

基于组合模型的交通事故严重程度预测方法

石雪怀¹, 戚湧^{1†}, 张伟斌², 李千目¹

(1. 南京理工大学计算机科学与工程学院, 2. 南京理工大学电子工程与光电技术学院 南京 210094)

摘要: 由于各个单一分类模型对道路交通事故严重程度预测的局限性, 本文致力于建立一种组合模型。本文结合卷积神经网络提取时空维度中的特征信息, 采用 stacking 方式将 CNNs 与 XGBoost 组合, 最终生成道路交通事故严重性的分类模型(多层提升算法), 实验结果表明, 此模型在测试集上预测精度为 91.51%, 组合模型比单一分类模型具有更好的分类结果。基于组合模型的分层结果, 对交通事故特征进行重要性排序, 开展特征相关性分析, 为减少道路交通事故及减轻道路交通事故严重等级的管理措施提供参考依据。

关键词: 交通安全; 交通事故严重程度; XGBoost; 卷积神经网络; 诱因分析

中图分类号: TP391.4

文献标志码: A

文章编号:

A Traffic Accident Prediction Approach Based on An Ensemble Approach

SHI Xuehuai¹, QI Yong^{1†}, ZHANG Weibin², LI Qianmu¹

(1. School of Computer Science & Engineering, 2. School of Electronic Engineering & Photoelectric Technology, Nanjing University of Science and Technology, Nanjing 210094, China)

【Abstract】As the limitations of the single classifier on traffic accident severity, an ensemble approach is proposed. Using CNNs to extract the features from the spatial dimension, getting an ensemble approach with XGBoost and CNNs by stacking (Multi-level Boosting Algorithm). The predicting precision of the approach is 91.51% on the validation set. In comparison with the single classification model, the result of the experiment shows a better performance. For providing useful information for reducing the number of traffic accidents and downgrading the severity of traffic accident, the paper gives out a correlation analysis by sorting the features based on the predictions.

【Key words】Traffic Safety, Traffic Accident Severity, XGBoost, CNNs, Inducement Analysis

1. 概述

当今世界, 道路交通事故已对社会造成严重的危害, 成为现代化交通运输中亟需解决的一个重要问题。在我国, 道路交通事故每年都会造成巨大的人员伤亡和经济损失, 2016 年共发生各类道路交通事故伤亡事故超过 15 万, 死伤超过 10 万余人^[1], 直接经济损失 50 亿余元。因此, 对道路交通事故严重程度进行科学预测, 并结合预测结果进行诱因分析, 掌握严重程度较高的道路交通事故的发生、发展、分

布规律与特征, 依此规律和特征为最大程度规避严重交通事故发生的风险提出行之有效的交通安全对策, 成为当前亟需解决的问题。

道路交通事故严重程度预测是以道路交通微观环境活动作为考察对象, 研究各个有关变量分布及其变动对于事故严重程度的影响, 从而提出行之有效的解决方案, 以减轻道路事故严重性和减少道路交通事故。近年来, 国内外学者已经对道路交通事故严重性预测问题进行了大量研究^[2-12], 国内道路交通事故严重程度预测研究如马壮林等采用有

1

收稿日期: 2018-03-04; **修回日期:** **基金项目:** 国家重点研发计划政府间国际科技创新合作重点专项(2016YFE0108000); 江苏省重点研发计划(产业前瞻与共性关键技术)项目(BE2017163); 中央高校基本科研业务费专项资金(No. 30916015104); 中兴合作研究项目(2016ZTE04-11)

作者简介: 石雪怀(1995-), 男, 江苏泰州人, 硕士研究生, 研究方向为数据挖掘(shixuehuai@163.com);

通讯作者: 戚湧(1970-), 男, 江苏泰州人, 博士, 教授、博士生导师, CCF 高级会员(22866S), 研究方向为数据挖掘、智能交通系统;

张伟斌(1975-), 男, 陕西西安人, 博士, 教授, 研究方向为智能交通系统;

李千目(1979-), 男, 江苏南京人, 博士, 教授、博士生导师, CCF 会员(06376S), 研究方向为数据挖掘。

序 Logit 模型和广义 Logit 模型, 建立公路隧道交通事故严重程度预测模型^[2]; 刘海珠等从宏观层面筛选道路交通事故的主要影响因素, 将事故严重程度划分为三个等级, 建立基于累积 Logistic 回归模型的道路交通事故严重程度预测模型^[6]。国外道路交通事故严重程度预测研究如 Manuel Fogue^[8], 通过选取事故严重程度相关数据特征(肇事车行驶速度、肇事车类型、安全气囊状态等)建立事故严重程度智能分类系统来为紧急服务机构提供指导意见; Maher Ibrahim Sameen 等^[9]则采用递归神经网络建立深度学习模型对事故严重性进行分类等。由于事故数据为序列化数据, 在时空上具有一定关联性, O'Donnell 等^[13]采用有序 Logit 模型和有序 probit 模型, 分析事故严重程度与驾驶人属性和车辆特征之间的关系; Maher Ibrahim Sameen 等^[9]则通过建立递归神经网络来捕获数据之间的时空关联, 通过添加多个长短时记忆模型最终通过两个全连接层和 softmax 多分类层实现道路交通事故严重性分类。

在目前交通事故严重程度预测模型中, 更多根据事故结果信息预测事故严重程度, 没有获取足够路段信息, 如马壮林等采用有序 Logit 模型和广义 Logit 模型建立公路隧道交通事故严重程度预测模型^[2]; 在 Manuel Fogue^[8]建立的事故严重程度智能分类系统, 其目的在于紧急服务机构提供指导意见, 没有通过相关数据特征分析事故发生主要原因; Maher Ibrahim Sameen 等^[9]采用递归神经网络建立深度学习模型对事故严重性进行分类, 通过照明情况、路表情况等道路信息特征以及事故发生时间等环境特征对事故严重性分类, 但此模型并没有获取邻近路段信息特征且在验证集上预测精度为 71.77%。因此, 提升道路交通事故严重程度预测精度并依此对事故发生原因进行诱因分析成为道路交通事故严重程度研究的主要内容。基于上述文献分析, 本文通过路号获取事故地点邻近道路信息, 结合卷积神经网络提取道路在空间维度中的信息, 采用 stacking 方式将 CNNs 与 XGBoost 组合, 最终生成道路交通事故严重性的分类模型(多层提升算法), 并依据分类模型的预测结果进行诱因分析, 为减少道路交通事故及减轻道路交通事故严重等级

的管理措施提供参考依据。

2 卷积神经网络与 XGBoost 简介

2.1 卷积神经网络

卷积神经网络(Convolutional Neural Network, CNN)最早由 Kunihiko Fukushima 在 1980 年提出^[15], Lecun 在 1998 年对 BP 算法优化进一步提升 CNN 性能^[16], 2012 年 CNN 在 ImageNet LSVRC-2012 contest^[18]中取得了极高的准确度最终将 CNN 研究推向高潮。

CNN 结构如图 1 所示。由图 1 可知, CNN 主要包含: 1) 卷积层(Convolutional layer); 2) 线性整流层(Rectified Linear Units layer, ReLU layer); 3) 池化(Pooling); 4) 全连接层; 5) 损失函数层。CNN 相比于传统神经网络具有以下优点: a) 能够获取非结构性数据的区域关联信息; b) CNN 实现了特征提取的封装, 简化训练过程; c) 同时学习和产生信息提取和分类; d) 权重共享可以减少网络的训练参数, 使神经网络结构变得更简单, 适应性更强^[16]。CNN 也有一些缺点: a) CNN 的特性使其难以泛化; b) 由于卷积层和池化层的计算量庞大, 使 CNN 的研究依赖于硬件性能的支持^[17]。卷积神经网络目前广泛应用于影像识别、视讯分析、自然语言处理、药物发现和围棋等多个领域。

2.2 XGBoost

XGBoost 由 Chen T 等^[22]于 2016 年提出, 它基于 GBDT(又名 MART(Multiple Additive Regression Tree)), 是一种迭代的决策树算法, 该算法由多棵决策树组成, 所有树的结论累加起来做最终答案^[24], 对 GBDT 进行优化, 提供缓存感知预读取技术、分布式外存计算技术、AllReduce 容错工具提高现有提升树增强算法运算速率, 解决当前提升树增强算法局限于百万级别数据量的问题, 提升算法运行速率。XGBoost 的结构如图 2 所示, 具体介绍详见文献[22][23]。

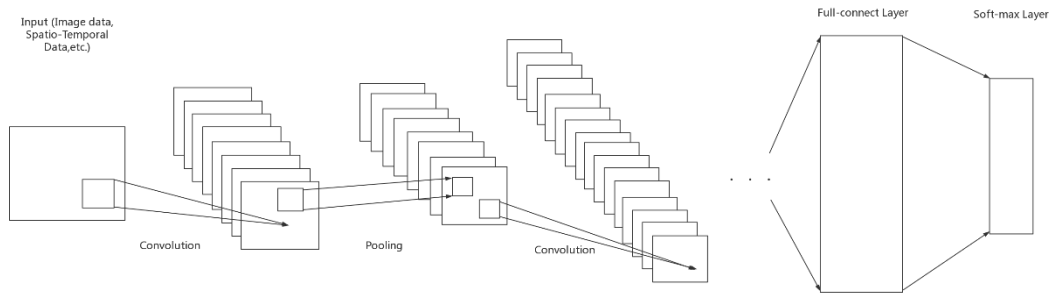


图 1 卷积神经网络结构图

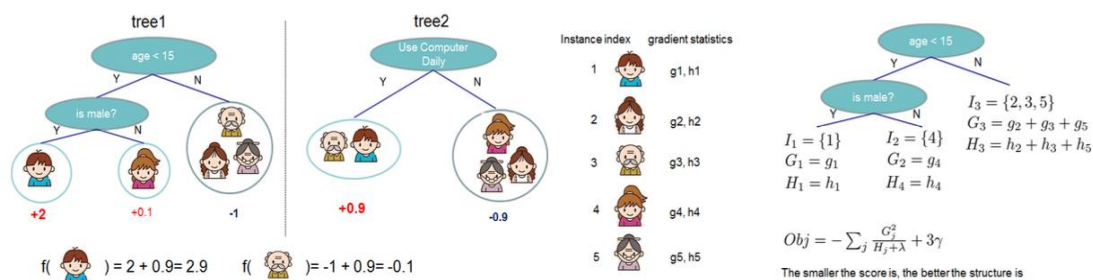


图 2 XGBoost 模型结构和评分方法

数据集对每起事故采集了 69 条特征，本文依照交通安全工程^[14]对特征分类，详情如表 1 所示。

3 模型建立

3.1 强关联规则特征选择方法

根据 GA1082 一般事故采集标准，深圳市事故

表 1 2014-2016 年深圳事故数据特征分类

大类	小类	特征
道路特征	公路	公路行政等级, 公里数, 单向路宽, 双向路宽, 路面结构, 路面附着系数, 道路类型, 道路线型, 长下坡路段, 在道路·横断面位置
	城市道路	路口路段类型
	交通管理设施	交通信号方式（控制）, 交通标志标线完善
	安全防护设施	中央隔离设施, 其他交通安全设施不全, 路侧防护设施类型, 道路物理隔离
环境特征	天气条件	天气, 能见度
	道路环境	事故多发点段, 地形, 照明条件, 路表情况, 路面状况, 道路安全属性, 道路安全隐患督办等级
事故特征	事故基本信息	所属中队, 文书状态, 所辖乡镇, 现场, 行人数量, 行政区划, 调解人, 米数
	事故时空信息	事故发生时间, 事故地点, 路号, 路名
	事故原因信息	事故认定原因, 事故认定原因分类小类
	事故类型信息	事故形态, 事故严重程度, 单车事故, 是否简易程序, 路外事故类型, 车辆间事故, 逃逸事故侦破
	事故伤亡信息	受伤人数, 当事人总数, 抢救死亡人数, 机动车数量, 死亡人数, 直接财产损失, 重伤人数, 轻伤人数, 行人数量, 非机动车数量

由表 1 可知，本文采用 {道路特征, 环境特征, 事故严重程度} 为数据集，以事故严重程度为预测

目标，使用道路特征和环境特征预测事故严重程度，寻求道路特征与环境特征对于事故严重程度的影

响。

本文使用一种强关联规则特征选择算法，通过关联规则进行特征选择，优先选取后件为标签的短规则，达到减小特征空间维度的目的，算法伪代码见算法 1。

基于强关联规则特征选择算法结果，获取预测路段及邻近三条道路特征集合，并结合本路段其他特征集合预测道路交通事故严重程度，特征选择算法结果详见表 2。

算法 1 强关联规则特征选择算法	
1: Input: data set X , feature set F , association rule set Y , true rule set Φ , target label y , iteration number T , support threshold λ , confidence threshold v , final feature set M	
2: Initialize $Y \leftarrow \text{Apriori}(X, F)$, $t \leftarrow 0$, $M \leftarrow \{\}$	
3: for ρ in Y then	
4: if y in ρ and $\text{lift}(y, \rho) > 1$ then	
5: Add ρ to Φ	
6: end if	
7: end for	
8: $\Phi \leftarrow \text{sort}(\Phi)$	
9: while $t < T$ do	
10: Get Q_{first} from Φ	
11: if $v_{Q_{\text{first}}} > v$ and $\lambda_{Q_{\text{first}}} > \lambda$ then	
12: $M_{Q_{\text{first}}} \leftarrow \text{getFeatures}(Q_{\text{first}})$	
13: Add $M_{Q_{\text{first}}}$ to M without repetition	
14: end if	
15: $t \leftarrow t + 1$	
16: end while	
17: Output: M	

表 2 强关联特征选择的特征集合

小类	特征
邻近三段道路特征	事故多发点段, 交通标志标线完善, 其他交通安全设施不全, 路侧防护设施类型, 路口路段类型, 路面状况, 路面结构, 交通信号方式(控制), 道路安全属性, 道路安全隐患督办等级, 道路物理隔离, 道路类型, 道路线型, 长下坡路段
本路段特征	路宽, 中央隔离设施, 是否双道, 事故多发点段, 交通信号方式(控制), 交通标志标线完善, 其他交通安全设施不全, 在道路横断面位置, 地形, 照明条件, 路口路段类型, 路表情况, 路面状况, 路面结构, 道路安全属性, 道路安全隐患督办等级, 道路物理隔离, 道路类型, 道路线型, 长下坡路段, 路侧防护设施类型
环境特征	是否节假日, 是否白天, 天气, 能见度

3.2 事故严重程度预测方法

在深度学习中，特征通过不同层之间进行信息传递，本文受此启发，采用 stacking 的方式将 CNNs 与 xgboost 组合，采用卷积神经网络提取预测路段及邻近的三条路段道路特征的空间信息，将提取的空间信息结合其他特征，采用 stacking 的方式与 xgboost 组合，在此命名为多层提升算法，算法伪代码见算法 2，算法流程图见图 3。

多层提升算法每层包括一个 CNNs 提取空间数据信息输出 3-dim 空间信息的增强特征，结合其他特征代入 xgboost，预测事故严重程度，输入 3-dim 结果向量，并作为增强特征代入下一层的 CNNs 和 xgboost。最终，使用 argmax 确定预测结果。多层提升算法的层数自动确定，受 deep forest^[26]启发，在增加下一层时使用 k 折交叉验证，若预测结果没有明显提升则停止拓展下一层。

算法 2 多层提升算法

```

1: Algorithm configurations: parameters in CNN, parameters in XGBoost
2: Input: data set  $\mathbf{X}$ (include spatial features  $\mathbf{X}_s$ , other features  $\mathbf{X}_o$ ), target label  $\mathbf{y}$ , cross validation parameter  $\mathbf{k}$ , performance estimation  $\lambda$ 
3: Initialize performance estimation  $\mathbf{v} \leftarrow 1$ , level  $\mathbf{l} \leftarrow 0$ 
4: while  $\mathbf{v} > \lambda$  then
5:   if  $\mathbf{l} == 0$  then
6:      $\hat{\mathbf{y}}_{c,l} \leftarrow h\_CNNs(\mathbf{X}_s)$ 
7:     fit CNNs( $\mathbf{X}_s, \mathbf{y}$ )
8:      $\hat{\mathbf{y}}_{XG,l} \leftarrow h\_XGBoost(\mathbf{X}_o, \hat{\mathbf{y}}_c)$ 
9:     fit XGBoost( $\mathbf{X}_o, \hat{\mathbf{y}}_{c,l}, \mathbf{y}$ )
10:     $\mathbf{v} \leftarrow kFold\_estimation(\mathbf{X}, \mathbf{y}, \mathbf{k})$ 
11:   end if
12:    $\mathbf{l} \leftarrow \mathbf{l} + 1$ 
13:    $\hat{\mathbf{y}}_{c,l} \leftarrow h\_CNNs(\mathbf{X}_s)$ 
14:   fit CNNs( $\mathbf{X}_s, \hat{\mathbf{y}}_{XG,l}, \mathbf{y}$ )
15:    $\hat{\mathbf{y}}_{XG,l} \leftarrow h\_XGBoost(\mathbf{X}_o, \hat{\mathbf{y}}_c)$ 
16:   fit XGBoost( $\mathbf{X}_o, \hat{\mathbf{y}}_{c,l}, \mathbf{y}$ )
17:    $\mathbf{v} \leftarrow kFold\_estimation(\mathbf{X}, \mathbf{y}, \mathbf{k})$ 
18: end while
19:  $\hat{\mathbf{y}} = \mathbf{argmax}(\hat{\mathbf{y}}_{XG,l})$ 
20: Output:  $\hat{\mathbf{y}}$ 

```

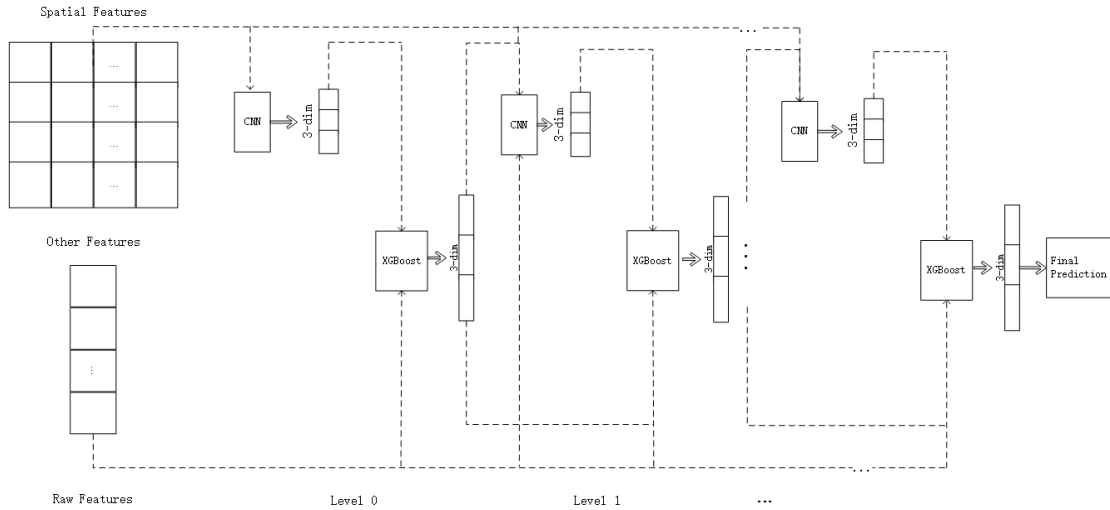


图 3 多层提升算法流程图

4 实验分析

4.1 基于组合模型的预测结果

本文实验数据由道路交通安全研究信息共享平台的深圳市道路交通事故数据分析竞赛提供，数据内容包含深圳市 2014 至 2016 年 20 余万条交通事故信息，包括伤亡事故和轻微事故。本文通过数

据提供的道路信息和环境信息，预测不同事故的严重程度，并采用 5 折交叉验证对预测结果进行评估。本文根据数据预处理结果使用对应分类模型对数据预测。

多层提升算法基于 5 折交叉验证结果对模型参数调整，确定层数为 2 层。最终，采用 5 折交叉验证对此模型预测精度进行评估，多层提升算法预测精度为 91.51%。

4.2 误差分析及比较

为了进一步验证多层提升算法的准确性,使用单一 XGBoost 模型、单一 RF 模型(random forest)

与多层提升算法分类预测结果进行对比,可以判断各模型的准确性,详情如表 3 所示。

表 3 不同模型预测精度比较

评价指标 \ 模型	RF	XGBoost	多层提升算法
acc-on-trainSet	97.29%	97.42%	99.57%
acc-on-cv	74.32%	71.69%	91.51%
test-merror-mean	25.68%	28.31%	8.49%
test-merror-std	16.17%	2.12%	3.90%
train-merror-mean	2.71%	2.05%	0.36%
train-merror-std	None	0.28%	0.11%

表 3 中 acc-on-trainSet 表示模型在用于模型拟合的训练集中的准确率, acc-on-cv 表示模型基于 5 折交叉验证对于测试集的预测准确率, test-merror-mean、test-merror-std、train-merror-mean 和 train-merror-std 是以交叉验证为基础的评估指标, test-merror-mean 表示在 cv 集合中多个测试集中多分类误差平均值, test-merror-std 表示多个测试集中多分类误差的标准差, train -merror-mean 表示多个训练集中多分类误差平均值, train-merror-std 表示多个训练集中多分类误差的标准差。从表 3 可以看出: (1) 无论是 RF、XGBoost 或组合模型, 模型总是在训练集上拟合效果最好, 模型都偏于过拟合。(2) 组合模型相对于 RF 或 XGBoost 单个模型, 在训练集和交叉验证的测试集上表现都更好, 且组合模型在交叉验证的测试集上的多分类误差的标准差为 3.90%, 模型表现较为稳定且良好。

4.3 诱因分析

根据基于组合模型的预测结果, 根据道路交通专家意见, 在组合模型特征重要性排序图的基础上, 本文寻求道路信息特征集合和环境信息特征集合中对道路交通事故严重性预测结果产生重要影响的特征, 并依此进行道路交通事故的诱因分析, 模型重要性排序信息如图 4 所示。由图 4 可得, 路宽、能见度、路侧防护设施类型、路口路段类型、道路物理隔离和交通信号方式(控制)都起到了比较重要的作用。

在图 6 和图 7 中, 对分类变量的不同分类进行了数值变换, 其中“交通信号方式(控制)”中 1-无控制, 2-民警指挥, 3-信号灯, 4-标志, 5-标线, 6-其他, 本实验由 1-6 组合为 11 种类别, 其他特征对应数值的解释详情如表 4 所示。

表 4 不同特征数值解释表

特征	数值	解释	特征	数值	解释
交通信号方式(控制)	1	其它安全设施	路侧防护设施类型	1	绿化带
	2	356		2	混凝土护栏
	3	56		3	防护墩(柱)
	4	34		4	行道树
	5	16		5	金属护栏
	6	3456		6	无防护
	7	456	道路物理隔离	1	中心隔离
	8	45		2	机非隔离
	9	345		3	中心隔离加机非隔离
	10	无信号			
	11	标线			

路口路段类型	1	三枝分叉口			
	2	四枝分叉口		4	无隔离
	3	多枝分叉口			
	4	环行交叉		1	50m 以内
	5	普通路段	能见度	2	50~100m
	6	变窄路段		3	100~200m
	7	路侧险要路段		4	200m 以上
	8	匝道口	事故类型(即事故严重程度)	1	财产损失事故
	9	路段进出处		2	伤人事故
	10	高架路段		3	死亡事故
	11	桥梁			
	12	隧道			
	13	其他特殊路段			

由图 6 和图 7 可得：

路宽：由图 6 中“路宽”与“事故类型”的双变量关系图可得，死亡事故和伤人事故更易出现于路宽较大的路段，相对伤人事故，死亡事故更集中于路宽较大的路段。由图 7 中“路宽”与“事故类型”的双变量关系图可得，两变量 pearsonr 相关系数为 0.23，且假设 pearsonr=0 的 p 值估计为 0.0088，说明“路宽”与“事故类型”有一定的相关性。

交通信号方式（控制）：由图 6 中“交通信号方式（控制）”与“事故类型”的双变量关系图可得，事故更易出现于“交通信号方式（控制）”为 8 的情况下，即由标志和标线组成的交通控制方式。由图 7 中“交通信号方式（控制）”与“事故类型”的双变量关系图可得，相对于财产损失事故，此种控制情况下更易出现伤人事故和死亡事故，且“交通信号方式（控制）”与“事故类型”两变量 pearsonr 相关系数为 0.32，假设 pearsonr=0 的 p 值估计为 0.00022，说明“交通信号方式（控制）”与“事故类型”有一定的相关性。

能见度：由图 6 中“能见度”与“事故类型”的双变量关系图可得，道路交通事故并非出现于能见度很低的情况下，如图可知，大多数道路交通事故出现于能见度为 4 的情况下，即能见度 200m 以上，相对于财产损失事故，伤人事故和死亡事故更易出现于能见度为 4 的情况下。由图 7 中“能见度”与“事故类型”的双变量关系图可得，相对于财产损失事故，能见度 200m 以上的情况下更易出现伤人事故和死亡事故，且“能见度”与“事故类型”

两变量 pearsonr 相关系数为 0.038，假设 pearsonr=0 的 p 值估计为 0.67，说明“能见度”与“事故类型”关联度不大。

路侧防护设施类型：由图 6 中“路侧防护设施类型”与“事故类型”的双变量关系图可得，相对于其他路侧防护设施，1、2 和 6 路侧防护设施类型下，即绿化带、混凝土护栏和无防护，道路交通事故发生概率更大，且相对于财产损失事故，路侧防护设施类型 1 和 6 更易导致伤人事故和死亡事故，即绿化带和无防护。由图 7 中“路侧防护设施类型”与“事故类型”的双变量关系图可得，“路侧防护设施类型”与“事故类型”两变量 pearsonr 相关系数为 0.087，假设 pearsonr=0 的 p 值估计为 0.33，说明“路侧防护设施类型”与“事故类型”关联度不大。

道路物理隔离：由图 6 中“道路物理隔离”与“事故类型”的双变量关系图可得，相对于其他路侧防护设施，道路物理隔离类型 1、4 下，即中心隔离和无隔离下，道路交通事故发生概率更大。由图 7 中“道路物理隔离”与“事故类型”的双变量关系图可得，“路侧防护设施类型”与“事故类型”两变量 pearsonr 相关系数为 0.17，假设 pearsonr=0 的 p 值估计为 0.053，说明“道路物理隔离”与“事故类型”有一定的关联性。

路口路段类型：由图 6 中“路口路段类型”与“事故类型”的双变量关系图可得，大部分道路交通事故都发生在路口路段类型 2 和 5 的情况下，即四枝分叉口和普通路段，且相对于财产损失事故和

死亡事故，路口路段类型 2 下更易发生伤人事故。由图 7 中“路口路段类型”与“事故类型”的双变量关系图可得，“路口路段类型”与“事故类型”两变量 pearsonr 相关系数为-0.3，假设 pearsonr=0 的 p 值估计为 0.00051，说明“路口路段类型”与“事故类型”的关联性很强。

综上所述，相对于财产损失事故，在道路较宽的路段易出现严重程度更高的事故；在由标志和标线组成的交通控制方式更易出现伤人事故和死亡事故；且不同于“恶劣天气更易发生严重道路交通事故”，数据显示能见度为 4 的情况下，即能见度 200m 以上；相对于财产损失事故，伤人事故和死亡事故出现次数更多，相对于其他路侧防护设施，绿化带、混凝土护栏和无防护，道路交通事故发生概率更大，且相对于财产损失事故，路侧防护设施类型为绿化带和无防护的情况下，伤人事故和死亡事故发生次数更多；相对于机非隔离和中心隔离加机非隔离，道路物理隔离类型为中心隔离或无隔离的情况下，道路交通事故发生概率更大；且在四枝分

叉口路段更易发生伤人事故。因此我们应根据上述数据描述结果采取合理手段减轻道路交通事故严重程度以及尽量避免道路交通事故。

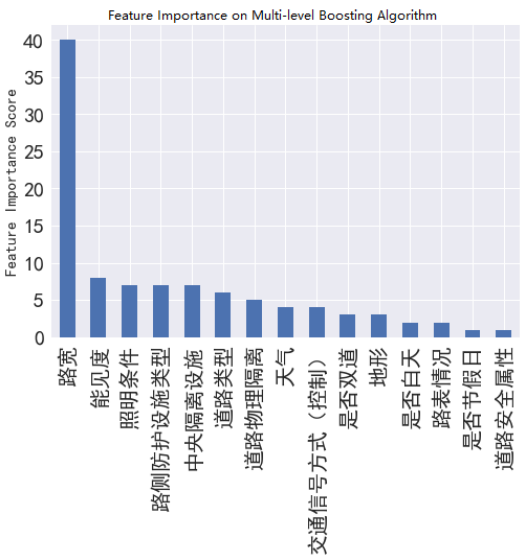


图 4 模型特征重要性排序

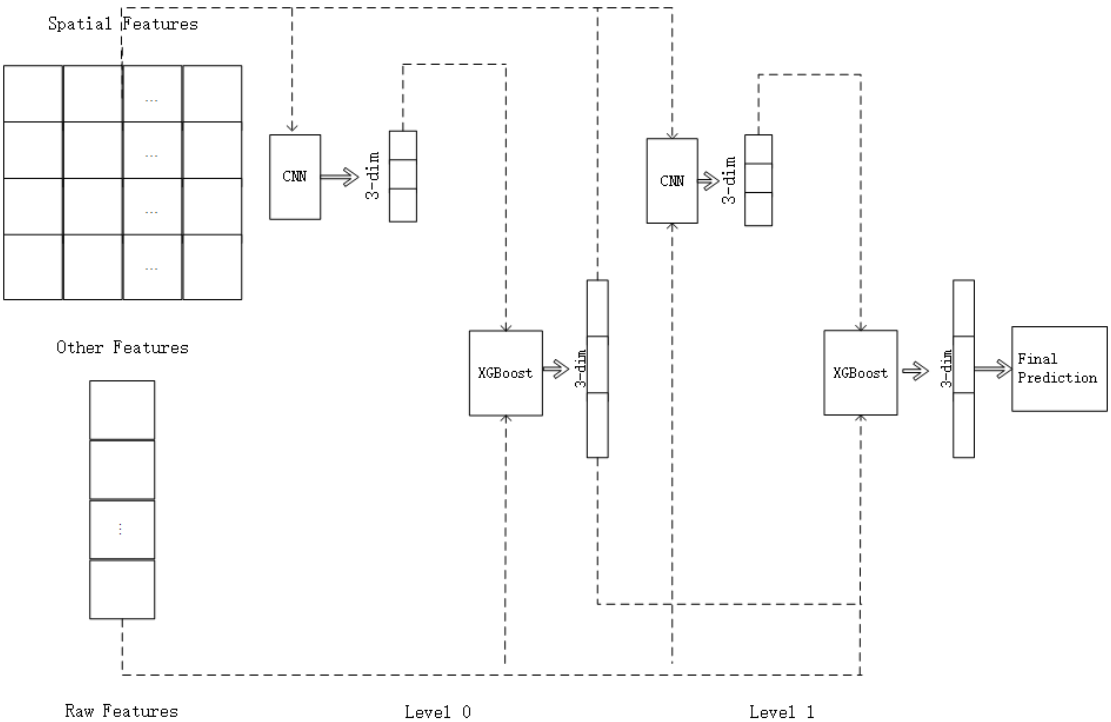


图 5 本文的多层提升算法架构

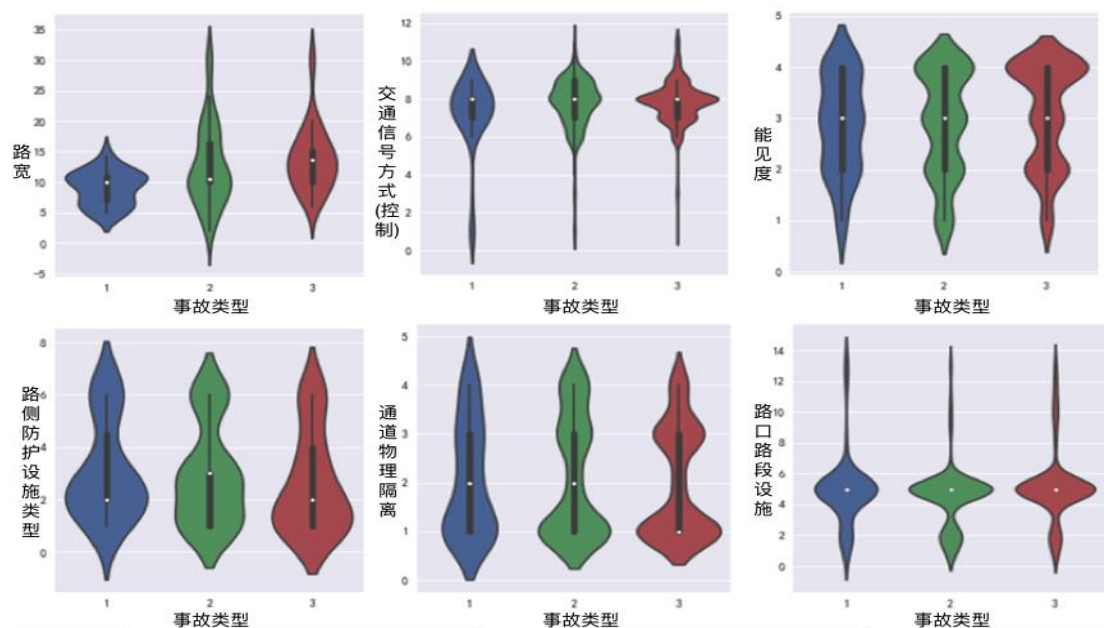


图 6 不同特征与“事故类型”关系图（小提琴图）

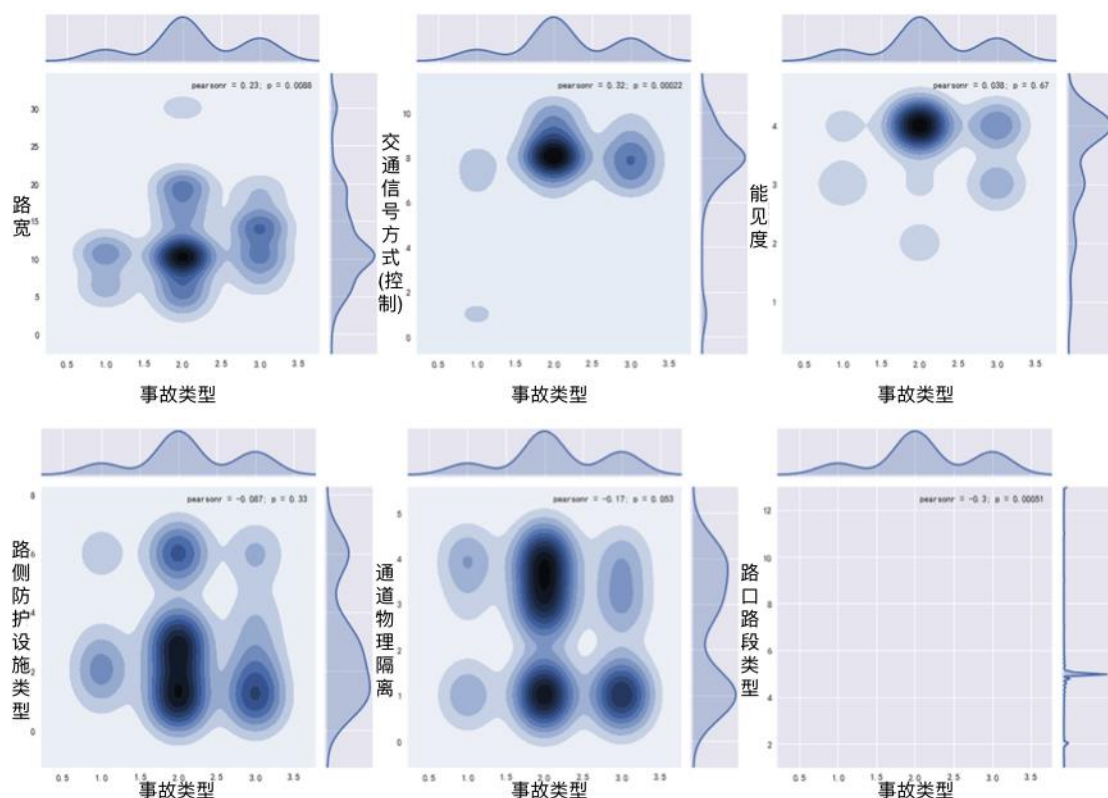


图 7 不同特征与“事故类型”关系图（KDE 图）

5 结论

本文根据 2014-2016 年深圳道路事故数据的实际情况，提出了基于组合模型的道路交通事故严重程度预测方法，结合卷积神经网络提取时空维度中的特征信息，采用 stacking 方式将 CNNs 与 XGBoost 组

合，最终生成道路交通事故严重性的分类模型（多层提升算法），实验结果表明，组合模型比任何单一分类模型具有更好的分类结果，是一种有效的道路交通事故严重程度预测方法。基于组合模型优秀的分类结果，对特征进行重要性排序，进行特征相关性分析，最终为减少道路交通事故及减轻道路交通事故严重等级提供参考意见。根据道路交通安全研

究信息共享平台的要求,本数据不得对外发布、传播,因此在此只公开本实验相关代码,详见<https://github.com/SXHSine/TrafficAccidentAnalysis>。

由于道路交通事故的发生具有时空关联性,在将来的研究中,通过获取事故的路段空间数据,可以根据空间位置更为准确地获取附近路段信息,获取更多邻近道路的空间信息,结合有关驾驶人的相关属性和操作行为数据,从而更准确地预测道路交通事故严重程度,并分析致灾原因,为避免事故的发生提供更多建设性意见。

参考文献:

- [1]. 2016 年全国道路交通事故数据统计[N/OL]. (2017-03-01) [2017-12-14].
<http://www.peichang.cn/detail/id18666.html>
- [2]. 马壮林, 张祎祎, 杨杨, 等. 公路隧道交通事故严重程度预测模型研究[J]. 中国安全科学学报, 2015, 25(5):75-79.
- [3]. 孙重静, 辛飞飞. 人车冲突风险度评价指标计算研究综述[J]. 交通信息与安全, 2016, 34(2):9-16.
- [4]. 马壮林, 邵春福, 董春娇, 等. 基于累积 Logistic 模型的交通事故严重程度时空分析[J]. 中国安全科学学报, 2011, 21(9):94-99.
- [5]. 侯树展, 孙小端, 贺玉龙, 等. 高速公路交通事故严重程度与交通流特征的关系研究[J]. 中国安全科学学报, 2011, 21(9):106-112.
- [6]. 刘海珠. 道路交通事故严重程度影响因素分析及预测模型建立[D]. 吉林大学, 2014.
- [7]. 马柱, 陈雨人, 张兰芳. 城市道路交通事故严重程度影响因素分析[J]. 重庆交通大学学报(自然科学版), 2014, 33(1):111-114.
- [8]. Fogue M, Garrido P, Martinez F J, et al. Using Data Mining and Vehicular Networks to Estimate the Severity of Traffic Accidents[M]// Management Intelligent Systems. Springer Berlin Heidelberg, 2012:37-46.
- [9]. Sameen M I, Pradhan B. Severity Prediction of Traffic Accidents with Recurrent Neural Networks[J]. Applied Sciences, 2017, 7(6).
- [10]. Khera D, Singh W. A Review on Injury Severity in Traffic System using Various Data Mining Techniques[J]. International Journal of Computer Applications, 2014, 100(3):17-22.
- [11]. Gupta M, Solanki V K, Singh V K. A Novel Framework to Use Association Rule Mining for classification of traffic accident severity[J]. Ingeniería solidaria, 2017, 13(21): 37-44.
- [12]. Baru A. Injury Severity Levels and Associated Factors among Road Traffic Accident Victims Referred To Emergency Departments of Selected Public Hospitals in Addis Ababa, Ethiopia: The Study Based On Haddon Matrix[J]. 2017.
- [13]. O'Donnell C J, Connor D H. Predicting the severity of motor vehicle accident injuries using models of ordered multiple choice. [J]. Accident Analysis & Prevention, 1996, 28(6):739-753.
- [14]. 肖贵平, 朱晓宁. 交通安全工程[M]. 中国铁道出版社, 2011.
- [15]. Fukushima K, Miyake S. Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition[M]//Competition and cooperation in neural nets. Springer, Berlin, Heidelberg, 1982: 267-285.
- [16]. LeCun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition[J]. Proceedings of the IEEE, 1998, 86(11): 2278-2324.
- [17]. Convolutional Neural Networks (LeNet) - DeepLearning 0.1 documentation. DeepLearning 0.1. LISA Lab. [31 August 2013].
- [18]. Olshausen B A. Emergence of simple-cell receptive field properties by learning a sparse code for natural images[J]. Nature, 1996, 381(6583): 607-609.
- [19]. LeCun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition[J]. Proceedings of the IEEE, 1998, 86(11): 2278-2324.
- [20]. Sabes P N, Jordan M I. Advances in neural information processing systems[C]//In G. Tesauro & D. Touretzky & T. Leed (Eds.), Advances in Neural Information Processing Systems. 1995.
- [21]. Krizhevsky A, Sutskever I, Hinton G E. ImageNet classification with deep convolutional neural networks[J]. Communications of the ACM, 2017, 60(2):2012.
- [22]. Chen T, Guestrin C. Xgboost: A scalable tree boosting system[C]//Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining. ACM, 2016: 785-794.

- [23]. Chen T, He T, Benesty M, et al. xgboost: Extreme Gradient Boosting[J]. 2017.
- [24]. Friedman J H. Greedy function approximation: a gradient boosting machine[J]. Annals of statistics, 2001: 1189-1232.
- [25]. Schapire R E, Singer Y. Improved Boosting Algorithms Using Confidence-rated Predictions[J]. Machine Learning, 1999, 37(3):297-336.
- [26]. Zhou Z H, Feng J. Deep forest: Towards an alternative to deep neural networks[J]. arXiv preprint arXiv:1702.08835, 2017.