

# 决策树算法应用于 MIMIC-III 数据库的 ICU 患者急性肾损伤预测研究

高文鹏<sup>1</sup> 吕海金<sup>2</sup> 周琅<sup>1</sup> 郭圣文<sup>3</sup>

**摘要** 目的 急性肾损伤(acute kidney injury, AKI)是重症监护病房(intensive care unit, ICU)最常见的并发症和致死因素之一。准确预测具 AKI 风险的患者,明确与 AKI 发生相关的关键因素,可为临床决策与风险患者干预提供有效指导。方法 采用公开的重症监护室数据库 MIMIC-III,提取 30 020 例患者记录(包括 AKI 患者 17 222 名,Non-AKI 患者 12 798 名),收集其住 ICU 期间基本信息、生理生化指标、药物使用、合并症等临床信息。将患者按 4:1 比例随机划分训练集和独立测试集,应用逻辑回归、随机森林与 LightGBM 3 种机器学习方法,分别建立 24 h、48 h 与 72 h 3 个时间点的 AKI 预测模型,采用十折交叉验证法,对各种模型进行训练与测试,预测患者是否发生 AKI,并获取重要特征。此外,利用 24 h 预测模型,在一周时间窗口内对 ICU 患者进行每隔 24 h 预测。结果 3 种学习模型中,LightGBM 性能最优,其 24 h、48 h 和 72 h 模型预测 AKI 的受试者工作特征曲线(receiver operator characteristic curve, ROC 曲线)下面积(area under curve, AUC)值分别为 0.90、0.88、0.87, F1 值分别为 0.91、0.88、0.86,在每隔 24 h 预测时,提前 1 d、2 d 和 3 d 预测 AKI 的成功率分别为 89%、83%、80%。已住院时长、体质量、白蛋白、收缩压、碳酸氢盐、葡萄糖、白细胞计数、体温、舒张压、血尿素氮等是预测 ICU 患者 AKI 的重要特征,仅使用 24 个重要特征,模型仍能取得良好的预测性能。结论 基于 ICU 患者的基本信息、生理生化指标、药物使用及合并症等临床信息,应用机器学习模型,可对其是否发生 AKI 进行多时间点的有效预测,并明确其关键风险因素。

**关键词** 急性肾损伤;重症监护室;机器学习;风险预测;重要特征

**DOI:** 10.3969/j.issn.1002-3208.2021.06.010.

**中图分类号** R318 **文献标志码** A **文章编号** 1002-3208(2021)06-0609-09

本文著录格式 高文鹏,吕海金,周琅,等.决策树算法应用于 MIMIC-III 数据库的 ICU 患者急性肾损伤预测研究[J].北京生物医学工程,2021,40(6):609-617. GAO Wenpeng, LYU Haijin, ZHOU Lang, et al. Decision tree algorithm applied to MIMIC-III database for the prediction of acute kidney injury in ICU patients[J]. Beijing Biomedical Engineering, 2021, 40(6): 609-617.

## Decision tree algorithm applied to MIMIC-III database for the prediction of acute kidney injury in ICU patients

GAO Wenpeng<sup>1</sup>, LYU Haijin<sup>2</sup>, ZHOU Lang<sup>1</sup>, GUO Shengwen<sup>3</sup>

<sup>1</sup> Department of Biomedical Engineering, School of Material Science and Engineering, South China University of Technology, Guangzhou 510006;

<sup>2</sup> SICU, The Third Affiliated Hospital, Sun Yat-sen University, Guangzhou 510630;

<sup>3</sup> School of Automation Science and Engineering, South China University of Technology, Guangzhou 510640

Corresponding author: GUO Shengwen (E-mail: shwguo@scut.edu.cn)

**【Abstract】 Objective** Acute kidney injury (AKI) is one of the most common complications and fatal factors in intensive care unit (ICU). Accurate prediction of AKI risk and identification of key factors related to AKI can provide effective guidance for clinical decision-making and intervention for patients with AKI risk.

**Methods** A total of 30 020 patients in ICU ( including 17 222 AKI patients and 12 798 Non-AKI patients) were selected from the public database MIMIC-III in this study and basic information ,physiological and biochemical indicators ,drug use and comorbidity during their stay in ICU were collected. All patients were randomly divided into training sets and independent testing sets according to the ratio of 4 : 1 ,and logistic regression ,random forest and LightGBM were applied to construct models for AKI predication in three time points including 24 h , 48 h and 72 h ,respectively. The 10-fold cross validation was used to train and validate various models to predict the occurrence of AKI and obtain important features. Furthermore ,24 h prediction models were used to predict AKI every 24 h during the 7-day window. **Results** LightGBM achieved the best performance with AUC values of 0.90 ,0.88 ,0.87 for 24 h ,48 h and 72 h prediction ,respectively and F1 values were 0.91 ,0.88 and 0.86. In prediction of every 24 h ,the success rates of identifying AKI patients were 89% ,83% and 80% in one day ,two days and three days in advance ,respectively. It was found that the length of stay in ICU ,body weight ,albumin , systolic blood pressure ,bicarbonate ,glucose ,white blood cell count ,body temperature ,diastolic blood pressure and blood urea nitrogen played vital roles in predicting AKI for ICU patients. Using only 24 important features , the models could still achieve prominent prediction performance. **Conclusions** Based on basic information , physiological and biochemical indicators ,drug use and comorbidity ,machine learning methods can be adopted to effectively predict AKI risk for ICU patients at several time points ,and determine the dominant factors relative to AKI.

**【Keywords】** acute kidney injury; intensive care unit; machine learning; risk prediction; important feature

## 0 引言

急性肾损伤(acute kidney injury ,AKI)在重症监护室(intensive care unit ,ICU)患者中很常见,具有较高的发病率与死亡率<sup>[1-2]</sup>。全球肾脏病预后组织 KDIGO(Kidney Disease: Improving Global Outcomes)发表的 AKI 临床实践指南<sup>[3]</sup>中对 AKI 的具体判别标准为:48 h 内血清肌酐升高 $\geq 26.5 \mu\text{mol/L}$ 或在 7 d 内血清肌酐增加到基线值的 1.5 倍以上,或尿量小于  $0.5 \text{ mL}/(\text{kg} \cdot \text{h})$ 且持续时间不少于 6 h。研究表明,AKI 导致 ICU 患者更高的治疗费用、不良的临床反应和慢性肾病的发展<sup>[3]</sup>,并且是 ICU 患者高死亡率的独立影响因素<sup>[4-5]</sup>。

由于血清肌酐是 AKI 的非特异性标志物,对 AKI 的诊断具有一定滞后性<sup>[6]</sup>,而临床中尿量不易监测且操作误差较大,故寻找影响 AKI 的重要临床因素并进行早期预测,是对 ICU 内 AKI 风险患者及时干预、指导治疗的关键。迄今为止,有关 AKI 的

早期预测通常有以下两种方法:一种是寻找具有特异性的生物标志物,另一种是基于统计学或机器学习方法建立风险预测模型。由于生物标志物的方法费用高、可纳入样本量少且受个体差异影响较大,其临床应用受限。

随着开源重症数据库的建立和医院电子病历(electronic health records ,EHR)的普及,ICU 患者临床数据的可用性不断提高<sup>[7]</sup>,从而为 AKI 预测研究提供充分的数据支持,相关研究逐渐增多。如 Haines 等<sup>[8]</sup>收集伦敦皇家医院 ICU 内 830 名患者的人口学信息和入院后 24 h 内的血液学指标,使用逻辑回归进行 AKI 预测,结果显示预测 AKI 1~3 期的受试者工作特征曲线(receiver operator characteristic curve ,ROC 曲线)下面积(area under curve ,AUC)值为 0.70,预测 AKI 2~3 期 AUC 值为 0.91。Malhotra 等<sup>[9]</sup>收集两家独立医院 ICU 共 2 017 名患者的人口学信息、合并症、生命体征、血液学指标及干预措施,使用多变量回归分析进行 AKI 预测,独立测试集的 AUC 值为 0.81。也有学者利用开源重症数据库进行 AKI 预测,如李千惠<sup>[10]</sup>从 Medical Information Mart for Intensive Care(MIMIC)重症数据库中提取 1 690 例患者(AKI 患者 840 例)的生命体征和血液学指标,使用逻辑回归、Adaboost 及多层感知机 3 种模

作者单位:1 华南理工大学材料科学与工程学院生物医学工程系(广州 510006)

2 中山大学附属第三医院外科 ICU(广州 510630)

3 华南理工大学自动化科学与工程学院(广州 510640)

通信作者:郭圣文。E-mail: shwguo@scut.edu.cn

型进行AKI早期预测,结果表明多层感知机性能最好, $F_{1.5}$ 分数为0.944。

张渊等<sup>[11]</sup>从MIMIC数据库中提取1166名患者(AKI患者884例)的人口学信息、生命体征和血液学指标,使用逻辑回归、随机森林及LightGBM 3种模型进行AKI发生前24 h预测,发现LightGBM模型最优,AUC值为0.92。Zimmerman等<sup>[12]</sup>提取MIMIC数据库23950例患者的人口学信息及入院后24 h内、生命体征、血液学指标、干预措施,使用逻辑回归、随机森林与多层感知机进行AKI早期预测,平均AUC值为0.783。

目前对ICU患者AKI预测的研究主要存在以下不足:(1)纳入样本量不足,特别是AKI患者的样本量普遍偏少,使得模型可靠性不足;(2)临床或数据库信息利用不充分,可能遗漏重要影响因素;(3)预测时间不及时,缺乏连续预警功能,导致临床医生没有足够的时间进行干预。

针对以上不足,本研究基于30020名ICU患者的人口学信息、入院信息、用药情况、生命体征、血液学指标、危重症评分、合并症、干预措施等8类临床信息,按4:1随机划分训练集与独立测试集,分别应用逻辑回归、随机森林和LightGBM 3种机器学习算法,建立24 h、48 h和72 h 3个时间点的AKI预测模型,对不同模型的性能进行评估、比较与分析,并使用最优模型进行连续24 h预测,明确与AKI事件发生相关的重要因素。

## 1 研究数据

### 1.1 数据来源

本研究所用数据来自MIMIC-III数据库<sup>[13]</sup>,MIMIC-III是美国麻省理工提供的一个公开免费的多参数重症监护数据库,包含了从2001年6月1日至2012年10月31日在波士顿的贝斯以色列女执事医疗中心ICU收治的46520名患者的住院记录,具有样本量大、临床信息丰富的特点。

### 1.2 数据筛选

入组标准:年龄>18岁;ICU住院时长>24 h。排除具有以下合并症的患者:肾结石;输尿管结石;肾癌;肾盂癌;尿路梗阻性疾病。首测肌酐值属于正常范围( $31.8 \sim 116.0 \mu\text{mol/L}$ )的患者,以首测肌酐值作为基线肌酐值;首测肌酐值不属于正常范围的患者,在排除慢性肾病的前提下,取 $116.0 \mu\text{mol/L}$

作为基线肌酐值;多次入住ICU的患者,若连续两次住ICU时间间隔超过48 h,则分别按照不同样本纳入。

根据以上标准,最终纳入30020名患者,其中AKI患者共17222名,占比57.4%。

### 1.3 变量纳入

(1) 人口学信息:年龄、性别、体质量、身高。

(2) 入院信息:入院方式、ICU类型。

(3) 药物使用:抗生素、利尿剂、他克莫司、利福平、两性霉素、顺铂。

(4) 生命体征:平均动脉压、心率、呼吸频率、收缩压、舒张压、体温。

(5) 危重症患者健康评分:Elixhauser comorbidity score、SAPS II、SOFA score、APSIH。

(6) 合并症:高血压、糖尿病、心肌梗死、心力衰竭、脓毒症、癌症。

(7) 血液学指标:肌酐、血红蛋白、白蛋白、pH、碳酸氢盐、碱剩余、乳酸、钾、氯、钠、白细胞计数、葡萄糖、血尿素氮、胆红素。

(8) 干预措施:机械通气、肾脏替代治疗。

其中,入院方式包括ELECTIVE(有计划的入院)、URGENT(急诊,不危及生命)、EMERGENCY(急诊,危及生命)3种;ICU类型包括SICU、MICU、CCU、TSICU、CSRU 5种;合并症与干预措施用“0/1”表示,0表示无,1表示有;危重症患者评分中,Elixhauser comorbidity score是针对患者的合并症评分,SAPS II是简化急性生理学评分,SOFA score是器官衰竭评分,APS III为急性生理学及慢性健康状况评分。

### 1.4 数据收集时间窗

本研究对AKI进行提前24 h、48 h、72 h预测,参考Peng等<sup>[14]</sup>的数据收集方法,对于AKI患者,其数据收集范围为入ICU到AKI发生前24 h、48 h、72 h;对于Non-AKI患者,其数据收集范围为入ICU到出ICU前24 h、48 h、72 h。数据收集窗口如图1所示。

### 1.5 特征构建

根据数据收集时间窗,获取特征参数后,对生命体征与血液学指标,分别计算它们在时间窗内的首检值、最小值、最大值、均值、标准差等,将统计特征纳入特征队列,以反映特征的统计分布特性。最后得到的特征维数为102。

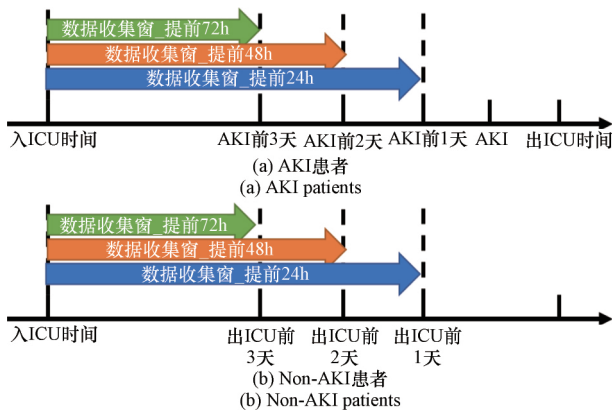


图 1 AKI 患者与 Non-AKI 患者的数据收集窗口

Figure 1 Data collection windows for AKI patients and Non-AKI patients

## 2 研究方法

### 2.1 逻辑回归

逻辑回归<sup>[15]</sup>是一种经典的广义的线性分析模型。

如果用  $x$ 、 $\theta$ 、 $h$ 、 $y$  分别表示训练数据、模型参数、预测输出函数和真实标签,则分类问题实际上是一个伯努利分布:

$$P(y = 1 | x; \theta) = h_{\theta}(x) \quad (1)$$

$$P' = P(y = 0 | x; \theta) = 1 - h_{\theta}(x) \quad (2)$$

式(1)和式(2)可以合并为:

$$P' = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y} \quad (3)$$

将式(3)的最大似然作为目标函数:

$$L(\theta) = \prod_{i=1}^n (h_{\theta}(x^i))^y (1 - h_{\theta}(x^i))^{1-y^i} \quad (4)$$

应用梯度下降法对  $L(\theta)$  求对数,再对模型参数  $\theta$  求偏导。

与线性回归不同,逻辑回归通过引入一个单调可微的 Sigmoid 函数作为输出函数,从而将线性方程预测的连续值映射成 1/0 两个离散值。

逻辑回归的计算代价不高,易于理解和实现,但对模型中自变量多重共线性较为敏感,容易欠拟合,分类精度较低,难以处理数据不平衡的问题。

### 2.2 随机森林

随机森林<sup>[16]</sup>以决策树为基本分类器,利用集成学习的 Bagging 思想,通过有放回的随机抽样,从原始数据集中随机选取数据子集与特征,构建多个决

策树进行分类,其输出类别是单个树输出类别的众数。

随机森林通过将多个分类器进行结合,常可获得比单一学习器更优越的泛化性能,但在噪声比较大的数据上容易过拟合。随机森林的算法流程图如图 2 所示。

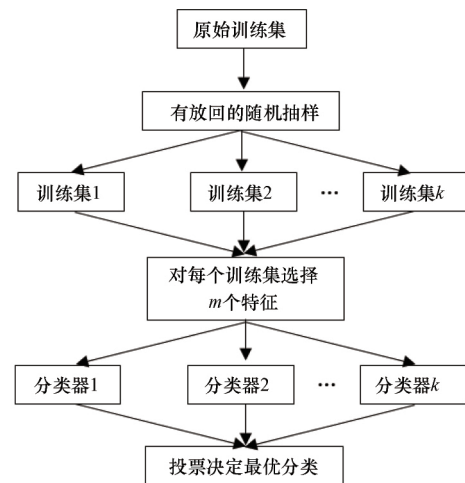


图 2 随机森林流程图

Figure 2 Flowchart of random forest

### 2.3 LightGBM

LightGBM (light gradient boosting machine) 是对梯度提升决策树算法 (gradient boosting decision tree, GBDT) 的优化<sup>[17]</sup>,主要包含两个算法:单边梯度采样 (gradient-based one-side sampling, GOSS) 和互斥特征绑定 (exclusive feature bundling, EFB)。

GOSS 算法通过区分不同梯度的训练数据,保留较大梯度数据的同时对较小梯度数据随机采样,从而达到减少计算量、提升运算效率的目的。定义  $O$  表示某个固定节点的训练集,训练集实例为  $x_1, x_2, \dots, x_n$ ,特征维度为  $s$ ,分割特征为  $j$ ,信息增益为  $V$ <sup>[17]</sup>。每次梯度迭时,模型数据变量的损失函数的负梯度方向表示为  $g_1, g_2, \dots, g_n$ ,则分割特征  $j$  在分割点  $d$  的信息增益为:

$$\hat{V}_{j|O}(d) = \frac{1}{n} \left( \frac{\left( \sum_{\{x_i \in O: x_{ij} \leq d\}} g_i \right)^2}{n_{L|O}^i(d)} + \frac{\left( \sum_{\{x_i \in O: x_{ij} > d\}} g_i \right)^2}{n_{R|O}^i(d)} \right) \quad (5)$$

GOSS 算法中,首先根据数据的梯度将训练数据降序。保留梯度最大的前  $\alpha\%$  的数据,作为数据子集

A,再从剩余的数据中进行随机采样得到数据子集B。

EFB是通过特征捆绑的方式减少特征维度,从而提升计算效率。原始特征个数为 $feature$ ,合并后特征个数为 $bundle$ ,这种方式的特征复杂度从 $O_1(data \times feature)$ 降到 $O_2(data \times bundle)$ ,由于 $bundle$ 远小于 $feature$ ,模型能够极大地加速GBDT的训练过程而不影响最后的精度。

此外,LightGBM将连续的浮点特征值离散化成 $k$ 个整数,构造出宽度为 $k$ 的直方图,当遍历完一次数据后,直方图累积了离散化需要的统计量,之后进行节点分裂时,可以根据直方图上的离散值寻找最佳分割点,减少对内存的消耗。LightGBM还摒弃了大部分GBDT使用的按层生长(level-wise)的决策树生长策略,使用带有深度限制的按叶子生长(leaf-wise)的策略,在分裂次数相同的情况下,可以降低更多的误差,得到更好的精度。

## 2.4 性能评价指标

评价指标采用精确性(precision)、敏感性(sensitivity)、F1值和AUC。

$$\text{precision} = \frac{TP}{TP + FP} \quad (6)$$

$$\text{sensitivity} = \frac{TP}{TP + FN} \quad (7)$$

$$F1 = 2 \times \frac{\text{precision} \times \text{sensitivity}}{\text{precision} + \text{sensitivity}} \quad (8)$$

式中: true negative(TN)称为真阴率,表明实际是负样本预测成负样本的样本数; false positive(FP)称为假阳率,表明实际是负样本预测成正样

本的样本数; false negative(FN)称为假阴率,表明实际是正样本预测成负样本的样本数; true positive(TP)称为真阳率,表明实际是正样本预测成正样本的样本数。

## 3 实验结果

### 3.1 预测结果

将30 020名患者按4:1随机划分训练集和独立测试集,训练集和独立测试集的人数分别为24 016(其中AKI患者13 778名,Non-AKI患者10 238名)和6 004(包括AKI患者3 444名,Non-AKI患者2 560名)。模型训练阶段采用十折交叉验证,训练结束后,使用独立测试集评估所训练模型的性能。

各模型十折交叉验证与独立测试集的结果如表1、表2所示,不同模型在对相同时间点预测时,除了逻辑回归模型预测24 h后的AKI敏感性最高之外,逻辑回归、随机森林、LightGBM模型的性能依次递增;随着预测时间点由24 h、48 h至72 h增加,预测难度上升,同一模型预测性能逐渐下降。LightGBM的精确性和敏感性均比较高,且相差不大,随机森林也类似,因此,二者的F1值和AUC值也较高,但逻辑回归的精确性和敏感性相差较大,即逻辑回归不能有效平衡查准率与查全率,其F1值和AUC值较低。综合比较,LightGBM在3个时间点的预测性能均最优,在独立测试集上24 h、48 h、72 h预测AKI的AUC值分别为0.90、0.88、0.87。不同模型不同时间点预测的ROC曲线如图3所示。

表1 不同时间点十折交叉验证的结果

Table 1 Results of 10-fold cross-validation of different time points

预测时间	模型	精确率(95%CI)	召回率(95%CI)	F1值(95%CI)	AUC(95%CI)
24 h	逻辑回归	0.69 (0.68-0.70)	0.93 (0.92-0.94)	0.79 (0.78-0.80)	0.82 (0.81-0.83)
	随机森林	0.88 (0.87-0.89)	0.88 (0.87-0.89)	0.88 (0.87-0.89)	0.94 (0.93-0.95)
	LightGBM	0.92 (0.91-0.93)	0.89 (0.88-0.90)	0.90 (0.89-0.91)	0.96 (0.95-0.97)
48 h	逻辑回归	0.71 (0.70-0.72)	0.82 (0.81-0.83)	0.76 (0.75-0.77)	0.80 (0.79-0.81)
	随机森林	0.88 (0.87-0.89)	0.85 (0.84-0.86)	0.86 (0.85-0.87)	0.93 (0.92-0.94)
	LightGBM	0.90 (0.89-0.91)	0.86 (0.75-0.87)	0.88 (0.87-0.89)	0.94 (0.93-0.95)
72 h	逻辑回归	0.82 (0.81-0.83)	0.40 (0.39-0.41)	0.54 (0.53-0.55)	0.80 (0.79-0.81)
	随机森林	0.86 (0.85-0.87)	0.81 (0.79-0.83)	0.83 (0.82-0.84)	0.92 (0.90-0.94)
	LightGBM	0.86 (0.78-0.94)	0.84 (0.79-0.89)	0.84 (0.80-0.88)	0.94 (0.93-0.95)

表 2 不同时间点独立测试集的结果

Table 2 Results of independent test sets at different time points

预测时间	模型	精确率	召回率	F1 值	AUC
24 h	逻辑回归	0.67	0.94	0.79	0.66
	随机森林	0.88	0.89	0.89	0.86
	LightGBM	0.92	0.89	0.91	0.90
48 h	逻辑回归	0.71	0.84	0.77	0.72
	随机森林	0.89	0.85	0.87	0.86
	LightGBM	0.91	0.86	0.88	0.88
72 h	逻辑回归	0.82	0.35	0.50	0.65
	随机森林	0.87	0.81	0.84	0.86
	LightGBM	0.87	0.84	0.86	0.87

应用性能最优的 LightGBM 对 AKI 患者进行连续 24 h 预测,即从入 ICU 第 1 天开始,预测 24 h 后发生 AKI 的风险,第 2 天则根据当天及当天之前的数据,进行 24 h 预测,直至转出 ICU 为止。将临床确诊为 AKI 的时间减去模型首次预测 AKI 成功的时间,作为提前时间(天数),统计预测成功的天数与比例(表 3)。表 3 说明,LightGBM 连续 24 h 预测的成功率较高,提前 1、2、3 天预测 AKI 风险患者的成功率分别为 89%、83% 和 80%。当对 1 587 名根据 KDIGO 标准确诊的 AKI 患者,模型能提前 3 天获知其中 1 272 名患者具有 AKI 风险,此时给临床医生预留了 3 天的时间进行干预治疗。

表 3 LightGBM 提前 1~3 天预测 AKI 的成功率

Table 3 Success rates of AKI prediction of 1~3 days in advance using LightGBM

提前预测天数	AKI 患者人数	预测成功人数	成功率
1	3 444	3 079	0.89
2	2 695	2 240	0.83
3	1 587	1 272	0.80

### 3.2 重要特征

特征重要性可以反映每个特征对模型预测能力的贡献程度。根据特征在模型训练过程中被使用的次数,得到 LightGBM 24 h、48 h 与 72 h 3 个时间点的所有特征权重列表,选择排列前 35 的特征,如图 4 所示。

3 个时间点预测模型得到的位居前列的特征大部分相同,均包括已住院时长、体质量、体温最小值、白细胞计数最大值/最小值、碳酸氢盐最大值、葡萄糖最小值、舒张压最大值、血尿素氮最小值/最大值、APS III 评分、体温最大值/首测值、血红蛋白首测值、心率首测值、血清肌酐最小值/最大值、收缩压最大值、心率首测值/最大值、葡萄糖最大值/首测值、收缩压首测值和氯最大值等 24 个特征。

值得注意的是,3 个时间点预测模型中:已住院时长、体质量均位列前 2,提示在早期预测 AKI 时,应首先观察这两个指标;白细胞计数最大值、碳酸氢盐最大值、体温最小值均位列前 10,需重点关注;而血清肌酐的重要性均未排入前 10,表明血清肌酐对模型的作用不够理想,仅以肌酐作为 AKI 诊断标准敏感性不足,具有一定滞后性。

为进一步验证重要特征的作用,并对特征进行降维,在仅使用 24 个重要特征的情况下,对 LightGBM 不同时间点预测模型进行训练和测试。独立测试集的结果如表 4 所示。

由表 4 结果可知,当仅使用 24 个重要特征时,LightGBM 预测模型 24 h、48 h、72 h 的 AUC 值分别为 0.88、0.86、0.85。与使用全部特征相比,AUC 值降低不超过 2 个百分点。说明在 102 维全部特征当中,24 个重要特征对 LightGBM 模型的预测性能贡献了绝大部分作用,仅使用重要特征也能对 AKI 进行连续、有效的预测。

表 4 LightGBM 仅使用 24 个重要特征的预测结果

Table 4 Prediction results of LightGBM using only 24 important features

预测时间/h	精确率	召回率	F1 值	AUC
24	0.89	0.87	0.89	0.88
48	0.89	0.85	0.87	0.86
72	0.86	0.81	0.85	0.85

## 4 讨论

实验结果表明,3 种机器学习模型中,LightGBM 模型性能最优,随机森林次之,逻辑回归较差。使用性能最佳的 LightGBM 可对 AKI 风险患者进行连续预测,仅使用重要特征时仍有较好的性能。

与以往的 AKI 预测研究相比<sup>[8-12]</sup>,本研究的主

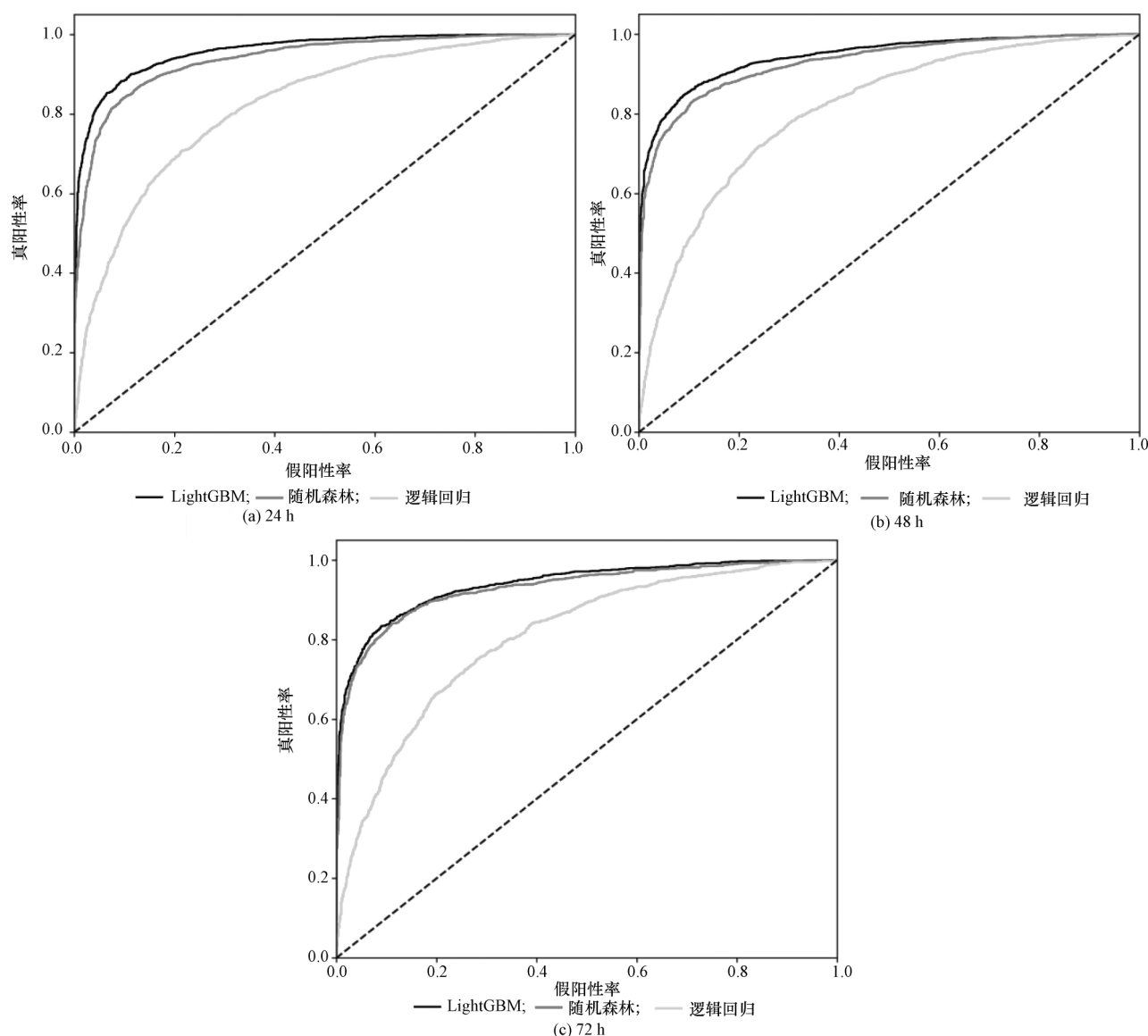


图3 3种模型AKI预测的ROC比较

Figure 3 Comparison of the ROC obtained by three prediction models

要优势为:(1) 纳入较大样本量。总样本量、AKI患者数均大于上文提及的相关研究,样本平衡性更好,结果可靠性更高。(2) 使用8大类临床信息构造特征,降低了遗漏关键因素的可能性,并对特征降维,使用重要特征进行预测,进一步验证了本文得到的重要特征的作用。(3) 构建了24 h、48 h、72 h 3个时间点预测模型,并对AKI风险患者进行连续预测,不仅能持续监测患者病情,也为临床医生进行干预与治疗预留了更多时间。

LightGBM 24 h 预测模型中,排列前10的主要

特征包括已住院时长、体质量、白蛋白首测值、白细胞计数最大值、收缩压最小值、碳酸氢盐最大值、葡萄糖最小值、白细胞计数最小值、体温最大值、血尿素氮最小值。3个时间点预测模型中,位居前30的特征绝大部分相同。

ICU病房中,患者住院时间长,说明患者健康状况严峻、疾病更为复杂,从而导致AKI的潜在风险升高。体质量作为最常见且易得到的指标,本实验数据中,患者体质量每增加5 kg,AKI发病率升高3.74%。

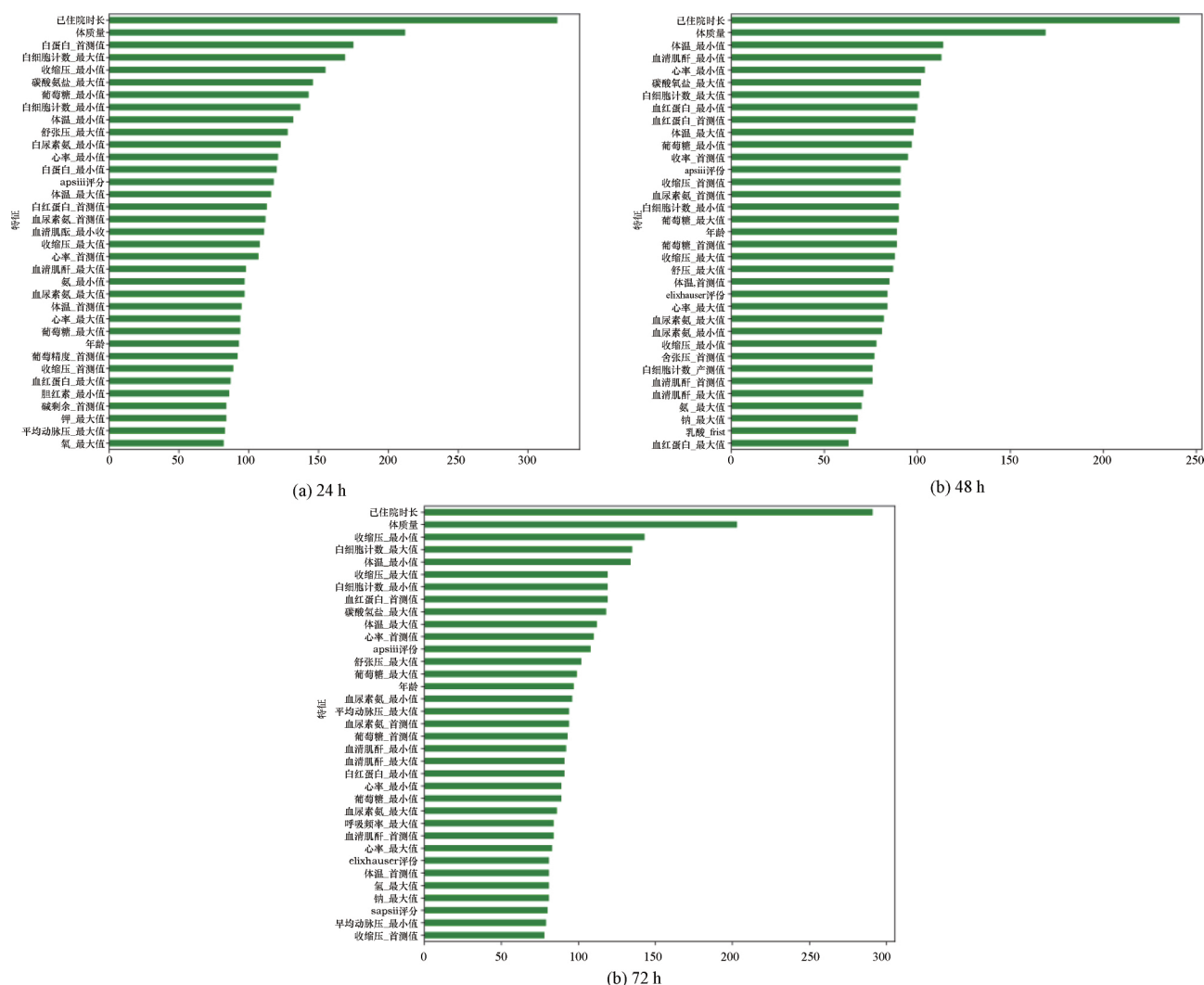


图 4 LightGBM 3 个时间点 AKI 预测的前 35 重要特征

Figure 4 Top 35 important features obtained by LightGBM at three time points

体温过低会导致肾血流量下降、肾小管功能受损,可引起酸中毒和碱中毒<sup>[18-19]</sup>。在 ICU 内,超过 40% 的体温过低患者会出现 AKI<sup>[18]</sup>。

白细胞在炎症反应、宿主防御和修复中具有重要作用,是大多数器官损伤过程中起关键作用的免疫学因素之一。白细胞计数与 AKI 风险之间呈现 U 型关系,白细胞计数降低导致的高 AKI 风险可能归因于淋巴细胞和单核细胞的减少,白细胞计数升高导致的高 AKI 风险可能归因于中性粒细胞的增多<sup>[20]</sup>。

血清中碳酸氢盐有助于增加向肾脏的氧输送,同时能够中和肾脏部分的酸中毒,而较低的碳酸氢盐水平会增加肾脏缺血性损伤的风险,特别是在危

重症情况下<sup>[21]</sup>。

本研究尚存在以下不足:(1) 样本来自单中心数据库,模型鲁棒性有待多中心数据的进一步验证;(2) 仅预测 AKI 与非 AKI,没有对 AKI 按 KIDGO 诊断标准进行分级(I-III 级)预测。

## 5 结论

本研究基于 MIMIC-III 数据库,提取 30 020 名患者的人口学信息、入院信息、生命体征、危重症评分、合并症、血液学指标、用药及干预措施等 8 类临床信息,分别使用逻辑回归、随机森林和 LightGBM 3 种机器学习算法,建立了 24 h、48 h、72 h 3 个时间点的 AKI 预测模型,比较不同模型的预测性能,获



取重要特征。结果表明 LightGBM 的预测性能最优,且在连续预测 AKI 风险患者时具有高达 80% 的识别率,在仅使用重要特征进行预测时仍有较高的性能。研究结果可对 ICU 患者发生 AKI 风险提供连续、有效的预测,明确重要影响因素,并为医护人员进行及时合理的干预提供重要指导。

#### 参考文献

- [1] Lameire NH, Bagga A, Cruz D, et al. Acute kidney injury: an increasing global concern [J]. *Lancet*, 2013, 382 ( 9887 ): 170-179.
- [2] Bagshaw SM, George C, Gibney RTN, et al. A multi-center evaluation of early acute kidney injury in critically ill trauma patients [J]. *Renal Failure*, 2008, 30( 6 ): 581-589.
- [3] Mizuno T, Sato W, Ishikawa K, et al. KDIGO ( kidney disease: improving global outcomes ) criteria could be a useful outcome predictor of cisplatin-induced acute kidney injury [J]. *Oncology*, 2012, 82( 6 ): 354-359.
- [4] Eriksson M, Brattström O, Mårtensson J, et al. Acute kidney injury following severe trauma: risk factors and long-term outcome [J]. *Journal of Trauma and Acute Care Surgery*, 2015, 79( 3 ): 407-412.
- [5] Ostermann M, Joannidis M. Acute kidney injury 2016: diagnosis and diagnostic workup [J]. *Critical Care*, 2016, 20: 299.
- [6] Panagidis D, Nanas S, Kokkoris S. Biomarkers of acute kidney injury in a mixed ICU population. a narrative review [J]. *Health & Research Journal*, 2019, 5( 4 ): 150.
- [7] Rojas JC, Carey KA, Edelson DP, et al. Predicting intensive care unit readmission with machine learning using electronic health record data [J]. *Annals of the American Thoracic Society*, 2018, 15( 7 ): 846-853.
- [8] Haines RW, Lin SP, Hewson R, et al. Acute kidney injury in trauma patients admitted to critical care: development and validation of a diagnostic prediction model [J]. *Scientific Reports*, 2018, 8: 3665.
- [9] Malhotra R, Kashani KB, Macedo E, et al. A risk prediction score for acute kidney injury in the intensive care unit [J]. *Nephrology Dialysis Transplantation*, 2017, 32( 5 ): 814-822.
- [10] 李千惠. 基于机器学习的急性肾损伤预测及临床应用优化 [D]. 北京: 北京交通大学, 2019.  
Li QH. Prediction of acute kidney injury and clinical application optimization based on machine learning [D]. Beijing: Beijing Jiaotong University, 2019.
- [11] 张渊, 冯聪, 李开源, 等. ICU 患者急性肾损伤发生风险的 LightGBM 预测模型 [J]. *解放军医学院学报*, 2019, 40( 4 ): 316-320.  
Zhang Y, Feng C, Li KY, et al. LightGBM model for predicting acute kidney injury risk in ICU patients [J]. *Academic Journal of Chinese PLA Medical School*, 2019, 40( 4 ): 316-320.
- [12] Zimmerman LP, Reifman PA, Smith A, et al. Early prediction of acute kidney injury following ICU admission using a multivariate panel of physiological measurements [J]. *BMC Medical Informatics and Decision Making*, 2019, 19( Suppl 1 ): 16.
- [13] Johnson A, Pollard TJ, Shen L, et al. MIMIC - III, a freely accessible critical care database [J]. *Scientific Data*, 2016( 3 ): 160035.
- [14] Peng C, Waitman LR, Yong H, et al. Predicting inpatient acute kidney injury over different time horizons: How early and accurate? [J]. *AMIA... Annual Symposium Proceedings/ AMIA Symposium*, 2018, 2017( 2017 ): 565-574.
- [15] van Houwelingen S. Ridge estimators in logistic regression [J]. *Journal of the Royal Statistical Society Series C*, 1992, 41( 1 ): 191-201.
- [16] Criminisi A, Shotton J, Konukoglu E. Decision forests: a unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning [J]. *Foundations and Trends in Computer Graphics and Vision*, 2011, 7( 2-3 ): 81-227.
- [17] Ke G, Meng Q, Finley T, et al. LightGBM: a highly efficient gradient boosting decision tree [C]//NIPS' 17: the 31st International Conference on Neural Information Processing Systems. Long Beach, USA: NIPS, 2017, 3149-3157.
- [18] Mallet ML. Pathophysiology of accidental hypothermia [J]. *QJM*, 2002, 95( 12 ): 775-785.
- [19] Mégarbane B, Axler O, Chary I, et al. Hypothermia with indoor occurrence is associated with a worse outcome [J]. *Intensive Care Medicine*, 2000, 26( 12 ): 1843-1849.
- [20] Han SS, Ahn SY, Ryu J, et al. U-shape relationship of white blood cells with acute kidney injury and mortality in critically ill patients [J]. *The Tohoku Journal of Experimental Medicine*, 2014, 232( 3 ): 177-185.
- [21] Gujadhur A, Tiruvoipati R, Cole E, et al. Serum bicarbonate may independently predict acute kidney injury in critically ill patients: an observational study [J]. *World Journal of Critical Care Medicine*, 2015, 4( 1 ): 71-76.

( 2021-01-20 收稿, 2021-04-08 修回 )