

文章编号: 1003-0077(2024)01-0156-10

## 非结构化数据表征增强的术后风险预测模型

王亚强<sup>1,2,3</sup>, 杨潇<sup>1,2,3</sup>, 郝学超<sup>4</sup>, 舒红平<sup>1,3</sup>, 陈果<sup>4</sup>, 朱涛<sup>4</sup>

(1. 成都信息工程大学 软件工程学院, 四川 成都 610225;

2. 成都信息工程大学 数据科学与工程研究所, 四川 成都 610225;

3. 成都信息工程大学 软件自动生成与智能服务四川省重点实验室, 四川 成都 610225;

4. 四川大学华西医院 麻醉手术中心, 四川省 成都市 621005)

**摘要:** 准确的术后风险预测对临床资源的规划、应急方案的准备以及患者术后风险和死亡率的降低具有积极的作用。目前, 术后风险预测主要基于患者的基本信息、术前的实验室检查及术中的生命体征等结构化数据, 蕴含着丰富语义信息的非结构化术前诊断的价值尚待验证。针对上述问题, 该文提出一种非结构化数据表征增强的术后风险预测模型, 利用自注意力机制, 将结构化数据与术前诊断进行信息加权融合。基于临床数据, 该文将所提出的模型与术后风险预测常用的统计机器学习模型以及最新的深度神经网络进行对比, 在肺部并发症风险预测、ICU 入室风险预测和心血管不良风险预测任务上的  $F_1$  值平均提升了 9.533%, 同时预测模型还具有良好的可解释性。

**关键词:** 术后风险预测; 自注意力机制; 数据表征; 信息融合

中图分类号: TP391

文献标识码: A

## An Unstructured Data Representation Enhanced Method for Postoperative Risk Prediction

WANG Yaqiang<sup>1,2,3</sup>, YANG Xiao<sup>1,2,3</sup>, HAO Xuechao<sup>4</sup>, SHU Hongping<sup>1,3</sup>, CHEN Guo<sup>4</sup>, ZHU Tao<sup>4</sup>

(1. College of Software Engineering, Chengdu University of Information Technology, Chengdu, Sichuan 610225, China;

2. Institute for Data Science and Engineering, Chengdu University of Information Technology, Chengdu, Sichuan 610225, China;

3. Sichuan Key Laboratory of Software Automatic Generation and Intelligent Service, Chengdu University of Information Technology, Chengdu, Sichuan 610225, China;

4. Department of Anesthesiology, Sichuan University, Chengdu, Sichuan 621005, China)

**Abstract:** Postoperative risk prediction has a positive effect on clinical resource plan, emergency plan preparation and postoperative risk and mortality reduction. To employ the unstructured preoperative diagnosis with rich semantic information, this paper proposes a postoperative risk prediction model via unstructured data representation enhancement. The model utilizes self-attention to fuse the structured data with unstructured preoperative diagnosis. Compared with the baseline methods, the proposed model improves  $F_1$ -Score by an average of 9.533% on the tasks of the pulmonary complication risk prediction, the ICU admission risk prediction and the cardiovascular adverse risk prediction.

**Keywords:** postoperative risk prediction; self-attention mechanism; data representation; information fusion

收稿日期: 2023-03-19 定稿日期: 2023-07-04

基金项目: 四川大学华西医院 1·3·5 项目(ZYJC21008); 国家重点研究与发展计划项目(2018YFC2001800)

## 0 引言

术后并发症(如肺部并发症<sup>[1]</sup>、心血管不良<sup>[2]</sup>、ICU 入室<sup>[3]</sup>等)风险(后文简称“术后风险”)所导致的术后 30 天内死亡,已成为全球排名第三位的人群死亡原因<sup>[4]</sup>。准确的术后风险预测对医生进行合理的临床资源规划、应急方案准备具有重要的辅助作用,对患者的术后风险发生和死亡率降低具有积极意义<sup>[5,6]</sup>。

目前,术后风险预测主要基于患者的基本信息(如体温、血压、体重等)、术前的实验室检查(如氧分压、氧饱和、蛋白等)、术中的生命体征(如出血量等)等结构化数据,利用极限梯度提升(eXtreme Gradient Boosting, XGBoost)、逻辑回归(Logistic Regression, LR)、随机森林、人工神经网络等模型实现<sup>[2,5]</sup>。

近年来,深度神经网络在各领域的预测任务中表现优秀,受到研究者的广泛关注,也被引入术后风险预测任务<sup>[6]</sup>。Fritz<sup>[7]</sup>等人构建了一种多路径卷积神经网络,提取和融合患者基本信息、共病情况、术前实验室检查和术中生命体征等结构化数据中的特

征,用于患者术后死亡风险预测。Barbieri<sup>[8]</sup>等人利用双向门控循环单元,将结构化数据之间的时间信息以拼接的方式融入数据表征,采用注意力机制提取重要特征,用于患者术后 ICU 入室风险预测。现有方法的核心是如何将结构化数据中的离散型和连续型特征向量化,形成基于深度神经网络的术后风险预测模型的数据表征。

在术前数据中,除结构化数据外,还包含语义丰富的非结构化术前诊断数据。术前诊断中不仅包含医生基于的医学知识,还包含根据局部的结构化数据,对患者病情的总结信息,以及医生以整体的结构化数据为依据,利用经验知识,对患者病况的推断信息。如图 1 中患者 1 的术前数据所示,根据结构化数据收缩压 156 mmHg(毫米汞柱)与舒张压 76 mmHg,基于医学知识“成人的收缩压和舒张压正常范围应在 90 mmHg 至 120 mmHg 之间”,因此,医生在术前诊断中总结该患者有“高血压病”,且属于“3 级很高危”。此外,依据目前患者整体的结构化数据,医生根据经验知识,推断患者是“肺部感染”。更进一步地,术前诊断的整体描述,反映了当前患者的全局状态。这些语义信息能够丰富术后风险预测的特征,有助于增强预测模型的性能。

编号	体温	是否使用活性药物	基于规范知识总结的患者状态信息		术前诊断	全局语义信息
			收缩压	舒张压		
患者1	36.5	1	156	76	1: 高血压病(3级 很高危) 2: 肺部感染	
患者2	36.4	0	113	70	直肠恶性肿瘤	
患者3	36.7	0	105	66	左膝重度关节炎	

图 1 结构化的患者基本信息和术前实验室检查数据

然而,术前诊断数据尚未在术后风险预测任务中被有效利用。如何充分地利用非结构化的术前诊断数据,形成有效的术后风险预测数据表征,尚有待进一步探索。

综上,本文围绕非结构化的术前诊断数据如何增强术后风险预测任务这一问题展开研究,主要的贡献包括以下三个方面:

(1) 与围术期医学专家合作,经过清洗、处理、转换和去隐私过程,构建了一份包含 12 240 个实例、面向术后风险预测任务的数据集。该数据集的结构化数据部分包含了 95 列离散型变量、61 列连续型变量、一列非结构化的术前诊断变量以及三列二元的术后风险标签变量,分别表示肺部并发症、心血管不良和 ICU 入室风险的发生情况。

(2) 为充分地利用非结构化的术前诊断数据,本文提出一种非结构化数据表征增强的术后风险预测模型,利用自注意力机制,将结构化数据与局部的细粒度实体信息及全局的粗粒度文本语义加权融合,有效地将非结构化数据用于增强术后风险预测性能。

(3) 本文提出的基于自注意力机制融合结构化与非结构化数据的模型结构,为术后风险预测带来了良好的可解释性。细节实验结果分析发现,利用自注意力机制获得的关系权重矩阵,可以解释和展示出非结构化数据,不仅增强了重要的结构化数据的贡献度,而且还补充了风险预测信息。

实验结果表明,本文提出的非结构化数据表征增强的术后风险预测模型明显优于所对比的常用统

计机器学习模型和最新的深度神经网络,在三种重要的术后风险预测(包括肺部并发症风险预测、ICU入室风险预测和心血管不良风险预测)任务上,本文提出的模型均取得了最优的结果, $F_1$ 值分别达到了66.909%、60.833%和55.888%。此外,通过消融实验,进一步验证了本文提出的模型有效地加权融合了局部的细粒度实体信息和全局的粗粒度文本语义信息。利用非结构化术前诊断数据表征增强术后风险预测模型后,肺部并发症风险预测的 $F_1$ 值提升了6.878%,ICU入室风险预测提升了7.641%,心血管不良风险预测提升了9.541%。

## 1 相关工作

术后风险预测是医学信息学领域的研究热点问题。当前的研究主要集中在验证统计机器学习模型在术后风险预测任务上的有效性,以及面向特定类型的术后并发症风险的特征分析两个层面。Canet<sup>[9]</sup>等人利用逻辑回归模型,确定了7个独立且具有良好的鉴别能力的危险因素后,构建了术后肺部并发症风险预测指标,用于评估和预测术后肺部并发症的个体风险。Hill<sup>[10]</sup>等人采用随机森林模型,自动地发现重要的术前特征,将结构化的美国麻醉医师协会身体状况特征与术前特征相结合,提升术后死亡风险的预测性能。与先前工作不同,本文提出了一种非结构化数据表征增强的术后风险预测模型,该模型基于自注意力机制,在预测中有效地融合结构化数据和非结构化语义信息,并提供良好的可解释性。

术后风险预测目前的主要研究对象是术前和术中的结构化数据,其中包含两种类型的变量,一种是离散型变量,另一种是连续型变量。其中连续型变量通常会被离散化后,与离散型变量一同构建特征向量,作为术后风险预测模型的输入<sup>[11]</sup>。本文的实验主要基于结构化的患者基本信息和术前的实验室检查数据。本文采用与先前工作相同的连续型变量的基本处理方法。差异在于本文借鉴Fritz<sup>[7]</sup>等人的思想,将离散型变量和离散化的连续型变量构建离散特征词典,并基于深度神经网络学习离散特征的嵌入表征。

术后风险预测除可利用术前和术中的结构化数据作为特征之外,通过观察发现,包含医学语义信息的非结构化术前诊断数据也可用于增强术后风险预测。Zhang<sup>[12]</sup>等人提出将英文临床文本利用Doc2Vec

模型<sup>[13]</sup>直接形成数据表征,然后与结构化数据合并的方式,将非结构化数据与结构化数据融合,应用于住院死亡率、住院时间长短和术后30天再入院的预测任务,该方法在英文临床数据MIMIC-III<sup>[14]</sup>上进行了实验验证。与该工作不同,本文首次探索了将中文非结构化临床文本引入术后风险预测的方法。

此外,本文通过观察还发现,在非结构化的术前诊断中,既包含全局的粗粒度文本语义信息,还包含局部的细粒度医学实体信息,它们均可作为术后风险预测提供医学语义特征(如图1所示)。为将这些信息与离散特征的嵌入表征相融合,本文首先基于常用的中文MedBERT<sup>①</sup>获得实体的嵌入表征,并将术前诊断视为句子后,采用词嵌入平均池化的方法将其向量化。然后利用自注意力机制<sup>[15]</sup>,将离散特征的嵌入表征与实体的嵌入表征以及向量化的术前诊断进行加权融合,在综合地利用全局和局部的文本语义信息的基础上,还为模型带来了良好的可解释性<sup>[16]</sup>。

## 2 术后风险预测

### 2.1 任务定义

本文将术后风险预测定义为一项二分类任务,采用有监督学习方法解决。定义 $(x, y)$ 为一个训练实例, $x$ 中包含 $x_{num}$ 、 $x_{cat}$ 和 $x_{PD}$ 三种类型的特征。其中, $x_{num}$ 表示表格数据中的连续型特征,共 $m$ 列, $x_{cat}$ 表示表格数据中的离散型特征,共 $n$ 列, $x_{PD}$ 表示非结构化的术前诊断文本数据, $y$ 表示术后风险发生的情况,用1或0分别表示风险的发生或未发生。

### 2.2 表格数据的向量表征

本文提出的术后风险预测模型主要利用结构化表格数据和非结构化术前诊断文本数据对术后风险进行预测(模型的结构如图2所示)。结构化表格数据由 $x_{num}$ 和 $x_{cat}$ 组成。本文采用分类与回归树算法<sup>[17]</sup>,先将连续型特征转换为离散型特征,在引入医学语义信息的同时,降低数据的复杂度。转换后的连续型特征不仅能够表达医学语义,还被统一成离散型特征。转换后的连续型变量表征被定义为

① URL: <https://code.ihub.org.cn/projects/1775>

$x_{n2cat}$  :

$$x_{n2cat} = discretize(x_{num}) \quad (1)$$

处理离散型变量表征的常用方式是采用实体嵌入<sup>[18]</sup>的方法,即为每一个离散型变量构建一个特征词表,词表大小为当前离散型变量的不同取值的数量。然而该方法在建模的过程中仅考虑了单一变量下的不同取值之间的语义关联,而不同的变量之间的相关性未被考虑其中。为引入全局不同变量之间的语义关联,本文改进了原始的实体嵌入方法,让所有的离散型变量共用特征词表。每一个离散型变量(包括  $x_{cat}$  和  $x_{n2cat}$ )的不同取值,都会被赋予唯一的

索引值  $x_i$ ,其中  $i \in [0, |V|]$ , $|V|$ 是所有的离散型变量的不同取值的数量总和,即共用的特征词表的词表大小。每个  $x_i$  都将通过学习过程被映射为一个维度为  $d$  的向量,定义为  $e_{tabular}$ ,其中  $d$  为超参数。通过构建全局共用的特征词表,原始的离散型变量转换为语义向量之后,不仅扩充了医学语义信息,并且不同的离散型变量之间也产生了语义关联。相比原始的实体嵌入方法,该方法解决了不同离散型变量之间语义关联缺失的问题。最后,将所有的  $e_{tabular}$  拼接形成表格数据的向量表征  $E_{tabular}$ 。

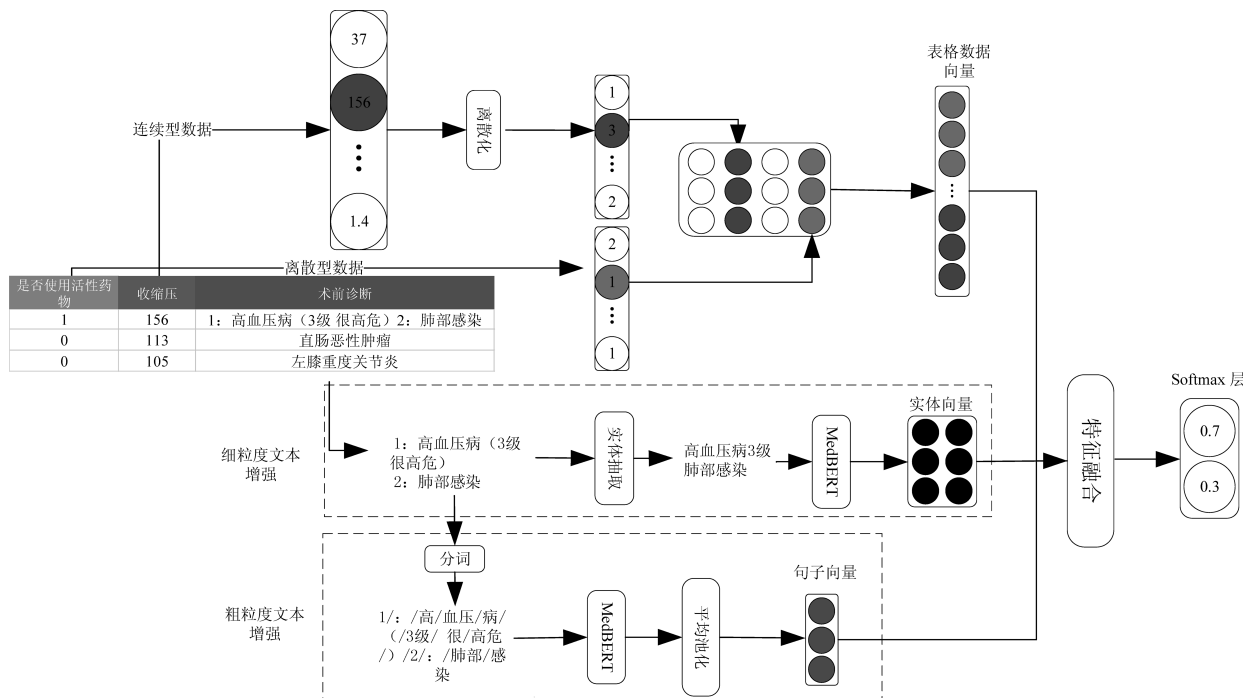


图2 模型结构图

### 2.3 文本数据的向量表征

术前诊断文本  $x_{PD}$  主要包含医生总结的病人身体症状和初步推断的病情描述,两者可统一定义为病症实体。因此,术前诊断文本可以归纳为由多个病症实体、连接词以及标点符号构成的集合,每个实例  $x_{PD}$  包含  $l_{max}$  项的病症实体, $l_{max}$  表示数据集中,  $x_{PD}$  中最多可包含的病症实体数量。

术前诊断文本可以有两种向量表征方法,一种是形如利用 Doc2Vec 模型<sup>[13]</sup>得到的全局语义向量,获取该类向量表征的方法我们称之为粗粒度文本的向量表征方法;另一种是直接将病症实体对应的语义向量拼接,形成细粒度文本的向量表征。后文将具体介绍它们获取术前诊断文本粗粒度语义信息和

细粒度语义信息的方法。

#### 2.3.1 粗粒度语义向量表征方法

为获取术前诊断文本的粗粒度语义向量表征,本文先将文本进行了分词<sup>①</sup>,得到分词列表  $\{token_0, token_1, \dots, token_p\}$ ,其中  $p$  表示文本分词后得到的词的数量。将分词列表输入领域微调后的预训练模型 MedBERT 中,生成维度为 768 的动态词向量列表  $\{e_0^{768}, e_1^{768}, \dots, e_p^{768}\}$ 。其中,768 是 MedBERT 的词向量维度。为进一步获取句子向量,本文采用平均池化的方法整合词向量的语义信息。对词向量矩阵中的每一列求均值,将词向量矩阵压缩为包含整

① 本文实验中直接采用了 <https://huggingface.co/hfl/chinese-macbert-base> 的内置分词工具

个术前诊断语义信息的粗粒度语义向量表征  $e_{\text{sentence}}$ , 如式(2)所示。

$$e_{\text{sentence}} = \text{MeanPooling}(\{e_0^{768}, \dots, e_p^{768}\}) \quad (2)$$

### 2.3.2 细粒度语义向量表征方法

将术前诊断文本分词后, 通过 MedBERT 生成的词向量被压缩为单一向量, 会导致局部语义信息的丢失, 且无法明确术前诊断文本中哪些信息在术

后风险预测过程中起到了关键作用。为保留术前诊断文本中的局部细粒度实体语义信息, 本文首先利用医学领域数据集, 基于 BERT+BiLSTM+CRF 模型训练得到实体抽取模型<sup>[19]</sup>, 然后利用该模型抽取  $x_{\text{PD}}$  中的病症实体, 形成病症实体集合  $\{\omega_0, \dots, \omega_k, \dots, \omega_K\}$ , 其中,  $K$  表示当前  $x_{\text{PD}}$  中抽取得到的病症实体数量。

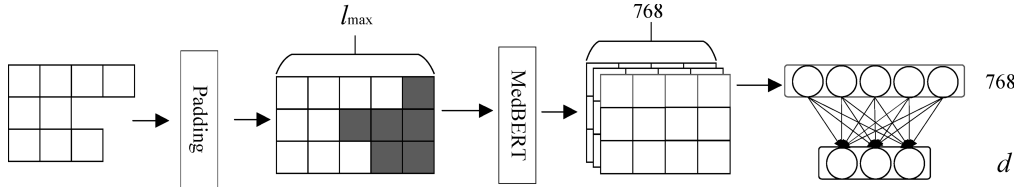


图3 细粒度语义向量的补全与降维转换方法

因为  $x_{\text{PD}}$  中包含的实体数量可能不一致, 为后续处理统一, 本文将病症实体数量未达到  $l_{\text{max}}$  的集合, 通过补全特殊字符[PAD]的方式, 形成数量均为  $l_{\text{max}}$  的实体集合(如图3所示)。然后, 每一个实体  $\omega_k$  将通过 MedBERT 转换为蕴含医学语义的向量  $e_k^{768}$ 。为后续与表格数据的向量表征进行融合, 细粒度语义向量进一步通过全连接层降维, 从 768 维降至  $d$  维, 得到降维后的细粒度语义向量集合  $\{e_0^d, e_1^d, \dots, e_{l_{\text{max}}}^d\}$ 。最后, 将含有全局语义信息的粗粒度向量表征和含有局部语义信息的细粒度向量表征组合, 得到最终的术前诊断文本的向量表征  $E_{\text{text}}$  如式(3)所示。

$$E_{\text{text}} = \{e_0^d, e_1^d, \dots, e_{l_{\text{max}}}^d, e_{\text{sentence}}\} \quad (3)$$

## 2.4 特征融合方法

在特征融合层, 本文选择采用 Self-Attention

机制<sup>[15]</sup>将表格数据表征  $E_{\text{tabular}}$  与文本数据的向量表征  $E_{\text{text}}$  进行特征融合(如图4所示)。首先, 将表示表格数据信息的数据表征  $E_{\text{tabular}}$  与表示文本语义信息的数据表征  $E_{\text{text}}$  拼接, 形成新的特征向量集合  $E_X$ , 并将  $E_X$  通过三个参数矩阵  $W^Q$ 、 $W^K$  和  $W^V$  映射为三个不同的矩阵  $Q$ 、 $K$  和  $V$ 。然后对  $Q$  和  $K^T$  执行点积并利用  $d_k$  放缩结果, 以保证训练过程中梯度的稳定性。其中,  $d_k$  是指矩阵  $K$  的维度, 计算方法如公式(4)的 Softmax 函数的输入所示。随后执行 Softmax 函数进行归一化, 得到不同的数据表征之间(包含表格数据表征和文本数据表征)的注意力权重  $W_{\text{weight}}$ , 其计算方法如式(4)所示。

$$W_{\text{weight}} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \quad (4)$$

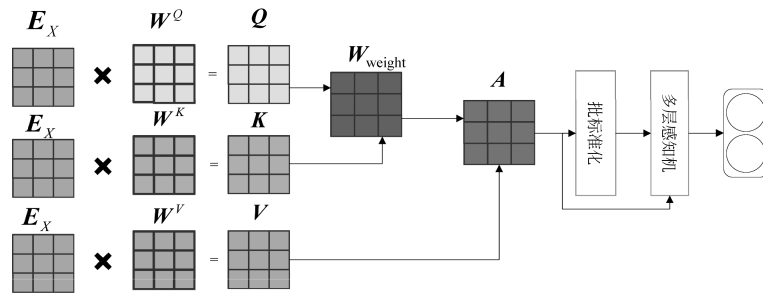


图4 特征融合方法

最后将  $W_{\text{weight}}$  与  $V$  相乘得到增强后的特征表示  $A$ 。具体计算过程如式(5)~式(7)所示。

$$E_X = E_{\text{tabular}} \oplus E_{\text{text}} \quad (5)$$

$$Q = E_X W^Q, \quad K = E_X W^K, \quad V = E_X W^V \quad (6)$$

$$A = \text{Attention}(Q, K, V) = W_{\text{weight}} V \quad (7)$$

通过注意力机制, 模型可以自动地学习到特征在推理过程中的重要性或贡献度。因此, 在模型推理过程中, 可以通过提取并分析注意力权重矩阵, 来

探究在模型预测过程中,各特征发挥作用的重要程度,从而为模型带来良好的可解释性。

为了解决梯度消失问题,受文献[20]和[21]的启发,表征矩阵  $\mathbf{A}$  在输入前馈神经网络之前,还经过了残差网络和层标准化操作。接着将向量输入到带有 sigmoid 激活函数的前馈神经网络中,计算预测术后风险的发生概率  $P$  如式(8)所示。

$$P = \text{sigmoid}(\mathbf{W}^T \mathbf{A} + \mathbf{b}) \quad (8)$$

在公式(8)中, $\mathbf{W}$  和  $\mathbf{b}$  是前馈神经网络将学习的参数。最后,模型的损失定义为如式(9)所示。

$$\text{loss} = -\frac{1}{M} \sum_{i=1}^M (y_i \log P_i + (1 - y_i) \log(1 - P_i)) \quad (9)$$

其中, $M$  指批量包含的实例数量。

### 3 实验

#### 3.1 实验数据

本文实验采用了从医院的临床管理系统中获取的数据,其中包含患者的基本信息、术前实验室检查数据和术前诊断,以及病人术后发生的肺部并发症风险、ICU 入室风险和心血管不良风险结局。该数据经过了如下基本的预处理过程:

- (1) 删除了有关患者身份的个人信息;
- (2) 删除了缺失率高于 50% 的变量。

最终得到包含 12 240 个实例的术后风险预测数据集,该数据集中包含 95 项离散型变量和 61 项连续型变量以及 1 项术前诊断变量。数据集中包含的三种术后风险的标签分布如图 5 所示,肺部并发症风险的阳性率为 15.93%,ICU 入室风险的阳性率为 6.25%,心血管不良风险的阳性率为 3.02%。实验中,本文将数据集按照 7 : 1 : 2 的比例划分得到训练集、验证集和测试集。

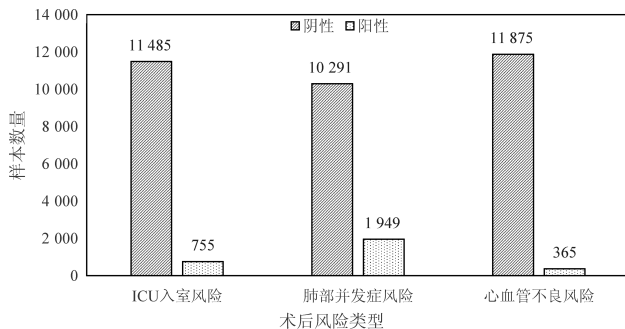


图 5 实验数据中的标签分布

#### 3.2 评估指标

为了评估模型的效果,本文采用精确率(Precision)、召回率(Recall)和  $F_1$  值作为主要的评估指标,具体的计算公式如下:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (10)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (11)$$

$$F_1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12)$$

其中,TP 表示在阳性实例中,模型预测为阳性的实例数量;FP 表示在阴性实例中,模型预测为阳性的实例数量;FN 表示在阳性实例中,模型预测为阴性的实例数量。

#### 3.3 参数设置

模型训练采用了 Adam 优化器,初始学习率设置为  $3e-4$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,批量大小设置为 128,训练轮次设置为 100,超参数  $d$  设置为 32,Dropout 比例参数设置为 0.5。在以上超参数设置条件下,本文提出的术后风险预测模型达到了收敛。

本文为验证所提出模型在术后风险预测任务上的有效性,在实验中选择了两种常用的统计机器学习模型 LR 和 XGBoost 以及两种最新的基于表格数据分类的神经网络 Wide&Deep<sup>[22]</sup> 和 Tabtransformer<sup>[23]</sup> 作为对比模型。LR 和 XGBoost 采用 scikit-learn 框架<sup>[24]</sup> 实现,Wide&Deep 和 Tabtransformer 采用开源的代码库<sup>①</sup> 实现。

#### 3.4 实验分析

首先,在三项术后风险预测任务上对比了模型的预测性能,实验结果如表 1 所示。

从表 1 所列结果可以观察到,Wide&Deep 和 Tabtransformer 在三项术后风险的预测任务上均优于 LR 和 XGBoost,特别是在阳性率较低的心血管不良风险预测任务上,Wide&Deep 和 Tabtransformer 的表现远优于 LR 和 XGBoost。该结果说明,神经网络在术后风险预测任务上的性能优于统计机器学习模型,这与文献[6]和文献[7]报告的结果保持一致。

① [https://github.com/jrzaurin/pytorch-widedeep/tree/pytorch\\_widedeep](https://github.com/jrzaurin/pytorch-widedeep/tree/pytorch_widedeep)

表1 实验整体结果

(单位: %)

模型	肺部并发症风险			ICU 入室风险			心血管不良风险		
	Precision	Recall	$F_1$	Precision	Recall	$F_1$	Precision	Recall	$F_1$
LR	71.508	31.068	43.316	72.840	36.875	48.963	45.833	13.415	20.755
XGBoost	69.965	48.293	57.143	74.118	39.375	51.429	41.176	8.537	14.141
Wide&-Deep	73.214	54.089	62.215	74.603	37.600	50.000	52.778	29.788	37.255
Tabtransformer	68.563	60.422	64.236	75.385	39.200	51.579	55.556	30.303	39.216
Our	68.378	65.723	<b>66.909</b>	65.088	57.664	<b>60.833</b>	77.395	44.260	<b>55.888</b>

此外,从表1中所列结果还可以看出,通过引入术前诊断文本数据表征,本文提出的模型在肺部并发症、心血管不良和ICU入室三个术后风险预测任务上均取得了最优的性能, $F_1$ 分别达到了66.909%、55.888%和60.833%。该结果证明,本文提出的文本数据表征增强的术后风险预测模型是有效的。

进一步观察表1中的结果发现,相比于其他模型,本文提出的模型是在保持了良好的精确率的条件下,大幅地提升了召回率,从而提升了 $F_1$ 。该结果说明,当模型引入非结构化的术前诊断数据表征后,进一步丰富了特征的医学语义信息,对阳性实例的预测带来了额外的医学语义信息补充,从而帮助

模型将之前无法判断的阳性实例准确地预测为阳性,进而提高了模型的召回率。

### 3.5 消融实验

为进一步验证文本数据表征对模型预测效果增强的作用,并探究文本中粗粒度语义信息和细粒度语义信息对预测任务的影响,本文还设计了不加入文本以及分别加入粗粒度和细粒度语义信息的对比消融实验,结果如表2所示,其中,“-E”表示模型中去除细粒度语义向量表征,“-S”表示模型中去除粗粒度语义向量表征,“-E-S”表示模型中去除所有的文本数据。

表2 消融实验结果

(单位: %)

模型	肺部并发症风险			ICU 入室风险			心血管不良风险		
	Precision	Recall	$F_1$	Precision	Recall	$F_1$	Precision	Recall	$F_1$
Our-E-S	60.842	59.444	60.031	53.878	52.598	53.192	52.854	42.687	46.347
Our-E	62.54	62.394	62.754	56.365	53.670	54.913	53.551	42.663	46.905
Our-S	68.089	66.010	66.883	61.129	58.152	59.570	79.697	43.029	55.577
Our	68.378	65.723	<b>66.909</b>	65.088	57.664	<b>60.833</b>	77.395	44.260	<b>55.888</b>

观察表2发现,在加入文本数据后,本文提出的模型在肺部并发症风险预测任务上的性能提高了6.878%,在ICU入室风险预测中提高了7.641%,在心血管不良风险预测中提高了9.541%,并且无论是单独加入粗粒度文本的语义向量还是细粒度文本的语义向量,模型的预测性能均得到明显改善。该结果说明,非结构化术前诊断中的信息对术后风险预测具有积极的作用,为术后风险预测提供了额外的决策信息,有效地增强了模型的预测能力。

此外,观察表2还可以发现,阳性率越低的术后风险,通过引入非结构化术前诊断数据表征后,模型的预测性能提升越高。该结果说明,对于阳性实例

更少的术后风险,模型需要更多的特征才能更准确地预测阳性病例,引入非结构化的术前诊断能够为模型带来更丰富的医学语义信息,从而使本文提出的模型在阳性率越低的术后风险预测中表现得越出色。

从表2结果还能够看出,相比于全局的粗粒度语义向量表征的缺失,模型对于局部的细粒度语义向量表征的缺失更加敏感。该结果说明,在术后风险预测的过程中引入围术期医学领域知识,对模型的预测性能提升具有重要的作用,这也进一步说明了本文提出的非结构化数据表征增强的术后风险预测模型的有效性和应用价值。

更进一步地,从表 2 中还可以看出,当模型同时引入粗粒度语义向量表征和细粒度语义向量表征时,模型的预测性能达到最优。该结果说明,当用非结构化数据表征增强术后风险预测模型时,既需要引入粗粒度语义向量表征携带的全局语义信息,又需要引入细粒度语义向量表征携带的局部语义信息。

### 3.6 细节分析

本文提出的模型通过自注意力机制为术后风险预测模型带来了可解释性。为验证和说明该效果,本文选取了一个发生了术后心血管不良的病人的案例,观察模型的注意力权重矩阵。该实例的术前诊断是“右肺上叶结节,高血压 3 级”。本文提出的模型准确地预测该实例的术后心血管不良风险结局。提取模型的注意力权重矩阵  $W_{weight}$ ,并画出其热力图(如图 6 所示)。在图 6 中,横轴上的“右肺上叶结节”和“高血压 3 级”是术前诊断中的实体病症,[PAD]是补全的字符,其余行的描述以及列的描述均是表格数据包含的变量。

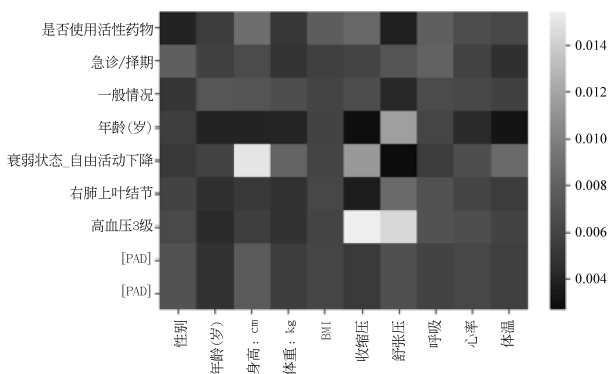


图 6 心血管不良实例的注意力权重热力图

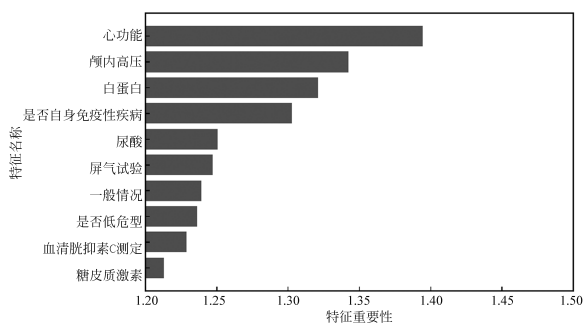
从图 6 可以看出,在术前诊断描述中,“高血压 3 级”显著地与表格数据中的收缩压和舒张压变量具有强关联。该强关联预示着模型通过训练,学习到了数据集中包含的医学领域知识关联信息,该关

联信息保存在了  $W_{weight}$  中,在术后风险预测中起到了重要的预示作用。另一个方面,该结果还说明,利用自注意力机制为术后风险预测模型带来了可解释性。总体地,实验结果验证了本文提出的模型在增强术后风险预测性能方面的鲁棒性和可解释性。

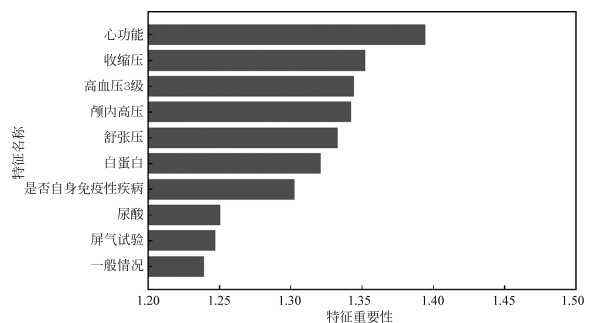
通过对比引入非结构化数据表征前后, $W_{weight}$  中包含的权重值按列求和后得到的每个变量在术后风险预测中的权重比率排序,进一步观察在术后风险预测中起重要作用的变量与术后风险结局是否存在医学语义的相关性,对比结果如图 7 所示。

从图 7 可以看出,权重比率越高,说明变量在预测中具有更高的重要性。从结果可以看出,在引入非结构化数据表征后,与术后心血管不良风险强相关的医学变量收缩压与舒张压的权重比率排序更加靠前。该结果进一步证实了本文提出的模型在提升术后风险预测性能的同时,还学习到了医学领域知识,具有更好的可解释性。

从图 7 中还可以看到,在引入非结构化数据表征后,临床医生根据经验总结或推断得到的额外的重要医学语义信息,也在风险预测中起到了重要的作用,例如临床医生基于收缩压和舒张压总结并记录下的“高血压 3 级”疾病。一方面,该结果证明了本文提出的模型学习到了医学领域知识,并对术后风险预测起到了积极的作用。另一方面,该结果还说明,本文通过直觉观察提出的模型是正确的,术前诊断中包含了大量的医学语义信息,这些信息既包含表格数据中已有的医学语义信息,还包含大量可用于丰富原始表格数据的额外的医学领域知识,这些信息会对模型的预测性能提升起到积极的作用。更进一步地,该结果也说明,本文提出的模型在提升了术后风险预测性能的同时,还具有良好的鲁棒性和结果可解释性。



(a) 未引入非结构化数据表征



(b) 引入非结构化数据表征

图 7 变量在模型术后风险预测中的重要性排序



## 4 结束语

术后风险预测在临床医学中具有重要意义,基于表格数据构建统计机器学习模型和深度神经网络,实现术后风险预测是常见的方式。非结构化术前诊断数据中蕴含了大量额外的医学领域知识,可为术后风险预测提供丰富的医学语义信息,然而它们尚未被有效利用。针对该问题,本文提出了一种新的模型,用非结构化数据表征增强术后风险预测,并在模型中引入自注意力机制,在有效融合表格数据和非结构化数据的同时,为模型带来良好的可解释性。实验结果表明,本文提出的非结构化数据表征增强的术后风险预测模型的性能显著高于其他比较的基线模型和先进模型。通过消融实验,验证了在术后风险预测中引入非结构化术前诊断数据的重要性,证明了本文提出的模型的有效性。此外,通过对模型的注意力权重的细节分析发现,利用自注意力机制将表格数据与非结构化的术前诊断融合用于术后风险预测,为模型带来了良好的可解释性。

## 参考文献

- [1] 魏娟, 邓惠民, 吕欣. 术后肺部并发症围手术期风险因素及防治策略[J]. 同济大学学报(医学版), 2021, 42(6): 736-743.
- [2] LUNDBERG S, NAIR B, VAVILALA M, et al. Explainable machine-learning predictions for the prevention of hypoxaemia during surgery[J]. Nature Biomedical Engineering, 2018, 2(10): 749-760.
- [3] CHIEW C, LIU N, WONG T, et al. Utilizing machine learning methods for preoperative prediction of postsurgical mortality and intensive care unit admission[J]. Annals of surgery, 2020, 272(6): 1133-1139.
- [4] LI P, LUO Y, YU X, et al. Readiness of healthcare providers for e-hospitals: a cross-sectional analysis in China before COVID-19[J]. BMJ Open, 2022, 12(2): e054169.
- [5] XUE B, LI D, LU C, et al. Use of machine learning to develop and evaluate models using preoperative and intraoperative data to identify risks of postoperative complications[J]. JAMA Network Open, 2021, 4(3): e212240-e212240.
- [6] BONDE A, VARADARAJAN K, BONDE N, et al. Assessing the utility of deep neural networks in predicting postoperative surgical complications: A retrospective study[J]. The Lancet Digital Health, 2021, 3(8): e471-e485.
- [7] FRITZ B, CUI Z, ZHANG M, et al. Deep-learning model for predicting 30-day postoperative mortality[J]. British Journal of Anaesthesia, 2019, 123(5): 688-695.
- [8] BARBIERI S, KEMP J, PEREZ-CONCHA O, et al. Benchmarking deep learning architectures for predicting readmission to the ICU and describing patients-at-risk[J]. Scientific Reports, 2020, 10(1): 1111.
- [9] CANET J, GALLART L, GOMAR C, et al. Prediction of postoperative pulmonary complications in a population-based surgical cohort[J]. Anesthesiology, 2010, 113(6): 1338-1350.
- [10] HILL B, BROWN R, GABEL E, et al. An automated machine learning-based model predicts postoperative mortality using readily-extractable preoperative electronic health record data[J]. British Journal of Anaesthesia, 2019, 123(6): 877-886.
- [11] ARIK S, PFISTER T. Tabnet: attentive interpretable tabular learning[C]//Proceedings of the AAAI Conference on Artificial Intelligence, 2021, 35(8): 6679-6687.
- [12] ZHANG D, YIN C, ZENG J, et al. Combining structured and unstructured data for predictive models: a deep learning approach[J]. BMC Medical Informatics and Decision Making, 2020, 20(1): 280.
- [13] LE Q, MIKOLOV T. Distributed representations of sentences and documents [C]//Proceedings of the 31st International Conference on Machine Learning, 2014: 1188-1196.
- [14] JOHNSON A, POLLARD T, SHEN L, et al. MIMIC-III, a freely accessible critical care database[J]. Scientific Data, 2016, 3: 160035.
- [15] BAHDANAU D, CHO K, BENGIO Y. Neural machine translation by jointly learning to align and translate[C]//Proceedings of the 3rd International Conference on Learning Representations, 2015.
- [16] HAO Y, DONG L, WEI F, et al. Self-attention attribution: interpreting information interactions inside transformer[C]//Proceedings of the AAAI Conference on Artificial Intelligence, 2021, 35(14): 12963-12971.
- [17] LOH W-Y. Classification and regression trees[J]. WIREs Data Mining and Knowledge Discovery, 2011, 1(1): 14-23.
- [18] GUO C, BERKHAHN F. Entity embeddings of categorical variables[J/OL]. arXiv preprint arXiv: 1604.06737, 2016.
- [19] DAI Z, WANG X, NI P, et al. Named entity recognition using BERT BiLSTM CRF for Chinese electronic health records[C]//Proceedings of the 12th In-

- ternational Congress on Image and Signal Processing, BioMedical Engineering and Informatics, 2019: 1-5.
- [20] HE K, ZHANG X, REN S, et al. Deep residual learning for image recognition [C]//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016: 770-778.
- [21] BA JL, KIROS JR, HINTON G. Layer normalization[J/OL]. arXiv preprint arXiv: 1607.06450, 2016.
- [22] CHENG H-T, KOC L, HARMSEN J, et al. Wide&deep learning for recommender systems[C]//Proceedings of the 1st Workshop on Deep Learning for Recommender Systems, 2016: 7-10.
- [23] HUANG X, KHETAN A, CVITKOVIC M, et al. Tabtransformer: Tabular data modeling using contextual embeddings [J/OL]. arXiv preprint arXiv: 2012.06678, 2020.
- [24] PEDREGOSA F, VAROQUAUX G, GRAMFORT A, et al. Scikit-learn: Machine learning in python[J]. Journal of Machine Learning Research, 2016, 12: 2825-2830.



王亚强(1984—), 博士, 副教授, 主要研究领域为机器学习、自然语言处理、医学信息学。  
E-mail: yaqwang@cuit.edu.cn



杨潇(1998—), 硕士, 主要研究领域为机器学习、自然语言处理、医学信息学。



朱涛(1969—), 通信作者, 博士, 教授, 主要研究领域为麻醉学、智慧麻醉、医学信息学。  
E-mail: 739501155@qq.com