

2023 年第八届“数维杯”大学生 数学建模挑战赛论文

基于 BP 神经网络对宫内节育器的生产策略

摘 要

节育环是我国女性使用的最普遍的节育方式，不同型号的节育环对于人群的适用情况不同。为了探究某公司的新生产的两种节育环是否适合投入生产，VCu260、VCu380 与已被临床应用的 MCu 节育器一起进行临床试验。

对于问题 1，首先从不同视角统计分析两院的主诉情况，使用 pandas 数据包深度挖掘了 Mcu、VCu260 和 VCu380 组的 8 种不适情况，各种分布结果均表明一院和二院的临床数据是具有显著性差异的。其次，对两个医院临床受试者的身体指标与节育器的理化指标进行统计对比分析，可得“既往应用节育器情况”和“宫颈扩张情况”是导致这种差异的主要因素。

对于问题 2，从清洗后的数据选出身体指标(年龄、初潮年龄、月经周期、月经经期、宫腔深度)，5 种指标分布结果表明每个医院使用三种节育器的人群身体状况相同；同时，对这 5 种指标进行相关性分析，相关系数矩阵表明选择的身体指标不存在强相关。为了量化不适情况，计算了 1、3、6、12 月和整体的不适情况，分别建立了一院、二院的多元线性回归模型，综合分析模型的显著性检验结果，表明初潮年龄和月经周期对出现不适状况的影响较大，年龄和月经经期不显著对出现不适情况的影响较小。

对于问题 3，根据受访者的主诉情况判断节育器质量的好坏，从而构建节育器质量好坏的标准。同时使用清洗后的数据选出身体指标和节育器的理化指标作为输入特征，节育器的适用情况指标作为输出特征；然后，分别使用一院和二院的对照组(MCu 节育器)数据建立 BP 神经网络预测模型对输入和输出特征进行训练，得到的节育器质量模型。根据训练好的模型预测实验组(VCu260、VCu380)的结果可知 VCu260 节育器在人群中的适应情况要优于 VCu380 节育器。

对于问题 4，基于问题 3 建立的节育器质量预测模型，通过可视化隐藏层权重系数得到了分析影响节育器质量的决定性因素。实验结果表明应用节育器情况、宫腔深度、使用节育器型号和宫颈扩张情况是影响宫内节育器质量的决定因素。

关键词：多元线性回归；BP 神经网络；相关性分析；主控因素

目 录

一、 问题重述	2
1.1 背景介绍	2
1.2 要解决的问题	2
二、 问题分析	2
2.1 问题 1 的分析	3
2.2 问题 2 的分析	3
2.3 问题 3 的分析	3
2.4 问题 4 的分析	4
三、 模型假设	4
四、 定义与符号说明	4
五、 模型的建立与求解	4
5.1 问题 1 的模型建立与求解	6
5.1.1 问题 1 模型的建立	6
5.1.2 问题 1 模型的求解	6
5.1.3 结果	10
5.2 问题 2 的模型建立与求解	10
5.2.1 问题 2 模型的建立	10
5.2.2 问题 2 模型的求解	13
5.2.3 结果	15
5.3 问题 3 的模型建立与求解	15
5.3.1 问题 3 模型的建立	15
5.3.2 问题 3 模型的求解	16
5.3.3 结果	19
5.4 问题 4 的模型建立与求解	19
六、 模型的评价及优化	20
6.1 误差分析	20
6.1.1 针对于问题 1 的误差分析	20
6.1.2 针对于问题 2 的误差分析	20
6.1.3 针对于问题 3 的误差分析	21
6.1.4 针对于问题 4 的误差分析	21
6.2 模型的优点	21
6.3 模型的缺点	21
6.4 模型的推广	21
参考文献	22
附 录	23

一、问题重述

1.1 背景介绍

目前我国育龄妇女的主要避孕措施是宫内节育器（IUD）。节育器，由于初期使用的装置多是环状的，又叫节育环，是一种放置在子宫腔内的避孕装置。这是一种相对安全、有效、经济、可逆、简便的节育器具。据悉，我国约 70% 妇女选用 IUD 作为避孕方法，占世界 IUD 避孕总人数的 80%。但是节育器的种类和形状有很多，它在给女性带来方便的同时，也可能会引起很多不适应的症状，故节育器的使用需要经过临床试验。

1.2 要解决的问题

问题一：根据附件 1 与附件 2，分析两个医院的临床数据有无显著性差异，若存在显著性差异，对导致这种差异的因素进行分析。

问题二：结合附件 1 与附件 2，分析受试者的身体指标与随访主诉情况的联系，并说明受试者的身体指标是否是受试者出现不适状况的主要因素。

问题三：根据受试者的身体指标、节育器的理化指标与随访时的主诉情况，建立节育器质量模型，并分析 VCu260 与 VCu380 记忆型宫内节育器的质量哪个更优，更适合生产。

问题四：结合问题三，根据建立的节育器质量模型，探究影响宫内节育器质量的决定因素。

二、问题分析

宫内节育器的种类很多，国内常用的有金属单环、麻花环、混合环、节育环、T 形环等。某公司研发了两种型号的 VCu 记忆型宫内节育器，采用镍钛记忆合金丝支架，除具有独特形状记忆功能外，还具有抗腐蚀、耐磨损、超弹性和对身体的副作用较小等优点。但是节育器也可能会引起疼痛、不适、脱落或者出血等症状。为了探究这两种型号的节育器是否适合投入生产，特与已被临床应用的 MCu 功能性宫内节育器一起临床试验。

2.1 问题 1 的分析

问题 1 要求根据附件一与附件二的数据,分析两个医院的临床数据有无显著性差异,若存在显著性差异,对导致这种差异的因素进行分析。由于人为记录数据可能会存在部分误差和异常,这会影响后续过程分析,所以需要对附件一和附件二中的数据进行清洗。利用清洗后的数据使用统计方法结合 Pandas 数据包进行分析,分别统计一院、二院中的床受试者和节育器的基本数据以及主诉情况,对比判断存在的显著性差,并分析得出导致这种差异的因素。

2.2 问题 2 的分析

问题 2 要求结合附件 1 与附件 2,分析受试者的身体指标与随访主诉情况的联系,并说明受试者的身体指标是否是受试者出现不适状况的主要因素。首先从附件一提取出身体指标(年龄、初潮年龄、月经周期、月经经期、宫腔深度),由于数据来源于一院、二院的实验组和对照组,需要判定一、二院人群的身体指标是否分布相同,相同可以整体分析,否则需要分组分析;其次,需要对选定的身体指标进行相关性分析,判断所选指标中是否有强相关的指标需要剔除;再者,统计分析附件二中主诉情况,量化不适情况指标;最后,需要建立多元统计分析模型,并检验模型的有效性,对模型诊断和评价,探究影响不适状况的主要因素。

2.3 问题 3 的分析

问题三是根据受试者的身体指标、节育器的理化指标与随访时的主诉情况,建立节育器质量模型并分析两种节育器(VCu260、VCu380)的质量哪个更优。首先需要构建节育器质量好坏的标准,根据受访者的主诉情况(8 种不适情况,即“有不适人数”指标)判断节育器质量的好坏;其次,需要对附件一和附件二中的失访数据进行筛选,对无法判断的数据删除,确定体指标和节育器的理化指标作为输入特征,1、3、6、12 月的“有不适人数”指标作为输出特征;最后,建立预测模型对输入和输出特征进行训练,使用对照组(MCu)的数据集进行建模,再对建立的数学模型进行评估,检查模型的拟合程度和预测效果。利用训练好的模型预测实验组(VCu260、VCu380)的“有不适人数”指标,

计算节育器是否好坏，并将实际结果与预测结果比较得出结论。

2.4 问题 4 的分析

问题 4 要求根据问题 3 建立的节育器质量模型，对影响宫内节育器质量的决定因素进行分析。理想状态下，问题 3 建立的节育器质量预测模型是最优的，能够根据决定性因素对节育器的质量做出判断，可以对模型输入特征的系数权重进行可视化，找出权重较大的特征，从而分析影响节育器质量的决定性因素。

三、模型假设

1. 假设受访人的不适情况不受外界其他因素的影响；
2. 假设节育器的使用情况只和使用者的身体指标和产品质量有关；
3. 假设节育器的质量问题依靠受试者的不适情况来反映的；
4. 一院、二院临床受试者除了身体指标外不存在其他差异；

四、定义与符号说明

符号定义	符号说明	符号定义	符号说明
n	样本数	H_i	神经元的输出
m	特征的个数	ω_{ij}	权重
F	模型	a_j	偏置
ε	误差向量	Y_k	输出结果
b	回归系数向量	b_k	偏置
X	设计矩阵	η	学习率
x_i	自变量	e_k	误差
y_i	因变量	f	网络层
\hat{b}	回归系数		

注：未列出符号及重复的符号以出现处为准。

五、模型的建立与求解

数据的预处理：

根据题目需求描述，需要对身体指标、节育器的理化指、两个医院的主诉情况记录进行分析与建模，说明影响不适状况的因素，并找出影响宫内节育器质量的决定因素。考虑到实际情况的不确定性，附件一和附件二的数据存在部分误差，以及失访情况下造成的数据缺失，需要对原始数据进一步处理方可使用。

（1）错误数据：在统计失访下的主诉情况时，二院的部分数据已失访但是主诉情况却有记录，无法确认是否有效，予以删除、

（2）疑惑数据：在统计附件一中的节育器型号情况时，有的受访人员使用了两种型号的节育器，便于后续建模选用较大型号进行分析；关于附件二的VCu260组受访人的主诉情况，“有不适人数”统计指标存在争议，为了和其他几个对照组统一，根据1、3、6、12月的8种不适情况进行计算，替代原始的数据。

（3）缺失数据：

1、附件一的处理：关于“既往应用节育器情况”这个指标，共包含IUD、无和其他3种情况，为了后续建模需要对这三种进行打标签，IUD对应1、无对应2、其他对应3，统一成单个指标；“使用节育器型号情况”指标的处理方式同上，小对应1、中对应2、大对应3，处理后的数据如下：

表 5-1 预处理后的数据

序号	组别	年龄	初潮年龄	月经周期	月经经期	节育器情况	宫腔深度	节育器型号
53	1	38	13	32	4	1	7	2
55	1	35	16	32	3	3	6	1
58	1	25	14	33	4	2	7	2
59	1	30	14	32	4	3	8	2

2、附件二的处理：对于单个受访人员，如果在某月没有失联，当月8种不适情况的空白处均填补为0，代表没有出现不适情况；失联对应当月8种不适情况的空白处全部替换为Nan，不参与计算。

经过预处理后的数据符合后续处理的需求，同时也很好的代表了原始数据的特征，使用处理后的数据进行建模较为合理。（完整数据详见支撑材料）

（4）节育器好坏指标：

根据附件二中的“有不适人数”进行判断处理，共有1、3、6、12这个四个月的不适指标，1代表这个月有不适状况，0代表这个月无不适状况。根据节育器的判别标准：前期出现不适症状，但后期症状消失，可认为受试者适应该节育器，即节育器质量没什么问题。

为了更合理的建立模型，这里我们判定超出 3 个月不会出现不适情况即认定该节育器安全(0/好，1/坏)，所以通过分析 1、3、6、12 月的不适情况可以得出：

表 5-2 节育器好坏评价指标

一月	三月	六月	十二月	是否安全
1	0	0	0	好
1	1	0	0	好
0	1	1	0	好
1	1	1	0	坏
1	0	1	0	坏
0	0	1	1	坏

5.1 问题 1 的模型建立与求解

5.1.1 问题 1 模型的建立

我们需要解决的问题是分析一院和二院的临床数据是否有显著性差异，如果有显著性差异，分析导致差异的原因。通过对附件一和附件二中数据的观察，两个医院的临床数据不仅包含了三组节育器的使用情况，同时还包含了不适情况的受试者随访记录。为了全面的分析两院的临床数据，具体分析步骤如下：

(1) 根据附件二中的数据，对一院和二院中分别使用三种节育器出现不适情况的受试者人数进行统计分析，观察两个医院使用这三种节育器的受试者出现不适情况是否一致；

(2) 统计附件二中的 8 种不适情况(脱落、因症取出、怀孕、非月经期出血、疼痛、经量多、分泌物增多、经期/周期异常)，分别分析一院和二院中使用三种节育器的受试者人数的相似规律；

(3) 再对分组后的数据进一步的挖掘，分别统计每个医院使用三种节育器的受试者在 1、3、6、12 月出现不适情况；

(4) 根据步骤 1、2、3 的统计分析可以得出一院和二院的临床数据是否有显著性差异，若有显著性差异，继续步骤 5；否则，两个医院临床数据不具有显著性差异。

(5) 对附件一中的数据与附件二中的数据进行相关性分析，观察是哪几种因素导致了两个医院的临床数据有显著性差异。

5.1.2 问题 1 模型的求解

通过对附件 2 中一、二院临床受试者和节育器的基本数据进行分析，统计出随访中 3 种节育器不适情况的人数，得出不同的层次数据的统计结果，如下所示：

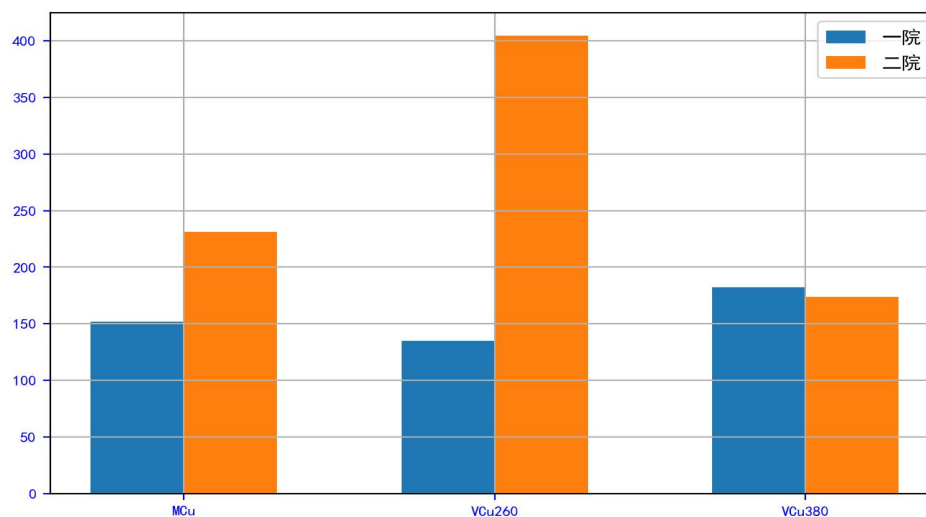


图 5-1 3 种节育器出现不适情况的统计结果

根据图 5-1 可以看出，一院和二院的临床数据在出现不适情况的统计中是存在一定差异的，因此可以进行进一步分析：

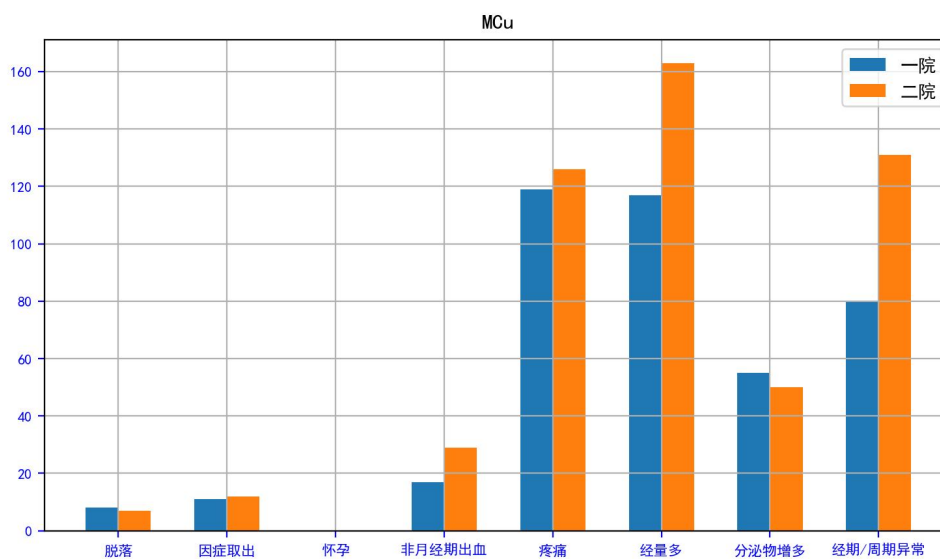


图 5-2 八种不适情况在第一组(MCu 功能性宫内节育器)中统计结果

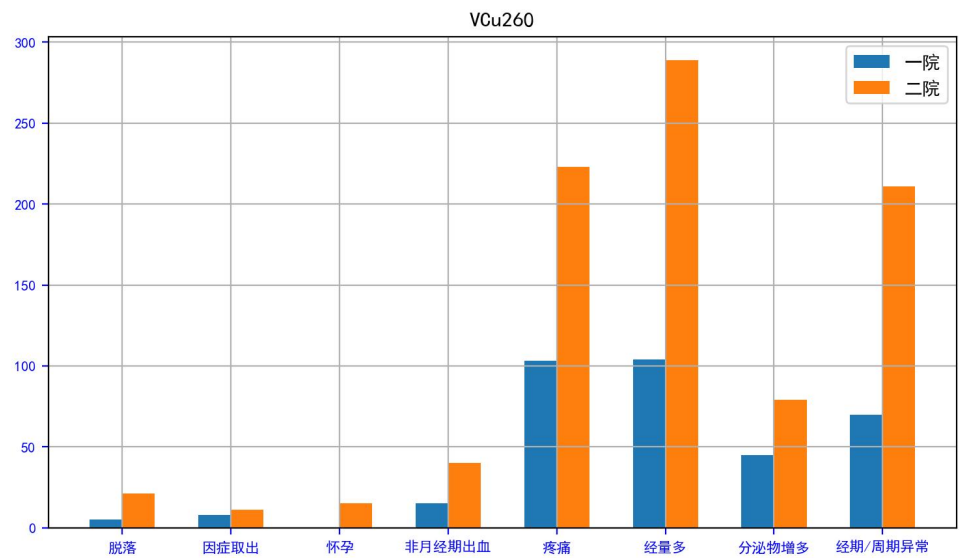


图 5-3 八种不适情况在第二组(VCu260 记忆型宫内节育器)中统计结果

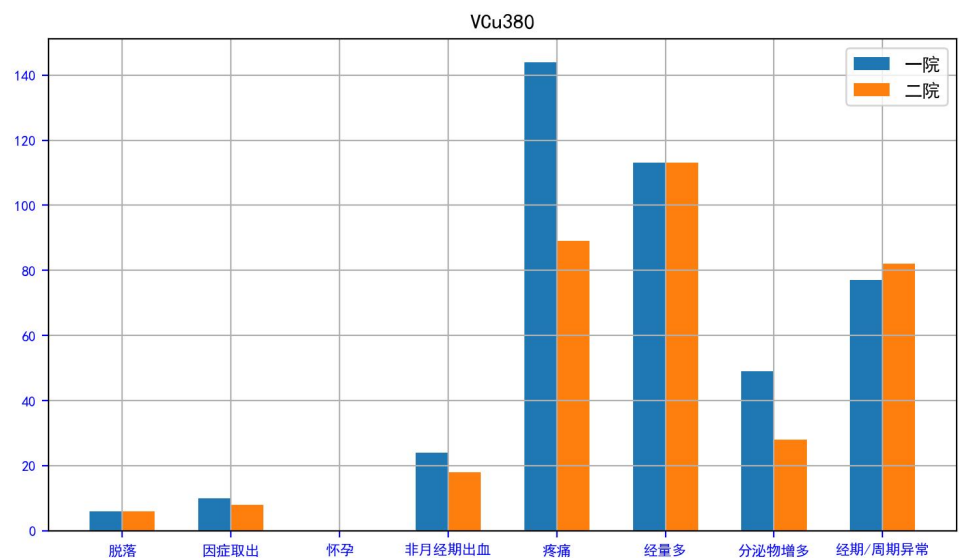


图 5-4 八种不适情况在第三组(VCu380 记忆型宫内节育器)中统计结果

通过图 5-2 到图 5-4 可以明显发现在疼痛、经量多、经期异常这几个特征一院、二院存在较大的差异，且在不同组别中也差距较大。

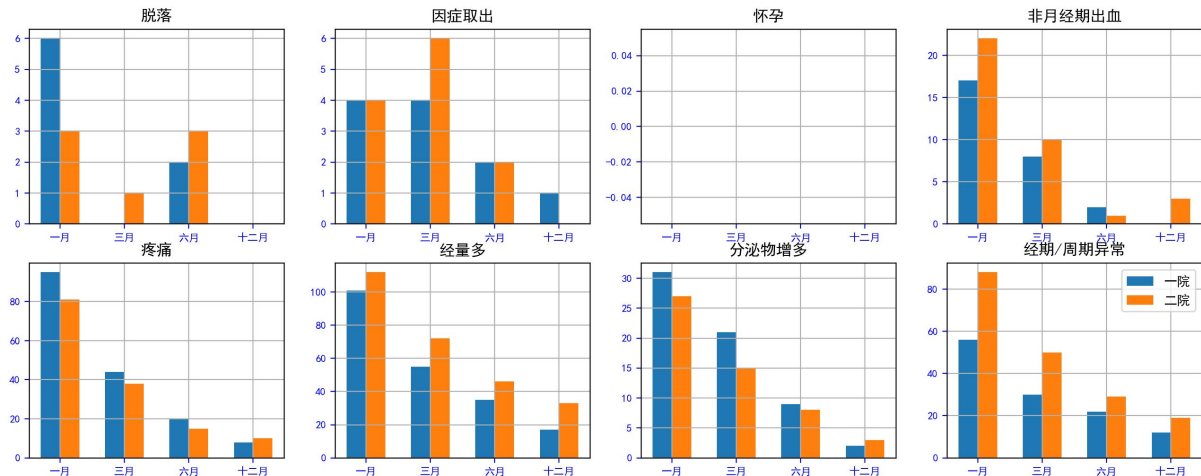


图 5-5 第一组(MCu 功能性宫内节育器)出现不适情况统计结果

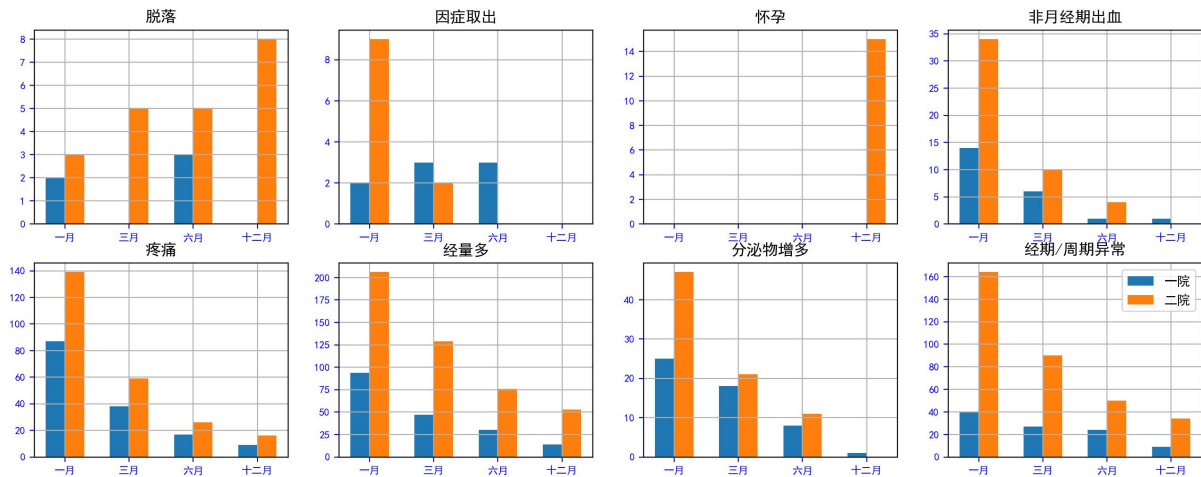


图 5-6 第二组(VCu260 记忆型宫内节育器)出现不适情况统计结果

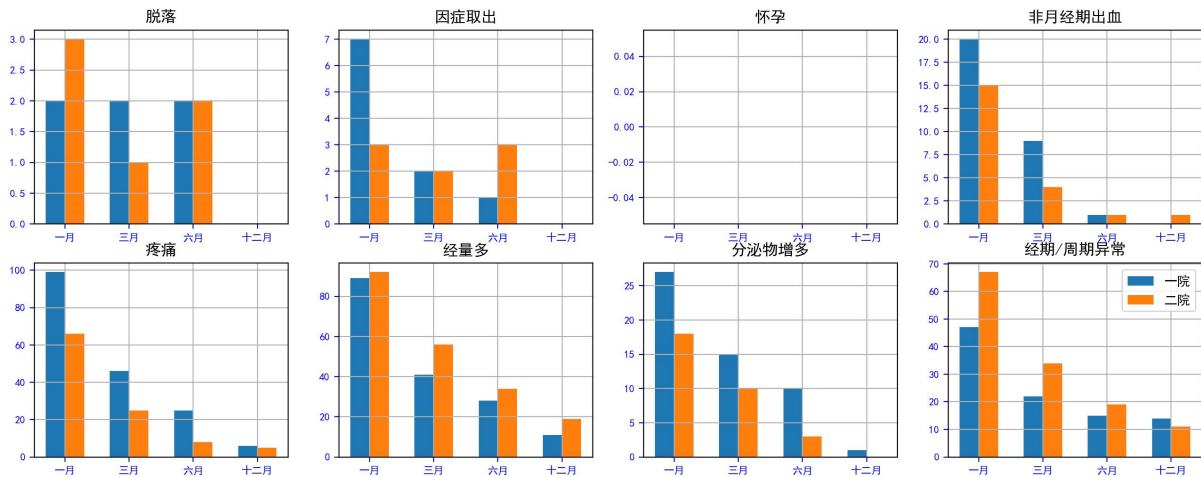


图 5-7 第三组(VCu380 记忆型宫内节育器)出现不适情况统计结果

对图 5-5 到图 5-7 分析，可以发现一院和二院的临床数据在时间月份上也存在显著性差异。同时，结合附件一中临床受试者的身体指标与节育器的理化指标对这种差异进一步分析，其统计指标如下：

表 5-3 一院的身体指标与节育器的理化指标分布

	年龄	初潮年龄	月经周期	月经经期	应用情况	宫腔深度	节育器型号	宫颈扩张情况
mean	31.8	14.2	30.3	4.3	1.8	7.5	2.03	0.47
std	5.2	1.26	1.57	0.79	0.88	0.79	0.43	0.49
min	20	12	26	3	1	6	1	0
max	40	18	34	7	3	9	3	1

表 5-4 二院的身体指标与节育器的理化指标分布

	年龄	初潮年龄	月经周期	月经经期	应用情况	宫腔深度	节育器型号	宫颈扩张情况
mean	30.6	14.14	30.23	4.63	1.86	7.30	2.8	0.04
std	4.07	1.21	1.32	0.82	0.96	0.69	0.58	0.24
min	22	12	28	3	1	6	0	0
max	40	18	35	7	3	9	3	1

由表 5-3 与 5-4 对比分析可知，两院的临床受试者身体指标整体分布差异不大，其年龄、初潮年龄、月经周期的分布基本一致，月经经期、宫腔深度、使用节育器型号的分布有点差距，既往应用节育器情况和宫颈扩张情况差距最大。因此，可以判断出既往应用节育器情况和宫颈扩张情况是引起两院之间数据出现显著性差异的因素。

5.1.3 结果

通过对模型的建立和求解，利用 pandas 数据包对筛选后的数据进行初步运算。首先通过 pandas 对数据筛选处理，得出数据的统计性结果。然后根据统计性的结果判断一院和二院的临床数据是具有显著性差异的，其次通过对两院的临床数据进行统计对比分析，可得“既往应用节育器情况”和“宫颈扩张情况”是导致这种差异的主要因素。

5.2 问题 2 的模型建立与求解

5.2.1 问题 2 模型的建立

(1) 数据的选择

为了探明受试者的身体指标是否是受试者出现不适状况的主要因素，首先从附件一中提取需要的数据，所选择的身体指标为年龄、初潮年龄、月经周期、月经经期、宫腔

深度。而受访人群来自于二院和四院，且对使用MCu、VCu260 和VCu380 三种节育器的人群进行分组，人群的身体指标分布情况如下：

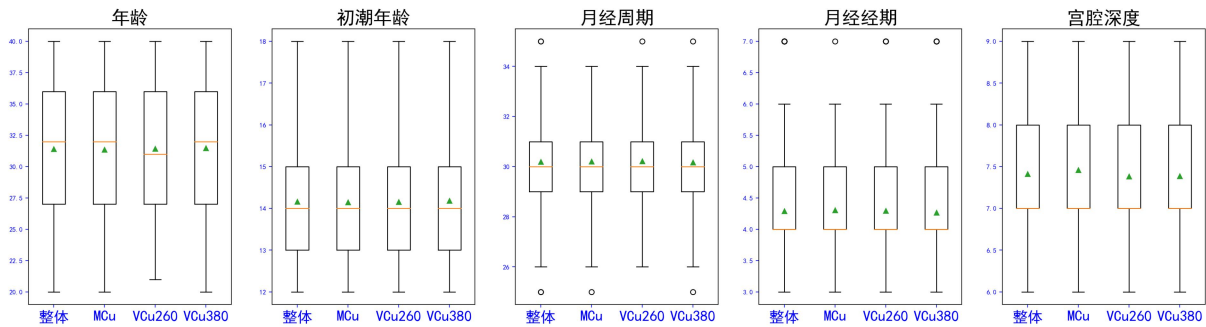


图 5-8 一院的各组人群的身体指标分布

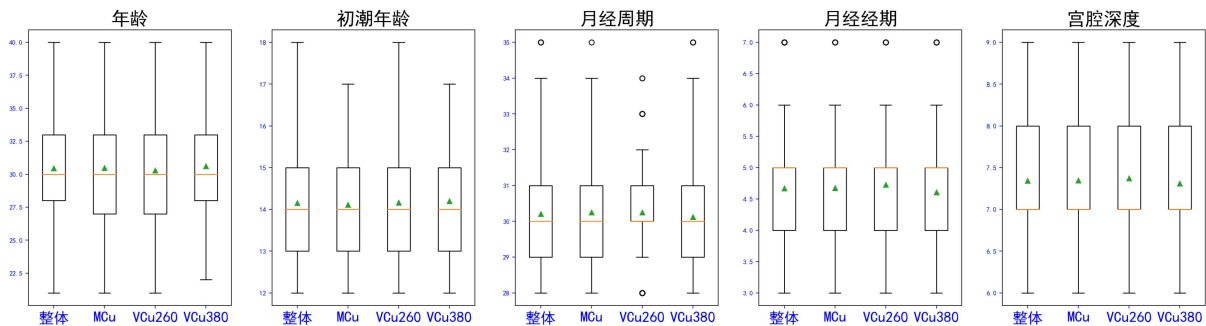


图 5-9 二院的各组人群的身体指标分布

从两个医院的各组人群的身体指标分布情况可以看出，整体的分布情况和各组的分布情况基本相同，固可以把一院、二院作为整体进行后续分析。同时，我们对年龄、初潮年龄、月经周期、月经经期、宫腔深度这 5 种指标做了相关性分析，如下所示：

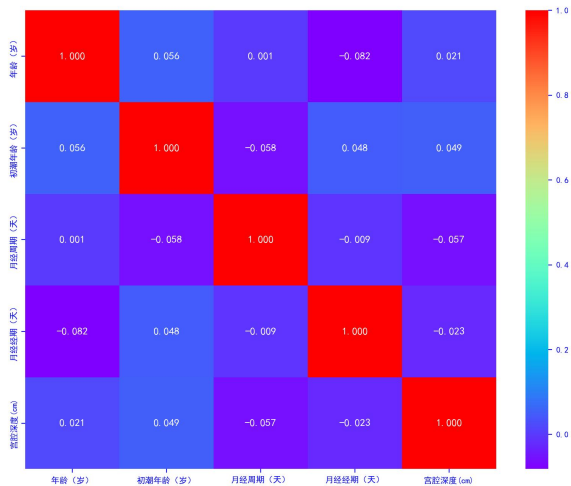


图 5-10 一院身体指标的相关系数矩阵

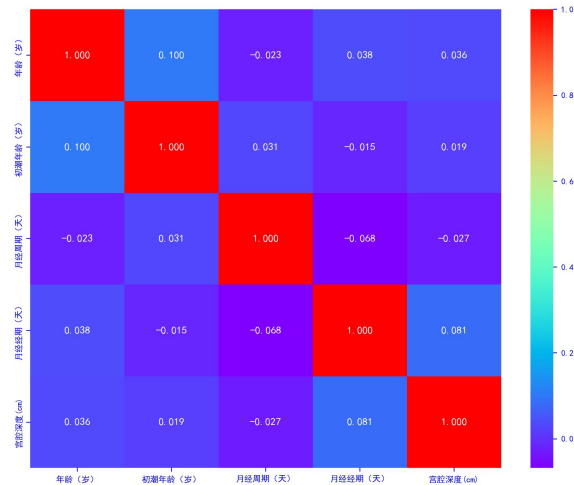


图 5-11 二院身体指标的相关系数矩阵

由一院和二院的身体指标的相关性矩阵热力图可以看出，这 5 种指标之间的相关性很小，两两之间没有相关系数都低于 0.1，不需要进一步处理。

同时，为了量化受试者的不适状况，这里我们统计了每个人的脱落、因症取出、怀孕、非月经期出血、疼痛、经量多、分泌物增多、经期/周期异常症状，计算了 1、3、6、12 月每个月内的症状个数，以及整体时间的症状个数，具体如下：

表 5-5 一院受试者的不适状况统计结果

序号	一月	三月	六月	十二月	总数
10	4	0	0	0	4
17	3	2	2	0	7
36	0	2	0	0	2
49	3	3	0	0	6

(2) 模型的建立

多元线性回归的是一元线性模型的拓展，其基本原理也是建立一个线性模型 F ，只不过使用多个自变量(特征)来预测一个因变量(目标)^[1]。假设有 n 个样本，其中第 i 个样本的自变量为 $x_i = (x_{i1}, x_{i2}, \dots, x_{im})$ ，这里 m 为特征的个数，因变量为 y_i ，模型为 F ，则多元线性回归模型可表达如下：

$$y_i = b_0 + b_1x_{i1} + b_2x_{i2} + \dots + b_mx_{im} + \varepsilon_i \tag{1}$$

其中 $b_0, b_1, b_2, \dots, b_m$ 代表回归系数， ε_i 是误差。

在多元线性回归中，回归系数表示对应自变量的影响大小，误差项表示模型无法解释的随机噪声^[2]。多元线性回归的目标是求解使得误差最小的回归系数，通常使用最小二乘法来求解回归系数，其中最小二乘法的思想是最小化所有样本误差的平方和^[3-5]。使

用矩阵形式表述，则可将多元线性回归表示为：

$$y = Xb + \varepsilon \quad (2)$$

其中， $y = (y_1, y_2, \dots, y_n)^T$ ， X 是 $n \times (m+1)$ 的设计矩阵，其中第一列全为常数 1，表示截距项。 $b = (b_0, b_1, b_2, \dots, b_m)^T$ 表示回归系数向量， $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$ 表示误差向量。

最小二乘法的目标是求解 b 使得样本误差 ε 的平方和最小：

$$\hat{b} = \arg \min_b \sum_{i=1}^n (y_i - x_i^T b)^2 \quad (3)$$

求解出回归系数 \hat{b} 后，可以利用多元线性回归模型进行分析。

5.2.2 问题 2 模型的求解

为了讨论身体指标是否和不适状况有关联，这里取身体指标年龄 x_1 、初潮年龄 x_2 、月经周期 x_3 、月经经期 x_4 、宫腔深度 x_5 这 5 个指标作为多元线性模型的输入，1 月 y_1 、3 月 y_2 、6 月 y_3 、12 月 y_4 和总月 y_5 的症状个数作为模型的输出。分别建立身体指标和每个月不适情况个数之间的线性回归模型，下面分别介绍一院和二院的求解情况：

(1) 一院的结果

根据预处理后的数据，建立身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和不适情况总数 y_5 之间的线性回归模型：

$$y_5 = 0.99 - 0.104x_1 + 0.064x_2 + 0.44x_3 - 0.078x_4 + 0.19x_5 \quad (4)$$

为了更加全面的评估模型，使用显著性检验判断身体指标和不适情况之间的关系是否显著，确认模型是否满足基本假设，是否存在问题，这里主要使用 P 值和 F 统计量，具体结果如下：

表 5-6 一院模型的评价表

组别	P					F
	x1	x2	x3	x4	x5	
总数	0.218	0.337	0.001	0.51	0.442	2.82

同时，为了探究身体指标和出现不适情况的时间之间联系，我们还分别建立了下面四种回归模型：

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和一月不适情况 y_1 之间的线性回归模型：

$$y_1 = 3.55 - 0.042x_1 - 0.048x_2 + 0.15x_3 - 0.04x_4 - 0.12x_5 \quad (5)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和三月不适情况 y_5 之间的线性回归模型:

$$y_2 = -0.903 - 0.044x_1 + 0.063x_2 + 0.138x_3 + 0.025x_4 - 0.04x_5 \quad (6)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和六月不适情况 y_3 之间的线性回归模型:

$$y_3 = -0.326 - 0.038x_1 + 0.016x_2 + 0.115x_3 - 0.023x_4 + 0.245x_5 \quad (7)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和十二月不适情况 y_4 之间的线性回归模型:

$$y_4 = -1.33 - 0.02x_1 + 0.032x_2 + 0.034x_3 - 0.039x_4 + 0.102x_5 \quad (8)$$

他们的显著性检验统计如下:

表 5-7 一院模型的评价表

组别	P					F
	x1	x2	x3	x4	x5	
一月	0.346	0.169	0.031	0.513	0.341	1.629
三月	0.302	0.063	0.041	0.676	0.744	1.908
六月	0.345	0.625	0.073	0.682	0.037	1.739
十二月	0.445	0.116	0.4	0.286	0.181	1.402

(2) 二院的结果

同理, 可建立二院的 身体指标和不适状况之间的回归模型, 具体如下:

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和不适情况总数 y_5 之间的线性回归模型:

$$y_5 = 5.1 - 0.029x_1 - 0.093x_2 - 0.316x_3 - 0.015x_4 + 0.044x_5 \quad (9)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和一月不适情况 y_1 之间的线性回归模型:

$$y_1 = 2.17 - 0.094x_1 - 0.013x_2 - 0.089x_3 - 0.013x_4 + 0.016x_5 \quad (10)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和三月不适情况 y_5 之间的线性回归模型:

$$y_2 = 2.51 - 0.021x_1 + 0.061x_2 + 0.063x_3 + 0.008x_4 - 0.059x_5 \quad (11)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和六月不适情况 y_3 之间的线性回归模型:

$$y_3 = 0.74 - 0.001x_1 - 0.019x_2 - 0.005x_3 + 0.026x_4 + 0.096x_5 \quad (12)$$

身体指标 $(x_1, x_2, x_3, x_4, x_5)$ 和十二月不适情况 y_4 之间的线性回归模型:

$$y_4 = -0.33 + 0.022x_1 + 0.010x_2 - 0.001x_3 - 0.006x_4 + 0.005x_5 \quad (13)$$

他们的显著性检验统计如下:

表 5-8 二院模型的评价表

组别	P					F
	x1	x2	x3	x4	x5	
一月	0.822	0.733	0.145	0.802	0.985	0.45

三月	0.957	0.083	0.263	0.87	0.455	1.038
六月	0.970	0.414	0.905	0.430	0.078	0.028
十二月	0.351	0.640	0.982	0.854	0.927	0.248
总数	0.660	0.173	0.747	0.854	0.751	0.437

5.2.3 结果

通过对上述一院、二院回归模型的结果分析可得，整体结果都表明受试者的身体指标不具有显著性，受试者出现不适状况的因素和身体指标有关联，但显著性不大。具体来说，综合对比两院的分析情况来讲，初潮年龄和月经周期对出现不适状况的影响较大，在实际生活中这两者对节育器的使用更为相关，也符合常识；其次宫腔深度也和出现不适状况有关联，宫腔深度的大小和节育器的型号相对应；而年龄和月经经期不显著，这两者和节育器的使用关联性不大，因此对出现不适情况的影响较小。

5.3 问题 3 的模型建立与求解

5.3.1 问题 3 模型的建立

我们需要解决的问题是根据受试者的身体指标、节育器的理化指标与随访时的主诉情况，建立节育器质量模型，并分析 VCu260 与 VCu380 记忆型宫内节育器的质量哪个更优，更适合生产。根据题目可知，已被临床应用的 MCu 功能性宫内节育器为对照组，因此对附件一二的数据进行预处理后，选择一组(MCu 功能性宫内节育器)的数据为训练数据，二组和三组(VCu260 与 VCu380)的数据为预测数据。采用 BP 神经网络进行建模^[6]，其简单拓扑结构如下：

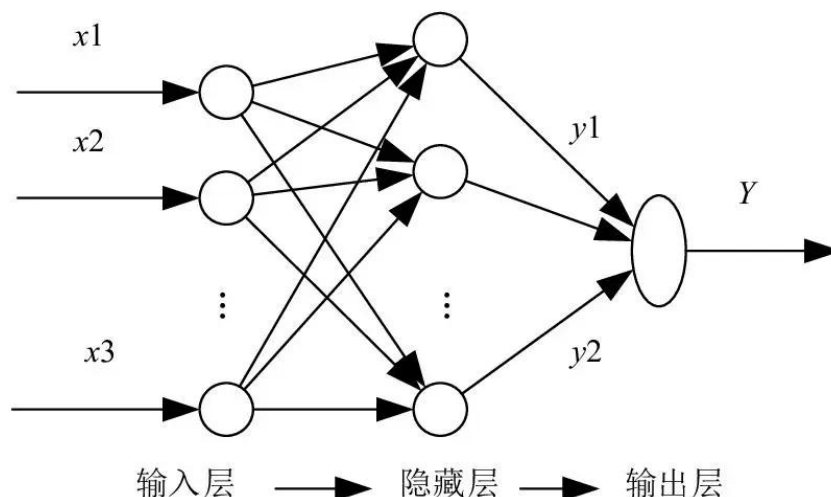


图 5-12 BP 网络的简单拓扑结构

BP(Back-Propagation)神经网络是一种常见的人工神经网络,作为一种非线性的数据建模和处理方法,具有广泛的应用前景,可以在不同领域中发挥重要作用。它可以通过学习样本中的特征,快速地完成对未知数据的分类或识别。通常由输入层、隐藏层和输出层组成^[7-9]。输入层是输入数据的读入端,除了神经元的个数之外隐藏层的其他参数均为待训练的,最后输出层就是输出最终想要的结果,训练步骤如下:

Step1: 初始化网络参数并设置适当的误差阈值、最大训练次数以及学习率和激活函数;

Step2: 计算隐藏层的输出 H_j , $H_j = f(\sum_{i=1}^n \omega_{ij}x_i + a_j)$, 其中 H_j 表示第 j 个神经元的输出, x_i 表示第 i 个输入数据, ω_{ij} 表示第 i 个输入到第 j 个神经元之间的权重, a_j 表示偏置。

Step3: 计算输出层结果 Y_k , $Y_k = \sum_{j=1}^l H_j \omega_{jk} + b_k$, 其中 Y_k 表示第 k 个输出结果, ω_{jk} 表示第 j 个神经元到第 k 个输出之间的权重, b_k 表示偏置。

Step4: 计算输出结果与真实结果之间的误差进行反向传播更新权重,权重更新公式为:
$$\begin{cases} \omega_{ij} = \omega_{ij} + \eta H_j (1 - H_j) x_i \sum_{k=1}^m \omega_{jk} e_k \\ \omega_{jk} = \omega_{jk} + \eta H_j e_k \end{cases}$$
 其中 η 表示提前设置的学习率, e_k 表示误差。

Step5: 判断网络输出的误差是否减小到预先设置的阈值以下,若误差减小到设置阈值以下或者训练次数超过预先设置的最大训练次数,则算法结束,否则返回第二步继续迭代。

5.3.2 问题 3 模型的求解

根据 python 编程对附件一二的数据处理后得出每个人的最终是否适应节育器的指标值,我们通过最终适应节育器的人数在总人数中的占比大小评价节育器的质量好坏。然后利用一组(MCu)的受试者的身体指标、节育器的理化指标数据与随访时的主诉情况作为模型的训练数据集,二组和三组(VCu260 与 VCu380)的受试者的身体指标、节育器的理化指标数据与随访时的主诉情况作为模型的预测数据集,进而对二组和三组(VCu260 与 VCu380)的节育器的质量好坏进行预测。

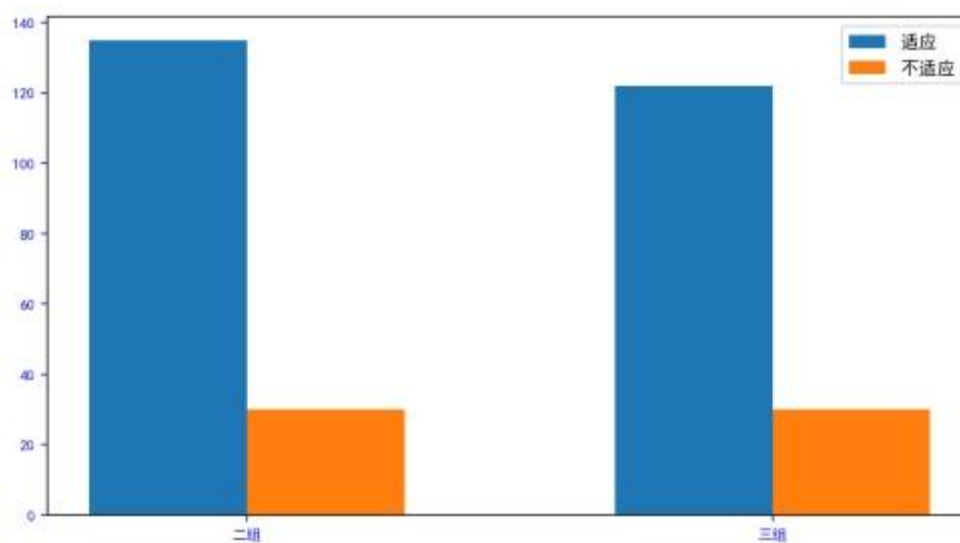


图 5-13 一院的人群适应情况的预测结果

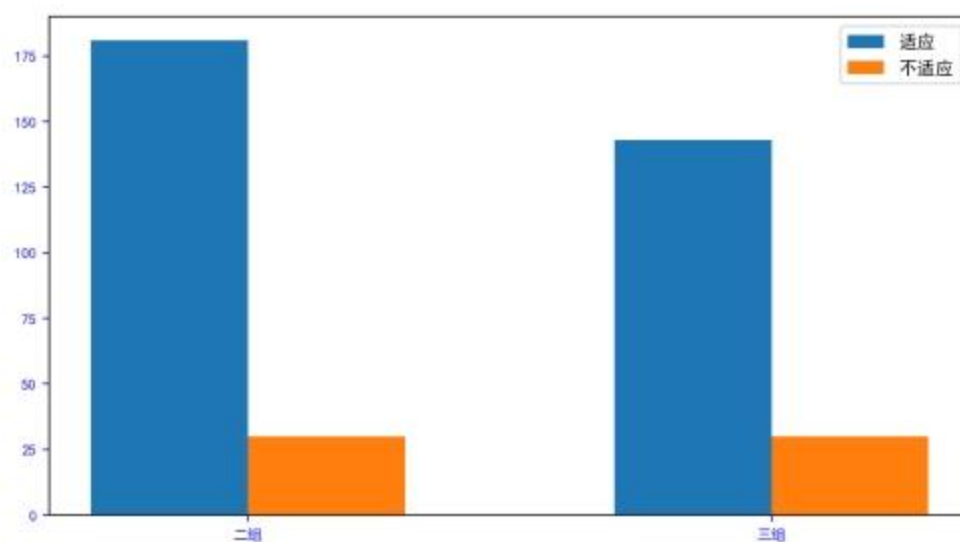


图 5-14 二院的人群适应情况的预测结果

由上图所示，根据受试者的身体指标、节育器的理化指标与随访时的主诉情况进行适应节育器预测分析，可以从预测分布图看出，一院和二院的分布情况大致相同。

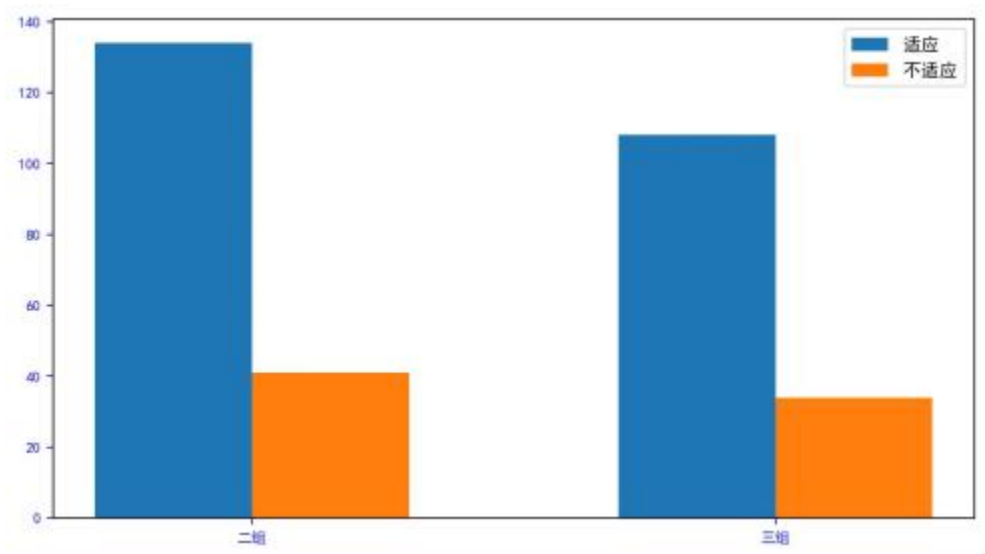


图 5-15 一院实际记录的人群适应情况

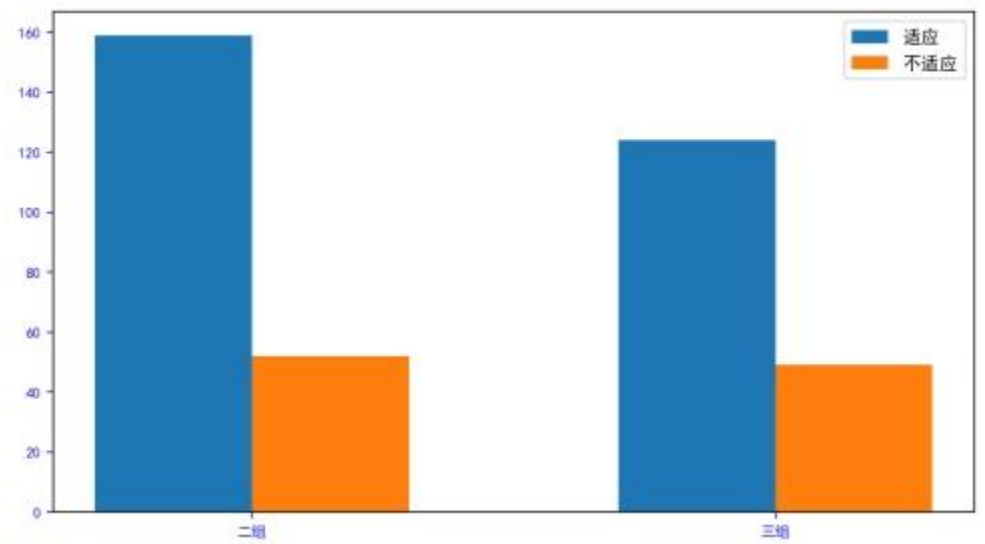


图 5-16 二院实际记录的人群适应情况

由上图所示，从一院和二院的实际记录数据分布与理想情况下（预测的分布情况）的分布情况中能够看出相较于三组(VCu380)，二组(VCu260)的分布更接近理想状态。同时我们统计了适应人数在总体人数中的占比情况，如下表。

表 5-9 受试者的适应人数占比情况				
	二组（预测）	三组（预测）	二组（记录）	三组（记录）
一院	0.818	0.803	0.766	0.761
二院	0.858	0.827	0.754	0.717

由上图和表所示，根据 BP 网络建模能够得出节育器质量模型，并从预测结果中可以看出，二组(VCu260)的节育器在人群中的适应情况要优于三组(VCu380)。

5.3.3 结果

将随访时的主诉情况数据进行处理筛选，得到每个人最终是否适应节育器。通过对一组(MCu)数据建立的质量评价模型可以对二组和三组(VCu260 与 VCu380)的数据进行预测，从而得到一个理想情况下二组与三组(VCu260 与 VCu380)使用节育环后的适应情况。通过质量评价模型得到的预测结果将与医院的实际记录的结果做对比。从而得出VCu260 节育器的质量优于 VCu380 节育器的结论。

5.4 问题 4 的模型建立与求解

问题 4 要求根据问题 3 建立的节育器质量模型，探究影响宫内节育器质量的决定因素。在理想状态下，可以认为问题 3 建立的节育器质量预测模型是最优的。因此，使用问题 3 建立的模型可以对决定节育器的质量好坏的决定性因素做出判断。

模型中隐藏层的权重系数在一定程度上可以衡量输入特征对输出结果的影响程度。因此，对模型输入特征的系数权重进行可视化，找出权重较大的特征，从而就能够分析出影响节育器质量的决定性因素。如下图是一院模型的 8 个输入特征的权重系数大于 0.9 的部分指标可视化结果，可以得出一院中节育器的质量好坏的决定性因素为：应用节育器情况、宫腔深度、使用节育器型号和宫颈扩张情况。

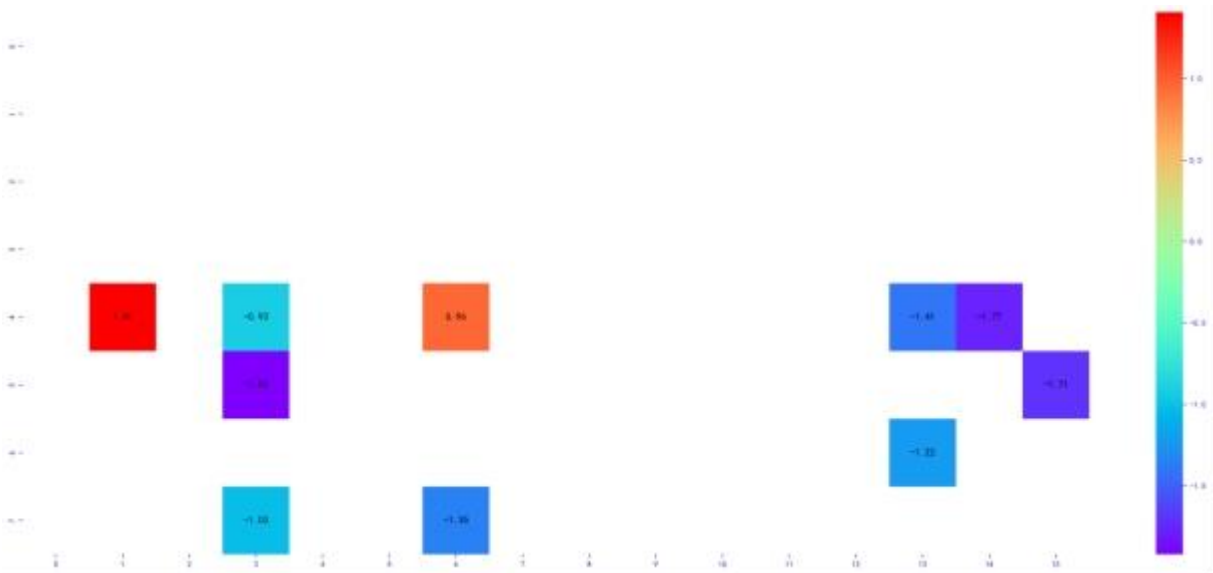


图 5-17 一院模型权重系数可视化结果

二院模型的 8 个输入特征的权重系数大于 0.9 的部分指标可视化结果如下，可以得

出二院中节育器的质量好坏的决定性因素为：应用节育器情况、宫腔深度、使用节育器型号和宫颈扩张情况。

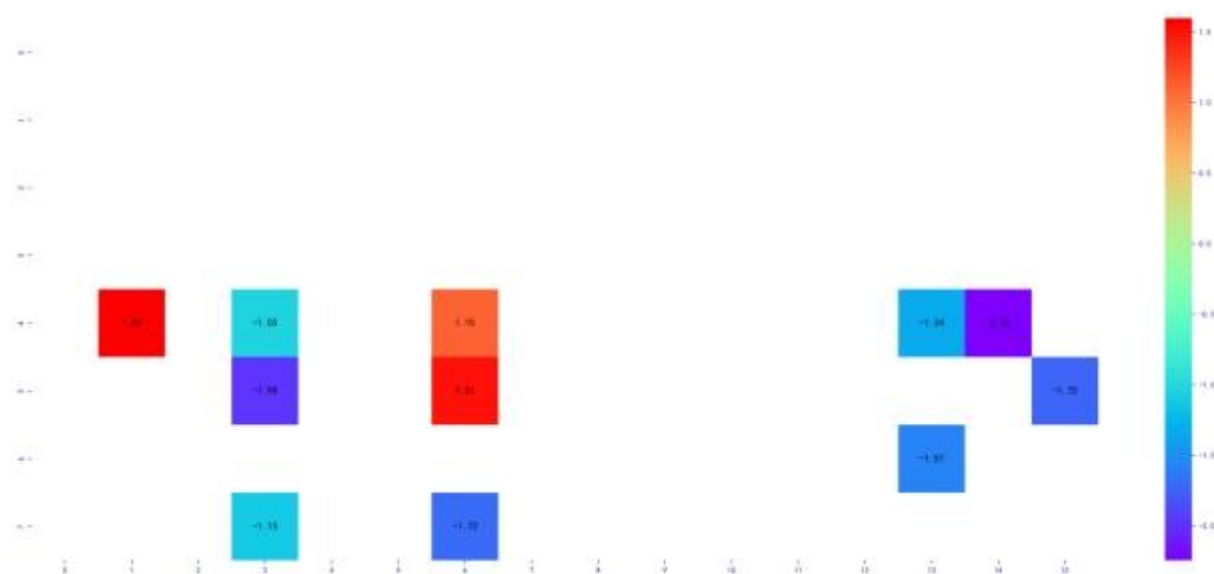


图 5-18 二院模型权重系数可视化结果

由上分析可知，两个医院的影响宫内节育器质量的决定因素完全一致。因此很容易得知影响宫内节育器质量的决定因素是：应用节育器情况、宫腔深度、使用节育器型号和宫颈扩张情况。

六、模型的评价及优化

6.1 误差分析

6.1.1 针对于问题 1 的误差分析

问题 1 中两个医院临床受试者主诉情况的分布差异性大，但是统计分析的身体指标与节育器的理化指标分布性差异显著性不大，可能是没有对失访数据进行处理，清洗后的数据可以看出导致这种差异的因素。

6.1.2 针对于问题 2 的误差分析

问题 2 中由于附件一和附件二中的数据量不匹配，建立的回归统计模型代表性不强，

通过填补附件二中感受良好的受访信息，模型的回归效果和显著性增强了。

6.1.3 针对于问题 3 的误差分析

问题 3 需要对两种节育器的好坏做出判断，需要根据使用者在 1、3、6、12 月的不适状况做出判断，通过构建判断指标，提高了模型的数据质量，也提高了模型的预测效果。

6.1.4 针对于问题 4 的误差分析

问题 4 中需要根据问题 3 的模型来探究影响宫内节育器质量的决定因素，没有构建质量判断指标前分析得到的决定因素不显著，构建后的权重热力图效果更明显。

6.2 模型的优点

通过构建节育器质量好坏标准，使用网络建立身体指标、节育器的理化指标与附件二中“有不适人数”之间的模型，没有直接预测节育器质量的好坏，而是预测 1、3、6、12 月的不适状况，再根据构建的评估标准计算质量好坏，模型的适用性更强，容错率更低。

6.3 模型的缺点

统计分析的数据之间存在微弱的不一致性，使得模型应用的时候存在偏差。

6.4 模型的推广

本文所建立的模型能够根据节育器使用者的情况做出 1、3、6、12 出现不适情况的判断，通过分析用户的身体指标、节育器的理化指标，对未来可能的一些不适状况作何判断和预测。针对预测存在的误差，后期可对主控因素进行挖掘，建立与 8 种不适状况之间的模型，更加详细的预测不同时间点出现的不同状况，做出更精确的判断。

参考文献

- [1] 刘严. 多元线性回归的数学模型[J]. 沈阳工程学院学报: 自然科学版, 2005, 1(2): 128-129.
- [2] Forrest D R, Hetland R D, DiMarco S F. Multivariable statistical regression models of the areal extent of hypoxia over the Texas–Louisiana continental shelf[J]. Environmental Research Letters, 2011, 6(4): 045002.
- [3] 戚德虎, 康继昌. BP 神经网络的设计[J]. 计算机工程与设计, 1998, 19(2): 48-50.
- [4] Li J, Cheng J, Shi J, et al. Brief introduction of back propagation (BP) neural network algorithm and its improvement[C]//Advances in Computer Science and Information Engineering: Volume 2. Springer Berlin Heidelberg, 2012: 553-558.
- [5] Eberly L E. Multiple linear regression[J]. Topics in Biostatistics, 2007: 165-187.
- [6] 邹乐强. 最小二乘法原理及其简单应用[J]. 科技信息, 2010 (23): 282-283.
- [7] Grant S W, Hickey G L, Head S J. Statistical primer: multivariable regression considerations and pitfalls[J]. European Journal of Cardio-Thoracic Surgery, 2019, 55(2): 179-185.
- [8] Björck Å. Least squares methods[J]. Handbook of numerical analysis, 1990, 1: 465-652.
- [9] Mansard E P D, Funke E R. The measurement of incident and reflected spectra using a least squares method[M]//Coastal Engineering 1980. 1980: 154-172.
- [10] Mansard E P D, Funke E R. The measurement of incident and reflected spectra using a least squares method[M]//Coastal Engineering 1980. 1980: 154-172.

附 录

问题 1:

```

from tensorflow.keras.layers import *
from tensorflow.keras import *
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import tensorflow as tf
import matplotlib.pyplot as plt
import matplotlib as mpl

mpl.rcParams['font.family'] = 'SimHei'
mpl.rcParams['axes.unicode_minus'] = False
mpl.rcParams['xtick.labelsize'] = 8
mpl.rcParams['ytick.labelsize'] = 8
mpl.rcParams['xtick.color'] = 'b'
mpl.rcParams['ytick.color'] = 'b'

name1 = ["一月", "三月", "六月", "十二月"]
name2 = ["脱落", "因症取出", "怀孕", "非月经期出血", "疼痛", "经量多", "分泌物增多", "经期/周期异常"]
name3 = ["MCu", "VCu260", "VCu380"]

#####
def prouduct(m1, m2):
    data2 = pd.read_excel('C:/Users/abc/Desktop/建模/附件 2：两个医院随访的节育器使用后主诉情况.xlsx', sheet_name=m1)
    data2 = data2[2:]
    data2 = data2[data2["组别"] == m2]
    data2 = data2[data2.columns[2:-4]]
    data2 = data2.fillna(value=0)

    n = data2.values.shape[0]
    D = np.zeros((4, 8))
    for i in range(4):
        for k in range(n):
            if data2.values[k][i] > 0:
                continue
            for j in range(1, 9):
                D[i][j-1] += data2.values[k][j*4+i]
    return D

for c in range(1, 4):
    D11 = prouduct(0, c)
    D21 = prouduct(1, c)

    x = np.arange(4)+0.7

    plt.subplots(figsize=(16, 6), dpi=200)

```



```

for i in range(8):
    plt.subplot(2, 4, i+1)

    plt.bar(x, D11.T[i], 0.3, label='一院')
    plt.bar(x+0.3, D21.T[i], 0.3, label='二院')

    # plt.bar(name1, D11.T[i])
    # plt.bar(name1 D21.T[i])
    # plt.barh(name1, D.T[5])
    plt.grid()
    plt.title(name2[i])
    plt.xticks(np.arange(4)+0.85, name1)
plt.legend()
plt.show()

#####
def prouduct1(m1, m2):
    data2 = pd.read_excel('C:/Users/abc/Desktop/建模/附件 2：两个医院随访的节育器使用后主诉情
况.xlsx', sheet_name=m1)
    data2 = data2[2:]
    data2 = data2[data2["组别"] == m2]
    data2 = data2[data2.columns[2+4:-4]]
    data2 = data2.fillna(value=0)

    n = data2.values.shape[0]
    D = np.zeros((n, 8))
    for i in range(n):
        for j in range(8):
            if sum(data2.values[i][j]*4:(j+1)*4) > 0:
                D[i][j] = 1

    return np.sum(D, axis=0), n
num1 = []
num2 = []
for c in range(1, 4):
    D11, a = prouduct1(0, c)
    D21, b = prouduct1(1, c)
    num1.append(a)
    num2.append(b)

x = np.arange(8)+0.7

plt.subplots(figsize=(9, 5), dpi=200)

plt.bar(x, D11, 0.3, label='一院')
plt.bar(x+0.3, D21, 0.3, label='二院')

plt.grid()
plt.title(name3[c-1])
plt.xticks(np.arange(8)+0.85, name2)
plt.legend()
plt.show()

#####
x = np.arange(3)+0.7

```

```
plt.subplots(figsize=(8, 5), dpi=200)
plt.bar(x, num1, 0.3, label='一院')
plt.bar(x+0.3, num2, 0.3, label='二院')

plt.grid()
plt.xticks(np.arange(3)+0.85, name3)
plt.legend()
plt.show()
```

问题 2:

```
from tensorflow.keras.layers import *
from tensorflow.keras import *
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import tensorflow as tf
import matplotlib.pyplot as plt
import matplotlib as mpl

mpl.rcParams['font.family'] = 'SimHei'
mpl.rcParams['axes.unicode_minus'] = False
mpl.rcParams['xtick.labelsize'] = 8
mpl.rcParams['ytick.labelsize'] = 8
mpl.rcParams['xtick.color'] = 'b'
mpl.rcParams['ytick.color'] = 'b'

name1 = ["一月", "三月", "六月", "十二月"]
name2 = ["脱落", "因症取出", "怀孕", "非月经期出血", "疼痛", "经量多", "分泌物增多", "经期/周期异常"]
name3 = ["MCu", "VCu260", "VCu380"]

data2 = pd.read_excel('C:/Users/abc/Desktop/建模/附件 2: 两个医院随访的节育器使用后主诉情况.xlsx',
sheet_name=1)
data2 = data2[2:]
data2 = data2.fillna(value=0)
data3 = data2[data2.columns[6:-4]]
data2["一月"] = np.sum(data3.values[:, 0::4], axis=1)
data2["三月"] = np.sum(data3.values[:, 1::4], axis=1)
data2["六月"] = np.sum(data3.values[:, 2::4], axis=1)
data2["十二月"] = np.sum(data3.values[:, 3::4], axis=1)
data2["总数"] = np.sum(data3.values[:, ::4], axis=1)

data2[data2.columns[-9]] = data2["一月"].map(lambda x: 1 if x > 0 else 0)
data2[data2.columns[-8]] = data2["三月"].map(lambda x: 1 if x > 0 else 0)
data2[data2.columns[-7]] = data2["六月"].map(lambda x: 1 if x > 0 else 0)
data2[data2.columns[-6]] = data2["十二月"].map(lambda x: 1 if x > 0 else 0)

data2 = data2.set_index('序号')
print(data2)
```

```

data1 = pd.read_excel('C:/Users/abc/Desktop/建模/附件 1: 两个院临床受试者和节育器的基本数据.xlsx',
sheet_name=1)
data1 = data1[1:]
data1 = data1.fillna(value=0)

data1["应用情况"] = data1[data1.columns[6]].map(lambda x: 1 if x > 0 else 0) +
data1[data1.columns[7]].map(lambda x: 2 if x > 0 else 0) + data1[data1.columns[8]].map(lambda x: 3 if
x > 0 else 0)
data1[data1.columns[10]] = data1[data1.columns[10]].map(lambda x: 1 if x > 0 else 0)
data1[data1.columns[11]] = data1[data1.columns[11]].map(lambda x: 2 if x > 0 else 0)
data1[data1.columns[12]] = data1[data1.columns[12]].map(lambda x: 3 if x > 0 else 0)

data1["型号情况"] = data1[data1.columns[10:13]].max(axis=1)

data1[data1.columns[10]] = data1["型号情况"]
del data1["型号情况"]
data1[data1.columns[6]] = data1["应用情况"]
del data1["应用情况"]

del data1[data1.columns[12]]
del data1[data1.columns[11]]
del data1[data1.columns[8]]
del data1[data1.columns[7]]

data1 = data1.set_index('序号')
print(data1)

Data = data1.merge(data2, how="inner", on="序号")
# print(Data[Data.columns[-9:-5]])
# print(Data[Data.columns[-5:]])

# Data = Data[Data["组别_x"] == 1]
Data = Data.fillna(value=0)

Data["组别_y"] = Data[Data.columns[-9]].map(lambda x: str(int(x))) +
Data[Data.columns[-8]].map(lambda x: str(int(x))) + Data[Data.columns[-7]].map(lambda x: str(int(x))) +
Data[Data.columns[-6]].map(lambda x: str(int(x)))

cla = ["0000", "1000", "0100", "0010", "1100"]
Data["组别_y"] = Data["组别_y"].map(lambda x: 0 if x in cla else 1)
Data = Data[~Data.index.duplicated()]

Data[Data.columns[2:8]].describe()

X = Data[Data.columns[1:7]]
del X[X.columns[-2]]
Y = Data[Data.columns[-5:]]

import statsmodels.api as sm

X = X.values

```

Y = Y.values

for i in range(5):

```

    X1 = sm.add_constant(X)
    Y1 = Y[:, i]
    model = sm.OLS(Y1, X1)
    results = model.fit()
    print(' 回归参数的估计值分别为 ----- :
\n',results.params)
    print('检验的结果为: \n',results.summary())

```

```

import seaborn as sns
#f, (ax1,ax2) = plt.subplots(figsize=(9, 9), nrow=2)
#sns.heatmap(A, annot=True, fmt='.2f', ax=ax1)
#sns.heatmap(A, mask = abs(A) < 0.5, fmt='.2f', annot=True, annot_kws={"weight": "bold"}, ax=ax2)

```

```

def prouduct2(m1):
    data1 = Data
    data1 = data1[1:]
    del data1[data1.columns[6]]
    del data1[data1.columns[6]]
    del data1[data1.columns[6]]
    data11 = data1[data1["组别"] == 1]
    data12 = data1[data1["组别"] == 2]
    data13 = data1[data1["组别"] == 3]

    disMat = data1[data1.columns[2:7]].corr()

    f, ax1 = plt.subplots(figsize=(10, 8), dpi=200)
    sns.heatmap(disMat, fmt='.3f', annot=True, annot_kws={"weight": "bold"}, cmap='rainbow')

    print(disMat)
    name33 = ["整体"] + name3
    plt.subplots(figsize=(25, 6), dpi=200)
    for i in range(5):
        ax = plt.subplot(1, 5, i+1)
        ax.boxplot([data1[data1.columns[2+i]].values, data11[data11.columns[2+i]].values,
data12[data12.columns[2+i]].values, data13[data13.columns[2+i]].values],showmeans=True)
        ax.set_xticklabels(name33, fontsize=18)
        plt.title(name4[i], fontsize=22)

    plt.show()

prouduct2(0)
prouduct2(1)

```

问题三:

```

from tensorflow.keras.layers import *
from tensorflow.keras import *
import numpy as np

```

```

import matplotlib.pyplot as plt
import pandas as pd
import tensorflow as tf
import matplotlib.pyplot as plt
import matplotlib as mpl

in_put = tf.keras.Input(shape=((8,)))
out = tf.keras.layers.Dense(16, activation='relu')(in_put)
out = tf.keras.layers.Dense(32, activation='relu')(out)
out = tf.keras.layers.Dense(10, activation='relu')(out)
out_put = tf.keras.layers.Dense(4, activation='sigmoid')(out)
model = tf.keras.Model(inputs = in_put, outputs = out_put)
model.summary()
model.compile('adam', loss='mse', metrics='mae')

# 用对照组数据进行建模
def data_split(m1,m2):
    # 上面的 Data
    data = pd.read_excel(r'C:\Users\78096\Desktop\2023 年第八届数维杯数学建模挑战赛题目\11.xlsx',
sheet_name=m1)
    # data = data[2:].fillna(0)
    data = data[data["组别_x"] == m2]
    X = data[data.columns[1:9]].values
    Y = data[data.columns[-9:-5]].values
    return X,Y

from sklearn import model_selection
# 一院建模
X11,Y11 = data_split(0,1)
x_train, x_test, y_train, y_test = model_selection.train_test_split(X11, Y11, test_size = 0.1, random_state =
1234)
history = model.fit(x_train, y_train, batch_size=1, epochs=500, verbose=1, validation_data=(x_test,
y_test))
Y11_pre = np.int64(model.predict(X11)>0.5)

# 预测一院两个实验组的适应程度，实际情况和理论上更接近的实验组质量更佳
X12,Y12_real = data_split(0,2)
Y12_pre = np.int64(model.predict(X12)>0.5)
X13,Y13_real = data_split(0,3)
Y13_pre = np.int64(model.predict(X13))

# 二院建模
X21,Y21 = data_split(1,1)
x_train, x_test, y_train, y_test = model_selection.train_test_split(X21, Y21, test_size = 0.1, random_state =
1234)
history = model.fit(x_train, y_train, batch_size=1, epochs=500, verbose=1, validation_data=(x_test,
y_test))
Y21_pre = np.int64(model.predict(X21)>0.5)

# 预测二院两个实验组的适应程度，实际情况和理论上更接近的实验组质量更佳
X22,Y22_real = data_split(1,2)
Y22_pre = np.int64(model.predict(X22))
X23,Y23_real = data_split(1,3)
Y23_pre = np.int64(model.predict(X23))

```

```

# 预测的最终结果
Data = pd.read_excel(r'C:\Users\78096\Desktop\2023 年第八届数维杯数学建模挑战赛题目\结果.xlsx',
sheet_name=1)
D = Data[Data.columns[-5]].map(lambda x: str(int(x))) + Data[Data.columns[-4]].map(lambda x: str(int(x)))
+ Data[Data.columns[-3]].map(lambda x: str(int(x))) + Data[Data.columns[-2]].map(lambda x: str(int(x)))

cla = ["0000", "1000", "0100", "0010", "1100"]
D = D.map(lambda x: 0 if x in cla else 1)
D = D.values

# 一院第一组计算正确率
m1 = 1
m2 = 3
Data = pd.read_excel(r'C:\Users\78096\Desktop\2023 年第八届数维杯数学建模挑战赛题目\结果.xlsx',
sheet_name=m1)

data1 = Data[Data["组别_x"] == m2]["组别_y"].values
data2 = Data[Data["组别_x"] == m2]["预测结果"].values
right = 0
for i in range(len(data1)):
    if data1[i] == data2[i]:
        right += 1
print(right/len(data1))

# 最终判定是否适应
def panduan(m1):
    Data = pd.read_excel(r'C:\Users\78096\Desktop\2023 年第八届数维杯数学建模挑战赛题目\结
果.xlsx', sheet_name=m1)
    D = Data[Data.columns[0]].map(lambda x: str(int(x))) + Data[Data.columns[1]].map(lambda x:
str(int(x))) + Data[Data.columns[2]].map(lambda x: str(int(x))) + Data[Data.columns[3]].map(lambda x:
str(int(x)))

    cla = ["0000", "1000", "0100", "0010", "1100"]
    D = D.map(lambda x: 0 if x in cla else 1)
    Y_pre = D.values
    return Y_pre

Y12_pre = panduan(0)
Y13_pre = panduan(1)
Y22_pre = panduan(2)
Y23_pre = panduan(3)
Y1_true = np.concatenate((YY12,YY13))
Y1_pre = np.concatenate((Y12_pre,Y13_pre))
Y2_true = np.concatenate((YY22,YY23))
Y2_pre = np.concatenate((Y22_pre,Y23_pre))

# 二组和三组 0 和 1 的个数
def plt_bar(D1, D2):
    x = np.arange(2)+0.7
    plt.subplots(figsize=(9, 5), dpi=200)
    plt.bar(x, D1, 0.3, label='适应')
    plt.bar(x+0.3, D2, 0.3, label='不适应')

```

```

plt.xticks(np.arange(2)+0.85, ["一组", "二组"])
plt.legend()
plt.show()
# 一院的预测图
D11 = [len(Y12_pre)-sum(Y12_pre), len(Y13_pre)-sum(Y13_pre)]
D12 = [sum(Y12_pre), sum(Y13_pre)]
plt_bar(D11, D12)

# 二院的预测图
D21 = [len(Y22_pre)-sum(Y22_pre), len(Y23_pre)-sum(Y23_pre)]
D22 = [sum(Y22_pre), sum(Y23_pre)]
plt_bar(D21, D22)

# 一院的真实图
D11 = [len(Y12)-sum(Y12), len(Y13)-sum(Y13)]
D12 = [sum(Y12), sum(Y13)]
plt_bar(D11, D12)

# 二院的真实图
D21 = [len(Y22)-sum(Y22), len(Y23)-sum(Y23)]
D22 = [sum(Y22), sum(Y23)]
plt_bar(D21, D22)

```

问题四：

```

import pandas as pd
import tensorflow as tf
import numpy as np
import matplotlib.pyplot as plt
from scipy.signal import savgol_filter
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
import tensorflow as tf
from tensorflow.keras.layers import Dense
plt.rcParams['axes.unicode_minus'] = False # 用来正常显示负号
plt.rcParams['font.sans-serif'] = ['SimHei'] # 用来正常显示中文标签

path = r'C:\Users\78096\Desktop\采收率\新建文件夹' # Recovery rate 采收率\2

# 加载模型

from tensorflow.keras.models import load_model
model = load_model(r'C:\Users\78096\Desktop\2023 年第八届数维杯数学建模挑战赛题目\model2.h5')
model.summary()

for i, layer in enumerate(model.layers):
    if i == 1:#8
        M = layer.variables[0].numpy()

import seaborn as sns
plt.subplots(figsize=(24, 10), dpi=500)
sns.heatmap(M, mask = abs(M)<=0.9, fmt='.2f', annot=True, annot_kws={"weight": "bold", "color": "k"},
linewidths=0, cmap='rainbow')

```