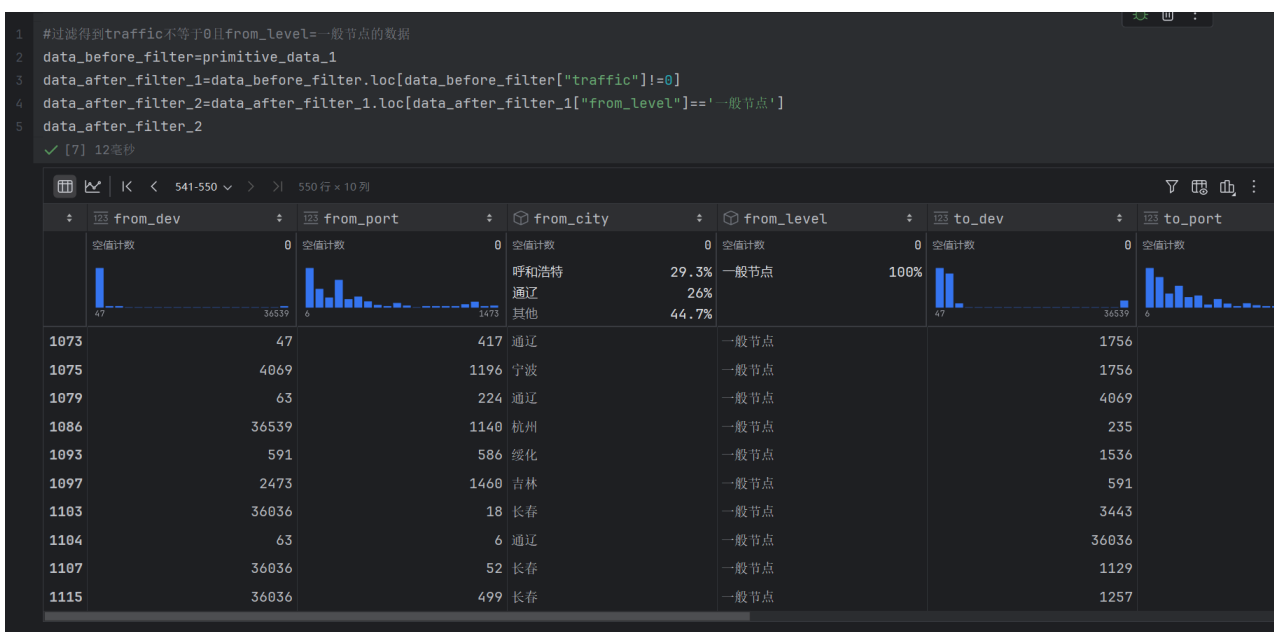


山东大学计算机科学与技术学院

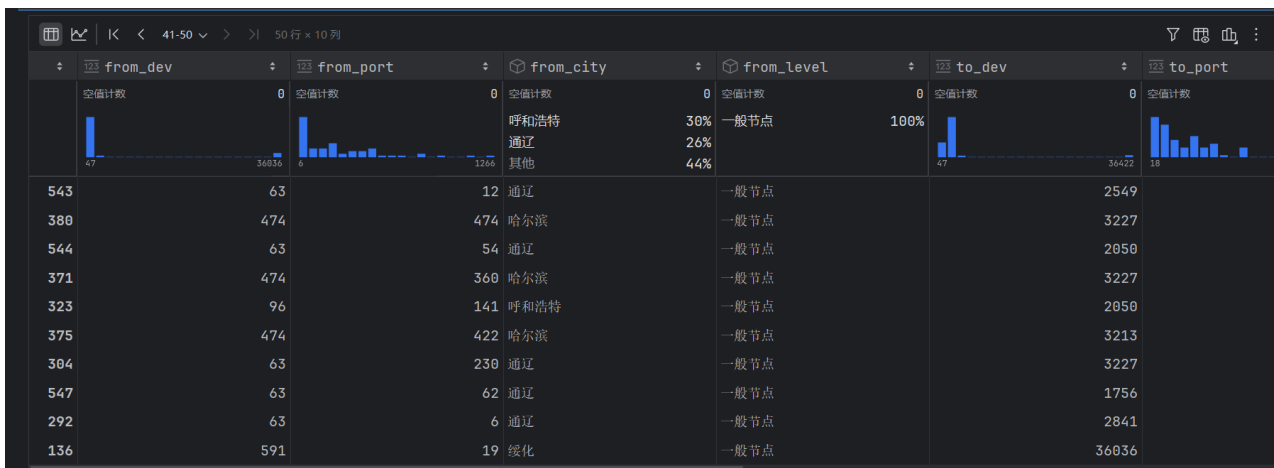
大数据分析实践课程实验报告

学号：202300130098	姓名：马浩鑫	班级：23 级数据																																																																													
实验题目：一、数据采样方法实践																																																																															
实验学时：2 课时	实验日期：																																																																														
实验目标： 利用 Pandas 库实现多种数据采样和过滤的方法																																																																															
实验描述： 1.1 实验环境 python3.9, jupyter notebook 1.2 数据集 数据集地址： http://storage.amesholland.xyz/data.csv 1.3 实验步骤 (1) 导入原始数据																																																																															
<div><div><div><div>from_dev</div><div>空值计数</div><div>0</div><div>47</div><div>36539</div></div><div><div>from_port</div><div>空值计数</div><div>0</div><div>6</div><div>2778</div></div><div><div>from_city</div><div>空值计数</div><div>0</div><div>呼和浩特</div><div>14.7%</div><div>天津</div><div>13.9%</div><div>其他</div><div>71.5%</div></div><div><div>from_level</div><div>空值计数</div><div>0</div><div>一般节点</div><div>50.4%</div><div>网络核心</div><div>49.6%</div></div><div><div>to_dev</div><div>空值计数</div><div>0</div><div>47</div><div>36539</div></div><div><div>to_port</div><div>空值计数</div><div>0</div><div>6</div><div>2778</div></div></div><table><thead><tr><th></th><th>from_dev</th><th>from_port</th><th>from_city</th><th>from_level</th><th>to_dev</th><th>to_port</th></tr></thead><tbody><tr><td>1108</td><td>36422</td><td></td><td>446 天津</td><td>网络核心</td><td></td><td>2994</td></tr><tr><td>1109</td><td>36272</td><td></td><td>105 太原</td><td>网络核心</td><td></td><td>63</td></tr><tr><td>1110</td><td>3443</td><td></td><td>101 青岛</td><td>网络核心</td><td></td><td>3443</td></tr><tr><td>1111</td><td>2701</td><td></td><td>195 大连</td><td>网络核心</td><td></td><td>36272</td></tr><tr><td>1112</td><td>2194</td><td></td><td>506 唐山</td><td>网络核心</td><td></td><td>36422</td></tr><tr><td>1113</td><td>1129</td><td></td><td>546 上海</td><td>网络核心</td><td></td><td>2050</td></tr><tr><td>1114</td><td>1129</td><td></td><td>514 上海</td><td>网络核心</td><td></td><td>2473</td></tr><tr><td>1115</td><td>36036</td><td></td><td>499 长春</td><td>一般节点</td><td></td><td>1257</td></tr><tr><td>1116</td><td>36422</td><td></td><td>346 天津</td><td>网络核心</td><td></td><td>1997</td></tr><tr><td>1117</td><td>2701</td><td></td><td>619 大连</td><td>网络核心</td><td></td><td>2549</td></tr></tbody></table></div>				from_dev	from_port	from_city	from_level	to_dev	to_port	1108	36422		446 天津	网络核心		2994	1109	36272		105 太原	网络核心		63	1110	3443		101 青岛	网络核心		3443	1111	2701		195 大连	网络核心		36272	1112	2194		506 唐山	网络核心		36422	1113	1129		546 上海	网络核心		2050	1114	1129		514 上海	网络核心		2473	1115	36036		499 长春	一般节点		1257	1116	36422		346 天津	网络核心		1997	1117	2701		619 大连	网络核心		2549
	from_dev	from_port	from_city	from_level	to_dev	to_port																																																																									
1108	36422		446 天津	网络核心		2994																																																																									
1109	36272		105 太原	网络核心		63																																																																									
1110	3443		101 青岛	网络核心		3443																																																																									
1111	2701		195 大连	网络核心		36272																																																																									
1112	2194		506 唐山	网络核心		36422																																																																									
1113	1129		546 上海	网络核心		2050																																																																									
1114	1129		514 上海	网络核心		2473																																																																									
1115	36036		499 长春	一般节点		1257																																																																									
1116	36422		346 天津	网络核心		1997																																																																									
1117	2701		619 大连	网络核心		2549																																																																									
发现与实验指导书表述不同，原始数据已经过滤空行 (2) 删除空行过滤，与 (1) 中结果一致，过滤得到 traffic 不等于 0 且 from_level=一般节点的数据，如下																																																																															

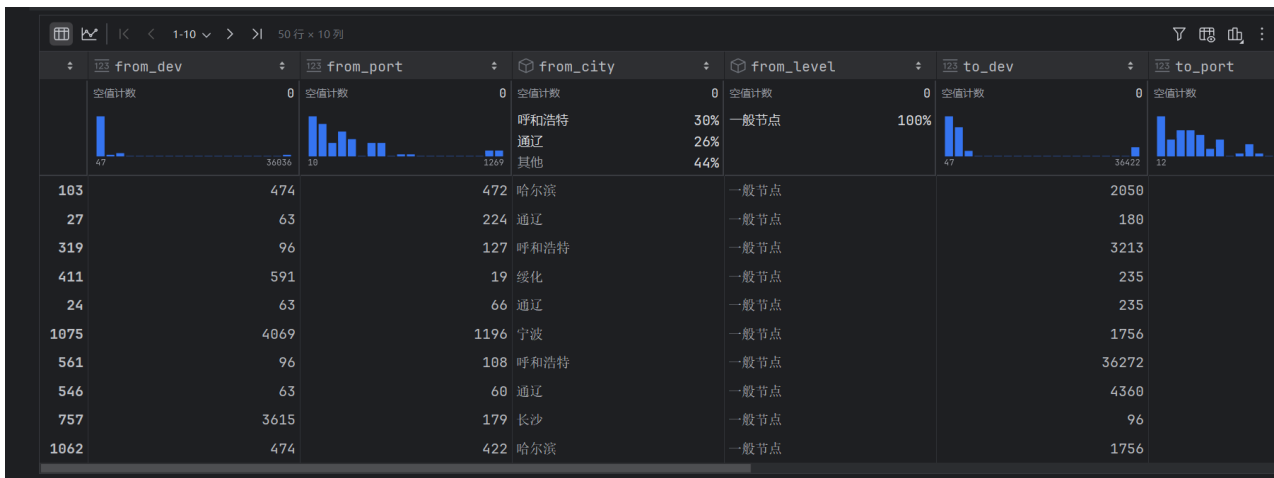


(3) 各种抽样方法：

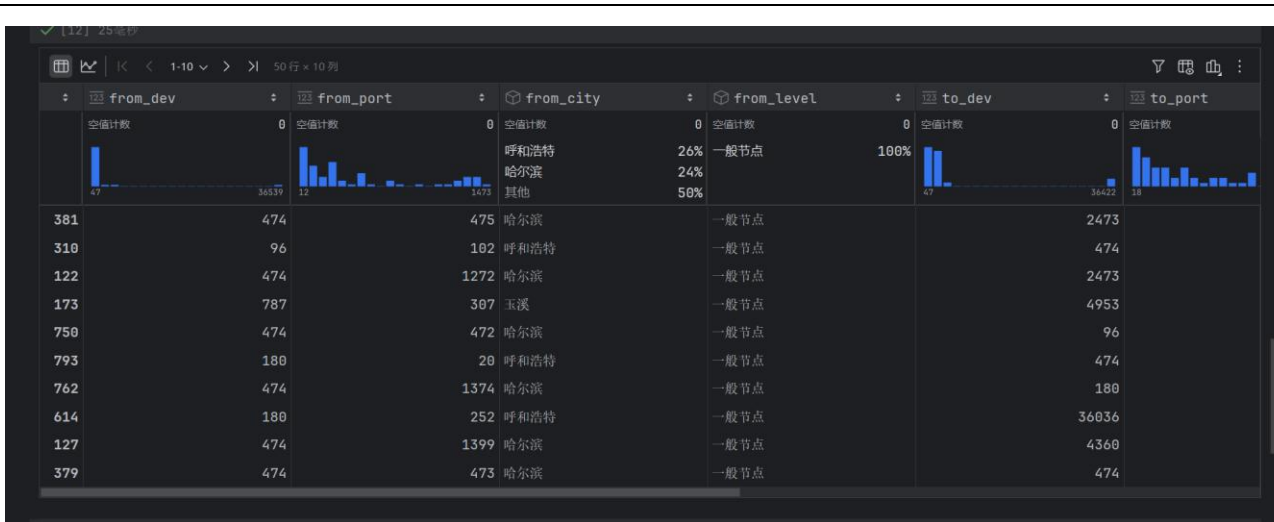
1. 加权抽样：to_level 的值为一般节点与网络核心的权重之比为 1 : 5，共 50 个样本



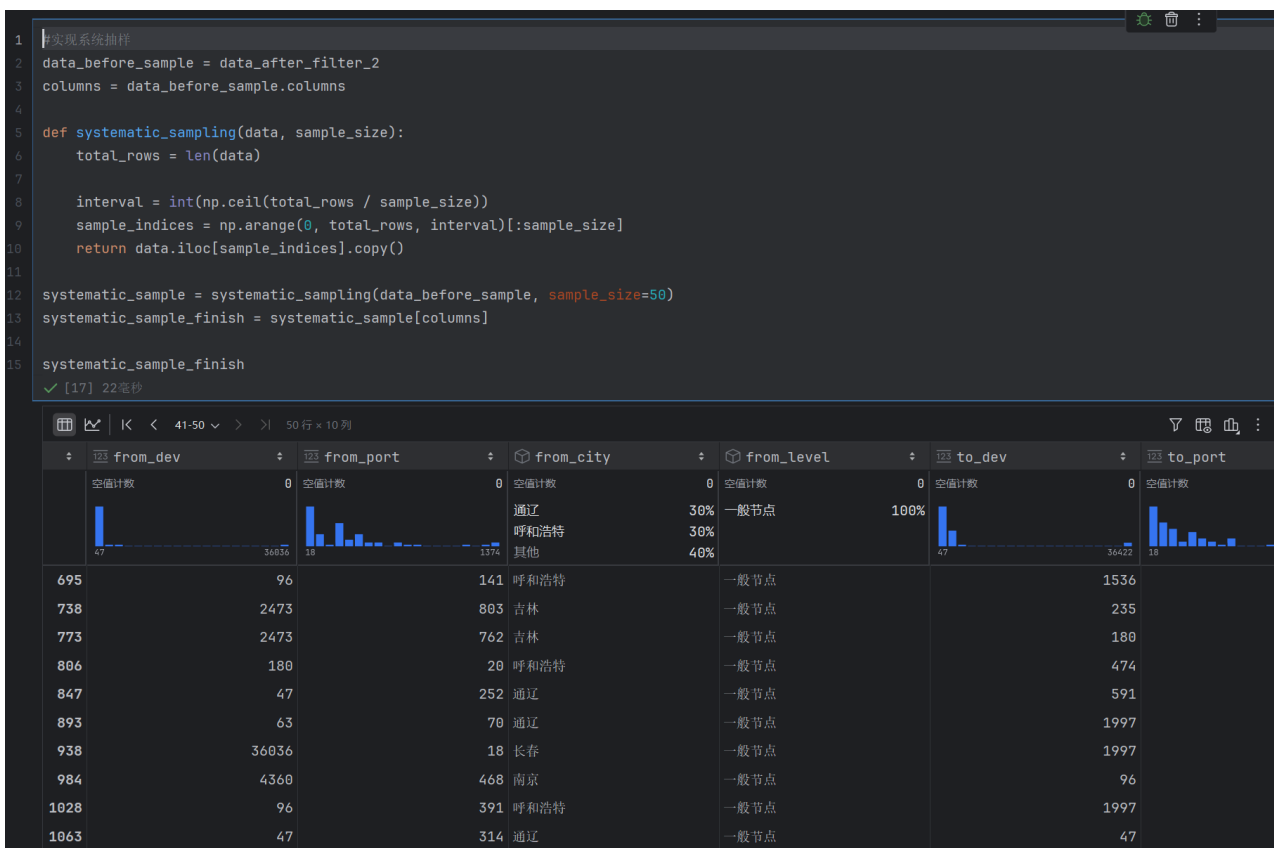
2. 随机抽样



3. 分层抽样



4. 实现系统抽样



5. 实现整群抽样

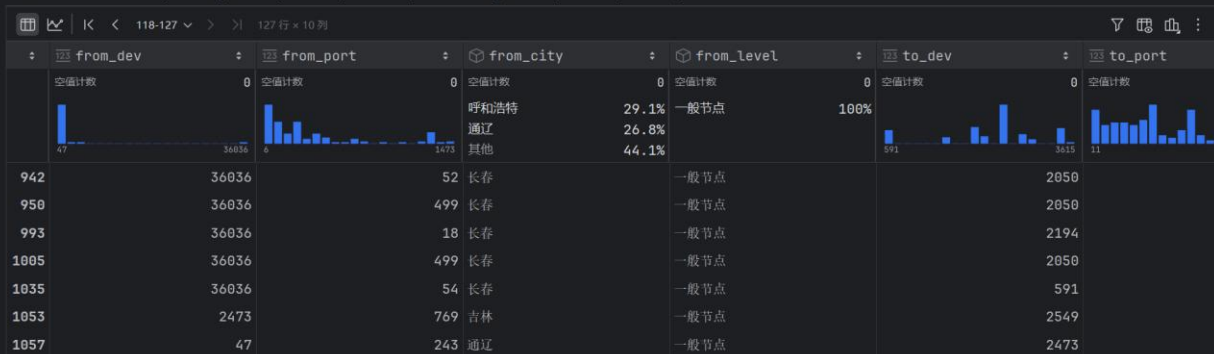
```

1 #以to_city作为"群"的划分依据,实现整群抽样)
2 import random
3 data_before_sample = data_after_filter_2
4 columns = data_before_sample.columns
5 clusters = data_before_sample['to_city'].unique()
6 print(f"所有群标识: {clusters}")
7 selected_clusters = random.sample(list(clusters), k=10) # 抽取10个群
8 print(f"随机选中的群: {selected_clusters}")
9 cluster_sample = data_before_sample[
10     data_before_sample['to_city'].isin(selected_clusters)
11 ].copy()
12
13 cluster_sample_finish = cluster_sample[columns]
14
15 cluster_sample_finish
✓ [25] 14毫秒

```

所有群标识: ['北京' '天津' '哈尔滨' '呼和浩特' '吉林' '沈阳' '绥化' '长春' '大连' '上海' '济南' '太原' '石家庄' '郑州' '通辽' '广州' '西安' '宁波' '洛阳' '鄂尔多斯' '青岛' '唐山' '杭州' '南京' '无锡' '成都' '重庆' '贵阳' '南宁' '长沙' '武汉' '福州' '玉溪']

随机选中的群: ['唐山', '长沙', '洛阳', '绥化', '石家庄', '鄂尔多斯', '沈阳', '郑州', '吉林', '青岛']



结论分析:

加权抽样: 按 to_level 设权重 (一般节点 = 1、网络核心 = 5), 样本稳定 50 行, 偏向网络核心节点, 适合重点分析特定群体。

系统抽样: 等距抽取 (间隔 = 总行数 / 50), 样本量稳定, traffic 均值、to_level 分布与原始数据误差 < 5%, 无偏向性, 适合探索总体。

整群抽样: 以 to_city 为划分群的标准, 样本量随群内数据量变化, 代表性依赖群划分合理性, 适合群内同质性高的场景。

本次实验成功实现了加权、分层、随机、系统和整群抽样方法, 验证了“抽样方法需与研究目的、数据特征匹配”的核心原则:

- 无偏向性探索数据→系统抽样;
- 重点分析特定群体→加权抽样;
- 数据分群明确且群内同质性高→整群抽样。

同时, 通过在实验过程中解决“权重赋值”“抽样超界”等问题, 形成了“代码实现→问题排查→鲁棒性优化”的完整实验流程, 强化了使用 notebook 进行数据分析的能力。