    绥化.
                                                    1536
         1079
                    63
                            224
                                    诵辽
                                           一般节点
                                                    4069
                                                           1196
                                                                   宁波 一般节点 50209459772 1.000000e+11
          494
                    47
                            252
                                    诵辽
                                           一般节点
                                                    1536
                                                            86 鄂尔多斯 网络核心 50478868327 1.000000e+11
          323
                    96
                            141
                                 呼和浩特
                                           一般节点
                                                    2050
                                                            391
                                                                 石家庄 网络核心 49814111100 1.000000e+11
          145
                   591
                            100
                                    绥化
                                          一般节点
                                                    2194
                                                            506
                                                                   唐山 网络核心 51437026945 1.000000e+11
          329
                            159
                                 呼和浩特
                                           一般节点
                                                           1088
                                                                              51159730271 1.000000e+11
          317
                                                            263
                                                                   太原 网络核心 50500915133 1.000000e+11
                    96
                            123
                                 呼和浩特
                                           一般节点
                                                    2360
          674
                   591
                            586
                                    绥化
                                           一般节点
                                                     47
                                                            243
                                                                   通辽 一般节点 50565152517 1.000000e+11
          322
                    96
                            136
                                 呼和浩特
                                           一般节点
                                                    3227
                                                            389
                                                                  济南 网络核心 50541979348 1.000000e+11
          1075
                           1196
                                    宁波
                                           一般节点
                                                    1756
                                                           1187
                                                                   北京 网络核心 50488255524 1.000000e+11
                                                    ._.
                                                           ---
                                                                 ..... .... ..... ..... ..... ...
3) 分层抽样
```

```
after_sample=pd.concat([ybjd.sample(17),wlhx.sample(33)])
[12]:
          from_dev from_port from_city from_level to_dev to_port to_city to_level
                                                                                traffic
                                                                                        bandwidth
      491
                                通辽
                                      一般节点
                                                               杭州
                                                                   一般节点 50888438116 1.000000e+11
      750
               474
                        472
                               哈尔滨
                                       一般节点
                                                96
                                                       157 呼和浩特 一般节点 50784665266
                                                                                     1.000000e+11
      556
                63
                        282
                                诵辽
                                      一般节点
                                                 63
                                                         6
                                                               通辽 一般节点 49489299594 1.000000e+11
      110
               474
                        672
                              哈尔滨
                                     一般节点
                                              47
                                                       242
                                                               通辽 一般节点 51555817613 1.000000e+11
      827
               474
                        422
                                      一般节点
                                                             哈尔滨 一般节点 49998657939 1.000000e+11
                               哈尔滨
                                                474
                                                       1410
      656
              4069
                       1196
                                宁波
                                       一般节点
                                                180
                                                       264 呼和浩特 一般节点 49766912004 1.000000e+11
        5
                47
                        243
                                通辽
                                       一般节点
                                                        124 呼和浩特
                                                                   一般节点 49942713747 1.000000e+11
     1097
              2473
                       1460
                                吉林
                                       一般节点
                                                591
                                                       586
                                                              绥化 一般节点 48409925693 1.000000e+11
      347
               180
                         42
                             呼和浩特
                                      一般节点
                                                4360
                                                       406
                                                               南京 一般节点 50178810628 1.000000e+11
      546
                63
                         60
                                通辽
                                      一般节点
                                                4360
                                                       468
                                                              南京 一般节点 47970715088 1.000000e+11
      354
               180
                        192
                             呼和浩特
                                       一般节点
                                                4360
                                                       271
                                                               南京 一般节点 51828297117 1.000000e+11
      804
               180
                                                474
                                                       475
                                                             哈尔滨 一般节点 49012460413 1.000000e+11
                        264
                             呼和浩特
                                       一般节点
      289
                47
                        417
                                诵订
                                       一般节点
                                                3615
                                                       191
                                                               长沙 一般节点 50099712071 1.000000e+11
              2473
                                                180
      764
                        941
                                吉林
                                       一般节点
                                                        26 呼和浩特 一般节点 49660872427 1.000000e+11
                47
                        241
                                                4953
      278
                                诵订
                                      一般节点
                                                       725
                                                               豊阳 一般节点 50008939996 1.000000e+11
      861
                47
                        417
                                通辽
                                      一般节点
                                                591
                                                       1284
                                                               绥化 一般节点 49276967001 1.000000e+11
      376
                        460
                                       一般节点
                                                               福州 一般节点
                                                                          48394911971 1.000000e+11
                47
                         74
                                通辽
                                      一般节点
                                                1756
                                                       776
                                                              北京 网络核心 50056871412 1.000000e+11
      531
                47
                        250
                                诵辽
                                      一般节点
                                                3227
                                                       195
                                                               济南 网络核心 48844966451 1.000000e+11
       26
                63
                         74
                                通辽
                                      一般节点
                                                2701 181 大连 网络核心 50364636480 1.000000e+11
```

□ ↑ ↓ 占 무 ▮

4) 补充: 系统抽样

[12]: ybjd=data_before_sample.loc[data_before_sample['to_level']=='一般节点']

wlhx=data_before_sample.loc[data_before_sample['to_level']=='网络核心']

```
[9]: systematic_sample = data_before_sample.copy()
                                                                                                  ⑥↑↓古♀▮
    n = 50
   r = np.random.randint(0, k) # 随机起点
    indices = np.arange(r, r + n*k, k)[:n]
   systematic_sample_finish = systematic_sample.iloc[indices]
systematic_sample_finish
                                                                            bandwidth
        from_dev from_port from_city from_level to_dev to_port to_city to_level
                                                                    traffic
            47
                    74
                           诵订
                                一般节点 1756
                                              776
                                                     北京 网络核心 50056871412 1.000000e+11
     1
                    260
     12
            47
                        通辽 一般节点 2549 835 沈阳 网络核心 50220958279 1.000000e+11
     23
                           诵订
                                一般节点 36422
                                                     天津 网络核心 50322780029 1.000000e+11
                                               560 上海 网络核心 49753614568 1.000000e+11
     34
            96
                    99
                        呼和浩特 一般节点 1257
     45
             96
                    134
                        呼和浩特
                                一般节点
                                        47
                                               252
                                                     通辽 一般节点 49416652053 1.000000e+11
                    346
                        呼和浩特
                                一般节点
                                        1257
                                                     上海 网络核心 47759033178 1.000000e+11
     67
            180
                                一般节点 1385
                    28
                        呼和浩特
                                               133
                                                     广州 网络核心 52798223188 1.000000e+11
                  188
     78
            180
                        呼和浩特 一般节点 36422
                                               350 天津 网络核心 49047066099 1.000000e+11
     89
            180
                               一般节点 1129
                                                     上海 网络核心 49512421445 1.000000e+11
                        呼和浩特
    102
            474
                    467
                       哈尔滨 一般节点 1257
                                               174 上海 网络核心 49987703744 1.000000e+11
    115
            474
                    683
                         哈尔定 一般节点 1997
                                               84
                                                     天津 网络核心 49446798762 1.000000e+11
    126
            474
                   1389
                          哈尔滨
                                一般节点
                                        1756
                                               1127
                                                     北京 网络核心 48259332712 1.000000e+11
    137
                           绥化
                                一般节点
                                               6
                                                     通辽 一般节点 49462640634 1.000000e+11
    148
           591
                   558
                          绥化 一般节点 36036
                                               499
                                                     长春 一般节点 49953028308 1.000000e+11
    162
            591
                   1274
                           绥化. 一般节点 1129
                                              203
                                                     上海 网络核心 49321244844 1.000000e+11
            787 307 王宮 一般节点 4953 686 無阳 一般节点 4939787960 1,000000+11
```

先统计数据总数和需要的样本数量, 计算出抽样间隔; 然后在前若干个数据里随机选一个起始点; 最后按照这个起点, 每隔固定数量依次抽取, 直到得到所需样本。

5) 补充:整群抽样



按照 to_level 分为若干个群,再随机抽取一个群,显示结果

结论分析:

本实验围绕数据采样方法展开,旨在通过实际操作掌握数据清洗与多种抽样技术的基本流程。实验中首先利用 Pandas 对原始数据进行预处理,删除空行并过滤掉无效记录,从而得到符合条件的数据集。在此基础上,依次完成了加权抽样、随机抽样和分层抽样等方法的实践,对比了不同抽样策略在样本分布上的差异。实验进一步扩展到系统抽样和整群抽样,以便全面理解各种抽样方式的适用场景与特点。通过本实验的完成,加深了对数据预处理和抽样在大数据分析中重要性的理解,为后续实验的深入学习奠定了基础。