

出血性脑卒中临床智能诊疗优化模型

摘要

在现代医学临床技术中，通过分析临床数据，我们可以对血性脑卒中患者的血肿扩张风险、血肿周围水肿的发生和演进规律进行评估，并结合临床和影像信息来预测患者的临床预后。本研究采用了逻辑回归模型，并结合决策树和多项式回归的方法，建立了随机森林模型，通过考虑临床实际情况对模型进行优化，并提出相应的优化方案，最终得到了适用于一般情况的普适性模型。

针对问题一，我们首先根据提供的数据判断患者发病后 48 小时内是否发生了血肿扩张事件，并同时记录血肿扩张的发生情况。我们根据发病到首次影像检查的时间间隔和后续影像检查的时间间隔来确定是否在 48 小时内发生了血肿扩张。考虑到数据的连续性和稳定性，我们采用贪心算法的思想来判断是否发生了血肿扩张，并通过建立逻辑回归模型来预测所有患者发生血肿扩张的概率。逻辑回归模型通过使用 Sigmoid 函数将线性预测转换为概率值。

针对问题二，为了计算真实值与拟合曲线之间的残差，我们根据患者的水肿体积和发病到各个随访时间点的时间间隔构建了进展曲线，以探索患者水肿体积随时间变化的个体差异，并建立不同人群的进展曲线。根据构建的水肿体积随时间进展曲线，我们计算患者的真实值与曲线之间的残差。为了分析水肿体积进展模式的影响，我们采用了多项式回归和随机森林模型。考虑到患者个体差异，患者脑补血肿体积随时间进展可能不符合简单的多项式关系，因此我们最终使用随机森林模型进行曲线拟合，并作为最终预测结果的模型。

针对问题三，为了平衡个人史、疾病史和发病相关因素以及首次影像结果之间的关系，我们分别使用逻辑回归模型、随机森林算法和 K 临近算法对患者的 mRS 评分进行预测。根据已有的临床和治疗结果，我们预测了所有含有随访影像检查的患者在 90 天时的 mRS 评分。通过逻辑回归模型和分析预测结果，我们考虑到随机森林模型和 K 临近算法之间的精度，建议在对预后相关结果和因素进行准确预测时，可以适当采用随机森林模型以提高预测精度。

关键词：随机森林模型 逻辑回归模型 平均相对误差 多项式回归

目 录

摘 要.....	2
一、问题重述.....	4
1.1 问题背景	4
1.2 问题提出	4
二、模型假设.....	6
2.1 模型的基本	6
三、模型的符号说明.....	7
四、问题一的模型建立与求解	8
4.1 任务 a 的分析与处理	8
4.1.1 任务 a 问题分析.....	8
4.1.2 任务 a 特征提取及数据分析.....	8
4.1.3 算法设计.....	9
4.2 任务 b 的分析与处理	9
4.2.1 任务 b 问题分析	9
4.2.3 任务 b 模型的建立	14
4.2.4 任务 b 模型结果	14
五、问题二的模型建立与求解	15
5.1 任务 a 的分析与处理	15
5.2 任务 b 的分析与处理	17
5.3 任务 c 的分析与处理.....	19
5.4 任务 d 的分析与处理	20
六、问题三的模型建立和求解	21
6.1 问题 a 的分析与处理	21
七、参考文献.....	23
附录（代码）	24

一、问题重述

1.1 问题背景

出血性脑卒中是指脑部血管破裂导致的脑内出血。它与缺血性脑卒中(脑血管阻塞)不同,属于一种较为严重的疾病。出血性脑卒中的常见原因包括高血压、脑动脉瘤、脑血管畸形和出血性疾病等。当血管破裂时,血液会进入脑组织,压迫周围健康脑组织并引起炎症反应,造成神经功能障碍。出血性脑卒中的症状可能会突然出现,如剧烈头痛、意识丧失、肢体无力、言语困难、视力改变等。如果怀疑发生了出血性脑卒中,应立即就医。治疗出血性脑卒中的方法包括控制高血压、手术修复破裂的血管、减轻颅内压力以及进行康复治疗等。同时,预防出血性脑卒中的关键是保持健康的生活方式,如保持血压正常、戒烟限酒、均衡饮食、适量运动、定期体检等。

出血性脑卒中指非外伤性脑实质内血管破裂引起的脑出血,占全部脑卒中发病率的10-15%。其病因复杂,通常因脑动脉瘤破裂、脑动脉异常等因素,导致血液从破裂的血管涌入脑组织,从而造成脑部机械性损伤,并引发一系列复杂的生理病理反应。出血性脑卒中起病急、进展快,预后较差,急性期内病死率高达45-50%,约80%的患者会遗留较严重的神经功能障碍,为社会及患者家庭带来沉重的健康和经济负担。因此,发掘出血性脑卒中的发病风险,整合影像学特征、患者临床信息及临床诊疗方案,精准预测患者预后,并据此优化临床决策具有重要的临床意义。

出血性脑卒中后,血肿范围扩大是预后不良的重要危险因素之一。在出血发生后的短时间内,血肿范围可能因脑组织受损、炎症反应等因素逐渐扩大,导致颅内压迅速增加,从而引发神经功能进一步恶化,甚至危及患者生命。因此,监测和控制血肿的扩张是临床关注的重点之一。此外,血肿周围的水肿作为脑出血后继发性损伤的标志,在近年来引起了临床广泛关注。血肿周围的水肿可能导致脑组织受压,进而影响神经元功能,使脑组织进一步受损,进而加重患者神经功能损伤。综上所述,针对出血性脑卒中后的两个重要关键事件,即血肿扩张和血肿周围水肿的发生及发展,进行早期识别和预测对于改善患者预后、提升其生活质量具有重要意义。

事实上,出血性脑卒中出现的危险因素很多,我们可以联合患者个人信息、治疗方案和预后等数据,构建智能诊疗模型,明确因素,实现精准预测。相关研究成果及科学依据将能够进一步应用于临床实践,为改善出血性脑卒中患者预后作出贡献。

1.2 问题提出

本文以上述分析为背景,主要考虑以下几个问题

(1) 血脉扩张风险相关因素,根据血肿体积前后变化,后续检查比首次检查绝对体积增加或者相对体积增加判断是否发生血肿扩张。结合发病到首次影像

时间间隔和后续影像检查时间间隔判断当前影像检查是否在发病 48 小时内。并以此为目标变量，结合相关信息预测患者发生血肿扩张的概率。

(2) 结合水肿体积和重复检查时间点，构建一条全体患者水肿体积随时间进展曲线，计算真实值和所拟合曲线之间存在的残差。探索患者之间的个体差异，构建水肿体积随时间进展曲线，计算残差。观察残差的分布情况，采用方根误差或平均绝对误差评估拟合曲线的精度，较小的残差表示模型拟合效果较好。

(3) 通过拟合曲线分析不同的治疗方法对水肿体积进展模式的影响，我们可以通过逻辑回归模型分析血肿体积、水肿体积以及治疗方法三者之间的关系进一步探索关联关系。

(4) 基于前述模型以及算法分析，可以预测 mRS 结果，结合所有已知临床、治疗等的影像结果，预测所有含随访影像检查患者的 mRS 评分，构建出预测模型，并分析患者预后的影像特征等关联关系。

二、模型假设

2.1 模型的基本

在本文的研究中心，采用均匀先验理论所需要进行的问题假设：

- （1）假设所有患者均发生血肿扩张，即所有患者的后续检查比首次检查绝对体积增加大于等于 6mL 或相对体积增加大于等于 33%。
- （2）假设患者水肿体积随时间进展模式的差异均相同。
- （3）假设所有含随访影像检查的患者 90 天 mRS 评分分布符合正态分布

三、模型的符号说明

表 3-1 模型符号说明

序号	符号	说明
1	$g, grad$	图像平均梯度
2	$c, cons$	图像对比度
3	μ	图像均值
4	ε	拟合残差
5	h	时间
6	σ	图像均方差
7	Reg	图像子窗格
8	s	信号强度特征
9	y	实际观测值
10	\hat{y}	预测值

四、问题一的模型建立与求解

4.1 任务 a 的分析与处理

4.1.1 任务 a 问题分析

任务 a 要求判断患者是否在 48 小时内发生血肿扩张事件，若发生，则记为 1，同时记录从首次随访时间到发生血肿的时间；否则记为 0，同时发生血肿的时间记为 0 小时。

考虑到需要表 2 给出的数据种类众多，因而先提取数据。除此之外，由于需要记录发生血肿扩张的时间，因此需要从最近的随访时间开始遍历，故使用贪心算法的思想进行求解。

4.1.2 任务 a 特征提取及数据分析

首先需要对表 1 的数据进行特征提取，得到两列原始数据，注意到后续需要对数据时间进行操作，因此将表中数据转为时间戳，如表 4-1 所示。

表 4-1 表 1 特征提取后的前 5 行内容

入院首次影像检查流水号	发病到首次影像检查时间间隔
020161212002136	0 days 02:30:00
120160406002131	0 days 03:00:00
220160413000006	0 days 02:00:00
320161215001667	0 days 01:00:00
420161222000978	0 days 05:00:00

接着判断表 1 空数据的基本情况如表 4-2 所示

表 4-2 表 1 空数据的基本情况

入院首次影像检查流水号	0
发病到首次影像检查时间间隔	0
dtype:	int64
Null columns:	Index([], dtype='object')

通过上表可知表 1 特征提取后不包含空数据，结合表 1 两列特征的实际意义，也可得出该结论。

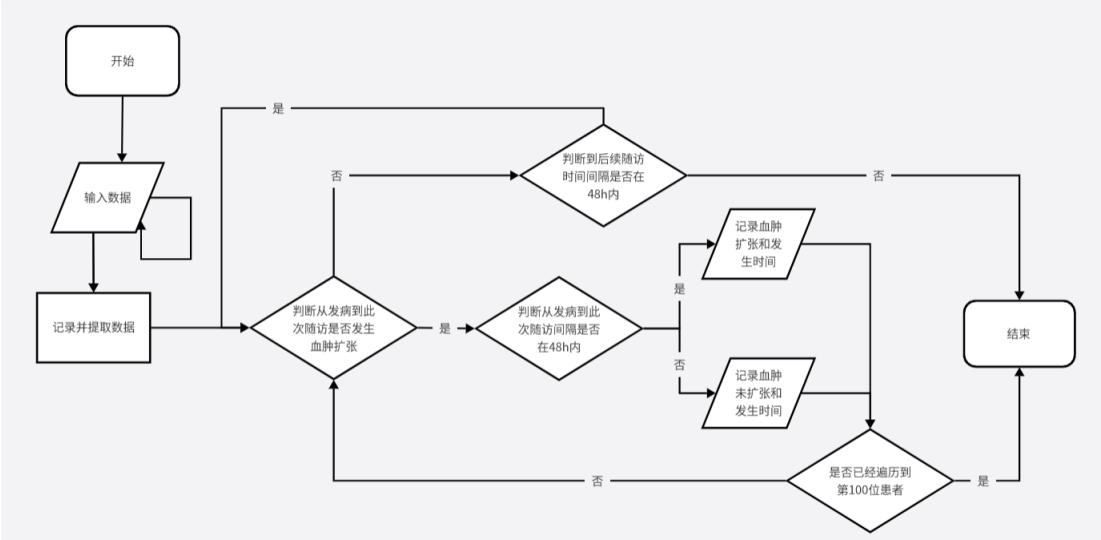
接着对表 2 的数据进行特征提取，得到问题所需的特征，提取结果如表 4-3 所示，不难发现，提取后的文件中包含大量空数据，考虑到空数据的实际意义和任务需求，故对空数据不做处理。

判断血肿发生的时间是否在 48h 内，需要查阅附表 1。在实现上，通过表 2 的流水号定位附表 1 对应的时间即可，因此问题 a 的数据处理已完成。

4.1.3 算法设计

任务 a 不仅需要判断是否发生血肿，还要记录血肿发生时间，因此考虑使用贪心算法思想：对于每位患者，优先从最近的随访时间开始判断是否发生血肿，如果发生血肿，再判断时间是否在 48 小时内。如果本次发生血肿的时间已经超过 48h，说明后续的随访必然会超过 48h。如果本次的随访时间在 48h 内，说明本次随访时间一定是最先发生随访的时间，即需要记录的时间，具体过程如图 4-1 所示。

图 4-1 任务 a 算法流程图



4.2 任务 b 的分析与处理

4.2.1 任务 b 问题分析

从问题 b 的描述中可知，我们要以是否发生血肿扩张事件为目标变量，基于患者的个人情况构建模型，来预测患者发生血肿扩张的概率。在这里我们通过逻辑回归的方法来预测概率。

4.2.2 任务 b 数据分析

首先对表 1 给出的原始数据进行数据分析，如表 4-4 所示：

表 4-3 表 1 原始数据的基本信息

#	列名	非空数据个数	数据类型
1	90 天 mRS	100	float64
2	数据集划分	160	object
3	入院首次影像检查流水号	160	int64
4	年龄	160	int64
5	性别	160	object

6	脑出血前 mRS 评分	160	int64
7	高血压病史	160	int64
8	卒中病史	160	int64
9	糖尿病史	160	int64
10	房颤史	160	int64
11	冠心病史	160	int64
12	吸烟史	160	int64
13	饮酒史	160	int64
14	发病到首次影像检查时间 间隔	160	float64
15	血压	160	object
16	脑室引流	160	int64
17	止血治疗	160	int64
18	降颅压治疗	160	int64
19	降压治疗	160	int64
20	镇静、镇痛治疗	160	int64
21	止吐护胃	160	int64
22	营养神经	160	int64

根据题意，前三行数据不需要作为模型的输入，因此需要处理的列为“性别”和“血压”。将“性别”替换为布尔型，“血压”扩充为“高压”和“低压”。接下来对处理后的表 1 进行数据分析，先对统计各字段数据的类别如下图 4-2。

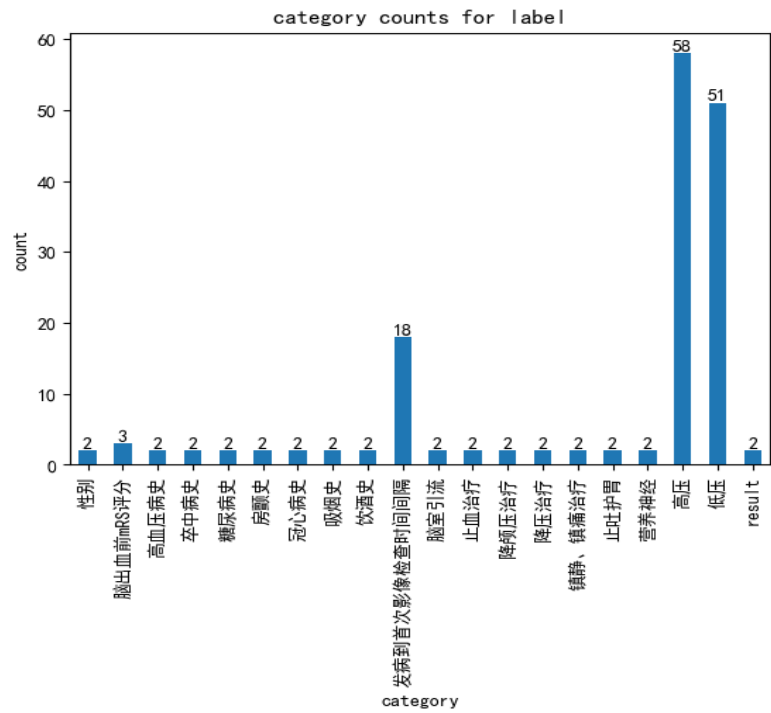


图 4-2 表 1 各字段类别统计图

“表 1”给出了较多的患者信息，通过上表得到不同字段的个数，下图 4-3 各字段的分布情况。

Scatter and Violin Plot

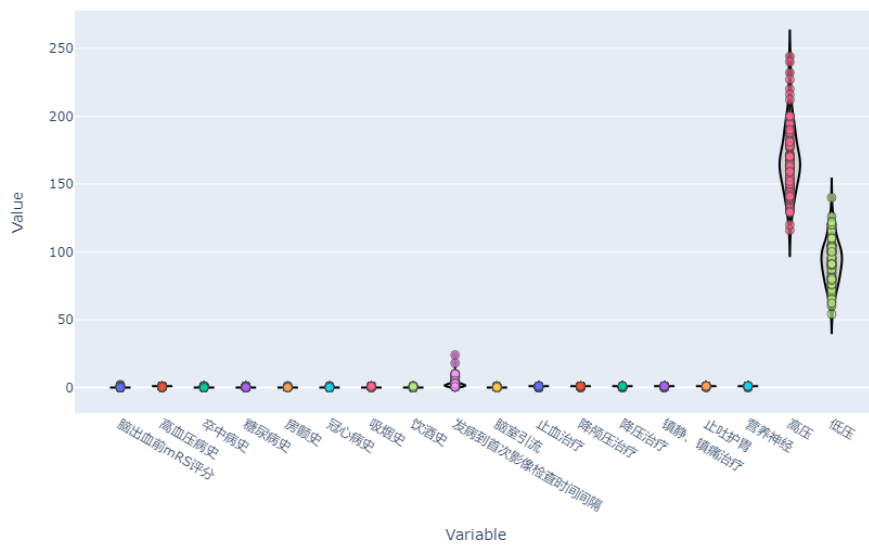


图 4-3 表 1 各字段数据分布图

Box Plot

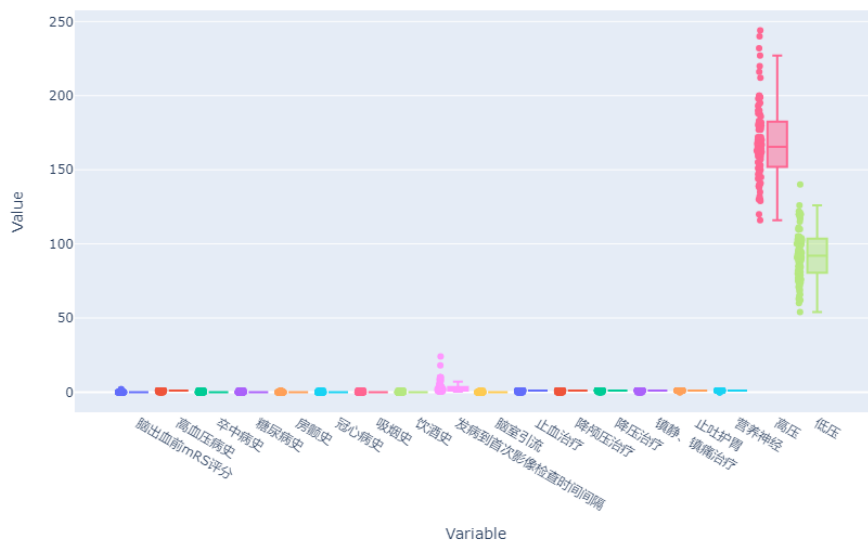


图 4-4 表 1 各字段盒图

根据图 4-4，可以轻易得知患者高压大多在 140-190 之间，低压大多在 60-110 之间。由于大多字段都为布尔型，因此需要进行进一步探索各个变量间的关系，如下图 4-5 所示。

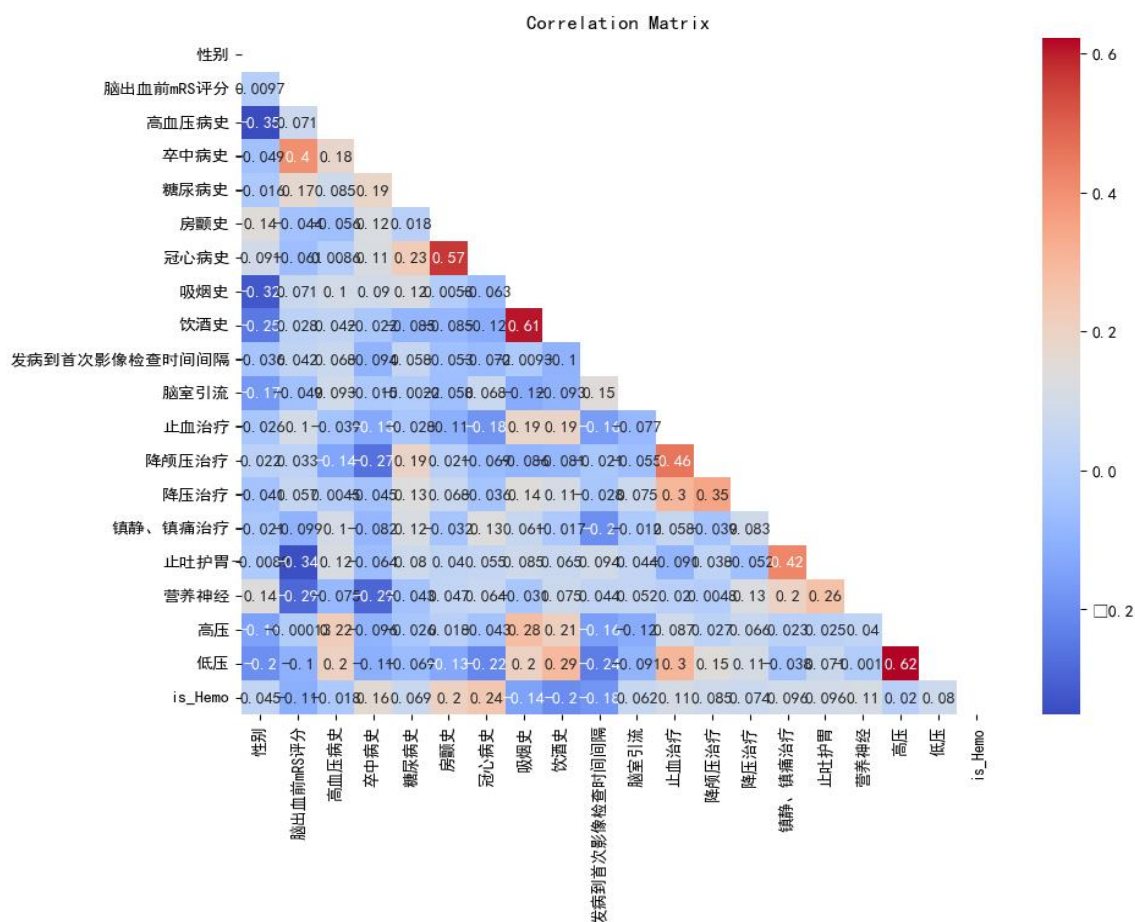


图 4-5 各影响因子对 Hemo 的相关性

根据上面的热力图，可以得到如下的相关性

表 4-5 与发生 Hemo 相关性表

影响因子	相关系数（绝对值）	正相关/负相关
脑出血前 mRs 评分	0.11	-
高血压病史	0.018	-
卒中病史	0.16	+
糖尿病史	0.069	+
房颤史	0.2	+
冠心病史	0.24	+
吸烟史	0.14	-
饮酒史	0.2	-
降颅压治疗	0.085	+
止血治疗	0.11	+

考虑到年龄对疾病的影响因素很大，因此将年龄作为横轴，进一步探索各变量和年龄关系图，如下图 4-5：

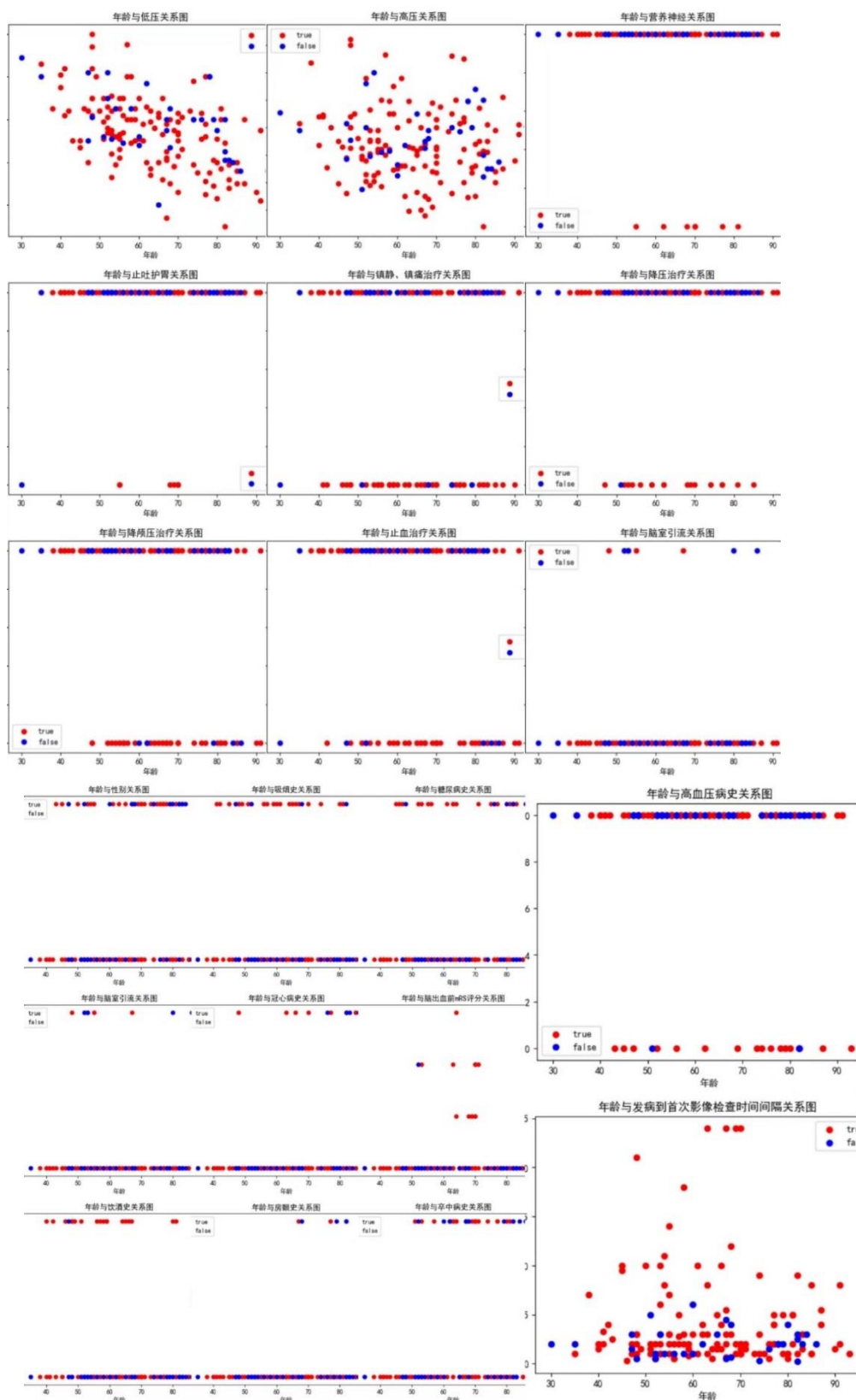


图 4-6 各指标与是否发生血肿对比

从图中不难看出，性别、营养神经与是否发生血肿扩张的关联性并不是那么密切，因此不将这两个影响因素作为模型的输入。

4.2.3 任务 b 模型的建立

任务 b 需要将 160 位患者分为数据集和测试集，这是一个小体量的监督分类问题，因此选择逻辑回归模型。在模型处理过程中，首先去掉输入中的“性别”和“营养神经”两个影响因子，作为模型一，再用所有字段作为逻辑回归模型的输入，记作模型二。

（数学表达式）逻辑回归模型通过使用 sigmoid 函数将线性预测转换为概率值，其数学公式如下：

$$P(Y = 1|X) = \frac{1}{(1 + \exp(-(\beta^0 + \beta^1X^1 + \beta^2X^2 + \dots + \beta_nX_n)))} \quad (4 - 1)$$

其中， $(P(Y=1|X))$ 表示在给定输入特征 (X) 的情况下，目标变量 (Y) 为类别 1 的概率。 $(\beta_0, \beta_1, \beta_2, \dots, \beta_n)$ 表示模型的系数（参数）， (X_1, X_2, \dots, X_n) 表示输入特征。

逻辑回归模型在训练阶段通过最大似然估计或正则化方法来学习参数值，从而得到最佳模型。在代码中，使用 GridSearchCV 进行网格搜索调优，并使用训练集数据 train_data 和 train_label 来训练模型并找到最佳的超参数组合 best_params_ 和最佳模型 best_model。

接下来，使用训练集数据 train_data 获取最佳模型的输出概率，通过调用 predict_proba 方法，得到对于每个样本属于各个类别的概率值。在代码中，train_prob 和 test_prob 分别存储了训练集和测试集的输出概率。

最后，使用列表解析和格式化字符串的方式，将输出概率值保留四位小数，并将其存储在 prob_issue1 列表中。

4.2.4 任务 b 模型结果

通过对两个模型的训练，得到如下的残差结果如表 4-6：

表 4-6 两类模型的总残差

模型名	总残差
模型一	1.2484214129536988
模型二	1.248439554901102

不难看出“性别”和“营养神经”会对模型的精度产影响，但是影响不大这与图 4-6，图 4-6 所分析的结果一致。

五、问题二的模型建立与求解

5.1 任务 a 的分析与处理

在预测患者脑部血肿体积随时间的进展时，患者脑部血肿体积的变化可能不符合简单的线性关系，而是会受到多种复杂因素的影响，因此我们倾向于建立多项式回归和随机森林模型来进行预测，使用最小二乘法来拟合数据来寻找最佳的拟合曲线。通过计算残差来评估模型的拟合效果和预测准确性，并绘制观测值与预测值之间的散点图，观察预测效果。

由于患者之间个体的差异，我们把患者分为 4 个亚组，对各个组进行多项式回归模型进行预测，并计算真实值和拟合曲线的残差。

5.2.1 任务 a 模型

基于上一小节所进行的数据处理，可以使用不同的方法建立回归模型，以期得到最优结果，包括多项式回归和随机森林模型。

考虑到患者个体差异因此患者脑补血肿体积随时间进展可能不符合简单的线性关系。首先使用多项式回归模型进行拟合曲线。

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_n X^n \quad (5-1)$$

在这个公式中，Y 是目标变量，X 是输入变量。 β_0 、 β_1 、 β_2 、...、 β_n 是模型的参数，分别对应截距和输入变量的权重。 X^2 、 X^3 、...、 X^n 表示输入变量的多项式项。 $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_n X^n$

在这个公式中，Y 是目标变量，X 是输入变量。 β_0 、 β_1 、 β_2 、...、 β_n 是模型的参数，分别对应截距和输入变量的权重。 X^2 、 X^3 、...、 X^n 表示输入变量的多项式项。

在本例中 X 对应时间间隔，Y 对应水肿体积。

通过对上一小节中对数据处理和分析知道，时间大于 500h 后，其数据不在变化，因此数据不具有参考行。考虑数据筛选，将时间间隔超过 500h 的数据过滤掉得到如下肿块随时间变化图（图 5-1）

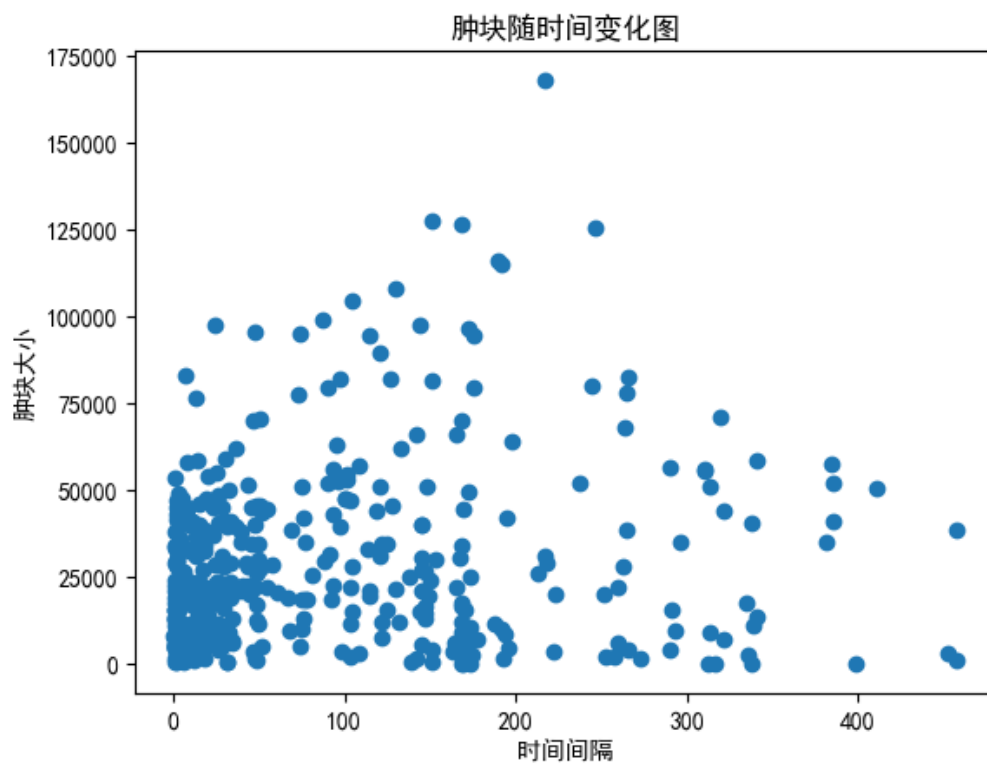


图 5-1 肿块随时间变化散点图

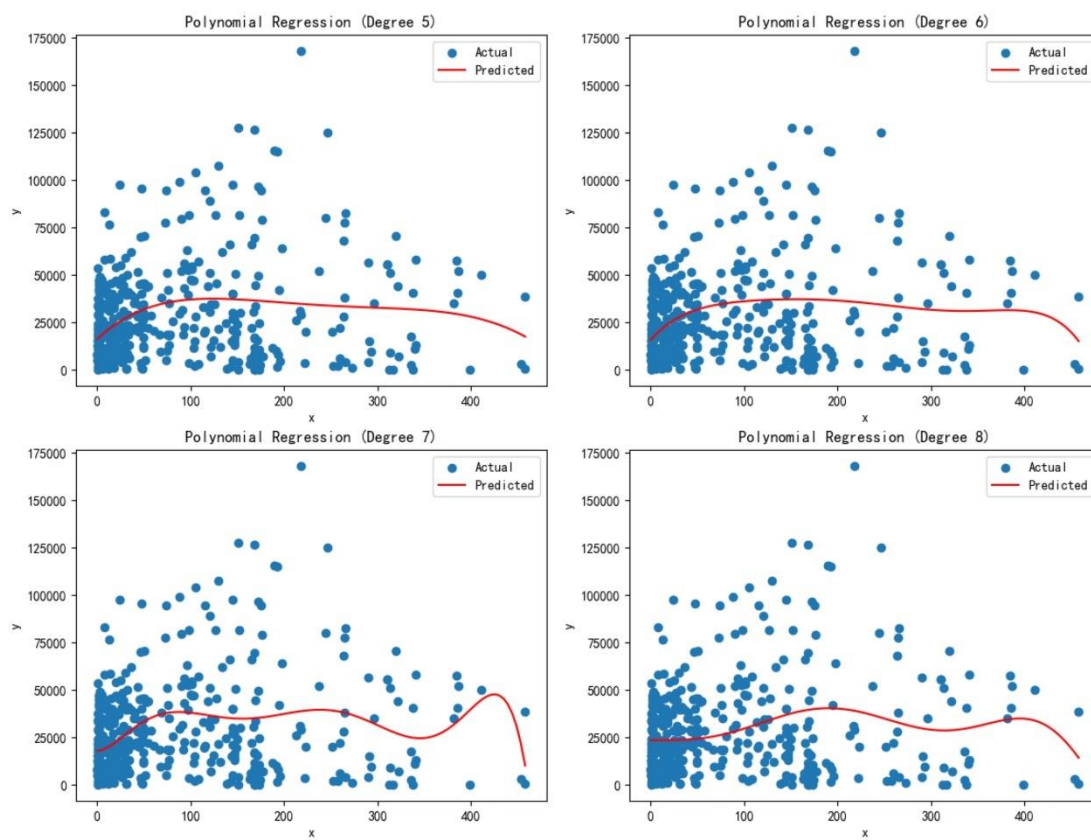


图 5-2 多项式阶数分别为 5,6,7,8 时拟合曲线和散点对比图

在尝试不同的多项式阶数后，通过模型输出结果得出在多项式阶数达到 7 阶时真实值和拟合值之间残差最小拟合效果最好。

接下来，使用随机森林模型进行拟合曲线。随机森林是一种集成学习方法，它通过构建多个决策树来进行预测和分类。决策树用于回归任务时，其预测公式如下：

$$y_{pred} = Mean(y_{leaf}) \quad (5-2)$$

决策树将每个样本分配到一个叶节点，并将该叶节点中样本的平均值作为预测结果。通过组合多个决策树的预测结果来进行最终的预测。对于回归任务，预测公式如下：

$$y_{pred} = Mean(y_{tree}) \quad (5-3)$$

随机森林将每个决策树的预测结果取平均，得到最终的预测结果拟合出如下图曲线。

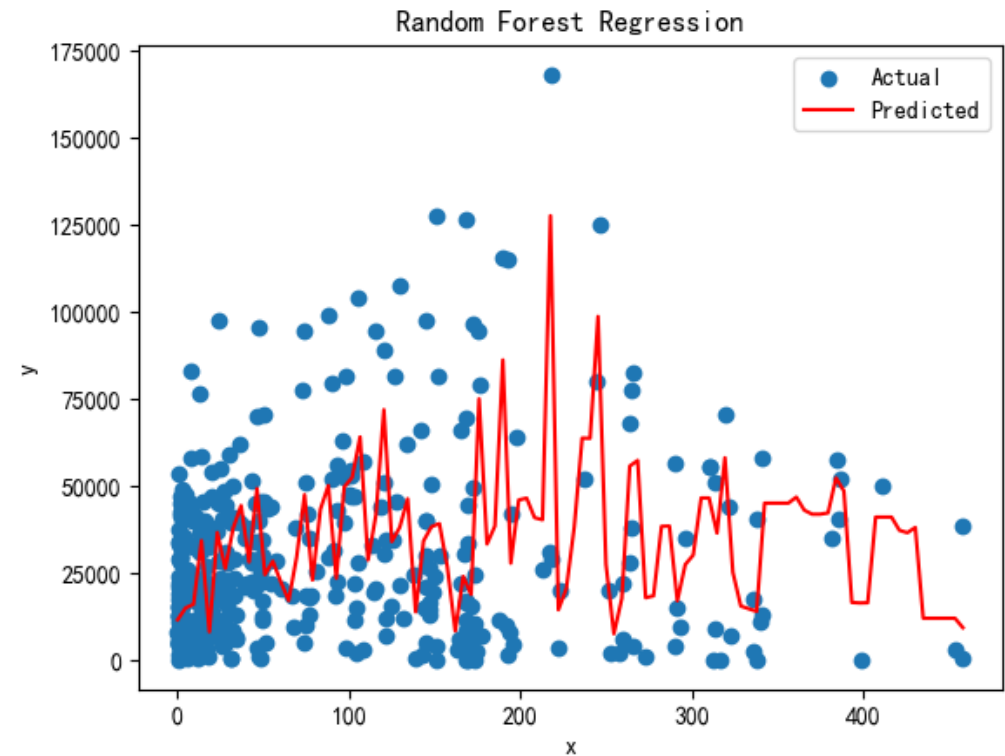


图 5-3 随机森林法构建拟合曲线

5.2 任务 b 的分析与处理

根据题意，我们将患者分为以下 4 个亚组，Class1：既抽烟又喝酒，Class2：只抽烟不喝酒，Class3：只喝酒不抽烟，Class4：不抽烟也不喝酒。绘制如下图所示的不同亚组水肿体积随时间变化散点图。

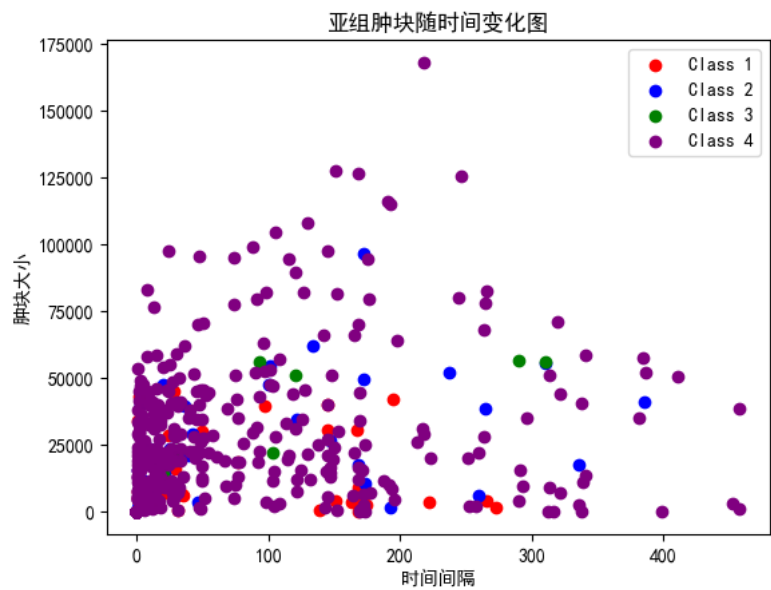


图 5-4 水肿体积随时间变化图

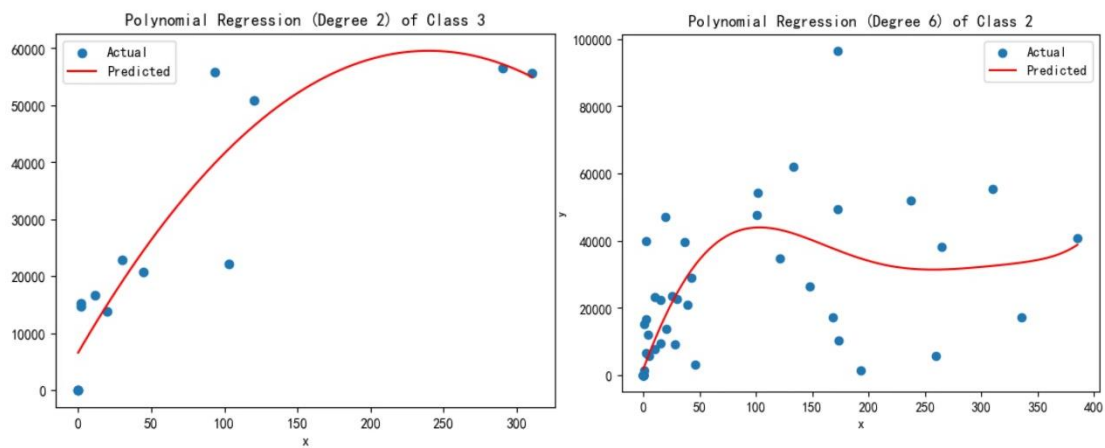


图 5-5 Class3 与 Class2 对比

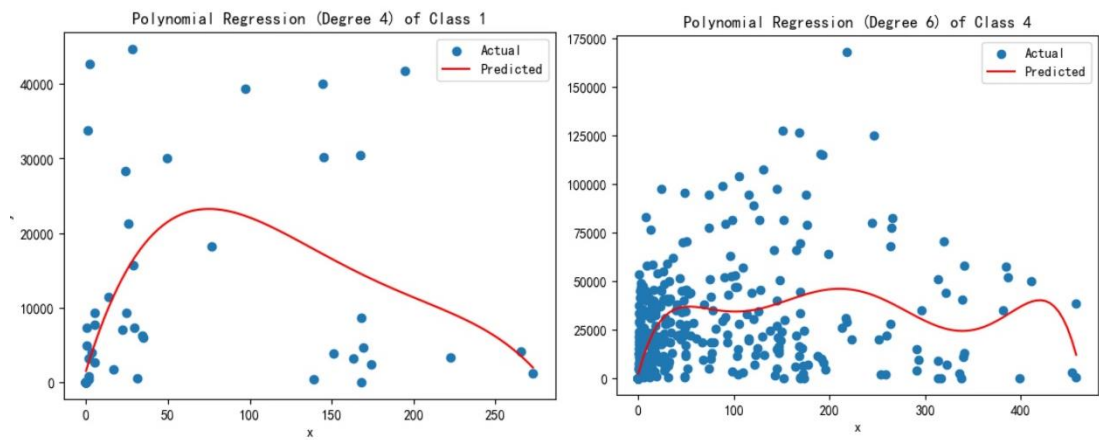


图 5-6 Class1 与 Class4 对比

通过图（5-5）和（5-6）不难看出抽烟喝酒对水肿体积影响较大。

5.3 任务 c 的分析与处理

通过结合“表 1”和“表 2”绘制不同治疗方法和水肿体积相关性热力图如下（图 5-7）

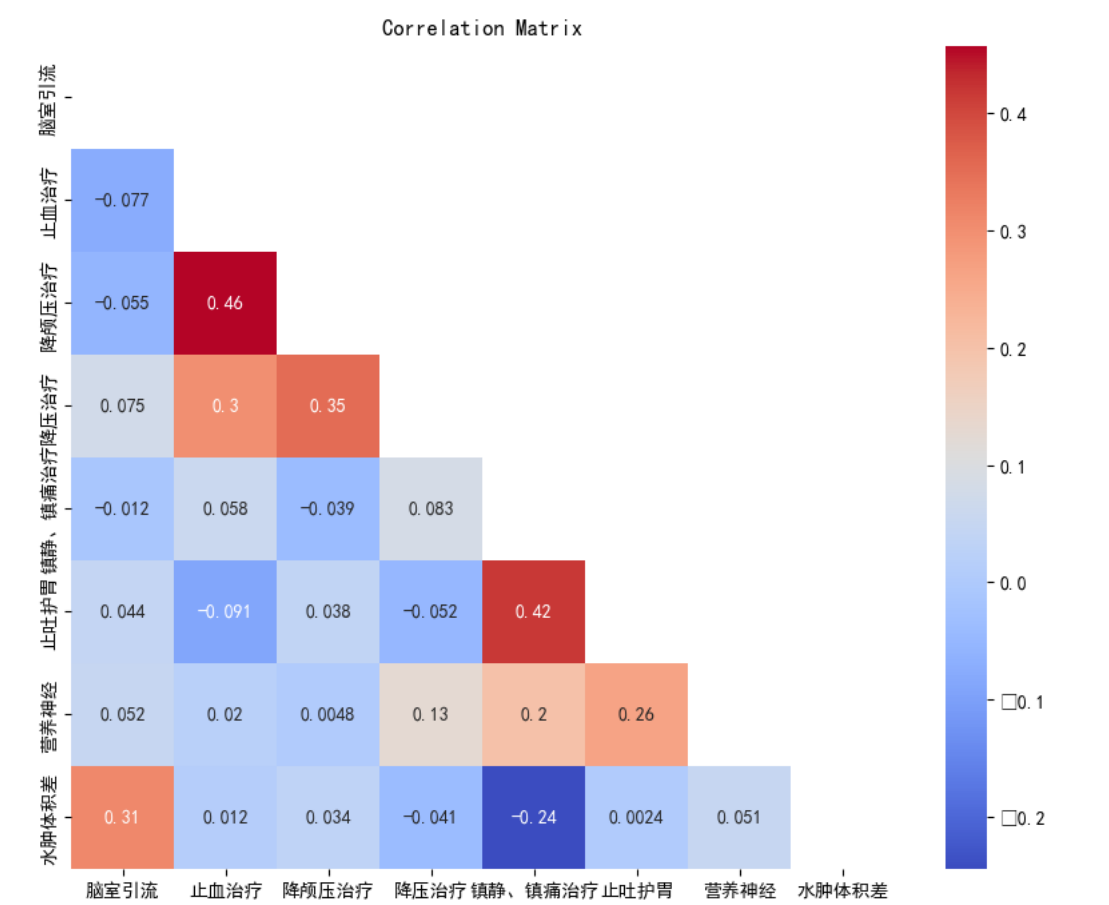


图 5-7 不同治疗方法和水肿体积相关性图

通过图 5-7 可以绘制出下表（表 5-1）

表 5-1 不同治疗方法和水肿体积相关性表

影响因子	相关系数（绝对值）	正负
脑室引流	0.31	+
止血治疗	0.012	+
降颅压治疗	0.034	+
降压治疗	0.041	-
镇静镇痛治疗	0.24	-
止吐护胃	0.0024	+
营养神经	0.051	+

通过表 5-1 可以看出降压治疗和镇静镇痛治疗对患者水肿体积存在副作用，而其他治疗手段均会对治疗起正向作用，有利于改善患者症状，其中脑室引流效果最好。

5.4 任务 d 的分析与处理

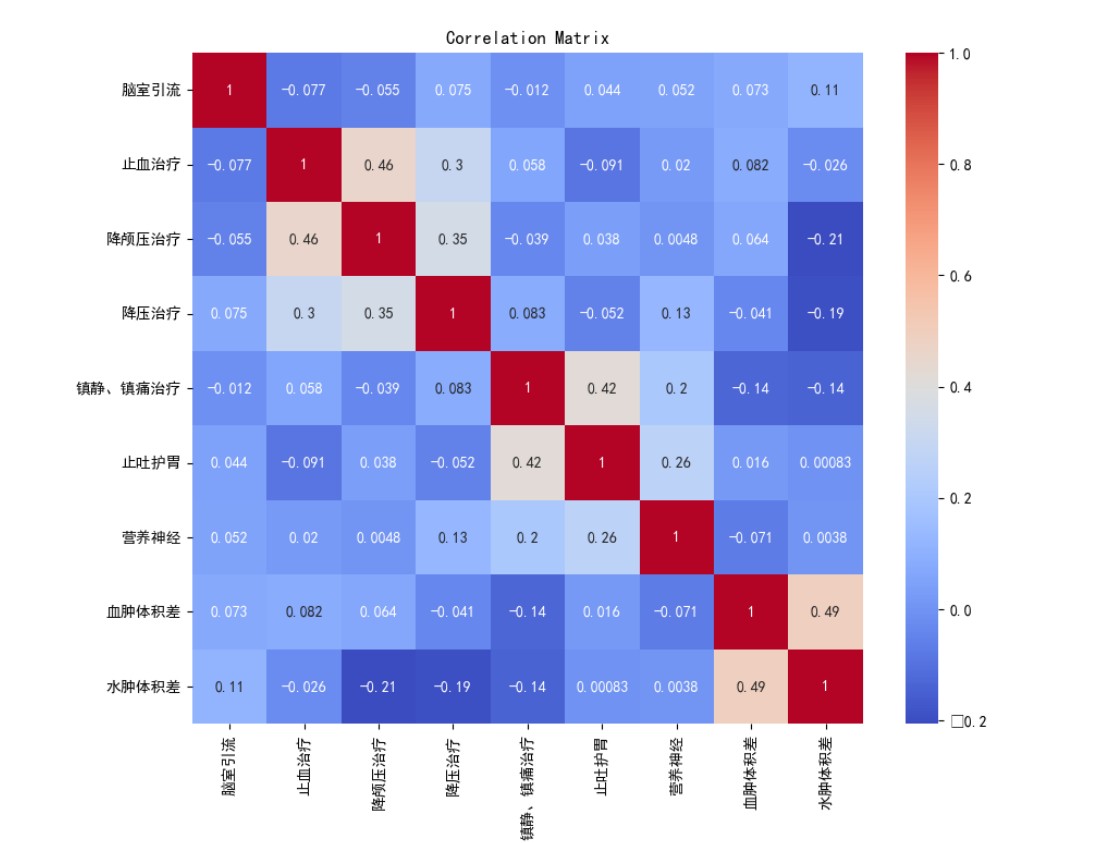


图 5-8 治疗手段与血肿水肿体积相关性

通过图 5-8 可以看出血肿体积和水肿体积成正相关，而且降压治疗和镇静、镇痛治疗对血肿和水肿体积起反作用，因此这两种治疗方案不适合对患者使用，推荐使用脑室引流行治疗。

六、问题三的模型建立和求解

6.1 问题 a 的分析与处理

由于患者个体差异，患者的发病历史和 mRs 评分可能不符合简单的线性关系，因此在预测 mRs 评分时选择以下三种回归模型并选择精度最好的模型进行预测。

- 1) 逻辑回归模型
- 2) 随机森林模型
- 3) K 临近算法模型

K 近邻 (K-Nearest Neighbors, 简称 KNN) 算法是一种基于距离度量的分类和回归算法。下面是 KNN 算法的公式表示：

对于两个样本 x_i 和 x_j ，常用的距离度量包括欧氏距离、曼哈顿距离等，表示为 $d(x_i, x_j)$ ；对于一个新样本 x_{new} ，KNN 算法通过以下公式进行分类预测：

$$y_{pred} = MajorityVote(y_{neighbors}) \quad (6-1)$$

其中 $y_{neighbors}$ 是 x_{new} 的 K 个最近邻居的类别标签。

KNN 算法通过计算样本之间的距离，并利用最近邻居的信息进行分类预测，采用多数表决的方式确定预测结果。

逻辑回归模型和随机森林模型的公式在上一小节有过介绍，在这里不再赘述。

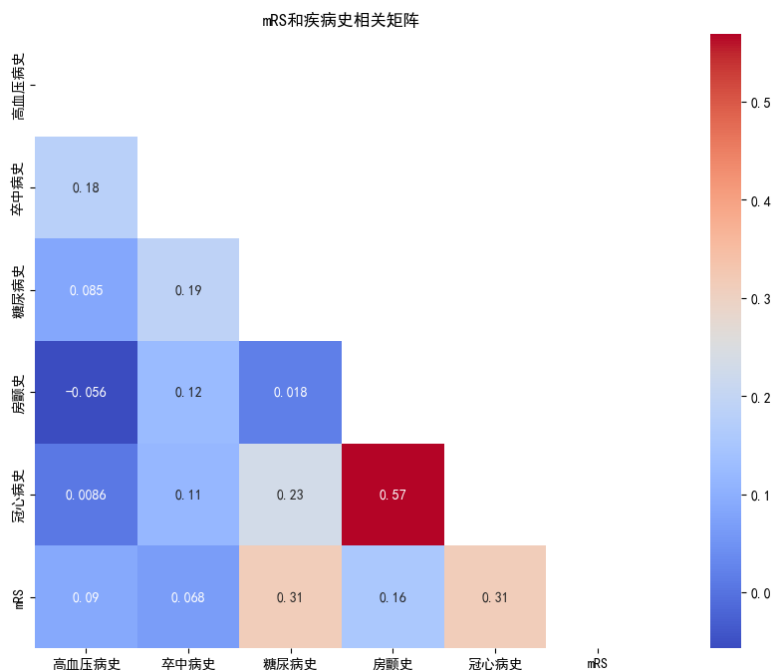


图 6-1 mRS 和疾病史相关性

从图 6-1 可以看出患者有相关病史会促使 mRS 评分升高。

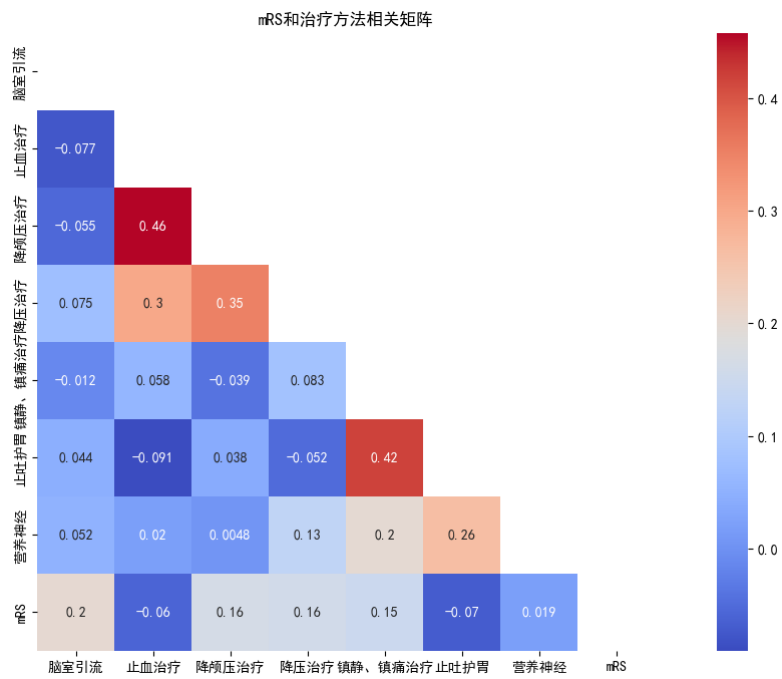


图 6-2 mRS 和治疗方法相关性

从图 6-2 可以看出患者进行相关治疗手段之后会使 mRS 评分降低或者减慢 mRS 评分上涨的速度。

七、参考文献

- [1] Yu J, Zheng J, Xu R, et al. Expansion of Intracerebral Hemorrhage: Incidence and Predictors. *World Neurosurg.* 2017;104:22-30.
- [2] Demchuk AM, Dowlatshahi D, Rodriguez-Luna D, et al. Prediction of Hematoma Growth and Outcome in Patients With Intracerebral Hemorrhage Using the CT-Angiography Spot Sign (PREDICT): A Prospective Observational Study. *Lancet Neurol.* 2012;11(4):307-14.
- [3] 刘娟, 王治国, 胡志坚. 脑出血后血肿扩张的危险因素和预测方法[J]. *临床神经病学杂志*, 2018, 31(6): 507-511.
- [4] 余骏, 李恭膺. 脑出血中的“Spot Sign”及扩散性血肿与病死率的关系[J]. *实用神经疾病杂志*, 2019, 22(11): 26-29.
- [5] Guo D, Wilkinson DA, Thompson BG, et al. Hematoma expansion predictors in intracerebral hemorrhage: Clinical and molecular perspectives. *Neurology.* 2016;86(11): 949-56.
- [6] van Asch CJ, Luitse MJ, Rinkel GJ, et al. Incidence, case fatality, and functional outcome of intracerebral haemorrhage over time, according to age, sex, and ethnic origin: a systematic review and meta-analysis. *Lancet Neurol.* 2010;9(2):167-76.
- [7] 杨秀芳, 楼晓琳, 吴荣华. 出血性脑卒中的临床特点及转归分析[J]. *医学综述*, 2019, 25(14): 2778-2781.
- [8] 陈萌, 彭红, 周锐. 出血性脑卒中患者的康复治疗研究进展[J]. *中国康复医学杂志*, 2020, 35(2): 230-233.
- [9] 张海滨, 赵智民, 李鹏. 出血性脑卒中的急性处理与护理进展[J]. *护理学杂志*, 2017, 32(16): 54-57.
- [10] 李宜春, 胡晓江, 范丹. 出血性脑卒中超早期治疗进展[J]. *中东国际医学杂志*, 2020, 20(3): 282-287.
- [11] Huynh TJ, Aviv RI, Dowlatshahi D, et al. Validation of the 9-point and 24-point hematoma expansion prediction scores and derivation of the PREDICT A/B scores. *Stroke.* 2015;46(11):3105-3110. doi:10.1161/STROKEAHA.115.010274
- [12] 王利娟, 陈晓通, 邓伟新. 出血性脑卒中的影像学诊断和评估新进展[J]. *中华放射学杂志*, 2018, 52(11): 878-882.

附录（代码）

```
1. import numpy as np
2. import pandas as pd
3. import matplotlib.pyplot as plt
4. from sklearn.preprocessing import LabelEncoder
5. from sklearn.metrics import accuracy_score
6. file_path = '../E 题/数据/表 1-患者列表及临床信息.xlsx'
7. table1 = pd.read_excel(file_path)
8. table1.info()
9. # 从 table_1 选择两列
10. table_1_issue1 = table_1[['入院首次影像检查流水号', '发病到首次影像检查时间间隔']]
11.
12. # 将时间间隔转换为时间间隔类型
13. table_1_issue1['发病到首次影像检查时间间隔'] = pd.to_timedelta(table_1_issue1['发病到首次影像检查时间间隔'], unit='h')
14. table_1_issue1.head()
15. # 统计空值的个数
16. null_count = table_1_issue1.isnull().sum()
17.
18. # 打印空值的个数
19. print('Null count:')
20. print(null_count)
21.
22. # 找到包含空值的列
23. null_columns = null_count[null_count > 0].index
24.
25. # 打印包含空值的列
26. print('Null columns:')
27. print(null_columns)
28. file_path = '../E 题/数据/表 2-患者影像信息血肿及水肿的体积及位置.xlsx'
29. table_2 = pd.read_excel(file_path)
30. table_2.head()
31. schedule_1 = pd.read_excel('../E 题/数据/附表 1-检索表格-流水号 vs 时间.xlsx')
32.
33. schedule_1.head()
34. # 筛选 table2 的列
```

```

35.index = ['首次检查流水号', '随访 1 流水号', '随访 2 流水号', '随访 3 流
    水号', '随访 4 流水号', '随访 5 流水号']
36.index_list = []
37.for i in index:
38.    index_list.append(table_2.columns.get_loc(i))
39.    index_list.append(table_2.columns.get_loc(i) + 1)
40.table_2_issue1 = table_2.iloc[:, index_list]
41.table_2_issue1.head()
42.from datetime import timedelta
43.# 判断血肿体积变化
44.result, time_list = [0 for _ in range(160)], ['0 小时
    ' for _ in range(160)]
45.latest_index = 0
46.for index, row in table_2_issue1.iterrows():
47.    for num in range(1, 6):
48.        latest_index = 2 * num + 1
49.        # 先判断是否发生血肿
50.        if row[latest_index] - row[1] >= 6000 or (row[latest_inde
            x] - row[1]) / row[1] > 0.33:
51.            # 再判断发生血肿是否在 48h 内
52.            # 计算率先达到血肿条件的时间戳索引, 判断这段时间是否超出
                48h
53.            row_index, col_index = np.where(schedule_1.values ==
                row[latest_index-1])
54.            # 如果血肿在 48h 内
55.            if schedule_1.iloc[row_index, col_index-
                1].values[0] != None and pd.Timedelta(days=2) >= schedule_1.iloc[
                row_index, col_index-
                1].values[0][0] + table_1_issue1.iloc[index]['发病到首次影像检查时
                间间隔'] - schedule_1.iloc[row_index, 2].values[0]:
56.                interval_time = schedule_1.iloc[row_index, col_in
                    dex-1].values[0][0] - schedule_1.iloc[row_index, 2].values[0]
57.                interval_time = timedelta(microseconds=int(interv
                    al_time / np.timedelta64(1, 'us')))
58.                result[index] = 1
59.                time_list[index] = f'{interval_time.total_seconds
                    () / 3600:.4f}小时'
60.                # print(index, num, result[index], time_list[inde
                    x])
61.                break
62.print(result)
63.print(time_list)

```



```
64.from openpyxl import load_workbook
65.
66.# 打开答案文件.xlsx
67.workbook = load_workbook('../E 题/数据/表 4-答案文件.xlsx')
68.
69.# 选择活动的工作表
70.worksheet = workbook.active
71.
72.# 获取起始单元格的列号和行号
73.result_start_column = 'C'
74.start_row = 4
75.time_start_column = 'D'
76.
77.# 逐行写入 result 列表的元素
78.for i in range(len(result)):
79.    # 定位
80.    result_current_column = result_start_column
81.    time_current_column = time_start_column
82.    result_cell_coordinate = f"{result_current_column}{start_row
+ i}"
83.    time_cell_coordinate = f"{time_current_column}{start_row + i}
"
84.    # 写入
85.    worksheet[result_cell_coordinate] = result[i]
86.    worksheet[time_cell_coordinate] = time_list[i]
87.
88.# 保存工作簿到文件
89.workbook.save('../E 题/数据/表 4-答案文件.xlsx')
90.# 处理血压
91.table_1[["高压", "低压"]] = table_1["血压
"].str.split("/", expand=True)
92.table_1[["高压", "低压"]] = table_1[["高压", "低压"]].astype(int)
93.table_1 = table_1.drop("血压", axis=1)
94.table_1.head()
95.# 处理性别 0-男 1-女
96.table_1['性别'] = table_1['性别'].replace({"男": 0, "女": 1})
97.table_1.head()
98.# 将 result 拼接过来
99.table_1 = pd.concat([table_1, pd.Series(result).rename('result')]
, axis=1)
100.# 统计空值的个数
101.null_count = table_1.isnull().sum()
```

```
102.
103. # 打印空值的个数
104. print('Null count:')
105. print(null_count)
106.
107. # 找到包含空值的列
108. null_columns = null_count[null_count > 0].index
109.
110. # 打印包含空值的列
111. print('Null columns:')
112. print(null_columns)
113. data = table_1.iloc[:, 2:]
114. train_data, test_data = data.iloc[:100, :], data.iloc[100:, :]
115. train_label, test_label = result[:100], result[100:]
116. train_data.shape, test_data.shape, len(train_label), len(test_label)
117. plt.rcParams['font.sans-serif'] = ['SimHei']
118. category_counts = train_data.nunique()
119.
120. plt.figure()
121. category_counts.plot(kind='bar')
122. plt.title('category counts for label')
123. plt.xlabel('category')
124. plt.ylabel('count')
125. # 显示数字
126. for i, v in enumerate(category_counts):
127.     plt.text(i, v, str(v), ha='center', va='bottom')
128.
129. plt.show()
130. train_data.info()
131. table_1_a = table_1.iloc[:, :12]
132. table_1_b = table_1.iloc[:, 12:]
133. import plotly.graph_objects as go
134. import pandas as pd
135. # 创建数据
136. data = train_data.iloc[:, 1:-1]
137.
138. # 创建图表布局
139. fig = go.Figure()
140.
141. # 添加每个变量的散点图
142. for column in data.columns:
```

```
143.     fig.add_trace(go.Scatter(  
144.         x=[column] * len(data),  
145.         y=data[column],  
146.         mode='markers',  
147.         name=column,  
148.         marker=dict(  
149.             size=8,  
150.             line=dict(width=1),  
151.             opacity=0.7  
152.         )  
153.     ))  
154.  
155. # 设置图表布局  
156. fig.update_layout(  
157.     title='Multiple Indicator Strip Plot',  
158.     xaxis=dict(title='Variable'),  
159.     yaxis=dict(title='Value'),  
160.     showlegend=False,  
161.     width=900,  
162.     height=600  
163. )  
164.  
165. # 显示图表  
166. fig.show()  
167. # 创建数据  
168. data = train_data.iloc[:, 1:-1]  
169.  
170. # 创建图表布局  
171. fig = go.Figure()  
172.  
173. # 添加每个变量的盒图  
174. for column in data.columns:  
175.     fig.add_trace(go.Box(  
176.         y=data[column],  
177.         name=column,  
178.         boxpoints='all',  
179.         jitter=0.3,  
180.         pointpos=-1.8  
181.     ))  
182.  
183. # 设置图表布局  
184. fig.update_layout(  

```

```
185.     title='Box Plot',
186.     xaxis=dict(title='Variable'),
187.     yaxis=dict(title='Value'),
188.     showlegend=False,
189.     width=900,
190.     height=600
191.)
192.
193.# 显示图表
194.fig.show()
195.# 创建数据
196.data = train_data.iloc[:, 1:-1]
197.
198.# 创建图表布局
199.fig = go.Figure()
200.
201.# 添加每个变量的散点图
202.for column in data.columns:
203.    fig.add_trace(go.Scatter(
204.        x=[column] * len(data),
205.        y=data[column],
206.        mode='markers',
207.        name=column,
208.        marker=dict(
209.            size=8,
210.            line=dict(width=1),
211.            opacity=0.7
212.        )
213.    ))
214.
215.# 添加每个变量的小提琴图
216.for column in data.columns:
217.    fig.add_trace(go.Violin(
218.        y=data[column],
219.        name=column,
220.        box_visible=True,
221.        meanline_visible=True,
222.        fillcolor='lightgray',
223.        line_color='black'
224.    ))
225.
226.# 设置图表布局
```

```
227.fig.update_layout(  
228.     title='Scatter and Violin Plot',  
229.     xaxis=dict(title='Variable'),  
230.     yaxis=dict(title='Value'),  
231.     showlegend=False,  
232.     width=900,  
233.     height=600  
234.)  
235.  
236.# 显示图表  
237.fig.show()  
238.train_data.shape  
239.# 数据分析  
240.# for i in range(2, len(table_1.columns)-1):  
241.#     subset1 = table_1[table_1['result'] == 0]  
242.#     subset2 = table_1[table_1['result'] == 1]  
243.  
244.#     plt.scatter(subset1['年龄'  
    '], subset1[table_1.columns[i]], color='red', label='true')  
245.#     plt.scatter(subset2['年龄'  
    '], subset2[table_1.columns[i]], color='blue', label='false')  
246.  
247.#     plt.xlabel('年龄')  
248.#     plt.ylabel(table_1.columns[i])  
249.#     plt.title(f'年龄与{table_1.columns[i]}关系图')  
250.#     plt.legend()  
251.#     plt.show()  
252.import pandas as pd  
253.import seaborn as sns  
254.import matplotlib.pyplot as plt  
255.import numpy as np  
256.  
257.# 创建数据  
258.data = train_data.iloc[:, :-1].assign(is_Hemo=train_label) # 假  
    设train_data 的标签位于第一列  
259.  
260.# 计算相关系数矩阵  
261.corr_matrix = data.corr()  
262.  
263.# 创建遮盖矩阵  
264.mask = np.triu(np.ones_like(corr_matrix, dtype=bool))  
265.
```

```
266.# 绘制相关矩阵
267.plt.figure(figsize=(16, 8))
268.sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=True,
                mask=mask)
269.plt.title('Correlation Matrix')
270.plt.show()
271.table_3 = pd.read_excel('../E 题/数据/表 3-患者影像信息水肿及水肿的
                形状及灰度分布.xlsx')
272.table_3
273.# 创建数据
274.data = table_2.iloc[:,100, 2:24].assign(is_Hemo=train_label) #
                假设train_data 的标签位于第一列
275.
276.# 计算相关系数矩阵
277.corr_matrix = data.corr()
278.
279.# 创建遮盖矩阵
280.mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
281.
282.# 绘制相关矩阵
283.plt.figure(figsize=(16, 8))
284.sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=True,
                mask=mask)
285.plt.title('Correlation Matrix')
286.plt.show()
287.# 创建数据
288.data = table_3.iloc[:,100, :].assign(is_Hemo=train_label) # 假设
                train_data 的标签位于第一列
289.
290.# 计算相关系数矩阵
291.corr_matrix = data.corr()
292.
293.# 创建遮盖矩阵
294.mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
295.
296.# 绘制相关矩阵
297.plt.figure(figsize=(16, 8))
298.sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=True,
                mask=mask)
299.plt.title('Correlation Matrix')
300.plt.show()
301.from sklearn.linear_model import LogisticRegression
```

```

302. from sklearn.model_selection import GridSearchCV
303.
304. # 创建逻辑回归模型对象
305. model = LogisticRegression()
306.
307. # 定义超参数的取值范围
308. param_grid = {
309.     'C': [0.1, 1, 10],
310.     'penalty': ['l1', 'l2'],
311.     'solver': ['liblinear']
312. }
313.
314. # 创建 GridSearchCV 对象
315. grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5)
316.
317. # 使用训练集进行网格搜索调优
318. grid_search.fit(train_data, train_label)
319.
320. # 输出最佳超参数组合
321. print(grid_search.best_params_)
322.
323. # 输出最佳模型
324. best_model = grid_search.best_estimator_
325.
326. # 使用训练集获取模型的输出概率
327. train_prob = best_model.predict_proba(train_data)
328. test_prob = best_model.predict_proba(test_data)
329. prob_issue1 = [f"{x:.4f}" for x in train_prob[:, 1:].reshape(-1)] + [f"{x:.4f}" for x in test_prob[:, 1:].reshape(-1)]
330. print(prob_issue1)
331. sum_resid = sum(np.abs(test_label-test_prob[:, 1:].reshape(-1)))
332.
333. sum_resid
334. # data = table_1.iloc[:, 2:]
335. # data1 = data.drop(columns=['性别', '营养神经'])
336. # data1.head()
337. # train_data, test_data = data1.iloc[:100, :], data1.iloc[100:, :]
338. # # 创建逻辑回归模型对象
339. # model1 = LogisticRegression()
340.

```

```
341.## # 定义超参数的取值范围
342.## param_grid = {
343.##     'C': [0.1, 1, 10],
344.##     'penalty': ['l1', 'l2'],
345.##     'solver': ['liblinear']
346.## }
347.
348.## # 创建 GridSearchCV 对象
349.## grid_search = GridSearchCV(estimator=model1, param_grid=param_
    grid, cv=5)
350.
351.## # 使用训练集进行网格搜索调优
352.## grid_search.fit(train_data, train_label)
353.
354.## # 输出最佳超参数组合
355.## print(grid_search.best_params_)
356.
357.## # 输出最佳模型
358.## best_model = grid_search.best_estimator_
359.
360.## # 使用训练集获取模型的输出概率
361.## train_prob = best_model.predict_proba(train_data)
362.## test_prob = best_model.predict_proba(test_data)
363.## prob_issue1 = [f"{x:.4f}" for x in train_prob[:, 1:].reshape(-
    1)] + [f"{x:.4f}" for x in test_prob[:, 1:].reshape(-1)]
364.## print(prob_issue1)
365.## sum_resid = sum(np.abs(test_label-test_prob[:, 1:].reshape(-
    1)))
366.
367.## sum_resid
368.from openpyxl import load_workbook
369.
370.## 打开答案文件.xlsx
371.workbook = load_workbook('../E 题/数据/表 4-答案文件.xlsx')
372.
373.## 选择活动的工作表
374.worksheet = workbook.active
375.
376.## 获取起始单元格的列号和行号
377.prob_start_column = 'E'
378.start_row = 4
379.
```



```

380. # 逐行写入 result 列表的元素
381. for i in range(len(prob_issue1)):
382.     # 定位
383.     prob_current_column = prob_start_column
384.     prob_cell_coordinate = f"{prob_current_column}{start_row + i
        }"
385.     # 写入
386.     worksheet[prob_cell_coordinate] = prob_issue1[i]
387.
388. # 保存工作簿到文件
389. workbook.save('../E 题/数据/表 4-答案文件.xlsx')
390. # 筛选 table2 的列
391. index = [f'随访{i}流水号' for i in range(1, 9)]
392. index.insert(0, '首次检查流水号')
393. table_indices, schedule_indices = [], [2*i for i in range(1, len
        (index)+1)]
394. for i in index:
395.     table_indices.append(table_2.columns.get_loc(i) + 12)
396.
397. Y_train = table_2.iloc[:100, table_indices]
398. schedule_1_issue2 = schedule_1.iloc[:100, schedule_indices]
399. Y_train.head()
400. # 创建空的 x_train DataFrame, 与 schedule_1_issue2 有相同的结构
401. X_train = pd.DataFrame(columns=schedule_1_issue2.columns)
402.
403. # 遍历 schedule_1_issue2 中的每一行
404. for i in range(len(schedule_1_issue2)):
405.     # 获取当前行的时间戳
406.     current_row = schedule_1_issue2.iloc[i]
407.
408.     # 创建空的 diff 列表, 用于存储时间差值
409.     diff = []
410.
411.     # 遍历当前行的每个单元格
412.     for j in range(len(current_row)):
413.         # 判断当前单元格是否为空
414.         if pd.isnull(current_row[j]):
415.             # 如果当前单元格为空, 则不进行处理, 将时间差值设置为 NaN
416.             diff.append(float('nan'))
417.         else:
418.             # 获取前一个单元格的时间戳
419.             if j == 0:

```

```
420.         previous_cell = pd.to_datetime(current_row[j])
421.     else:
422.         previous_cell = schedule_1_issue2.iloc[i, j - 1]
423.
424.     # 将当前单元格时间戳转换为 datetime 类型
425.     current_cell = pd.to_datetime(current_row[j])
426.
427.     # 计算时间差，并将结果添加到 diff 列表
428.     time_diff = current_cell - previous_cell + table_1_i
        ssue1.iloc[i]['发病到首次影像检查时间间隔']
429.     diff.append(round(time_diff.total_seconds() / 3600,
        2))
430.
431.     # 将 diff 列表添加到 x_train
432.     X_train.loc[i] = diff
433.
434. # 打印处理后的 x_train
435. X_train.head()
436. import numpy as np
437. import matplotlib.pyplot as plt
438.
439. # 获取对应位置的 x 和 y 值
440. x_values = X_train.values.flatten().tolist()
441. y_values = Y_train.values.flatten().tolist()
442.
443. # 创建筛选后的数据列表
444. filtered_x = []
445. filtered_y = []
446.
447. # 根据条件筛选数据
448. for x, y in zip(x_values, y_values):
449.     if x <= 500:
450.         filtered_x.append(x)
451.         filtered_y.append(y)
452.
453. # 绘制散点图
454. plt.scatter(filtered_x, filtered_y)
455.
456. # 设置横纵坐标标签
457. plt.xlabel('时间间隔')
458. plt.ylabel('肿块大小')
459. plt.title('肿块随时间变化图')
```

```
460.
461. # 显示图像
462. plt.show()
463. # 找到空数据所在的索引
464. nan_indices = np.isnan(filtered_x) | np.isnan(filtered_y)
465.
466. # 获取空数据所在的索引
467. indices = np.where(nan_indices)[0]
468.
469. # 删除空数据所在的索引对应的元素
470. x = np.delete(filtered_x, indices)
471. y = np.delete(filtered_y, indices)
472. len(x), len(y)
473. import numpy as np
474. from sklearn.preprocessing import PolynomialFeatures
475. from sklearn.linear_model import LinearRegression
476.
477. # 将 train_data 转换为二维数组
478. x_train = np.array(x).reshape(-1, 1)
479.
480. # 将 train_label 转换为二维数组
481. y_train = np.array(y).reshape(-1, 1)
482.
483. best_degree = None
484. best_residual_sum = float('inf')
485.
486. # 尝试不同的多项式阶数
487. for degree in range(1, 10):
488.     # 使用多项式特征转换
489.     poly_features = PolynomialFeatures(degree=degree)
490.     x_train_poly = poly_features.fit_transform(x_train)
491.
492.     # 创建多项式回归模型并进行训练
493.     model = LinearRegression()
494.     model.fit(x_train_poly, y_train)
495.
496.     # 使用模型对训练数据进行预测
497.     y_train_pred = model.predict(x_train_poly)
498.
499.     # 计算残差和
500.     residual_sum = np.sum(np.abs(y_train - y_train_pred))
501.
```

```
502.     # 更新最小残差和和对应的阶数
503.     if residual_sum < best_residual_sum:
504.         best_residual_sum = residual_sum
505.         best_degree = degree
506.
507. # 使用最佳阶数重新训练模型
508. poly_features = PolynomialFeatures(degree=best_degree)
509. x_train_poly = poly_features.fit_transform(x_train)
510. model = LinearRegression()
511. model.fit(x_train_poly, y_train)
512.
513. # 生成测试数据
514. x_test = np.linspace(np.min(x_train), np.max(x_train), 100).reshape(-1, 1)
515. x_test_poly = poly_features.transform(x_test)
516.
517. # 使用模型进行预测
518. y_pred = model.predict(x_test_poly)
519.
520. # 绘制真实值和预测值的图像
521. plt.scatter(x_train, y_train, label='Actual')
522. plt.plot(x_test, y_pred, color='red', label='Predicted')
523. plt.xlabel('x')
524. plt.ylabel('y')
525. plt.title('Polynomial Regression (Degree {})'.format(best_degree))
526. plt.legend()
527. plt.show()
528.
529. # 使用模型对训练数据进行预测
530. y_train_pred = model.predict(x_train_poly)
531.
532. # 计算最终的残差和
533. residual_sum = np.sum(np.abs(y_train - y_train_pred))
534.
535. # 打印最终的残差和
536. print(residual_sum, best_degree)
537. import numpy as np
538. from sklearn.ensemble import RandomForestRegressor
539.
540. # 将 train_data 转换为二维数组
541. x_train = np.array(x).reshape(-1, 1)
```

```
542.
543. # 将 train_label 转换为一维数组
544. y_train = np.array(y)
545.
546. # 创建随机森林回归模型并进行训练
547. model = RandomForestRegressor(n_estimators=100)
548. model.fit(x_train, y_train)
549.
550. # 生成测试数据
551. x_test = np.linspace(np.min(x_train), np.max(x_train), 100).reshape(-1, 1)
552.
553. # 使用模型进行预测
554. y_pred = model.predict(x_test)
555.
556. # 绘制真实值和预测值的图像
557. plt.scatter(x_train, y_train, label='Actual')
558. plt.plot(x_test, y_pred, color='red', label='Predicted')
559. plt.xlabel('x')
560. plt.ylabel('y')
561. plt.title('Random Forest Regression')
562. plt.legend()
563. plt.show()
564. # 使用模型对训练数据进行预测
565. y_train_pred = model.predict(x_train)
566.
567. # 计算最终的残差和
568. residual_sum = np.sum(np.abs(y_train - y_train_pred))
569.
570. # 打印最终的残差和
571. print(residual_sum)
572. x_values_without_nan = np.nan_to_num(x_values)
573. y_values_without_nan = np.nan_to_num(y_values)
574. y_values_pred = model.predict(x_values_without_nan.reshape(-1, 1))
575. # 获取 x_values_without_nan 中数值为 0 的索引位置
576. zero_indices = np.where(x_values_without_nan == 0)
577.
578. # 将 y_values_pred 中对应索引位置的值替换为 0
579. y_values_pred[zero_indices] = 0
580. # 计算 y_values_pred 和 y_values_without_nan 之间的绝对值差
581. abs_diff = y_values_pred - y_values_without_nan
```

```
582.
583. # 每 9 个数据求和并存储在 resid_list 中
584. resid_list = []
585. for i in range(0, len(abs_diff), 9):
586.     resid_sum = np.sum(abs_diff[i:i+5])
587.     resid_list.append(resid_sum / 9)
588.
589. # resid_list
590. # resid = y_train.reshape(-1) - y_train_pred
591. # 打开答案文件.xlsx
592. workbook = load_workbook('../E 题/数据/表 4-答案文件.xlsx')
593.
594. # 选择活动的工作表
595. worksheet = workbook.active
596.
597. # 获取起始单元格的列号和行号
598. start_column = 'F'
599. start_row = 4
600.
601. # 逐行写入 result 列表的元素
602. for i in range(len(resid_list)):
603.     # 定位
604.     current_column = start_column
605.     current_column = start_column
606.     cell_coordinate = f"{current_column}{start_row + i}"
607.     # 写入
608.     worksheet[cell_coordinate] = resid_list[i]
609.
610. # 保存工作簿到文件
611. workbook.save('../E 题/数据/表 4-答案文件.xlsx')
612. train_table_1 = table_1.iloc[:100, :]
613. idx_1 = train_table_1[(train_table_1["吸烟史"] == 1) & (train_table_1["饮酒史"] == 1)].index.tolist()
614. idx_2 = train_table_1[(train_table_1["吸烟史"] == 1) & (train_table_1["饮酒史"] == 0)].index.tolist()
615. idx_3 = train_table_1[(train_table_1["吸烟史"] == 0) & (train_table_1["饮酒史"] == 1)].index.tolist()
616. idx_4 = train_table_1[(train_table_1["吸烟史"] == 0) & (train_table_1["饮酒史"] == 0)].index.tolist()
617. len(idx_1) + len(idx_2) + len(idx_3) + len(idx_4)
618. # 根据索引分类 x_train 的行
619. x_train_1 = X_train.loc[idx_1]
```

```
620.x_train_2 = X_train.loc[idx_2]
621.x_train_3 = X_train.loc[idx_3]
622.x_train_4 = X_train.loc[idx_4]
623.
624.x_train_1.shape, x_train_2.shape, x_train_3.shape, x_train_4.shape
625.# 根据索引分类x_train 的行
626.y_train_1 = Y_train.loc[idx_1]
627.y_train_2 = Y_train.loc[idx_2]
628.y_train_3 = Y_train.loc[idx_3]
629.y_train_4 = Y_train.loc[idx_4]
630.
631.y_train_1.shape, y_train_2.shape, y_train_3.shape, y_train_4.shape
632.# 填充缺失值为均值
633.x_train_1 = x_train_1.fillna(0)
634.x_train_2 = x_train_2.fillna(0)
635.x_train_3 = x_train_3.fillna(0)
636.x_train_4 = x_train_4.fillna(0)
637.
638.y_train_1 = y_train_1.fillna(0)
639.y_train_2 = y_train_2.fillna(0)
640.y_train_3 = y_train_3.fillna(0)
641.y_train_4 = y_train_4.fillna(0)
642.import numpy as np
643.import matplotlib.pyplot as plt
644.
645.# 获取对应位置的 x 和 y 值
646.x_values_1 = x_train_1.values.flatten().tolist()
647.y_values_1 = y_train_1.values.flatten().tolist()
648.x_values_2 = x_train_2.values.flatten().tolist()
649.y_values_2 = y_train_2.values.flatten().tolist()
650.x_values_3 = x_train_3.values.flatten().tolist()
651.y_values_3 = y_train_3.values.flatten().tolist()
652.x_values_4 = x_train_4.values.flatten().tolist()
653.y_values_4 = y_train_4.values.flatten().tolist()
654.
655.# 创建筛选后的数据列表
656.filtered_x1 = []
657.filtered_y1 = []
658.filtered_x2 = []
659.filtered_y2 = []
```

```
660.filtered_x3 = []
661.filtered_y3 = []
662.filtered_x4 = []
663.filtered_y4 = []
664.
665.# 根据条件筛选数据
666.for x, y in zip(x_values_1, y_values_1):
667.    if x <= 500:
668.        filtered_x1.append(x)
669.        filtered_y1.append(y)
670.
671.for x, y in zip(x_values_2, y_values_2):
672.    if x <= 500:
673.        filtered_x2.append(x)
674.        filtered_y2.append(y)
675.
676.for x, y in zip(x_values_3, y_values_3):
677.    if x <= 500:
678.        filtered_x3.append(x)
679.        filtered_y3.append(y)
680.
681.for x, y in zip(x_values_4, y_values_4):
682.    if x <= 500:
683.        filtered_x4.append(x)
684.        filtered_y4.append(y)
685.
686.# 绘制散点图
687.plt.scatter(filtered_x1, filtered_y1, c='red', label='Class 1')
688.plt.scatter(filtered_x2, filtered_y2, c='blue', label='Class 2')
689.plt.scatter(filtered_x3, filtered_y3, c='green', label='Class 3'
690.    )
691.plt.scatter(filtered_x4, filtered_y4, c='purple', label='Class 4
692.    ')
693.
694.# 设置横纵坐标标签
695.plt.xlabel('时间间隔')
696.plt.ylabel('肿块大小')
697.plt.title('亚组肿块随时间变化图')
698.
699.# 添加图例
700.plt.legend()
701.
```



```
700.# 显示图像
701.plt.show()
702.# 找到空数据所在的索引
703.nan_indices = np.isnan(filtered_x1) | np.isnan(filtered_y1)
704.
705.# 获取空数据所在的索引
706.indices = np.where(nan_indices)[0]
707.
708.# 删除空数据所在的索引对应的元素
709.x1 = np.delete(filtered_x1, indices)
710.y1 = np.delete(filtered_y1, indices)
711.import numpy as np
712.from sklearn.preprocessing import PolynomialFeatures
713.from sklearn.linear_model import LinearRegression
714.
715.# 将 train_data 转换为二维数组
716.x_train = np.array(x1).reshape(-1, 1)
717.
718.# 将 train_label 转换为二维数组
719.y_train = np.array(y1).reshape(-1, 1)
720.
721.best_degree = None
722.best_residual_sum = float('inf')
723.
724.# 尝试不同的多项式阶数
725.for degree in range(1, 10):
726.    # 使用多项式特征转换
727.    poly_features = PolynomialFeatures(degree=degree)
728.    x_train_poly = poly_features.fit_transform(x_train)
729.
730.    # 创建多项式回归模型并进行训练
731.    model = LinearRegression()
732.    model.fit(x_train_poly, y_train)
733.
734.    # 使用模型对训练数据进行预测
735.    y_train_pred = model.predict(x_train_poly)
736.
737.    # 计算残差和
738.    residual_sum = np.sum(np.abs(y_train - y_train_pred))
739.
740.    # 更新最小残差和和对应的阶数
741.    if residual_sum < best_residual_sum:
```

```
742.         best_residual_sum = residual_sum
743.         best_degree = degree
744.
745. # 使用最佳阶数重新训练模型
746. poly_features = PolynomialFeatures(degree=best_degree)
747. x_train_poly = poly_features.fit_transform(x_train)
748. model = LinearRegression()
749. model.fit(x_train_poly, y_train)
750.
751. # 生成测试数据
752. x_test = np.linspace(np.min(x_train), np.max(x_train), 100).reshape(-1, 1)
753. x_test_poly = poly_features.transform(x_test)
754.
755. # 使用模型进行预测
756. y_pred = model.predict(x_test_poly)
757.
758. # 绘制真实值和预测值的图像
759. plt.scatter(x_train, y_train, label='Actual')
760. plt.plot(x_test, y_pred, color='red', label='Predicted')
761. plt.xlabel('x')
762. plt.ylabel('y')
763. plt.title('Polynomial Regression (Degree {}) of Class 1'.format(
    best_degree))
764. plt.legend()
765. plt.show()
766.
767. # 使用模型对训练数据进行预测
768. y_train_pred = model.predict(x_train_poly)
769.
770. # 计算最终的残差和
771. residual_sum = np.sum(np.abs(y_train - y_train_pred))
772.
773. # 打印最终的残差和
774. print(residual_sum, best_degree)
775. resid_list2 = [0 for _ in range(100)]
776. i = 0
777. for idx in idx_1:
778.     resid = np.mean((y_train - y_train_pred)[i:i+8])
779.     resid_list2[idx] = resid
780.     i += 1
781. # 找到空数据所在的索引
```

```
782.nan_indices = np.isnan(filtered_x2) | np.isnan(filtered_y2)
783.
784.# 获取空数据所在的索引
785.indices = np.where(nan_indices)[0]
786.
787.# 删除空数据所在的索引对应的元素
788.x2 = np.delete(filtered_x2, indices)
789.y2 = np.delete(filtered_y2, indices)
790.import numpy as np
791.from sklearn.preprocessing import PolynomialFeatures
792.from sklearn.linear_model import LinearRegression
793.
794.# 将 train_data 转换为二维数组
795.x_train = np.array(x2).reshape(-1, 1)
796.
797.# 将 train_label 转换为二维数组
798.y_train = np.array(y2).reshape(-1, 1)
799.
800.best_degree = None
801.best_residual_sum = float('inf')
802.
803.# 尝试不同的多项式阶数
804.for degree in range(1, 10):
805.    # 使用多项式特征转换
806.    poly_features = PolynomialFeatures(degree=degree)
807.    x_train_poly = poly_features.fit_transform(x_train)
808.
809.    # 创建多项式回归模型并进行训练
810.    model = LinearRegression()
811.    model.fit(x_train_poly, y_train)
812.
813.    # 使用模型对训练数据进行预测
814.    y_train_pred = model.predict(x_train_poly)
815.
816.    # 计算残差和
817.    residual_sum = np.sum(np.abs(y_train - y_train_pred))
818.
819.    # 更新最小残差和和对应的阶数
820.    if residual_sum < best_residual_sum:
821.        best_residual_sum = residual_sum
822.        best_degree = degree
823.
```

```
824. # 使用最佳阶数重新训练模型
825. poly_features = PolynomialFeatures(degree=best_degree)
826. x_train_poly = poly_features.fit_transform(x_train)
827. model = LinearRegression()
828. model.fit(x_train_poly, y_train)
829.
830. # 生成测试数据
831. x_test = np.linspace(np.min(x_train), np.max(x_train), 100).reshape(-1, 1)
832. x_test_poly = poly_features.transform(x_test)
833.
834. # 使用模型进行预测
835. y_pred = model.predict(x_test_poly)
836.
837. # 绘制真实值和预测值的图像
838. plt.scatter(x_train, y_train, label='Actual')
839. plt.plot(x_test, y_pred, color='red', label='Predicted')
840. plt.xlabel('x')
841. plt.ylabel('y')
842. plt.title('Polynomial Regression (Degree {}) of Class 2'.format(
    best_degree))
843. plt.legend()
844. plt.show()
845.
846. # 使用模型对训练数据进行预测
847. y_train_pred = model.predict(x_train_poly)
848.
849. # 计算最终的残差和
850. residual_sum = np.sum(np.abs(y_train - y_train_pred))
851.
852. # 打印最终的残差和
853. print(residual_sum, best_degree)
854. i = 0
855. for idx in idx_2:
856.     resid = np.mean((y_train - y_train_pred)[i:i+9])
857.     resid_list2[idx] = resid
858.     i += 1
859. # 找到空数据所在的索引
860. nan_indices = np.isnan(filtered_x3) | np.isnan(filtered_y3)
861.
862. # 获取空数据所在的索引
863. indices = np.where(nan_indices)[0]
```

```
864.
865. # 删除空数据所在的索引对应的元素
866. x3 = np.delete(filtered_x3, indices)
867. y3 = np.delete(filtered_y3, indices)
868. import numpy as np
869. from sklearn.preprocessing import PolynomialFeatures
870. from sklearn.linear_model import LinearRegression
871.
872. # 将 train_data 转换为二维数组
873. x_train = np.array(x3).reshape(-1, 1)
874.
875. # 将 train_label 转换为二维数组
876. y_train = np.array(y3).reshape(-1, 1)
877.
878. best_degree = None
879. best_residual_sum = float('inf')
880.
881. # 尝试不同的多项式阶数
882. for degree in range(1, 4):
883.     # 使用多项式特征转换
884.     poly_features = PolynomialFeatures(degree=degree)
885.     x_train_poly = poly_features.fit_transform(x_train)
886.
887.     # 创建多项式回归模型并进行训练
888.     model = LinearRegression()
889.     model.fit(x_train_poly, y_train)
890.
891.     # 使用模型对训练数据进行预测
892.     y_train_pred = model.predict(x_train_poly)
893.
894.     # 计算残差和
895.     residual_sum = np.sum(np.abs(y_train - y_train_pred))
896.
897.     # 更新最小残差和和对应的阶数
898.     if residual_sum < best_residual_sum:
899.         best_residual_sum = residual_sum
900.         best_degree = degree
901.
902. # 使用最佳阶数重新训练模型
903. poly_features = PolynomialFeatures(degree=best_degree)
904. x_train_poly = poly_features.fit_transform(x_train)
905. model = LinearRegression()
```

```
906.model.fit(x_train_poly, y_train)
907.
908.# 生成测试数据
909.x_test = np.linspace(np.min(x_train), np.max(x_train), 100).resh
ape(-1, 1)
910.x_test_poly = poly_features.transform(x_test)
911.
912.# 使用模型进行预测
913.y_pred = model.predict(x_test_poly)
914.
915.# 绘制真实值和预测值的图像
916.plt.scatter(x_train, y_train, label='Actual')
917.plt.plot(x_test, y_pred, color='red', label='Predicted')
918.plt.xlabel('x')
919.plt.ylabel('y')
920.plt.title('Polynomial Regression (Degree {}) of Class 3'.format(
best_degree))
921.plt.legend()
922.plt.show()
923.
924.# 使用模型对训练数据进行预测
925.y_train_pred = model.predict(x_train_poly)
926.
927.# 计算最终的残差和
928.residual_sum = np.sum(np.abs(y_train - y_train_pred))
929.
930.# 打印最终的残差和
931.print(residual_sum, best_degree)
932.
933.i = 0
934.for idx in idx_3:
935.    resid = np.mean((y_train - y_train_pred)[i:i+7])
936.    resid_list2[idx] = resid
937.    i += 1
938.# 找到空数据所在的索引
939.nan_indices = np.isnan(filtered_x4) | np.isnan(filtered_y4)
940.
941.# 获取空数据所在的索引
942.indices = np.where(nan_indices)[0]
943.
944.# 删除空数据所在的索引对应的元素
945.x4 = np.delete(filtered_x4, indices)
```

```
946.y4 = np.delete(filtered_y4, indices)
947.import numpy as np
948.from sklearn.preprocessing import PolynomialFeatures
949.from sklearn.linear_model import LinearRegression
950.
951.# 将 train_data 转换为二维数组
952.x_train = np.array(x4).reshape(-1, 1)
953.
954.# 将 train_label 转换为二维数组
955.y_train = np.array(y4).reshape(-1, 1)
956.
957.best_degree = None
958.best_residual_sum = float('inf')
959.
960.# 尝试不同的多项式阶数
961.for degree in range(1, 10):
962.    # 使用多项式特征转换
963.    poly_features = PolynomialFeatures(degree=degree)
964.    x_train_poly = poly_features.fit_transform(x_train)
965.
966.    # 创建多项式回归模型并进行训练
967.    model = LinearRegression()
968.    model.fit(x_train_poly, y_train)
969.
970.    # 使用模型对训练数据进行预测
971.    y_train_pred = model.predict(x_train_poly)
972.
973.    # 计算残差和
974.    residual_sum = np.sum(np.abs(y_train - y_train_pred))
975.
976.    # 更新最小残差和和对应的阶数
977.    if residual_sum < best_residual_sum:
978.        best_residual_sum = residual_sum
979.        best_degree = degree
980.
981.# 使用最佳阶数重新训练模型
982.poly_features = PolynomialFeatures(degree=best_degree)
983.x_train_poly = poly_features.fit_transform(x_train)
984.model = LinearRegression()
985.model.fit(x_train_poly, y_train)
986.
987.# 生成测试数据
```

```

988.x_test = np.linspace(np.min(x_train), np.max(x_train), 100).reshape(-1, 1)
989.x_test_poly = poly_features.transform(x_test)
990.
991.# 使用模型进行预测
992.y_pred = model.predict(x_test_poly)
993.
994.# 绘制真实值和预测值的图像
995.plt.scatter(x_train, y_train, label='Actual')
996.plt.plot(x_test, y_pred, color='red', label='Predicted')
997.plt.xlabel('x')
998.plt.ylabel('y')
999.plt.title('Polynomial Regression (Degree {}) of Class 4'.format(
    best_degree))
1000.    plt.legend()
1001.    plt.show()
1002.
1003.    # 使用模型对训练数据进行预测
1004.    y_train_pred = model.predict(x_train_poly)
1005.
1006.    # 计算最终的残差和
1007.    residual_sum = np.sum(np.abs(y_train - y_train_pred))
1008.
1009.    # 打印最终的残差和
1010.    print(residual_sum, best_degree)
1011.    i = 0
1012.    for idx in idx_4:
1013.        resid = np.mean((y_train - y_train_pred)[i:i+9])
1014.        resid_list2[idx] = resid
1015.        i += 1
1016.    from openpyxl import load_workbook
1017.
1018.    # 打开答案文件.xlsx
1019.    workbook = load_workbook('../E 题/数据/表 4-答案文件.xlsx')
1020.
1021.    # 选择活动的工作表
1022.    worksheet = workbook.active
1023.
1024.    # 获取起始单元格的列号和行号
1025.    prob_start_column = 'G'
1026.    start_row = 4
1027.

```



```
1028. # 逐行写入 result 列表的元素
1029. for i in range(len(resid_list2)):
1030.     # 定位
1031.     prob_current_column = prob_start_column
1032.     prob_cell_coordinate = f"{prob_current_column}{start_row
        + i}"
1033.     # 写入
1034.     worksheet[prob_cell_coordinate] = resid_list2[i]
1035.
1036. # 写入内容 1 到指定单元格
1037. for i, value in enumerate(idx_1):
1038.     # 计算目标单元格的行和列索引
1039.     row = 'H'
1040.     column = value + 4
1041.
1042.     # 将内容 1 写入目标单元格
1043.     worksheet[f"{row}{column}"] = 1
1044.
1045. for i, value in enumerate(idx_2):
1046.     # 计算目标单元格的行和列索引
1047.     row = 'H'
1048.     column = value + 4
1049.
1050.     # 将内容 1 写入目标单元格
1051.     worksheet[f"{row}{column}"] = 2
1052.
1053. for i, value in enumerate(idx_3):
1054.     # 计算目标单元格的行和列索引
1055.     row = 'H'
1056.     column = value + 4
1057.
1058.     # 将内容 1 写入目标单元格
1059.     worksheet[f"{row}{column}"] = 3
1060.
1061. for i, value in enumerate(idx_4):
1062.     # 计算目标单元格的行和列索引
1063.     row = 'H'
1064.     column = value + 4
1065.
1066.     # 将内容 1 写入目标单元格
1067.     worksheet[f"{row}{column}"] = 4
1068.
```

```

1069. # 保存工作簿到文件
1070. workbook.save('../E 题/数据/表 4-答案文件.xlsx')
1071. def calculate_difference(row):
1072.     first_valid = row.first_valid_index() # 获取第一个非 NaN
        值的索引
1073.     last_valid = row.last_valid_index() # 获取最后一个非 NaN
        值的索引
1074.     if first_valid is None or last_valid is None: # 如果没有
        非 NaN 值, 则返回 NaN
1075.         return np.nan
1076.     else:
1077.         return row[first_valid] - row[last_valid] # 计算差值
1078. Y_train['水肿体积差
    '] = Y_train.apply(calculate_difference, axis=1)
1079. table_1_issue2 = pd.concat([table_1.iloc[:100, -10:-
    3], Y_train['水肿体积差']], axis=1)
1080. table_1_issue2
1081. # 计算相关系数矩阵
1082. corr_matrix = table_1_issue2.corr()
1083. # 创建遮盖矩阵
1084. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1085. # 绘制相关系数矩阵的热力图
1086. plt.figure(figsize=(10, 8))
1087. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1088. plt.title('Correlation Matrix')
1089. plt.show()
1090. table_2_issue2 = table_2_issue1.iloc[:100, 1::2]
1091. table_2_issue2['血肿体积差
    '] = table_2_issue2.apply(calculate_difference, axis=1)
1092. table_issue2 = pd.concat([table_1.iloc[:100, -10:-
    3], table_2_issue2['血肿体积差'], Y_train['水肿体积差']], axis=1)
1093. table_issue2
1094. # 计算相关系数矩阵
1095. corr_matrix = table_issue2.corr()
1096.
1097. # 绘制相关系数矩阵的热力图
1098. plt.figure(figsize=(10, 8))
1099. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True)
1100. plt.title('Correlation Matrix')
1101. plt.show()

```

```
1102. train_data
1103. table_1.shape
1104. table_11 = pd.concat([table_1.iloc[:, 1:-3], table1['90 天
    mRS']], axis=1)
1105. import pandas as pd
1106. import seaborn as sns
1107. import matplotlib.pyplot as plt
1108. import numpy as np
1109.
1110. # 创建数据
1111. data = table_11.iloc[:, :]
1112.
1113. # 计算相关系数矩阵
1114. corr_matrix = data.corr()
1115.
1116. # 创建遮盖矩阵
1117. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1118.
1119. # 绘制相关矩阵
1120. plt.figure(figsize=(16, 8))
1121. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1122. plt.title('Correlation Matrix')
1123. plt.show()
1124. # 读表1 的数据
1125. table_1_issue3 = table_11[['年龄', '糖尿病史', '冠心病史', '吸
    烟史', '饮酒史', '脑室引流', '90 天 mRS']].iloc[:, :]
1126. table_1_issue3
1127. table_22 = pd.concat([table_2.iloc[:, 2:24], table_1_issue3['
    90 天 mRS']], axis=1)
1128. table_22
1129.
1130. # 计算相关系数矩阵
1131. corr_matrix = table_22.corr()
1132. # 创建遮盖矩阵
1133. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1134. # 绘制相关系数矩阵的热力图
1135. plt.figure(figsize=(10, 8))
1136. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1137. plt.title('Correlation Matrix')
1138. plt.show()
```

```

1139. # 读表1 的数据
1140. table_2_issue3 = table_22[['HM_volume', 'ED_volume']]
1141. table_2_issue3
1142. # 计算相关系数矩阵
1143. corr_matrix = table_3.iloc[:160, 2:].corr()
1144. # 创建遮盖矩阵
1145. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1146. # 绘制相关系数矩阵的热力图
1147. plt.figure(figsize=(10, 8))
1148. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1149. plt.title('Correlation Matrix')
1150. plt.show()
1151. table_3_issue3 = table_3[['NCCT_original_firstorder_90Percent
    ile', 'NCCT_original_firstorder_InterquartileRange', 'NCCT_origin
    al_firstorder_MeanAbsoluteDeviation', 'NCCT_original_firstorder_R
    ange', 'NCCT_original_firstorder_RobustMeanAbsoluteDeviation']]
    .iloc[:160, :]
1152. table_3_issue3.shape
1153. train_data3, test_data3 = table_issue3.iloc[:100, :-
    1], table_issue3.iloc[100:, :-1]
1154. train_label3, test_label3 = table_issue3.iloc[:100, -
    1:], table_issue3.iloc[100:, -1:]
1155. train_data3.shape, test_data3.shape, train_label3.shape, test
    _label3.shape
1156. from sklearn.linear_model import LogisticRegression
1157. from sklearn.metrics import accuracy_score
1158.
1159. # 创建逻辑回归模型对象
1160. logreg = LogisticRegression(multi_class='ovr')
1161.
1162. # 训练逻辑回归模型
1163. logreg.fit(train_data3, train_label3)
1164.
1165. # 在测试集上进行预测
1166. logreg_pred = logreg.predict(train_data3)
1167.
1168. # 计算逻辑回归模型的精度
1169. logreg_accuracy = sum(train_label3.iloc[:, 0] == logreg_pred)
    / len(logreg_pred)
1170. print("逻辑回归模型精度: ", logreg_accuracy)
1171. from sklearn.ensemble import RandomForestClassifier

```

```
1172.
1173. # 创建随机森林模型
1174. rf = RandomForestClassifier(n_estimators=8)
1175.
1176. # 训练随机森林模型
1177. rf.fit(train_data3, train_label3)
1178.
1179. # 在测试集上进行预测
1180. rf_pred = rf.predict(train_data3)
1181. # 计算随机森林模型的精度
1182. rf_accuracy = sum(train_label3.iloc[:, 0] == rf_pred) / len(r
    f_pred)
1183. print("随机森林模型精度: ", rf_accuracy)
1184. from sklearn.neighbors import KNeighborsClassifier
1185. from sklearn.metrics import accuracy_score
1186.
1187. # 创建K 临近算法模型对象
1188. knn = KNeighborsClassifier(n_neighbors=2)
1189.
1190. # 训练K 临近算法模型
1191. knn.fit(train_data3, train_label3)
1192.
1193. # 在测试集上进行预测
1194. knn_pred = knn.predict(train_data3)
1195.
1196. # 计算K 临近算法模型的精度
1197. knn_accuracy = sum(train_label3.iloc[:, 0] == knn_pred) / len
    (knn_pred)
1198. print(f"K 临近算法模型精度: ", knn_accuracy)
1199. mRS_list = rf.predict(table_issue3.iloc[:, :-1])
1200. mRS_list
1201. from openpyxl import load_workbook
1202.
1203. # 打开答案文件.xlsx
1204. workbook = load_workbook('../E 题/数据/表 4-答案文件.xlsx')
1205.
1206. # 选择活动的工作表
1207. worksheet = workbook.active
1208.
1209. # 获取起始单元格的列号和行号
1210. prob_start_column = 'I'
1211. start_row = 4
```

```
1212.
1213.     # 逐行写入 result 列表的元素
1214.     for i in range(len(mRS_list)):
1215.         # 定位
1216.         prob_current_column = prob_start_column
1217.         prob_cell_coordinate = f"{prob_current_column}{start_row
+ i}"
1218.         # 写入
1219.         worksheet[prob_cell_coordinate] = mRS_list[i]
1220.
1221.
1222.     # 保存工作簿到文件
1223.     workbook.save('../E 题/数据/表 4-答案文件.xlsx')
1224.     table = pd.concat([table_1.iloc[:100, :], table_2.iloc[:100,
1:], table_3.iloc[:100, :]], axis=1)
1225.     table.shape
1226.     table = table.fillna(table.mean())
1227.     # 数据清洗
1228.     # 使用均值进行缺失值填充
1229.     table.fillna(table.mean(), inplace=True)
1230.     # 使用 Z-score 方法检测和处理异常值
1231.     from scipy import stats
1232.     z_scores = stats.zscore(table)
1233.     df = table[(z_scores < 3).all(axis=1)]
1234.
1235.     null_counts = df.isnull().sum()
1236.     print(null_counts)
1237.     print(table.dtypes)
1238.     table_1.head()
1239.     # 创建数据
1240.     data = table_1[['年龄', '性别', '吸烟史', '饮酒史
']].assign(mRS=train_label3) # 假设 train_data 的标签位于第一列
1241.
1242.     # 计算相关系数矩阵
1243.     corr_matrix = data.corr()
1244.
1245.     # 创建遮盖矩阵
1246.     mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1247.
1248.     # 绘制相关矩阵
1249.     plt.figure(figsize=(16, 8))
```

```
1250. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1251. plt.title('mRS 和个人史相关矩阵')
1252. plt.show()
1253. table_1 = table_1.drop(columns=['吸烟史', '饮酒史'])
1254. table_1
1255. # 创建数据
1256. data = table_1.iloc[:100, 4:9].assign(mRS=train_label3) # 假
    设train_data 的标签位于第一列
1257.
1258. # 计算相关系数矩阵
1259. corr_matrix = data.corr()
1260.
1261. # 创建遮盖矩阵
1262. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1263.
1264. # 绘制相关矩阵
1265. plt.figure(figsize=(16, 8))
1266. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1267. plt.title('mRS 和疾病史相关矩阵')
1268. plt.show()
1269. # 创建数据
1270. data = table_1.iloc[:100, 10:-
    3].assign(mRS=train_label3) # 假设train_data 的标签位于第一列
1271.
1272. # 计算相关系数矩阵
1273. corr_matrix = data.corr()
1274.
1275. # 创建遮盖矩阵
1276. mask = np.triu(np.ones_like(corr_matrix, dtype=bool))
1277.
1278. # 绘制相关矩阵
1279. plt.figure(figsize=(16, 8))
1280. sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', square=
    True, mask=mask)
1281. plt.title('mRS 和治疗方法相关矩阵')
1282. plt.show()
1283. table_2
```