

基于Cluster-GCN的医疗保险欺诈检测研究

李淑锦^{1,2}, 梅浩²

¹广东培正学院数字经济研究所, 广东 广州

²杭州电子科技大学经济学院, 浙江 杭州

收稿日期: 2023年7月18日; 录用日期: 2023年7月31日; 发布日期: 2023年9月6日

摘要

医疗保险欺诈是现阶段医疗保险机构面临的重大挑战, 由于保险基金规模的扩张, 欺诈导致的基金损失也在增加。随着大数据技术的发展, 医疗保险欺诈检测技术在不断提高。本文针对个人的多条就诊记录与多个就诊医院, 使用基于图聚类(社群发现)改进GCN的Cluster-GCN方法来搭建被保险人与医院、被保险人之间的关系网络, 并设定行为信息更新与聚合的方式。检测结果表明, Cluster-GCN在医疗保险欺诈检测中有较好的分类性能, AUC达到0.86以上, 总体识别的准确精度达到83.37%, 反欺诈的识别率为77.25%。

关键词

保险欺诈, 图卷积网络(GCN), Cluster-GCN, 行为信息

Research on Detection of Medical Insurance Fraud Based on Cluster-GCN

Shujin Li^{1,2}, Hao Mei²

¹The Institute of Digital Economics, Guangdong Peizheng College, Guangzhou Guangdong

²School of Economics, Hangzhou Dianzi University, Hangzhou Zhejiang

Received: Jul. 18th, 2023; accepted: Jul. 31st, 2023; published: Sep. 6th, 2023

Abstract

Medical Insurance fraud is a major challenge faced by medical insurance institutions at this stage. Due to the expansion of the insurance fund scale, the fund loss caused by fraud is also increasing. With the development of Big data technology, medical Insurance fraud detection technology is constantly improving. This article focuses on multiple individual medical records and multiple hospital visits and uses the improved Cluster-GCN method based on graph clustering (community

discovery) to build a relationship network between insured individuals and hospitals, as well as between insured individuals. It also sets a way to update and aggregate behavioral information. The detection results show that Cluster-GCN has good classification performance in the detection of medical insurance fraud, with an AUC of more than 0.86, overall recognition accuracy of 83.37%, and anti-fraud accuracy of 77.25%.

Keywords

Insurance Fraud, Graph Convolutional Networks (GCN), Cluster-GCN, Behavioral Information

Copyright © 2023 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

随着基本医疗保险覆盖面进一步扩大, 中国的医疗保险基金支出也在不断增加。根据国家统计局的数据, 截至 2022 年年末, 全国城镇居民参加基本医疗保险人数达到 9.8 亿人, 医疗保险基金支出从 2016 年的 1.1 万亿元增长到 2022 年的 3 万亿元, 年均复合增速为 18%。医疗保险欺诈是指医疗机构或个人通过各种手段骗取医疗保险基金的行为。随着医疗保险基金的扩张, 医疗保险欺诈问题也日益突出。中国社会保障学会指出, 全球因欺诈导致的医疗保险基金损失占医疗保险基金支出的 4.57%。以此为依据, 可以计算出 2022 年全国医保基金因欺诈的损失要高达 1371 亿元, 医疗保险欺诈问题日益突出, 已经成为各国政府和医疗机构关注的焦点。

医疗保险欺诈是指个人或组织利用医疗保险制度中的漏洞、优惠政策或者非法手段, 骗取、盗取医疗保险基金或者其他医疗资源的行为。这种行为通常包括虚假申报、虚开发票、虚构病情、假冒身份等等, 旨在非法获得医疗保险待遇。医疗保险反欺诈则是对医疗保险欺诈进行有效的判断, 目的是识别和防止医疗保险欺诈行为, 保障医疗保险基金的安全和有效利用。目前医疗保险检测已然成为医疗保险领域的一项重要任务。

2. 文献综述

国内关于医疗保险反欺诈已有不少的研究成果。从研究内容来看, 保险反欺诈工作主要分为事前监管和事中监管, 事前监管是指保险机构在开展业务之前, 通过申请人提交的申请信息来判断是否给出保险以及保费如何定价; 事中监管是指保险机构在业务进行的过程中, 对被保险人的行为进行监管与规范, 从而识别其行为是否存在欺诈的事实。陈清华等(2023) [1]通过对 31 个省市共 111 起骗保数据进行分析, 采用 fsQCA 方法概括出参保人群妨碍、政治压力缺位、政策过程艰难和社会意识缺乏 4 种典型骗保模式。郑文珍等[2] [3]采用深圳市某家医院一个月的医保数据记录, 在对原始数据进行预处理后, 分别建立 RF 模型和 XGB 模型对医保正常数据和欺诈数据进行分类预测, 最终 F1 和 AUC 的结果显示 XGB 模型具有更好地评估识别能力。李杰等[4]认为, XGBoost 模型能够有效建立医保反欺诈系统, 识别出潜在的医保欺诈现象, 能够对医疗欺诈领域起到拨乱反正的作用。林源(2015) [5]采用 PCA 和 BP 神经网络模型对新农合医疗机构的住院服务的滥用行为进行识别, 研究显示 BP 神经网络模型能够有效检测出医疗欺诈行为, 模型结果显著优于 logistic 回归模型。周杰辉等(2021) [6]设计了医保反欺诈的可视化系统 Medicare Vis, 通过可视化对实例研究进行分析, 结果显示该方法能够有效检测出欺诈行为的关联性。曹鲁慧等(2020) [7]

采用 TLSTM 的方法来挖掘用户的就医时间序列信息, 判断用户存在就医欺诈的可能性。李金灿等(2021) [8]介绍了运用在医保欺诈检测领域的有监督学习方法: 神经网络和决策树, 以及无监督学习方法: 聚类分析、离群检测和关联法则挖掘, 并分析了不同模型的优势与不足。有监督学习方法需要较多的参数与迭代, 无监督学习法对异常值敏感, 并且处理海量数据的效率低下。易东义[9]和易东义等[10]提出了病人与医生关系网络, 采用 GCN 算法来挖掘欺诈信息, 并结合主动学习方法来标注医保欺诈标签。同时在保险欺诈检测中, 遇到的一个主要问题是欺诈数据的不平衡, 原因是申请信息并不包含参保人在参保后的行为信息, 通过行为信息(如消费信息)对参保人进行行为评分来实时识别欺诈的办法开始得到重视。吴文龙等(2021) [11]采用生成对抗网络将得到的仿真数据加入到医保数据中, 克服类别不平衡的影响。

事前监管的弊端不仅体现在无法监管参保人的事中行为, 还无法对出现的团伙欺诈和医疗机构违规行为进行监督。本文将从以下三个方面进行研究: 第一, 观测样本存在多条行为信息, 包含了在一个医院的多次就诊信息与在多个医院的就诊信息, 研究如何通过知识图谱和行为信息来搭建医院与被保险人、被保险人与被保险人之间的关系。第二, 根据医院的就诊项目与欺诈人群的就诊信息来识别易发生欺诈的处方、病情信息; 第三, 研究如何对就诊的行为信息进行更新与补充。在识别医疗保险欺诈的过程中可能的贡献有: 1) 是将医院作为连接关系网络的判断依据, 被保险人作为图的节点, 为两个去过同一家医院的被保险人建立联系, 增强反欺诈的识别能力; 2) 是通过 GCN 来学习自身节点和邻接节点的信息, 通过聚合-更新-聚合的方式来生成节点表征; 3) 是通过在损失函数中对不同标签数据附加权重的方式来平衡类别权重, 提升欺诈样本的影响, 并结合 Cluster-GCN 的算法来减少算法的复杂度。

3. 评估方法介绍

3.1. GCN 模型

图卷积网络(GCN)是图神经网络(Graph Neural Networks, GNN)的一种。GNN 是一种深度学习模型, 它用于处理图数据, 例如社交网络、化学分子和交通网络等。GNN 通过对节点特征进行聚合和信息传递, 从而实现对整个图的建模和分析。GNN 的基本思想是通过聚合每个节点周围的信息来更新节点的特征表示, 并利用这些更新的节点特征来预测节点的标签或执行其他任务。GNN 主要包括两个主要组成部分: 节点嵌入(node embedding)和图嵌入(graph embedding)。

节点嵌入主要是将每个节点的特征表示为低维向量, 同时考虑节点周围的邻居节点。节点嵌入可以通过各种各样的方式来计算, 例如 GCN (Graph Convolutional Network)、Graph SAGE、GAT (Graph Attention Network)等。图嵌入则是将整个图表示为一个低维向量, 可以用于图分类或图生成等任务。图嵌入可以通过池化(pooling)或图卷积等方式实现除了一些众所周知的英文缩写, 如 IP、CPU、FDA, 所有的英文缩写文中第一次出现时都应该给出其全称。文章标题中尽量避免使用生僻的英文缩写。

图 1 为节点信息更新的方式。以 A 点为例, A 点的邻接节点是 B、C、D, B 节点通过将自身节点的信息与聚合过来的 A、C 节点的信息进行变换, 得到新的 B 节点信息。同样的过程发生在 C 和 D 节点上。A 再将聚合得到的 B、C、D 节点信息与自身节点信息进行变换得到新的节点信息。多次更新后的节点信息就作为节点的特征表示。

在 GCN 模型中, 记一个图的节点个数为 N , 图的邻接矩阵为 A , $X^{(l)}$ 是第 l 层的节点表征。一个 L 层的图卷积神经网络由 L 个图卷积层组成, 每一层都通过聚合邻接节点的上一层的表征来生成中心节点的当前层的表征:

$$z^{(l+1)} = A'X^{(l)}W^{(l)} \quad (1)$$

$$X^{(l+1)} = \sigma(z^{(l+1)}) \quad (2)$$

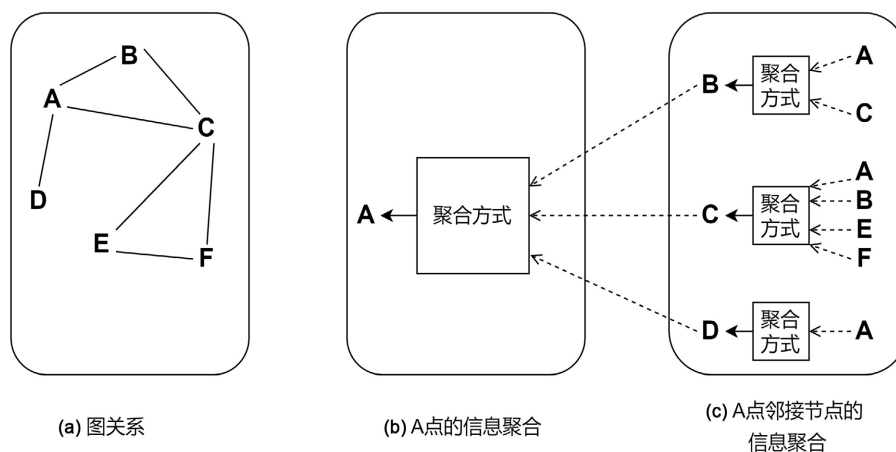


Figure 1. Mode of the node information updated
图 1. 节点信息更新的方式

公式(1)中, $X^{(l)}$ 表示 l 层的 N 个节点表征; $X^{(0)} = X$, X 是输入的初始节点属性; A' 是归一化和规范化后的邻接矩阵, 用于聚合其他节点的信息; $W^{(l)}$ 是权重矩阵, 包含了要训练的参数; 激活函数 $\sigma(\cdot)$ 设定为 ReLU, 用于传递当前层的节点信息。

训练的目标是通过最小化损失函数来学习公式(1)中的权重矩阵:

$$R = w_{y_k} \sum_{i \in y_k} \text{loss}(y_{k_i}, z_i^L), \quad k = 0, 1 \quad (3)$$

由于样本存在类别不平衡, 为了消除影响, 通过附加权重的方法来提升欺诈样本的影响。式(3)中, y_0 表示未发生欺诈的用户, y_1 表示欺诈用户。 w_{y_k} 表示 y_k 的类别权重, y_{k_i} 表示第 i 个节点的真实类别, z_i^L 表示第 i 个节点的预测类别。

3.2. Cluster-GCN 模型

为了减少计算的复杂度, 在实际训练阶段, 本文采用了 Cluster-GCN 的方法, 通过图聚类算法将节点分为多个簇, 在每个 epoch 中, 不放回地抽取随机的 q 个簇来构成一个 batch。GCN 和 Cluster-GCN 的复杂度如表 1 所示, 其中 b 表示每个 batch 的节点数。

Table 1. Time complexity and spatial complexity
表 1. 时间复杂度与空间复杂度

	GCN	Cluster-GCN
时间复杂度	$O(LAF + LNF^2)$	$O(LAF + LNF^2)$
空间复杂度	$O(LNF + LF^2)$	$O(bLF + LF^2)$

4. 模型建立与结果分析

4.1. 数据预处理

本文一共收集了 2016 年 20,000 个医保用户的医疗数据, 共包含了 183 万条就诊记录。数据集包含了诊断病种名称与处方中的药物名称, 还根据每次就诊中的处方信息对就诊记录进行了细分, 共记录了 530 多万条就诊的具体项目信息。在 20,000 个医保用户中, 欺诈人数 1000, 占总人数的 5%。由于数据

过多, 研究采用文本识别的方法来对诊断病种和处方药物进行分类, 同时对出现频率较高的文本信息进行了 one-hot 编码, 对编码取值为 1 的处方药物用相应的费用数据来代替, 具体的特征编码见表 2。

Table 2. Feature coding

表 2. 特征编码

编码特征	诊断频次较高的项目
诊断病种名称	糖尿病, 偏瘫, 肾, 高血压, 冠心病, 心脏病, 血症, 血糖, 视网膜, 神经病, 肺, 脑梗, 血管, 贫血, 尿毒症, 血瘀, 骨, 肝, 癌, 其他
药物名称	挂号, 脑心通胶囊, 肝素钠, 心电图, 通心络胶囊, 肾炎康复片, 肝素钙, 肺力咳合剂, 丹芪偏瘫胶囊, 乙肝, 癌胚抗原, 紫外线消毒, 稳心颗粒, 血栓心脉宁片, 参松养心胶囊, 通脉养心丸, 麝香保心丸, 其他

4.2. 建模过程

在生成节点信息的过程中, 本文以个人编码为主键, 分别采用平均和加总两种方式对各用户的处方项目费用特征与诊断病种进行处理, 最终获得 20,000 条动态就诊数据和 82 个特征。在这基础上, 记录参保人曾经就诊过的医院, 并为去过同一个医院的参保人建立关系, 最终得到一个包含了 20,000 个节点和 3200 多万条边的图数据集, 并根据 6:2:2 的比例将数据集分为训练集、验证集和测试集。在实际应用中, 节点的特征会在每一次就诊后都发生更新, 节点之间也可能会产生新的边关系, 由此实现实时识别医疗欺诈。

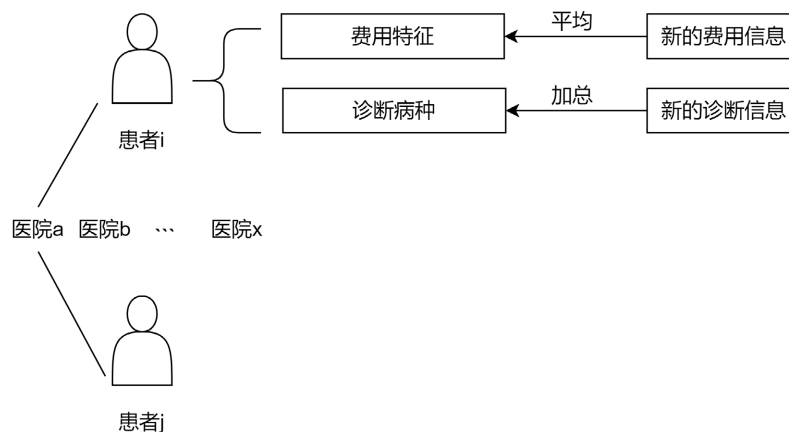


Figure 2. The way to update the doctor-patient information in real time

图 2. 实时更新医患信息的方式

实时更新医患信息的方式如图 2 所示, 患者 i 和 j 由于都去过同一家医院, 就为两个节点添上边, i 和 j 互为邻接节点, 通过图 2 的方式对节点信息进行更新。节点的信息包含了费用信息和诊断病种信息。为了保留过去的就诊记录, 并对未来的就诊行为信息进行更新, 本文采用平均法对每个项目的费用信息进行记录并在未来进行更新, 平均法得到的费用能够用于检验患者该项费用是否异常; 采用加总的方式对诊断病种信息进行记录。由于每条就诊记录的诊断病种进行过 one-hot 编码, 所以将个人的诊断病种特征进行加总就能得到患者被诊断出某一病种的次数。

本文采用图卷积神经网络(GCN)来学习节点表征, 并对节点的类别做出预测。通过设置三层卷积层

来对患者的行为特征进行处理, 并分别将隐藏层特征节点设置为 60、40、20, 每一层患者的节点表征都受到上一层与他去过同一家医院的患者节点表征的影响。本文将学习率设置为 0.001, 分别把图分为了 400、700 和 1000 簇, 并将 batch_size 设置为 10。最终输出欺诈或非欺诈的二分类预测。因此文章的算法考虑了医疗机构监管缺位对患者行为可能存在的引导作用, 以及患者之间的行为影响机制。

4.3. 结果分析

本文分别训练了 400、700 和 1000 簇的 Cluster-GCN 模型, 并比较了模型在测试集上的 AUC、Recall、F1 值和 Accuracy, 具体的结果见表 3。Accuracy 是总体的预测准确率, Accuracy⁰ 和 Accuracy¹ 分别代表非欺诈样本的预测准确率和欺诈样本的预测准确率。

Table 3. Classification performance of the models on the test set

表 3. 模型在测试集上的分类性能

模型	AUC	Recall	F1	Accuracy	Accuracy ⁰	Accuracy ¹	簇数
Cluster-GCN	0.8614	0.9859	0.8774	0.7895	0.7903	0.7725	400
	0.8605	0.9824	0.8886	0.8073	0.8111	0.7381	700
	0.8648	0.9812	0.8851	0.8337	0.8406	0.7095	1000
MLP	0.8024			0.891	0.9116	0.5120	

从 AUC、Recall 和 F1 指标来看, 簇数对模型的分性能影响不大。从 Accuracy 来看, 簇数为 400 的训练策略的总体预测准确率和在非欺诈样本上的预测准确率较低, 但在欺诈样本上的预测准确率最高, 这是因为非欺诈样本数量较多, 所以对总体预测准确率的影响权重更大。随着簇数的增加, 模型在欺诈样本上的预测准确率反而下降。

除了搭建 Cluster-GCN 以外, 本文还训练了多层感知机(MLP)模型。比较二者的分类结果发现, Cluster-GCN 拥有更高的 AUC 和 Accuracy¹。虽然 MLP 有着较高的总体预测准确率 89.1%, 但反欺诈的预测能力即 Accuracy¹ 只有 51.2%。在实际应用中, 医保基金更侧重于识别欺诈用户, 所以 Cluster-GCN 具有更强的医疗保险反欺诈的能力。

5. 结论

随着医疗保险基金规模的日益壮大与就诊信息的爆炸式增长, 医疗保险系统亟需一个能够高效、实时处理大数据的模型来应对欺诈带来的挑战, 减少欺诈损失。本文通过图神经网络来对医疗保险欺诈行为进行检测, 并利用文本分析的方法对诊断病种和处方药物两类特征进行分类。为了实现实时识别欺诈, 设计了个人就诊行为信息的处理方式, 以去过同一家医院为标准建立被保险人与医院、被保险人之间的关系, 并通过 GCN 图神经网络来更新节点之间的信息。研究表明, Cluster-GCN 模型在检测医疗保险欺诈行为上具有更好的分类性能, 其 AUC 达到 0.86 以上, 总体识别欺诈的准确程度达到 83.37%, 反欺诈的识别率为 77.25%。

基金项目

国家社会科学基金项目(17BJY233)资助。

参考文献

- [1] 陈清华, 陈永成, 刘青, 明爱恋, 吴海波. 基于 fsQCA 组态视角的基本医疗保险欺诈骗保的影响因素研究[J]. 医

- 学与社会, 2023, 36(4): 116-121.
- [2] 郑文珍, 翁雯璇, 董火柴. 基于 Apriori-XGB 的医保反欺诈预警模型研究[J]. 产业与科技论坛, 2021, 20(14): 34-36.
 - [3] 郑文珍, 赖俊龙, 翁雯璇. 基于医保欺诈检测的 RF 与 XGB 模型比较[J]. 金融科技时代, 2021(6): 68-73.
 - [4] 李杰, 兰巧玲, 马士豪. 基于大数据的基本医疗保险参保人欺诈风险评估[J]. 中国卫生政策研究, 2018, 11(10): 43-50.
 - [5] 林源. 基于 BP 神经网络的新农合欺诈识别实证研究——以定点医疗机构欺诈滥用为中心[J]. 云南师范大学学报(哲学社会科学版), 2015, 47(3): 117-128.
 - [6] 周杰辉, 朱融晨, 张玮, 陆俊华, 应豪超, 吴健, 陈为. Medicare Vis: 面向医保反欺诈的联合可视分析方法[J]. 医学与社会, 2021, 33(9): 1311-1317.
 - [7] 曹鲁慧, 秦丰林, 闫中敏. 基于 TLSTM 的医疗保险欺诈检测[J]. 计算机工程与应用, 2020, 56(21): 237-241.
 - [8] 李金灿, 徐珂琳, 於州, 魏艳, 仇春涓, 秦国友, 汪荣明, 徐望红. 大数据技术在医保反欺诈中的应用[J]. 中国医疗保险, 2021(1): 48-52.
 - [9] 易东义. 基于关系网与主动学习的医保欺诈识别系统研究及应用[D]: [硕士学位论文]. 深圳: 深圳大学, 2019.
 - [10] 易东义, 邓根强, 董超雄. 基于图卷积神经网络的医保欺诈检测算法[J]. 计算机应用, 2020, 40(5): 1272-1277.
 - [11] 吴文龙, 周喜, 王轶, 王保全. WKAG: 一种针对不平衡医保数据的欺诈检测方法[J]. 计算机工程与应用, 2021, 57(9): 247-254.