

融合医学本体知识的药物推荐算法

洪高枫, 黄 杰, 万 健

(浙江科技学院 信息与工程学院, 杭州 310023)

摘 要: 针对传统的药物推荐技术忽略了医疗编码体系中蕴含的医学本体分类信息, 导致推荐效率不高的问题, 提出了一种融合图神经网络技术和注意力机制的药物推荐算法。首先通过图神经网络学习医学本体图中的分类关系, 对医疗代码进行嵌入表示; 其次将医疗表示输入结合注意力机制的循环神经网络中捕捉患者病历信息特征, 同时引入药物相互作用知识, 通过图神经网络学习药物相互作用关系; 最后进行多标签分类来输出推荐药物。算法在电子病历数据集上进行试验, 结果如下: F_1 值为 63.09%, 杰卡德相似系数为 47.43%, 精确度调用曲线值为 71.64%, 优于对比的其他方法。试验结果表明本算法能够丰富医疗代码表示, 提升药物推荐的准确率, 并且降低推荐药物组合中相互作用概率。本算法能够向医生推荐合适的药物, 对医疗智能化发展具有一定的参考价值。

关键词: 药物推荐; 图神经网络; 注意力机制

中图分类号: TP301.6

文献标志码: A

文章编号: 1671-8798(2022)03-0233-09

Drug recommendation algorithm integrating medical ontology representation

HONG Gaofeng, HUANG Jie, WAN Jian

(School of Information and Electronic Engineering, Zhejiang University of
Science and Technology, Hangzhou 310012, Zhejiang, China)

Abstract: In response to the low efficiency of recommendation arising from the problem that the traditional drug recommendation technology tends to ignore the medical ontology information contained in the medical code, a new drug recommendation algorithm was proposed by integrating the graph neural network with the attention mechanism. Firstly, the algorithm applied the graph neural network technology to learn the classification relationship in the medical ontology and conducted the embedded representation of medical code. Secondly, the patient's medical record information characteristics were captured through recurrent neural

收稿日期: 2021-06-23

基金项目: 国家自然科学基金项目(61972358); 浙江省重点研发计划项目(2020C03071)

通信作者: 万 健(1969—), 男, 福建省泉州人, 教授, 博士, 主要从事云计算及大数据研究。E-mail: wanjian@zust.edu.cn。

networks in combination with the attention mechanism. In addition, the knowledge of drug interaction was introduced to learn the drug interaction relationship through the graph neural network. Finally, the recommended drugs were output through the multi-label classification. The algorithm was tested on an EMR data set, whose experimental results show that the F_1 value is 63.09%, the Jaccard similarity coefficient is 47.43%, and the accurate call curve is 71.64%, which proves this method is superior to the contrast methods. The algorithm can enrich the representation of medical code, improve the accuracy of drug recommendations, and reduce the interaction rate in the recommended drug combination. The algorithm can recommend the appropriate drugs for doctors, which is of reference value for the development of medical intelligence.

Keywords: drug recommendation; graph neural network; attention mechanism

随着现代医学技术的发展,电子病历广泛使用,从而积累了大量的临床数据,如生命体征^[1]、疾病诊断^[2]、处方药物^[3]、医疗费用^[4]等数据。同时,深度学习技术为医疗数据的挖掘和利用提供了新的技术手段,是目前的一个热点研究方向^[5-6]。其中,基于电子病历的组合药物推荐算法,能够根据患者病情的变化特征、药物属性及大量药物之间的作用关系,辅助医生制定安全有效的处方^[7],具有重要的研究价值。

早期的药物推荐系统多基于规则。李枫林等^[8]基于描述逻辑与语义推理向用户推荐合适的抗高血压药物,实现了简单的药物推荐算法。Chen 等^[9]从患者的诊断、疾病分类、症状、检测结果等医学信息中得出用药规则以推荐药物。深度学习运用到药物推荐系统之后,Gong 等^[10]将患者的体征、医学诊断、既往用药等信息嵌入一个低维空间,并使用该嵌入表示进行推荐。

注意力机制是在深度学习模型中嵌入的一种特殊结构,用来自动学习和计算输入数据对输出数据的贡献大小。将注意力机制运用到药物推荐算法中能够捕捉患者病历中对当前推荐药物有贡献的内容进行学习。Choi 等^[11]使用注意力机制学习患者历史就诊的权重,对每次就诊之间的关系进行建模,获得以往就诊对当前药物推荐的贡献矩阵,帮助医生理解模型提供的推荐依据。Ma 等^[12]使用注意力机制来解决不同就诊时间间隔的临床检测数据对患者诊断的影响,从而提升模型的准确率。

图神经网络能够处理医疗实体间的复杂关系,它在医疗领域的使用方法就是在医疗图数据上进行卷积,包含基于谱方法和基于空间的方法。谱方法是在谱域上进行卷积,首先将图数据上的信号转换到谱域,然后在谱域上进行卷积的定义,再变换到空间域,代表模型为图卷积神经网络(graph convolutional networks, GCN)。空间方法是直接在空间上进行卷积,通过将相邻节点的信息聚合来更新结点的表示,代表模型为图注意力网络(graph attention networks, GAT)。Shang 等^[13]结合动态记忆网络和循环神经网络来模拟患者电子病历中的时序依赖,使用图卷积网络对药物之间的相互作用进行建模,降低推荐药物间的相互作用率,提升了推荐药物的安全性。Bhoi 等^[14]在推荐系统中使用图注意力网络对药物间相互作用进行建模,给不同模块量化得分,使得推荐药物具有可解释性。

以往研究未能充分利用药物相互作用知识,忽略医疗本体关系,导致医疗表示稀疏,推荐质量不够高。针对这些问题,本研究提出一种基于图神经网络结合注意力技术的药物推荐算法(graph augmented neural network with attention for medication recommendation, GRAD)。利用该算法构造医疗代码的本体关系图,使用图神经网络学习医疗代码的本体语义信息,丰富医疗代码的表示;将医疗表示输入结合注意力机制的循环神经网络中捕捉患者历史病历特征;同时引入药物相互作用知识,学习药物相互作用知识;并在公开的电子病历数据集上进行多个试验并评估。

1 药物推荐定义

1.1 患者病历定义

患者的电子病历可以表示为 t 次就诊的序列,如图1所示。一位患者的病历表示为 $\mathbf{P}=(x_1, x_2, \dots, x_t)$, t 为该患者就诊次数。患者的第 i 次就诊表示为 $x_i=(\mathbf{d}_i, \mathbf{p}_i, \mathbf{m}_i)$, $i=1, 2, \dots, t$, 其中 \mathbf{d}_i 为患者第 i 次就诊的国际疾病分类(international classification of diseases, icd-9)^[15] 诊断代码,如{'4373', '43820', 'V452'}; \mathbf{p}_i 为一个患者病历中第 i 次就诊的 icd-9 手术代码,如{'3731', '8872', '893'}; \mathbf{m}_i 为第 i 次就诊的药物解剖学、治疗学及化学分类法(anatomical therapeutic chemical, ATC)药物代码,如{'A12C', 'A03B', 'C01C'}。病历中出现的所有诊断代码、手术代码、药物代码的数量分别用 N_d, N_p, N_m 来表示。

	<div>电子病历</div> <div></div>	<div>电子病历</div> <div></div>	<div>电子病历</div> <div></div>		
诊断	4373,43820,V452	442.84,578.1,V12.5	442.84,559.0,285.1		
诊断	3731,8872,893	4444,3893,9904	444,8847,9904		
诊断	A12C,A03C,C01C	A06C,A07A,N02A	A06A,N02B,A12B		
	1	...	t-1	t	时间

图1 电子病历的表示

Fig. 1 Graphical illustration of EMR

1.2 推荐任务定义

以电子病历数据为基础,给定一名患者过去病历记录中的手术代码 $\mathbf{P}_t, \mathbf{P}_t=(p_1, p_2, \dots, p_t)$, 诊断代码 $\mathbf{D}_t, \mathbf{D}_t=(d_1, d_2, \dots, d_t)$ 和历史用药代码 $\mathbf{M}_{t-1}, \mathbf{M}_{t-1}=(m_1, m_2, \dots, m_{t-1})$ 。引入药物知识和医疗代码的 本体信息,将药物推荐当作序列预测任务,通过生成多标签输出 $\hat{\mathbf{y}}_t \in \{0, 1\}^{N_m}$ 来进行药物推荐。

2 GRAD 算法

2.1 算法框架描述

GRAD 算法框架如图2所示。它主要由图嵌入表示模块、患者病历表示模块、历史用药检索模块、药物相互作用模块和输出模块组成。图嵌入表示模块使用图神经网络对手术、诊断、药物本体图进行学习,对医疗代码使用向量表示并输出。患者病历表示模块对患者的历史诊断和手术代码向量使用注意力机制结合循环神经网络进行学习,输出患者历史病历向量 \mathbf{q}_t 。历史用药检索模块使用患者病历向量 \mathbf{q}_t 与历史用药向量相结合,得到历史用药向量 \mathbf{v}_t 。药物相互作用模块使用图卷积神经网络学习药物相互作用知识,与患者病历向量 \mathbf{q}_t 相结合,得到包含药物相互作用信息的药物向量 \mathbf{d}_t 。输出模块对 $\mathbf{q}_t, \mathbf{d}_t, \mathbf{v}_t$ 3 个向量进行拼接,通过多分类输出得到推荐的药物向量。

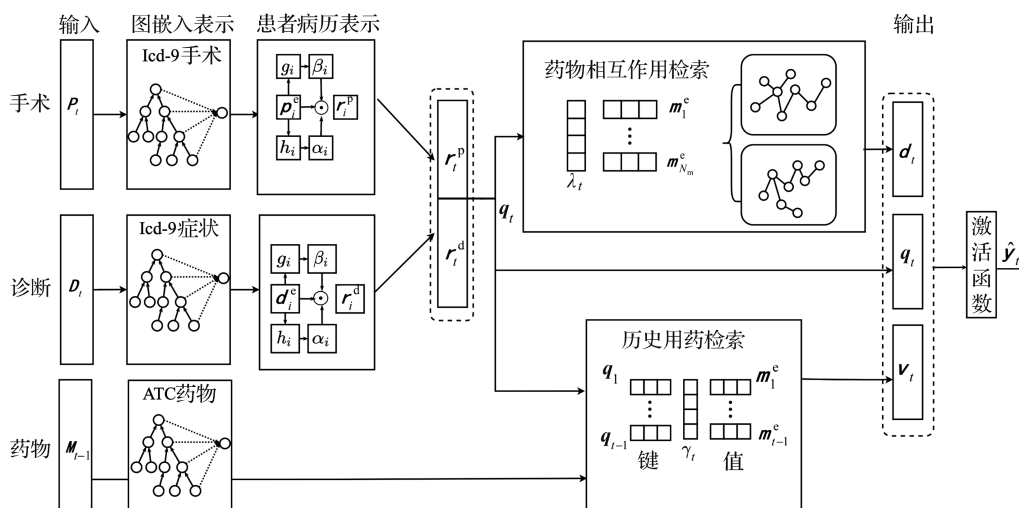


图2 GRAD 算法框架

Fig. 2 Framework of GRAD algorithm

2.2 医疗代码的本体图嵌入表示

医疗代码的编码体系结构可以描述为一个树状的结构, icd-9 编码部分结构表示如图 3 所示。叶子结点代表 icd-9 诊断代码, 其祖先结点代表具有医学分类概念的医学本体。如 C_9 413.0(卧位心绞痛)和 C_{10} 413.1(Prinzmetal 氏心绞痛)同属于 413 下的心绞痛分类, 属于 410~414 缺血性心脏病总分类。通过电子病历中出现的医疗代码, 根据编码知识来构造相应的医学本体关系图。使用图神经网络对医学本体树的各结点特征进行学习, 使得本体树中的结点特征包含不同关联度的其他结点信息, 获得更加全面的表示。最终得到融合医学本体语义信息的医疗代码表示, 以用于后续的推荐任务。

使用经过设计的图注意力网络实现消息在 icd-9 本体图中的传递(图 3)。对所有结点初始化, 对图中每个结点用一个初始化的向量 e_* 表示, $e_* \in \mathbb{R}^{64 \times 1}$ 。第一步将所有的非叶子结点 $e_i \in \bar{C}_*$ 表示为其自身嵌入表示和其所有子结点的嵌入表示之和, 采用图注意力机制的方式进行结合:

$$e_i = \sum_{k=1}^K \sigma \left(\sum_{j \in c(i) \cup \{i\}} \alpha_{ij}^k e_j \right). \quad (1)$$

式(1)中: σ 为一个非线性激活函数(这里使用的是 ReLU 函数); $c(i)$ 为结点 i 的所有孩子结点集合; α_{ij}^k 为第 k 个注意力下各个结点的系数,

$$\alpha_{ij}^k = \frac{\exp(\sigma_{\text{LeakyReLU}}(a^T [e_i \| e_j]))}{\sum_{l \in c(i) \cup \{i\}} \exp(\sigma_{\text{LeakyReLU}}(a^T [e_i \| e_l]))}. \quad (2)$$

式(2)中: a 为可学习的矩阵向量, $a \in \mathbb{R}^{2 \times 64}$; $\sigma_{\text{LeakyReLU}}$ 为 LeakyReLU 激活函数; $\|$ 为矩阵的拼接。

第二步对所有叶子结点 $e_i \in C_*$ 进行嵌入学习, 利用图注意力网络将其所有的祖先结点的嵌入表示和自身结点的嵌入表示相结合,

$$e_i = \sum_{k=1}^K \sigma \left(\sum_{j \in a(i) \cup \{i\}} \alpha_{ij}^k e_j \right). \quad (3)$$

第 k 个注意力下各个结点的系数

$$\alpha_{ij}^k = \frac{\exp(\sigma_{\text{LeakyReLU}}(a^T [e_i \| e_j]))}{\sum_{l \in a(i) \cup \{i\}} \exp(\sigma_{\text{LeakyReLU}}(a^T [e_i \| e_l]))}. \quad (4)$$

式(4)中: $a(i)$ 为叶子结点 i 的所有祖先结点; $\sigma_{\text{LeakyReLU}}$ 为 LeakyReLU 激活函数。图 3 展示了叶子节点 C_{10} 通过图注意力网络结合其所属分类本体信息的过程。最终得到了融合本体信息的所有诊断代码表示。

对电子病历中的诊断代码、手术代码和药物代码分别构造 icd-9 疾病、icd-9 手术和 ATC 药物的医学本体图, 再通过图神经网络学习诊断代码、手术代码和药物代码的结构信息。输入患者的医疗代码, 得到其医疗代码对应的嵌入表示。于是一位患者的第 i 次就诊可以表示为 (d_i^e, p_i^e, m_i^e) 。

2.3 病历学习及药物相互作用检索

2.3.1 历史诊断和手术的学习

采用结合注意力机制的循环神经网络学习患者历史病历, 将患者病历中的诊断、手术表示视作序列, 输入循环神经网络。使用注意力机制将历史诊断的隐藏层输出信息结合到当前诊断表示中。将历史诊断信息注意力权重表示为 $\alpha^d, \alpha^d[i] (1 \leq i \leq t)$ 为第 i 次就诊权重的量值,

$$\begin{aligned} h_1, h_2, \dots, h_t &= G_a(d_1^e, d_2^e, \dots, d_t^e); \\ \alpha^d &= \text{softmax}(F_a(h_1), F_a(h_2), \dots, F_a(h_t)). \end{aligned} \quad (5)$$

式(5)中: G_a 为一个方法, 这个方法是将诊断代码的时序向量输入循环神经网络(gate recurrent unit,

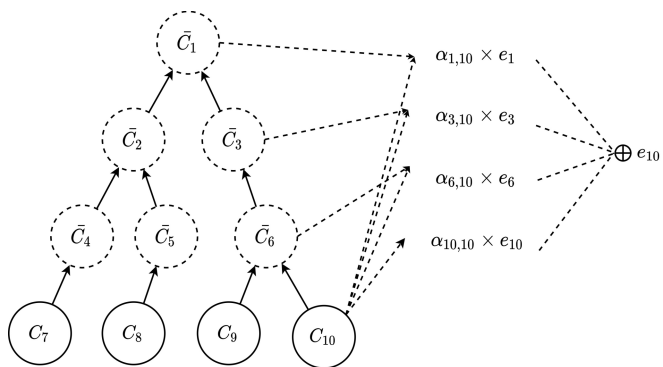


图 3 icd-9 编码部分结构表示

Fig. 3 Graphical illustration of icd-9 coding part

GRU)并输出; h_t 为GRU中的一个隐藏层输出; F_a 为权重向量的 $\mathbb{R}^{64 \times 1}$ 线性变换函数。同样将诊断代码的时序向量输入循环神经网络 G_β 中,并用 \tanh 作为激活函数获得权重 β^d ,

$$\begin{aligned} \mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_t &= G_\beta(\mathbf{d}_1^e, \mathbf{d}_2^e, \dots, \mathbf{d}_t^e); \\ \beta^d[j] &= \tanh(F_\beta(\mathbf{g}_j)), j = 1, 2, \dots, t. \end{aligned} \quad (6)$$

式(6)中: \mathbf{g}_i 为GRU中的隐藏层输出。结合2个诊断权重得到包含历史诊断信息的表示

$$\mathbf{r}_t^d = \sum_{i=1}^t \alpha^d[i] \beta^d[i] \otimes \mathbf{d}_i^e. \quad (7)$$

式(7)中: \otimes 为逐元素相乘; $\mathbf{r}_t^d \in \mathbb{R}^{64 \times 1}$ 。

将手术表示序列输入结合注意力机制的循环神经网络中得到具有历史手术信息的手术表示 \mathbf{r}_t^p ,连接2个向量通过线性变化输出得到患者历史病历的表示 $\mathbf{q}_t, \mathbf{q}_t \in \mathbb{R}^{64 \times 1}$,

$$\mathbf{q}_t = f([\mathbf{r}_t^p, \mathbf{r}_t^d]). \quad (8)$$

2.3.2 历史用药检索

采用记忆神经网络存储历史用药信息。历史用药以键值对的形式存储,键为结合患者过去诊断和手术的表示向量 \mathbf{q}_i ,值是对应使用药物的表示 \mathbf{m}_i^e 。键值对由一对元组 $\langle \mathbf{q}_i, \mathbf{m}_i^e \rangle (i \in [1, t-1])$ 组成。给定当前就诊结合历史病历的向量 \mathbf{q}_t ,计算之前就诊信息的权重

$$\gamma_t = \text{softmax}(\mathbf{q}_t \cdot \mathbf{q}_1, \mathbf{q}_t \cdot \mathbf{q}_2, \dots, \mathbf{q}_t \cdot \mathbf{q}_{t-1}). \quad (9)$$

通过注意力 $\gamma_t \in \mathbb{R}^{t-1}$ 得到历史用药向量 \mathbf{v}_t ,第 i 次历史用药向量:

$$\begin{cases} \mathbf{v}_t[i] = \gamma_t[i] \cdot \mathbf{m}_i^e; \\ \mathbf{v}_t = \sum_{i=1}^{t-1} \mathbf{v}_t[i]. \end{cases} \quad (10)$$

2.3.3 药物相互作用检索和输出

引入药物相互作用知识,用邻接矩阵 \mathbf{A}_C 表示电子病历中的药物共存关系,用邻接矩阵 \mathbf{A}_D 表示药物相互作用关系。采用图卷积神经网络学习药物共现关系和药物相互作用关系,以表示药物相互作用知识。用 $\mathbf{A}_* \in \mathbb{R}^{N_m \times N_m}$ 来统一表示邻接矩阵 \mathbf{A}_C 和 \mathbf{A}_D ,每个 \mathbf{A}_* 可以表示为

$$\mathbf{A}_* = \mathbf{D}^{-1}(\mathbf{A}_* + \mathbf{I})\mathbf{D}^{-1}. \quad (11)$$

式(11)中: \mathbf{A}_* 为一个对称归一化拉普拉斯矩阵; \mathbf{D} 为一个对角矩阵; \mathbf{I} 为一个单位矩阵。然后使用图卷积神经网络来学习每个图中药物之间的关系,将药物的相互作用和共现关系结合到嵌入表示中:

$$\begin{cases} \mathbf{Z}_C = \tilde{\mathbf{A}}_C \tanh(\tilde{\mathbf{A}}_C \mathbf{m}^e \mathbf{W}_C^1) \mathbf{W}_C^2; \\ \mathbf{Z}_D = \tilde{\mathbf{A}}_D \tanh(\tilde{\mathbf{A}}_D \mathbf{m}^e \mathbf{W}_D^1) \mathbf{W}_D^2. \end{cases} \quad (12)$$

式(12)中: $\mathbf{W}_C^1, \mathbf{W}_C^2, \mathbf{W}_D^1, \mathbf{W}_D^2 \in \mathbb{R}^{64 \times 64}$,为参数矩阵。与模型GAMENET^[13]不同,本文方法使用从药物本体图输出的药物表示矩阵作为输入。经过学习可以得到药物共现关系图的表示矩阵 $\mathbf{Z}_C \in \mathbb{R}^{64 \times N_m}$ 和药物相互作用图的表示矩阵 $\mathbf{Z}_D \in \mathbb{R}^{64 \times N_m}$ 。基于这些药物表示和查询向量 \mathbf{q}_t 计算注意力

$$\lambda_t = \text{softmax}((\mathbf{Z}_C + \omega \mathbf{Z}_D)^T \cdot \mathbf{q}_t). \quad (13)$$

式(13)中: $(\mathbf{Z}_C + \omega \mathbf{Z}_D)^T$ 为 $(\mathbf{Z}_C + \omega \mathbf{Z}_D)$ 的转置矩阵; $\omega \in \mathbb{R}$,为药物相互作用的权重。得到检索药物相互作用信息的向量

$$\mathbf{d}_t = (\mathbf{Z}_C + \omega \mathbf{Z}_D) \cdot \lambda_t. \quad (14)$$

将患者的历史诊断及手术表示 \mathbf{q}_t ,历史用药的向量 \mathbf{v}_t ,药物相互作用的检索结果 \mathbf{d}_t 相连接,使用非线性激活函数变换多分类输出,得到推荐的药物向量 $\hat{\mathbf{y}}_t, \hat{\mathbf{y}}_t \in \mathbb{R}^{N_m}$,当 $\hat{\mathbf{y}}_t[i] > 0.5$ 时则推荐药物 i 。

$$\hat{\mathbf{y}}_t = \sigma([\mathbf{q}_t, \mathbf{v}_t, \mathbf{d}_t]). \quad (15)$$

2.4 损失函数设置

将药物推荐任务视作序列预测的多标签分类问题,使用二元交叉熵损失函数、多标签分类间距损失

函数来控制推荐药物的准确率。定义药物相互作用损失函数来衡量推荐药物的相互作用率。

二元交叉熵损失函数定义如下:

$$L_{\text{bce}} = - \sum_t^T \sum_i^{N_m} \mathbf{y}_t^i \lg \sigma(\hat{\mathbf{y}}_t^i) + (1 - \mathbf{y}_t^i) \lg(1 - \sigma(\hat{\mathbf{y}}_t^i)). \quad (16)$$

式(16)中: T 为患者就诊的总次数; \mathbf{y}_t^i 和 $\hat{\mathbf{y}}_t^i$ 分别为第 t 次就诊的真实用药和算法推荐的药物。定义多标签间距损失函数

$$L_{\text{multi}} = \sum_t^T \sum_i^{N_m} \sum_j^U \frac{\max(0, 1 - (\hat{\mathbf{y}}_t[j] - \hat{\mathbf{y}}_t[i]))}{|U|}. \quad (17)$$

式(17)中: U 为真实病历使用药物的集合。药物与药物相互作用的损失函数

$$L_{\text{DDI}} = \sum_t^T \sum_{i,j} (\mathbf{A}_D \odot (\hat{\mathbf{y}}_t^T \hat{\mathbf{y}}_t)) [i, j]. \quad (18)$$

使用 3 个超参数 γ_1 、 γ_2 、 γ_3 来联合 3 个损失函数,

$$L_{\text{combine}} = \gamma_1 L_{\text{bce}} + \gamma_2 L_{\text{multi}} + \gamma_3 L_{\text{DDI}}. \quad (19)$$

式(19)中: $\gamma_1 + \gamma_2 + \gamma_3 = 1$ 。

3 试验设计与结果分析

3.1 数据集介绍

本试验使用重症监护医学信息数据库(Medical Information Mark for Intensive Care, MIMIC-III)^[16] 电子病历(Electronic Medical Record, EMR)数据,该数据集是由麻省理工学院计算生理学实验室发布的免费公开的重症监护室数据集。本试验使用该数据库中的诊断单、手术单和处方单数据,筛选患者在进入重症监护室(Intensive Care Unit, ICU)后的 24 h 内接受的药物,并将药物编码从国家药物代码(National Drug Code, NDC)编码转换为 ATC 三级编码,并使用 icd-9 编码实现诊断代码和手术代码的统一。数据集的统计见表 1。

3.2 评价指标及训练设置

为衡量推荐准确度,使用杰卡德相似系数(Jaccard)(真实药物与推荐药物交集的大小除以并集大小)、 F_1 值(精确率和召回率的调和平均)、精确度调用曲线面积(precision recall area under curve, PRAUC)作为准确率的衡量指标。同时为衡量推荐药物的安全性,定义药物相互作用率 R_{DDI} ,即推荐组合药物中含有药物相互作用(drug-drug interaction, DDI)药物的比例

$$R_{\text{DDI}} = \frac{\sum_k^N \sum_t^{T_k} \sum_{i,j} |(\mathbf{m}_i, \mathbf{m}_j) \in \hat{\mathbf{y}}_t^k | (\mathbf{m}_i, \mathbf{m}_j) \in \mathbf{A}_D |}{\sum_k^N \sum_t^{T_k} \sum_{i,j} 1}. \quad (20)$$

式(20)中: N 为测试集中患者的数量; T_k 为第 k 个患者的就诊次数。定义相对药物相互作用率 Δ_{DDI} ,即推荐药物的 DDI 与电子病历(EMR)数据集中真实 DDI 的相对百分比

$$\Delta_{\text{DDI}} = \frac{R_{\text{DDI}} - R_{\text{DDI,EMR}}}{R_{\text{DDI,EMR}}} \times 100\%. \quad (21)$$

试验将以 4 : 1 : 1 的比例将数据集随机分为训练集、验证集和测试集。超参数在验证集中进行调整,取值分别为 $\gamma_1 = 0.98$ 、 $\gamma_2 = 0.01$ 、 $\gamma_3 = 0.01$ 。GRU 和图神经网络中随机失活(dropout)设置为 0.4,图注意力网络中注意力个数 K 为 2。使用自适应估计优化器来优化模型的参数,学习率为 0.000 4,根据 40 个训练周期后验证集上的结果选择性能最佳的模型,试验的最终结果为测试集上的结果。所有方法均使用 PyTorch 1.7.0 软件实现,系统为 Ubuntu 18.04,显卡为具有 12 GB 内存的 NVIDIA GEFORCE

表1 数据集的统计

Table 1 Statistics of data set

属性	数量
患者数量/个	6 350
病历记录/条	15 016
诊断代码/条	1 958
手术代码/条	1 426
药物代码/条	145
患者平均用药/条	8.8
患者平均诊断/条	10.51
患者平均手术/条	3.84

RTX 2080Ti。

3.3 对比方法介绍

逻辑回归(logistic regression, LR)方法使用 L_2 正则化逻辑回归,将医疗代码用多热向量来表示并作为数据输入,使用二分类用于处理多标签输出。RETAIN 方法^[11]基于两层注意力网络模型学习的药物组合序列,选择过去就诊中重要的临床变量来进行推荐。LEAP 方法^[17]使用循环神经网络来学习药物间依赖关系,使用基于内容的注意力机制来捕获标签到实例间的映射。GAMENET 方法^[13]通过存储模块将历史用药和药物相互作用 DDI 使用图卷积网络集成来降低推荐药物间的相互作用率。PREMIER 方法^[14]使用注意力机制学习患者历史表示,结合图注意力机制学习药物相互作用来进行安全的药物推荐。

3.4 试验结果及分析

3.4.1 不同模型对比试验

不同模型对比试验见表2。试验结果表明,在所有方法中,本文方法可以达到最好的效果。本研究提出的方法在 Jaccard、PRAUC 和 F_1 方面比 PREMIER 分别高出 1.02%、1.09% 和 1.23%。同时本文方法兼顾药物相互作用,在取前 40 种药物相互作用情况下和同类深度学习方法对比,本文方法的 R_{DDI} 较低,为 0.070 5。此外本文方法推荐的平均药物数量为 14.98 种,在与各个深度学习方法对比中最接近病历真实的平均药物数量 14.68 种。这表明本文算法比其他算法更能有效提高药物推荐的准确率和安全性。LR 为机器学习的方法,在 F_1 、PRAUC、Jaccard 得分方面不如其他深度学习方法。深度学习方法中 LEAP 为基于实例的方法,在各项性能指标中低于 RETAIN、GAMENET 等基于患者历史病历时序序列的方法。这也证实了患者的历史就诊信息对当前药物推荐非常重要。GAMENET 和 PREMIER 方法引入了药物相互作用知识,相比 LR、LEAP 方法提升了药物推荐的准确率,降低了药物的相互作用率,但将这些方法将医疗代码看作独立个体,忽视了医疗代码间的医学知识,各项指标低于本文方法。这表明了结合医疗本体图知识后的算法对药物推荐的准确率和药物相互作用率的控制具有提升作用。

表2 不同模型对比试验

Table 2 Comparison experiment of different models

模型	$R_{DDI}40$	Jaccard	PRAUC	F_1	推荐药物数量/种
LR	0.078 2	0.408 7	0.673 9	0.566 9	11.37
LEAP ^[17]	0.063 3	0.391 1	0.569 9	0.548 6	15.96
RETAIN ^[11]	0.086 3	0.414 0	0.661 2	0.574 6	18.46
GAMENET ^[13]	0.078 0	0.447 9	0.688 3	0.605 3	13.93
PREMIER ^[14]	0.075 4	0.464 1	0.705 5	0.618 6	13.96
本文方法去掉 DDI 图	0.076 7	0.473 8	0.715 7	0.630 6	15.03
本文方法去掉本体图	0.071 6	0.469 1	0.709 8	0.622 0	14.34
本文方法	0.070 5	0.474 3	0.716 4	0.630 9	14.98

去掉 DDI 图,不进行药物相互作用检索的模型,在没有药物相互作用知识图的情况下, Jaccard、PRAUC、 F_1 值变动不大,但模型推荐药物中的 R_{DDI} 变高,达到 0.076 7,这说明本文方法中把药物间相互作用关系的知识和具有历史就诊信息的查询向量相结合,降低了推荐药物中的相互作用率,能够有效提高用药安全性。去掉本体图模块,直接用简单的多热编码嵌入矩阵代替医学本体图嵌入进行试验,在没有医疗代码本体结构的嵌入后, F_1 值下降 0.89%, PRAUC 下降 0.66%, 药物推荐的准确率出现明显下降,这表明本研究使用的图神经网络具有对高阶结构特征的编码能力,能够丰富医学本体的嵌入表示,在一定程度上弥补了训练数据稀疏的问题,提高了药物推荐的准确率。

3.4.2 不同就诊次数的试验

由于每位患者的就诊次数不同,故应考虑以往就诊次数对药物推荐质量的影响。不同就诊次数下 F_1 值对比如图4所示, LEAP 为基于实例的方法,就诊次数的增加对之后的药物推荐没有影响, F_1 值一直在较低的水平波动。GAMENET 为基于历史就诊时序的方法,能够学习患者之前的用药,相比 LEAP

方法有明显的提升,但随着就诊次数的增加, F_1 值开始下降。PREMIER 相比 GAMENET 增加了时序注意力机制,提升了 F_1 值,但随着就诊次数的增加, F_1 值仍然会下降。对于不同的时序长度,本研究使用的算法优于其他比较的方法。本文算法在以就诊次数为分类依据的所有试验中 F_1 值均为最高。特别是对于就诊次数多的患者,与其他方法相比仍能保持较高的准确率,这表明本文方法对患者病历中的长时序依赖具有更好的建模能力。

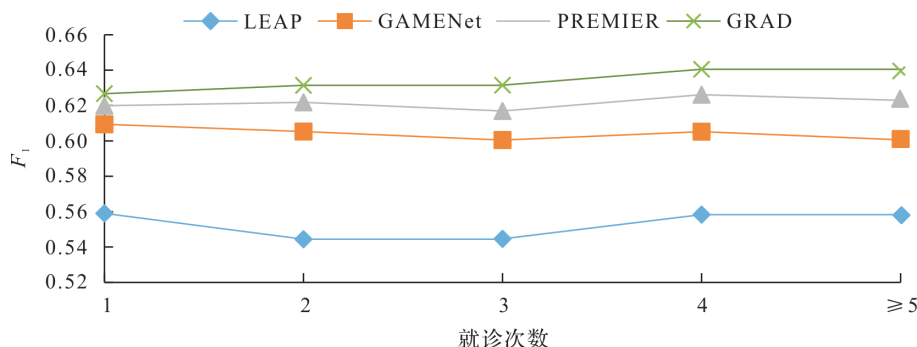


图 4 不同就诊次数下 F_1 值对比

Fig. 4 Comparison of F_1 value under different number of visits

3.4.3 不同 DDI 数量试验

关于算法对药物相互作用的效果,本研究还做了进一步的试验。分别使用前 40、60、80、100 种 DDI 类型,以探讨在使用不同数量的 DDI 情况下,本算法及各对比方法受到的影响。不同 DDI 数量试验结果见表 3,RETAIN 方法没有引入药物相互作用知识,使得 R_{DDI} 一直处于较高水平;GAMENET 和 PREMIER 引入了药物相互作用知识,在 DDI 数量较少时能够控制推荐药物的相互作用率,但随着 DDI 数量的增加推荐药物的相互作用率快速增加;LEAP 是基于实例的方法,推荐药物按照之前的处方进行推荐使得推荐药物的相互作用能控制在较低水平,但在 DDI 数量取前 100 时 Δ_{DDI} 仍会大于 0;本文方法优于对比方法,尽管 Δ_{DDI} 从 -18.48% 上升到 -0.26%,但考虑到 DDI 的数量从 40 变为 100,GRAD 是唯一能够实现 R_{DDI} 降低的算法,并且无论 DDI 类型有多少, Δ_{DDI} 始终大于零。这表明在引入药物相互作用知识后本算法能有效降低推荐药物的相互作用率,更具安全性。

表 3 不同 DDI 数量试验结果

Table 3 Results of experiment at different number of DDI

方法	DDI 数量为 40		DDI 数量为 60		DDI 数量为 80		DDI 数量为 100	
	R_{DDI}	$\Delta_{DDI}/\%$	R_{DDI}	$\Delta_{DDI}/\%$	R_{DDI}	$\Delta_{DDI}/\%$	R_{DDI}	$\Delta_{DDI}/\%$
GAMENET ^[13]	0.078 0	-9.81	0.198 5	-0.03	0.300 0	2.28	0.403 5	3.73
RETAIN ^[11]	0.086 3	-0.21	0.206 5	4.00	0.298 5	1.77	0.397 9	2.29
LEAP ^[17]	0.063 3	-26.81	0.192 3	-3.15	0.293 1	-0.08	0.396 4	1.90
PREMIER ^[14]	0.075 4	-12.82	0.198 9	0.17	0.309 5	5.52	0.420 1	7.99
GRAD	0.070 5	-18.48	0.186 9	-5.87	0.287 3	-2.05	0.388 0	-0.26

4 结 语

本研究提出一种基于图神经网络和注意力机制的药物推荐算法,将每个医学本体当作一个结点,采用图神经网络捕捉其中各个医学本体之间的关系,学习包含医学本体知识的高阶特征;同时使用注意力来对患者的历史病历进行更高效的建模,引入药物相互作用知识,提供兼顾用药安全的药物推荐。以 MIMIC 电子病历为数据源的综合试验结果表明,本算法能够提升推荐药物的准确率,并降低推荐药物间的相互作用率。尽管本研究提出的算法取得不错效果,但研究中仍存在问题,例如由于数据规模的限制在医学表示时仍存在稀疏性等问题,后期的研究将引入更多的知识以增强医学表示。

参考文献:

- [1] 李鹏飞,鲁法明,包云霞,等. 基于医疗过程挖掘与患者体征的药物推荐方法[J]. 计算机集成制造系统,2020,26(6):1668.
- [2] 张晓博,杨燕,李天瑞,等. 基于医疗文本数据聚类的帕金森病早期诊断预测[J]. 计算机应用,2020,40(10):3088.
- [3] 刘杰,金柳颀,景波. 基于药物和疾病特征关联的药物重定位混合推荐算法[J]. 计算机应用研究,2020,37(3):672.
- [4] 芮晨,李杰,郭栋伟,等. 基于 LIME-BP 神经网络的医疗费用预测研究[J]. 中国卫生统计,2020,37(5):698.
- [5] 王露潼,王红,宋永强,等. 基于 FT-LSTM 模型的临床事件诊断序列预测研究[J]. 计算机应用研究,2020,37(10):2961.
- [6] 吴宗友,白昆龙,杨林蕊,等. 电子病历文本挖掘研究综述[J]. 计算机研究与发展,2021,58(3):513.
- [7] 周虎,于跃,张正宇,等. 基于电子病历的 ADR 知识发现与应用模型研究[J]. 中国卫生事业管理,2020,37(2):81.
- [8] 李枫林,李娜. 基于情景的医药信息服务本体建模及规则推理研究[J]. 情报理论与实践,2016,39(5):120.
- [9] CHEN Z, MARPLE K, SALAZAR E, et al. A physician advisory system for chronic heart failure management based on knowledge patterns[J]. Theory & Practice of Logic Programming,2016,16(5):604.
- [10] GONG F, WANG M, WANG H, et al. Smr: medical knowledge graph embedding for safe medicine recommendation[J]. Big Data Research,2021,23:100174.
- [11] CHOI E, BAHADORI M T, KULAS J A, et al. RETAIN: an interpretable predictive model for healthcare using reverse time attention mechanism[C]//Proceedings of the 30th International Conference on Neural Information Processing Systems. Barcelona: CAI,2016:3512.
- [12] MA L, ZHANG C, WANG Y, et al. Concare: personalized clinical feature embedding via capturing the healthcare context[C]//Proceedings of the AAAI Conference on Artificial Intelligence. New York: AAAI,2020:833.
- [13] SHANGE J, XIAO C, MA T, et al. GAMENET: graph augmented memory networks for recommending medication combination[C]//Proceedings of the AAAI Conference on Artificial Intelligence. Hawaii: AAAI,2019:1126.
- [14] BHOI S, LI L M, HSU W. PREMIER: personalized recommendation for medical prescriptions from electronic records[EB/OL]. (2020-08-28)[2022-05-12]. <https://arxiv.org/abs/2008>.
- [15] SLEE V. The international classification of diseases: ninth revision (icd-9)[J]. Annals of Internal Medicine,1978,88(3):424.
- [16] JOHNSON A E W, POLLARD T J, SHEN L, et al. MIMIC-III, a freely accessible critical care database[J]. Scientific Data,2016,3(1):1.
- [17] ZHANG Y, CHEN R, TANG J, et al. LEAP: learning to prescribe effective and safe treatment combinations for multimorbidity[C]//Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM,2017:1315.