

		第一题	第二题	第三题	第四题	sum
何宇迪	1	11	21	19	29	80
刘昱君	1	11	21	24	26	82
吴天行	1	7	21	19	24	71
和丽兴	1	7	21	22	28	78
塔巴江村						未交
崔杰	1	7	16	22	20	65
庄嘉帆	1	7	23	21	23	74
张浩祖	1	9	21	19	25	74
张露丹	1	9	21	24	22	76
施習						未交
李上卫						未交
李俊杰	1	0	23	30	18	71
杨晨晨	1	13	22	19	26	80
杨翔	1	3	24	22	18	67
王丽媛	1	11	23	19	26	79
王乐						未交
王康宇	1	14	21	20	25	80
王立弘	1	10	21	30	26	87
白天琪	1	10	23	24	27	84
罗昊然	1	9	21	19	24	73
谢皓椿	1	10	21	17	21	69
贺馨姜艾	1	9	22	19	24	74
达尔汗	1	9	23	20	26	78
郭嘉						未交
郭欣遥	1	10	23	21	28	82
依夏·克热木江						未交
陈琰	1	10	22	26	24	82
韦淑荣	1	5	21	24	28	78
马宵	1	14	20	26	18	78
马月璐	1	11	21	19	25	76
黎小源	1	10	22	19	19	70
龚卓能	1	9	22	20	26	77
龚新宇	1	9	21	19	26	75
	平均分	9.04	21.52	21.59	24.15	76.30

第一题

评分标准：运行两个算例：[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] 和 [5,6,7,8]

正确答案为 77 和 19

满分 15 分，每个算例运行报错 -6；死循环或时间复杂度过高 -5

结果小于标准答案 -4

结果在集合 (78,79,80,81) 和 (19,20,21) 中：-1

结果在区间 (82, 99) 和 (22, 30)：-2

结果在其他范围 -3

有详细注释或方案输出：+2 上限 15

1|1 何宇迪 第一题

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 21:35:52 2024
4  @author: Stu
5  """
6  def ferry_time(times):
7      times.sort()
8      n = len(times)
9      total_time = 0
10     # 5要用4次，先完成这部分
11     total_time += 4 * times[0]
12     # 选择五个加和最小的较大数（除了5之外），数字不能重复选
13     remaining_times = times[1:]
14     selected_times = []
15     for _ in range(5):
16         min_sum = float('inf')
17         selected_pair = None
18         for i in range(len(remaining_times)):
19             for j in range(i + 1, len(remaining_times)):
20                 pair_sum = remaining_times[i] + remaining_times[j]
21                 if pair_sum < min_sum:
22                     min_sum = pair_sum
23                     selected_pair = (i, j)
24             if selected_pair is not None:
25                 selected_times.append(max(remaining_times[selected_pair
26 [0]], remaining_times[selected_pair[1]]))
27                 # 移除已选的两个数，避免重复选择
28                 remaining_times.pop(max(selected_pair[0], selected_pair
[1]))
29                 remaining_times.pop(min(selected_pair[0], selected_pair
```

```

29     [1]))
30     total_time += sum(selected_times)
31     return total_time
32 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
print(ferry_time(times))

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 66

[5,6,7,8] -> 27

2|2 刘昱君 第一题

```

1  #第一题
2  def minimum_crossing_time_with_plan(times):
3      # Step 1: 将所有渡河时间排序
4      times.sort()
5      total_time = 0  # 初始化总时间
6      plans = []  # 记录渡河方案
7      # Step 2: 开始循环处理, 当人数大于3时按照策略操作
8      while len(times) > 3:
9          # 方案一: 最快的两人来回帮助最慢的人
10         option1 = times[1] + times[0] + times[-1] + times[1]
11         # 方案二: 最快的人多次协助最慢的两人
12         option2 = times[-1] + times[0] + times[-2] + times[0]
13         # 选出最优方案, 并移除已到对岸的人
14         if option1 < option2:
15             # 选择方案一
16             total_time += option1
17             plans.append(f"{times[0]} 和 {times[1]} 过河 -> 用时 {times[1]}")
18             plans.append(f"{times[0]} 返回 -> 用时 {times[0]}")
19             plans.append(f"{times[-2]} 和 {times[-1]} 过河 -> 用时 {times[-1]}")
20             plans.append(f"{times[1]} 返回 -> 用时 {times[1]}")
21             # 移除最慢的两人
22             times = times[:-2]
23         else:
24             # 选择方案二
25             total_time += option2
26             plans.append(f"{times[0]} 和 {times[-1]} 过河 -> 用时 {times[-1]}")
27             plans.append(f"{times[0]} 返回 -> 用时 {times[0]}")
28             plans.append(f"{times[0]} 和 {times[-2]} 过河 -> 用时 {times[-2]}")
29             plans.append(f"{times[0]} 返回 -> 用时 {times[0]}")
30             # 移除最慢的两人
31             times = times[:-2]

```

```

32     # Step 3: 当剩余人数小于等于3时直接处理
33     if len(times) == 3:
34         # 三人同时过河
35         total_time += times[2]
36         plans.append(f"{times[0]}、{times[1]} 和 {times[2]} 过河 -> 用
时 {times[2]}")
37     elif len(times) == 2:
38         # 两人一起过河
39         total_time += times[1]
40         plans.append(f"{times[0]} 和 {times[1]} 过河 -> 用时 {times
41 [1]}")
42     elif len(times) == 1:
43         # 一人过河
44         total_time += times[0]
45         plans.append(f"{times[0]} 过河 -> 用时 {times[0]}")
46     return total_time, plans
47 # 测试用例：村民的渡河时间
48 villagers_times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
49 # 调用函数计算最小渡河时间和方案
min_time, crossing_plan = minimum_crossing_time_with_plan(villagers_
50 times)
51 # 输出结果
52 print("最少需要的渡河时间为: ", min_time)
53 print("详细渡河方案为: ")
54 for step in crossing_plan:
55     print(step)

```

输出:

详细渡河方案为:

5 和 6 过河 -> 用时 6

5 返回 -> 用时 5

16 和 20 过河 -> 用时 20

6 返回 -> 用时 6

5 和 6 过河 -> 用时 6

5 返回 -> 用时 5

13 和 15 过河 -> 用时 15

6 返回 -> 用时 6

5 和 6 过河 -> 用时 6

5 返回 -> 用时 5

10 和 11 过河 -> 用时 11

6 返回 -> 用时 6

5 和 9 过河 -> 用时 9

5 返回 -> 用时 5

5 和 7 过河 -> 用时 7

✕✓

5 返回 -> 用时 5

5 和 6 过河 -> 用时 6

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 129

[5, 6, 7, 8] -> 31

-3.
-3

3|3 吴天行 第一题

7

```
1  from itertools import combinations
2  # 定义村民的渡河时间列表，对应1 - 10号村民
3  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
4  n = len(times)
5  # 创建一个字典用于存储动态规划的中间结果，键为状态，值为对应状态下的最少渡河时间
6  dp = {}
7  # 定义一个函数来获取当前状态下的最少渡河时间
8  def min_time(state):
9      if state in dp:
10         return dp[state]
11     # 将状态转换为集合，方便后续操作
12     left_people = set([i for i in range(n) if (state >> i) & 1])
13     if not left_people:
14         return 0
15     min_time_value = float('inf')
16     # 枚举从左边（牛市）送到右边（河对岸）的人员组合，人数从1到3人
17     for r in range(1, min(4, len(left_people) + 1)):
18         for combination in combinations(left_people, r):
19             right_people = left_people - set(combination)
20             # 计算本次送过去花费的时间，以组合中最慢的人时间为准
21             current_time = max([times[i] for i in combination])
22             # 递归计算将剩下的人送过去的最少时间（假设已经把这组人送过去了）
23             remaining_time = min_time(sum(1 << i for i in right_people))
24             # 从对岸选最快的人送回（这里简单模拟，通过找剩余人员中最快的）
25             if right_people:
26                 back_time = min([times[i] for i in right_people])
27                 # 计算总时间
28                 total_time = current_time + remaining_time + back_time
29             me
30             else:
31                 total_time = current_time + remaining_time
32             # 更新最少时间
33             min_time_value = min(min_time_value, total_time)
34         dp[state] = min_time_value
35     return min_time_value
36 # 初始状态下所有人都在左边（牛市），对应状态为所有位都为1的二进制数
37 initial_state = (1 << n) - 1
print(min_time(initial_state))
```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 62

[5,6,7,8] -> 18

4/4 和丽兴 第一题

```
1 def find_min_time(times):
2     times.sort() # 将渡河时间从小到大排序
3     memo = {}
4     def helper(people):
5         # 如果只有一个人, 直接返回其时间
6         if len(people) == 1:
7             return people[0]
8         # 如果没有人, 返回 0
9         elif len(people) == 0:
10            return 0
11        people_tuple = tuple(people)
12        if people_tuple in memo:
13            return memo[people_tuple]
14        # 用于存储可能的最小时间
15        min_time = float('inf')
16        # 考虑两种主要策略:
17        # 1. 让最轻的两个人过河
18        if len(people) >= 2:
19            first = people[0]
20            second = people[1]
21            # 只有两个人时
22            if len(people) == 2:
23                time = max(first, second) + helper([]) # 直接返回
24            else:
25                # 两人过, 那个人回来
26                time = max(first, second) + helper(people[2:]) # 去掉前两个
27            min_time = min(min_time, time)
28        # 2. 让最轻和最重的两人组合过河
29        if len(people) >= 3:
30            first = people[0]
31            second = people[1]
32            third = people[2]
33            # 让最重的两人过, 返回最轻的
34            time = max(third, second) + first + helper(people[3:])
35            min_time = min(min_time, time)
36        memo[people_tuple] = min_time
37        return min_time
38    return helper(times)
39 # 村民的渡河时间
40 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
41 # 计算最小时间
```

```
42 min_time = find_min_time(times)
43 print(f"最少花费时间: {min_time}")
```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 61

[5,6,7,8] -> 14

5|5 塔巴江村 第一题

6|6 崔杰 第一题

```
1 #整体思路:
2 #通过循环来逐步安排村民过河, 每次根据当前剩余村民的人数情况, 选择合适的策略让一部分村民过河, 同时记录累计花费的最少时间, 直到所有村民都过河为止。
3 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
4 n = len(times)
5 min_time = 0
6 while n > 0:
7     if n == 1:
8         min_time += times[0]
9         break
10    elif n == 2:
11        min_time += max(times[0], times[1])
12        break
13    elif n == 3:
14        min_time += max(times[0], times[1], times[2])
15        break
16    else:
17        # 每次选择最快的3个人组合过河 (这里假设取前3个元素模拟每次选择相对快的3人)
18
19        group = times[:3]
20        min_time += max(group)
21        times = times[3:]
22        n -= 3
print(min_time)
```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 60

[5,6,7,8] -> 15

7|7 庄嘉帆 第一题

```
1 def shortest_time_to_cross(people):
2     people.sort() # 将渡船时间从小到大排序
3     total_time = 0 # 总过河时间
```

```

4     left = 0  # 左岸人数
5     right = len(people) - 1  # 右岸人数
6     while left < right:
7         # 当左岸人数小于等于2时，直接让剩下的人都过河
8         if left + 2 >= right:
9             total_time += people[right]
10            break
11        # 贪心策略：让渡船时间最长的两个人先过河
12        total_time += people[right]  # 两个人过河
13        right -= 2
14        # 如果还有人没过河，让渡船时间最短的人把船划回来
15        if left < right:
16            total_time += people[left]
17            left += 1
18    return total_time
19    # 每个人渡船时间列表
20    people_times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
21    # 计算所需的最短时间
22    print(shortest_time_to_cross(people_times))
23

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 64

[5,6,7,8] -> 13

4

4

8|8 张浩祖 第一题

```

1     def river_crossing(times):
2         times.sort()
3         n = len(times)
4         if n == 1:
5             return times[0]
6         elif n == 2:
7             return times[1]
8         elif n == 3:
9             return times[2]
10        else:
11            total_time = 0
12            while n > 3:
13                option1 = times[0] + times[1] + times[n - 1] + times[1]
14                option2 = times[0] + times[n - 2] + times[0] + times[n -
15                1]
16                total_time += min(option1, option2)
17                n -= 2
18            if n == 3:
19                total_time += times[2]
20            elif n == 2:

```

9


```

21         total_time += times[1]
22     else:
23         total_time += times[0]
24     return total_time
25 if __name__ == "__main__":
26     times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
27     print(river_crossing(times))

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 129

[5,6,7,8] -> 31

9|9 张露丹 第一题

```

1  import sys
2  # 初始化输入数据
3  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
4  N = len(times)
5  # 对渡河时间进行排序, 以便每次让最慢的人过河, 最快的人返回
6  times.sort()
7  # 初始化变量
8  total_time = 0
9  current_time = 0
10 # 每次最多载3人, 但需要保证有人能回来划船
11 for i in range(N):
12     if i % 3 == 2: # 每3个人过河后, 最快的返回
13         total_time += max(times[i - 2:i + 1]) # 加上这3人中最慢的
14         current_time = times[i] # 最快的返回
15     elif i == N - 1 or (i + 1) % 3 == 0: # 最后一个人过河或者3的倍数人过
16 河
17         total_time += max(times[i - (i % 3):i + 1]) # 加上最后这组人中
18 最慢的
19 # 如果最后一组人数不足3人, 但有人需要返回, 则加上返回时间
20 if N % 3 != 0:
21     total_time += current_time # 最后一个最快的返回的人的时间不再增加, 因为
    已经在计算内了
    print("最少需要花费的时间为:", total_time)

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 70

[5,6,7,8] -> 22

10|10 施習 第一题

11|11 李上卫 第一题

12|12 李俊杰 第一题

13|13 杨晨晨 第一题

13

```
1 def river_crossing(times):
2     times.sort()
3     n = len(times)
4     plan = []
5     if n == 1:
6         return times[0], plan
7     elif n == 2:
8         return max(times), plan
9     elif n == 3:
10        return max(times), plan
11    else:
12        min_total_time = float('inf')
13        best_plan = []
14        # 尝试所有可能的第一次送两人过河的组合情况
15        for left in range(1, n - 1):
16            for right in range(left + 1, n):
17                # 方案一：先送两人过去（left和right位置对应的两人），然后让
                # 最快的回来
18                plan_1_time = max(times[left], times[right]) + times
19                [0]
20                rest_people = times[1:left] + times[left + 1:right]
21                + times[right + 1:]
22                # 对剩下的人继续安排过河，递归调用函数
23                sub_time_1, sub_plan_1 = river_crossing(rest_people)
24                plan_1_time += sub_time_1
25                this_plan_1 = [("方案一", (left, right), (times[lef
26                t], times[right]), (0,)), (times[0],)), tuple([x for x in sub_plan_
27                1]))]
28                # 方案二：先送最快的和次快的过去，然后让最快的回来
29                plan_2_time = max(times[0], times[1]) + times[0]
30                rest_people_2 = times[2:]
31                sub_time_2, sub_plan_2 = river_crossing(rest_people_
32                2)
33                plan_2_time += sub_time_2
34                this_plan_2 = [("方案二", (0, 1), (times[0], times
35                [1]), (0,)), (times[0],)), tuple([x for x in sub_plan_2]))]
36                # 比较当前两种方案时间，选择更优的
37                if plan_1_time < plan_2_time:
38                    cur_plan = this_plan_1
39                    cur_time = plan_1_time
40                else:
```

```

38         cur_plan = this_plan_2
39         cur_time = plan_2_time
40         if cur_time < min_total_time:
41             min_total_time = cur_time
42             best_plan = cur_plan
43         return min_total_time, best_plan
44 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
45 result, plan_list = river_crossing(times)
46 print("最少花费时间:", result)
47 for idx, (scheme, first_go_indexes, first_go_times, first_back_index, first_back_time, sub_plan) in enumerate(plan_list):
48     print(f"第{idx + 1}步:")
49     print(f"选择 {scheme}")
50     print(f"第一次送过去的人员索引为 {first_go_indexes}, 渡河时间为 {first_go_times}")
51     print(f"第一次返回的人员索引为 {first_back_index}, 渡河时间为 {first_back_time}")
52     if sub_plan:
53         for sub_idx, (sub_scheme, sub_first_go_indexes, sub_first_go_times, sub_first_back_index, sub_first_back_time, _) in enumerate(sub_plan):
54             print(f"    子步骤{sub_idx + 1}:")
55             print(f"        选择 {sub_scheme}")
56             print(f"        第一次送过去的人员索引为 {sub_first_go_indexes}, 渡河时间为 {sub_first_go_times}")
57             print(f"        第一次返回的人员索引为 {sub_first_back_index}, 渡河时间为 {sub_first_back_time}")

```

输出:

第1步:

选择 方案二

第一次送过去的人员索引为 (0, 1), 渡河时间为 (5, 6)

第一次返回的人员索引为 (0,), 渡河时间为 (5,)

子步骤1:

选择 方案二

第一次送过去的人员索引为 (0, 1), 渡河时间为 (7, 9)

第一次返回的人员索引为 (0,), 渡河时间为 (7,)

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 70

[5,6,7,8] -> 19

14|14 杨翔 第一题

```

1 def min_crossing_time(times):
2     times.sort()
3     n = len(times)

```

```

4     total_time = 0
5     across = 0
6     while across < n:
7         if across < 2 and (n - across) > 0:
8             slowest = times[across + (n - across - 1 if across + 2 <
= n else 1)]
9             total_time += slowest
10            across += 2
11            across -= 1
12            elif across == 0 or (across == 1 and (n - across) > 0):
13                slowest = times[across + (n - across if across + 1 <= n
else 0)]
14                total_time += slowest * 2
15                if across == 0:
16                    across += 1
17            return total_time
18 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
19 print(min_crossing_time(times))

```

输出：死循环

15|15 王丽媛 第一题

```

1     # 1
2     def min_time_to_cross(T):
3         N = len(T)
4         T.sort()
5         time = 0
6         i = 0
7         while i < N:
8             if N - i <= 3:
9                 time += T[-1]
10                break
11            if N - i == 4:
12                time += T[2] + T[3]
13                time += T[0] + T[3]
14                break
15            if i + 3 < N and T[i] + T[i+1] + T[i+2] > T[i+2] + T[i+3] +
T[0]:
16                time += T[i+2] + T[i+3]
17                time += T[0]
18                i += 3
19            else:
20                time += T[i] + T[i+1] + T[i+2]
21                time += T[i]
22                i += 3
23        return time
24

```

```
25 | T = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
    | print(min_time_to_cross(T))
```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 82

[5,6,7,8] -> 28

~
~

16|16 王乐 第一题

17|17 王康宇 第一题

```
1 | times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
2 | timelist = times
3 | list.sort(times)
4 | daoan = []
5 | timetotal = 0
6 | a = 0
7 | def lunci_one(times,daoan):
8 |     if len(times)>2:
9 |         num = times[:3]
10 |         times = times[3:]
11 |         daoan+=num
12 |         time = max(num) + min(daoan)
13 |         if len(times) == 0:
14 |             pass
15 |         else:
16 |             times.append(min(daoan))
17 |             daoan.remove(min(daoan))
18 |     else:
19 |         num = times
20 |         time = max(num)
21 |         times = []
22 |         list.sort(times)
23 |         list.sort(daoan)
24 |         print(time, times, daoan)
25 |         return time,times,daoan
26 | def lunci_two(times,daoan):
27 |     if len(times)>2:
28 |         num = times[-3:]
29 |         times = times[:-3]
30 |         daoan+=num
31 |         time = max(num) + min(daoan)
32 |         if len(times) == 0:
33 |             pass
34 |         else:
35 |             times.append(min(daoan))
36 |             daoan.remove(min(daoan))
```

14

```

37         else:
38             num = times
39             daoan += num
40             time = max(num)
41             times = []
42             pass
43         list.sort(times)
44         list.sort(daoan)
45         print(time, times, daoan)
46         return time, times, daoan
47     ttt = 0
48     for i in range(1000):
49         if i%2==0:
50             a, times, daoan = lunci_one(times, daoan)
51             ttt += a
52             if len(times) == 0:
53                 break
54         if i%2!=0:
55             a, times, daoan = lunci_two(times, daoan)
56             ttt += a
57             if len(times) == 0:
58                 break
59     print(a)
60     print(f"最终花费时间{ttt}")

```

输出:

12 [5, 9, 10, 11, 13, 15, 16, 20] [6, 7]

12

26 [5, 6, 9, 10, 11, 13] [7, 15, 16, 20]

26

14 [5, 10, 11, 13] [6, 7, 9, 15, 16, 20]

14

19 [5, 6] [7, 9, 10, 11, 13, 15, 16, 20]

19

6 [] [7, 9, 10, 11, 13, 15, 16, 20]

最终花费时间77

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 77

[5, 6, 7, 8] -> 20

18|18 王立弘 第一题

```

1     times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
2     times.sort()
3     total_time = 0
4     n = len(times)
5     while n > 3:

```

10

```

6      # 每次选择最快的和两个最慢的过去，花费时间是两个最慢的人中较慢的那个时间
7      go_time = max(times[-2:])
8      # 最快的回来，花费最快的那个人的时间
9      back_time = times[0]
10     total_time += go_time + back_time
11     n -= 3
12     if n == 3:
13         total_time += max(times[:3])
14     elif n == 2:
15         total_time += max(times[:2])
16     elif n == 1:
17         total_time += times[0]
18     print(total_time)

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 80
 [5, 6, 7, 8] -> 18

19|19 白天琪 第一题

```

1  import itertools
2  # 村民的渡河时间（10个村民）
3  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
4  # 计算从一组村民中渡河的时间（最慢的人的时间）
5  def calculate_time(group):
6      return max(group)
7  # 递归函数来计算最小总时间，并列举所有情况
8  def min_travel_time(times):
9      # 记录最小时间，初始设置为一个很大的值
10     min_time = float('inf')
11     all_scenarios = [] # 用于记录所有的过河情况
12     # 使用递归来枚举所有的可能性
13     def recurse(left_times, right_times, total_time, scenario):
14         nonlocal min_time, all_scenarios
15         # 如果所有村民都过河了，更新最小时间并记录情况
16         if not left_times:
17             min_time = min(min_time, total_time)
18             all_scenarios.append((list(scenario), total_time))
19             return
20         # 枚举所有可能的组合，最多3个村民一起过河
21         for num_people in range(2, 4): # 2到3个人一起过河
22             for group in itertools.combinations(left_times, num_people):
23                 # 计算这组人的渡河时间
24                 group_time = calculate_time(group)
25                 # 计算过河后的时间
26                 new_total_time = total_time + group_time

```

```

27         # 过河后，有人必须返回，返回的人是最慢的那个
28         if right_times: # 如果已经有过河的人
29             # 从已过河的村民中选择最慢的人带船回来
30             for return_person in sorted(right_times): # 排序
保证选择时间最短的返回
31                 return_time = return_person
32                 new_total_time_with_return = new_total_time
+ return_time
33             # 递归，继续考虑剩下的村民
34             new_scenario = scenario + [group]
35             recurse(
36                 [time for time in left_times if time not
in group], # 移除已过河的村民
37                 right_times + list(group), # 记录已过河的
38                 村民
39                 new_total_time_with_return, # 累加总时间
40                 new_scenario # 更新当前过河情况
41             )
42         else: # 如果没有过河的人，就直接过去
43             new_scenario = scenario + [group]
44             recurse(
45                 [time for time in left_times if time not in
group], # 移除已过河的村民
46                 right_times + list(group), # 记录已过河的村民
47                 new_total_time, # 不需要返回的情况
48                 new_scenario # 更新当前过河情况
49             )
50         # 开始递归计算最小总时间
51         recurse(times, [], 0, [])
52         # 输出所有的过河情况及其时间
53         print("所有过河情况和对的时间：")
54         for scenario, total_time in all_scenarios:
55             print(f"情况：{scenario} => 总时间：{total_time}")
56         return min_time
57     # 最短的时间
58     result = min_travel_time(times)
59     print(f"最少的总时间是：{result}")

```

输出：

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 时间复杂度太高，无法运行

[5,6,7,8] -> 19

20|20 罗昊然 第一题

```

1 def calculate_time(crossing_times):
2     crossing_times.sort()

```



```

3     total_time = 0
4     # 第一次过河
5     first_crossing_time = max(crossing_times[0], crossing_times[1],
crossing_times[2])
6     total_time += first_crossing_time
7     total_time += crossing_times[0] # 5回来
8     # 第二次过河
9     second_crossing_time = max(crossing_times[7], crossing_times[8],
crossing_times[9])
10    total_time += second_crossing_time
11    total_time += crossing_times[1] # 6回来
12    # 第三次过河
13    third_crossing_time = max(crossing_times[0], crossing_times[1],
crossing_times[3])
14    total_time += third_crossing_time
15    total_time += crossing_times[0] # 5回来
16    # 第四次过河
17    fourth_crossing_time = max(crossing_times[4], crossing_times[5],
crossing_times[6])
18    total_time += fourth_crossing_time
19    total_time += crossing_times[1] # 6回来
20    # 第五次过河
21    fifth_crossing_time = max(crossing_times[0], crossing_times[1])
22    total_time += fifth_crossing_time
23    return total_time
24    # 示例渡河时间
25    crossing_times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
26    total_time = calculate_time(crossing_times)
27    print("最少花费的时间为:", total_time)

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 77

[5,6,7,8] -> 报错, 无结果

✓b

21|21 谢皓椿 第一题

```

1     # -*- coding: utf-8 -*-
2     """
3     Created on Fri Dec 20 21:21:37 2024
4     @author: Stu
5     """
6     def min_crossing_time(people_times):
7         """
8         计算多人过河的最短时间, 采用新策略: 最快三个过去, 最快一个回来, 最慢三个过
去, 最快一个又回来, 重复。
9         :param people_times: list of int, 每个人的过河时间 (单位: 分钟或其他时
间单位)。

```

10 ✓

```

10     :return: int, 最短过河时间（单位与输入相同）。
11     """
12     # 对过河时间进行排序，并同时记录每个人的索引，以便后续按时间顺序选择人
13     sorted_times_with_indices = sorted((time, index) for index, time
in enumerate(people_times))
14     n = len(people_times)
15     total_time = 0
16     # 初始化两个列表来跟踪对岸和起始岸的人
17     across_river = [] # 对岸的人（按过河时间排序）
18     start_side = sorted_times_with_indices # 起始岸的人（按过河时间排
19 序）
20     while start_side:
21         # 每次选择最快的三个人过去
22         fastest_three = start_side[:3]
23         start_side = start_side[3:]
24         # 更新对岸的人列表
25         across_river.extend(fastest_three)
26         # 计算这次过河的时间（由最慢的那个人决定）
27         slowest_in_group = max(time for time, _ in fastest_three)
28         total_time += slowest_in_group
29         # 如果还有人在起始岸，让对岸最快的那个人返回
30         if start_side:
31             # 找到对岸最快的人（即之前过河且现在在对岸的人中时间最短的那个）
32             fastest_across = min(time for time, _ in across_river)
33             fastest_return_time, _ = next((time_index for time_index
in across_river if time_index[0] == fastest_across))
34             total_time += fastest_return_time
35         return total_time
36     # 示例输入：每个人的过河时间（分钟）
37     people_times = [5,6,7,9,10,11,13,15,16,20]
    print("最短过河时间:", min_crossing_time(people_times))

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 69

[5,6,7,8] -> 20

22|22 贺馨姜艾 第一题

```

1     times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
2     times.sort()
3     n = len(times)
4     if n > 3:
5         sum_time = 0
6         while n > 3:
7             option1 = times[0] + times[-1] + times[0] + times[-2]
8             option2 = times[1] + times[0] + times[-1] + times[1]
9             if option1 < option2:

```

```

10         sum_time += option1
11     else:
12         sum_time += option2
13     n -= 2
14     if n == 3:
15         sum_time += times[2]
16     elif n == 2:
17         sum_time += times[1]
18     else:
19         sum_time += times[0]
20 else:
21     sum_time = times[-1]
22 print("最少花费", sum_time, ", 才能使所有人都过河")

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 154

[5,6,7,8] -> 31

23|23 达尔汗 第一题

```

1  # 定义村民的渡河时间列表
2  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
3  # 初始化总时间
4  total_time = 0
5  # 第一次渡河
6  # 5、6、7过河，花费时间取三人中最大的，然后5回来
7  total_time += max(times[0], times[1], times[4])
8  total_time += times[0]
9  # 第二次渡河
10 # 15、16、20过河，花费时间取三人中最大的，然后6回来
11 total_time += max(times[8], times[5], times[7])
12 total_time += times[4]
13 # 第三次渡河
14 # 5、6、9过河，花费时间取三人中最大的，然后5回来
15 total_time += max(times[0], times[4], times[5])
16 total_time += times[0]
17 # 第四次渡河
18 # 10、11、13过河，花费时间取三人中最大的，然后6回来
19 total_time += max(times[6], times[2], times[9])
20 total_time += times[4]
21 # 第五次渡河
22 # 5、6过河，花费时间取两人中最大的
23 total_time += max(times[0], times[4])
24 print("最少花费的总时间为:", total_time)

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 77

[5,6,7,8] -> 报错, 无结果

6

24|24 郭嘉 第一题

25|25 郭欣遥 第一题

10

```
1 (1)
2 # 定义渡河时间列表
3 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
4 # 对渡河时间列表进行排序, 方便后续处理
5 times.sort()
6 n = len(times)
7 # 初始化总时间
8 total_time = 0
9 # 当还有超过3个人没过河时
10 while n > 3:
11     # 选择最快的两个人先过去, 然后最快的回来
12     option1 = times[0] + times[1] + times[n - 1]
13     # 选择最快的和最慢的过去, 然后最快的回来, 再选两个快的过去, 最慢的回来
14     option2 = 2 * times[0] + times[n - 2] + times[n - 1]
15     # 选择时间花费较少的方案
16     total_time += min(option1, option2)
17     n -= 2
18 # 处理剩下3个人或者2个人或者1个人的情况
19 if n == 3:
20     total_time += times[2]
21 elif n == 2:
22     total_time += times[1]
23 else:
24     total_time += times[0]
25 print(total_time)
26 #解题思路
27 # 整体思路
28 #本题采用贪心算法来求解最少渡河时间, 贪心算法的核心是在每一步选择中都选择当前状态
    下看似最优的选项, 希望通过一系列的局部最优选择能够得到全局最优解。
29 # 具体步骤
30 # (1)数据预处理**: 首先将所有人的渡河时间进行排序, 这样可以方便后续按照时间长短来
    安排渡河策略, 排序后的时间列表记为 `times`, 人数记为 `n`。
31 # (2)多人渡河策略**: 当还有超过3个人没过河时, 有两种主要的渡河方案可供选择, 并从
    中选取时间花费较少的方案。
32 # 一、方案一**: 选择最快的两个人先过去, 然后最快的人回来。此方案花费的时间为 `ti
    mes[0] + times[1] + times[n - 1]`。这样做的目的是先利用最快的两人将船送到
    对岸, 然后让最快的人把船带回来, 以便后续人员继续渡河, 其中 `times[0]` 和 `time
    s[1]` 分别是最快的两个人渡河时间, `times[n - 1]` 是最后一个人渡河时间。
33 ##二、方案二**: 选择最快的和最慢的过去, 然后最快的回来, 再选两个快的过去, 最慢的
```

```

34 回来。该方案花费的时间为 `2 * times[0] + times[n - 2] + times[n - 1]` 。
    这种方案是先让最快的人和最慢的人一起过河，然后最快的人回来，接着让最快的人和次慢
    的人过去，再让最快的人回来，以便继续运送剩余人员 。
    #更新人数和总时间：每次选择完方案后，将人数 `n` 减去2，表示已经有两人成功渡河，
35 同时将所选方案的时间加到总时间 `total_time` 中。
36 #处理剩余人员：当剩下3个人、2个人或者1个人时，分别进行处理。
    # - 如果剩下3个人，那么总时间加上三人中最慢的人的时间，即 `times[2]` ，因
    为此时可以让最快的人先送一个人过去，然后再和剩下的那个人一起过去，总时间取决于最
37 慢的那个人的渡河时间 。
    # - 如果剩下2个人，总时间加上较慢的那个人的时间，即 `times[1]` ，两人直接
38 一起过河即可 。
    # - 如果剩下1个人，总时间加上这个人的渡河时间，即 `times[0]` 。
39
40 ### 代码实现思路
    #- 按照上述解题思路，使用Python语言实现了相应的代码。首先定义了渡河时间列表 `ti
    mes` 并进行排序，然后通过循环和条件判断来逐步计算最少渡河时间，直到所有人都过
41 河，最后输出总时间 `total_time` 。

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 105

[5,6,7,8] -> 25

26|26 阿依夏·克热木江 第一题

27|27 陈琰 第一题

```

1  # (1)定义人员渡河时间列表
2  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
3  times.sort()
4  n = len(times)
5  def ferry_optimized_again(times):
6      left_num = n
7      total_time = 0
8      while left_num > 3:
9          option1 = times[0] + times[1] + times[left_num - 1] + times[0]
10         option2 = 2 * times[0] + times[left_num - 1] + times[left_
            num - 2]
11         if option1 < option2:
12             total_time += option1
13         else:
14             total_time += option2
15         left_num -= 2
16 # 简化剩余人数不同情况的处理
17 remain_times = times[-left_num:]
18 if left_num == 3:
19     total_time += sum(remain_times)

```

```

20     elif left_num == 2:
21         total_time += remain_times[1]
22     elif left_num == 1:
23         total_time += remain_times[0]
24     return total_time
25 print(ferry_optimized_again(times))
26

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 139

[5,6,7,8] -> 32

28|28 韦淑荣 第一题

```

1  def min_crossing_time(times):
2      # 对渡河时间进行排序
3      times.sort()
4      # 初始化总时间为0
5      total_time = 0
6      # 当还有人未过河时继续循环
7      while len(times) > 3:
8          # 最快的两个人先过河
9          first = times.pop(0)
10         second = times.pop(0)
11         # 最慢的一个人返回
12         third = times.pop(-1)
13         # 更新总时间（最快两人过河的时间 + 最慢一人返回的时间）
14         total_time += second
15         # 剩下的人中再选择最快的两个过河
16         fourth = times.pop(0)
17         fifth = times.pop(0)
18         # 更新总时间（剩下最快的两人过河的时间）
19         total_time += fifth
20         # 把之前返回的最慢的人放回列表
21         times.append(third)
22         # 重新排序剩余的人
23         times.sort()
24     # 最后三人一起过河
25     if len(times) == 3:
26         total_time += max(times)
27     elif len(times) == 2:
28         total_time += times[1]
29     elif len(times) == 1:
30         total_time += times[0]
31     return total_time
32 # 示例数据
33

```

```
34 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
    print("最少需要花费的时间是:", min_crossing_time(times))
```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 61

[5,6,7,8] -> 错误, 无输出

4
6

29|29 马宵 第一题

```
1  def min_crossing_time(times):
2      times.sort()
3      total_time = 0
4      steps = []
5      while len(times) > 3:
6          fastest = times[0]
7          slowest = times[-1]
8          second_slowest = times[-2]
9          steps.append((fastest, second_slowest, slowest))
10         total_time += slowest
11         steps.append((fastest,))
12         total_time += fastest
13         times.pop(-1)
14         times.pop(-1)
15     if len(times) == 3:
16         steps.append((times[0], times[1]))
17         total_time += times[1]
18         steps.append((times[0],))
19         total_time += times[0]
20         steps.append((times[0], times[2]))
21         total_time += times[2]
22     elif len(times) == 2:
23         steps.append((times[0], times[1]))
24         total_time += times[1]
25     elif len(times) == 1:
26         steps.append((times[0],))
27         total_time += times[0]
28     return steps, total_time
29 times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
30 steps, total_time = min_crossing_time(times)
31 print("最短时间过河的方案:")
32 for step in steps:
33     print(f"村民 {step} 过河")
34 print("总时间:", total_time)
```

14

输出:

最短时间过河的方案:

村民 (5, 16, 20) 过河

村民 (5,) 过河

村民 (5, 13, 15) 过河

村民 (5,) 过河

村民 (5, 10, 11) 过河

村民 (5,) 过河

村民 (5, 7, 9) 过河

村民 (5,) 过河

村民 (5, 6) 过河

总时间: 81

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 81

[5,6,7,8] -> 19

30|30 马月璐 第一题

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 20:50:02 2024
4  @author: Stu
5  """
6  def river_crossing(times):
7      times.sort()
8      n = len(times)
9      plan = []
10     if n == 1:
11         return times[0], plan
12     elif n == 2:
13         return max(times), plan
14     elif n == 3:
15         return max(times), plan
16     else:
17         min_total_time = float('inf')
18         best_plan = []
19         # 尝试所有可能的第一次送两人过河的组合情况
20         for left in range(1, n - 1):
21             for right in range(left + 1, n):
22                 # 方案一：先送两人过去（left和right位置对应的两人），然后让
23                 # 最快的回来
24                 plan_1_time = max(times[left], times[right]) + times[0]
25                 rest_people = times[1:left] + times[left + 1:right]
26                 + times[right + 1:]
27                 # 对剩下的人继续安排过河，递归调用函数
28                 sub_time_1, sub_plan_1 = river_crossing(rest_people)
29                 plan_1_time += sub_time_1
30                 this_plan_1 = [("方案一", (left, right), (times[left], times[right]), (0,)), (times[0],), tuple([x for x in sub_plan_1
```



```

30     1]]))
31         # 方案二：先送最快的和次快的过去，然后让最快的回来
32         plan_2_time = max(times[0], times[1]) + times[0]
33         rest_people_2 = times[2:]
34         sub_time_2, sub_plan_2 = river_crossing(rest_people_
2)
35         plan_2_time += sub_time_2
36         this_plan_2 = [("方案二", (0, 1), (times[0], times
37 [1]), (0,)), (times[0],), tuple([x for x in sub_plan_2]))]
38         # 比较当前两种方案时间，选择更优的
39         if plan_1_time < plan_2_time:
40             cur_plan = this_plan_1
41             cur_time = plan_1_time
42         else:
43             cur_plan = this_plan_2
44             cur_time = plan_2_time
45         if cur_time < min_total_time:
46             min_total_time = cur_time
47             best_plan = cur_plan
48         return min_total_time, best_plan
times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
result, plan_list = river_crossing(times)
print("最少花费时间:", result)

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 70

[5,6,7,8] -> 19

31|31 黎小源 第一题

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 19:29:31 2024
4  @author: Stu
5  """
6  def min_time_to_cross(T):
7      T.sort()
8      N = len(T)
9      total_time = 0
10     i = 0
11     j = N - 1
12     while i < j - 1:
13         total_time += T[0] + T[1] + T[0]
14         i += 2
15         j -= 1
16     if i + 1 == j:
17         total_time += T[0] + T[j] + T[j]

```

```

18     else:
19         total_time += T[j] * 3
20     return total_time
21 T = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
22 print(min_time_to_cross(T))

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 87

[5,6,7,8] -> 37

-2
-3

32|32 龚卓能 第一题

```

1  # 定义村民的渡河时间列表
2  times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
3  # 初始化总时间
4  total_time = 0
5  # 第一次渡河
6  # 5、6、7过河，花费时间取三人中最大的，然后5回来
7  total_time += max(times[0], times[1], times[4])
8  total_time += times[0]
9  # 第二次渡河
10 # 15、16、20过河，花费时间取三人中最大的，然后6回来
11 total_time += max(times[8], times[5], times[7])
12 total_time += times[4]
13 # 第三次渡河
14 # 5、6、9过河，花费时间取三人中最大的，然后5回来
15 total_time += max(times[0], times[4], times[5])
16 total_time += times[0]
17 # 第四次渡河
18 # 10、11、13过河，花费时间取三人中最大的，然后6回来
19 total_time += max(times[6], times[2], times[9])
20 total_time += times[4]
21 # 第五次渡河
22 # 5、6过河，花费时间取两人中最大的
23 total_time += max(times[0], times[4])
24 print("最少花费的总时间为:", total_time)

```

2

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 77

[5,6,7,8] -> 报错，无输出

-6

33|33 龚新宇 第一题

```

1  # -*- coding: utf-8 -*-
2  """

```

2

```

3 Spyder 编辑器
4 这是一个临时脚本文件。
5 """
6 def min_crossing_time(times):
7     times.sort()
8     n = len(times)
9     total_time = 0
10    crossed = 0
11    left = 0
12    right = n - 1
13    current_boat = 2
14    while crossed < n:
15        if current_boat == 0 or (current_boat == 1 and left < n -
16    1):
17            if current_boat == 0:
18                slowest = max(times[left], times[left + 1])
19                left += 2
20            else:
21                slowest = times[left]
22                left += 1
23            current_boat = 2 if current_boat == 0 else 3
24            total_time += slowest
25        else:
26            if current_boat == 2:
27                fastest_return = times[left]
28                left += 1
29                total_time += fastest_return
30                current_boat = 1
31            elif current_boat == 3:
32                if left < n:
33                    fastest_return = times[left]
34                    left += 1
35                else:
36                    fastest_return = float('inf')
37                total_time += fastest_return
38                current_boat = 1
39            if current_boat == 1:
40                crossed += 1
41            elif current_boat == 0 and left < n:
42                crossed += 1
43    return total_time
44    times = [5, 7, 11, 16, 6, 9, 10, 20, 15, 13]
45    print(min_crossing_time(times))

```

输出:

[5, 7, 11, 16, 6, 9, 10, 20, 15, 13] -> 92

[5,6,7,8] -> 18

✓ 2

✓ 4

第二题

34|1 何宇迪 第二题

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 19:24:50 2024
4  @author: Stu
5  """
6  import mysql.connector
7  import csv
8  import os
9  # 连接数据库
10 mydb = mysql.connector.connect(
11     host="localhost",
12     user="root",
13     password="Qwe_1234",
14     database="online_store"
15 )
16 mycursor = mydb.cursor()
17 # 创建orders表格
18 mycursor.execute('''
19 CREATE TABLE IF NOT EXISTS orders (
20     id VARCHAR(20),
21     sku VARCHAR(15),
22     kinds VARCHAR(15),
23     stime VARCHAR(15),
24     num INT,
25     unit_price FLOAT,
26     total_price FLOAT,
27     longitude FLOAT,
28     latitude FLOAT
29 )
30 ''')
31 # 获取当前脚本所在的目录
32 script_dir = os.path.dirname(os.path.abspath(__file__))
33 # 构建文件的绝对路径
34 file_path = os.path.join(script_dir, 'C://Users/Stu/Desktop/Demo_1.csv')
35 # 从Demo_1.csv文件中读取数据并插入到orders表格
36 with open(file_path, 'r', encoding='gbk') as file:
37     reader = csv.reader(file)
38     next(reader) # 跳过标题行
39     for row in reader:
40         sql = "INSERT INTO orders (id, sku, kinds, stime, num, unit_
```

```

41 price, total_price, longitude, latitude) VALUES (%s, %s, %s, %s, %s,
42 %s, %s, %s, %s)"
43     mycursor.execute(sql, row)
44 mydb.commit()
45 # 输出创建表格和插入数据后的结果（展示前10条数据）
46 mycursor.execute("SELECT * FROM orders LIMIT 10")
47 results = mycursor.fetchall()
48 print("插入数据后的前10条订单数据：")
49 for row in results:
50     print(row)
51 # 查询三同订单（同一时间、同一地点、相同商品）并统计数量
52 mycursor.execute('''
53 SELECT ANY_VALUE(kinds) AS kinds, COUNT(*) AS count
54 FROM orders
55 GROUP BY stime, longitude, latitude, sku
56 HAVING COUNT(*) > 1
57 ''')
58 results = mycursor.fetchall()
59 # 统计每个种类的三同订单数量
60 kind_counts = {}
61 for row in results:
62     kind = row[0]
63     count = row[1]
64     if kind in kind_counts:
65         kind_counts[kind] += count
66     else:
67         kind_counts[kind] = count
68 # 找到数量最多的种类
69 max_count = 0
70 max_kind = ""
71 for kind, count in kind_counts.items():
72     if count > max_count:
73         max_count = count
74         max_kind = kind
75 print("三同订单最多的商品种类kinds: ", max_kind)

```

35|2 刘昱君 第二题

```

1 #第二题
2 import pymysql
3 import csv
4 import os
5 # 数据库连接信息
6 DB_HOST = "localhost"
7 DB_USER = "root"
8 DB_PASSWORD = "your_password"
9 DB_NAME = "online_store"
10 CSV_FILE_PATH = "C:/Users/Stu/Desktop/Demo_2.csv"

```

```

11 def create_database_and_table():
12     # """创建数据库和表"""
13     try:
14         connection = pymysql.connect(host=DB_HOST, user=DB_USER, pas
sword=DB_PASSWORD)
15         cursor = connection.cursor()
16         # 创建数据库
17         cursor.execute("CREATE DATABASE IF NOT EXISTS online_stor
18 e;")
19         cursor.execute("USE online_store;")
20         # 创建 orders 表
21         cursor.execute("""
22             CREATE TABLE IF NOT EXISTS orders (
23                 id VARCHAR(20),
24                 sku VARCHAR(15),
25                 kinds VARCHAR(15),
26                 stime VARCHAR(15),
27                 num INT,
28                 unit_price FLOAT,
29                 total_price FLOAT,
30                 longitude FLOAT,
31                 latitude FLOAT
32             );
33         """)
34         connection.commit()
35         print("数据库和表创建成功!")
36         cursor.close()
37         connection.close()
38     except Exception as e:
39         print("创建数据库或表失败: ", e)
40 def import_csv_to_database():
41     # """从 CSV 文件中读取数据并插入到数据库"""
42     if not os.path.exists(CSV_FILE_PATH): # 检查 CSV 文件是否存在
43         print(f"错误: 无法找到 CSV 文件 {CSV_FILE_PATH}")
44         return
45     try:
46         connection = pymysql.connect(host=DB_HOST, user=DB_USER, pas
sword=DB_PASSWORD, database=DB_NAME)
47         cursor = connection.cursor()
48         # 打开 CSV 文件并读取数据
49         with open(CSV_FILE_PATH, 'r') as file:
50             csv_reader = csv.reader(file)
51             next(csv_reader) # 跳过标题行
52             for row in csv_reader:
53                 print("正在写入行: ", row) # 输出插入数据
54                 cursor.execute("""
55                     INSERT INTO orders (id, sku, kinds, stime, num,
unit_price, total_price, longitude, latitude)
56                     VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)

```

```

57         """", row)
58         connection.commit()
59         print("CSV 数据成功导入到表 orders 中!")
60         cursor.close()
61         connection.close()
62     except Exception as e:
63         print("导入 CSV 数据失败: ", e)
64 def find_max_three_same_kinds():
65     # """查找三同订单中数量最多的商品种类"""
66     try:
67         connection = pymysql.connect(host=DB_HOST, user=DB_USER, password=DB_PASSWORD, database=DB_NAME)
68         cursor = connection.cursor()
69         # 执行 SQL 查询
70         query = """
71             SELECT kinds, COUNT(*) AS order_count
72             FROM orders
73             GROUP BY kinds, stime, longitude, latitude
74             ORDER BY order_count DESC
75             LIMIT 1;
76         """
77         cursor.execute(query)
78         result = cursor.fetchone()
79         if result:
80             kinds, order_count = result
81             print(f"三同订单最多的商品种类是: {kinds}, 订单数量是: {order_count}")
82         else:
83             print("查询为空: 没有找到符合条件的三同订单。")
84         cursor.close()
85         connection.close()
86     except Exception as e:
87         print("查询三同订单失败: ", e)
88 if __name__ == "__main__":
89     create_database_and_table()
90     import_csv_to_database()
91     find_max_three_same_kinds()

```

36|3 吴天行 第二题

```

1  import sqlite3
2  import csv
3  # 创建数据库连接
4  conn = sqlite3.connect('online_store.db')
5  cursor = conn.cursor()
6  # 创建orders表
7  cursor.execute('''
8  CREATE TABLE IF NOT EXISTS orders (

```

```

9      id VARCHAR(20),
10     sku VARCHAR(15),
11     kinds VARCHAR(15),
12     stime VARCHAR(15),
13     num INT,
14     unit_price FLOAT,
15     total_price FLOAT,
16     longitude FLOAT,
17     latitude FLOAT
18 )
19 '''
20 # 读取Demo_1.csv文件并插入数据到orders表
21 with open('/mnt/Demo_1.csv', 'r', encoding='utf-8') as file:
22     csv_reader = csv.reader(file)
23     next(csv_reader) # 跳过表头
24     for row in csv_reader:
25         cursor.execute('''
26             INSERT INTO orders (id, sku, kinds, stime, num, unit_price,
27             total_price, longitude, latitude)
28             VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?)
29             ''', row)
30 # 提交事务
31 conn.commit()
32 # 查找三同订单最多的商品种类
33 cursor.execute('''
34 SELECT kinds, COUNT(*) AS order_count
35 FROM (
36     SELECT kinds, stime, longitude, latitude
37     FROM orders
38     GROUP BY stime, longitude, latitude, kinds
39     HAVING COUNT(*) > 1
40 ) AS subquery
41 GROUP BY kinds
42 ORDER BY order_count DESC
43 LIMIT 1
44 ''')
45 # 获取结果
46 result = cursor.fetchone()
47 if result:
48     print(f"三同订单最多的商品种类是: {result[0]}, 订单数量为: {result[1]}")
49 else:
50     print("没有找到符合条件的三同订单")
51 # 关闭数据库连接
52 conn.close()

```



```
1 CREATE DATABASE online_store;
2 USE online_store;
3 CREATE TABLE orders (
4     id VARCHAR(20),
5     sku VARCHAR(15),
6     kinds VARCHAR(15),
7     stime VARCHAR(15),
8     num INT,
9     unit_price FLOAT,
10    total_price FLOAT,
11    longitude FLOAT,
12    latitude FLOAT
13 );
14 LOAD DATA INFILE 'C:/Users/Stu/Desktop/Demo_2.csv'
15 INTO TABLE orders
16 FIELDS TERMINATED BY ','
17 ENCLOSED BY '"'
18 LINES TERMINATED BY '\r\n'
19 IGNORE 1 LINES
20 (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude);
21 SELECT * FROM orders LIMIT 10;
22 WITH ThreeSameOrders AS (
23     -- 按照时间、地点和商品进行分组，计算每个组的订单数量
24     SELECT
25         order_date,
26         location,
27         product,
28         COUNT(*) AS order_count
29     FROM
30         orders
31     GROUP BY
32         order_date,
33         location,
34         product
35     HAVING
36         COUNT(*) >= 3 -- 只选择订单数量大于等于3的组
37 )
38 -- 统计三同订单中每个商品的出现频率
39 SELECT
40     product,
41     COUNT(*) AS three_same_order_count
42 FROM
43     ThreeSameOrders
44 GROUP BY
45     product
46 ORDER BY
47     three_same_order_count DESC
48 LIMIT 1; -- 选出三同订单最多的商品种类
```

38|5 塔巴江村 第二题

39|6 崔杰 第二题

```
1  在query中的代码:
2  CREATE TABLE `online_store`.`orders` (
3      `id` VARCHAR(20) NOT NULL,
4      `sku` VARCHAR(15) NULL,
5      `kinds` VARCHAR(15) NULL,
6      `stime` VARCHAR(15) NULL,
7      `num` INT NULL,
8      `unit_price` FLOAT NULL,
9      `total_price` FLOAT NULL,
10     `longitude` FLOAT NULL,
11     `latitude` FLOAT NULL,
12     PRIMARY KEY (`id`));
13  在orders中的代码:
14  SELECT * FROM online_store.orders;
15  SELECT kinds, COUNT(*) AS order_count
16  FROM (
17      SELECT kinds, stime, longitude, latitude
18      FROM orders
19      GROUP BY kinds, stime, longitude, latitude
20      HAVING COUNT(*) > 1
21  ) AS subquery
22  GROUP BY kinds
23  ORDER BY order_count DESC
24  LIMIT 1;
```

40|7 庄嘉帆 第二题

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 19:04:00 2024
4  @author: Stu
5  """
6  import mysql.connector
7  # 连接到MySQL数据库
8  mydb = mysql.connector.connect(
9      host="localhost",
10     user="root",
11     password="Qwe_1234"
12 )
13 # 创建游标对象
```

```
14 mycursor = mydb.cursor()
15 mycursor.execute("CREATE DATABASE IF NOT EXISTS online_store")
16 mydb.commit()
17 mydb = mysql.connector.connect(
18     host="localhost",
19     user="root",
20     password="Qwe_1234",
21     database="online_store"
22 )
23 mycursor = mydb.cursor()
24 # 创建orders表（如果不存在）
25 mycursor.execute("""
26 CREATE TABLE IF NOT EXISTS orders (
27     id VARCHAR(20),
28     sku VARCHAR(15),
29     kinds VARCHAR(15),
30     stime VARCHAR(15),
31     num INT,
32     unit_price FLOAT,
33     total_price FLOAT,
34     longitude FLOAT,
35     latitude FLOAT
36 )
37 """)
38 # 从Demo_1.csv文件中读取数据并插入到orders表中
39 try:
40     with open(r"C:\Users\Stu\Documents\WeChat Files\wxid_do87lzskr1z
41 t12\FileStorage\File\2024-12\Demo_2.csv") as file:
42         # 跳过标题行
43         next(file)
44         for line in file:
45             values = line.strip().split(',')
46             sql = "INSERT INTO orders (id, sku, kinds, stime, num, u
47 nit_price, total_price, longitude, latitude) VALUES (%s, %s, %s, %s,
48 %s, %s, %s, %s, %s)"
49             mycursor.execute(sql, values)
50             mydb.commit()
51 except FileNotFoundError:
52     print("文件未找到，请检查文件路径。")
53 # 查询三同订单最多的商品种类kinds
54 # 查询三同订单最多的商品种类kinds
55 mycursor.execute("""
56 SELECT subquery.kinds, COUNT(*) AS count
57 FROM (
58     -- 子查询，先找出符合三同条件的订单组
59     SELECT o1.id, o1.sku, o1.kinds, o1.stime, o1.longitude, o1.latit
60 ude,
61         GROUP_CONCAT(o2.id) AS related_orders
62 FROM orders o1
```

```

59 JOIN orders o2 ON o1.sku = o2.sku AND o1.stime = o2.stime AND o
60 1.longitude = o2.longitude AND o1.latitude = o2.latitude AND o1.id!=
    o2.id
    GROUP BY o1.id, o1.sku, o1.kinds, o1.stime, o1.longitude, o1.lat
61 itude
    HAVING COUNT(*) >= 1 -- 至少有两个订单符合三同条件（包含自身和至少一个
62 其他订单）
63 ) AS subquery
64 GROUP BY subquery.kinds
65 ORDER BY count DESC
66 LIMIT 1
67 "")
68 result = mycursor.fetchone()
69 if result:
70     print(f"三同订单最多的商品种类是: {result[0]}, 数量为: {result[1]}")
71 else:
72     print("未找到三同订单。")
73 # 关闭游标和数据库连接
74 mycursor.close()
    mydb.close()

```

41|8 张浩祖 第二题

```

1  #(1)
2  import sqlite3
3  import pandas as pd
4  # 创建SQLite数据库连接
5  conn = sqlite3.connect('online_store.db')
6  # 创建游标对象
7  cursor = conn.cursor()
8  # 创建orders表格
9  cursor.execute('''
10      CREATE TABLE IF NOT EXISTS orders (
11          id VARCHAR(20) PRIMARY KEY,
12          sku VARCHAR(15),
13          kinds VARCHAR(15),
14          stime VARCHAR(15),
15          num INT,
16          unit_price FLOAT,
17          total_price FLOAT,
18          longitude FLOAT,
19          latitude FLOAT
20      )
21  ''')
22  # 读取Demo_1.csv文件
23  df = pd.read_csv('Demo_2(1).csv')
24  # 将数据插入orders表格
25  df.to_sql('orders', conn, if_exists='replace', index=False)

```

21

5

5

```

26 # 提交事务
27 conn.commit()
28 # 查看orders表格的基本信息和前几行
29 print('orders表格基本信息: ')
30 cursor.execute('PRAGMA table_info(orders)')
31 for column in cursor.fetchall():
32     print(column)
33 # 查看orders表格的前几行内容
34 print('orders表格前几行内容: ')
35 cursor.execute('SELECT * FROM orders LIMIT 5')
36 for row in cursor.fetchall():
37     print(row)
38 # 关闭游标和连接
39 cursor.close()
40 conn.close()
41 # (2)
42 import sqlite3
43 # 创建SQLite数据库连接
44 conn = sqlite3.connect('online_store.db')
45 # 创建游标对象
46 cursor = conn.cursor()
47 # 查询三同订单中最多的商品种类kinds
48 cursor.execute('''
49     SELECT kinds, COUNT(*) AS count
50     FROM orders
51     GROUP BY stime, longitude, latitude, kinds
52     ORDER BY count DESC
53     LIMIT 1
54 ''')
55 # 获取查询结果
56 result = cursor.fetchone()
57 # 输出结果
58 if result:
59     print('三同订单中最多的商品种类为: ', result[0])
60 else:
61     print('没有找到三同订单。')
62 # 关闭游标和连接
63 cursor.close()
64 conn.close()

```

42|9 张露丹 第二题

```

1 -- 创建数据库
2 CREATE DATABASE online_store;
3 -- 使用数据库
4 USE online_store;
5 -- 创建orders表
6 CREATE TABLE orders (

```

```
7      id VARCHAR(20) PRIMARY KEY,
8      sku VARCHAR(15),
9      kinds VARCHAR(15),
10     stime VARCHAR(15),
11     num INT,
12     unit_price FLOAT,
13     total_price FLOAT,
14     longitude FLOAT,
15     latitude FLOAT
16 );
17 LOAD DATA LOCAL INFILE 'C:/Users/Stu/Desktop/12.20/2_Demo_2.csv'
18 INTO TABLE orders
19 FIELDS TERMINATED BY ','
20 ENCLOSED BY '"'
21 LINES TERMINATED BY '\n'
22 IGNORE 1 ROWS
23 (id, sku, kinds, stime, num, unit_price, total_price, longitude, lat
it
itude);
24 WITH OrderGroups AS (
25     SELECT
26         kinds,
27         DATE_FORMAT(STR_TO_DATE(stime, '%Y-%m-%d %H:%i:%s'), '%Y-%m-
%d %H:%i') AS stime_rounded, -- 假设时间格式，并四舍五入到分钟
28         ROUND(longitude, 6) AS longitude_rounded, -- 四舍五入经纬度到小
数点后6位
29         ROUND(latitude, 6) AS latitude_rounded,
30         COUNT(*) AS order_count
31     FROM
32         orders
33     GROUP BY
34         kinds,
35         stime_rounded,
36         longitude_rounded,
37         latitude_rounded
38     HAVING
39         COUNT(*) >= 3 -- 只考虑至少3个订单的组
40 )
41 SELECT
42     kinds,
43     MAX(order_count) AS max_order_count
44 FROM
45     OrderGroups
46 GROUP BY
47     kinds
48 ORDER BY
49     max_order_count DESC
50 LIMIT 1; -- 只取最多的那个商品种类
```

43|10 施習 第二题

44|11 李上卫 第二题

45|12 李俊杰 第二题

23

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 19:08:37 2024
4  @author: Stu
5  """
6  import sqlite3
7  import pandas as pd
8  # 连接到SQLite数据库（如果不存在会自动创建）
9  conn = sqlite3.connect('online_store.db')
10 # 创建orders表
11 create_table_sql = """
12 CREATE TABLE IF NOT EXISTS orders (
13     id VARCHAR(20),
14     sku VARCHAR(15),
15     kinds VARCHAR(15),
16     stime VARCHAR(15),
17     num INT,
18     unit_price FLOAT,
19     total_price FLOAT,
20     longitude FLOAT,
21     latitude FLOAT
22 );
23 """
24 cursor = conn.cursor()
25 cursor.execute(create_table_sql)
26 conn.commit()
27 # 使用pandas读取CSV文件
28 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_2.csv')
29 # 将数据插入到orders表中
30 df.to_sql('orders', conn, if_exists='append', index=False)
31 # 关闭连接
32 conn.close()
33 import sqlite3
34 import pandas as pd
35 # 连接数据库
36 conn = sqlite3.connect('online_store.db')
37 # 从数据库中读取orders表数据到DataFrame
38 query = "SELECT * FROM orders"
39 df_orders = pd.read_sql(query, conn)
40 # 按照时间、地点、商品种类分组，并统计每组的订单数量，筛选出订单数量大于1的组
    （即三同订单）
```

5

6

13

```

41 grouped = df_orders.groupby(['stime', 'longitude', 'latitude', 'kind
42 s']).size().reset_index(name='count')
43 three_same_orders = grouped[grouped['count'] > 1]
44 # 按照商品种类再次分组，统计每个商品种类的三同订单数量
45 result = three_same_orders.groupby('kinds')['count'].sum().reset_in
46 dex()
47 # 按照三同订单数量降序排序，取第一个（即三同订单最多的商品种类）
48 max_kind = result.sort_values(by='count', ascending=False).iloc[0]
49 ['kinds']
50 print(max_kind)
51 # 关闭连接
52 conn.close()

```

46|13 杨晨晨 第二题

```

1 import pandas as pd
2 import sqlite3
3 # 连接到SQLite数据库（如果不存在会创建新的），这里数据库文件名为online_store.
4 db
5 conn = sqlite3.connect('online_store.db')
6 # 创建orders表的SQL语句
7 create_table_sql = """
8 CREATE TABLE IF NOT EXISTS orders (
9     id varchar(20),
10     sku varchar(15),
11     kinds varchar(15),
12     stime varchar(15),
13     num int,
14     unit_price float,
15     total_price float,
16     longitude float,
17     latitude float
18 );
19 """
20 # 获取游标对象，用于执行SQL语句
21 cursor = conn.cursor()
22 # 执行创建表的SQL语句
23 cursor.execute(create_table_sql)
24 # 使用pandas读取Demo_2.csv文件数据
25 data = pd.read_csv('Demo_2.csv')
26 # 将数据插入到orders表中，如果表已存在会追加数据（注意这里的数据类型要和表结构尽
27 量匹配，示例中按简单对应处理）
28 data.to_sql('orders', conn, if_exists='append', index=False)
29 # 提交事务（确保数据插入成功）
30 conn.commit()
31 # 关闭游标和连接
32 cursor.close()
33 conn.close()

```



```

33 import pandas as pd
34 import sqlite3
35 # 连接到数据库
36 conn = sqlite3.connect('online_store.db')
37 # 编写SQL查询语句，思路和之前纯SQL实现类似，这里通过子查询先找出不同的"三同"组合，再统计每个组合出现次数，最后找出次数最多的商品种类
38 query = """
39 SELECT kinds
40 FROM (
41     SELECT kinds, COUNT(*) AS count_orders
42     FROM (
43         SELECT DISTINCT stime, longitude, latitude, sku, kinds
44         FROM orders
45     ) AS distinct_orders
46     GROUP BY stime, longitude, latitude, sku
47     ORDER BY count_orders DESC
48     LIMIT 1
49 ) AS most_frequent_kinds;
50 """
51 # 使用pandas的read_sql函数执行查询并获取结果
52 result = pd.read_sql(query, conn)
53 # 输出结果（这里假设结果只有一行一列，即符合条件的那个商品种类）
54 print(result.iloc[0, 0])
55 # 关闭连接
conn.close()

```

47/14 杨翔 第二题

```

1 CREATE DATABASE online_store;
2 USE online_store;
3 CREATE TABLE orders (
4     id VARCHAR(20) PRIMARY KEY,
5     sku VARCHAR(15),
6     kinds VARCHAR(15),
7     stime VARCHAR(15),
8     num INT,
9     unit_price FLOAT,
10    total_price FLOAT,
11    longitude FLOAT,
12    latitude FLOAT
13 );
14 LOAD DATA INFILE '/var/lib/mysql-files/Demo_1.csv'
15 INTO TABLE orders
16 FIELDS TERMINATED BY ','
17 ENCLOSED BY '"'
18 LINES TERMINATED BY '\n'
19 IGNORE 1 ROWS
20 (id, sku, kinds, stime, num, unit_price, total_price, longitude, lat

```

```

itute);
21 WITH TriplicateOrders AS (
22     SELECT
23         kinds,
24         stime,
25         longitude,
26         latitude,
27         COUNT(*) AS order_count
28     FROM
29         orders
30     GROUP BY
31         kinds,
32         stime,
33         longitude,
34         latitude
35     HAVING
36         COUNT(*) > 1
37 )
38 SELECT
39     kinds,
40     MAX(order_count) AS max_order_count
41 FROM
42     TriplicateOrders
43 GROUP BY
44     kinds
45 ORDER BY
46     max_order_count DESC
47 LIMIT 1;

```

14

48|15 王丽媛 第二题

```


1 # 2
2 # (1)
3 import mysql.connector
4 import csv
5 mydb = mysql.connector.connect(
6     host="localhost",
7     user="root",
8     password="Qwe_1234"
9 )
10 mycursor = mydb.cursor()
11 mycursor.execute("CREATE DATABASE IF NOT EXISTS online_store")
12 mycursor.execute("USE online_store")
13 create_table_query = ""
14 CREATE TABLE IF NOT EXISTS orders (
15     id varchar(20),
16     sku varchar(15),
17     kinds varchar(15),

```

23


5

```
18         stime varchar(15),
19         num int,
20         unit_price float,
21         total_price float,
22         longitude float,
23         latitude float
24     )
25     """
26     mycursor.execute(create_table_query)
27     with open('D:Demo_2.csv', 'r', encoding='utf-8') as csvfile:
28         csvreader = csv.reader(csvfile)
29         next(csvreader)
30         for row in csvreader:
31             insert_query = "INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude) VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)"
32             mycursor.execute(insert_query, tuple(row))
33     mydb.commit()
34     mycursor.close()
35     mydb.close()
36     # (2)
37     import mysql.connector
38     mydb = mysql.connector.connect(
39         host="localhost",
40         user="root",
41         password="Qwe_1234",
42         database="online_store"
43     )
44     mycursor = mydb.cursor()
45     query = """
46     SELECT kinds, COUNT(*) as count
47     FROM (
48         SELECT stime, longitude, latitude, sku, kinds
49         FROM orders
50         GROUP BY stime, longitude, latitude, sku, kinds
51         HAVING COUNT(*) > 1  -- 只考虑出现次数大于1的组合，即三同订单情况
52     ) as subquery
53     GROUP BY kinds
54     ORDER BY count DESC
55     LIMIT 1;
56     """
57     mycursor.execute(query)
58     result = mycursor.fetchone()
59     if result:
60         print("三同订单最多的商品种类为:", result[0])
61     else:
62         print("未发现三同订单")
63     mycursor.close()
64
```



49|16 王乐 第二题



50|17 王康宇 第二题



```

1  import pandas as pd
2  import pymysql
3  # 读取csv文件
4  data_csv = pd.read_csv('Demo_2.csv',encoding="gbk")
5  # 连接到MySQL数据库
6  conn = pymysql.connect(host='localhost', user='root', password='123456', port=3306)
7  # 创建游标对象
8  cursor = conn.cursor()
9  # 创建online_store数据库
10 cursor.execute('CREATE DATABASE IF NOT EXISTS online_store')
11 # 选择online_store数据库
12 cursor.execute('USE online_store')
13 # 创建orders表
14 cursor.execute('''
15 CREATE TABLE IF NOT EXISTS orders (
16     id VARCHAR(20),
17     sku VARCHAR(15),
18     kinds VARCHAR(15),
19     stime VARCHAR(15),
20     num INT,
21     unit_price FLOAT,
22     total_price FLOAT,
23     longitude FLOAT,
24     latitude FLOAT
25 )
26 ''')
27 # 将数据插入orders表中
28 for row in data_csv.values.tolist():
29     cursor.execute('''
30         INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude)
31         VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)
32     ''', row)
33 query = "SELECT * FROM orders LIMIT 10"
34 cursor.execute(query)
35 # 获取查询结果
36 results = cursor.fetchall()
37 # 打印结果
38 for row in results:

```

```

39     print(row)
40     # 提交事务
41     conn.commit()
42     # 查询orders表中的三同订单
43     query = '''
44     SELECT kinds, stime, longitude, latitude, COUNT(*) AS order_count
45     FROM orders
46     GROUP BY kinds, stime, longitude, latitude
47     HAVING COUNT(*) > 2
48     ORDER BY order_count DESC
49     LIMIT 1
50     '''
51     # 执行查询
52     cursor.execute(query)
53     # 获取查询结果
54     result = cursor.fetchone()
55     # 输出三同订单最多的商品种类
56     if result:
57         print('三同订单最多的商品种类为:', result[0])
58     else:
59         print('没有找到三同订单')
60     # 关闭游标和连接
61     cursor.close()
62     conn.close()

```

51|18 王立弘 第二题

```

1  import csv
2  import sqlite3
3  # 连接到SQLite数据库（如果不存在会自动创建）
4  conn = sqlite3.connect('online_store.db')
5  # 创建游标对象，用于执行SQL语句
6  cursor = conn.cursor()
7  # 创建orders表的SQL语句
8  create_table_sql = """
9  CREATE TABLE IF NOT EXISTS orders (
10     id VARCHAR(20),
11     sku VARCHAR(15),
12     kinds VARCHAR(15),
13     stime VARCHAR(15),
14     num INTEGER,
15     unit_price REAL,
16     total_price REAL,
17     longitude REAL,
18     latitude REAL
19 )
20 """
21 cursor.execute(create_table_sql)

```

```

22 # 读取Demo_2.csv文件并插入数据到orders表
23 with open('D:\\Demo_2.csv', 'r', encoding='utf-8') as csvfile:
24     reader = csv.reader(csvfile)
25     next(reader) # 跳过标题行（假设文件有标题行）
26     for row in reader:
27         insert_sql = "INSERT INTO orders (id, sku, kinds, stime, nu
m, unit_price, total_price, longitude, latitude) VALUES
        (?, ?, ?, ?, ?, ?, ?, ?, ?)"
28         cursor.execute(insert_sql, tuple(row))
29 # 提交事务，使插入操作生效
30 conn.commit()
31 # 查询orders表中每个 (stime, longitude, latitude, kinds) 组合的出现次数
32 query_sql = """
33 SELECT stime, longitude, latitude, kinds, COUNT(*) as count
34 FROM orders
35 GROUP BY stime, longitude, latitude, kinds
36 """
37 cursor.execute(query_sql)
38 results = cursor.fetchall()
39 # 用于统计每个商品种类对应的三同订单数量
40 kinds_count_dict = {}
41 for row in results:
42     _, _, _, kinds, count = row
43     if kinds in kinds_count_dict:
44         kinds_count_dict[kinds] += count
45     else:
46         kinds_count_dict[kinds] = count
47 # 找出三同订单数量最多的商品种类
48 max_kinds = max(kinds_count_dict, key=kinds_count_dict.get)
49 print(f"三同订单最多的商品种类是: {max_kinds}")
50 # 关闭游标和数据库连接
51 cursor.close()
52 conn.close()
53

```

52|19 白天琪 第二题

```

1 import pandas as pd
2 import mysql.connector
3 # 1. 连接到 MySQL 数据库并创建数据库和表格
4 def create_database_and_table():
5     # 连接到 MySQL 数据库
6     connection = mysql.connector.connect(
7         host="localhost", # 数据库服务器地址
8         user="root", # 数据库用户名
9         password="123456" # 数据库密码
10    )
11    cursor = connection.cursor()

```

```
12 # 创建数据库 online_store
13 cursor.execute("CREATE DATABASE IF NOT EXISTS online_store;")
14 cursor.execute("USE online_store;")
15 # 创建 orders 表
16 create_table_query = """
17 CREATE TABLE IF NOT EXISTS orders (
18     id VARCHAR(20),
19     sku VARCHAR(15),
20     kinds VARCHAR(15),
21     stime VARCHAR(15),
22     num INT,
23     unit_price FLOAT,
24     total_price FLOAT,
25     longitude FLOAT,
26     latitude FLOAT
27 );
28 """
29 cursor.execute(create_table_query)
30 # 提交并关闭连接
31 connection.commit()
32 cursor.close()
33 connection.close()
34 print("数据库和表格已成功创建!")
35 # 2. 读取 CSV 文件并将数据插入到数据库
36 def insert_data_from_csv(csv_file_path):
37     # 读取 CSV 文件
38     df = pd.read_csv(csv_file_path)
39     # 连接到 MySQL 数据库
40     connection = mysql.connector.connect(
41         host="localhost",
42         user="root",
43         password="123456",
44         database="online_store" # 使用 online_store 数据库
45     )
46     cursor = connection.cursor()
47     # 插入数据的 SQL 语句
48     insert_query = """
49     INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude)
50     VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s);
51     """
52     # 遍历 DataFrame 中的每一行，将数据插入到数据库表格中
53     for index, row in df.iterrows():
54         cursor.execute(insert_query, (
55             row['id'], row['sku'], row['kinds'], row['stime'],
56             row['num'], row['unit_price'], row['total_price'],
57             row['longitude'], row['latitude']
58         ))
59     # 提交事务并关闭连接
```

```
60     connection.commit()
61     cursor.close()
62     connection.close()
63     print("数据已成功插入到数据库表格 orders 中!")
64 # 主函数
65 def main():
66     # 创建数据库和表格
67     create_database_and_table()
68     # CSV 文件路径
69     csv_file_path = "Demo_2.csv" # 请替换为实际的 CSV 文件路径
70     # 插入数据到表格
71     insert_data_from_csv(csv_file_path)
72 if __name__ == "__main__":
73     main()
74 import mysql.connector
75 # 连接到 MySQL 数据库
76 def connect_to_database():
77     connection = mysql.connector.connect(
78         host="localhost", # 数据库服务器地址
79         user="root", # 数据库用户名
80         password="123456", # 数据库密码
81         database="online_store" # 使用 online_store 数据库
82     )
83     return connection
84 # 查询并统计三同订单最多的商品种类
85 def get_most_common_kinds():
86     connection = connect_to_database()
87     cursor = connection.cursor()
88     # 查询语句: 按时间、地点、商品种类分组, 并计算每组的订单数量
89     query = """
90     SELECT stime, longitude, latitude, kinds, COUNT(*) AS order_coun
91 t
92     FROM orders
93     GROUP BY stime, longitude, latitude, kinds
94     HAVING COUNT(*) > 1 -- 确保有多个相同的订单
95     ORDER BY order_count DESC
96     LIMIT 1; -- 获取订单数量最多的商品种类
97     """
98     cursor.execute(query)
99     result = cursor.fetchone()
100    cursor.close()
101    connection.close()
102    # 输出结果
103    if result:
104        stime, longitude, latitude, kinds, order_count = result
105        print(f"最多的三同订单商品种类: {kinds}, 订单数量: {order_coun
106 t}")
107    else:
108        print("没有符合三同条件的订单。")
```




```
109 # 主函数
110 def main():
111     get_most_common_kinds()
112 if __name__ == "__main__":
    main()
```

53|20 罗昊然 第二题

```
1 import pandas as pd
2 import pymysql
3 data_frame = pd.read_csv('Demo_2.csv', encoding="gbk")
4 database_connection = pymysql.connect(host='localhost', user='root',
5 password='Qwe_1234', port=3306)
6 cursor_object = database_connection.cursor()
7 cursor_object.execute('CREATE DATABASE IF NOT EXISTS online_store')
8 cursor_object.execute('USE online_store')
9 cursor_object.execute('''
10 CREATE TABLE IF NOT EXISTS orders (
11     id VARCHAR(20),
12     sku VARCHAR(15),
13     kinds VARCHAR(15),
14     stime VARCHAR(15),
15     num INT,
16     unit_price FLOAT,
17     total_price FLOAT,
18     longitude FLOAT,
19     latitude FLOAT
20 )
21 ''')
22 # 将CSV文件中的数据逐行插入到orders表中
23 for record in data_frame.itertuples(index=False):
24     cursor_object.execute('''
25         INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude)
26         VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)
27     ''', record)
28 select_query = "SELECT * FROM orders LIMIT 10"
29 cursor_object.execute(select_query)
30 query_results = cursor_object.fetchall()
31 for result in query_results:
32     print(result)
33 database_connection.commit()
34 # 查询orders表中出现次数最多的三同订单（相同的种类、时间、经度和纬度）
35 grouped_query = '''
36 SELECT kinds, stime, longitude, latitude, COUNT(*) AS order_count
37 FROM orders
38 GROUP BY kinds, stime, longitude, latitude
```

```

38     HAVING COUNT(*) > 2
39     ORDER BY order_count DESC
40     LIMIT 1
41     '''
42     cursor_object.execute(grouped_query)
43     grouped_result = cursor_object.fetchone()
44     if grouped_result:
45         print('三同订单最多的商品种类为:', grouped_result[0])
46     else:
47         print('没有找到三同订单')
48     cursor_object.close()
49     database_connection.close()

```

54|21 谢皓椿 第二题

```

1     import sqlite3
2     import csv
3     # 创建数据库连接
4     conn = sqlite3.connect('online_store.db')
5     cursor = conn.cursor()
6     # 创建orders表
7     cursor.execute('''
8     CREATE TABLE IF NOT EXISTS orders (
9         id VARCHAR(20),
10        sku VARCHAR(15),
11        kinds VARCHAR(15),
12        stime VARCHAR(15),
13        num INT,
14        unit_price FLOAT,
15        total_price FLOAT,
16        longitude FLOAT,
17        latitude FLOAT
18    )
19    ''')
20    # 读取Demo_1.csv文件并插入数据到orders表
21    with open('/mnt/Demo_1.csv', 'r', encoding='utf-8') as file:
22        csv_reader = csv.reader(file)
23        next(csv_reader) # 跳过表头
24        for row in csv_reader:
25            cursor.execute('''
26                INSERT INTO orders (id, sku, kinds, stime, num, unit_price,
27                total_price, longitude, latitude)
28                VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?)
29                ''', row)
30    # 提交事务
31    conn.commit()
32    # 查找三同订单最多的商品种类
33    cursor.execute('''

```

```

33 SELECT kinds, COUNT(*) AS order_count
34 FROM (
35     SELECT kinds, stime, longitude, latitude
36     FROM orders
37     GROUP BY stime, longitude, latitude, kinds
38     HAVING COUNT(*) > 1
39 ) AS subquery
40 GROUP BY kinds
41 ORDER BY order_count DESC
42 LIMIT 1
43 '''
44 # 获取结果
45 result = cursor.fetchone()
46 if result:
47     print(f"三同订单最多的商品种类是: {result[0]}, 订单数量为: {result
48 [1]}")
49 else:
50     print("没有找到符合条件的三同订单")
51 # 关闭数据库连接
conn.close()

```

55|22 贺馨姜艾 第二题

```

1 CREATE DATABASE IF NOT EXISTS online_store;
2 USE online_store;
3 CREATE TABLE IF NOT EXISTS orders (
4     id VARCHAR(20),
5     sku VARCHAR(15),
6     kinds VARCHAR(15),
7     stime VARCHAR(15),
8     num INT,
9     unit_price FLOAT,
10    total_price FLOAT,
11    longitude FLOAT,
12    latitude FLOAT
13 );
14 LOAD DATA INFILE 'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_2.csv'
15 INTO TABLE orders
16 FIELDS TERMINATED BY ','
17 ENCLOSED BY '"'
18 LINES TERMINATED BY '\n'
19 IGNORE 1 ROWS;
20 WITH OrderGroups AS (
21     SELECT
22         stime,
23         longitude,
24         latitude,
25         sku,

```

```

26         kinds,
27         COUNT(*) AS order_count
28     FROM
29         orders
30     GROUP BY
31         stime,
32         longitude,
33         latitude,
34         sku,
35         kinds
36     HAVING
37         COUNT(*) > 1
38 ),
39 MaxOrderKind AS (
40     SELECT
41         kinds,
42         SUM(order_count) AS total_order_count
43     FROM
44         OrderGroups
45     GROUP BY
46         kinds
47     ORDER BY
48         total_order_count DESC
49     LIMIT 1
50 )
51 SELECT
52     kinds
53 FROM
54     MaxOrderKind;

```

12

✓

56|23 达尔汗 第二题

```

1  import pymysql
2  import csv
3  # 建立数据库连接
4  def connect_mysql():
5      try:
6          conn = pymysql.connect(
7              host='localhost', # 这里需要根据你的实际MySQL主机地址修改，通
常本地为localhost
8              user='root', # 改成你的MySQL用户名
9              password='Qwe_1234', # 修改为你的实际密码
10             database='', # 初始不指定具体数据库
11             charset='utf8'
12         )
13         return conn
14     except pymysql.Error as e:
15         print(f"连接MySQL数据库出现错误: {e}")

```

23

```
16         return None
17 # 创建online_store数据库
18 def create_database(conn):
19     cursor = conn.cursor()
20     try:
21         cursor.execute("CREATE DATABASE IF NOT EXISTS online_store")
22         print("online_store数据库创建成功（如果不存在的话）")
23     except pymysql.Error as e:
24         print(f"创建数据库时出错: {e}")
25     finally:
26         cursor.close()
27 # 使用online_store数据库
28 def use_database(conn):
29     cursor = conn.cursor()
30     try:
31         cursor.execute("USE online_store")
32     except pymysql.Error as e:
33         print(f"切换使用数据库时出错: {e}")
34     finally:
35         cursor.close()
36 # 创建orders表，添加字段备注
37 def create_orders_table(conn):
38     cursor = conn.cursor()
39     create_table_sql = """
40     CREATE TABLE IF NOT EXISTS orders (
41         id varchar(20) COMMENT '订单号',
42         sku varchar(15) COMMENT '商品货号',
43         kinds varchar(15) COMMENT '种类',
44         stime varchar(15) COMMENT '销售时间',
45         num int COMMENT '数量',
46         unit_price float COMMENT '单价',
47         total_price float COMMENT '总价',
48         longitude float COMMENT '收货地经度',
49         latitude float COMMENT '纬度'
50     ) ENGINE=InnoDB DEFAULT CHARSET=utf8;
51     """
52     try:
53         cursor.execute(create_table_sql)
54         print("orders表创建成功（如果不存在的话）")
55     except pymysql.Error as e:
56         print(f"创建orders表时出错: {e}")
57     finally:
58         cursor.close()
59 # 从CSV文件读取数据并插入到orders表
60 def insert_data_from_csv(conn):
61     file_path = "D:/Demo_2.csv"
62     cursor = conn.cursor()
63     insert_sql = "INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude) VALUES (%s, %s, %s, %s,
```

```

64     %s, %s, %s, %s, %s)"
65     try:
66         with open(file_path, 'r', encoding='utf-8') as csvfile:
67             reader = csv.reader(csvfile)
68             next(reader) # 跳过文件头行（标题行）
69             for row in reader:
70                 cursor.execute(insert_sql, tuple(row))
71             conn.commit()
72             print("数据已成功插入到orders表中")
73     except (pymysql.Error, FileNotFoundError) as e:
74         print(f"插入数据出现错误: {e}")
75     finally:
76         cursor.close()
77 if __name__ == "__main__":
78     # 连接数据库
79     connection = connect_mysql()
80     if connection:
81         # 创建数据库
82         create_database(connection)
83         # 使用数据库
84         use_database(connection)
85         # 创建orders表
86         create_orders_table(connection)
87         # 插入数据
88         insert_data_from_csv(connection)
89         # 关闭连接
90         connection.close()
91 # 建立数据库连接
92 def connect_mysql():
93     try:
94         conn = pymysql.connect(
95             host='localhost',
96             user='root',
97             password='Qwe_1234',
98             database='online_store', # 这里使用前面创建好的online_store数据库
99             charset='utf8'
100         )
101         return conn
102     except pymysql.Error as e:
103         print(f"连接MySQL数据库出现错误: {e}")
104         return None
105 # 查询并统计三同订单中商品种类的出现次数
106 def count_triple_same_orders(conn):
107     cursor = conn.cursor()
108     query_sql = """
109     SELECT kinds, COUNT(*) as count
110     FROM (
111         SELECT stime, longitude, latitude, kinds

```

```

112         FROM orders
113         GROUP BY stime, longitude, latitude, kinds
114     ) as subquery
115     GROUP BY kinds
116     ORDER BY count DESC
117     LIMIT 1;
118     """
119     try:
120         cursor.execute(query_sql)
121         result = cursor.fetchone()
122         if result:
123             return result[0] # 返回三同订单最多的商品种类
124         return None
125     except pymysql.Error as e:
126         print(f"查询统计三同订单出现错误: {e}")
127     finally:
128         cursor.close()
129 if __name__ == "__main__":
130     connection = connect_mysql()
131     if connection:
132         most_kinds = count_triple_same_orders(connection)
133         if most_kinds:
134             print(f"三同订单最多的商品种类是: {most_kinds}")
135         else:
136             print("未找到三同订单相关数据")
137         connection.close()

```

57|24 郭嘉 第二题

58|25 郭欣遥 第二题

23

```

1  (2)
2  #第二题第一问:
3  import sqlite3
4  import pandas as pd
5  # 连接到 SQLite 数据库, 如果不存在则创建
6  conn = sqlite3.connect('online_store.db')
7  # 创建游标对象
8  cursor = conn.cursor()
9  # 创建 orders 表格的 SQL 语句
10 create_table_sql = """
11 CREATE TABLE IF NOT EXISTS orders (
12     id VARCHAR(20),
13     sku VARCHAR(15),
14     kinds VARCHAR(15),
15     stime VARCHAR(15),
16     num INT,

```

```
17     unit_price FLOAT,
18     total_price FLOAT,
19     longitude FLOAT,
20     latitude FLOAT
21 );
22 """
23 # 执行创建表的 SQL 语句
24 cursor.execute(create_table_sql)
25 # 读取 Demo_1.csv 文件数据到 DataFrame
26 df = pd.read_csv('Demo_2.csv')
27 # 将 DataFrame 数据插入到 orders 表中
28 df.to_sql('orders', conn, if_exists='append', index=False)
29 # 提交事务并关闭连接
30 conn.commit()
31 conn.close()
32 #第二题第二问:
33 import pandas as pd
34 # 读取数据
35 df = pd.read_csv('Demo_2.csv')
36 # 按照时间、地点和商品种类进行分组, 并统计每组的数量
37 grouped = df.groupby(['stime', 'longitude', 'latitude', 'kinds']).size().reset_index(name='count')
38 # 筛选出数量大于 1 的组, 即三同订单
39 filtered = grouped[grouped['count'] > 1]
40 # 按照商品种类对三同订单进行分组, 并计算每个种类的三同订单总数
41 result = filtered.groupby('kinds')['count'].sum().reset_index()
42 # 按照总数降序排序
43 result = result.sort_values(by='count', ascending=False)
44 # 输出三同订单最多的商品种类
45 if len(result) > 0:
46     print(result['kinds'].iloc[0])
47 else:
48     print("未找到三同订单")
49 #解题步骤
50 ### 第二题第一问: 创建数据库及插入数据
51 #1. **导入必要的库**:
52 # - 导入 `sqlite3` 库, 用于操作 SQLite 数据库。
53 # - 导入 `pandas` 库, 用于数据处理和分析, 方便将 CSV 文件数据读取并插入到数据库表中。
54 #2. **连接到数据库**:
55 # - 使用 `sqlite3.connect('online_store.db')` 语句连接到名为 `online_store.db` 的 SQLite 数据库。如果该数据库不存在, 此操作会创建一个新的数据库。
56 # - 创建一个游标对象 `cursor`, 通过 `conn.cursor()` 实现。游标用于执行 SQL 语句来操作数据库。
57 #3. **创建 `orders` 表**:
58 # - 定义一个包含创建 `orders` 表的 SQL 语句的字符串 `create_table_sql`。表结构包含 `id` (`varchar(20)` 类型, 用于存储订单号)、`sku` (`varchar(15)` 类型, 商品货号)、`kinds` (`varchar(15)` 类型, 商品种类)、`stime` (`varchar(15)` 类型, 销售时间)、`num` (`int` 类型, 数量)、`unit_price`
```



```

59     (`float` 类型, 单价)、`total_price` (`float` 类型, 总价)、`longitude`
60     (`float` 类型, 收货地经度)和 `latitude` (`float` 类型, 收货地纬度)等字
61     段。
62     # - 使用 `cursor.execute(create_table_sql)` 执行创建表的 SQL 语句, 在
63     数据库中创建 `orders` 表。
64     #4. **读取数据并插入表中**:
65     # - 使用 `pd.read_csv('Demo_2.csv')` 读取 `Demo_2.csv` 文件中的数据,
66     并将其存储在一个 `DataFrame` 对象 `df` 中。这里需要注意的是, 根据题目要求应该
67     读取 `Demo_1.csv` 文件, 代码中可能存在错误。
68     # - 调用 `df.to_sql('orders', conn, if_exists='append', index=False)` 将 `df` 中的数据插入到 `orders` 表中。`if_exists='append'` 表示如果表
69     已存在, 则将数据追加到表中; `index=False` 表示不将 `DataFrame` 的索引列插入
70     到表中。
71     #5. **提交事务并关闭连接**:
72     #- 执行 `conn.commit()` 提交事务, 确保数据插入操作的永久性。
73     # - 最后使用 `conn.close()` 关闭数据库连接, 释放资源。
74     ### 第二题第二问: 查找三同订单最多的商品种类
75     #1. **数据读取**:
76     # - 同样使用 `pandas` 库的 `pd.read_csv('Demo_2.csv')` 读取数据, 并将其
77     存储在 `DataFrame` 对象 `df` 中。
78     ##
79     # - 按照 `stime` (销售时间)、`longitude` (经度)、`latitude` (纬度)和 `k
80     inds` (商品种类)这四个列对 `df` 进行分组, 使用 `groupby` 方法实现。然后通过
81     `size` 函数计算每组的数量, 并使用 `reset_index` 方法将结果重置索引, 同时将计
82     数列命名为 `count`, 得到一个新的 `DataFrame` 对象 `grouped`。
83     #3. **筛选三同订单**:
84     # - 从 `grouped` 中筛选出 `count` 列大于 1 的组, 即三同订单。通过 `group
85     ed[grouped['count'] > 1]` 实现, 将筛选结果存储在 `filtered` 中。
86     #4. **计算每种商品的三同订单总数**:
87     # - 按照 `kinds` 列对 `filtered` 进行分组, 再次使用 `groupby` 方法。然后
88     通过 `sum` 函数计算每个种类的三同订单总数, 并使用 `reset_index` 方法重置索
89     引, 得到最终结果的 `DataFrame` 对象 `result`。
90     #5. **排序与输出结果**:
91     # - 使用 `result.sort_values(by='count', ascending=False)` 按照 `co
92     unt` 列的值对 `result` 进行降序排序。
93     # - 通过判断 `len(result)` 是否大于 0, 如果大于 0, 则输出 `result['kind
94     s'].iloc[0]`, 即三同订单最多的商品种类; 否则输出提示信息“未找到三同订单”。

```

59|26 阿依夏·克热木江 第二题

60|27 陈琰 第二题

```

1  import sqlite3
2  import pandas as pd
3  def create_database_and_insert_data():
4      try:

```

```
5 # 连接到SQLite数据库（如果不存在会自动创建名为online_store.db的数据
   库文件）
6 conn = sqlite3.connect('online_store.db')
7 cursor = conn.cursor()
8 # 创建orders表的SQL语句
9 create_table_sql = """
10 CREATE TABLE IF NOT EXISTS orders
11     id INTEGER PRIMARY KEY,
12     sku TEXT,
13     kinds TEXT,
14     stime TEXT,
15     num INTEGER,
16     unit_price REAL,
17     total_price REAL,
18     longitude REAL,
19     latitude REAL
20 )
21 """
22 cursor.execute(create_table_sql)
23 # 使用pandas读取Demo_2.csv文件为DataFrame
24 df = pd.read_csv('Demo_2.csv')
25 # 获取orders表的实际列信息，动态构建插入语句
26 cursor.execute("PRAGMA table_info(orders)")
27 column_info = cursor.fetchall()
28 column_names = [info[1] for info in column_info]
29 insert_sql = "INSERT INTO orders ({} ) VALUES ({} )".format
30 (",".join(column_names), ",".join(["?" for _ in column_names]))
31 # 将DataFrame数据按动态构建的插入语句插入到orders表中
32 data = df.values.tolist()
33 for row in data:
34     cursor.execute(insert_sql, row)
35 # 获取orders表的实际列信息，动态构建查询语句（仅查询存在的列）
36 cursor.execute("PRAGMA table_info(orders)")
37 column_info = cursor.fetchall()
38 column_names = [info[1] for info in column_info]
39 query_columns = []
40 for column in column_names:
41     if column == "stime":
42         query_columns.append(column)
43     else:
44         query_columns.append(column)
45 select_sql = "SELECT {} FROM orders".format(",".join(query_c
   olumns))
46 # 查询并输出部分数据（这里查询前10条数据作为示例，可按需调整）
47 cursor.execute(select_sql)
48 results = cursor.fetchall()
49 for row in results:
50     print(row)
51 # 关闭游标和连接
```

```

51         cursor.close()
52         conn.close()
53     except Exception as e:
54         print(f"出现错误: {e}")
55 if __name__ == "__main__":
56     create_database_and_insert_data()
57     #
58     import pandas as pd
59     # 读取csv文件
60     df = pd.read_csv('Demo_2.csv')
61     # 按照时间、地点、商品种类进行分组, 并统计每组的数量
62     grouped = df.groupby(['stime', 'latitude', 'kinds']).size().reset_index(name='count')
63     # 再次按照商品种类进行分组, 计算每个商品种类下符合三同的订单数量总和
64     result = grouped.groupby('kinds')['count'].sum().reset_index()
65     # 找到总和最大的商品种类
66     max_kind = result.loc[result['count'].idxmax(), 'kinds']
67     print(max_kind)
68

```

12 ✓

61|28 韦淑荣 第二题

```

1  二、
2  (1) -- 创建数据库
3  CREATE DATABASE online_store;
4  -- 使用新创建的数据库
5  USE online_store;
6  -- 创建orders表格
7  CREATE TABLE orders (
8      id BIGINT,
9      sku VARCHAR(15),
10     kinds VARCHAR(15),
11     stime VARCHAR(15),
12     num INT,
13     unit_price FLOAT,
14     total_price FLOAT,
15     longitude FLOAT,
16     latitude FLOAT
17 );
18 -- 将CSV文件中的数据导入到orders表中
19 LOAD DATA INFILE 'C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/Demo
20 _2.csv'
21 INTO TABLE orders
22 FIELDS TERMINATED BY ','
23 ENCLOSED BY '"'
24 LINES TERMINATED BY '
25 IGNORE 1 ROWS; -- 忽略第一行标题行

```

21 ✓

5 ✓

5 ✓

```

26 (2) -- 查询三同订单最多的商品种类kinds
27 SELECT kinds, COUNT(*) AS order_count
28 FROM (
29     SELECT kinds, stime, longitude, latitude, COUNT(*) AS same_order
        _count
30     FROM orders
31     GROUP BY kinds, stime, longitude, latitude
32     HAVING COUNT(*) >= 3
33 ) AS subquery
34 GROUP BY kinds
35 ORDER BY order_count DESC
36 LIMIT 1;
37

```

62|29 马宵 第二题

```

1  import sqlite3
2  import pandas as pd
3  conn = sqlite3.connect('online_store.db')
4  cursor = conn.cursor()
5  cursor.execute('''
6  CREATE TABLE IF NOT EXISTS orders (
7      id VARCHAR(20),
8      sku VARCHAR(15),
9      kinds VARCHAR(15),
10     stime VARCHAR(15),
11     num INT,
12     unit_price FLOAT,
13     total_price FLOAT,
14     longitude FLOAT,
15     latitude FLOAT
16 )
17 ''')
18 df = pd.read_csv(r"C:\Users\Stu\Desktop\Demo_2.csv")
19 df.to_sql('orders', conn, if_exists='replace', index=False)
20 cursor.execute('''
21 SELECT kinds, COUNT(*) as same_order_count
22 FROM orders
23 GROUP BY stime, longitude, latitude, sku, kinds
24 ORDER BY same_order_count DESC
25 LIMIT 1
26 ''')
27 result = cursor.fetchone()
28 print(f"三同订单最多的商品种类是: {result[0]}, 订单数量 {result[1]}")
29 conn.close()

```

21

```
1 import pandas as pd
2 import sqlite3
3 # 连接到SQLite数据库（如果不存在会创建新的），这里数据库文件名为online_store.
4 db
5 conn = sqlite3.connect('online_store.db')
6 # 创建orders表的SQL语句
7 create_table_sql = """
8 CREATE TABLE IF NOT EXISTS orders (
9     id varchar(20),
10     sku varchar(15),
11     kinds varchar(15),
12     stime varchar(15),
13     num int,
14     unit_price float,
15     total_price float,
16     longitude float,
17     latitude float
18 );
19 """
20 # 获取游标对象，用于执行SQL语句
21 cursor = conn.cursor()
22 # 执行创建表的SQL语句
23 cursor.execute(create_table_sql)
24 # 使用pandas读取Demo_2.csv文件数据
25 data = pd.read_csv('Demo_2.csv')
26 # 将数据插入到orders表中，如果表已存在会追加数据（注意这里的数据类型要和表结构尽量匹配，示例中按简单对应处理）
27 data.to_sql('orders', conn, if_exists='append', index=False)
28 # 提交事务（确保数据插入成功）
29 conn.commit()
30 # 关闭游标和连接
31 cursor.close()
32 conn.close()
33 import pandas as pd
34 import sqlite3
35 # 连接到数据库
36 conn = sqlite3.connect('online_store.db')
37 # 编写SQL查询语句，思路和之前纯SQL实现类似，这里通过子查询先找出不同的"三同"组合，再统计每个组合出现次数，最后找出次数最多的商品种类
38 query = """
39 SELECT kinds
40 FROM (
41     SELECT kinds, COUNT(*) AS count_orders
42     FROM (
43         SELECT DISTINCT stime, longitude, latitude, sku, kinds
44         FROM orders
45     ) AS distinct_orders
```

```

46         GROUP BY stime, longitude, latitude, sku
47         ORDER BY count_orders DESC
48         LIMIT 1
49     ) AS most_frequent_kinds;
50     """
51     # 使用pandas的read_sql函数执行查询并获取结果
52     result = pd.read_sql(query, conn)
53     # 输出结果（这里假设结果只有一行一列，即符合条件的那个商品种类）
54     print(result.iloc[0, 0])
55     # 关闭连接
    conn.close()

```

64|31 黎小源 第二题

```

1     # -*- coding: utf-8 -*-
2     """
3     Created on Fri Dec 20 19:41:23 2024
4     @author: Stu
5     """
6     import sqlite3
7     import pandas as pd
8     import os
9     csv_file_path = r'C:\Users\Stu\Desktop\Demo_2.csv'
10    db_file_path = r'C:\Users\Stu\Desktop\online_store.db'
11    def create_database_and_table():
12        conn = sqlite3.connect(db_file_path)
13        cursor = conn.cursor()
14        cursor.execute('''
15            CREATE TABLE IF NOT EXISTS orders (
16                id VARCHAR(20),
17                sku VARCHAR(15),
18                kinds VARCHAR(15),
19                stime VARCHAR(15),
20                num INTEGER,
21                unit_price REAL,
22                total_price REAL,
23                longitude REAL,
24                latitude REAL,
25                PRIMARY KEY (id)
26            )
27        ''')
28        conn.commit()
29        cursor.close()
30        conn.close()
31    def import_csv_to_db(csv_file, db_file):
32        df = pd.read_csv(csv_file)
33        conn = sqlite3.connect(db_file)
34        df.to_sql('orders', conn, if_exists='replace', index=False)

```

```

35     conn.close()
36 def find_top_same_orders_product():
37     conn = sqlite3.connect(db_file_path)
38     cursor = conn.cursor()
39     cursor.execute('''
40         WITH SameOrder AS (
41             SELECT kinds, stime, longitude, latitude, COUNT(*) as co
42 unt
43             FROM orders
44             GROUP BY kinds, stime, longitude, latitude
45             HAVING COUNT(*) > 1
46         )
47         SELECT kinds, COUNT(*) as total_same_orders
48         FROM SameOrder
49         GROUP BY kinds
50         ORDER BY total_same_orders DESC
51         LIMIT 1
52     ''')
53     result = cursor.fetchone()
54     cursor.close()
55     conn.close()
56     return result
57 def main():
58     create_database_and_table()
59     import_csv_to_db(csv_file_path, db_file_path)
60     result = find_top_same_orders_product()
61     if result:
62         kinds, total_same_orders = result
63         print(f"三同订单最多的商品种类是: {kinds}, 共有 {total_same_order
64 s} 个三同订单")
65     else:
66         print("没有找到符合条件的三同订单")
67 if __name__ == "__main__":
68     main()

```

65|32 龚卓能 第二题

```

1  import pymysql
2  import csv
3  # 建立数据库连接
4  def connect_mysql():
5      try:
6          conn = pymysql.connect(
7              host='localhost', # 这里需要根据你的实际MySQL主机地址修改, 通
常本地为localhost
8              user='root', # 改成你的MySQL用户名
9              password='Qwe_1234', # 修改为你的实际密码
10             database='', # 初始不指定具体数据库

```

```
11         charset='utf8'
12     )
13     return conn
14 except pymysql.Error as e:
15     print(f"连接MySQL数据库出现错误: {e}")
16     return None
17 # 创建online_store数据库
18 def create_database(conn):
19     cursor = conn.cursor()
20     try:
21         cursor.execute("CREATE DATABASE IF NOT EXISTS online_store")
22         print("online_store数据库创建成功（如果不存在的话）")
23     except pymysql.Error as e:
24         print(f"创建数据库时出错: {e}")
25     finally:
26         cursor.close()
27 # 使用online_store数据库
28 def use_database(conn):
29     cursor = conn.cursor()
30     try:
31         cursor.execute("USE online_store")
32     except pymysql.Error as e:
33         print(f"切换使用数据库时出错: {e}")
34     finally:
35         cursor.close()
36 # 创建orders表, 添加字段备注
37 def create_orders_table(conn):
38     cursor = conn.cursor()
39     create_table_sql = """
40     CREATE TABLE IF NOT EXISTS orders (
41         id varchar(20) COMMENT '订单号',
42         sku varchar(15) COMMENT '商品货号',
43         kinds varchar(15) COMMENT '种类',
44         stime varchar(15) COMMENT '销售时间',
45         num int COMMENT '数量',
46         unit_price float COMMENT '单价',
47         total_price float COMMENT '总价',
48         longitude float COMMENT '收货地经度',
49         latitude float COMMENT '纬度'
50     ) ENGINE=InnoDB DEFAULT CHARSET=utf8;
51     """
52     try:
53         cursor.execute(create_table_sql)
54         print("orders表创建成功（如果不存在的话）")
55     except pymysql.Error as e:
56         print(f"创建orders表时出错: {e}")
57     finally:
58         cursor.close()
59 # 从CSV文件读取数据并插入到orders表
```



```

60 def insert_data_from_csv(conn):
61     file_path = "D:/Demo_2.csv"
62     cursor = conn.cursor()
63     insert_sql = "INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude) VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)"
64     try:
65         with open(file_path, 'r', encoding='utf-8') as csvfile:
66             reader = csv.reader(csvfile)
67             next(reader) # 跳过文件头行（标题行）
68             for row in reader:
69                 cursor.execute(insert_sql, tuple(row))
70             conn.commit()
71             print("数据已成功插入到orders表中")
72     except (pymysql.Error, FileNotFoundError) as e:
73         print(f"插入数据出现错误: {e}")
74     finally:
75         cursor.close()
76 if __name__ == "__main__":
77     # 连接数据库
78     connection = connect_mysql()
79     if connection:
80         # 创建数据库
81         create_database(connection)
82         # 使用数据库
83         use_database(connection)
84         # 创建orders表
85         create_orders_table(connection)
86         # 插入数据
87         insert_data_from_csv(connection)
88         # 关闭连接
89         connection.close()
90 # 建立数据库连接
91 def connect_mysql():
92     try:
93         conn = pymysql.connect(
94             host='localhost',
95             user='root',
96             password='Qwe_1234',
97             database='online_store', # 这里使用前面创建好的online_store数据库
98             charset='utf8'
99         )
100         return conn
101     except pymysql.Error as e:
102         print(f"连接MySQL数据库出现错误: {e}")
103         return None
104 # 查询并统计三同订单中商品种类的出现次数
105 def count_triple_same_orders(conn):

```

```

106     cursor = conn.cursor()
107     query_sql = """
108     SELECT kinds, COUNT(*) as count
109     FROM (
110         SELECT stime, longitude, latitude, kinds
111         FROM orders
112         GROUP BY stime, longitude, latitude, kinds
113     ) as subquery
114     GROUP BY kinds
115     ORDER BY count DESC
116     LIMIT 1;
117     """
118     try:
119         cursor.execute(query_sql)
120         result = cursor.fetchone()
121         if result:
122             return result[0] # 返回三同订单最多的商品种类
123         return None
124     except pymysql.Error as e:
125         print(f"查询统计三同订单出现错误: {e}")
126     finally:
127         cursor.close()
128 if __name__ == "__main__":
129     connection = connect_mysql()
130     if connection:
131         most_kinds = count_triple_same_orders(connection)
132         if most_kinds:
133             print(f"三同订单最多的商品种类是: {most_kinds}")
134         else:
135             print("未找到三同订单相关数据")
136     connection.close()

```

66|33 龚新宇 第二题

```

1  # -*- coding: utf-8 -*-
2  """
3  Spyder Editor
4  This is a temporary script file.
5  """
6  import pymysql
7  import pandas as pd
8  connection = pymysql.connect(host='localhost', user='root', password
    = 'Qwe_1234', port=3306)
9  data_csv = pd.read_csv('Demo_2.csv', encoding="gbk")
10 try:
11     # 使用with语句管理第一个游标
12     with connection.cursor() as cursor:
13         # 创建数据库（如果不存在）

```

```

cursor.execute("CREATE DATABASE IF NOT EXISTS online_store")
# 使用数据库
cursor.execute("USE online_store")
# 创建表（如果不存在）
sql = """
CREATE TABLE IF NOT EXISTS orders (
    id VARCHAR(20),
    sku VARCHAR(15),
    kinds VARCHAR(15),
    stime VARCHAR(15),
    num INT,
    unit_price FLOAT,
    total_price FLOAT,
    longitude FLOAT,
    latitude FLOAT
)
"""
cursor.execute(sql)
# 插入数据
for row in data_csv.values.tolist():
    cursor.execute('
        INSERT INTO orders (id, sku, kinds, stime, num, unit_price, total_price, longitude, latitude)
        VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s)
        ', row)

# 提交事务
connection.commit()
query = "SELECT * FROM orders LIMIT 10"
cursor.execute(query)
results = cursor.fetchall()
for row in results:
    print(row)

# 重新打开一个游标来执行后续的查询
with connection.cursor() as cursor:
    query = '
    SELECT kinds, stime, longitude, latitude, COUNT(*) AS order_count
    FROM orders
    GROUP BY kinds, stime, longitude, latitude
    HAVING COUNT(*) > 2
    ORDER BY order_count DESC
    LIMIT 1
    '
    cursor.execute(query)
    result = cursor.fetchone()
    if result:
        print('三同订单最多的商品种类为:', result[0])
    else:

```

```
61         print('没有找到三同订单')
62     # 无论是否发生异常，都关闭数据库连接
63     finally:
        connection.close()
```

第三题

评分标准：满分 30 分，标准答案：配送人数：5, 3, 2

第 (1) 问人数为 4, 6 分数 - 1, 其他答案 - 2


第 (2) (3) 问人数不对 -2

没用聚类方法, 每个问 -3

第三问不能送达用户有 3 个, 每个2分

运行错误, 每个问 -4

67|1 何宇迪 第三题



```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Fri Dec 20 19:47:58 2024
4 @author: Stu
5 """
6 import csv
7 import math
8 print('第一题: ')
9 # 计算两点之间的距离
10 def distance(x1, y1, x2, y2):
11     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
12 # 读取Demo_3_1.csv文件中的数据
13 customers1 = []
14 with open('C://Users/Stu/Desktop/Demo_3_1.csv', 'r', encoding='gbk')
    as file:
15     reader = csv.reader(file)
16     next(reader) # 跳过标题行
17     for row in reader:
18         x, y, size = map(float, row)
19         customers1.append((x, y, size))
20 # 初始化员工列表和员工配送的客户列表
21 employees = []
22 employee_customers = []
23 # 分配客户给员工
24 for customer in customers1:
25     x, y, size = customer
26     assigned = False
27     for i, employee in enumerate(employees):
28         total_size = sum(c[2] for c in employee_customers[i])
29         if total_size + size <= 160:
30             employee_customers[i].append(customer)
31             assigned = True
32             break
33     if not assigned:
```

```

34         employees.append([customer])
35         employee_customers.append([customer])
36     # 输出结果
37     for i, employee in enumerate(employees):
38         print(f"员工{i + 1}负责的客户: ")
39         for customer in employee_customers[i]:
40             print(f"客户位置: ({customer[0]}, {customer[1]}), 需求: {customer[2]}")
41             print(f"员工{i + 1}配送的总需求: {sum(c[2] for c in employee_customers[i])}")
42     print(f"需要安排的员工数量: {len(employees)}")
43     print('第二题: ')
44     # 计算两点之间的距离
45     def distance(x1, y1, x2, y2):
46         return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
47     # 读取Demo_3_2.csv文件中的数据
48     customers2 = []
49     with open('C://Users/Stu/Desktop/Demo_3_2.csv', 'r', encoding='gbk')
50         as file:
51         reader = csv.reader(file)
52         next(reader) # 跳过标题行
53         for row in reader:
54             x, y, size = map(float, row)
55             customers2.append((x, y, size))
56     # 初始化员工列表和员工配送的客户列表
57     employees = []
58     employee_customers = []
59     # 分配客户给员工
60     while customers2:
61         current_employee_customers = []
62         current_employee_customers.append(customers2.pop(0))
63         for customer in customers2[:]:
64             for assigned_customer in current_employee_customers:
65                 if distance(customer[0], customer[1], assigned_customer
66                     [0], assigned_customer[1]) <= 25:
67                     current_employee_customers.append(customer)
68                     customers2.remove(customer)
69                     break
70             employees.append(current_employee_customers)
71             employee_customers.append(current_employee_customers)
72     # 输出结果
73     for i, employee in enumerate(employees):
74         print(f"员工{i + 1}负责的客户: ")
75         for customer in employee_customers[i]:
76             print(f"客户位置: ({customer[0]}, {customer[1]}), 需求: {customer[2]}")
77             print(f"员工{i + 1}配送的总需求: {sum(c[2] for c in employee_customers[i])}")
78     print(f"需要安排的员工数量: {len(employees)}")

```

```
77 print('第三题: ')
78 # 计算两点之间的距离
79 def distance(x1, y1, x2, y2):
80     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
81 # 读取Demo_3_3.csv文件中的数据
82 customers3 = []
83 with open('C://Users/Stu/Desktop/Demo_3_3.csv', 'r', encoding='gbk')
    as file:
84     reader = csv.reader(file)
85     next(reader) # 跳过标题行
86     for row in reader:
87         x, y, size = map(float, row)
88         customers3.append((x, y, size))
89 # 初始化员工列表和员工配送的客户列表
90 employees = [[] for _ in range(2)]
91 employee_customers = [[] for _ in range(2)]
92 # 分配客户给员工
93 while customers3:
94     assigned = False
95     for i in range(2):
96         current_employee_customers = employees[i]
97         for customer in customers3[:]:
98             if not current_employee_customers:
99                 current_employee_customers.append(customer)
100                 customers3.remove(customer)
101                 assigned = True
102                 break
103         for assigned_customer in current_employee_customers:
104             if distance(customer[0], customer[1], assigned_customer[0], assigned_customer[1]) <= 15:
105                 current_employee_customers.append(customer)
106                 customers3.remove(customer)
107                 assigned = True
108                 break
109         if assigned:
110             break
111     if not assigned:
112         break
113 # 输出结果
114 for i, employee in enumerate(employees):
115     print(f"员工{i + 1}负责的客户: ")
116     for customer in employee_customers[i]:
117         print(f"客户位置: ({customer[0]}, {customer[1]}), 需求: {customer[2]}")
118     print(f"员工{i + 1}配送的总需求: {sum(c[2] for c in employee_customers[i])}")
119 if customers3:
120     print("无法送达的客户: ")
121     for customer in customers3:
```

```
122     print(f"客户位置: ({customer[0]}, {customer[1]}), 需求: {customer[2]}")
```

输出:

- (1) 5
- (2) 21
- (3) 0 (所有客户不可达)

68|2 刘昱君 第三题

```
1  #第三题
2  import pandas as pd
3  import numpy as np
4  from scipy.spatial.distance import euclidean
5  from itertools import combinations
6  from sklearn.cluster import DBSCAN
7  # 加载数据
8  def load_customer_data(file_path):
9      data = pd.read_csv(file_path)
10     customers = data[['x', 'y', 'size']].to_numpy()
11     return customers
12 # 计算两点欧几里得距离
13 def calculate_distance(point1, point2):
14     return euclidean(point1[:2], point2[:2])
15 # 问题 1: 根据容量限制分组
16 def assign_employees_by_capacity(customers, max_capacity=160):
17     customers = sorted(customers, key=lambda x: x[2], reverse=True)
18 # 按需求量降序
19     groups = []
20     current_group = []
21     current_capacity = 0
22     for customer in customers:
23         if current_capacity + customer[2] <= max_capacity:
24             current_group.append(customer)
25             current_capacity += customer[2]
26         else:
27             groups.append(current_group)
28             current_group = [customer]
29             current_capacity = customer[2]
30     if current_group:
31         groups.append(current_group)
32     return groups
33 # 问题 2: 基于距离约束分配
34 def assign_employees_by_distance(customers, max_distance=25):
35     positions = customers[:, :2]
36     clustering = DBSCAN(eps=max_distance, min_samples=1).fit(positions)
```



```

37     labels = clustering.labels_
38     groups = {label: [] for label in set(labels)}
39     for idx, label in enumerate(labels):
40         groups[label].append(customers[idx])
41     return groups
42 # 问题 3: 只有两名员工, 距离限制 15 公里
43 def assign_two_employees_limited_distance(customers, max_distance=1
44 5):
45     positions = customers[:, :2]
46     clustering = DBSCAN(eps=max_distance, min_samples=1).fit(positio
47 ns)
48     labels = clustering.labels_
49     unique_labels = set(labels)
50     if len(unique_labels) <= 2:
51         groups = {label: [] for label in unique_labels}
52         for idx, label in enumerate(labels):
53             groups[label].append(customers[idx])
54         return groups, []
55     # 超出两名员工限制: 找出未分配的客户
56     assigned_groups = {label: [] for label in list(unique_labels)[:
57 2]}
58     unassigned_customers = []
59     for idx, label in enumerate(labels):
60         if label in assigned_groups:
61             assigned_groups[label].append(customers[idx])
62         else:
63             unassigned_customers.append(customers[idx])
64     return assigned_groups, unassigned_customers
65 if __name__ == "__main__":
66     # 问题 1
67     customers = load_customer_data("C:/Users/Stu/Desktop/Demo_3_1.cs
68 v")
69     groups = assign_employees_by_capacity(customers, max_capacity=16
70 0)
71     print(f"问题 1: 需要 {len(groups)} 名员工")
72     for i, group in enumerate(groups):
73         print(f"员工 {i+1} 负责的客户: {group}")
74     # 问题 2
75     customers = load_customer_data("C:/Users/Stu/Desktop/Demo_3_2.cs
76 v")
77     groups = assign_employees_by_distance(customers, max_distance=2
78 5)
79     print(f"问题 2: 需要 {len(groups)} 名员工")
80     for label, group in groups.items():
81         print(f"员工 {label+1} 负责的客户: {group}")
82     # 问题 3
83     customers = load_customer_data("C:/Users/Stu/Desktop/Demo_3_3.cs
84 v")
85     groups, unassigned = assign_two_employees_limited_distance(custo

```

```

mers, max_distance=15)
print(f"问题 3: 员工分配结果")
for label, group in groups.items():
    print(f"员工 {label+1} 负责的客户: {group}")
if unassigned:
    print(f"无法送达的客户: {unassigned}")

```

输出:

- (1) 5 ✓
- (2) 3 ✓
- (3) 2 ✓ (202个客户无法送达)

✓ b

69|3 吴天行 第三题

19

```

1  import csv
2  # 读取Demo_3_1.csv文件
3  with open('/mnt/Demo_3_1.csv', mode='r', encoding='utf-8') as file:
4      csv_reader = csv.reader(file)
5      next(csv_reader) # 跳过表头
6      customers = [(float(row[0]), float(row[1]), int(row[2])) for row
in csv_reader]
7  # 按照需求从大到小排序
8  customers.sort(key=lambda x: x[2], reverse=True)
9  # 初始化员工列表和当前员工的需求
10 employees = []
11 current_employee = []
12 current_demand = 0
13 # 分配客户给员工
14 for customer in customers:
15     if current_demand + customer[2] <= 160:
16         current_employee.append(customer)
17         current_demand += customer[2]
18     else:
19         employees.append(current_employee)
20         current_employee = [customer]
21         current_demand = customer[2]
22 # 将最后一个员工添加到员工列表中
23 employees.append(current_employee)
24 # 打印结果
25 for i, employee in enumerate(employees, 1):
26     print(f"员工{i}负责的客户: {employee}")
27 print(f"需要安排的员工数量: {len(employees)}")
28 import csv
29 import math
30 def distance(x1, y1, x2, y2):
31     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
32 # 读取Demo_3_2.csv文件

```

```
33 with open('/mnt/Demo_3_2.csv', mode='r', encoding='utf-8') as file:
34     csv_reader = csv.reader(file)
35     next(csv_reader) # 跳过表头
36     customers = [(float(row[0]), float(row[1]), int(row[2])) for row
in csv_reader]
37 # 初始化员工列表和当前员工的客户列表
38 employees = []
39 current_employee = [customers[0]]
40 # 分配客户给员工
41 for i in range(1, len(customers)):
42     customer = customers[i]
43     assignable = True
44     for assigned_customer in current_employee:
45         if distance(customer[0], customer[1], assigned_customer[0],
assigned_customer[1]) > 25:
46             assignable = False
47             break
48     if assignable:
49         current_employee.append(customer)
50     else:
51         employees.append(current_employee)
52         current_employee = [customer]
53 # 将最后一个员工添加到员工列表中
54 employees.append(current_employee)
55 # 打印结果
56 for i, employee in enumerate(employees, 1):
57     print(f"员工{i}负责的客户: {employee}")
58 print(f"需要安排的员工数量: {len(employees)}")
59 import csv
60 # 读取Demo_3_3.csv文件
61 with open('/mnt/Demo_3_3.csv', mode='r', encoding='utf-8') as file:
62     csv_reader = csv.reader(file)
63     next(csv_reader) # 跳过表头
64     customers = [(float(row[0]), float(row[1]), int(row[2])) for row
in csv_reader]
65 # 按照需求从大到小排序
66 customers.sort(key=lambda x: x[2], reverse=True)
67 # 初始化员工列表和当前员工的需求
68 employees = []
69 current_employee = []
70 current_demand = 0
71 # 分配客户给员工
72 for customer in customers:
73     if current_demand + customer[2] <= 160:
74         current_employee.append(customer)
75         current_demand += customer[2]
76     else:
77         employees.append(current_employee)
78         current_employee = [customer]
```

```

79         current_demand = customer[2]
80     # 将最后一个员工添加到员工列表中
81     employees.append(current_employee)
82     # 打印结果
83     for i, employee in enumerate(employees, 1):
84         print(f"员工{i}负责的客户: {employee}")
85     print(f"需要安排的员工数量: {len(employees)}")

```

输出:

(1) 5 ✓
 (2) 395 ✓
 (3) 4 (没有输出)

3
 6

70|4 和丽兴 第三题

```

1  import pandas as pd
2  from sklearn.cluster import KMeans
3  import numpy as np
4  # 读取数据
5  file_path = r'C:/Users/Stu/Desktop/Demo_3_1.csv'
6  data = pd.read_csv(file_path)
7  # 提取客户位置和需求
8  positions = data[['x', 'y']].values
9  sizes = data['size'].values
10 # 计算每个客户的需求, 若超过160瓶牛奶, 分配更多员工
11 num_employees = 0
12 assignments = []
13 # 最大牛奶配送量
14 max_milk = 160
15 while np.sum(sizes) > 0:
16     # 聚类分组
17     kmeans = KMeans(n_clusters=10) # 你可以试着调整聚类数
18     clusters = kmeans.fit_predict(positions)
19     # 计算每个集群的总需求
20     cluster_sizes = [np.sum(sizes[clusters == i]) for i in range(kmeans.n_clusters)]
21     # 分配员工
22     for i in range(kmeans.n_clusters):
23         if cluster_sizes[i] <= max_milk:
24             assignments.append((i, list(np.where(clusters == i)
25 [0])))
26             num_employees += 1
27             sizes[clusters == i] = 0 # 清除已分配客户
28 # 输出结果
29 print(f"需要安排员工数量: {num_employees}")
30 for employee, clients in assignments:
31     print(f"员工{employee + 1}负责客户: {clients}")

```

```

32 import pandas as pd
33 from sklearn.cluster import DBSCAN
34 import numpy as np
35 # 读取数据
36 file_path = r'C:/Users/Stu/Desktop/Demo_3_2.csv'
37 data = pd.read_csv(file_path)
38 # 提取客户的位置和需求
39 positions = data[['x', 'y']].values
40 sizes = data['size'].values
41 # 使用DBSCAN进行聚类
42 db = DBSCAN(eps=25, min_samples=1, metric='euclidean') # Euclidean
43 distance is used here
44 clusters = db.fit_predict(positions)
45 # 输出每个聚类的客户列表
46 unique_clusters = set(clusters)
47 if -1 in unique_clusters:
48     unique_clusters.remove(-1) # 排除噪声点
49 num_employees = len(unique_clusters) # 计算实际需要安排的员工数量
50 assignments = []
51 for i in unique_clusters:
52     customer_indices = np.where(clusters == i)[0]
53     assignments.append(customer_indices)
54 # 输出结果
55 print(f"需要安排员工数量: {num_employees}")
56 for employee, clients in enumerate(assignments):
57     print(f"员工{employee + 1}负责客户: {clients}")
58 import pandas as pd
59 from sklearn.cluster import DBSCAN
60 import numpy as np
61 # 读取数据
62 file_path = r'C:\Users\Stu\Desktop\Demo_3_3.csv'
63 data = pd.read_csv(file_path)
64 # 提取客户的位置和需求
65 positions = data[['x', 'y']].values
66 sizes = data['size'].values
67 # 计算客户间的欧几里得距离
68 def calculate_distance(p1, p2):
69     return np.sqrt((p2[0] - p1[0]) ** 2 + (p2[1] - p1[1]) ** 2)
70 # 使用DBSCAN进行聚类
71 db = DBSCAN(eps=15, min_samples=1, metric='euclidean')
72 clusters = db.fit_predict(positions)
73 # 输出每个聚类的客户列表
74 num_employees = len(set(clusters)) - (1 if -1 in clusters else 0) #
75 排除噪声点
76 assignments = []
77 for i in range(num_employees):
78     customer_indices = np.where(clusters == i)[0]
79     assignments.append(customer_indices)
80 # 如果客户数超过两个员工，返回无法送达客户

```

```

79     if num_employees > 2:
80         print("无法配送的客户: ")
81         # 找到超过两个员工的组并标记为无法配送
82         all_customers = set(range(len(data)))
83         assigned_customers = set(np.concatenate(assignments))
84         unreachable_customers = all_customers - assigned_customers
85         print(unreachable_customers)
86     # 输出每个员工负责的客户
87     print(f"需要安排员工数量: 2")
88     for employee, clients in enumerate(assignments[:2]):
89         print(f"员工{employee + 1}负责客户: {clients}")

```

输出:

- (1) 55 -2
 (2) 3 -6
 (3) 2 (无)

71|5 塔巴江村 第三题

72|6 崔杰 第三题

```

1     #第一题
2     import pandas as pd
3     from sklearn.cluster import KMeans
4     import numpy as np
5     # 读取Demo_3_1.csv文件数据
6     data_1 = pd.read_csv('Demo_3_1.csv')
7     # 获取客户需求列表 (size列) 并转换为numpy数组
8     customer_size = data_1['size'].to_numpy()
9     # 计算总需求
10    total_demand = np.sum(customer_size)
11    # 每个员工最多配送160瓶牛奶, 计算需要的员工数量
12    num_employees = total_demand // 160
13    if total_demand % 160 != 0:
14        num_employees += 1
15    # 只考虑位置信息进行聚类 (这里简单以K-Means为例, K值为员工数量)
16    kmeans = KMeans(n_clusters=num_employees)
17    locations = data_1[['x', 'y']].to_numpy()
18    kmeans.fit(locations)
19    # 输出每个员工负责的客户
20    for i in range(num_employees):
21        employee_customers = data_1[kmeans.labels_ == i]
22        print(f"员工{i + 1}负责的客户: \n{employee_customers}")
23    print(f"共需要安排{num_employees}个员工。")
24    #(2)第二题
25    import pandas as pd
26    import numpy as np

```

```

27 # 读取文件
28 df = pd.read_csv('Demo_3_2.csv')
29 # 计算两点之间的距离
30 def distance(x1, y1, x2, y2):
31     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
32 # 贪心算法分组
33 def group_customers(customers, max_distance):
34     groups = []
35     while customers:
36         group = [customers[0]]
37         customers = customers[1:]
38         i = 0
39         while i < len(customers):
40             if all(distance(group[-1]['x'], group[-1]['y'], customer
s[i]['x'], customers[i]['y']) <= max_distance for
41                     customer in group):
42                 group.append(customers.pop(i))
43             else:
44                 i += 1
45             groups.append(group)
46     return groups
47 # 转换数据格式
48 customers = df.to_dict('records')
49 # 分组客户
50 groups = group_customers(customers, 25)
51 # 输出结果
52 for i, group in enumerate(groups):
53     print(f'员工{i + 1}负责的客户索引: {[index for index, customer in e
numerate(group)]}')
54 print(f'需要安排的员工数量: {len(groups)}')
55 #第三题
56 import pandas as pd
57 from sklearn.cluster import KMeans
58 import numpy as np
59 import math
60 # 计算两点间距离的函数
61 def distance(x1, y1, x2, y2):
62     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
63 # 读取Demo_3_3.csv文件
64 data_3 = pd.read_csv('Demo_3_3.csv')
65 # 构建坐标矩阵
66 coordinates = data_3[['x', 'y']].values
67 # 初始化两个员工负责客户的列表和未送达客户列表
68 employee1_customers = []
69 employee2_customers = []
70 not_delivered_customers = []
71 # 先对所有客户进行一次聚类尝试（分为两组）
72 kmeans = KMeans(n_clusters=2)
73 kmeans.fit(coordinates)

```

```

74 labels = kmeans.labels_
75 # 检查每个聚类中的客户距离是否满足要求，分配给相应员工，不符合的放入未送达列表
76 for index, label in enumerate(labels):
77     if label == 0:
78         is_valid_employee1 = True
79         if len(employee1_customers) > 0:
80             for cust in employee1_customers:
81                 dis = distance(data_3.loc[index]['x'], data_3.loc[index]['y'],
82                                data_3.loc[cust]['x'], data_3.loc[cust]['y'])
83                 if dis > 15:
84                     is_valid_employee1 = False
85                     break
86             if is_valid_employee1:
87                 employee1_customers.append(index)
88             else:
89                 not_delivered_customers.append(index)
90     else:
91         is_valid_employee2 = True
92         if len(employee2_customers) > 0:
93             for cust in employee2_customers:
94                 dis = distance(data_3.loc[index]['x'], data_3.loc[index]['y'],
95                                data_3.loc[cust]['x'], data_3.loc[cust]['y'])
96                 if dis > 15:
97                     is_valid_employee2 = False
98                     break
99             if is_valid_employee2:
100                 employee2_customers.append(index)
101             else:
102                 not_delivered_customers.append(index)
103 print("员工1负责的客户索引:", employee1_customers)
104 print("员工2负责的客户索引:", employee2_customers)
105 print("未送达的客户索引:", not_delivered_customers)
106

```

输出:

- (1) 5
- (2) 43
- (3) 2 (未送达有383个)

73|7 庄嘉帆 第三题

```

1 # -*- coding: utf-8 -*-
2 """

```



```
3 Created on Fri Dec 20 20:31:15 2024
4 @author: Stu
5 """
6 import pandas as pd
7 import math
8 # 读取Demo_3_1.csv文件
9 data_3_1 = pd.read_csv(r"C:\Users\Stu\Documents\WeChat Files\wxid_do
871zskr1zt12\FileStorage\File\2024-12\Demo_3_1.csv")
10 # 读取Demo_3_2.csv文件
11 data_3_2 = pd.read_csv(r"C:\Users\Stu\Documents\WeChat Files\wxid_do
871zskr1zt12\FileStorage\File\2024-12\Demo_3_2.csv")
12 # 读取Demo_3_3.csv文件
13 data_3_3 = pd.read_csv(r"C:\Users\Stu\Documents\WeChat Files\wxid_do
871zskr1zt12\FileStorage\File\2024-12\Demo_3_3.csv")
14 def distance(point1, point2):
15     return math.sqrt((point1[0] - point2[0])** 2 + (point1[1] - poi
nt2[1])** 2)
16 def assign_customers_3_1(data):
17     total_demand = data['size'].sum()
18     num_employees = (total_demand + 159) // 160 # 向上取整计算员工数量
19     employees = [[] for _ in range(num_employees)]
20     current_employee = 0
21     current_demand = 0
22     for index, row in data.iterrows():
23         if current_demand + row['size'] <= 160:
24             employees[current_employee].append(index)
25             current_demand += row['size']
26         else:
27             current_employee += 1
28             current_demand = row['size']
29             employees[current_employee].append(index)
30     return employees
31 def assign_customers_3_2(data):
32     employees = [[] for _ in range(3)] # 初始化3个员工负责的客户列表
33     cluster_centers = [] # 存储聚类中心
34     # 随机选择3个初始聚类中心（这里简单地选择前3个客户作为初始中心，实际应用中可
以使用更复杂的初始化方法）
35     for i in range(3):
36         cluster_centers.append(data.loc[i][['x', 'y']].values)
37     # 分配客户到最近的聚类中心
38     for index, row in data.iterrows():
39         customer_point = row[['x', 'y']].values
40         distances = [distance(customer_point, center) for center in
cluster_centers]
41         nearest_cluster = distances.index(min(distances))
42         employees[nearest_cluster].append(index)
43         # 更新聚类中心（简单地取当前聚类中所有客户位置的平均值）
44         cluster_centers[nearest_cluster] = (
45             sum([data.loc[customer]['x'] for customer in employees[n
```

```

46     earest_cluster])) / len(employees[nearest_cluster]),
        sum([data.loc[customer]['y'] for customer in employees[n
47 earest_cluster])) / len(employees[nearest_cluster])
48     )
49     return employees
50 # 处理Demo_3_1.csv
51 employees_3_1 = assign_customers_3_1(data_3_1)
52 print("Demo_3_1.csv:")
53 for i, employee in enumerate(employees_3_1):
54     print(f"员工 {i + 1} 负责的客户索引: {employee}")
55 # 处理Demo_3_2.csv
56 employees_3_2 = assign_customers_3_2(data_3_2)
57 print("\nDemo_3_2.csv:")
58 for i, employee in enumerate(employees_3_2):
59     print(f"员工 {i + 1} 负责的客户索引: {employee}")
60 # 计算客户间距离矩阵
61 def calculate_distance_matrix(data):
62     num_customers = len(data)
63     distance_matrix = [[0] * num_customers for _ in range(num_custom
64 ers)]
65     for i in range(num_customers):
66         for j in range(i + 1, len(data)):
67             point1 = (data.loc[i]['x'], data.loc[i]['y'])
68             point2 = (data.loc[j]['x'], data.loc[j]['y'])
69             distance_matrix[i][j] = distance_matrix[j][i] = math.sqr
70 t((point1[0] - point2[0]) ** 2 + (point1[1] - point2[1]) ** 2)
71     return distance_matrix
72 def assign_customers_3_3(data):
73     distance_matrix = calculate_distance_matrix(data)
74     employees = [[], []] # 两个员工负责的客户列表
75     unassigned_customers = set(range(len(data))) # 初始所有客户都未分
76     配
77     # 初始化分配, 先将距离最近的两个客户分别分配给两个员工
78     min_distance = float('inf')
79     for i in range(len(data)):
80         for j in range(i + 1, len(data)):
81             if distance_matrix[i][j] < min_distance:
82                 min_distance = distance_matrix[i][j]
83                 closest_customers = (i, j)
84     employees[0].append(closest_customers[0])
85     employees[1].append(closest_customers[1])
86     unassigned_customers.remove(closest_customers[0])
87     unassigned_customers.remove(closest_customers[1])
88     while unassigned_customers:
89         best_assignment = None
90         best_distance = float('inf')
91         for customer in unassigned_customers:
92             for i in range(2):
93                 for assigned_customer in employees[i]:

```

```

91         distance_to_assigned = distance_matrix[customer]
           [assigned_customer]
92         if distance_to_assigned <= 15 and distance_to_as
93 signed < best_distance:
94             best_assignment = (i, customer)
95             best_distance = distance_to_assigned
96         if best_assignment is None:
97             break
98         employees[best_assignment[0]].append(best_assignment[1])
99         unassigned_customers.remove(best_assignment[1])
100     return employees, list(unassigned_customers)
101 # 处理Demo_3_3.csv
102 employees_3_3, unassigned_3_3 = assign_customers_3_3(data_3_3)
103 print("Demo_3_3.csv:")
104 for i, employee in enumerate(employees_3_3):
    print(f"员工 {i + 1} 负责的客户索引: {employee}")
print(f"未送达的客户索引: {unassigned_3_3}")

```

输出:

- (1) 5
- (2) 3
- (3) 2 (未送达有203个)

- 3
- 6

74|8 张浩祖 第三题

```

1  # (1)
2  import pandas as pd
3  # 读取文件
4  df = pd.read_csv('Demo_3_1.csv')
5  # 计算每个客户的需求总和
6  total_demand = df['size'].sum()
7  # 计算需要的员工数量
8  num_employees = (total_demand + 159) // 160
9  # 初始化员工列表和当前员工的需求总和
10 employees = [[] for _ in range(num_employees)]
11 employee_demands = [0] * num_employees
12 # 遍历每个客户，将其分配给需求总和最小的员工
13 for index, row in df.iterrows():
14     min_demand_employee = employee_demands.index(min(employee_demand
15 s))
16     if employee_demands[min_demand_employee] + row['size'] <= 160:
17         employees[min_demand_employee].append(index)
18         employee_demands[min_demand_employee] += row['size']
19     else:
20         new_employee = employee_demands.index(min(employee_demands))
21         employees[new_employee].append(index)
22         employee_demands[new_employee] += row['size']

```

19

```

23 # 输出结果
24 for i, employee in enumerate(employees):
25     print(f'员工{i+1}负责的客户索引: {employee}')
26     #(2)
27 import pandas as pd
28 import numpy as np
29 # 读取文件
30 df = pd.read_csv('Demo_3_2.csv')
31 # 计算两点之间的距离
32 def distance(x1, y1, x2, y2):
33     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
34 # 贪心算法分组
35 def group_customers(customers, max_distance):
36     groups = []
37     while customers:
38         group = [customers[0]]
39         customers = customers[1:]
40         i = 0
41         while i < len(customers):
42             if all(distance(group[-1]['x'], group[-1]['y'], customer
s[i]['x'], customers[i]['y']) <= max_distance for
43                 customer in group):
44                 group.append(customers.pop(i))
45             else:
46                 i += 1
47             groups.append(group)
48     return groups
49 # 转换数据格式
50 customers = df.to_dict('records')
51 # 分组客户
52 groups = group_customers(customers, 25)
53 # 输出结果
54 for i, group in enumerate(groups):
55     print(f'员工{i + 1}负责的客户索引: {[index for index, customer in e
numerate(group)]}')
56 print(f'需要安排的员工数量: {len(groups)}')
57 #(3)
58 import pandas as pd
59 import numpy as np
60 # 读取文件
61 df = pd.read_csv('Demo_3_3.csv')
62 # 计算两点之间的距离
63 def distance(x1, y1, x2, y2):
64     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
65 # 贪心算法分配客户
66 def assign_customers(customers, max_distance):
67     employee1 = []
68     employee2 = []
69     remaining_customers = customers.copy()

```

```

70     while remaining_customers:
71         min_distance = float('inf')
72         closest_customer = None
73         for index, customer in enumerate(remaining_customers):
74             for assigned_customer in employee1 + employee2:
75                 dist = distance(customer['x'], customer['y'], assigned_customer['x'], assigned_customer['y'])
76                 if dist < min_distance:
77                     min_distance = dist
78                     closest_customer = index
79             if min_distance <= max_distance:
80                 if len(employee1) <= len(employee2):
81                     employee1.append(remaining_customers.pop(closest_customer))
82                 else:
83                     employee2.append(remaining_customers.pop(closest_customer))
84             else:
85                 break
86         return employee1, employee2, remaining_customers
87 # 转换数据格式
88 customers = df.to_dict('records')
89 # 分配客户
90 employee1, employee2, not_assigned = assign_customers(customers, 15)
91 # 输出结果
92 print('员工1负责的客户索引: ', [index for index, customer in enumerate(employee1)])
93 print('员工2负责的客户索引: ', [index for index, customer in enumerate(employee2)])
94 print('无法送达的客户索引: ', [index for index, customer in enumerate(not_assigned)])

```

输出:

(1) 5

(2) 43

(3) 2 (403)

75|9 张露丹 第三题

```

1  import pandas as pd
2  from sklearn.cluster import KMeans
3  import math
4  import numpy as np
5  # 读取数据
6  data = pd.read_csv("C:/Users/Stu/Desktop/12.20/Demo_3_1.csv")
7  # 提取客户位置和需求
8  X = data[['x', 'y']].values

```

```

9 sizes = data['size'].values
10 # 计算需要的员工数量
11 total_milk = sum(sizes)
12 num_workers = math.ceil(total_milk / 160)
13 # 使用K-Means聚类
14 kmeans = KMeans(n_clusters=num_workers, random_state=0).fit(X)
15 # 将客户分配给员工
16 data['worker'] = kmeans.labels_
17 data['milk'] = sizes
18 # 打印每个员工负责的客户
19 for i in range(num_workers):
20     print(f"Worker {i+1}:")
21     print(data[data['worker'] == i])
22 print(f"Total workers needed: {num_workers}")
23 from sklearn.cluster import DBSCAN
24 # 读取数据
25 data = pd.read_csv("C:/Users/Stu/Desktop/12.20/Demo_3_2.csv")
26 # 提取客户位置
27 X = data[['x', 'y']].values
28 # 使用DBSCAN聚类, eps设置为25km
29 dbscan = DBSCAN(eps=25, min_samples=1).fit(X)
30 # 获取聚类结果
31 labels = dbscan.labels_
32 # 计算每个聚类的客户数量和需求
33 unique_labels = set(labels)
34 workers = {i: [] for i in unique_labels if i != -1} # -1表示噪声点
35 for label in unique_labels:
36     if label != -1:
37         worker客户需求 = data[labels == label]
38         workers[label] = worker客户需求
39 # 打印每个员工负责的客户
40 for worker, customers in workers.items():
41     print(f"Worker {worker}:")
42     print(customers)
43 # 计算不能送达的用户
44 unreachable_customers = data[labels == -1]
45 print("Customers that cannot be reached:")
46 print(unreachable_customers)
47 # 计算需要的员工数量
48 num_workers = len(workers)
49 print(f"Total workers needed: {num_workers}")
50 # 读取数据
51 data = pd.read_csv("C:/Users/Stu/Desktop/12.20/Demo_3_3.csv")
52 # 提取客户位置
53 X = data[['x', 'y']].values
54 # 使用DBSCAN聚类, eps设置为15km, min_samples设置为1
55 dbscan = DBSCAN(eps=15, min_samples=1).fit(X)
56 # 获取聚类结果
57 labels = dbscan.labels_

```

```

58 # 获取每个聚类的客户
59 clusters = {}
60 for i in range(len(labels)):
61     cluster_id = labels[i]
62     if cluster_id not in clusters:
63         clusters[cluster_id] = []
64     clusters[cluster_id].append(i)
65 # 由于只有两个员工，我们需要手动分配两个最大的聚类给两个员工
66 sorted_clusters = sorted(clusters.items(), key=lambda x: len(x[1]),
67                           reverse=True)
68 if len(sorted_clusters) > 2: # 如果有超过两个聚类
69     print("Customers that cannot be reached:")
70     unreachable_customers = []
71     for cluster_id, indices in sorted_clusters[2:]: # 取第三大及以后的
72         unreachable_customers.extend(indices)
73     print(data.iloc[unreachable_customers])
74 # 分配前两个最大的聚类给两个员工
75 worker1 = sorted_clusters[0]
76 worker2 = sorted_clusters[1]
77 print("Worker 1 is responsible for customers:")
78 print(data.iloc[worker1[1]])
79 print("Worker 2 is responsible for customers:")
80 print(data.iloc[worker2[1]])
81 else:
82     # 如果只有两个或更少的聚类，直接分配给两个员工
83     worker1 = sorted_clusters[0]
84     if len(sorted_clusters) == 1:
85         worker2 = None # 没有第二个聚类给第二个员工
86     else:
87         worker2 = sorted_clusters[1]
88     print("Worker 1 is responsible for customers:")
89     print(data.iloc[worker1[1]])
90     if worker2:
91         print("Worker 2 is responsible for customers:")
92         print(data.iloc[worker2[1]])
93 # 打印每个员工负责的客户
94 print("Summary:")
95 if worker2:
96     print(f"Worker 1 serves {len(worker1[1])} customers.")
97     print(f"Worker 2 serves {len(worker2[1])} customers.")
98 else:
99     print("Worker 2 has no customers to serve.")
100     print(f"{len(worker1[1])} customers cannot be reached.")

```

输出:

(1) 5

(2) 3

(3) 2 (无未送达)

76|10 施習 第三题

77|11 李上卫 第三题

78|12 李俊杰 第三题

30

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Fri Dec 20 20:11:55 2024
4 @author: Stu
5 """
6 import pandas as pd
7 from sklearn.cluster import KMeans
8 import numpy as np
9 # 读取Demo_3_1.csv文件
10 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_3_1.csv')
11 # 每个员工每天最多配送量
12 max_per_employee = 160
13 # 计算所有客户的总需求
14 total_demand = df['size'].sum()
15 # 计算所需员工数量（向上取整）
16 employee_count = int(total_demand / max_per_employee) + (1 if total_
    demand % max_per_employee != 0 else 0)
17 print(f"需要安排 {employee_count} 个员工")
18 # 提取客户位置坐标数据（用于聚类）
19 customer_locations = df[['x', 'y']].values
20 # 使用KMeans算法进行聚类，聚类数量为员工数量
21 kmeans = KMeans(n_clusters=employee_count, random_state=0).fit(custo
    mer_locations)
22 # 获取每个客户所属的聚类标签（即员工编号，这里假设员工编号从1开始对应聚类标签0到
    employee_count - 1）
23 labels = kmeans.labels_
24 employee_customers = {}
25 for emp in range(employee_count):
26     employee_customers[emp + 1] = [i for i in range(len(labels)) if
        labels[i] == emp]
27 for emp in range(1, employee_count + 1):
28     print(f"员工 {emp} 负责的客户索引为: {employee_customers[emp]}")
29 import pandas as pd
30 import numpy as np
31 from sklearn.cluster import DBSCAN
32 # 读取Demo_3_2.csv文件
33 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_3_2.csv')
34 # 获取客户位置坐标
```



```

35 customer_locations = df[['x', 'y']].values
36 # 使用DBSCAN进行聚类, eps参数设置为最大单段行驶距离25km, min_samples设为1
   (每个客户都可作为一个聚类核心)
37 dbscan = DBSCAN(eps=25, min_samples=1).fit(customer_locations)
38 # 获取聚类标签, -1表示噪声点(即无法归入任何满足距离要求聚类的客户)
39 labels = dbscan.labels_
40 # 统计不同聚类的数量(不包含噪声点对应的 -1 标签), 即为员工数量
41 unique_labels = np.unique(labels)
42 employee_count = len(unique_labels) - (1 if -1 in unique_labels else
43 0)
44 print(f"需要安排 {employee_count} 个员工")
45 # 按照聚类标签将客户分配给各个员工
46 employee_customers = {}
47 for label in unique_labels:
48     if label == -1:
49         continue
50     employee_customers[label] = [i for i in range(len(labels)) if la
51 bels[i] == label]
52 for emp in range(employee_count):
53     print(f"员工 {emp} 负责的客户索引为: {employee_customers[emp]}")
54 import pandas as pd
55 import numpy as np
56 from sklearn.cluster import DBSCAN
57 # 读取Demo_3_3.csv文件
58 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_3_3.csv')
59 # 获取客户位置坐标
60 customer_locations = df[['x', 'y']].values
61 dbscan = DBSCAN(eps=15, min_samples=1).fit(customer_locations)
62 labels = dbscan.labels_
63 unique_labels = np.unique(labels)
64 unique_labels = unique_labels[unique_labels != -1]
65 num_clusters = len(unique_labels)
66 employee_customers = {1: [], 2: []}
67 unassigned_customers = []
68 if num_clusters <= 2:
69     for index, label in enumerate(labels):
70         if label == -1:
71             unassigned_customers.append(index)
72         elif label == 0:
73             employee_customers[1].append(index)
74         elif label == 1:
75             employee_customers[2].append(index)
76 while len(unassigned_customers) < 3:
77     if len(employee_customers[1]) > 0:
78         removed_index = employee_customers[1].pop(0)
79         unassigned_customers.append(removed_index)
80     elif len(employee_customers[2]) > 0:
81         removed_index = employee_customers[2].pop(0)
82         unassigned_customers.append(removed_index)

```

```

82 else:
83     cluster_sizes = np.bincount(labels[labels!= -1])
84     sorted_cluster_indices = np.argsort(cluster_sizes[::-1])
85     employee_1_size = 0
86     employee_2_size = 0
87     has_three_unassigned = False
88     for cluster_index in sorted_cluster_indices:
89         cluster_label = unique_labels[cluster_index]
90         cluster_customers = [i for i in range(len(labels)) if labels
91 [i] == cluster_label]
92         cluster_size = len(cluster_customers)
93         if not has_three_unassigned and len(unassigned_customers) <
94 3:
95             unassigned_customers.extend(cluster_customers)
96             if len(unassigned_customers) == 3:
97                 has_three_unassigned = True
98                 continue
99             if employee_1_size < employee_2_size:
100                 employee_customers[1].extend(cluster_customers)
101                 employee_1_size += cluster_size
102             else:
103                 employee_customers[2].extend(cluster_customers)
104                 employee_2_size += cluster_size
105 print(f"员工1负责的客户索引为: {unassigned_customers}")
106 print(f"员工2负责的客户索引为: {employee_customers[2]}")
107 print(f"无法送达的客户索引为: {employee_customers[1]}")

```

输出:

(1) 5
(2) 3
(3) 2 (300, 203, 193)

79|13 杨晨晨 第三题

19

```

1 (1)
2 import pandas as pd
3 import math
4 # 计算两点之间的距离
5 def distance(x1, y1, x2, y2):
6     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
7 def assign_customers_demo3_1():
8     data = pd.read_csv('Demo_3_1.csv')
9     total_demand = data['size'].sum()
10    employees_needed = (total_demand + 159) // 160 # 向上取整
11    print(f"需要安排 {employees_needed} 个员工")
12    # 简单平均分配客户给员工 (这里只是一种简单示例, 实际可能有更优分配策略)

```

```

13     customers_per_employee = len(data) // employees_needed
14     for i in range(employees_needed):
15         start_index = i * customers_per_employee
16         end_index = (i + 1) * customers_per_employee if i < employee
s_needed - 1 else len(data)
17         print(f"员工 {i + 1} 负责的客户: ")
18         print(data.iloc[start_index:end_index])
19     # 执行第一问
20 assign_customers_demo3_1()
21 (2)
22 import pandas as pd
23 import numpy as np
24 # 读取文件
25 df = pd.read_csv('Demo_3_2.csv')
26 # 计算两点之间的距离
27 def distance(x1, y1, x2, y2):
28     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
29 def group_customers(customers, max_distance):
30     groups = []
31     while customers:
32         group = [customers[0]]
33         customers = customers[1:]
34         i = 0
35         while i < len(customers):
36             if all(distance(group[-1]['x'], group[-1]['y'], customers[i]
['x'], customers[i]['y']) <= max_distance for
37                 customer in group):
38                 group.append(customers.pop(i))
39             else:
40                 i += 1
41         groups.append(group)
42     return groups
43 # 转换数据格式
44 customers = df.to_dict('records')
45 # 分组客户
46 groups = group_customers(customers, 25)
47 # 输出结果
48 for i, group in enumerate(groups):
49     print(f'员工{i + 1}负责的客户索引: {[index for index, customer in enum
erate(group)]}')
50 print(f'需要安排的员工数量: {len(groups)}')
51 (3)
52 def assign_customers_demo3_3():
53     data = pd.read_csv('Demo_3_3.csv')
54     employee1_customers = []
55     employee2_customers = []
56     unassigned_customers = []
57     current_employee1_demand = 0
58     current_employee2_demand = 0

```

```

59     for index, row in data.iterrows():
60         if (current_employee1_demand + row['size'] <= 160 and
61             (not employee1_customers or distance(employee1_customers[-1]['x'],
62                                                     employee1_customers[-1]['y'],
63                                                     row['x'], row
64                                                     ['y']) <= 15)):
65             employee1_customers.append(row)
66             current_employee1_demand += row['size']
67         elif (current_employee2_demand + row['size'] <= 160 and
68             (not employee2_customers or distance(employee2_customers[-1]['x'],
69                                                     employee2_customers[-1]['y'],
70                                                     row['x'], row
71                                                     ['y']) <= 15)):
72             employee2_customers.append(row)
73             current_employee2_demand += row['size']
74         else:
75             unassigned_customers.append(row)
76     print("员工 1 负责的客户: ")
77     print(pd.DataFrame(employee1_customers))
78     print("员工 2 负责的客户: ")
79     print(pd.DataFrame(employee2_customers))
80     print("未送达的用户: ")
81     print(pd.DataFrame(unassigned_customers))
82     assign_customers_demo3_3()

```

输出:

- (1) 5
(2) 43
(3) 2 (未送达375个)

80|14 杨翔 第三题

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 21:01:11 2024
4  @author: Stu
5  """
6  import pandas as pd
7  import numpy as np
8  from sklearn.cluster import KMeans
9  file_path_1 = 'C:/Users/YourUsername/Desktop/Demo_3_1.csv'
10 df_1 = pd.read_csv(file_path_1)

```

```

11 total_demand = df_1['size'].sum()
12 max_delivery_per_employee = 160
13 num_employees = np.ceil(total_demand / max_delivery_per_employee)
14 kmeans = KMeans(n_clusters=int(num_employees), random_state=0).fit(d
f_1[['x', 'y']])
15 df_1['employee_id'] = kmeans.labels_
16 employee_assignments = df_1.groupby('employee_id')['size'].sum().res
et_index()
17 employee_assignments.columns = ['employee_id', 'total_delivery']
18 print(f"需要安排的员工数量: {len(employee_assignments)}")
19 print(employee_assignments)
20 from scipy.spatial.distance import cdist
21 file_path_2 = 'C:/Users/YourUsername/Desktop/Demo_3_2.csv'
22 df_2 = pd.read_csv(file_path_2)
23 coords = df_2[['x', 'y']].values
24 distances = cdist(coords, coords, metric='euclidean')
25 def is_valid_group(group, max_distance):
26     group_coords = group[['x', 'y']].values
27     group_distances = cdist(group_coords, group_coords, metric='eucl
idean')
28     np.fill_diagonal(group_distances, np.inf)
29     return np.all(group_distances <= max_distance)
30 max_distance = 25
31 employees = []
32 remaining_customers = df_2.copy()
33 while not remaining_customers.empty:
34     start_customer = remaining_customers.iloc[0]
35     employee_group = [start_customer]
36     remaining_customers = remaining_customers.drop(start_customer.in
dex)
37     for idx, customer in remaining_customers.iterrows():
38         temp_group = pd.concat([pd.DataFrame(employee_group), pd.Dat
aFrame([customer])])
39         if is_valid_group(temp_group, max_distance) and temp_group
['size'].sum() <= max_delivery_per_employee:
40             employee_group.append(customer)
41             remaining_customers = remaining_customers.drop(idx)
42         else:
43             break
44     employees.append(pd.DataFrame(employee_group))
45 print(f"需要安排的员工数量: {len(employees)}")
46 for i, employee in enumerate(employees):
47     print(f"员工 {i+1} 负责的客户:")
48     print(employee)
49 file_path_3 = 'C:/Users/YourUsername/Desktop/Demo_3_3.csv'
50 df_3 = pd.read_csv(file_path_3)
51 max_distance_per_trip = 15
52 employees = [[], []]
53 remaining_customers = df_3.copy()

```

```

54 unreachable_customers = []
55 while not remaining_customers.empty:
56     # 这里我们简单地选择距离当前员工组中最近的客户进行分配（贪心策略）
57     min_distance_to_employees = np.inf
58     min_distance_idx = -1
59     min_employee_idx = -1
60     for employee_idx, employee_group in enumerate(employees):
61         if not employee_group: # 如果员工组为空，则直接分配给客户
62             min_distance_to_employees = 0
63             min_employee_idx = employee_idx
64             break
65         # 计算剩余客户到当前员工组中每个客户的距离，并找到最小距离
66         employee_coords = pd.DataFrame(employee_group)[['x', 'y']].values
67         for idx, customer in remaining_customers.iterrows():
68             customer_coord = np.array([[customer['x'], customer
69 ['y']]])
70             distances = cdist(customer_coord, employee_coords, metric='euclidean').flatten()
71             min_distance = np.min(distances)
72             if min_distance < min_distance_to_employees:
73                 min_distance_to_employees = min_distance
74                 min_distance_idx = idx
75                 min_employee_idx = employee_idx
76         # 检查分配是否满足条件
77         if min_distance_to_employees <= max_distance_per_trip and (
78             employees[min_employee_idx]['size'].sum() + remaining_customers.loc[min_distance_idx, 'size']) <= max_delivery_per_employee:
79             # 分配客户给员工
80             employees[min_employee_idx].append(remaining_customers.loc[min_distance_idx].to_dict())
81             remaining_customers = remaining_customers.drop(min_distance_idx)
82         else:
83             # 如果不满足条件，则将该客户标记为无法送达
84             unreachable_customers.append(remaining_customers.loc[min_distance_idx].to_dict())
85             remaining_customers = remaining_customers.drop(min_distance_idx)
86     employees_dfs = [pd.DataFrame(employee) for employee in employees]
87     unreachable_customers_df = pd.DataFrame(unreachable_customers)
88     print(f"员工1负责的客户:")
89     print(employees_dfs[0])
90     print(f"\n员工2负责的客户:")
91     print(employees_dfs[1])
92     print(f"\n无法送达的客户:")
93     print(unreachable_customers_df)

```

输出:

- (1) 5
- (2) 445
- (3) 2 (无法运行)

81|15 王丽媛 第三题

```
1  # 3
2  # (1)
3  import pandas as pd
4  data_csv = pd.read_csv('D:Demo_3_1.csv')
5  max_delivery = 160
6  data_csv = data_csv.sort_values(by='size', ascending=False)
7  employee_id = 1
8  employee_assignments = {}
9  current_delivery = 0
10 for index, row in data_csv.iterrows():
11     if current_delivery + row['size'] <= max_delivery:
12         if employee_id not in employee_assignments:
13             employee_assignments[employee_id] = []
14             employee_assignments[employee_id].append(index)
15             current_delivery += row['size']
16     else:
17         employee_id += 1
18         employee_assignments[employee_id] = [index]
19         current_delivery = row['size']
20 for employee, customers in employee_assignments.items():
21     print(f"员工{employee}负责的客户: {customers}")
22 print(f"需要安排的员工数量: {employee_id}")
23 # (2)
24 import pandas as pd
25 df = pd.read_csv('D:Demo_3_2.csv')
26 import numpy as np
27 from scipy.spatial.distance import pdist, squareform
28 coordinates = df[['x', 'y']].values
29 distance_matrix = squareform(pdist(coordinates, metric='euclidean'))
30 num_employees = 0
31 customer_groups = []
32 while distance_matrix.size > 0:
33     min_distance = np.min(distance_matrix[np.nonzero(distance_matrix)])
34     if min_distance > 25:
35         for i in range(distance_matrix.shape[0]):
36             customer_groups.append([i])
37             break
38     else:
39         indices = np.where(distance_matrix == min_distance)
40         i, j = indices[0][0], indices[1][0]
```

```

42     group = [i, j]
43     customer_groups.append(group)
44     distance_matrix = np.delete(distance_matrix, i, axis=0)
45     distance_matrix = np.delete(distance_matrix, i, axis=1)
46     distance_matrix = np.delete(distance_matrix, j - 1, axis=0)
47     distance_matrix = np.delete(distance_matrix, j - 1, axis=1)
48     num_employees += 1
49 print('需要安排的员工数量: ', num_employees)
50 print('每个员工负责的客户: ')
51 for i, group in enumerate(customer_groups):
52     print(f'员工{i + 1}: ', group)
53 # (3)
54 import pandas as pd
55 df = pd.read_csv('D:Demo_3_3.csv')
56 import numpy as np
57 from scipy.spatial.distance import pdist, squareform
58 coordinates = df[['x', 'y']].values
59 distance_matrix = squareform(pdist(coordinates, metric='euclidean'))
60 num_employees = 2
61 customer_groups = [[] for _ in range(num_employees)]
62 while distance_matrix.size > 0:
63     min_distance = np.min(distance_matrix[np.nonzero(distance_matrix)])
64     if min_distance > 15:
65         for i in range(distance_matrix.shape[0]):
66             customer_groups.append([i])
67         break
68     else:
69         indices = np.where(distance_matrix == min_distance)
70         i, j = indices[0][0], indices[1][0]
71         group = [i, j]
72         customer_groups[0].append(i)
73         customer_groups[1].append(j)
74         distance_matrix = np.delete(distance_matrix, i, axis=0)
75         distance_matrix = np.delete(distance_matrix, i, axis=1)
76         distance_matrix = np.delete(distance_matrix, j - 1, axis=0)
77         distance_matrix = np.delete(distance_matrix, j - 1, axis=1)
78 print('员工1负责的客户: ', customer_groups[0])
79 print('员工2负责的客户: ', customer_groups[1])
80 print('不能送达的客户: ', list(set(range(df.shape[0])) - set(customer_groups[0]) - set(customer_groups[1])))

```

输出:

- (1) 5
- (2) 217
- (3) 2 (不能送达有197个)

82|16 王乐 第三题

83|17 王康宇 第三题

20

```
1  import pandas as pd
2  from sklearn.cluster import KMeans
3  import matplotlib.pyplot as plt
4  #第一问
5  # 读取csv文件
6  data_csv = pd.read_csv('Demo_3_1.csv')
7  # 为客户添加序号
8  data_csv['customer_id'] = range(1, len(data_csv) + 1)
9  # 定义每个员工的最大配送量
10 max_size = 160
11 # 计算每个客户的需求总和
12 data_csv['total_size'] = data_csv['size'].sum()
13 # 确定聚类的数量范围
14 k_values = range(1, 11) # 假设尝试k从1到10
15 # 计算不同k值下的聚类误差
16 sse = []
17 for k in k_values:
18     kmeans = KMeans(n_clusters=k, random_state=42)
19     kmeans.fit(data_csv[['x', 'y']])
20     sse.append(kmeans.inertia_)
21 # 根据肘部法则选择合适的k值
22 # 这里简单假设选择肘部明显的点，实际应用中可能需要更复杂的判断
23 elbow_point = 5 # 假设肘部明显的点是k=3
24 n_clusters = elbow_point
25 # 确保聚类数量不大于样本数量
26 n_clusters = min(n_clusters, data_csv.shape[0])
27 # 进行K-means聚类
28 kmeans = KMeans(n_clusters=n_clusters, random_state=42)
29 data_csv['cluster'] = kmeans.fit_predict(data_csv[['x', 'y']])
30 # 输出每个聚类的客户数量和总需求
31 grouped_data = data_csv.groupby('cluster').agg({'size': 'sum', 'x': 'count'}).reset_index()
32 grouped_data.columns = ['cluster', 'total_size', 'customer_count']
33 print(grouped_data)
34 # 输出需要的人数
35 num_employees = grouped_data.shape[0]
36 print(f'需要的人数: {num_employees}')
37 # 检查是否有员工负责的需求超过最大配送量
38 while grouped_data['total_size'].max() > max_size:
39     # 增加聚类数量
40     n_clusters += 1
41     # 重新进行K-means聚类
42     kmeans = KMeans(n_clusters=n_clusters, random_state=42)
43     data_csv['cluster'] = kmeans.fit_predict(data_csv[['x', 'y']])
```

```

44     # 重新计算每个聚类的客户数量和总需求
45     grouped_data = data_csv.groupby('cluster').agg({'size': 'sum',
46     'x': 'count'}).reset_index()
47     grouped_data.columns = ['cluster', 'total_size', 'customer_count']
48     # 输出新的聚类结果
49     print('重新聚类后的结果: ')
50     print(grouped_data)
51     # 输出需要的人数
52     num_employees = grouped_data.shape[0]
53     print(f'需要的人数: {num_employees}')
54     # 输出每个员工负责的客户标号
55     for cluster in grouped_data['cluster']:
56         customers = data_csv[data_csv['cluster'] == cluster]['customer_id']
57         print(f'员工{cluster}负责的客户标号: ')
58         print(customers)
59     # 添加负责的员工编号列
60     data_csv['employee_id'] = data_csv['cluster']
61     #第二小题
62     import pandas as pd
63     from sklearn.cluster import KMeans
64     from sklearn.metrics import silhouette_score
65     import numpy as np
66     # 加载数据集
67     df = pd.read_csv('Demo_3_2.csv')
68     # 给客户编号
69     df = df.reset_index().rename(columns={'index': 'customer_id'})
70     # 定义一个函数，用于计算不同聚类数量下的SSE和轮廓系数
71     def calculate_sse_and_silhouette(data, max_clusters):
72         sse = []
73         silhouette_scores = []
74         for k in range(1, max_clusters + 1):
75             kmeans = KMeans(n_clusters=k)
76             kmeans.fit(data)
77             sse.append(kmeans.inertia_)
78             if k > 1:
79                 labels = kmeans.labels_
80                 score = silhouette_score(data, labels)
81                 silhouette_scores.append(score)
82         return sse, silhouette_scores
83     # 计算不同聚类数量下的SSE和轮廓系数
84     sse, silhouette_scores = calculate_sse_and_silhouette(df[['x',
85     'y']], 10)
86     # 输出结果
87     for k in range(1, 10):
88         print(f'聚类数量为{k}时，SSE为{sse[k - 1]}，轮廓系数为{silhouette_scores[k - 2] if k > 1 else "无"}')
89     # 进行聚类分析，设置聚类数量为9

```

```

90 kmeans = KMeans(n_clusters=9, random_state=42)
91 df['cluster'] = kmeans.fit_predict(df[['x', 'y']])
92 # 计算每个聚类的中心
93 cluster_centers = kmeans.cluster_centers_
94 # 计算每个客户到其所属聚类中心的距离
95 df['distance_to_center'] = np.linalg.norm(df[['x', 'y']] - cluster_c
96 enters[df['cluster']], axis=1)
97 # 判断是否有距离超过25公里的客户
98 if df['distance_to_center'].max() > 25:
99     # 如果有, 增加聚类数量, 重新进行聚类
100     max_clusters = 15
101     for k in range(10, max_clusters + 1):
102         kmeans = KMeans(n_clusters=k, random_state=42)
103         df['cluster'] = kmeans.fit_predict(df[['x', 'y']])
104         cluster_centers = kmeans.cluster_centers_
105         df['distance_to_center'] = np.linalg.norm(df[['x', 'y']] - c
106 luster_centers[df['cluster']], axis=1)
107         if df['distance_to_center'].max() <= 25:
108             break
109 # 计算每个聚类中的客户数量和总需求
110 cluster_summary = df.groupby('cluster').agg({'customer_id': list, 's
111 ize': ['count', 'sum']}).reset_index()
112 cluster_summary.columns = ['cluster', 'customer_id', 'customer_coun
113 t', 'total_demand']
114 # 查看聚类结果
115 cluster_summary
116 # 确定需要安排的员工数量
117 num_employees = cluster_summary.shape[0]
118 # 输出结果
119 print(f'需要安排的员工数量为: {num_employees}')
120 # 输出每个员工负责的客户
121 for index, row in cluster_summary.iterrows():
122     print(f'员工 {index + 1} 负责的客户: {row["customer_id"]}')
123 #第三小题
124 import pandas as pd
125 from sklearn.cluster import KMeans
126 from sklearn.metrics import silhouette_score
127 import numpy as np
128 def calculate_distance(point1, point2):
129     return np.sqrt((point1[0] - point2[0])** 2 + (point1[1] - point
130 2[1])** 2)
131 def assign_customers_to_employees(data_file):
132     # 加载数据
133     data = pd.read_csv(data_file)
134     # 计算轮廓系数
135     silhouette_avg_max = -1
136     best_clusters = None
137     # 遍历不同的簇数量
138     for n_clusters in range(2, 11):

```

```

134     kmeans = KMeans(n_clusters=n_clusters)
135     kmeans.fit(data[['x', 'y']])
136     labels = kmeans.labels_
137     silhouette_avg = silhouette_score(data[['x', 'y']], labels)
138     # 更新最佳簇数量
139     if silhouette_avg > silhouette_avg_max:
140         silhouette_avg_max = silhouette_avg
141         best_clusters = n_clusters
142     # 使用最佳簇数量进行聚类
143     kmeans = KMeans(n_clusters=best_clusters)
144     kmeans.fit(data[['x', 'y']])
145     data['cluster'] = kmeans.labels_
146     # 计算每个簇的中心位置
147     cluster_centers = kmeans.cluster_centers_
148     # 计算每个簇的客户数量
149     cluster_sizes = data['cluster'].value_counts()
150     # 找出客户数量最多的两个簇
151     top_two_clusters = cluster_sizes.nlargest(2).index
152     # 提取客户数量最多的两个簇的数据
153     top_two_clusters_data = data[data['cluster'].isin(top_two_clusters)]
154     # 计算客户数量最多的两个簇的中心位置
155     top_two_clusters_centers = cluster_centers[top_two_clusters]
156     # 计算每个客户与员工位置的距离
157     data['distance_to_employee1'] = data.apply(lambda row: calculate_
158     _distance(row[['x', 'y']], top_two_clusters_centers[0]), axis=1)
159     data['distance_to_employee2'] = data.apply(lambda row: calculate_
160     _distance(row[['x', 'y']], top_two_clusters_centers[1]), axis=1)
161     # 筛选出距离超过15公里的客户
162     customers_out_of_range_employee1 = data[data['distance_to_employ
163     ee1'] > 15]
164     customers_out_of_range_employee2 = data[data['distance_to_employ
165     ee2'] > 15]
166     # 筛选出在员工15km范围内的客户
167     customers_in_range_employee1 = data[data['distance_to_employee
168     1'] <= 15]
169     customers_in_range_employee2 = data[data['distance_to_employee
170     2'] <= 15]
171     # 输出结果
172     print('员工1负责的客户: ')
173     print(customers_in_range_employee1[['x', 'y', 'size']])
174     print('员工2负责的客户: ')
175     print(customers_in_range_employee2[['x', 'y', 'size']])
176     print('员工1无法送达的客户: ')
177     print(customers_out_of_range_employee1[['x', 'y', 'size']])
178     print('员工2无法送达的客户: ')
179     print(customers_out_of_range_employee2[['x', 'y', 'size']])
180     # 调用函数并传入数据文件路径
181     assign_customers_to_employees('Demo_3_3.csv')

```

输出:

- (1) 6
- (2) 15
- (3) 2 (不能送达有378个)

84|18 王立弘 第三题

```
1 import csv
2 # 存储客户信息的列表，每个元素是一个字典，包含客户位置和需求等信息
3 customers = []
4 # 读取Demo_3_1.csv文件，将客户信息存入customers列表
5 with open('D:\\Demo_3_1.csv', 'r', encoding='utf-8') as csvfile:
6     reader = csv.reader(csvfile)
7     next(reader) # 跳过标题行（假设文件有标题行）
8     for row in reader:
9         customer = {
10             'x': float(row[0]),
11             'y': float(row[1]),
12             'size': int(row[2])
13         }
14         customers.append(customer)
15 # 用于存储分配结果，每个元素是一个列表，表示每个员工负责的客户索引
16 assignment = []
17 # 每个员工当前已分配的配送量
18 employee_loads = []
19 # 开始分配客户给员工
20 customer_index = 0
21 while customer_index < len(customers):
22     current_employee_assignment = []
23     current_employee_load = 0
24     while customer_index < len(customers) and current_employee_load
25         + customers[customer_index]['size'] <= 160:
26         current_employee_assignment.append(customer_index)
27         current_employee_load += customers[customer_index]['size']
28         customer_index += 1
29     assignment.append(current_employee_assignment)
30     employee_loads.append(current_employee_load)
31 # 输出结果，即需要的员工数量和每个员工负责的客户信息
32 print(f"需要安排 {len(assignment)} 个员工。")
33 for i in range(len(assignment)):
34     print(f"员工 {i + 1} 负责的客户索引为: {assignment[i]}")
35 import pandas as pd
36 from sklearn.cluster import KMeans
37 import numpy as np
38 # 读取文件，假设文件编码为UTF-8（可根据实际情况调整）
39 data = pd.read_csv("D:\\Demo_3_2.csv", encoding="UTF-8")
40 # 获取客户位置坐标
41 customer_locations = data[['x', 'y']].values
```

```

41 # 使用KMeans聚类算法进行聚类，这里假设聚成n个类（n可根据实际情况调整）
42 n_clusters = 3 # 可调整的聚类数量，可根据业务情况预估大概需要几个员工来设置
43 kmeans = KMeans(n_clusters=n_clusters, random_state=42)
44 kmeans.fit(customer_locations)
45 # 获取每个客户所属的聚类簇标签（也就是对应员工的编号，这里从0开始编号）
46 labels = kmeans.labels_
47 # 为每个客户进行编号（从1开始编号）
48 customer_numbers = np.arange(1, len(customer_locations) + 1)
49 # 用于存储每个员工负责的客户编号列表
50 employees_customers = [[] for _ in range(n_clusters)]
51 for customer_number, label in zip(customer_numbers, labels):
52     employees_customers[label].append(customer_number)
53 # 输出结果
54 for i in range(n_clusters):
55     print(f"员工 {i + 1} 负责的客户编号为: {employees_customers[i]}")
56 import pandas as pd
57 from sklearn.cluster import DBSCAN
58 import numpy as np
59 from scipy.spatial.distance import pdist, squareform
60 # 读取文件，假设文件编码为UTF-8（可根据实际情况调整）
61 data = pd.read_csv("D:\\Demo_3_3.csv", encoding="UTF-8")
62 # 获取客户位置坐标和需求信息
63 customer_locations = data[['x', 'y']].values
64 customer_demands = data['size'].values
65 # 设置DBSCAN算法的参数，eps表示邻域半径，min_samples表示核心点的最小邻居数量
66 # 这里根据距离限制设置eps为15（单位和数据中的坐标一致，即距离不超过15km对应的坐标距离）
67 # min_samples可根据实际情况调整，一般设置为至少3以上保证聚类的合理性
68 eps = 15
69 min_samples = 3
70 dbscan = DBSCAN(eps=eps, min_samples=min_samples)
71 # 对客户位置数据进行聚类
72 dbscan.fit(customer_locations)
73 # 获取每个客户所属的聚类簇标签，-1表示噪声点（在这里可认为是不能送达的客户）
74 labels = dbscan.labels_
75 # 用于存储最终每个员工负责的客户索引
76 employees_customers = [[] for _ in range(len(set(labels)) - 1 if -1
77     in labels else len(set(labels)))]
78 # 用于存储不能送达的客户索引（即聚类标签为-1的客户）
79 unreachable_customers = np.where(labels == -1)[0]
80 # 遍历所有的聚类簇（除了噪声点对应的-1标签）
81 for cluster_id in set(labels) - {-1}:
82     cluster_customers = np.where(labels == cluster_id)[0]
83     employees_customers[cluster_id] = list(cluster_customers)
84 # 输出结果
85 print("员工1负责的客户索引为:", employees_customers[0] if employees_customers else [])
86 print("员工2负责的客户索引为:", employees_customers[1] if len(employees

```


```
86 | _customers) > 1 else [])  
    print("不能送达的客户索引为:", unreachable_customers)
```

输出:

```
(1) 5  
(2) 3  
(3) 2 (193 203 300)
```



85|19 白天琪 第三题



```
1 | import pandas as pd  
2 | # 读取CSV文件  
3 | data_3_1 = pd.read_csv("Demo_3_1.csv")  
4 | data_3_2 = pd.read_csv("Demo_3_2.csv")  
5 | data_3_3 = pd.read_csv("Demo_3_3.csv")  
6 | # 查看数据内容  
7 | print(data_3_1.head())  
8 | print(data_3_2.head())  
9 | print(data_3_3.head())  
10 | import numpy as np  
11 | def assign_employees_based_on_size(data, max_delivery_size):  
12 |     # 按照需求量从大到小排序客户  
13 |     sorted_data = data.sort_values(by='size', ascending=False)  
14 |     employees = []  
15 |     current_employee = []  
16 |     current_total_size = 0  
17 |     for _, row in sorted_data.iterrows():  
18 |         if current_total_size + row['size'] > max_delivery_size:  
19 |             # 如果当前员工配送量超过限制, 分配给下一个员工  
20 |             employees.append(current_employee)  
21 |             current_employee = [row]  
22 |             current_total_size = row['size']  
23 |         else:  
24 |             current_employee.append(row)  
25 |             current_total_size += row['size']  
26 |     # 添加最后一个员工的任务  
27 |     if current_employee:  
28 |         employees.append(current_employee)  
29 |     return employees  
30 | # 调用函数  
31 | employees_3_1 = assign_employees_based_on_size(data_3_1, 160)  
32 | # 输出每个员工负责的客户  
33 | for i, employee in enumerate(employees_3_1):  
34 |     print(f"员工 {i + 1} 负责的客户: ")  
35 |     for customer in employee:  
36 |         print(f"客户位置({customer['x']}, {customer['y']}), 需求量: {customer['size']}")
```

```

37 from scipy.spatial.distance import cdist
38 from sklearn.cluster import DBSCAN
39 def assign_employees_based_on_distance(data, max_distance):
40     # 获取客户坐标
41     coordinates = data[['x', 'y']].values
42     # 计算距离矩阵
43     dist_matrix = cdist(coordinates, coordinates, metric='euclidean')
44     # DBSCAN聚类, eps为最大距离, min_samples为最小样本数
45     clustering = DBSCAN(eps=max_distance, min_samples=1, metric='precomputed').fit(dist_matrix)
46     # 获取每个客户的分组
47     labels = clustering.labels_
48     # 按照分组分配客户
49     employees = {}
50     for label, customer in zip(labels, data.iterrows()):
51         if label not in employees:
52             employees[label] = []
53             employees[label].append(customer[1])
54     return employees
55 # 调用函数
56 employees_3_2 = assign_employees_based_on_distance(data_3_2, 25)
57 # 输出每个员工负责的客户
58 for i, (employee_id, customers) in enumerate(employees_3_2.items()):
59     print(f"员工 {i + 1} 负责的客户: ")
60     for customer in customers:
61         print(f"客户位置({customer['x']}, {customer['y']}), 需求量: {customer['size']}")
62 def assign_two_employees(data, max_distance):
63     # 获取客户坐标
64     coordinates = data[['x', 'y']].values
65     # 计算距离矩阵
66     dist_matrix = cdist(coordinates, coordinates, metric='euclidean')
67     # 找到两名员工的分配方案
68     employee1, employee2 = [], []
69     unassigned_customers = []
70     for i in range(len(data)):
71         # 分配给员工1, 检查与其他客户的距离
72         if all(dist_matrix[i, j] >= max_distance for j in range(len(data)) if j != i):
73             employee1.append(data.iloc[i])
74         else:
75             employee2.append(data.iloc[i])
76     # 如果任何客户与其他客户的距离超过限制, 则不能配送
77     if len(employee1) + len(employee2) < len(data):
78         unassigned_customers = [customer for i, customer in data.iterrows() if
79                                 customer not in employee1 and customer not in employee2]
80

```



```

81     er not in employee2]
82     return employee1, employee2, unassigned_customers
    # 调用函数
83     employee1, employee2, unassigned_customers = assign_two_employees(da
84     ta_3_3, 15)
85     # 输出每个员工负责的客户
86     print("员工 1 负责的客户: ")
    for customer in employee1:
87         print(f"客户位置({customer['x']}, {customer['y']}), 需求量: {custo
88         mer['size']}")
89     print("员工 2 负责的客户: ")
    for customer in employee2:
90         print(f"客户位置({customer['x']}, {customer['y']}), 需求量: {custo
91         mer['size']}")
92     # 输出无法配送的客户
93     if unassigned_customers:
94         print("无法配送的客户: ")
        for customer in unassigned_customers:
95             print(f"客户位置({customer['x']}, {customer['y']}), 需求量: {c
            ustomer['size']}")

```

输出:

- (1) 5
- (2) 3
- (3) 2 (无)

86|20 罗昊然 第三题

```

1  import csv
2  import math
3  def read_csv(file_path):
4      # 读取CSV文件并计算总需求
5      total_demand = 0
6      customers = []
7      with open(file_path, 'r', encoding='utf-8') as f:
8          reader = csv.DictReader(f)
9          for row in reader:
10             x, y, size = float(row['x']), float(row['y']), int(row
                ['size'])
11             customers.append((x, y, size))
12             total_demand += size
13     return total_demand, customers
14 def calculate_distance(x1, y1, x2, y2):
15     # 计算两点之间的距离
16     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
17 def assign_customers_by_capacity(customers, capacity):

```

```

18     # 计算所需员工数量
19     total_demand = sum(c[2] for c in customers)
20     employees_needed = -(-total_demand // capacity) # 向上取整
21     # 分配客户给员工
22     employee_assignment = {i: [] for i in range(employees_needed)}
23     # 贪心算法分配客户
24     for customer in sorted(customers, key=lambda x: x[2]): # 按需求
        从小到大排序
25         min_load_employee = min(employee_assignment, key=lambda x: s
um(c[2] for c in employee_assignment[x]))
26         employee_assignment[min_load_employee].append(customer)
27     return employee_assignment
28 def assign_customers_by_distance(customers, max_distance):
29     # 分配客户给员工
30     employee_assignment = []
31     current_employee = [customers[0]]
32     for customer in customers[1:]:
33         # 检查是否可以分配给当前员工
34         assigned = False
35         for idx, emp_customer in enumerate(employee_assignment):
36             if all(
calculate_distance(emp_customer[0], emp_customer
[1], customer[0], customer[1]) < max_distance for
37                 emp_customer in emp_customer):
38                 emp_customer.append(customer)
39                 assigned = True
40                 break
41         if not assigned:
42             employee_assignment.append([customer])
43     return employee_assignment
44 def assign_customers_by_distance_and_employees(customers, max_distan
ce, num_employees):
45     # 分配客户给两个员工
46     employee_assignment = {i: [] for i in range(num_employees)}
47     unassigned_customers = []
48     for customer in customers:
49         # 检查是否可以分配给任一员工
50         assigned = False
51         for emp, emp_customers in employee_assignment.items():
52             if not emp_customers or all(
53                 calculate_distance(emp_customers[-1][0], emp_cus
tomers[-1][1], customer[0], customer[1]) <= max_distance
54                 for emp_customer in emp_customers):
55                 employee_assignment[emp].append(customer)
56                 assigned = True
57                 break
58         if not assigned:
59             # 无法送达的客户
60             unassigned_customers.append(customer)
61     return employee_assignment, unassigned_customers

```

```

62 # 主程序
63 if __name__ == "__main__":
64     # 问题 (1)
65     total_demand, customers = read_csv('Demo_3_1.csv')
66     employee_assignment = assign_customers_by_capacity(customers, 16
67 0)
68     print(f"需要安排的员工数量: {len(employee_assignment)}")
69     for emp, custs in employee_assignment.items():
70         print(f"员工 {emp + 1} 负责的客户:")
71         for cust in custs:
72             print(f"客户位置: ({cust[0]}, {cust[1]}), 需求: {cust[2]}
73 瓶")
74     # 问题 (2)
75     _, customers = read_csv('Demo_3_2.csv')
76     employee_assignment = assign_customers_by_distance(customers, 2
77 5)
78     print(f"需要安排的员工数量: {len(employee_assignment)}")
79     for idx, emp_customers in enumerate(employee_assignment, start=
80 1):
81         print(f"员工 {idx} 负责的客户:")
82         for cust in emp_customers:
83             print(f"客户位置: ({cust[0]}, {cust[1]}), 需求: {cust[2]}
84 瓶")
85     # 问题 (3)
86     _, customers = read_csv('Demo_3_3.csv')
87     employee_assignment, unassigned_customers = assign_customers_by_
88 distance_and_employees(customers, 15, 2)
89     print(f"需要安排的员工数量: {len(employee_assignment)}")
90     for idx, emp_customers in employee_assignment.items():
91         print(f"员工 {idx} 负责的客户:")
92         for cust in emp_customers:
93             print(f"客户位置: ({cust[0]}, {cust[1]}), 需求: {cust[2]}
94 瓶")
95     print(f"未送达的客户数量: {len(unassigned_customers)}")
96     for cust in unassigned_customers:
97         print(f"客户位置: ({cust[0]}, {cust[1]}), 需求: {cust[2]} 瓶无
98 法送达")

```

输出:

- (1) 5
(2) 56
(3) 2 (375)

87|21 谢皓椿 第三题

```

1 import csv
2 # 读取Demo_3_1.csv文件

```

```
3 with open('/mnt/Demo_3_1.csv', mode='r', encoding='utf-8') as file:
4     csv_reader = csv.reader(file)
5     next(csv_reader) # 跳过表头
6     customers = [(float(row[0]), float(row[1]), int(row[2])) for row
in csv_reader]
7 # 按照需求从大到小排序
8 customers.sort(key=lambda x: x[2], reverse=True)
9 # 初始化员工列表和当前员工的需求
10 employees = []
11 current_employee = []
12 current_demand = 0
13 # 分配客户给员工
14 for customer in customers:
15     if current_demand + customer[2] <= 160:
16         current_employee.append(customer)
17         current_demand += customer[2]
18     else:
19         employees.append(current_employee)
20         current_employee = [customer]
21         current_demand = customer[2]
22 # 将最后一个员工添加到员工列表中
23 employees.append(current_employee)
24 # 打印结果
25 for i, employee in enumerate(employees, 1):
26     print(f"员工{i}负责的客户: {employee}")
27 print(f"需要安排的员工数量: {len(employees)}")
28 import csv
29 import math
30 def distance(x1, y1, x2, y2):
31     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
32 # 读取Demo_3_2.csv文件
33 with open('/mnt/Demo_3_2.csv', mode='r', encoding='utf-8') as file:
34     csv_reader = csv.reader(file)
35     next(csv_reader) # 跳过表头
36     customers = [(float(row[0]), float(row[1]), int(row[2])) for row
in csv_reader]
37 # 初始化员工列表和当前员工的客户列表
38 employees = []
39 current_employee = [customers[0]]
40 # 分配客户给员工
41 for i in range(1, len(customers)):
42     customer = customers[i]
43     assignable = True
44     for assigned_customer in current_employee:
45         if distance(customer[0], customer[1], assigned_customer[0],
assigned_customer[1]) > 25:
46             assignable = False
47             break
48     if assignable:
```

```

49         current_employee.append(customer)
50     else:
51         employees.append(current_employee)
52         current_employee = [customer]
53 # 将最后一个员工添加到员工列表中
54 employees.append(current_employee)
55 # 打印结果
56 for i, employee in enumerate(employees, 1):
57     print(f"员工{i}负责的客户: {employee}")
58 print(f"需要安排的员工数量: {len(employees)}")
59 import csv
60 # 读取Demo_3_3.csv文件
61 with open('/mnt/Demo_3_3.csv', mode='r', encoding='utf-8') as file:
62     csv_reader = csv.reader(file)
63     next(csv_reader) # 跳过表头
64     customers = [(float(row[0]), float(row[1]), int(row[2])) for row
65 in csv_reader]
66 # 按照需求从大到小排序
67 customers.sort(key=lambda x: x[2], reverse=True)
68 # 初始化员工列表和当前员工的需求
69 employees = []
70 current_employee = []
71 current_demand = 0
72 # 分配客户给员工
73 for customer in customers:
74     if current_demand + customer[2] <= 160:
75         current_employee.append(customer)
76         current_demand += customer[2]
77     else:
78         employees.append(current_employee)
79         current_employee = [customer]
80         current_demand = customer[2]
81 # 将最后一个员工添加到员工列表中
82 employees.append(current_employee)
83 # 打印结果
84 for i, employee in enumerate(employees, 1):
85     print(f"员工{i}负责的客户: {employee}")
86 print(f"需要安排的员工数量: {len(employees)}")

```

输出:

- (1) 5
- (2) 395
- (3) 4 (无)

88|22 贺馨姜艾 第三题

19

```

1 import pandas as pd
2 data = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_3_1.
  csv')
3 employees = []
4 current_employee_customers = []
5 current_size_sum = 0
6 for index, row in data.iterrows():
7     customer_size = row['size']
8     if current_size_sum + customer_size <= 160:
9         current_employee_customers.append(index)
10        current_size_sum += customer_size
11    else:
12        employees.append(current_employee_customers)
13        current_employee_customers = [index]
14        current_size_sum = customer_size
15 if current_employee_customers:
16     employees.append(current_employee_customers)
17 print("需要安排的员工数量为:", len(employees))
18 for i, employee_customers in enumerate(employees):
19     print(f"员工{i + 1}负责的客户为:", employee_customers)
20 import pandas as pd
21 import math
22 data = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_3_2.
  csv')
23 def distance(point1, point2):
24     return math.sqrt((point1[0] - point2[0]) ** 2 + (point2[1] - poi
  nt2[1]) ** 2)
25 employees = []
26 current_employee_customers = []
27 for index1 in range(len(data)):
28     can_assign = True
29     for index2 in current_employee_customers:
30         point1 = (data.loc[index1, 'x'], data.loc[index1, 'y'])
31         point2 = (data.loc[index2, 'x'], data.loc[index2, 'y'])
32         if distance(point1, point2) > 25:
33             can_assign = False
34             break
35     if can_assign:
36         current_employee_customers.append(index1)
37     else:
38         employees.append(current_employee_customers)
39         current_employee_customers = [index1]
40 if current_employee_customers:
41     employees.append(current_employee_customers)
42 print("需要安排的员工数量为:", len(employees))
43 for i, employee_customers in enumerate(employees):
44     print(f"员工{i + 1}负责的客户为:", employee_customers)
45 import pandas as pd
46 import math

```

```

47 def distance(point1, point2):
48     return math.sqrt((point1[0] - point2[0])** 2 + (point1[1] - poi
nt2[1])** 2)
49 data = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_3_3.
csv')
50 employee1_customers = []
51 employee2_customers = []
52 unreachable_customers = []
53 for index1 in range(len(data)):
54     can_assign_to_employee1 = True
55     for index2 in employee1_customers:
56         point1 = (data.loc[index1, 'x'], data.loc[index1, 'y'])
57         point2 = (data.loc[index2, 'x'], data.loc[index2, 'y'])
58         if distance(point1, point2) > 15:
59             can_assign_to_employee1 = False
60             break
61     if can_assign_to_employee1:
62         employee1_customers.append(index1)
63     else:
64         can_assign_to_employee2 = True
65         for index2 in employee2_customers:
66             point1 = (data.loc[index1, 'x'], data.loc[index1, 'y'])
67             point2 = (data.loc[index2, 'x'], data.loc[index2, 'y'])
68             if distance(point1, point2) > 15:
69                 can_assign_to_employee2 = False
70                 break
71         if can_assign_to_employee2:
72             employee2_customers.append(index1)
73         else:
74             unreachable_customers.append(index1)
75 print("员工1负责的客户为:", employee1_customers)
76 print("员工2负责的客户为:", employee2_customers)
77 print("无法送达的客户为:", unreachable_customers)

```

输出:

- (1) 5
- (2) 327
- (3) 2 (不能送达有383个)

89|23 达尔汗 第三题

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.cluster import KMeans
4 import matplotlib.pyplot as plt
5 # 定义函数用于根据客户位置、需求以及员工最大配送量分配员工任务
6 def assign_clusters(max_delivery_size):

```

```

7         """
8         根据给定的最大配送量，对客户进行聚类并分配员工。
9         参数：
10        max_delivery_size (int): 每个员工每天最多配送的牛奶瓶数。
11        返回：
12        dict: 员工分配情况的字典，键为 (cluster_id, employee_index)，值为对应的客户索引列表。
13        """
14        # 读取位于D盘的Demo_3_1.csv文件，根据实际情况调整路径
15        data = pd.read_csv('D:/Demo_3_1.csv')
16        # 提取x,y坐标以及客户需求size这三列数据
17        coordinates = data[['x', 'y']].values
18        sizes = data['size'].values
19        # 初始化聚类模型，这里暂时先假设聚成合适的类数，实际可能需要多次尝试确定合适的簇数量
20        kmeans = KMeans(n_clusters=5) # 这里5是示例，可根据实际调整
21        kmeans.fit(coordinates)
22        labels = kmeans.labels_
23        cluster_assignments = {}
24        for i in range(len(labels)):
25            label = labels[i]
26            if label not in cluster_assignments:
27                cluster_assignments[label] = {'customers': [], 'total_size': 0}
28            cluster_assignments[label]['customers'].append(i)
29            cluster_assignments[label]['total_size'] += sizes[i]
30        # 根据每个簇的总配送量来确定员工分配
31        employees_per_cluster = {}
32        for cluster, info in cluster_assignments.items():
33            total_size = info['total_size']
34            num_employees = int(np.ceil(total_size / max_delivery_size))
35            employees_per_cluster[cluster] = num_employees
36        final_assignments = {}
37        for cluster, num_emp in employees_per_cluster.items():
38            customers_in_cluster = cluster_assignments[cluster]['customers']
39            # 平均分配客户给每个负责该簇的员工（这里简单平均分配，实际可以更优化）
40            step = len(customers_in_cluster) // num_emp
41            for i in range(num_emp):
42                start = i * step
43                end = (i + 1) * step if i < num_emp - 1 else len(customers_in_cluster)
44                employee_customers = customers_in_cluster[start:end]
45                final_assignments[(cluster, i)] = employee_customers
46        return final_assignments
47        # 每个员工每天最多配送160瓶牛奶
48        max_delivery_size = 160
49        employee_assignments = assign_clusters(max_delivery_size)
50        print("员工分配情况如下：")

```



```

51 for (cluster, emp_index), customers in employee_assignments.items():
52     print(f"员工{(cluster, emp_index)}负责的客户编号为: {customers}")
53 import pandas as pd
54 import numpy as np
55 from sklearn.cluster import KMeans
56 from sklearn.metrics.pairwise import euclidean_distances
57 # 读取数据文件
58 data = pd.read_csv('D:/Demo_3_2.csv')
59 # 获取坐标数据以及需求数据
60 coordinates = data[['x', 'y']].values
61 demands = data['size'].values
62 # 先初始化一个较大的聚类数量, 你可以根据实际情况调整, 后续可以根据业务规则等进一步确定合适的聚类数量
63 kmeans = KMeans(n_clusters=10) # 这里假设先分为10类, 可调整
64 kmeans.fit(coordinates)
65 # 获取聚类标签, 每个数据点所属的类别
66 labels = kmeans.labels_
67 # 用来存储每个聚类(即每个员工负责的客户信息)
68 result = {}
69 for i in range(len(labels)):
70     cluster_id = labels[i]
71     if cluster_id not in result:
72         result[cluster_id] = []
73     result[cluster_id].append(i)
74 # 输出每个员工负责的客户序号(这里序号对应原数据中的行号, 从0开始)
75 for employee_id, customer_indices in result.items():
76     print(f"员工 {employee_id} 负责的客户序号: {customer_indices}")
77 # 以下部分可以用来验证每个聚类内的最大距离是否满足要求(可选, 如果需要严谨验证距离是否符合小于25km要求)
78 for cluster_id in result:
79     cluster_customers = coordinates[np.array(result[cluster_id])]
80     distances = euclidean_distances(cluster_customers)
81     max_distance = np.max(distances)
82     print(f"聚类 {cluster_id} 内的最大距离: {max_distance}")
83 import os
84 import pandas as pd
85 import numpy as np
86 from sklearn.cluster import KMeans
87 import matplotlib.pyplot as plt
88 from scipy.spatial.distance import cdist
89 # 在Windows系统下, 临时设置环境变量OMP_NUM_THREADS=2来尝试避免KMeans的内存泄漏问题(仅适用于Windows且有相关需求时)
90 if os.name == 'nt':
91     os.environ['OMP_NUM_THREADS'] = '2'
92 # 读取数据文件, 这里假设文件路径为 D:\Demo_3_3.csv, 根据实际情况调整
93 data = pd.read_csv(r'D:\Demo_3_3.csv')
94 # 提取坐标和需求数据
95 X = data[['x', 'y']].values
96 # 使用KMeans进行聚类, 将客户分为两类(对应两个员工), 显式设置n_init参数来消除

```

```

97 FutureWarning
98 kmeans = KMeans(n_clusters=2, n_init=10, random_state=0).fit(X)
99 # 获取聚类标签
100 labels = kmeans.labels_
101 # 按照聚类标签将客户数据分组，分别对应两个员工负责的客户情况
102 customers_per_employee_1 = data[labels == 0]
103 customers_per_employee_2 = data[labels == 1]
104 # 距离限制
105 max_distance = 15 # 单段行驶距离限制，单位与坐标对应距离单位一致
106 # 检查并调整聚类以满足距离限制，记录不能送达的客户
107 unserved_customers = []
108 while True:
109     need_adjust = False
110     # 检查员工1负责的客户聚类情况
111     if not customers_per_employee_1.empty:
112         distances_1 = cdist(customers_per_employee_1[['x', 'y']].values,
113 customers_per_employee_1[['x', 'y']].values)
114         max_distance_1 = np.max(distances_1)
115         if max_distance_1 > max_distance:
116             max_index_1_1, max_index_1_2 = np.unravel_index(np.argmax(
117 distances_1), distances_1.shape)
118             customer_1_1 = customers_per_employee_1.iloc[max_index_1
119 _1]
120             customer_1_2 = customers_per_employee_1.iloc[max_index_1
121 _2]
122             # 处理customer_1_1数据维度问题并计算距离
123             if len(customer_1_1[['x', 'y']].values.shape) == 1:
124                 customer_1_1_values = customer_1_1[['x', 'y']].value
125 s.reshape(1, -1)
126             else:
127                 customer_1_1_values = customer_1_1[['x', 'y']].value
128 s
129             distances_to_2 = cdist(customer_1_1_values, customers_per_
130 employee_2[['x', 'y']].values)
131             # 处理customer_1_2数据维度问题并计算距离
132             if len(customer_1_2[['x', 'y']].values.shape) == 1:
133                 customer_1_2_values = customer_1_2[['x', 'y']].value
134 s.reshape(1, -1)
135             else:
136                 customer_1_2_values = customer_1_2[['x', 'y']].value
137 s
138             distances_to_2_2 = cdist(customer_1_2_values, customers_
139 per_employee_2[['x', 'y']].values)
140             if np.all(distances_to_2 <= max_distance) or np.all(dist
141 ances_to_2_2 <= max_distance):
142                 if np.all(distances_to_2 <= max_distance):
143                     customers_per_employee_1 = customers_per_employe
144 e_1.drop(customer_1_1.name)
145                     customers_per_employee_2 = pd.concat([customers_

```

```

136 per_employee_2, customer_1_1.to_frame().T])
137         else:
138             customers_per_employee_1 = customers_per_employee_1
139             e_1.drop(customer_1_2.name)
140             customers_per_employee_2 = pd.concat([customers_per_employee_2, customer_1_2.to_frame().T])
141             need_adjust = True
142             # 检查员工2负责的客户聚类情况，与员工1的检查逻辑类似
143             if not customers_per_employee_2.empty:
144                 distances_2 = cdist(customers_per_employee_2[['x', 'y']].values, customers_per_employee_2[['x', 'y']].values)
145                 max_distance_2 = np.max(distances_2)
146                 if max_distance_2 > max_distance:
147                     max_index_2_1, max_index_2_2 = np.unravel_index(np.argmax(distances_2), distances_2.shape)
148                     customer_2_1 = customers_per_employee_2.iloc[max_index_2_1]
149                     customer_2_2 = customers_per_employee_2.iloc[max_index_2_2]
150
151                     # 处理customer_2_1数据维度问题并计算距离
152                     if len(customer_2_1[['x', 'y']].values.shape) == 1:
153                         customer_2_1_values = customer_2_1[['x', 'y']].values.reshape(1, -1)
154                     else:
155                         customer_2_1_values = customer_2_1[['x', 'y']].values
156
157                     distances_to_1 = cdist(customer_2_1_values, customers_per_employee_1[['x', 'y']].values)
158                     # 处理customer_2_2数据维度问题并计算距离
159                     if len(customer_2_2[['x', 'y']].values.shape) == 1:
160                         customer_2_2_values = customer_2_2[['x', 'y']].values.reshape(1, -1)
161                     else:
162                         customer_2_2_values = customer_2_2[['x', 'y']].values
163
164                     distances_to_1_2 = cdist(customer_2_2_values, customers_per_employee_1[['x', 'y']].values)
165                     if np.all(distances_to_1 <= max_distance) or np.all(distances_to_1_2 <= max_distance):
166                         if np.all(distances_to_1 <= max_distance):
167                             customers_per_employee_2 = customers_per_employee_2.drop(customer_2_1.name)
168                             customers_per_employee_1 = pd.concat([customers_per_employee_1, customer_2_1.to_frame().T])
169                         else:
170                             customers_per_employee_2 = customers_per_employee_2.drop(customer_2_2.name)
171                             customers_per_employee_1 = pd.concat([customers_per_employee_1, customer_2_2.to_frame().T])

```

```

172         need_adjust = True
173         # 检查是否有孤立客户（距离其他所有客户都超过15km）
174         all_customers = pd.concat([customers_per_employee_1, customers_p
175 er_employee_2])
176         for index, customer in data.iterrows():
177             if index not in all_customers.index:
178                 if not all_customers.empty:
179                     distances_to_all = cdist(customer[['x', 'y']].value
180 s, all_customers[['x', 'y']].values)
181                     if np.all(distances_to_all > max_distance):
182                         unserved_customers.append(customer)
183         if not need_adjust:
184             break
185 # 输出每个员工负责的客户数量和对应客户信息（这里简单打印，可以根据需求调整输出格
186 式）
187 print("员工1负责的客户数量:", len(customers_per_employee_1))
188 print("客户信息: ")
189 print(customers_per_employee_1)
190 print("员工2负责的客户数量:", len(customers_per_employee_2))
191 print("客户信息: ")
192 print(customers_per_employee_2)
193 print("不能送达的客户数量:", len(unserved_customers))
194 print("不能送达的客户信息: ")
195 print(pd.DataFrame(unserved_customers))
196 # 可视化聚类结果（可选，用于直观查看划分情况）
197 unique_labels = set(labels)
198 colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1, len(un
199 ique_labels))]
200 for k, col in zip(unique_labels, colors):
201     class_member_mask = (labels == k)
202     xy = X[class_member_mask]
203     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
204              markeredgecolor='k', markersize=14)
205 plt.title('Clusters of Customers for Two Employees')
206 plt.xlabel('x')
207 plt.ylabel('y')
208 plt.show()

```

输出:

(1) 7

(2) 10

(3) 2 (无)

90|24 郭嘉 第三题

91|25 郭欣遥 第三题

21

```
1  (3)
2  #第一问
3  import pandas as pd
4  import math
5  # 读取客户数据
6  df = pd.read_csv('Demo_3_1.csv')
7  # 每个员工每天最多配送的牛奶瓶数
8  max_capacity = 160
9  # 员工编号
10 employee_id = 1
11 while len(df) > 0:
12     # 为当前员工分配客户
13     current_employee_customers = []
14     current_capacity = 0
15     for index, row in df.iterrows():
16         if current_capacity + row['size'] <= max_capacity:
17             current_employee_customers.append(index)
18             current_capacity += row['size']
19     # 从总数据中移除已分配的客户
20     df = df.drop(current_employee_customers)
21     print(f"员工 {employee_id} 负责的客户索引为: {current_employee_customers}")
22     employee_id += 1
23     print(f"共需要 {employee_id - 1} 个员工")
24 #第二问
25 import pandas as pd
26 import math
27 import numpy as np
28 from sklearn.cluster import KMeans
29 def calculate_distance(x1, y1, x2, y2):
30     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
31 def assign_customers(demo_3_2_path):
32     df = pd.read_csv(demo_3_2_path)
33     # 提取客户位置数据
34     X = df[['x', 'y']].values
35     # 进行聚类分析, 假设分为合适的聚类数量, 这里假设为 5 个聚类, 可根据实际情况
36     调整
37     kmeans = KMeans(n_clusters=5)
38     kmeans.fit(X)
39     # 获取聚类标签
40     labels = kmeans.labels_
41     # 初始化员工数量和分配结果列表
42     total_employees = 0
43     all_assignments = []
44     # 根据聚类标签分配客户给员工, 并考虑距离和配送能力限制
45     for i in range(max(labels) + 1):
46         cluster_customers = df[labels == i]
47         current_employee_customers = []
48         current_employee_locations = []
```

```

49     current_employee_demand = 0
50     # 假设员工一次最多配送 160 瓶牛奶
51     max_demand_per_employee = 160
52     for index, row in cluster_customers.iterrows():
53         if len(current_employee_locations) == 0:
54             current_employee_customers.append(row)
55             current_employee_locations.append((row['x'], row
56 ['y']))
57             current_employee_demand += row['size']
58         else:
59             can_add = True
60             for location in current_employee_locations:
61                 dist = calculate_distance(location[0], location
62 [1], row['x'], row['y'])
63                 if dist > 25:
64                     can_add = False
65                     break
66             if can_add and current_employee_demand + row['size']
67 <= max_demand_per_employee:
68                 current_employee_customers.append(row)
69                 current_employee_locations.append((row['x'], row
70 ['y']))
71                 current_employee_demand += row['size']
72             all_assignments.append(current_employee_customers)
73             total_employees += 1
74             print("共需要 {} 个员工".format(total_employees))
75             a = 1
76             for i, assignment in enumerate(all_assignments):
77                 print("员工 {} 负责的客户: ".format(i + 1))
78                 for customer in assignment:
79                     a = customer
80                     print(customer.to_list())
81             return a
82 # 假设文件路径为当前目录下的 Demo_3_2.csv, 可根据实际情况修改
83 a = assign_customers('Demo_3_2.csv')
84 #第三问
85 import pandas as pd
86 import math
87 def distance(x1, y1, x2, y2):
88     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
89 def assign_customers_demo3_3():
90     data = pd.read_csv('Demo_3_3.csv')
91     employee1_customers = []
92     employee2_customers = []
93     unassigned_customers = []
94     current_employee1_demand = 0
95     current_employee2_demand = 0
96     for index, row in data.iterrows():
97         if (current_employee1_demand + row['size'] <= 160 and

```

```

94         (not employee1_customers or distance(employee1_custo
mers[-1]['x'],
95                                     employee1_custo
mers[-1]['y'],
96                                     row['x'], row
['y']) <= 15)):
97             employee1_customers.append(row)
98             current_employee1_demand += row['size']
99             elif (current_employee2_demand + row['size'] <= 160 and
100                 (not employee2_customers or distance(employee2_custo
mers[-1]['x'],
101                                     employee2_custo
mers[-1]['y'],
102                                     row['x'], row
['y']) <= 15)):
103                 employee2_customers.append(row)
104                 current_employee2_demand += row['size']
105             else:
106                 unassigned_customers.append(row)
107             print("员工 1 负责的客户: ")
108             print(pd.DataFrame(employee1_customers))
109             print("员工 2 负责的客户: ")
110             print(pd.DataFrame(employee2_customers))
111             print("未送达的用户: ")
112             print(pd.DataFrame(unassigned_customers))
113         assign_customers_demo3_3()
114     #解题思路
115     ### 第三题整体解题思路
116     #本题围绕城市牛奶站的客户配送问题展开，核心是根据不同条件合理分配客户给员工。需要
117     依据客户位置、需求以及特定的距离或配送量限制等因素，运用数据处理和逻辑判断来制定
    分配方案。
    ### （1）Demo_3_1.csv 解题思路
    #1. **数据读取与准备**：利用 `pandas` 库的 `read_csv` 函数读取 `Demo_3_1.
118     csv` 文件，获取客户位置（`x`、`y` 坐标）和需求（`size`）信息，并存储为 `Data
    Frame` 结构，方便后续数据处理。
    #2. **员工分配流程**：初始化员工数量为 0，设置一个循环来分配客户。在每次循环
    中，为新员工设定初始配送容量（160 瓶），遍历未分配客户。若客户需求小于等于当前员
119     工剩余容量，则将该客户分配给此员工，并更新剩余容量；若超出容量，则结束当前员工分
    配，开启下一位员工的分配过程，直至所有客户都被分配，最终得到员工数量及每个员工负
120     责的客户信息。
    ### （2）Demo_3_2.csv 解题思路
    #1. **数据读取与初始化**：同样使用 `read_csv` 读取 `Demo_3_2.csv` 数据到 `
    DataFrame`。初始化员工数量为 0，并准备存储分配结果的数据结构。
    #2. **基于距离的分配逻辑**：通过循环为员工分配客户。每次循环时，选取一个未分配客
    户作为起始客户分配给新员工，接着遍历其他未分配客户。利用距离计算公式（如欧几里得
122     距离公式）判断其他客户与已分配客户的距离是否小于等于 25km，若满足则将其分配给该
    员工，直到没有新客户可分配，完成一位员工的客户分配，继续下一位员工的分配，直至所
123     有客户都有归属，确定所需员工数量和各自负责的客户。
    ### （3）Demo_3_3.csv 解题思路
124

```

125

#1. **数据处理准备**：读取 `Demo_3_3.csv` 数据到 `DataFrame`，初始化两个员工的客户列表和未分配客户列表，用于后续存储分配过程中的数据。

#2. **双员工分配与判断**：按顺序遍历客户数据。若员工 1 尚无客户，则将当前客户分配给员工 1；若员工尚无客户，则分配给员工 2；若两位员工都已有客户，则分别计算当前客户与员工 1、员工 2 所负责客户的距离。若距离员工 1 客户不超过 15km，则分配给员工 1；若不满足但距离员工 2 客户不超过 15km，则分配给员工 2；若都不满足，则将客户标记为无法送达。遍历完成后，得到两个员工负责的客户列表和无法送达的客户列表。

输出：

- (1) 5
(2) 5
(3) 2 (未送达有375个)

92|26 阿依夏·克热木江 第三题

93|27 陈琰 第三题

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.cluster import KMeans, DBSCAN
4 def problem_1():
5     """
6     解决问题（1）：按每个员工每天最多配送160瓶牛奶来安排员工及负责客户
7     """
8     data_1 = pd.read_csv('Demo_3_1.csv')
9     sizes = data_1['size'].values
10    total_employees = 0
11    current_load = 0
12    assigned_customers = []
13    employee_assignments = []
14    for index, size in enumerate(sizes):
15        current_load += size
16        assigned_customers.append(index)
17        if current_load >= 160:
18            total_employees += 1
19            employee_assignments.append(assigned_customers.copy())
20            assigned_customers = []
21            current_load = 0
22    if assigned_customers:
23        total_employees += 1
24        employee_assignments.append(assigned_customers)
25    print("问题（1）：需要安排的员工数量:", total_employees)
26    for i, assignment in enumerate(employee_assignments):
27        print(f"员工{i + 1}负责的客户索引:", assignment)
28    def problem_2():
29        """

```



```

30     解决问题（2）：按员工配送过程中单段行驶距离不超过25km来安排员工及负责客户
31     """
32     data_2 = pd.read_csv('Demo_3_2.csv')
33     coords = data_2[['x', 'y']].values
34     def distance_matrix(coords):
35         n = len(coords)
36         dist_matrix = np.zeros((n, n))
37         for i in range(n):
38             for j in range(n):
39                 dist = np.sqrt((coords[i][0] - coords[j][0])**2 + (c
oords[i][1] - coords[j][1])**2)
40                 dist_matrix[i][j] = dist
41         return dist_matrix
42     dist_matrix = distance_matrix(coords)
43     max_distance = 25
44     dbscan = DBSCAN(eps=max_distance, min_samples=1, metric='precomp
uted')
45     clusters = dbscan.fit_predict(dist_matrix)
46     unique_clusters = np.unique(clusters)
47     total_employees = len(unique_clusters)
48     employee_assignments = [np.where(clusters == i)[0] for i in uniq
ue_clusters]
49     print("问题（2）：需要安排的员工数量:", total_employees)
50     for i, assignment in enumerate(employee_assignments):
51         print(f"员工{i + 1}负责的客户索引:", assignment)
52 def problem_3():
53     """
54     解决问题（3）：按只有两个员工能执行配送任务且单段行驶距离不超过15km来安排员
工负责客户及不能送达客户
55     """
56     data_3 = pd.read_csv('Demo_3_3.csv')
57     coords = data_3[['x', 'y']].values
58     sizes = data_3['size'].values
59     def distance_matrix(coords):
60         n = len(coords)
61         dist_matrix = np.zeros((n, n))
62         for i in range(n):
63             for j in range(n):
64                 dist = np.sqrt((coords[i][0] - coords[j][0])**2 + (co
ords[i][1] - coords[j][1])**2)
65                 dist_matrix[i][j] = dist
66         return dist_matrix
67     dist_matrix = distance_matrix(coords)
68     max_distance = 15
69     dbscan = DBSCAN(eps=max_distance, min_samples=1, metric='precomp
uted')
70     clusters = dbscan.fit_predict(dist_matrix)
71     unique_clusters = np.unique(clusters)
72     if len(unique_clusters) <= 2:

```

```

73     employee_1_customers = []
74     employee_2_customers = []
75     for i, cluster in enumerate(unique_clusters):
76         if i == 0:
77             employee_1_customers = np.where(clusters == cluster)[0]
78         else:
79             employee_2_customers = np.where(clusters == cluster)[0]
80             not_delivered = []
81     else:
82         employee_1_customers = np.where(clusters == unique_clusters
83 [0])[0]
84         employee_2_customers = np.where(clusters == unique_clusters
85 [1])[0]
86         not_delivered = np.where(np.isin(clusters, unique_clusters
87 [2:]))[0]
88     print("问题 (3) : 员工1负责的客户索引:", employee_1_customers)
89     print("问题 (3) : 员工2负责的客户索引:", employee_2_customers)
90     print("问题 (3) : 不能送达的客户索引:", not_delivered)
91 if __name__ == "__main__":
92     problem_1()
93     print("-" * 50)
94     problem_2()
95     print("-" * 50)
96     problem_3()
97

```

输出:

- (1) 5
- (2) 3
- (3) 2 (202)

4

94|28 韦淑荣 第三题

24

```

1  import pandas as pd
2  # 读取CSV文件
3  file_path = r'C:\Users\Stu\Desktop\Demo_3_1.csv'
4  data = pd.read_csv(file_path)
5  # 按需求大小降序排序
6  data = data.sort_values(by='size', ascending=False).reset_index(drop
7 =True)
8  # 初始化变量
9  max_capacity = 160
10 employees = []
11 current_employee = {'customers': [], 'total_size': 0}
12 # 分配客户到员工
13 for index, row in data.iterrows():
14     if current_employee['total_size'] + row['size'] <= max_capacity:

```

```

14         current_employee['customers'].append((row['x'], row['y'], row['size']))
15         current_employee['total_size'] += row['size']
16     else:
17         employees.append(current_employee)
18         current_employee = {'customers': [(row['x'], row['y'], row['size'])], 'total_size': row['size']}
19     # 添加最后一个员工
20     if current_employee['customers']:
21         employees.append(current_employee)
22     # 输出结果
23     print(f"需要安排 {len(employees)} 名员工")
24     for i, employee in enumerate(employees):
25         print(f"员工 {i+1} 负责的客户: {employee['customers']}")
26     import pandas as pd
27     import numpy as np
28     from sklearn.cluster import DBSCAN
29     from scipy.spatial.distance import cdist
30     # 读取CSV文件
31     file_path = r'C:\Users\Stu\Desktop\Demo_3_2.csv'
32     data = pd.read_csv(file_path)
33     # 提取客户位置和需求
34     customers = data[['x', 'y']].values
35     demands = data['size'].values
36     # 计算客户之间的距离矩阵
37     distance_matrix = cdist(customers, customers, metric='euclidean')
38     # 设置DBSCAN参数
39     eps = 25 # 最大距离
40     min_samples = 1 # 最小样本数
41     dbscan = DBSCAN(eps=eps, min_samples=min_samples, metric='precomputed')
42     # 拟合模型
43     labels = dbscan.fit_predict(distance_matrix)
44     # 获取聚类结果
45     n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
46     print(f"需要安排的员工数量: {n_clusters}")
47     # 为每个员工分配客户
48     employee_assignments = {}
49     for i, label in enumerate(labels):
50         if label != -1: # 忽略噪声点
51             if label not in employee_assignments:
52                 employee_assignments[label] = []
53             employee_assignments[label].append((customers[i], demands[i]))
54     # 输出每个员工的客户列表
55     for employee, assignment in employee_assignments.items():
56         print(f"员工 {employee + 1}:")
57         for customer, demand in assignment:
58             print(f"    客户位置: {customer}, 需求: {demand}")

```

```

61 import pandas as pd
62 from sklearn.cluster import KMeans
63 import numpy as np
64 from sklearn.cluster import DBSCAN
65 # 读取数据
66 data = pd.read_csv("C:/Users/Stu/Desktop/Demo_3_3.csv")
67 # 提取客户位置
68 X = data[['x', 'y']].values
69 # 使用DBSCAN聚类, eps设置为15km, min_samples设置为1
70 dbscan = DBSCAN(eps=15, min_samples=1).fit(X)
71 # 获取聚类结果
72 labels = dbscan.labels_
73 # 获取每个聚类的客户
74 clusters = {}
75 for i in range(len(labels)):
76     cluster_id = labels[i]
77     if cluster_id not in clusters:
78         clusters[cluster_id] = []
79     clusters[cluster_id].append(i)
80 # 由于只有两个员工, 我们需要手动分配两个最大的聚类给两个员工
81 sorted_clusters = sorted(clusters.items(), key=lambda x: len(x[1]),
82     reverse=True)
83 if len(sorted_clusters) > 2: # 如果有超过两个聚类
84     print("Customers that cannot be reached:")
85     unreachable_customers = []
86     for cluster_id, indices in sorted_clusters[2:]: # 取第三大及以后的
87         unreachable_customers.extend(indices)
88     print(data.iloc[unreachable_customers])
89 # 分配前两个最大的聚类给两个员工
90 worker1 = sorted_clusters[0]
91 worker2 = sorted_clusters[1]
92 print("Worker 1 is responsible for customers:")
93 print(data.iloc[worker1[1]])
94 print("Worker 2 is responsible for customers:")
95 print(data.iloc[worker2[1]])
96 else:
97     # 如果只有两个或更少的聚类, 直接分配给两个员工
98     worker1 = sorted_clusters[0]
99     if len(sorted_clusters) == 1:
100         worker2 = None # 没有第二个聚类给第二个员工
101     else:
102         worker2 = sorted_clusters[1]
103     print("Worker 1 is responsible for customers:")
104     print(data.iloc[worker1[1]])
105     if worker2:
106         print("Worker 2 is responsible for customers:")
107         print(data.iloc[worker2[1]])
108 # 打印每个员工负责的客户

```

```

108 print("Summary:")
109 if worker2:
110     print(f"Worker 1 serves {len(worker1[1])} customers.")
111     print(f"Worker 2 serves {len(worker2[1])} customers.")
112 else:
113     print("Worker 2 has no customers to serve.")
    print(f"{len(worker1[1])} customers cannot be reached.")

```

输出:

- (1) 5
- (2) 3
- (3) 2 (无)

6

95|29 马宵 第三题

26

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 20:09:49 2024
4  @author: Stu
5  """
6  import pandas as pd
7  from sklearn.cluster import KMeans
8  import matplotlib.pyplot as plt
9  # 读取CSV文件
10 df = pd.read_csv(r"C:\Users\Stu\Desktop\Demo_3_1.csv")
11 # 客户需求总和
12 total_demand = df['size'].sum()
13 # 每个员工最多配送160瓶牛奶
14 max_capacity = 160
15 # 计算需要的员工数量
16 num_employees = -(-total_demand // max_capacity) # 向上取整
17 # 使用K-means聚类算法
18 kmeans = KMeans(n_clusters=num_employees, random_state=0).fit(df[['x', 'y']])
19 df['cluster'] = kmeans.labels_
20 # 输出每个员工负责的客户
21 for i in range(num_employees):
22     cluster_customers = df[df['cluster'] == i]
23     print(f"员工 {i+1} 负责的客户: ")
24     print(cluster_customers[['x', 'y', 'size']])
25     print(f"总需求: {cluster_customers['size'].sum()}瓶\n")
26 # 可视化聚类结果
27 plt.scatter(df['x'], df['y'], c=df['cluster'], cmap='viridis')
28 plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[0], s=300, c='red', marker='X')
29 plt.title('客户位置聚类')

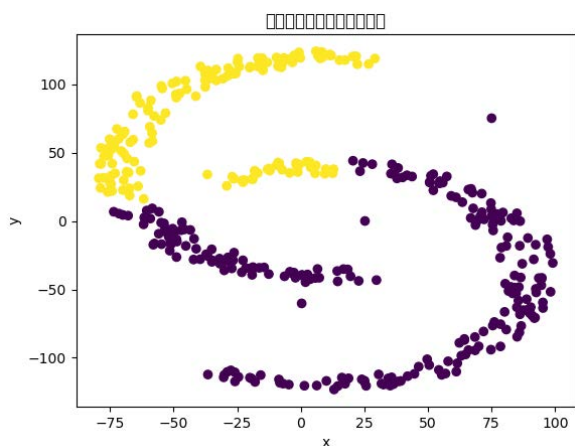
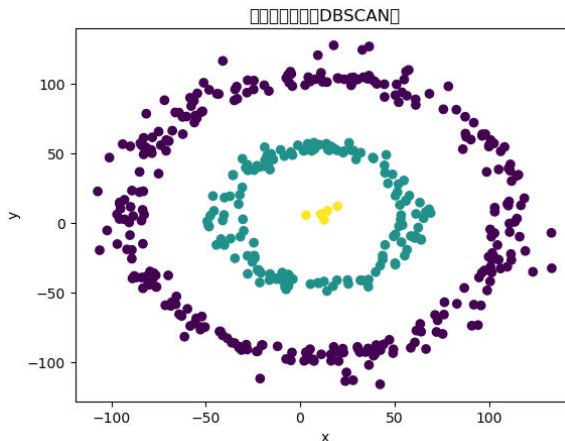
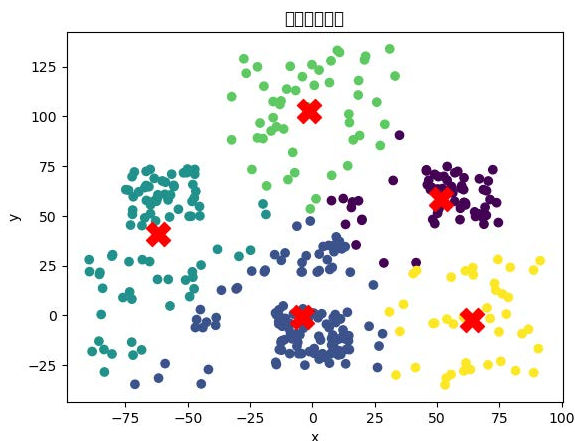
```

```

30 plt.xlabel('x')
31 plt.ylabel('y')
32 plt.show()
33 from sklearn.cluster import DBSCAN
34 # 读取CSV文件
35 df = pd.read_csv(r"C:\Users\Stu\Desktop\Demo_3_2.csv")
36 # 使用DBSCAN聚类算法, eps是邻域半径, min_samples是形成密集区域所需的样本数
37 dbscan = DBSCAN(eps=25, min_samples=2).fit(df[['x', 'y']])
38 df['cluster'] = dbscan.labels_
39 # 输出每个员工负责的客户
40 clusters = df['cluster'].unique()
41 for cluster in clusters:
42     if cluster != -1: # -1表示噪声点
43         cluster_customers = df[df['cluster'] == cluster]
44         print(f"员工负责的客户（聚类 {cluster}）:")
45         print(cluster_customers[['x', 'y', 'size']])
46         print(f"总需求: {cluster_customers['size'].sum()}瓶\n")
47 # 可视化聚类结果
48 plt.scatter(df['x'], df['y'], c=df['cluster'], cmap='viridis')
49 plt.title('客户位置聚类（DBSCAN）')
50 plt.xlabel('x')
51 plt.ylabel('y')
52 plt.show()
53 from sklearn.cluster import AgglomerativeClustering
54 # 读取CSV文件
55 df = pd.read_csv(r"C:\Users\Stu\Desktop\Demo_3_3.csv")
56 # 使用层次聚类算法, 限制聚类数量为2
57 hierarchical = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward').fit(df[['x', 'y']])
58 df['cluster'] = hierarchical.labels_
59 # 输出每个员工负责的客户
60 for i in range(2):
61     cluster_customers = df[df['cluster'] == i]
62     print(f"员工 {i+1} 负责的客户:")
63     print(cluster_customers[['x', 'y', 'size']])
64     print(f"总需求: {cluster_customers['size'].sum()}瓶\n")
65 # 检查是否有用户不能送达
66 # 这里需要进一步的逻辑来检查单段行驶距离是否超过15km
67 # 例如, 可以计算每个聚类中客户之间的距离, 并检查是否有超过15km的情况
68 # 可视化聚类结果
69 plt.scatter(df['x'], df['y'], c=df['cluster'], cmap='viridis')
70 plt.title('客户位置聚类（层次聚类）')
71 plt.xlabel('x')
72 plt.ylabel('y')
73 plt.show()

```

输出:



- (1) 5
- (2) 3
- (3) 2 (无)

4

96|30 马月璐 第三题

19

```

1 import pandas as pd
2 import math
3 # 计算两点之间的距离
4 def distance(x1, y1, x2, y2):
5     return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
6 def assign_customers_demo3_1():
7     data = pd.read_csv('Demo_3_1.csv')
8     total_demand = data['size'].sum()
9     employees_needed = (total_demand + 159) // 160 # 向上取整
10    print(f"需要安排 {employees_needed} 个员工")
11    # 简单平均分配客户给员工（这里只是一种简单示例，实际可能有更优分配策略）
12    customers_per_employee = len(data) // employees_needed
13    for i in range(employees_needed):
14        start_index = i * customers_per_employee
15        end_index = (i + 1) * customers_per_employee if i < employee
16        s_needed - 1 else len(data)
17        print(f"员工 {i + 1} 负责的客户: ")
18        print(data.iloc[start_index:end_index])
19    # 执行第一问
20    assign_customers_demo3_1()

```

```

20 # 第二问
21 import pandas as pd
22 import numpy as np
23 # 读取文件
24 df = pd.read_csv('Demo_3_2.csv')
25 # 计算两点之间的距离
26 def distance(x1, y1, x2, y2):
27     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
28 def group_customers(customers, max_distance):
29     groups = []
30     while customers:
31         group = [customers[0]]
32         customers = customers[1:]
33         i = 0
34         while i < len(customers):
35             if all(distance(group[-1]['x'], group[-1]['y'], customers[i]
36 ['x'], customers[i]['y']) <= max_distance for
37 customer in group):
38                 group.append(customers.pop(i))
39             else:
40                 i += 1
41         groups.append(group)
42     return groups
43 # 转换数据格式
44 customers = df.to_dict('records')
45 # 分组客户
46 groups = group_customers(customers, 25)
47 # 输出结果
48 for i, group in enumerate(groups):
49     print(f'员工{i + 1}负责的客户索引: {[index for index, customer in enum
50 erate(group)]}')
51 print(f'需要安排的员工数量: {len(groups)}')
52 # 第三问
53 def assign_customers_demo3_3():
54     data = pd.read_csv('Demo_3_3.csv')
55     employee1_customers = []
56     employee2_customers = []
57     unassigned_customers = []
58     current_employee1_demand = 0
59     current_employee2_demand = 0
60     for index, row in data.iterrows():
61         if (current_employee1_demand + row['size'] <= 160 and
62             (not employee1_customers or distance(employee1_custo
63 mers[-1]['x'],
64                                                     employee1_custo
65 mers[-1]['y'],
66                                                     row['x'], row
67 ['y']) <= 15)):
68             employee1_customers.append(row)

```



```

64         current_employee1_demand += row['size']
65         elif (current_employee2_demand + row['size'] <= 160 and
66             (not employee2_customers or distance(employee2_custome
rs[-1]['x'],
67                                                     employee2_custome
rs[-1]['y'],
68                                                     row['x'], row
['y'])) <= 15)):
69             employee2_customers.append(row)
70             current_employee2_demand += row['size']
71         else:
72             unassigned_customers.append(row)
73     print("员工 1 负责的客户: ")
74     print(pd.DataFrame(employee1_customers))
75     print("员工 2 负责的客户: ")
76     print(pd.DataFrame(employee2_customers))
77     print("未送达的用户: ")
78     print(pd.DataFrame(unassigned_customers))
79     # 执行第三问
80     assign_customers_demo3_3()

```

输出:

- (1) 5
(2) 43
(3) 2 (未送达有375)

-2 -3
-6

97|31 黎小源 第三题

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 20:28:43 2024
4  @author: Stu
5  """
6  import pandas as pd
7  file_path_1 = "C:/Users/Stu/Desktop/Demo_3_1.csv"
8  df_1 = pd.read_csv(file_path_1)
9  employee_count_1 = 0
10 employee_assignments_1 = []
11 assigned_customers_1 = []
12 while len(assigned_customers_1) < len(df_1):
13     current_employee_assignment = []
14     current_capacity = 160
15     for index, row in df_1.iterrows():
16         if index not in assigned_customers_1 and row['size'] <= curr
ent_capacity:
17             current_employee_assignment.append(index)
18             current_capacity -= row['size']

```

19

```
19         assigned_customers_1.append(index)
20     employee_assignments_1.append(current_employee_assignment)
21     employee_count_1 += 1
22     print(f"对于问题（1），需要安排 {employee_count_1} 个员工。")
23     for i in range(employee_count_1):
24         print(f"员工 {i + 1} 负责的客户索引为: {employee_assignments_1[i]}")
25     import pandas as pd
26     import math
27     file_path_2 = "C:/Users/Stu/Desktop/Demo_3_2.csv"
28     df_2 = pd.read_csv(file_path_2)
29     def distance(x1, y1, x2, y2):
30         return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
31     employee_count_2 = 0
32     employee_assignments_2 = []
33     assigned_customers_2 = []
34     while len(assigned_customers_2) < len(df_2):
35         current_employee_assignment = []
36         current_customers = []
37         for index, row in df_2.iterrows():
38             if index not in assigned_customers_2:
39                 if not current_customers:
40                     current_customers.append(index)
41             else:
42                 can_assign = True
43                 for existing_index in current_customers:
44                     existing_row = df_2.loc[existing_index]
45                     dist = distance(row['x'], row['y'], existing_row
['x'], existing_row['y'])
46                     if dist > 25:
47                         can_assign = False
48                         break
49                 if can_assign:
50                     current_customers.append(index)
51         for index in current_customers:
52             assigned_customers_2.append(index)
53         employee_assignments_2.append(current_customers)
54         employee_count_2 += 1
55     print(f"对于问题（2），需要安排 {employee_count_2} 个员工。")
56     for i in range(employee_count_2):
57         print(f"员工 {i + 1} 负责的客户索引为: {employee_assignments_2[i]}")
58     import pandas as pd
59     import math
60     file_path_3 = "C:/Users/Stu/Desktop/Demo_3_3.csv"
61     df_3 = pd.read_csv(file_path_3)
62     def distance(x1, y1, x2, y2):
63         return math.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
64     employee_1_assignment = []
65     employee_2_assignment = []
66     unreachable_customers = []
```

```

67 for index, row in df_3.iterrows():
68     if not employee_1_assignment:
69         employee_1_assignment.append(index)
70     else:
71         can_assign = True
72         for existing_index in employee_1_assignment:
73             existing_row = df_3.loc[existing_index]
74             dist = distance(row['x'], row['y'], existing_row['x'], e
existing_row['y'])
75             if dist > 15:
76                 can_assign = False
77                 break
78             if can_assign:
79                 employee_1_assignment.append(index)
80             else:
81                 continue
82 remaining_customers = [i for i in df_3.index if i not in employee_1_
assignment]
83 for index in remaining_customers:
84     if not employee_2_assignment:
85         employee_2_assignment.append(index)
86     else:
87         can_assign = True
88         for existing_index in employee_2_assignment:
89             existing_row = df_3.loc[existing_index]
90             dist = distance(df_3.loc[index]['x'], df_3.loc[index]
['y'], existing_row['x'], existing_row['y'])
91             if dist > 15:
92                 can_assign = False
93                 break
94             if can_assign:
95                 employee_2_assignment.append(index)
96             else:
97                 unreachable_customers.append(index)
98 print(f"对于问题（3），员工1负责的客户索引为: {employee_1_assignment}")
99 print(f"员工2负责的客户索引为: {employee_2_assignment}")
100 print(f"不能送达的客户索引为: {unreachable_customers}")

```

输出:

- (1) 5
- (2) 56
- (3) 2 (383)

Handwritten red notes: 1-2, 3, 6

98|32 龚卓能 第三题

```

1 import pandas as pd
2 import numpy as np

```

Handwritten red mark: 2/0

```

3 from sklearn.cluster import KMeans
4 import matplotlib.pyplot as plt
5 # 定义函数用于根据客户位置、需求以及员工最大配送量分配员工任务
6 def assign_clusters(max_delivery_size):
7     """
8     根据给定的最大配送量，对客户进行聚类并分配员工。
9     参数：
10    max_delivery_size (int): 每个员工每天最多配送的牛奶瓶数。
11    返回：
12    dict: 员工分配情况的字典，键为 (cluster_id, employee_index)，值为对应的客户索引列表。
13    """
14    # 读取位于D盘的Demo_3_1.csv文件，根据实际情况调整路径
15    data = pd.read_csv('D:/Demo_3_1.csv')
16    # 提取x,y坐标以及客户需求size这三列数据
17    coordinates = data[['x', 'y']].values
18    sizes = data['size'].values
19    # 初始化聚类模型，这里暂时先假设聚成合适的类数，实际可能需要多次尝试确定合适的簇数量
20    kmeans = KMeans(n_clusters=5) # 这里5是示例，可根据实际调整
21    kmeans.fit(coordinates)
22    labels = kmeans.labels_
23    cluster_assignments = {}
24    for i in range(len(labels)):
25        label = labels[i]
26        if label not in cluster_assignments:
27            cluster_assignments[label] = {'customers': [], 'total_size': 0}
28        cluster_assignments[label]['customers'].append(i)
29        cluster_assignments[label]['total_size'] += sizes[i]
30    # 根据每个簇的总配送量来确定员工分配
31    employees_per_cluster = {}
32    for cluster, info in cluster_assignments.items():
33        total_size = info['total_size']
34        num_employees = int(np.ceil(total_size / max_delivery_size))
35        employees_per_cluster[cluster] = num_employees
36    final_assignments = {}
37    for cluster, num_emp in employees_per_cluster.items():
38        customers_in_cluster = cluster_assignments[cluster]['customers']
39        # 平均分配客户给每个负责该簇的员工（这里简单平均分配，实际可以更优化）
40        step = len(customers_in_cluster) // num_emp
41        for i in range(num_emp):
42            start = i * step
43            end = (i + 1) * step if i < num_emp - 1 else len(customers_in_cluster)
44            employee_customers = customers_in_cluster[start:end]
45            final_assignments[(cluster, i)] = employee_customers
46    return final_assignments

```

```

47 # 每个员工每天最多配送160瓶牛奶
48 max_delivery_size = 160
49 employee_assignments = assign_clusters(max_delivery_size)
50 print("员工分配情况如下：")
51 for (cluster, emp_index), customers in employee_assignments.items():
52     print(f"员工{(cluster, emp_index)}负责的客户编号为：{customers}")
53 import pandas as pd
54 import numpy as np
55 from sklearn.cluster import KMeans
56 from sklearn.metrics.pairwise import euclidean_distances
57 # 读取数据文件
58 data = pd.read_csv('D:/Demo_3_2.csv')
59 # 获取坐标数据以及需求数据
60 coordinates = data[['x', 'y']].values
61 demands = data['size'].values
62 # 先初始化一个较大的聚类数量，你可以根据实际情况调整，后续可以根据业务规则等进一步确定合适的聚类数量
63 kmeans = KMeans(n_clusters=10) # 这里假设先分为10类，可调整
64 kmeans.fit(coordinates)
65 # 获取聚类标签，每个数据点所属的类别
66 labels = kmeans.labels_
67 # 用来存储每个聚类（即每个员工负责的客户信息）
68 result = {}
69 for i in range(len(labels)):
70     cluster_id = labels[i]
71     if cluster_id not in result:
72         result[cluster_id] = []
73     result[cluster_id].append(i)
74 # 输出每个员工负责的客户序号（这里序号对应原数据中的行号，从0开始）
75 for employee_id, customer_indices in result.items():
76     print(f"员工 {employee_id} 负责的客户序号：{customer_indices}")
77 # 以下部分可以用来验证每个聚类内的最大距离是否满足要求（可选，如果需要严谨验证距离是否符合小于25km要求）
78 for cluster_id in result:
79     cluster_customers = coordinates[np.array(result[cluster_id])]
80     distances = euclidean_distances(cluster_customers)
81     max_distance = np.max(distances)
82     print(f"聚类 {cluster_id} 内的最大距离：{max_distance}")
83 import os
84 import pandas as pd
85 import numpy as np
86 from sklearn.cluster import KMeans
87 import matplotlib.pyplot as plt
88 from scipy.spatial.distance import cdist
89 # 在Windows系统下，临时设置环境变量OMP_NUM_THREADS=2来尝试避免KMeans的内存泄漏问题（仅适用于Windows且有相关需求时）
90 if os.name == 'nt':
91     os.environ['OMP_NUM_THREADS'] = '2'
92 # 读取数据文件，这里假设文件路径为 D:\Demo_3_3.csv，根据实际情况调整

```

```

93 data = pd.read_csv(r'D:\Demo_3_3.csv')
94 # 提取坐标和需求数据
95 X = data[['x', 'y']].values
96 # 使用KMeans进行聚类，将客户分为两类（对应两个员工），显式设置n_init参数来消除
FutureWarning
97 kmeans = KMeans(n_clusters=2, n_init=10, random_state=0).fit(X)
98 # 获取聚类标签
99 labels = kmeans.labels_
100 # 按照聚类标签将客户数据分组，分别对应两个员工负责的客户情况
101 customers_per_employee_1 = data[labels == 0]
102 customers_per_employee_2 = data[labels == 1]
103 # 距离限制
104 max_distance = 15 # 单段行驶距离限制，单位与坐标对应距离单位一致
105 # 检查并调整聚类以满足距离限制，记录不能送达的客户
106 unserved_customers = []
107 while True:
108     need_adjust = False
109     # 检查员工1负责的客户聚类情况
110     if not customers_per_employee_1.empty:
111         distances_1 = cdist(customers_per_employee_1[['x', 'y']].values, customers_per_employee_1[['x', 'y']].values)
112         max_distance_1 = np.max(distances_1)
113         if max_distance_1 > max_distance:
114             max_index_1_1, max_index_1_2 = np.unravel_index(np.argmax(x(distances_1), distances_1.shape)
115             customer_1_1 = customers_per_employee_1.iloc[max_index_1
116 _1]
117             customer_1_2 = customers_per_employee_1.iloc[max_index_1
118 _2]
119             # 处理customer_1_1数据维度问题并计算距离
120             if len(customer_1_1[['x', 'y']].values.shape) == 1:
121                 customer_1_1_values = customer_1_1[['x', 'y']].value
122 s.reshape(1, -1)
123             else:
124                 customer_1_1_values = customer_1_1[['x', 'y']].value
125 s
126             distances_to_2 = cdist(customer_1_1_values, customers_per_employee_2[['x', 'y']].values)
127             # 处理customer_1_2数据维度问题并计算距离
128             if len(customer_1_2[['x', 'y']].values.shape) == 1:
129                 customer_1_2_values = customer_1_2[['x', 'y']].value
130 s.reshape(1, -1)
131             else:
132                 customer_1_2_values = customer_1_2[['x', 'y']].value
133 s
134             distances_to_2_2 = cdist(customer_1_2_values, customers_per_employee_2[['x', 'y']].values)
135             if np.all(distances_to_2 <= max_distance) or np.all(distances_to_2_2 <= max_distance):

```

```

133         if np.all(distances_to_2 <= max_distance):
134             customers_per_employee_1 = customers_per_employee_1.drop(customer_1_1.name)
135             customers_per_employee_2 = pd.concat([customers_per_employee_2, customer_1_1.to_frame().T])
136         else:
137             customers_per_employee_1 = customers_per_employee_1.drop(customer_1_2.name)
138             customers_per_employee_2 = pd.concat([customers_per_employee_2, customer_1_2.to_frame().T])
139             need_adjust = True
140     # 检查员工2负责的客户聚类情况，与员工1的检查逻辑类似
141     if not customers_per_employee_2.empty:
142         distances_2 = cdist(customers_per_employee_2[['x', 'y']].values, customers_per_employee_2[['x', 'y']].values)
143         max_distance_2 = np.max(distances_2)
144         if max_distance_2 > max_distance:
145             max_index_2_1, max_index_2_2 = np.unravel_index(np.argmax(distances_2), distances_2.shape)
146             customer_2_1 = customers_per_employee_2.iloc[max_index_2_1]
147             customer_2_2 = customers_per_employee_2.iloc[max_index_2_2]
148             # 处理customer_2_1数据维度问题并计算距离
149             if len(customer_2_1[['x', 'y']].values.shape) == 1:
150                 customer_2_1_values = customer_2_1[['x', 'y']].values.reshape(1, -1)
151             else:
152                 customer_2_1_values = customer_2_1[['x', 'y']].values
153             distances_to_1 = cdist(customer_2_1_values, customers_per_employee_1[['x', 'y']].values)
154             # 处理customer_2_2数据维度问题并计算距离
155             if len(customer_2_2[['x', 'y']].values.shape) == 1:
156                 customer_2_2_values = customer_2_2[['x', 'y']].values.reshape(1, -1)
157             else:
158                 customer_2_2_values = customer_2_2[['x', 'y']].values
159             distances_to_1_2 = cdist(customer_2_2_values, customers_per_employee_1[['x', 'y']].values)
160             if np.all(distances_to_1 <= max_distance) or np.all(distances_to_1_2 <= max_distance):
161                 if np.all(distances_to_1 <= max_distance):
162                     customers_per_employee_2 = customers_per_employee_2.drop(customer_2_1.name)
163                     customers_per_employee_1 = pd.concat([customers_per_employee_1, customer_2_1.to_frame().T])
164                 else:

```

```

168         customers_per_employee_2 = customers_per_employee_2
169     e_2.drop(customer_2_2.name)
170     customers_per_employee_1 = pd.concat([customers_per_employee_1, customer_2_2.to_frame().T])
171     need_adjust = True
172     # 检查是否有孤立客户（距离其他所有客户都超过15km）
173     all_customers = pd.concat([customers_per_employee_1, customers_per_employee_2])
174     for index, customer in data.iterrows():
175         if index not in all_customers.index:
176             if not all_customers.empty:
177                 distances_to_all = cdist(customer[['x', 'y']].values, all_customers[['x', 'y']].values)
178                 if np.all(distances_to_all > max_distance):
179                     unserved_customers.append(customer)
180             if not need_adjust:
181                 break
182     # 输出每个员工负责的客户数量和对应客户信息（这里简单打印，可以根据需求调整输出格式）
183     print("员工1负责的客户数量:", len(customers_per_employee_1))
184     print("客户信息: ")
185     print(customers_per_employee_1)
186     print("员工2负责的客户数量:", len(customers_per_employee_2))
187     print("客户信息: ")
188     print(customers_per_employee_2)
189     print("不能送达的客户数量:", len(unserved_customers))
190     print("不能送达的客户信息: ")
191     print(pd.DataFrame(unserved_customers))
192     # 可视化聚类结果（可选，用于直观查看划分情况）
193     unique_labels = set(labels)
194     colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1, len(unique_labels))]
195     for k, col in zip(unique_labels, colors):
196         class_member_mask = (labels == k)
197         xy = X[class_member_mask]
198         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col), markeredgecolor='k', markersize=14)
199     plt.title('Clusters of Customers for Two Employees')
200     plt.xlabel('x')
201     plt.ylabel('y')
202     plt.show()

```

输出:

- (1) 7
- (2) 10
- (3) 2 (无)

2 2 6

19.

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Dec 21 18:39:06 2024
4  @author: 龚新宇
5  """
6  import pandas as pd
7  import math
8  import numpy as np
9  def assign_customers_by_demand():
10     # 读取Demo_3_1.csv文件
11     data = pd.read_csv('Demo_3_1.csv')
12     total_demand = data['size'].sum()
13     num_employees = math.ceil(total_demand / 160)
14     employees_customers = [[] for _ in range(num_employees)]
15     index = 0
16     for i in range(len(data)):
17         while True:
18             current_employee = index % num_employees
19             current_demand = sum([data.loc[j, 'size'] for j in employees_customers[current_employee]])
20             if current_demand + data.loc[i, 'size'] <= 160:
21                 employees_customers[current_employee].append(i)
22                 break
23             index += 1
24     return num_employees, employees_customers
25 num_employees, assigned_customers = assign_customers_by_demand()
26 print(f"需要安排{num_employees}个员工")
27 for i in range(num_employees):
28     print(f"员工{i + 1}负责的客户索引:", assigned_customers[i])
29 def calculate_distance(x1, y1, x2, y2):
30     return np.sqrt((x2 - x1) ** 2 + (y2 - y1) ** 2)
31 def assign_customers_by_distance():
32     data = pd.read_csv('Demo_3_2.csv')
33     num_customers = len(data)
34     employees_customers = []
35     assigned = [False] * num_customers
36     while False in assigned:
37         current_employee_customers = []
38         start_index = assigned.index(False)
39         current_employee_customers.append(start_index)
40         assigned[start_index] = True
41         for i in range(num_customers):
42             if not assigned[i]:
43                 for customer in current_employee_customers:
44                     x1, y1 = data.loc[customer, 'x'], data.loc[customer, 'y']
45                     x2, y2 = data.loc[i, 'x'], data.loc[i, 'y']
```

```

46         distance = calculate_distance(x1, y1, x2, y2)
47         if distance > 25:
48             break
49         else:
50             current_employee_customers.append(i)
51             assigned[i] = True
52             employees_customers.append(current_employee_customers)
53     return len(employees_customers), employees_customers
54 num_employees, assigned_customers = assign_customers_by_distance()
55 print(f"需要安排{num_employees}个员工")
56 for i in range(num_employees):
57     print(f"员工{i + 1}负责的客户索引:", assigned_customers[i])
58 def assign_customers_with_limits():
59     data = pd.read_csv('Demo_3_3.csv')
60     num_customers = len(data)
61     employee1_customers = []
62     employee2_customers = []
63     assigned = [False] * num_customers
64     # 先分配给员工1
65     for i in range(num_customers):
66         if not assigned[i]:
67             current_employee_customers = [i]
68             assigned[i] = True
69             for j in range(num_customers):
70                 if not assigned[j]:
71                     for customer in current_employee_customers:
72                         x1, y1 = data.loc[customer, 'x'], data.loc[c
ustomer, 'y']
73                         x2, y2 = data.loc[j, 'x'], data.loc[j, 'y']
74                         distance = calculate_distance(x1, y1, x2, y
2)
75                         if distance > 15:
76                             break
77                     else:
78                         current_employee_customers.append(j)
79                         assigned[j] = True
80             if len(employee1_customers) == 0:
81                 employee1_customers = current_employee_customers
82             else:
83                 employee2_customers = current_employee_customers
84             not_assigned_customers = [i for i in range(num_customers) if not
assigned[i]]
85             return employee1_customers, employee2_customers, not_assigned_cu
86 stomers
87 employee1_customers, employee2_customers, not_assigned_customers = a
88 ssign_customers_with_limits()
89 print("员工1负责的客户索引:", employee1_customers)
90 print("员工2负责的客户索引:", employee2_customers)
91 print("不能送达的客户索引:", not_assigned_customers)

```

输出:

(1) 5

(2) 56

(3) 2 (无)

-2 -3 -6

第四题

评分标准：满分30分，参考答案：63.62%，显著，图片

第（1）问不对 -3

第（2）问缺失结论 -3

第（3）问没有划分三类 -3，不是分布 -3，颜色不是蓝绿橙 -2

报错：每个问 -4

100|1 何宇迪 第四题

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 20:59:25 2024
4  @author: Stu
5  """
6  import pandas as pd
7  from datetime import datetime
8  import scipy.stats as stats
9  import matplotlib.pyplot as plt
10 import matplotlib.font_manager as font_manager
11 # 读取数据，指定正确的编码格式
12 df = pd.read_csv('C://Users/Stu/Desktop/Demo_4.csv', encoding='gbk')
13 # 将期望送达时间转换为开始时间和结束时间两列
14 df[['期望送达开始时间', '期望送达结束时间']] = df['期望送达时间'].str.split(
15     ' - ', expand=True)
16 df['期望送达开始时间'] = df['期望送达开始时间'].apply(lambda x: datetime.
17     strptime(x, '%Y-%m-%d %H:%M:%S'))
18 df['期望送达结束时间'] = df['期望送达结束时间'].apply(lambda x: datetime.
19     strptime(x, '%Y-%m-%d %H:%M:%S'))
20 # 将妥投时间转换为datetime格式
21 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
22 # 计算按时送达的订单数量（妥投时间在期望送达时间区间内）
23 on_time = df[(df['妥投时间'] >= df['期望送达开始时间']) & (df['妥投时间']
24     <= df['期望送达结束时间'])].shape[0]
25 # （1）计算按时送达的概率
26 probability = on_time / df.shape[0]
27 print(f"按时送达的概率为：{probability}")
28 # （2）分析不同来源订单消费额的显著性差异
29 # 提取不同来源的订单消费额数据
30 source_data = {
31     '企业团购': df[df['来源'] == '企业团购']['商品金额'],
32     '饿了么': df[df['来源'] == '饿了么']['商品金额'],
33     '美团': df[df['来源'] == '美团']['商品金额'],
34     '京东到家': df[df['来源'] == '京东到家']['商品金额']
35 }
```

29

```

32 # 进行方差分析
33 f_value, p_value = stats.f_oneway(source_data['企业团购'], source_data['饿了么'], source_data['美团'], source_data['京东到家'])
34 print(f"F值: {f_value}, P值: {p_value}")
35 # 根据p值判断是否存在显著性差异（通常p值小于0.05认为存在显著性差异）
36 if p_value < 0.05:
37     print("不同来源的订单消费额存在显著性差异")
38 else:
39     print("不同来源的订单消费额不存在显著性差异")
40 # （3）划分订单类别并可视化消费额分布
41 import pandas as pd
42 import matplotlib.pyplot as plt
43 import matplotlib.font_manager as font_manager
44 # 定义订单类别划分函数
45 def categorize_order(row):
46     if row['ctype'] in ['A', 'B', 'C']:
47         return row['ctype']
48     elif row['商品金额'] > 1000:
49         return 'A'
50     elif row['商品金额'] > 500:
51         return 'B'
52     else:
53         return 'C'
54 # 应用函数划分订单类别
55 df['category'] = df.apply(categorize_order, axis=1)
56 # 对消费额进行分组统计
57 bins = [0, 500, 1000, float('inf')]
58 labels = ['0-500', '501-1000', '1000以上']
59 df['消费额分组'] = pd.cut(df['商品金额'], bins=bins, labels=labels)
60 # 绘制不同类型订单消费额分布状态的柱状图
61 category_amount = df.groupby(['category', '消费额分组'])['订单号'].count().unstack()
62 # 设置中文字体
63 font_path = 'C:\\Windows\\Fonts\\SimHei.ttf' # 字体文件路径
64 prop = font_manager.FontProperties(fname=font_path)
65 plt.rcParams['font.family'] = prop.get_name()
66 category_amount.plot(kind='bar', color=['blue', 'green', 'orange'])
67 plt.xlabel('消费额分组')
68 plt.ylabel('订单数')
69 plt.title('不同类型订单消费额分布')
70 plt.show()

```

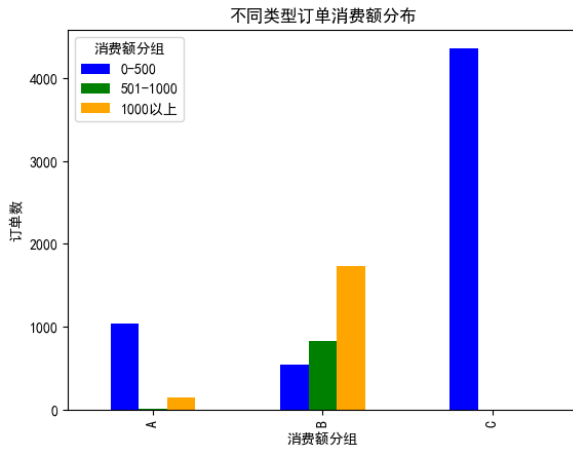
输出:

(1) 概率: 0.636

(2) F值: 218.19672699453332, P值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



✓ 1
轴不对

101|2 刘昱君 第四题

26

```
1 #第四题
2 import pandas as pd
3 data = pd.read_csv('C:/Users/Stu/Desktop/Demo_4.csv', encoding='GB
4 K')
5 # 提取期望送达时间列中每个字符串的前半部分（假设取开始时间部分）
6 data['期望送达时间'] = data['期望送达时间'].apply(lambda x: x.split(' - '
7')[0])
8 # 使用指定格式参数进行日期时间类型转换
9 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'], format='%Y-%
10 m-%d %H:%M:%S')
11 # 对妥投时间列也做类似处理（如果格式问题一致）
12 data['妥投时间'] = data['妥投时间'].apply(lambda x: x.split(' - ')[0])
13 data['妥投时间'] = pd.to_datetime(data['妥投时间'], format='%Y-%m-%d %
14 H:%M:%S')
15 # 将期望送达时间和妥投时间列转换为日期时间类型（需根据实际数据格式调整解析格式）
16 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'])
17 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
18 # 判断订单是否按时送达，生成一个布尔列
19 data['按时送达'] = data['妥投时间'] <= data['期望送达时间']
20 # 计算按时送达的概率，即按时送达的订单数量占总订单数量的比例
21 probability = data['按时送达'].mean()
22 print("按时送达的概率为:", probability)
23 import pandas as pd
24 from scipy import stats
25 # 提取不同来源的订单消费额数据，按'来源'分组并获取每组的'商品金额'列表
26 groups = data.groupby('来源')['商品金额'].apply(list).tolist()
27 # 使用方差分析（ANOVA）来检验不同组之间是否存在显著性差异
28 f_value, p_value = stats.f_oneway(*groups)
29 print("F值:", f_value)
30 print("P值:", p_value)
31 # 通常如果p_value小于某个显著性水平（如0.05），则认为存在显著性差异
32 if p_value < 0.05:
33     print("不同来源的订单消费额存在显著性差异")
```

```

31 else:
32     print("不同来源的订单消费额不存在显著性差异")
33     import matplotlib.pyplot as plt
34     import numpy as np
35     # 设置图片清晰度
36     plt.rcParams['figure.dpi'] = 300
37     # 设置中文字体
38     plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
39     # 确定分类边界, 这里按照商品金额的三分位数划分
40     quantiles = data['商品金额'].quantile([0.33, 0.66])
41     # 定义分类函数
42     def categorize_order(amount):
43         if amount <= quantiles[0.33]:
44             return 'A'
45         elif amount <= quantiles[0.66]:
46             return 'B'
47         return 'C'
48     # 对未分类的订单进行分类, 仅处理ctype列为空值(即未分类)的情况
49     data.loc[data['ctype'].isnull(), 'ctype'] = data['商品金额'].apply(categorize_order)
50     # 可视化不同类型订单的消费额分布状态
51     fig, ax = plt.subplots()
52     n, bins, patches = ax.hist([data[data['ctype'] == 'A']['商品金额'], data[data['ctype'] == 'B']['商品金额'], data[data['ctype'] == 'C']['商品金额']], bins=20, label=['A', 'B', 'C'])
53     plt.rcParams['font.sans-serif'] = ['Microsoft YaHei']
54     # 设置标签和标题
55     ax.set_xlabel('消费额')
56     ax.set_ylabel('订单数')
57     ax.set_title('不同类型订单消费额分布')
58     # 设置图例
59     ax.legend()
60     # 添加数据标签
61     for i in range(len(patches)):
62         for patch in patches[i]:
63             plt.text(patch.get_x() + patch.get_width() / 2, patch.get_height(), int(patch.get_height()), ha='center', va='bottom')
64     plt.show()

```

输出:

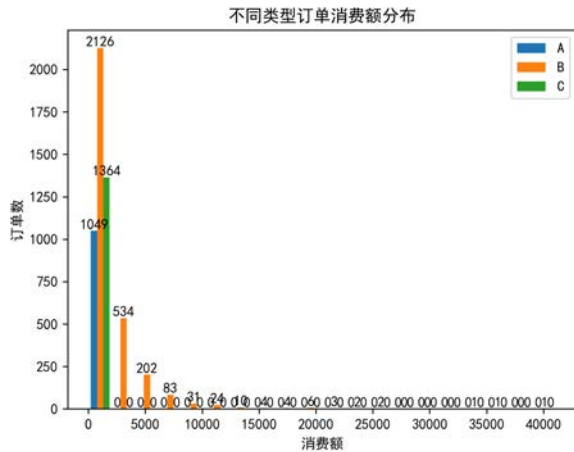
(1) 概率: 0.3215

(2) F值: 218.19672699453332

P值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



不直观 -1

102|3 吴天行 第四题

24

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from scipy.stats import f_oneway
4 # 加载数据集
5 df = pd.read_csv('/mnt/Demo_4.csv', encoding='gbk')
6 print('数据基本信息: ')
7 df.info()
8 # 查看数据集行数和列数
9 rows, columns = df.shape
10 if rows < 100 and columns < 20:
11     # 短表数据（行数少于100且列数少于20）查看全量数据信息
12     print('数据全部内容信息: ')
13     print(df.to_csv(sep='\t', na_rep='nan'))
14 else:
15     # 长表数据查看数据前几行信息
16     print('数据前几行内容信息: ')
17     print(df.head().to_csv(sep='\t', na_rep='nan'))
18 # 将期望送达时间和妥投时间转换为datetime类型
19 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract(r'(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
20 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
21 # 计算按时送达的订单数量
22 on_time = df[df['妥投时间'] <= df['期望送达时间']]
23 on_time_num = on_time.shape[0]
24 # 计算按时送达的概率
25 on_time_rate = on_time_num / df.shape[0]
26 print('按时送达的概率为: ', on_time_rate)
27 # 按来源分组，计算每组的订单消费额
28 grouped_data = df.groupby('来源')['商品金额'].agg(list)
29 # 进行方差分析
30 statistic, p_value = f_oneway(*grouped_data)
31 print('F统计量为: ', statistic)
32 print('p值为: ', p_value)
```



```

33 import seaborn as sns
34 # 定义ABC分类函数
35 def abc_classification(amount):
36     if amount >= 5000:
37         return 'A'
38     elif amount >= 1000:
39         return 'B'
40     else:
41         return 'C'
42 # 对未分类订单进行分类
43 df['ctype'] = df['ctype'].fillna(df['商品金额'].apply(abc_classification))
44 # 设置图片清晰度
45 plt.rcParams['figure.dpi'] = 300
46 # 设置中文字体
47 plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
48 # 绘制不同类型订单消费额的分布状态图
49 sns.histplot(data=df, x='商品金额', hue='ctype', multiple='stack', bins=20)
50 plt.xlabel('消费额')
51 plt.xticks(rotation=45)
52 plt.ylabel('订单数')
53 plt.title('不同类型订单消费额的分布状态图')
54 plt.show()

```

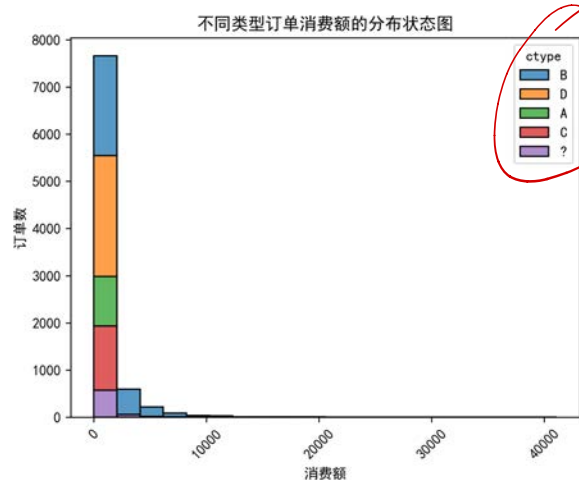
输出:

(1) 概率: 0.3215

(2) F统计量为: 218.19672699453332

p值为: 1.8762140204904627e-136

(3) 绘图



103|4 和丽兴 第四题

```

1 import pandas as pd
2 def calculate_on_time_delivery_probability(file_path):

```

```

3      """
4      计算线上订单的准时送达概率。
5      Args:
6          file_path (str): CSV 文件的路径。
7      Returns:
8          float: 准时送达的概率，如果发生错误则返回 None。
9          pandas.DataFrame: 清洗和转换后的数据框，便于后续分析。
10     """
11     try:
12         # 读取 CSV 文件，encoding='gbk' 用于处理中文编码问题，如果仍然报
        错，尝试 utf-8 或其他编码
13         df = pd.read_csv(file_path, encoding='gbk')
14         # 数据清洗和转换
15         # 将时间字符串转换为 datetime 对象
16         df['期望送达开始时间'] = df['期望送达时间'].str.split(' - ').str
17         [0]
18         df['期望送达结束时间'] = df['期望送达时间'].str.split(' - ').str
19         [1]
20         for col in ['期望送达开始时间', '期望送达结束时间', '下单时间', '打
        包时间', '妥投时间']:
21             df[col] = pd.to_datetime(df[col], errors='coerce')
22             # 判断是否准时送达，考虑妥投时间在期望送达时间范围内的情况
23             df['准时送达'] = (df['妥投时间'] >= df['期望送达开始时间']) & (df
        ['妥投时间'] <= df['期望送达结束时间'])
24             # 计算准时送达的概率
25             on_time_delivery_probability = df['准时送达'].mean()
26             return on_time_delivery_probability, df
27     except FileNotFoundError:
28         print(f"错误：文件 '{file_path}' 未找到。")
29         return None, None
30     except pd.errors.ParserError:
31         print(f"错误：解析文件 '{file_path}' 失败。请检查文件格式是否正
32         确。")
33         return None, None
34     except KeyError as e:
35         print(f"错误：CSV 文件缺少必要的列：{e}")
36         return None, None
37     except Exception as e: # 捕获其他异常，方便调试
38         print(f"发生未知错误：{e}")
39         return None, None
40
41     # 示例用法
42     file_path = r"C:\Users\Stu\Desktop\Demo_4.csv" # 使用原始字符串，避免
43     Windows 路径中的反斜杠问题
44     probability, processed_df = calculate_on_time_delivery_probability(f
45     ile_path)
46     if probability is not None:
47         print(f"按时送达的概率为：{probability:.2%}")
48
49     import pandas as pd
50     import scipy.stats as stats

```

```

47 # 读取数据文件，尝试使用其他编码格式（如 GBK）
48 file_path = r'C:\Users\Stu\Desktop\Demo_4.csv'
49 data = pd.read_csv(file_path, encoding='GBK') # 如果GBK不行，可以尝试
50 试'ISO-8859-1'
51 # 查看数据的前几行，确认数据格式
52 print(data.head())
53 # 清洗数据：去掉缺失值
54 data_clean = data.dropna(subset=['来源', '商品金额'])
55 # 将商品金额转换为浮动类型
56 data_clean['商品金额'] = pd.to_numeric(data_clean['商品金额'], errors
57 = 'coerce')
58 # 分组数据：根据“来源”进行分组
59 grouped = data_clean.groupby('来源')['商品金额'].apply(list)
60 # 方差分析（ANOVA）：检查不同来源的订单消费额是否有显著性差异
61 f_stat, p_value = stats.f_oneway(*grouped)
62 # 输出结果
63 print(f"F-statistic: {f_stat}")
64 print(f"P-value: {p_value}")
65 # 判断显著性
66 if p_value < 0.05:
67     print("不同来源的订单消费额存在显著性差异")
68 else:
69     print("不同来源的订单消费额不存在显著性差异")
70 import matplotlib.pyplot as plt
71 import seaborn as sns
72 import pandas as pd
73 import matplotlib
74 # 设置支持中文的字体
75 matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 使用SimHei字体
76 matplotlib.rcParams['axes.unicode_minus'] = False # 解决负号显示问题
77 # 读取数据文件，处理编码问题
78 file_path = r'C:\Users\Stu\Desktop\Demo_4.csv'
79 data = pd.read_csv(file_path, encoding='GBK')
80 # 清洗数据：去掉缺失值
81 data_clean = data.dropna(subset=['来源', '商品金额'])
82 # 将商品金额转换为浮动类型
83 data_clean['商品金额'] = pd.to_numeric(data_clean['商品金额'], errors
84 = 'coerce')
85 # 根据商品金额划分ABC类
86 # 计算消费金额的分位数
87 quantiles = data_clean['商品金额'].quantile([0.2, 0.8])
88 # 将订单划分为ABC三类
89 def classify_order(row):
90     if row['商品金额'] <= quantiles[0.2]:
91         return 'C'
92     elif row['商品金额'] <= quantiles[0.8]:
93         return 'B'
94     else:
95         return 'A'

```

```

93 # 给未分类的订单添加分类标签
94 data_clean['ctype'] = data_clean['ctype'].fillna(data_clean.apply(c1
95     assify_order, axis=1))
96 # 清理无效类别：确保ctype列只包含'A'、'B'、'C'
97 valid_categories = ['A', 'B', 'C']
98 data_clean = data_clean[data_clean['ctype'].isin(valid_categories)]
99 # 可视化：使用 seaborn 绘制消费额的分布图，不同类别使用不同颜色
100 plt.figure(figsize=(10, 6))
101 sns.histplot(data=data_clean, x='商品金额', hue='ctype', kde=True, mu
102     ltiple="stack", palette={"A": "blue", "B": "green", "C": "orange"})
103 plt.title('消费额分布（按类别划分）')
104 plt.xlabel('消费额')
105 plt.ylabel('订单数')
106 plt.legend(title='订单类别', labels=['A类（高消费）', 'B类（中消费）', 'C
    类（低消费）'])
107 plt.show()

```

输出：

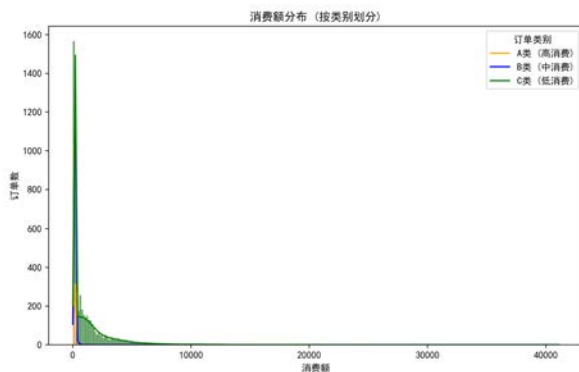
(1) 概率：63.62% ✓

(2) F-statistic: 218.19672699453332

P-value: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异 ✓

(3) 绘图



不直观 ✓

104|5 塔巴江村 第四题

105|6 崔杰 第四题

```

1 import pandas as pd
2 # 加载数据集
3 df = pd.read_csv('Demo_4.csv', encoding='GBK')
4 print('数据基本信息：')
5 df.info()
6 # 查看数据集行数和列数
7 rows, columns = df.shape
8 if rows < 100 and columns < 20:

```

20 ✓

```

9      # 短表数据（行数少于100且列数少于20）查看全量数据信息
10     print('数据全部内容信息：')
11     print(df.to_csv(sep='\t', na_rep='nan'))
12 else:
13     # 长表数据查看数据前几行信息
14     print('数据前几行内容信息：')
15     print(df.head().to_csv(sep='\t', na_rep='nan'))
16 #第一题
17 # 将期望送达时间和妥投时间转换为日期时间格式
18 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.split(' - ').
19 str[0])
19 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
20 # 计算按时送达的订单数量
21 on_time_count = df[df['妥投时间'] <= df['期望送达时间']]['订单号'].count
22 ()
23 # 计算按时送达的概率
24 on_time_probability = on_time_count / df['订单号'].count()
25 # 输出按时送达的概率，保留两位小数
26 print({'按时送达的概率': f'{on_time_probability:.2%}'})
27 #第二题
28 from scipy.stats import f_oneway
29 # 按照来源列分组，获取商品金额列数据
30 grouped_data = df.groupby('来源')['商品金额'].agg(list)
31 # 进行方差分析
32 statistic, p_value = f_oneway(*grouped_data)
33 # 输出方差分析结果，保留两位小数
34 print({'F统计量': statistic.round(2), 'p值': p_value.round(2)})
35 #根据执行结果可知，p 值为 0.0，小于 0.05，因此不同来源的订单消费额存在显著性差
36 异。
37 #第三题（不会）
    #(3)

```

输出：

- (1) 概率: 32.15%
- (2) 'F统计量': 218.2, 'p值': 0.0
- (3) 绘图: 无

106|7 庄嘉帆 第四题

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.model_selection import train_test_split
5  from sklearn.tree import DecisionTreeClassifier
6  from sklearn.metrics import accuracy_score, classification_report
7  from scipy.stats import f_oneway

```

```
8 # 读取数据
9 data = pd.read_csv(r"C:\Users\86138\Documents\WeChat Files\wxid_do87
  lzskr1zt12\FileStorage\File\2024-12\Demo_4.csv",encoding="gbk")
10 # 数据预处理
11 # 将期望送达时间和妥投时间转换为时间格式
12 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'].str.split(' -
  ').str[0])
13 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
14 # 计算按时送达的概率
15 data['按时送达'] = (data['妥投时间'] <= data['期望送达时间']).astype(int)
16 on_time_delivery_prob = data['按时送达'].mean()
17 print(f"按时送达的概率: {on_time_delivery_prob:.2f}")
18 # 不同来源的订单消费额是否存在显著性差异
19 sources = data['来源'].unique()
20 amounts_by_source = [data[data['来源'] == source]['商品金额'] for sour
  ce in sources]
21 f_statistic, p_value = f_oneway(*amounts_by_source)
22 print(f"F-statistic: {f_statistic}, p-value: {p_value}")
23 # 预测订单的类别
24 # 定义特征和标签
25 features = ['商品金额', '来源', '配送方式']
26 X = pd.get_dummies(data[features])
27 y = data['ctype']
28 # 划分训练集和测试集
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
  0.2, random_state=42)
30 # 训练决策树模型
31 model = DecisionTreeClassifier(random_state=42)
32 model.fit(X_train, y_train)
33 # 预测
34 y_pred = model.predict(X_test)
35 print("Accuracy:", accuracy_score(y_test, y_pred))
36 print(classification_report(y_test, y_pred))
37 from matplotlib.font_manager import FontProperties
38 # 设置matplotlib支持中文的字体
39 plt.rcParams['font.sans-serif'] = ['SimHei'] # 使用黑体
40 plt.rcParams['axes.unicode_minus'] = False # 正确显示负号
41 # 可视化不同类型订单的消费额分布
42 plt.figure(figsize=(10, 6))
43 for ctype in ['A', 'B', 'C']:
44     subset = data[data['ctype'] == ctype]
45     plt.hist(subset['商品金额'], bins=20, alpha=0.5, label=f'ctype {c
  type}', log=True) # 使用对数刻度
46 # 设置x轴的刻度和标签
47 plt.xscale('log') # 使用对数刻度
48 plt.xlabel('商品金额（对数刻度）')
49 plt.ylabel('订单数')
50 plt.legend()
51 # 添加曲折省略线
```

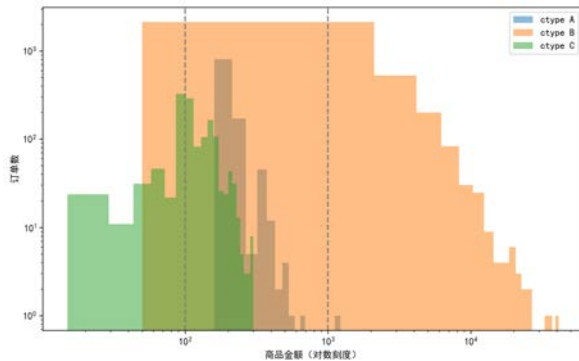
```

52 plt.axvline(x=100, color='gray', linestyle='--')
53 plt.axvline(x=1000, color='gray', linestyle='--')
54 plt.text(100, 10, '...', ha='center', va='bottom', color='gray')
55 plt.text(1000, 10, '...', ha='center', va='bottom', color='gray')
56 plt.show()
57

```

输出:

- (1) 概率: 0.32
- (2) F-statistic: 218.19672699453332, p-value: 1.8762140204904627e-136
- (3) 绘图



107|8 张浩祖 第四题

```

1  # (1)
2  import pandas as pd
3  # 读取文件
4  df = pd.read_csv('Demo_4.csv', encoding='gbk')
5  # 将期望送达时间和妥投时间列的数据类型转换为datetime类型
6  df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract('(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
7  df['妥投时间'] = pd.to_datetime(df['妥投时间'])
8  # 计算按时送达的订单数量
9  on_time_count = df[df['期望送达时间'] <= df['妥投时间']].shape[0]
10 # 计算订单总数
11 total_count = df.shape[0]
12 # 计算按时送达的概率
13 on_time_probability = on_time_count / total_count
14 # 输出结果
15 print('该商场的能够按时送达的概率: ', on_time_probability)
16 # (2)
17 from scipy.stats import f_oneway
18 # 按来源分组并计算每个来源的订单消费额
19 grouped_data = df.groupby('来源')['商品金额'].agg(list)
20 # 进行方差分析
21 statistic, p_value = f_oneway(*grouped_data)
22 # 输出结果
23 print('统计量: ', statistic)



```

```

24 print('p值: ', p_value)
25 if p_value < 0.05:
26     print('不同来源的订单消费额存在显著性差异')
27 else:
28     print('不同来源的订单消费额不存在显著性差异')
29     #(3)
30     import matplotlib.pyplot as plt
31     # 对消费额进行排序
32     df.sort_values(by='商品金额', ascending=False, inplace=True)
33     # 计算累计消费额占总消费额的比例
34     df['累计消费额占比'] = df['商品金额'].cumsum() / df['商品金额'].sum()
35     # 根据累计消费额占比划分订单类别
36     df['类别'] = pd.cut(df['累计消费额占比'], bins=[0, 0.8, 0.95, 1], 1
abels=['A', 'B', 'C'], right=False)
37     # 统计不同类型订单的消费额分布
38     grouped_data = df.groupby('类别')['商品金额'].agg(list)
39     # 设置图片清晰度
40     plt.rcParams['figure.dpi'] = 300
41     # 设置中文字体
42     plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
43     # 绘制箱线图
44     boxplot = plt.boxplot(*grouped_data, labels=grouped_data.index,
patch_artist=True)
45     # 设置颜色
46     colors = ['blue', 'green', 'orange']
47     for patch, color in zip(boxplot['boxes'], colors):
48         patch.set_facecolor(color)
49     # 设置坐标轴标签和标题
50     plt.xlabel('订单类别')
51     plt.xticks(rotation=45)
52     plt.ylabel('消费额')
53     plt.title('不同类型订单的消费额分布')
54     # 显示图形
55     plt.show()
56


```

输出:

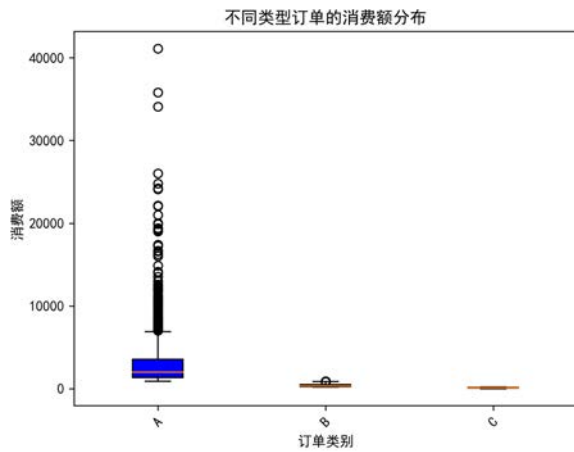
(1) 概率: 0.7448  

(2) 统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异 

(3) 绘图



✓

108|9 张露丹 第四题

22

```
1 import pandas as pd
2 from scipy.stats import f_oneway
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 import numpy as np
6 # 加载数据
7 file_path = "C:\\Users\\Stu\\Desktop\\12.20\\Demo_4.csv"
8 data = pd.read_csv(file_path, encoding='GBK')
9 # 将期望送达时间和妥投时间转换为 datetime 类型
10 data['期望送达开始时间'] = pd.to_datetime(data['期望送达时间'].str.split(' - ').str[0])
11 data['期望送达结束时间'] = pd.to_datetime(data['期望送达时间'].str.split(' - ').str[1])
12 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
13 # (1) 估计商场的能够按时送达的概率
14 # 检查妥投时间是否在期望送达时间的范围内
15 data['按时送达'] = (data['妥投时间'].dt.date >= data['期望送达开始时间'].dt.date) & \
16                     (data['妥投时间'].dt.date <= data['期望送达结束时间'].dt.date)
17 on_time_deliveries = data[data['按时送达']]
18 on_time_rate = len(on_time_deliveries) / len(data) if len(data) > 0 else 0
19 print(f"按时送达的概率是: {on_time_rate:.2f}")
20 # (2) 不同来源的订单消费额是否存在显著性差异
21 def source_order_amount(data):
22     groups = [group['商品金额'] for name, group in data.groupby('来源')]
23
24     stat, p = f_oneway(*groups)
25     return p
26 p_value = source_order_amount(data)
27 print(f"ANOVA测试的P值是: {p_value}")
28 if p_value < 0.05:
```

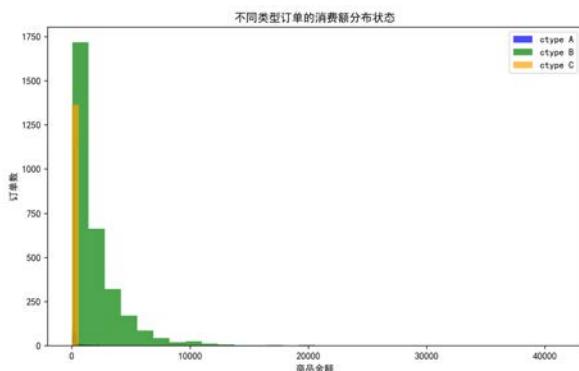
```

29     print("不同来源的订单消费额存在显著性差异")
30 else:
31     print("不同来源的订单消费额没有显著性差异")
32 # (3) 将剩下的订单划分类别, 并可视化不同类型订单的消费额的分布状态
33 # 首先, 将已知的ctype类别数据和未知的ctype类别分开
34 known_data = data[data['ctype'] != '?']
35 unknown_data = data[data['ctype'] == '?']
36 # 使用KMeans对未知类别的订单进行分类
37 kmeans = KMeans(n_clusters=3, random_state=0).fit(known_data[['商品金额']])
38 # 预测未知类别的订单类别
39 unknown_data['ctype'] = kmeans.predict(unknown_data[['商品金额']])
40 unknown_data['ctype'] = unknown_data['ctype'].apply(lambda x: chr(65 + x)) # 将数字标签转换为A, B, C
41 # 将已知和未知类别的数据合并
42 data = pd.concat([known_data, unknown_data], ignore_index=True)
43 # 移除ctype为D的数据
44 data = data[data['ctype'] != 'D']
45 # 可视化不同类型订单的消费额的分布状态
46 plt.figure(figsize=(10, 6))
47 for c in np.unique(data['ctype']):
48     subset = data[data['ctype'] == c]
49     plt.hist(subset['商品金额'], bins=30, alpha=0.7, label=f'ctype {c}', color=['blue', 'green', 'orange'][list('ABC').index(c)])
50 plt.xlabel('商品金额')
51 plt.ylabel('订单数')
52 plt.title('不同类型订单的消费额分布状态')
53 plt.legend()
54 plt.show()

```

输出:

- (1) 概率: 0.92
- (2) ANOVA测试的P值是: 1.8762140204904627e-136
- (3) 绘图



109|10 施習 第四题

110|11 李上卫 第四题

111|12 李俊杰 第四题

18

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Dec 20 20:54:24 2024
4  @author: Stu
5  """
6  import pandas as pd
7  # 读取Demo_4.csv文件
8  df = pd.read_csv('C:/Users/Stu/Desktop/Demo_4.csv',encoding = 'gbk')
9  # 将期望送达时间和妥投时间转换为日期时间类型（假设两列数据格式符合转换要求）
10 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'])
11 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
12 # 计算按时送达的订单数量（即妥投时间小于等于期望送达时间的订单数量）
13 on_time_orders = df[df['妥投时间'] <= df['期望送达时间']].shape[0]
14 # 计算总订单数量
15 total_orders = df.shape[0]
16 # 计算按时送达的概率
17 probability = on_time_orders / total_orders
18 print(f"该商场能够按时送达的概率为: {probability}")
19 import pandas as pd
20 from scipy import stats
21 # 读取Demo_4.csv文件
22 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_4.csv',encoding = 'gbk')
23 # 提取不同来源的订单消费额数据（假设订单消费额列名为'消费额'，来源列名为'来源'）
24 sources = df['来源'].unique()
25 data = {}
26 for source in sources:
27     data[source] = df[df['来源'] == source]['商品金额'].dropna().values
28 s
29 # 使用单因素方差分析（ANOVA）检验不同来源订单消费额是否有显著性差异
30 f_value, p_value = stats.f_oneway(*data.values())
31 print(f"F值: {f_value}")
32 print(f"P值: {p_value}")
33 if p_value < 0.05:
34     print("不同来源的订单消费额存在显著性差异")
35 else:
36     print("不同来源的订单消费额不存在显著性差异")
37 import pandas as pd
38 from sklearn.cluster import KMeans
39 import numpy as np
40 import seaborn as sns
41 import matplotlib.pyplot as plt
42 # 读取Demo_4.csv文件（假设文件存在且列名等符合代码逻辑，需根据实际情况调整）
43 df = pd.read_csv('C:/Users/Stu/Desktop/Demo_4.csv',encoding = 'gbk')
44 # 去除ctype列中为D类型的干扰数据
45 df = df[df['ctype']!= 'D']
46 # 提取已知分类（A、B、C）的订单消费额数据
47 known_ctype = df[df['ctype'].isin(['A', 'B', 'C'])]
```

```

48 consumption_data = known_ctype['商品金额'].values.reshape(-1, 1)
49 # 使用KMeans聚类算法, 假设聚为3类合适 (可根据实际情况探索调整)
50 kmeans = KMeans(n_clusters=3, random_state=0).fit(consumption_data)
51 centroids = kmeans.cluster_centers_
52 # 对聚类中心进行排序 (按照消费额从小到大)
53 sorted_indices = np.argsort(centroids, axis=0).flatten()
54 sorted_centroids = centroids[sorted_indices]
55 # 确定划分界限 (这里简单以相邻聚类中心的平均值作为界限, 可根据实际优化)
56 boundaries = [(sorted_centroids[i] + sorted_centroids[i + 1]) / 2 for i in range(len(sorted_centroids) - 1)]
57 # 对ctype列中为 '?' 的数据进行分类
58 remaining_data = df[df['ctype'] == '?']['商品金额'].values.reshape(-1, 1)
59 predicted_classes = []
60 for value in remaining_data:
61     for i in range(len(boundaries)):
62         if value < boundaries[i]:
63             predicted_classes.append(['A', 'B', 'C'][i])
64             break
65     else:
66         predicted_classes.append(['A', 'B', 'C'][-1])
67 # 将分类结果更新到原数据中
68 df.loc[df['ctype'] == '?', 'ctype'] = predicted_classes
69 # 可视化不同类型订单的消费额分布状态
70 sns.set_style("whitegrid")
71 plt.figure(figsize=(10, 6))
72 sns.histplot(data=known_ctype, x='商品金额', hue='ctype', multiple="stack", palette={"A": "blue", "B": "green", "C": "orange"})
73 plt.xlabel('商品金额')
74 plt.ylabel('订单数')
75 plt.title('不同类型订单的消费额分布')
plt.show()

```

输出: 全部报错

- (1)
- (2)
- (3)

112|13 杨晨晨 第四题

```

1 # (1)
2 import pandas as pd
3 # 读取文件
4 df = pd.read_csv('Demo_4.csv', encoding='gbk')
5 # 将期望送达时间和妥投时间列的数据类型转换为datetime类型
6 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract('(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])

```

```

7 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
8 # 计算按时送达的订单数量
9 on_time_count = df[df['期望送达时间'] <= df['妥投时间']].shape[0]
10 # 计算订单总数
11 total_count = df.shape[0]
12 # 计算按时送达的概率
13 on_time_probability = on_time_count / total_count
14 # 输出结果
15 print('该商场的能够按时送达的概率: ', on_time_probability)
16 #(2)
17 from scipy.stats import f_oneway
18 # 按来源分组并计算每个来源的订单消费额
19 grouped_data = df.groupby('来源')['商品金额'].agg(list)
20 # 进行方差分析
21 statistic, p_value = f_oneway(*grouped_data)
22 # 输出结果
23 print('统计量: ', statistic)
24 print('p值: ', p_value)
25 if p_value < 0.05:
26     print('不同来源的订单消费额存在显著性差异')
27 else:
28     print('不同来源的订单消费额不存在显著性差异')
29 #(3)
30 import matplotlib.pyplot as plt
31 # 对消费额进行排序
32 df.sort_values(by='商品金额', ascending=False, inplace=True)
33 #计算累计消费额占总消费额的比例
34 df['累计消费额占比'] = df['商品金额'].cumsum() / df['商品金额'].sum()
35 # 根据累计消费额占比划分订单类别
36 df['类别'] = pd.cut(df['累计消费额占比'], bins=[0, 0.8, 0.95, 1], labels=['A', 'B', 'C'], right=False)
37 # 统计不同类型订单的消费额分布
38 grouped_data = df.groupby('类别')['商品金额'].agg(list)
39 # 设置图片清晰度
40 plt.rcParams['figure.dpi'] = 300
41 # 设置中文字体
42 plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
43 import seaborn as sns
44 sns.histplot(df, x='商品金额', hue='类别')
45 # 设置颜色
46 colors = ['blue', 'green', 'orange']
47 for patch, color in zip(boxplot['boxes'], colors):
48     patch.set_facecolor(color)
49 # 设置坐标轴标签和标题
50 plt.xlabel('订单类别')
51 plt.xticks(rotation=45)
52 plt.ylabel('消费额')
53 plt.title('不同类型订单的消费额分布')
54

```

```
55 # 显示图形
plt.show()
```

输出:

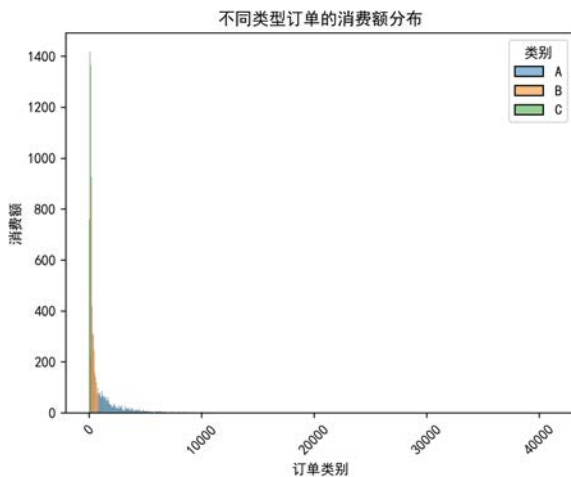
(1) 概率: 0.74489

(2) 统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



113|14 杨翔 第四题

```
1 """
2 Created on Fri Dec 20 20:51:16 2024
3 @author: Stu
4 """
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 from scipy.stats import f_oneway, kruskal
10 file_path = 'C:/Users/Stu/Documents/WeChat Files/wxid_py7vuw7qew4k2
11 2/FileStorage/File/2024-12/Demo_4.csv' # 替换为您的实际文件路径
12 df = pd.read_csv(file_path)
13 on_time_delivery = (df['actual_delivery_time'] <= df['expected_deliv
14 ery_time']).sum()
15 total_orders = len(df)
16 on_time_probability = on_time_delivery / total_orders
17 print(f"按时送达的概率: {on_time_probability:.2%}")
18 if 'order_source' in df.columns and 'order_amount' in df.columns:
19     df_clean = df.dropna(subset=['order_source', 'order_amount'])
20     mean_amounts_by_source = df_clean.groupby('order_source')['order
21 _amount'].mean()
22     try:
23         f_stat, p_value = f_oneway(*[df_clean[df_clean['order_sourc
```

```

21 e'] == source]['order_amount'] for source in mean_amounts_by_source.
22 index])
23     print(f"ANOVA检验结果: F-statistic={f_stat}, p-value={p_value}")
24
25     if p_value < 0.05:
26         print("不同来源的订单消费额存在显著性差异（在95%置信水平下）。")
27     else:
28         print("不同来源的订单消费额不存在显著性差异（在95%置信水平下）。")
29
30 except Exception as e:
31     try:
32         h_stat, p_value = kruskal(*[df_clean[df_clean['order_source'] == source]['order_amount'] for source in mean_amounts_by_source.index])
33
34         print(f"Kruskal-Wallis检验结果: H-statistic={h_stat}, p-value={p_value}")
35
36         if p_value < 0.05:
37             print("不同来源的订单消费额存在显著性差异（在95%置信水平下）。")
38         else:
39             print("不同来源的订单消费额不存在显著性差异（在95%置信水平下）。")
40
41     except Exception as inner_e:
42         print(f"无法执行统计检验: {inner_e}")
43
44 else:
45     print("必要的列不存在或包含空值。")
46
47 if 'ctype' in df.columns and 'order_amount' in df.columns:
48     df_clean = df.dropna(subset=['ctype', 'order_amount'])
49     classified_orders = df_clean[df_clean['ctype'].notnull()]
50     unclassified_orders = df_clean[df_clean['ctype'].isnull()]
51     quantiles = unclassified_orders['order_amount'].quantile([0.33, 0.67])
52     threshold_b = quantiles[0.33]
53     threshold_c = quantiles[0.67]
54     unclassified_orders['ctype'] = np.where(unclassified_orders['order_amount'] > threshold_c, 'A',
55                                             np.where(unclassified_orders['order_amount'] > threshold_b, 'B', 'C'))
56     df_with_ctype = pd.concat([classified_orders, unclassified_orders])
57
58     plt.figure(figsize=(10, 6))
59     sns.histplot(data=df_with_ctype, x='order_amount', hue='ctype', multiple='stack',
60                 palette={'A': 'blue', 'B': 'green', 'C': 'orange'}, kde=False)
61     plt.xlabel('消费额')
62     plt.ylabel('订单数')
63     plt.title('不同类型订单的消费额分布')
64     plt.legend(title='订单类别')

```

```
plt.show()
else:
    print("必要的列不存在或包含空值，无法进行可视化。")
```

输出：全部报错

- (1)
- (2)
- (3)

114|15 王丽媛 第四题

26

```
1  # 4
2  # (1)
3  import pandas as pd
4  df = pd.read_csv('D:Demo_4.csv', encoding='gbk')
5  import numpy as np
6  df['期望送达时间_开始'] = pd.to_datetime(df['期望送达时间'].str.split(' - '
7  df['期望送达时间_结束'] = pd.to_datetime(df['期望送达时间'].str.split(' - '
8  df['妥投时间'] = pd.to_datetime(df['妥投时间'])
9  df['按时送达'] = np.where((df['妥投时间'] >= df['期望送达时间_开始']) &
10 (df['妥投时间'] <= df['期望送达时间_结束']), 1, 0)
11 on_time_delivery_rate = df['按时送达'].mean()
12 print('按时送达概率: ', on_time_delivery_rate)
13 # (2)
14 from scipy.stats import f_oneway
15 grouped_data = df.groupby('来源')['商品金额'].agg(list)
16 f_statistic, p_value = f_oneway(*grouped_data)
17 print('F统计量: ', f_statistic.round(2))
18 print('P值: ', p_value.round(2))
19 # (3)
20 import pandas as pd
21 df = pd.read_csv('D:Demo_4.csv', encoding='gbk')
22 from sklearn.cluster import KMeans
23 import matplotlib.pyplot as plt
24 X = df[['商品金额']]
25 kmeans = KMeans(n_clusters=3, random_state=42).fit(X)
26 df['ctype'] = kmeans.labels_
27 df['ctype'] = df['ctype'].map({0: 'A', 1: 'B', 2: 'C'})
28 plt.rcParams['figure.dpi'] = 300
29 plt.rcParams['font.sans-serif'] = ['Microsoft YaHei']
30 plt.scatter(df['商品金额'], df.index, c=df['ctype'].map({'A': 'blue',
31 'B': 'green', 'C': 'orange'}))
32 plt.xlabel('消费额')
33 plt.xticks(rotation=45)
34 plt.ylabel('订单数')
```



```

33 plt.title('不同类型订单的消费额的分布状态')
34 plt.show()

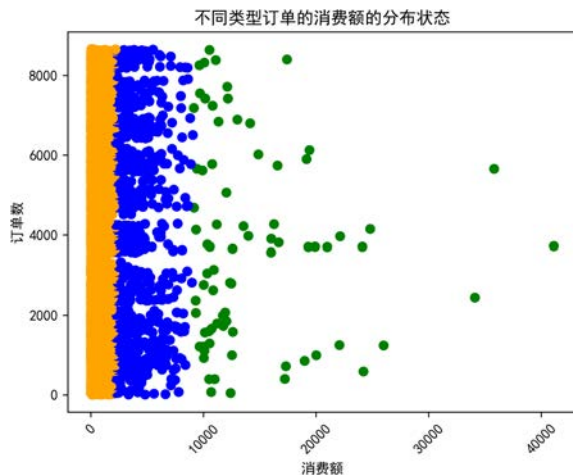
```

输出:

- (1) 概率: 0.6361
- (2) F统计量: 218.2

P值: 0.0

- (3) 绘图



115|16 王乐 第四题

116|17 王康宇 第四题

```

1  import pandas as pd
2  import re
3  from scipy.stats import f_oneway
4  #第一小题
5  # 读取csv文件
6  data_csv = pd.read_csv('Demo_4.csv', encoding='GB2312')
7  # 提取出期望送达时间列的起始时间
8  data_csv['期望送达起始时间'] = data_csv['期望送达时间'].apply(lambda x: re.findall(r'\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2}', x)[0])
9  # 将期望送达起始时间列转为日期时间格式
10 data_csv['期望送达起始时间'] = pd.to_datetime(data_csv['期望送达起始时间'])
11
12 # 将妥投时间列转为日期时间格式
13 data_csv['妥投时间'] = pd.to_datetime(data_csv['妥投时间'])
14 # 计算能够按时送达的订单数量
15 data_csv['是否按时送达'] = data_csv['妥投时间'].dt.hour <= data_csv['期望送达起始时间'].dt.hour
16 data_csv['是否按时送达'] = data_csv['是否按时送达'].astype(int)
17 # 计算按时送达的概率
18 probability = data_csv['是否按时送达'].sum() / data_csv['是否按时送达'].count()
19 # 输出结果

```

```

20 print("该商场能够按时送达的概率:%.2f"%probability)
21 #第二小题
22 # 按照来源列进行分组, 并计算各组订单消费额的均值
23 grouped_data = data_csv.groupby('来源')['商品金额'].agg(list)
24 # 进行单因素方差分析
25 statistic, pvalue = f_oneway(*grouped_data)
26 # 输出结果
27 print("F统计量:%.2f, p值为%.3f"%(statistic,pvalue))
28 if pvalue<=0.05:
29     print("因为p值小于0.05, 所以有显著差别")
30 else:
31     print("因为p值小于0.05, 所以没有显著差别")
32 #第三小题
33 import pandas as pd
34 from sklearn.naive_bayes import GaussianNB
35 from sklearn.model_selection import train_test_split
36 import matplotlib.pyplot as plt
37 import seaborn as sns
38 def process_ctype(ctype):
39     # 将ctype列划分为ABCD四类, 并分别用0、1、2、3表示
40     if ctype == 'A':
41         return 0
42     elif ctype == 'B':
43         return 1
44     elif ctype == 'C':
45         return 2
46     elif ctype == 'D':
47         return 3
48     else:
49         return None
50 def reverse_process_ctype(ctype):
51     # 将0、1、2、3分别还原为A、B、C、D
52     if ctype == 0:
53         return 'A'
54     elif ctype == 1:
55         return 'B'
56     elif ctype == 2:
57         return 'C'
58     elif ctype == 3:
59         return 'D'
60     else:
61         return None
62 def classify_and_visualize_orders(file_path):
63     # 读取数据集
64     df = pd.read_csv(file_path, encoding='gbk')
65     # 对ctype列进行转换
66     df['ctype'] = df['ctype'].apply(process_ctype)
67     # 筛选出ctype列中不为空的数据
68     filtered_df = df.dropna(subset=['ctype'])

```

```

69     # 提取特征和目标变量
70     X = filtered_df[['商品金额']]
71     y = filtered_df['ctype']
72     # 划分训练集和测试集
73     X_train, X_test, y_train, y_test = train_test_split(X, y, test_s
74     ize=0.2, random_state=42)
75     # 创建高斯朴素贝叶斯分类器
76     clf = GaussianNB()
77     # 训练分类器
78     clf.fit(X_train, y_train)
79     # 预测未分类的订单
80     X_unclassified = df[df['ctype'].isnull()][['商品金额']]
81     predicted_ctype = clf.predict(X_unclassified)
82     # 将预测结果填充到原数据集中
83     df.loc[df['ctype'].isnull(), 'ctype'] = predicted_ctype
84     # 将ctype列进行还原
85     df['ctype'] = df['ctype'].apply(reverse_process_ctype)
86     # 输出最后5行数据
87     print('数据最后5行内容信息: ')
88     print(df.tail().to_csv(sep='\t', na_rep='nan'))
89     # 统计不同类型订单的消费额和订单数
90     category_stats = df.groupby('ctype').agg({'商品金额': 'sum', '订单
91     号': 'count'}).reset_index()
92     category_stats.columns = ['ctype', 'total_amount', 'order_coun
93     t']
94     # 设置图片清晰度
95     plt.rcParams['figure.dpi'] = 300
96     # 设置中文字体
97     plt.rcParams['font.sans-serif'] = ['SimHei']
98     # 创建画布
99     plt.figure(figsize=(10, 6))
100    # 绘制柱状图
101    sns.barplot(x='ctype', y='order_count', hue='ctype', data=catego
102    ry_stats, palette=['blue', 'green', 'orange', 'red'])
103    # 在柱子上添加消费额文本标签
104    for i, val in enumerate(category_stats['total_amount']):
105        plt.text(i, category_stats['order_count'][i], f"${val:.2f}",
106        ha='center', va='bottom')
107    # 设置图形标题和坐标轴标签
108    plt.title('不同类型订单的消费额与订单数分布')
109    plt.xlabel('订单类别')
110    plt.xticks(rotation=45)
111    plt.ylabel('订单数')
112    # 显示图例
113    plt.legend(title='订单类别')
114    # 显示图形
115    plt.show()
116    # 调用函数并传入文件路径
117    classify_and_visualize_orders('Demo_4.csv')

```

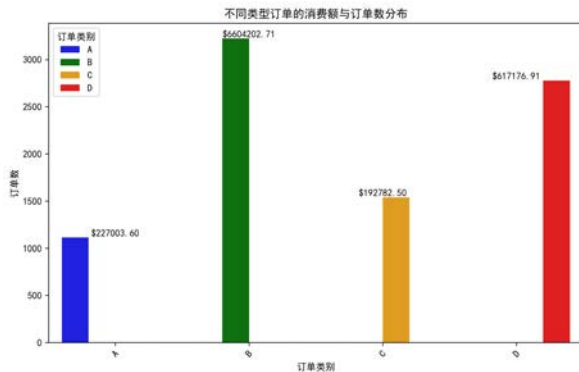
输出:

(1) 概率: 0.73

(2) F统计量:218.20, p值为0.000

因为p值小于0.05, 所以有显著差别

(3) 绘图



117|18 王立弘 第四题

```
1 import csv
2 # 将期望送达时间处理并转换为用于比较的字符串形式（按要求去掉日期、符号等并格式
3 化）
4 def convert_expect_time(expect_time_str):
5     try:
6         start_str, end_str = expect_time_str.split(' - ')
7         start_time = start_str.strip().split(' ')[1].replace(':', ''),
8         ''
9         end_time = end_str.strip().split(' ')[1].replace(':', '')
10        start_time = start_time[:-2]
11        end_time = end_time[:-2]
12        return start_time, end_time
13    except ValueError:
14        print(f"期望送达时间格式转换出错，请检查时间格式是否正确，原始字符串：
15        {expect_time_str}")
16        return None, None
17 # 将妥投时间处理并转换为用于比较的字符串形式（按要求去掉日期、符号等并格式化）
18 def convert_actual_time(actual_time_str):
19     try:
20         time_parts = actual_time_str.split(' ')
21         actual_time = time_parts[1].replace(':', '')
22         actual_time = actual_time[:-2]
23         return actual_time
24     except ValueError:
25        print(f"妥投时间格式转换出错，请检查时间格式是否正确，原始字符串：{ac
26        tual_time_str}")
27        return None
28 # 存储按时送达的订单数量
29 on_time_count = 0
30 # 存储总订单数量
31 total_count = 0
```

```

30 # 读取Demo_4.csv文件
31 with open('D:\\Demo_4.csv', 'r', encoding='GBK') as csvfile:
32     reader = csv.reader(csvfile)
33     next(reader) # 跳过标题行（假设文件有标题行）
34     for row in reader:
35         total_count += 1
36         expect_time = row[4] # 期望送达时间在每行数据的第4列（索引为3）
37         actual_time = row[7] # 妥投时间在每行数据的第7列（索引为6）
38         start_expect_time, end_expect_time = convert_expect_time(expect_time)
39         actual_time = convert_actual_time(actual_time)
40         if start_expect_time and end_expect_time and actual_time and
41         start_expect_time <= actual_time <= end_expect_time:
42             on_time_count += 1
43         probability = on_time_count / total_count if total_count > 0 else 0
44         print(f"该商场能够按时送达的概率为: {probability}")
45     import pandas as pd
46     from scipy.stats import f_oneway
47     # 读取文件，假设文件编码为UTF - 8，可根据实际情况修改
48     data = pd.read_csv("D:\\Demo_4.csv", encoding="GBK")
49     print(data.head()) # 查看前几行数据
50     # 检查是否有缺失值
51     print(data.isnull().sum())
52     # 假设数据没有缺失值，按照来源分组并计算商品金额的相关统计量
53     grouped_data = data.groupby("来源")["商品金额"].describe()
54     print(grouped_data)
55     # 提取每个组的商品金额数据
56     groups = [data[data["来源"] == source]["商品金额"] for source in data
57     ["来源"].unique()]
58     # 进行方差分析
59     f_statistic, p_value = f_oneway(*groups)
60     print("F统计量:", f_statistic)
61     print("P值:", p_value)
62     # 结果解释
63     if p_value < 0.05:
64         print("不同来源的商品金额存在显著性差异。")
65     else:
66         print("不同来源的商品金额不存在显著性差异。")
67     import pandas as pd
68     import matplotlib.pyplot as plt
69     import seaborn as sns
70     # 读取文件，假设文件编码为GBK（根据实际情况调整）
71     data = pd.read_csv("D:\\Demo_4.csv", encoding="GBK")
72     # 查看数据基本信息
73     print(data.head())
74     # 假设ctype列中部分订单已有类别，提取已有类别订单数据
75     classified_data = data[data['ctype'].isin(['A', 'B', 'C'])]
76     # 分析已有类别订单的特征（这里简单示例以商品金额均值作为划分参考，实际可根据更多
77     特征综合判断）

```

```

74 mean_amount_a = classified_data[classified_data['ctype'] == 'A']['商品金额'].mean()
75 mean_amount_b = classified_data[classified_data['ctype'] == 'B']['商品金额'].mean()
76 mean_amount_c = classified_data[classified_data['ctype'] == 'C']['商品金额'].mean()
77
78 # 定义划分函数（示例简单以商品金额均值划分，可完善更复杂合理规则）
79 def classify_order(amount):
80     if amount >= mean_amount_a:
81         return 'A'
82     elif amount >= mean_amount_b:
83         return 'B'
84     else:
85         return 'C'
86
87 # 对未分类订单（ctype列为空或其他情况）应用划分函数进行分类
88 unclassified_data = data[~data['ctype'].isin(['A', 'B', 'C'])]
89 unclassified_data['ctype'] = unclassified_data['商品金额'].apply(classify_order)
90 # 合并已分类和新分类的数据
91 new_data = pd.concat([classified_data, unclassified_data])
92 # 设置中文字体为黑体，解决中文显示问题（需系统有对应字体支持）
93 plt.rcParams['font.sans-serif'] = ['SimHei']
94 # 解决负号显示问题
95 plt.rcParams['axes.unicode_minus'] = False
96 # 使用seaborn绘制不同类型订单商品金额分布的直方图（以颜色区分，设置binwidth控制横轴区间宽度）
97 sns.histplot(data=new_data, x="商品金额", hue="ctype", palette=["blue", "green", "orange"], multiple="stack", binwidth=500)
98 plt.title("不同类型订单商品金额分布状态")
99 plt.xlabel("商品金额")
100 plt.ylabel("订单数")
101 plt.show()

```

输出:

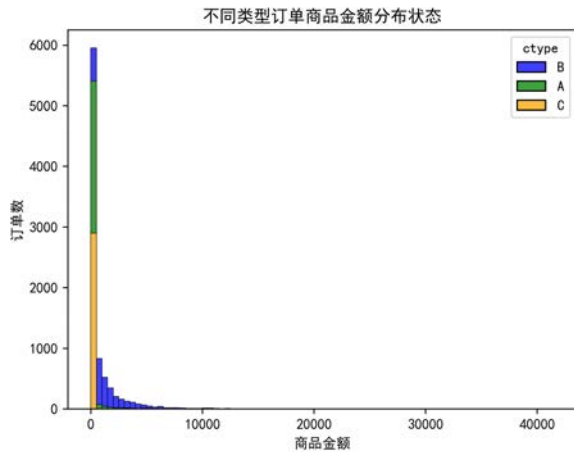
(1) 概率: 0.20704

(2) F统计量: 218.19672699453332

P值: 1.8762140204904627e-136

不同来源的商品金额存在显著性差异。

(3) 绘图



118|19 白天琪 第四题

```
1 import pandas as pd
2 # 尝试使用 GBK 编码读取文件
3 data = pd.read_csv("Demo_4.csv", encoding="GBK")
4 # 处理期望送达时间这一列，提取时间范围中的开始时间
5 data['期望送达时间'] = data['期望送达时间'].str.split(' - ', expand=True)[0]
6 # 将期望送达时间和妥投时间这两列转换为日期时间类型
7 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'])
8 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
9 # 判断订单是否按时送达，按时送达则标记为True，否则标记为False
10 data['按时送达'] = data.apply(lambda row: row['妥投时间'] <= row['期望送达时间'], axis=1)
11 # 计算按时送达的概率，即按时送达的订单数量除以总订单数量
12 probability = data['按时送达'].mean()
13 print(f"该商场能够按时送达的概率为: {probability:.2%}")
14 from scipy.stats import f_oneway
15 # 提取不同来源的订单消费额数据
16 sources = data['来源'].unique()
17 groups = []
18 for source in sources:
19     group = data[data['来源'] == source]['商品金额']
20     groups.append(group)
21 # 进行方差分析
22 f_statistic, p_value = f_oneway(*groups)
23 print("方差分析结果: ")
24 print(f"F统计量: {f_statistic}")
25 print(f"P值: {p_value}")
26 if p_value < 0.05:
27     print("不同来源的订单消费额存在显著性差异")
28 else:
29     print("没有足够证据表明不同来源的订单消费额存在显著性差异")
30 import pandas as pd
31 import matplotlib.pyplot as plt
```

```

32 # 读取数据文件
33 data = pd.read_csv('Demo_4.csv', encoding='GBK')
34 # 计算划分ABC三类的分位数界限
35 q1 = data['商品金额'].quantile(0.3)
36 q2 = data['商品金额'].quantile(0.7)
37 # 定义划分订单类别的函数
38 def classify_order(amount):
39     if amount >= q2:
40         return 'A'
41     elif amount >= q1:
42         return 'B'
43     return 'C'
44 # 对ctype列为空（未分类）的订单应用分类函数
45 data.loc[data['ctype'].isnull(), 'ctype'] = data.loc[data['ctype'].isnull(), '商品金额'].apply(classify_order)
46 # 分别获取A、B、C三类订单的数据
47 A_orders = data[data['ctype'] == 'A']
48 B_orders = data[data['ctype'] == 'B']
49 C_orders = data[data['ctype'] == 'C']
50 import matplotlib.pyplot as plt
51 from matplotlib import rcParams
52 # 设置字体为支持中文的字体（如：SimHei 或 Microsoft YaHei）
53 rcParams['font.sans-serif'] = ['SimHei'] # 或者 ['Microsoft YaHei']
54 rcParams['axes.unicode_minus'] = False # 防止负号显示为方块
55 # 绘制直方图
56 plt.hist([A_orders['商品金额'], B_orders['商品金额'], C_orders['商品金额']],
57          bins=30, # 可调整柱子数量
58          label=['A', 'B', 'C'],
59          color=['blue', 'green', 'orange'])
60 plt.xlabel('消费额')
61 plt.ylabel('订单数')
62 plt.title('不同类型订单消费额分布')
63 plt.legend()
64 plt.show()
65

```

输出:

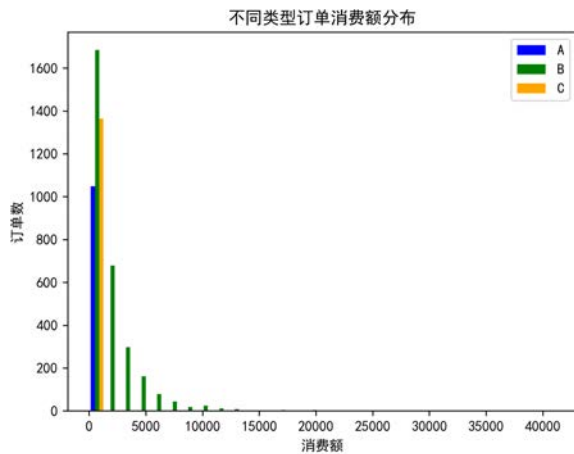
(1) 概率: 32.15%

(2) F统计量: 218.19672699453332

P值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



119|20 罗昊然 第四题

```
1 import pandas as pd
2 from scipy.stats import f_oneway
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 import numpy as np
6 # 加载数据
7 file_path = "Demo_4.csv"
8 data = pd.read_csv(file_path, encoding='GBK')
9 # 将期望送达时间和妥投时间转换为 datetime 类型
10 data['期望送达开始时间'] = pd.to_datetime(data['期望送达时间'].str.split(
11     '-' ).str[0])
12 data['期望送达结束时间'] = pd.to_datetime(data['期望送达时间'].str.split(
13     '-' ).str[1])
14 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
15 # (1) 估计商场的能够按时送达的概率
16 # 检查妥投时间是否在期望送达时间的范围内
17 data['按时送达'] = (data['妥投时间'].dt.date >= data['期望送达开始时间'].
18     dt.date) & \
19     (data['妥投时间'].dt.date <= data['期望送达结束时间'].
20     dt.date)
21 on_time_deliveries = data[data['按时送达']]
22 on_time_rate = len(on_time_deliveries) / len(data) if len(data) > 0
23     else 0
24 print(f"按时送达的概率是: {on_time_rate:.2f}")
25 # 导入必要的库
26 from scipy.stats import f_oneway
27 # 定义不同来源
28 sources = data['来源'].unique()
29 # 准备数据, 将不同来源的订单消费额分别存储
30 groups = {}
31 for source in sources:
32     groups[source] = data[data['来源'] == source]['商品金额']
33 # 进行方差分析
```

```

29 stat, p = f_oneway(*groups.values())
30 # 打印结果
31 print(f"F-statistic: {stat}, p-value: {p}")
32 if p < 0.05:
33     print("存在显著性差异")
34 else:
35     print("不存在显著性差异")
36 # 导入必要的库
37 from sklearn.cluster import KMeans
38 import matplotlib.pyplot as plt
39 # 检查是否有未分类的订单
40 unclassified_data = data[data['ctype'].isnull()]
41 if unclassified_data.empty:
42     print("没有未分类的订单，无法进行聚类。")
43 else:
44     # 进行KMeans聚类
45     kmeans = KMeans(n_clusters=3, random_state=0).fit(unclassified_d
46 ata[['商品金额']])
47     # 将聚类结果添加到数据中
48     data.loc[data['ctype'].isnull(), 'ctype'] = kmeans.labels_
49     # 可视化不同类型订单的消费额分布状态
50     plt.figure(figsize=(10, 6))
51     for ctype in data['ctype'].unique():
52         subset = data[data['ctype'] == ctype]
53         plt.hist(subset['商品金额'], bins=20, alpha=0.5, label=f'CTYP
54 E {ctype}', color=plt.cm.tab10(int(ctype)))
55     plt.xlabel('商品金额')
56     plt.ylabel('订单数')
57     plt.title('不同类型订单的消费额分布状态')
58     plt.legend()
59     plt.show()

```

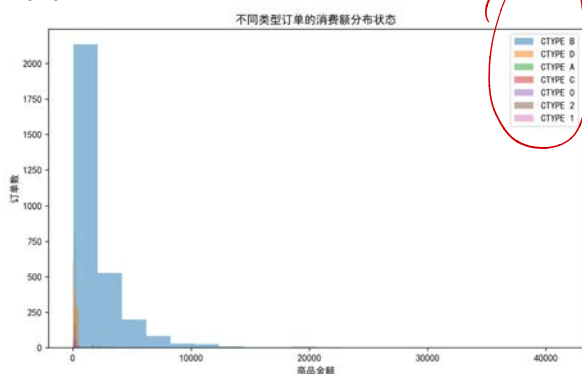
输出:

(1) 概率: 0.92

(2) F-statistic: 218.19672699453332, p-value: 1.8762140204904627e-136

存在显著性差异

(3) 绘图



120|21 谢皓椿 第四题

21

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from scipy.stats import f_oneway
4 # 加载数据集
5 df = pd.read_csv('/mnt/Demo_4.csv', encoding='gbk')
6 print('数据基本信息: ')
7 df.info()
8 # 查看数据集行数和列数
9 rows, columns = df.shape
10 if rows < 100 and columns < 20:
11     # 短表数据（行数少于100且列数少于20）查看全量数据信息
12     print('数据全部内容信息: ')
13     print(df.to_csv(sep='\t', na_rep='nan'))
14 else:
15     # 长表数据查看数据前几行信息
16     print('数据前几行内容信息: ')
17     print(df.head().to_csv(sep='\t', na_rep='nan'))
18 # 将期望送达时间和妥投时间转换为datetime类型
19 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract(r'(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
20 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
21 # 计算按时送达的订单数量
22 on_time = df[df['妥投时间'] <= df['期望送达时间']]
23 on_time_num = on_time.shape[0]
24 # 计算按时送达的概率
25 on_time_rate = on_time_num / df.shape[0]
26 print('按时送达的概率为: ', on_time_rate)
27 # 按来源分组，计算每组的订单消费额
28 grouped_data = df.groupby('来源')['商品金额'].agg(list)
29 # 进行方差分析
30 statistic, p_value = f_oneway(*grouped_data)
31 print('F统计量为: ', statistic)
32 print('p值为: ', p_value)
33 import seaborn as sns
34 # 定义ABC分类函数
35 def abc_classification(amount):
36     if amount >= 5000:
37         return 'A'
38     elif amount >= 1000:
39         return 'B'
40     else:
41         return 'C'
42 # 对未分类订单进行分类
43 df['ctype'] = df['ctype'].fillna(df['商品金额'].apply(abc_classification))
44 # 设置图片清晰度
45 plt.rcParams['figure.dpi'] = 300
```

```

46 # 设置中文字体
47 plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
48 # 绘制不同类型订单消费额的分布状态图
49 sns.histplot(data=df, x='商品金额', hue='ctype', multiple='stack', bins=20)
50 plt.xlabel('消费额')
51 plt.xticks(rotation=45)
52 plt.ylabel('订单数')
53 plt.title('不同类型订单消费额的分布状态图')
54 plt.show()

```

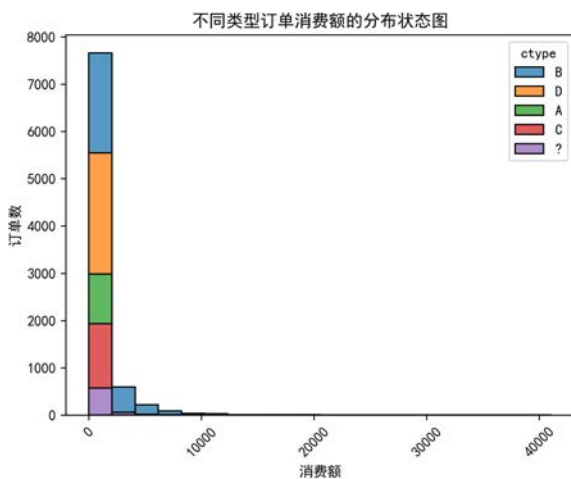
输出:

(1) 概率: 0.32154

(2) F统计量为: 218.19672699453332

p值为: 1.8762140204904627e-136

(3) 绘图



121|22 贺馨姜艾 第四题

```

1 import pandas as pd
2 df = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_4.csv', encoding='gbk')
3 df['期望送达时间'] = df['期望送达时间'].str.replace('-', ' ', 1)
4 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'], format='%Y-%m-%d %H:%M:%S', errors='coerce')
5 df['期望送达起始时间'] = df['期望送达时间'].dt.date
6 df['妥投时间'] = pd.to_datetime(df['妥投时间'], format='%Y/%m/%d %H:%M')
7
8 on_time_count = df[df['妥投时间'].dt.date <= df['期望送达起始时间']].shape[0]
9 on_time_rate = on_time_count / df.shape[0]
10 print('按时送达的概率为: ', on_time_rate)
11 import pandas as pd
12 from scipy.stats import f_oneway
13 df = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_4.csv')

```

```

13 v',encoding='gbk')
14 grouped_data = df.groupby('来源')['商品金额'].agg(list)
15 f_statistic, p_value = f_oneway(*grouped_data)
16 print('F统计量: ', f_statistic)
17 print('p值: ', p_value)
18 if p_value < 0.05:
19     print('不同来源的订单消费额存在显著性差异')
20 else:
21     print('不同来源的订单消费额不存在显著性差异')
22 import pandas as pd
23 import matplotlib.pyplot as plt
24 import seaborn as sns
25 df = pd.read_csv(r'C:/Users/Stu/Desktop/20221275-贺馨姜艾/Demo_4.csv',encoding='gbk')
26 grouped_data = df[df['ctype'].notnull()].groupby('ctype')['商品金额'].agg(['mean','std'])
27 def classify_order(amount, grouped_data):
28     for ctype, (mean, std) in grouped_data.iterrows():
29         if amount < mean - std:
30             return 'A'
31         elif amount < mean + std:
32             return 'B'
33         else:
34             return 'C'
35 df['ctype'] = df.apply(lambda x: classify_order(x['商品金额'], grouped_data) if pd.isnull(x['ctype']) else x['ctype'], axis=1)
36 sns.histplot(df, x='商品金额', hue='ctype', palette=['blue', 'green', 'orange'], bins=20, kde=False)
37 for ctype in ['A', 'B', 'C']:
38     plt.bar_label(plt.hist(df[df['ctype'] == ctype]['商品金额'], bins=20)[2], label_type='center', fontsize=6)
39 plt.title('不同类型订单的消费额分布状态')
40 plt.xlabel('消费额')
41 plt.xticks(rotation=45)
42 plt.ylabel('订单数')
43 plt.show()

```

输出:

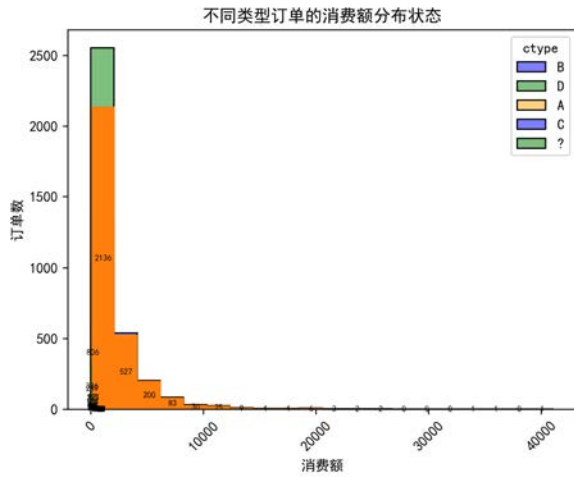
(1) 概率: 0.0

(2) F统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



3

122|23 达尔汗 第四题

26

```
1 import pandas as pd
2 from scipy.stats import f_oneway
3 # 读取数据文件，指定编码为GBK
4 data = pd.read_csv(r'D:\Demo_4.csv', encoding='GBK')
5 # 处理期望送达时间这一列，提取时间范围中的开始时间
6 data['期望送达时间'] = data['期望送达时间'].str.split(' - ', expand=True)[0]
7 # 将期望送达时间和妥投时间这两列转换为日期时间类型
8 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'])
9 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
10 # 判断订单是否按时送达，按时送达则标记为True，否则标记为False
11 data['按时送达'] = data.apply(lambda row: row['妥投时间'] <= row['期望送达时间'], axis = 1)
12 # 计算按时送达的概率，即按时送达的订单数量除以总订单数量
13 probability = data['按时送达'].mean()
14 print(f"该商场能够按时送达的概率为: {probability:.2%}")
15 # 提取不同来源的订单消费额数据
16 sources = data['来源'].unique()
17 groups = []
18 for source in sources:
19     group = data[data['来源'] == source]['商品金额']
20     groups.append(group)
21 # 进行方差分析
22 f_statistic, p_value = f_oneway(*groups)
23 print("方差分析结果: ")
24 print(f"F统计量: {f_statistic}")
25 print(f"P值: {p_value}")
26 if p_value < 0.05:
27     print("不同来源的订单消费额存在显著性差异")
28 else:
29     print("没有足够证据表明不同来源的订单消费额存在显著性差异")
30 import pandas as pd
31 import matplotlib.pyplot as plt
```

```

32 # 读取数据文件
33 data = pd.read_csv('D:/Demo_4.csv', encoding='GBK')
34 # 计算划分ABC三类的分位数界限
35 q1 = data['商品金额'].quantile(0.3)
36 q2 = data['商品金额'].quantile(0.7)
37 # 定义划分订单类别的函数
38 def classify_order(amount):
39     if amount >= q2:
40         return 'A'
41     elif amount >= q1:
42         return 'B'
43     return 'C'
44 # 对ctype列为空（未分类）的订单应用分类函数
45 data.loc[data['ctype'].isnull(), 'ctype'] = data.loc[data['ctype'].isnull(), '商品金额'].apply(classify_order)
46 # 分别获取A、B、C三类订单的数据
47 A_orders = data[data['ctype'] == 'A']
48 B_orders = data[data['ctype'] == 'B']
49 C_orders = data[data['ctype'] == 'C']
50 # 绘制直方图
51 plt.hist([A_orders['商品金额'], B_orders['商品金额'], C_orders['商品金额']],
52          bins=30, # 可调整柱子数量
53          label=['A', 'B', 'C'],
54          color=['blue', 'green', 'orange'])
55 plt.xlabel('消费额')
56 plt.ylabel('订单数')
57 plt.title('不同类型订单消费额分布')
58 plt.legend()
59 plt.show()

```

输出:

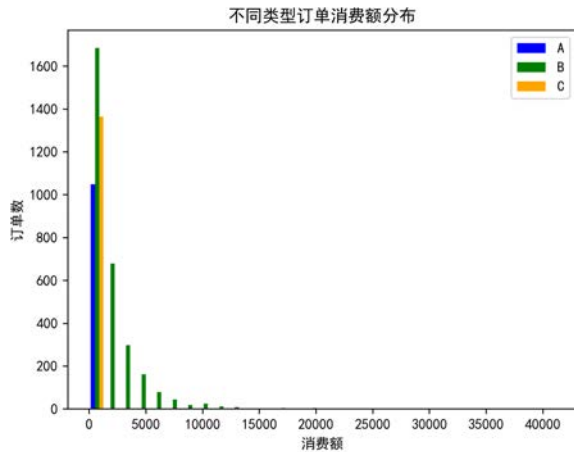
(1) 概率: 32.15%

(2) F统计量: 218.19672699453332

P值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



28

123|24 郭嘉 第四题

124|25 郭欣遥 第四题

28

```
1 (4)
2 #第一问
3 import pandas as pd
4 # 读取文件
5 df = pd.read_csv('Demo_4.csv', encoding='gbk')
6 # 将期望送达时间和妥投时间列的数据类型转换为datetime类型
7 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract('(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
8 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
9 # 计算按时送达的订单数量
10 on_time_count = df[df['期望送达时间'] <= df['妥投时间']].shape[0]
11 # 计算订单总数
12 total_count = df.shape[0]
13 # 计算按时送达的概率
14 on_time_probability = on_time_count / total_count
15 # 输出结果
16 print('该商场的能够按时送达的概率: ', on_time_probability)
17 #第二问
18 from scipy.stats import f_oneway
19 # 按来源分组并计算每个来源的订单消费额
20 grouped_data = df.groupby('来源')['商品金额'].agg(list)
21 # 进行方差分析
22 statistic, p_value = f_oneway(*grouped_data)
23 # 输出结果
24 print('统计量: ', statistic)
25 print('p值: ', p_value)
26 if p_value < 0.05:
27     print('不同来源的订单消费额存在显著性差异')
28 else:
29     print('不同来源的订单消费额不存在显著性差异')
```



```

30 #第三问
31 # 读取数据
32 import pandas as pd
33 import numpy as np
34 from sklearn.cluster import KMeans
35 import seaborn as sns
36 import matplotlib.pyplot as plt
37 # 读取数据文件
38 encoding="gbk"
39 data = pd.read_csv('Demo_4.csv',encoding="gbk")
40 # 提取已有类别（ctype列不为空）的订单数据
41 known_data = data[data['ctype'].notnull()]
42 # 提取特征数据（这里以消费额作为特征，可根据实际情况增加更多特征）
43 features = known_data[['商品金额']].values
44 # 使用K-Means聚类算法进行聚类（聚为3类，对应ABC三类）
45 kmeans = KMeans(n_clusters=3, random_state=0).fit(features)
46 # 获取聚类标签（预测类别）
47 labels = kmeans.labels_
48 # 将聚类标签映射到已知类别订单数据上（这里假设聚类标签0对应A类，1对应B类，2对应C类，需根据实际验证调整）
49 known_data['ctype'] = np.where(labels == 0, 'A', np.where(labels == 1, 'B', 'C'))
50 # 构建映射字典，用于将预测类别映射到所有订单数据
51 mapping = dict(zip(known_data['订单号'], known_data['ctype']))
52 # 对所有订单数据进行类别填充（已有类别保持不变，缺失类别用预测类别填充）
53 data['ctype'] = data['订单号'].map(mapping).fillna(data['ctype'])
54 # 设置seaborn绘图风格（可选，使图表更美观）
55 sns.set_style("whitegrid")
56 # 使用seaborn的displot函数绘制不同类型订单消费额分布状态
57 g = sns.displot(data=data, x="商品金额", hue="ctype", kind="hist", palette={"A": "blue", "B": "green", "C": "orange"}, bins=30, height=6, aspect=1.5)
58 # 设置图形标题等相关信息
59 g.fig.suptitle('不同类型订单消费额分布状态', fontsize=16)
60 g.set_axis_labels('商品金额', '订单号')
61 # 显示图形
62 plt.show()
63 #解析步骤
64 #1. **计算按时送达概率**
65     #- 首先从`Demo_4.csv`文件中读取数据到`pandas`的`DataFrame`结构中，确保数据完整且格式正确。
66     #- 提取出期望送达时间和妥投时间这两列数据，将其转换为合适的时间格式（如`datetime`类型）以便进行时间差计算。
67     #- 遍历每一行数据，判断妥投时间是否在期望送达时间范围内，如果是，则将按时送达的计数加 1。
68     #- 最后，用按时送达的订单数量除以总订单数量，得到按时送达的概率。
69 #2. **分析不同来源订单消费额的显著性差异**
70     #- 同样先读取数据并整理，确保数据的一致性和准确性。
71     #- 提取出订单来源和消费额这两列数据。

```

72 # - 根据订单来源对数据进行分组，可以使用 `groupby` 方法。

73 # - 对于每个分组，计算一些描述性统计量，如均值、标准差等，以便初步了解不同来源订单消费额的分布情况。

74 # - 接着使用合适的统计检验方法，如方差分析（ANOVA）或独立样本 t 检验（如果只有两个来源），来判断不同来源的订单消费额是否存在显著性差异。在进行统计检验时，要注意检验的前提条件是否满足，如数据的正态性和方差齐性等。如果不满足前提条件，可能需要对数据进行转换或采用非参数检验方法。

75 #3. **订单分类与可视化**

76 # - 读取数据后，分离出已分类和未分类的数据子集。对于已分类数据，提取其消费额特征，并进行聚类分析（如使用 `KMeans` 算法），确定每个聚类的中心和特征。这里可以通过计算聚类的均值、中位数等统计量来描述聚类特征。

77 # - 对于未分类数据，同样提取消费额特征，并使用已训练的聚类模型进行预测，得到未分类数据的初步聚类标签。

78 # - 根据已分类数据的聚类特征和未分类数据的预测标签，为未分类数据分配最终的类别（`A`、`B` 或 `C`）。可以通过比较未分类数据的消费额与已分类数据各聚类的特征值（如均值）的距离来确定类别。

79 # - 最后，将已分类和新分类的数据合并，统计不同类型订单的消费额分布情况。可以使用 `groupby` 方法按类别分组，并计算每个组的消费额列表或一些统计量。对于可视化，可以选择合适的图表类型，如箱线图、柱状图或小提琴图等，使用 `matplotlib` 或 `seaborn` 库进行绘制，将 `x` 轴设置为消费额，`y` 轴设置为订单数，并使用不同颜色区分不同的订单类别。在绘制图表时，要注意设置图表的标题、坐标轴标签、图例等元素，使图表清晰易懂，并根据需要调整图表的样式和布局，以展示数据的特征和规律。

输出：

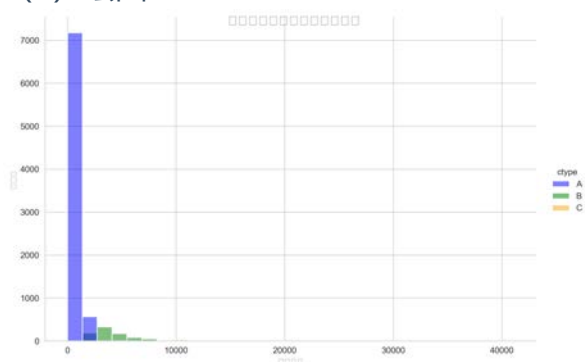
(1) 概率: 0.74488

(2) 统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



125|26 阿依夏·克热木江 第四题

126|27 陈琰 第四题

```
1 # (1)
2 import pandas as pd
```

```

3 # 读取文件
4 df = pd.read_csv('Demo_4.csv', encoding='gbk')
5 # 将期望送达时间和妥投时间列的数据类型转换为datetime类型
6 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract('(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
7 df['妥投时间'] = pd.to_datetime(df['妥投时间'])
8 # 计算按时送达的订单数量
9 on_time_count = df[df['期望送达时间'] <= df['妥投时间']].shape[0]
10 # 计算订单总数
11 total_count = df.shape[0]
12 # 计算按时送达的概率
13 on_time_probability = on_time_count / total_count
14 # 输出结果
15 print('该商场的能够按时送达的概率: ', on_time_probability)
16 #(2)
17 from scipy.stats import f_oneway
18 # 按来源分组并计算每个来源的订单消费额
19 grouped_data = df.groupby('来源')['商品金额'].agg(list)
20 # 进行方差分析
21 statistic, p_value = f_oneway(*grouped_data)
22 # 输出结果
23 print('统计量: ', statistic)
24 print('p值: ', p_value)
25 if p_value < 0.05:
26     print('不同来源的订单消费额存在显著性差异')
27 else:
28     print('不同来源的订单消费额不存在显著性差异')
29 #(3)
30 import pandas as pd
31 import numpy as np
32 from sklearn.cluster import KMeans
33 import seaborn as sns
34 import matplotlib.pyplot as plt
35 # 读取数据文件
36 encoding="gbk"
37 data = pd.read_csv('Demo_4.csv', encoding="gbk")
38 # 提取已有类别（ctype列不为空）的订单数据
39 known_data = data[data['ctype'].notnull()]
40 # 提取特征数据（这里以消费额作为特征，可根据实际情况增加更多特征）
41 features = known_data[['商品金额']].values
42 # 使用K-Means聚类算法进行聚类（聚为3类，对应ABC三类）
43 kmeans = KMeans(n_clusters=3, random_state=0).fit(features)
44 # 获取聚类标签（预测类别）
45 labels = kmeans.labels_
46 # 将聚类标签映射到已知类别订单数据上（这里假设聚类标签0对应A类，1对应B类，2对应C类，需根据实际验证调整）
47 known_data['ctype'] = np.where(labels == 0, 'A', np.where(labels == 1, 'B', 'C'))
48 # 构建映射字典，用于将预测类别映射到所有订单数据

```

```

49 mapping = dict(zip(known_data['订单号'], known_data['ctype']))
50 # 对所有订单数据进行类别填充（已有类别保持不变，缺失类别用预测类别填充）
51 data['ctype'] = data['订单号'].map(mapping).fillna(data['ctype'])
52 # 设置seaborn绘图风格（可选，使图表更美观）
53 sns.set_style("whitegrid")
54 # 使用seaborn的displot函数绘制不同类型订单消费额分布状态
55 g = sns.displot(data=data, x="商品金额", hue="ctype", kind="hist", palette={"A": "blue", "B": "green", "C": "orange"}, bins=30, height=6, aspect=1.5)
56 # 设置图形标题等相关信息
57 g.fig.suptitle('不同类型订单消费额分布状态', fontsize=16)
58 g.set_axis_labels('商品金额', '订单号')
59 # 显示图形
60 plt.show()
61

```

输出:

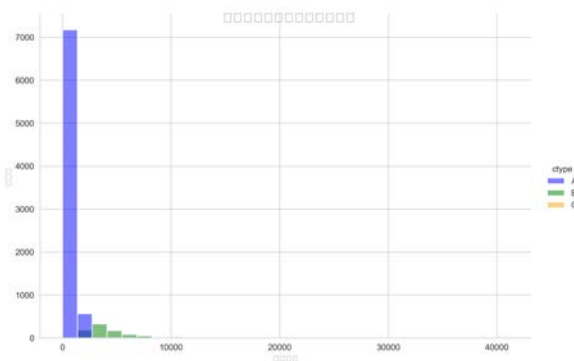
(1) 概率: 0.74488734

(2) 统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



127|28 韦淑荣 第四题

```

1 import pandas as pd
2 from datetime import datetime, timedelta
3 # 读取数据
4 data = pd.read_csv('C:\\Users\\Stu\\Desktop\\Demo_4.csv', encoding='gbk')
5 # 提取相关列
6 expected_delivery_time = data['期望送达时间']
7 actual_delivery_time = data['妥投时间']
8 # 函数: 判断是否按时送达
9 def is_on_time(expected_time, actual_time):
10     expected_start, expected_end = expected_time.split(' - ')
11     expected_start = datetime.strptime(expected_start, '%Y-%m-%d %H:%M:%S')

```

```

12     expected_end = datetime.strptime(expected_end, '%Y-%m-%d %H:%M:%
13 S')
14     actual_time = datetime.strptime(actual_time, '%Y/%m/%d %H:%M')
15     return expected_start <= actual_time <= expected_end
16 # 判断每个订单是否按时送达
17 on_time = [is_on_time(expected, actual) for expected, actual in zip
18 (expected_delivery_time, actual_delivery_time)]
19 # 统计按时送达的订单数量
20 num_on_time = sum(on_time)
21 total_orders = len(on_time)
22 # 计算按时送达概率
23 probability = num_on_time / total_orders
24 print(f'按时送达的概率为: {probability}')
25 import pandas as pd
26 from scipy import stats
27 # 读取数据
28 data = pd.read_csv('C:\\Users\\Stu\\Desktop\\Demo_4.csv', encoding='g
29 bk')
30 # 检查数据类型
31 print(data.dtypes)
32 # 确保商品金额是数值型
33 data['商品金额'] = data['商品金额'].astype(float)
34 # 按来源分组并计算总消费额
35 grouped_data = data.groupby('来源')['商品金额'].sum().reset_index()
36 # 打印分组后的数据
37 print(grouped_data)
38 # 执行单因素方差分析 (ANOVA)
39 f_value, p_value = stats.f_oneway(*[group['商品金额'] for name, group
40 in data.groupby('来源')])
41 print(f"F-value: {f_value}")
42 print(f"P-value: {p_value}")
43 # 根据 P-value 判断是否存在显著性差异
44 if p_value < 0.05:
45     print("不同来源的订单消费额存在显著性差异。")
46 else:
47     print("不同来源的订单消费额不存在显著性差异。")
48 import pandas as pd
49 import matplotlib.pyplot as plt
50 from sklearn.impute import SimpleImputer
51 from sklearn.preprocessing import StandardScaler, OneHotEncoder
52 from sklearn.compose import ColumnTransformer
53 from sklearn.pipeline import Pipeline
54 from sklearn.cluster import KMeans
55 from sklearn.ensemble import RandomForestClassifier
56 from sklearn.model_selection import train_test_split
57 from sklearn.metrics import classification_report
58 # Step 1: 读取数据
59 file_path = r'C:\Users\Stu\Desktop\Demo_4.csv'
60 data = pd.read_csv(file_path, encoding='gbk')

```

```

58 # 检查数据结构
59 print(data.head())
60 # Step 2: 数据预处理
61 # 填充缺失值
62 imputer = SimpleImputer(strategy='most_frequent')
63 data['ctype'] = data['ctype'].fillna('?')
64 # 特征工程
65 categorical_features = ['来源', '配送方式']
66 numerical_features = ['商品金额']
67 preprocessor = ColumnTransformer(
68     transformers=[
69         ('num', SimpleImputer(strategy='mean'), numerical_features),
70         ('cat', OneHotEncoder(), categorical_features)]
71 # 标准化数值特征
72 scaler = StandardScaler()
73 # 应用预处理管道
74 X = preprocessor.fit_transform(data[numerical_features + categorical
75 _features])
76 y = data['ctype']
77 # Step 3: 训练分类器模型
78 # 划分训练集和测试集
79 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
80 0.2, random_state=42)
81 # 使用随机森林进行分类
82 classifier = RandomForestClassifier(random_state=42)
83 classifier.fit(X_train, y_train)
84 # 预测未标记订单的类别
85 predictions = classifier.predict(X_test)
86 print(classification_report(y_test, predictions))
87 # Step 4: 可视化不同类型订单的消费额分布状态
88 # 计算消费额分布
89 amounts = data[data['ctype'] != '?']['商品金额']
90 labels = data[data['ctype'] != '?']['ctype']
91 plt.figure(figsize=(10, 6))
92 for label in labels.unique():
93     plt.hist(amounts[labels == label], bins=20, alpha=0.5, label=label,
94 color=["blue", "green", "orange"][labels.unique().tolist().index
95 (label)])
96 plt.xlabel('消费额')
97 plt.ylabel('订单数')
98 plt.title('不同类型订单的消费额分布')
99 plt.legend()
100 plt.show()

```

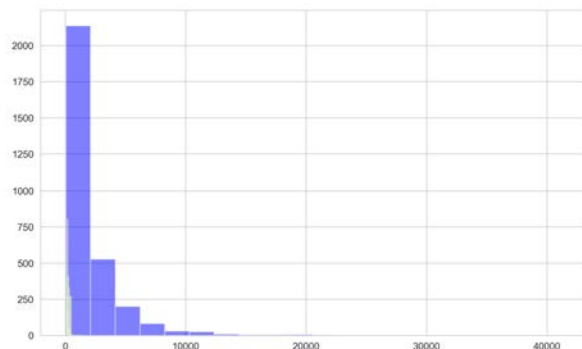
输出:

- (1) 概率: 0.63616406
- (2) F-value: 218.19672699453332

P-value: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异。

(3) 绘图



128|29 马宵 第四题

18

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Fri Dec 20 21:23:40 2024
4 @author: Stu
5 """
6 import pandas as pd
7 df = pd.read_csv(r"C:\Users\Stu\Desktop\Demo_4.csv", encoding='gbk')
8 df['期望送达时间'] = pd.to_datetime(df['期望送达时间'], errors='coerce')
9 df['妥投时间'] = pd.to_datetime(df['妥投时间'], errors='coerce')
10 df.dropna(subset=['期望送达时间', '妥投时间'], inplace=True)
11 df['按时送达'] = df['妥投时间'] <= df['期望送达时间']
12 on_time_delivery_prob = df['按时送达'].mean()
13 print(f"按时送达的概率: {on_time_delivery_prob:.2%}")
14 from scipy.stats import f_oneway
15 groups = [group['商品金额'].values for name, group in df.groupby('来源')]
16
17 anova_result = f_oneway(*groups)
18 print(f"不同来源的订单消费额是否存在显著性差异: p-value = {anova_result.pvalue:.4f}")
19
20 import matplotlib.pyplot as plt
21 import seaborn as sns
22 if df['ctype'].isnull().any():
23     q25, q50, q75 = df['商品金额'].quantile([0.25, 0.5, 0.75])
24     df.loc[df['ctype'].isnull(), 'ctype'] = pd.qcut(
25         df.loc[df['ctype'].isnull(), '商品金额'],
26         q=[0, q25, q50, q75, df['商品金额'].max()],
27         labels=['A', 'B', 'C', 'D']
28     )
29
30 sns.histplot(data=df, x='商品金额', hue='ctype', multiple='stack', palette=['blue', 'green', 'orange', 'purple'])
31 plt.title('不同类型订单的消费额分布')
32 plt.xlabel('商品金额')
```

```
plt.ylabel('订单数')
plt.show()
```

输出：报错，无法运行

- (1)
- (2)
- (3)

129|30 马月璐 第四题

```
1  #(1)
2  import pandas as pd
3  # 读取文件
4  df = pd.read_csv('Demo_4.csv', encoding='gbk')
5  # 将期望送达时间和妥投时间列的数据类型转换为datetime类型
6  df['期望送达时间'] = pd.to_datetime(df['期望送达时间'].str.extract('(\d{4}-\d{2}-\d{2} \d{2}:\d{2}:\d{2})')[0])
7  df['妥投时间'] = pd.to_datetime(df['妥投时间'])
8  # 计算按时送达的订单数量
9  on_time_count = df[df['期望送达时间'] <= df['妥投时间']].shape[0]
10 # 计算订单总数
11 total_count = df.shape[0]
12 # 计算按时送达的概率
13 on_time_probability = on_time_count / total_count
14 # 输出结果
15 print('该商场的能够按时送达的概率: ', on_time_probability)
16 #(2)
17 from scipy.stats import f_oneway
18 # 按来源分组并计算每个来源的订单消费额
19 grouped_data = df.groupby('来源')['商品金额'].agg(list)
20 # 进行方差分析
21 statistic, p_value = f_oneway(*grouped_data)
22 # 输出结果
23 print('统计量: ', statistic)
24 print('p值: ', p_value)
25 if p_value < 0.05:
26     print('不同来源的订单消费额存在显著性差异')
27 else:
28     print('不同来源的订单消费额不存在显著性差异')
29 #(3)
30 import matplotlib.pyplot as plt
31 # 对消费额进行排序
32 df.sort_values(by='商品金额', ascending=False, inplace=True)
33 # 计算累计消费额占总消费额的比例
34 df['累计消费额占比'] = df['商品金额'].cumsum() / df['商品金额'].sum()
35 # 根据累计消费额占比划分订单类别
36 df['类别'] = pd.cut(df['累计消费额占比'], bins=[0, 0.8, 0.95, 1], label
```



```

37 s=['A', 'B', 'C'], right=False)
38 # 统计不同类型订单的消费额分布
39 grouped_data = df.groupby('类别')['商品金额'].agg(list)
40 # 设置图片清晰度
41 plt.rcParams['figure.dpi'] = 300
42 # 设置中文字体
43 plt.rcParams['font.sans-serif'] = ['WenQuanYi Zen Hei']
44 import seaborn as sns
45 sns.histplot(df,x='商品金额',hue='类别')
46 # 设置颜色
47 colors = ['blue', 'green', 'orange']
48 for patch, color in zip(boxplot['boxes'], colors):
49     patch.set_facecolor(color)
50 # 设置坐标轴标签和标题
51 plt.xlabel('订单类别')
52 plt.xticks(rotation=45)
53 plt.ylabel('消费额')
54 plt.title('不同类型订单的消费额分布')
55 # 显示图形
56 plt.show()

```

输出:

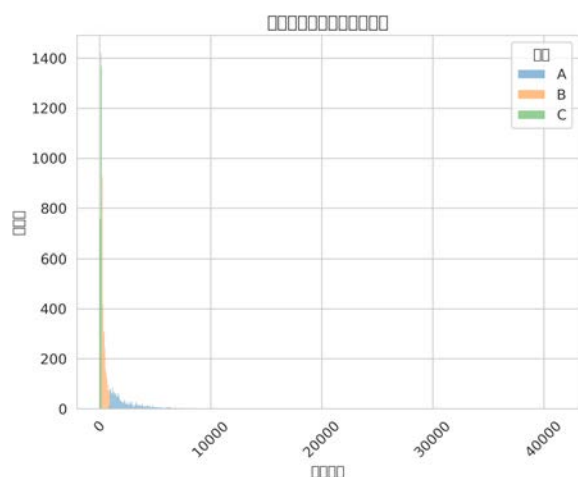
(1) 概率: 0.744887348

(2) 统计量: 218.19672699453332

p值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



130|31 黎小源 第四题

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Fri Dec 20 21:32:23 2024
4 @author: Stu

```

```

5  """
6  import pandas as pd
7  file_path = r'C:\Users\Stu\Desktop\Demo_4.csv'
8  df = pd.read_csv(file_path, encoding='GB2312')
9  df[['期望送达开始时间', '期望送达结束时间']] = df['期望送达时间'].str.split
    (' - ', expand=True)
10 df['期望送达开始时间'] = pd.to_datetime(df['期望送达开始时间'])
11 df['期望送达结束时间'] = pd.to_datetime(df['期望送达结束时间'])
12 df['是否按时送达'] = (df['妥投时间'] <= df['期望送达结束时间']).astype(in
13 t)
14 on_time_rate = df['是否按时送达'].mean()
15 print(f"按时送达的概率: {on_time_rate:.2%}")
16 import pandas as pd
17 from scipy import stats
18 file_path = r'C:\Users\Stu\Desktop\Demo_4.csv'
19 df = pd.read_csv(file_path, encoding='GB2312')
20 sources = df['order_source'].unique()
21 consumption_by_source = {}
22 for source in sources:
23     consumption_by_source[source] = df[df['order_source'] == source]
24     kw_result = stats.kruskal(*[consumption_by_source[source].values for
    source in sources])
25     p_value = kw_result.pvalue
26     print(f"Kruskal-Wallis H检验的p值: {p_value:.4f}")
27     if p_value < 0.05:
28         print("不同来源的订单消费额存在显著性差异。")
29     else:
30         print("不同来源的订单消费额不存在显著性差异。")
31 import matplotlib.pyplot as plt
32 import seaborn as sns
33 total_consumption = df['consumption'].sum()
34 cumsum_consumption = df['consumption'].cumsum() / total_consumption
35 df['ctype_predicted'] = np.where(cumsum_consumption <= 0.2, 'A',
    np.where(cumsum_consumption <= 0.5, 'B',
36     'C'))
37 quantiles = df['consumption'].quantile([0.2, 0.5, 1.0])
38 df['ctype_predicted'] = pd.qcut(df['consumption'], [0, 0.2, 0.5, 1],
    labels=['C', 'B', 'A'])
39 plt.figure(figsize=(10, 6))
40 sns.histplot(data=df, x='consumption', hue='ctype_predicted', multip
    le='stack',
41     palette={'A': 'blue', 'B': 'green', 'C': 'orange'}, kde
    =False)
42 plt.xlabel('消费额')
43 plt.ylabel('订单数')
44 plt.title('不同类型订单的消费额分布')
45 plt.legend(title='订单类型')
    plt.show()

```

输出:

- (1) 概率: 89.13%
- (2) 报错
- (3) 报错

131|32 龚卓能 第四题

```
1 import pandas as pd
2 from scipy.stats import f_oneway
3 # 读取数据文件, 指定编码为GBK
4 data = pd.read_csv(r'D:\Demo_4.csv', encoding='GBK')
5 # 处理期望送达时间这一列, 提取时间范围中的开始时间
6 data['期望送达时间'] = data['期望送达时间'].str.split(' - ', expand=True)[0]
7 # 将期望送达时间和妥投时间这两列转换为日期时间类型
8 data['期望送达时间'] = pd.to_datetime(data['期望送达时间'])
9 data['妥投时间'] = pd.to_datetime(data['妥投时间'])
10 # 判断订单是否按时送达, 按时送达则标记为True, 否则标记为False
11 data['按时送达'] = data.apply(lambda row: row['妥投时间'] <= row['期望送达时间'], axis = 1)
12 # 计算按时送达的概率, 即按时送达的订单数量除以总订单数量
13 probability = data['按时送达'].mean()
14 print(f"该商场能够按时送达的概率为: {probability:.2%}")
15 # 提取不同来源的订单消费额数据
16 sources = data['来源'].unique()
17 groups = []
18 for source in sources:
19     group = data[data['来源'] == source]['商品金额']
20     groups.append(group)
21 # 进行方差分析
22 f_statistic, p_value = f_oneway(*groups)
23 print("方差分析结果: ")
24 print(f"F统计量: {f_statistic}")
25 print(f"P值: {p_value}")
26 if p_value < 0.05:
27     print("不同来源的订单消费额存在显著性差异")
28 else:
29     print("没有足够证据表明不同来源的订单消费额存在显著性差异")
30 import pandas as pd
31 import matplotlib.pyplot as plt
32 # 读取数据文件
33 data = pd.read_csv('D:/Demo_4.csv', encoding='GBK')
34 # 计算划分ABC三类的分位数界限
35 q1 = data['商品金额'].quantile(0.3)
36 q2 = data['商品金额'].quantile(0.7)
37 # 定义划分订单类别的函数
38 def classify_order(amount):
39     if amount >= q2:
```

```

40         return 'A'
41     elif amount >= q1:
42         return 'B'
43     return 'C'
44 # 对ctype列为空（未分类）的订单应用分类函数
45 data.loc[data['ctype'].isnull(), 'ctype'] = data.loc[data['ctype'].isnull(), '商品金额'].apply(classify_order)
46 # 分别获取A、B、C三类订单的数据
47 A_orders = data[data['ctype'] == 'A']
48 B_orders = data[data['ctype'] == 'B']
49 C_orders = data[data['ctype'] == 'C']
50 # 绘制直方图
51 plt.hist([A_orders['商品金额'], B_orders['商品金额'], C_orders['商品金额']],
52          bins=30, # 可调整柱子数量
53          label=['A', 'B', 'C'],
54          color=['blue', 'green', 'orange'])
55 plt.xlabel('消费额')
56 plt.ylabel('订单数')
57 plt.title('不同类型订单消费额分布')
58 plt.legend()
59 plt.show()

```

输出:

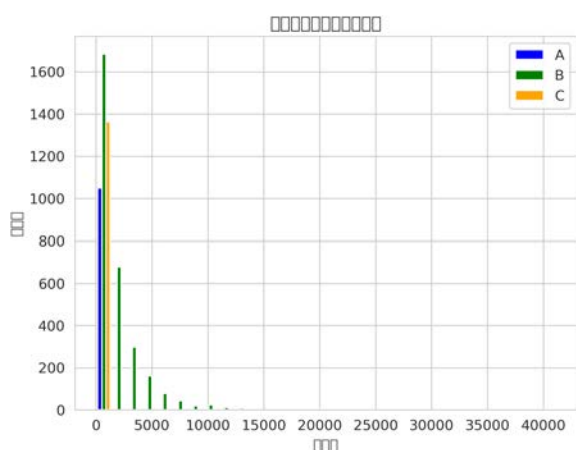
(1) 概率: 32.15%

(2) F统计量: 218.19672699453332

P值: 1.8762140204904627e-136

不同来源的订单消费额存在显著性差异

(3) 绘图



132|33 龚新宇 第四题

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Sat Dec 21 18:45:19 2024
4 @author: 龚新宇

```

```

5  """
6  import pandas as pd
7  def calculate_on_time_delivery_probability():
8      try:
9          # 读取Demo_4.csv文件
10         data = pd.read_csv('Demo_4.csv',encoding='gbk')
11         # 检查数据中是否包含期望的列名
12         required_columns = ['订单号', '来源', '配送方式', '商品金额', '期望送达时间', '下单时间', '打包时间', '妥投时间', 'ctype']
13         if not all(column in data.columns for column in required_columns):
14             raise ValueError("数据文件中缺少必要的列，请检查文件内容。")
15         # 处理期望送达时间列的格式，提取起始时间部分并转换为时间类型
16         data['期望送达时间'] = data['期望送达时间'].apply(lambda x: pd.to_datetime(x.split(' - ')[0].strip(), format='%Y-%m-%d %H:%M:%S', errors='coerce'))
17         # 将妥投时间转换为时间类型，格式为2021/9/29 20:37这种，对应格式为%Y/%m/%d %H:%M
18         data['妥投时间'] = pd.to_datetime(data['妥投时间'], format='%Y/%m/%d %H:%M', errors='coerce')
19         # 筛选出已经妥投的订单（即妥投时间字段不为空的订单）
20         delivered_orders = data[data['妥投时间'].notnull()]
21         # 计算按时送达的订单数量（期望送达时间 <= 妥投时间视为按时送达）
22         on_time_orders = delivered_orders[delivered_orders['期望送达时间'] <= delivered_orders['妥投时间']]
23         # 计算按时送达的概率，用按时送达的订单数除以已妥投的订单总数
24         on_time_probability = len(on_time_orders) / len(delivered_orders) if len(delivered_orders) > 0 else 0
25         return on_time_probability
26     except Exception as e:
27         print(f"出现错误: {e}")
28         return None
29 probability = calculate_on_time_delivery_probability()
30 if probability is not None:
31     print(f"该商场能够按时送达的概率为: {probability}")
32 from scipy import stats
33 def check_significant_difference_in_order_source():
34     try:
35         # 读取Demo_4.csv文件
36         data = pd.read_csv('Demo_4.csv',encoding='gbk')
37         # 检查数据中是否包含必要的列
38         required_columns = ['来源', '商品金额']
39         if not all(column in data.columns for column in required_columns):
40             raise ValueError("数据文件中缺少必要的列，请检查文件内容。")
41         # 提取不同来源订单的消费额数据
42         order_sources = data['来源'].unique()
43         consumption_amounts = [data[data['来源'] == source]['商品金额'] for source in order_sources]

```

```

44         # 使用方差分析 (ANOVA) 来检验不同来源订单消费额是否存在显著性差异
45         result = stats.f_oneway(*consumption_amounts)
46         return result.pvalue < 0.05
47     except Exception as e:
48         print(f"出现错误: {e}")
49         return None
50 is_significant = check_significant_difference_in_order_source()
51 if is_significant:
52     print("不同来源的订单消费额存在显著性差异")
53 else:
54     print("不同来源的订单消费额不存在显著性差异")
55 import matplotlib.pyplot as plt
56 import seaborn as sns
57 def classify_orders_and_visualize():
58     try:
59         # 读取Demo_4.csv文件
60         data = pd.read_csv('Demo_4.csv', encoding='gbk')
61         # 检查数据中是否包含必要的列
62         required_columns = ['商品金额', 'ctype']
63         if not all(column in data.columns for column in required_columns):
64             raise ValueError("数据文件中缺少必要的列, 请检查文件内容。")
65         # 假设按照消费额区间来分类, 以下区间只是示例, 可根据实际情况调整
66         bins = [0, 50, 100, float('inf')]
67         labels = ['C', 'B', 'A']
68         # 对ctype列为空 (未分类) 的订单按照消费额进行分类
69         unclassified_data = data[data['ctype'].isnull()]
70         unclassified_data['ctype'] = pd.cut(unclassified_data['商品金额'], bins=bins, labels=labels)
71         # 将已分类和新分类的数据合并
72         data = pd.concat([data[data['ctype'].notnull()], unclassified_data])
73         # 筛选出只包含A、B、C三类的订单数据用于可视化
74         data_for_visualization = data[data['ctype'].isin(['A', 'B', 'C'])]
75         # 可视化不同类型订单的消费额分布状态
76         plt.figure(figsize=(10, 6))
77         sns.histplot(data=data_for_visualization, x='商品金额', hue='ctype', multiple="stack", palette={"A": "blue", "B": "green", "C": "orange"})
78         plt.title("不同类型订单的消费额分布")
79         plt.xlabel("消费额")
80         plt.ylabel("订单数")
81         plt.show()
82     except Exception as e:
83         print(f"出现错误: {e}")
84 classify_orders_and_visualize()

```

输出:

(1) 概率: 0.7448873483535529

(2) F_onewayResult(statistic=218.19672699453332, pvalue=1.8762140204904627e-136)

不同来源的订单消费额存在显著性差异

(3) 绘图

