

基于Stacking多模型融合的电信用户 满意度影响因素研究

覃林, 谢本亮

贵州大学大数据与信息工程学院, 贵州 贵阳

收稿日期: 2024年6月28日; 录用日期: 2024年7月8日; 发布日期: 2024年8月23日

摘要

随着移动通信技术的飞速发展, 电信用户群体数量不断的攀升, 也越要求运营商重视用户使用体验, 不断提升网络使用满意度。本文基于2022年北京移动提供的客户语音和上网业务数据, 首先使用灰色关联分析筛选出满意度重要影响因素, 然后采用Stacking多模型融合策略, 结合使用随机森林、逻辑回归、K近邻、Adaboost、XGBoost共5种算法, 对客户满意度打分进行预测研究, 融合后模型在语音业务数据中预测准确率为0.613, 在上网业务数据中预测准确率为0.606, 为电信用户满意度评分预测和分析研究提供了一定的理论参考。

关键词

用户体验, Stacking, Adaboost, 满意度预测

Research on Influencing Factors of Telecommunication User Satisfaction Based on Stacking Multi-Model Fusion

Lin Qin, Benliang Xie

College of Big Data and Information Engineering, Guizhou University, Guiyang Guizhou

Received: Jun. 28th, 2024; accepted: Jul. 8th, 2024; published: Aug. 23rd, 2024

Abstract

With the rapid development of mobile communication technology, the number of telecommunication user groups continues to increase, and the more operators are required to pay attention to

user experience and continuously improve the satisfaction of network usage. This paper is based on the customer voice and Internet service data provided by Beijing Mobile in 2022, first we use grey correlation analysis to filter out the important factors affecting satisfaction, and then we adopt the stacking multi-model fusion strategy, combined with the use of Random Forest, Logistic Regression, K Nearest Neighbours, Adaboost and XGBoost in a total of five algorithms, to carry out prediction research on customer satisfaction scoring. The prediction accuracy of the model is 0.613 in voice service data and 0.606 in Internet service data, which provides certain theoretical reference for the prediction and analysis research of telecom customer satisfaction scoring.

Keywords

User Experience, Stacking, Adaboost, Satisfaction Prediction

Copyright © 2024 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

随着我国互联网不断发展壮大, 电信用户群体数量不断攀升, 三大运营商的用户数量也不断增长。随着工信部“提速降费”的政策施行, 也越来越要求运营商不断调整产品、服务、技术, 重视用户的使用体验[1], 提升网络服务质量, 提供给客户更好的使用体验。

在运营商的服务过程中, 客户满意度往往是运营商市场运营状况的重要体现[2]。客户满意度是客户对运营商产品服务的满意程度, 反映了客户期望与实际感知的产品服务之间的差异[3]。在数字化时代的今天, 随着用户数量的大幅增加, 移动产品的种类越来越丰富, 用户的意见数据也是海量和复杂的。因此, 利用海量的用户投诉和建议数据, 通过大数据分析手段, 建立客户满意度预测算法, 可以为运营商决策赋能, 有效提升客户的使用体验, 推动移动网络高质量可持续发展。

2. 特征处理分析

2.1. 特征预处理

特征工程是在应用机器学习训练前, 预先从原数据提取特征或者筛选特征的过程[4]。本文先对语音业务数据和上网业务数据做数据预处理, 剔除异常的数据, 删除缺失值过多的属性值数据, 然后对缺失值不太多的数据, 使用平均值填充数值型数据, 使用众数填充类别数据, 来达到对数据整体清洗。对于语音业务数据和上网业务的预测数据, 采用相同的数据预处理方法。

2.2. 灰色关联分析

灰色关联度分析(Grey Relation Analysis, GRA) [5], 是一种多因素统计分析的方法。在一个灰色系统中, 先假定关注的项目是受其他的因素影响或者说是相关的, 然后分析各因素的相关程度的方法。灰色关联度分析实质上是对不同序列两点间距离的反映, 使用数值的指标来描述事物或因素之间的关联程度。灰度关联度系数计算公式如下:

$$\zeta_i(k) = \frac{\min_i \min_k |x_0(k) - x_i(k)| + \rho \cdot \max_i \max_k |x_0(k) - x_i(k)|}{|x_0(k) - x_i(k)| + \rho \cdot \max_i \max_k |x_0(k) - x_i(k)|} \quad (1)$$

式中: x_0 为参考序列, x_i 为待比序列, ρ 为分辨系数, 通常取 0.5。

2.3. 特征选择

在机器学习中, 特征选择通常是重要的一步。一般数据的维度可能较大, 从数据集中选择出有效的特征, 可以降低数据的维度, 减少模型的学习成本, 提高模型的泛化性能。

3. Stacking 融合预测

3.1. 随机森林

随机森林[6]运用了 Bagging 集成学习策略, 它属于一个基于树的集合体, 每棵树都与一组随机变量存在关联。从输入数据里随机且有放回地选取特定数量的样本数据集, 接着为每个样本数据集配置一棵决策树, 最终对所有决策树的判别类结果实施投票操作。正是由于数据和属性都具备随机性, 极大地提升了随机森林的集成泛化能力。

3.2. 逻辑回归

逻辑回归[7]的基本思想是建立一条线(判定边界)把数据以分类目标划分开来。逻辑回归通常使用 Sigmoid 函数, 学习得到用于拟合的系数项参数集, 并用于组成线性回归的判定边界。逻辑回归具有较好的可解释性。但对输入数据的特征数量较为敏感。

3.3. K 近邻

K 最近邻(K-Nearest Neighbor, KNN) [8]是邻近算法采用 K 个邻居的特例, 其思想为: 若一个特征空间中大多数的样本属于某一个类别, 则在这个特征空间中, K 个最相似的样本也属于这个类别, 同时有这个类别上样本的特性。

3.4. Adaboost

Adaboost [9]依托于 Boost 集成学习方式, 具备将预测精准度仅稍高于随机猜测的弱学习器强化为预测精准度颇高的强学习器的能力。在训练过程中, 每一轮的迭代里, 都会增添一个新的弱分类器, 一直到错误率达到预先设定的足够小的值, 或者达到预先规定的最大迭代次数, 然后确定最终的强分类器。

3.5. XGBoost

XGBoost [10]是 boosting 算法的一种实现方式, 其采用多个简单的基学习器, 通过不断降低模型值和实际值的差, 不断生成新的树, 每棵树都是基于上一棵树和目标值的差值来进行学习, 从而有效降低模型的偏差, 在各类数据分析场景中都被广泛应用。XGBoost 的目标函数如式(2)所示, 其由模型预测误差和模型结构误差组成。

$$Obj(\theta) = L(\theta) + \Omega(\theta) = L(y_i, y_i') + \sum_{k=1}^l \Omega(f_k(x_i)) \quad (2)$$

3.6. Stacking

Stacking [11]算法的思想是基于预测的层次融合, 第一层先分别训练好各个单模型, 形成初级学习器。然后采用交叉验证法, 对预测结果进行交叉预测结果取平均, 最后默认使用逻辑回归分类器输出结果, 可以有效提升模型的泛化能力。

4. 实验环境与结果分析

4.1. 实验环境

实验在一台 Intel(R) Core(TM) i5-12400 CPU、RTX 3060 GPU、Windows64 位操作系统的计算机上进行实验。数据集使用的是北京移动在 2022 年 MathorCup 高校数学建模挑战赛——大数据竞赛上提供的数
据, 其中包含了语音业务用户满意度数据及其供预测的数据, 涉及了整体满意度、网络覆盖、语音清晰
度、场景描述、通话问题、资费问题等方面数据; 上网业务用户满意度数据及其供预测的数据, 其中涉
及了整体满意度、上网稳定性、场景描述、异常问题、APP 场景等方面数据。训练各算法时, 对语音业
务和上网业务数据分别以 7:3 划分训练集和测试集。

4.2. 评价标准

本文使用准确率作为主要评价指标。准确率是正确分类的样本数与总样本数之比, 衡量的是算法总
体分类能力。准确率计算公式如下式所示:

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

4.3. 特征筛选实验

由于语音业务和上网业务数据属性过多, 语音业务除去语音通话整体满意度属性, 共有 54 条属性。
上网业务除去手机上网整体满意度属性, 共有 124 条属性。若将全部属性输入预测, 则数据的维度是极
大的, 不利于算法进行学习。通过分析, 决定用户对语音业务和上网业务体验满意度的因素是复杂的,
但是可以尝试找出影响满意度的主要影响因素, 然后进行建模预测, 这样一方面降低了模型的学习成本,
另一方面模型可解释性也更好。

本文将经过特征预处理的语音业务数据进行灰色关联分析, 以语音通话整体满意度为参考对象, 其
他各因素为评价对象, 计算其关联度, 根据其关联度数据绘制出如图 1 所示的可视化相关性热图。

从图 1 中可以看出, 语音通话整体满意度为参考对象, 相关性为 1。语音通话清晰度、网络覆盖与
信号强度、语音通话稳定性、是否遇到过网络问题、居民小区、手机没有信号、有信号无法拨通、通话
过程突然中断、办公室等 9 个因素对语音通话整体满意度的影响由强至弱。

同理, 以手机上网整体满意度为参考对象, 经过对特征预处理的上网业务数据进行灰色关联分析,
根据其关联度数据绘制出如图 2 所示的可视化相关性热图。

从图 2 中不难看出, 手机上网整体满意度为参考对象, 相关性为 1。手机上网速度、网络覆盖与信号强
度、手机上网稳定性、网络信号差/没有信号、手机上网速度慢、上网过程中网络时断时续或时快时慢、打
开网页或 APP 图片慢、居民小区、显示有信号上不了网等 9 个因素对手机上网整体满意度的影响由强到弱排序。

在特征选择方面, 使用灰色关联度分析得到的结果, 选择重要的影响因素作为数据输入。在语音业
务方面, 选择语音通话清晰度、网络覆盖与信号强度、语音通话稳定性、是否遇到过网络问题、居民小
区、手机没有信号、有信号无法拨通、通话过程突然中断、办公室共 9 个影响显著因素数据; 在上网业
务方面, 选择手机上网速度、网络覆盖与信号强度、手机上网稳定性、网络信号差/没有信号、手机上网
速度慢、上网过程中网络时断时续或时快时慢、打开网页或 APP 图片慢、居民小区、显示有信号上不了
网共 9 个影响显著因素数据。

4.4. 单模型预测实验

在使用 Stacking 融合模型之前, 需要先对各个模型进行训练调参, 以确保各个单模型具有较好的

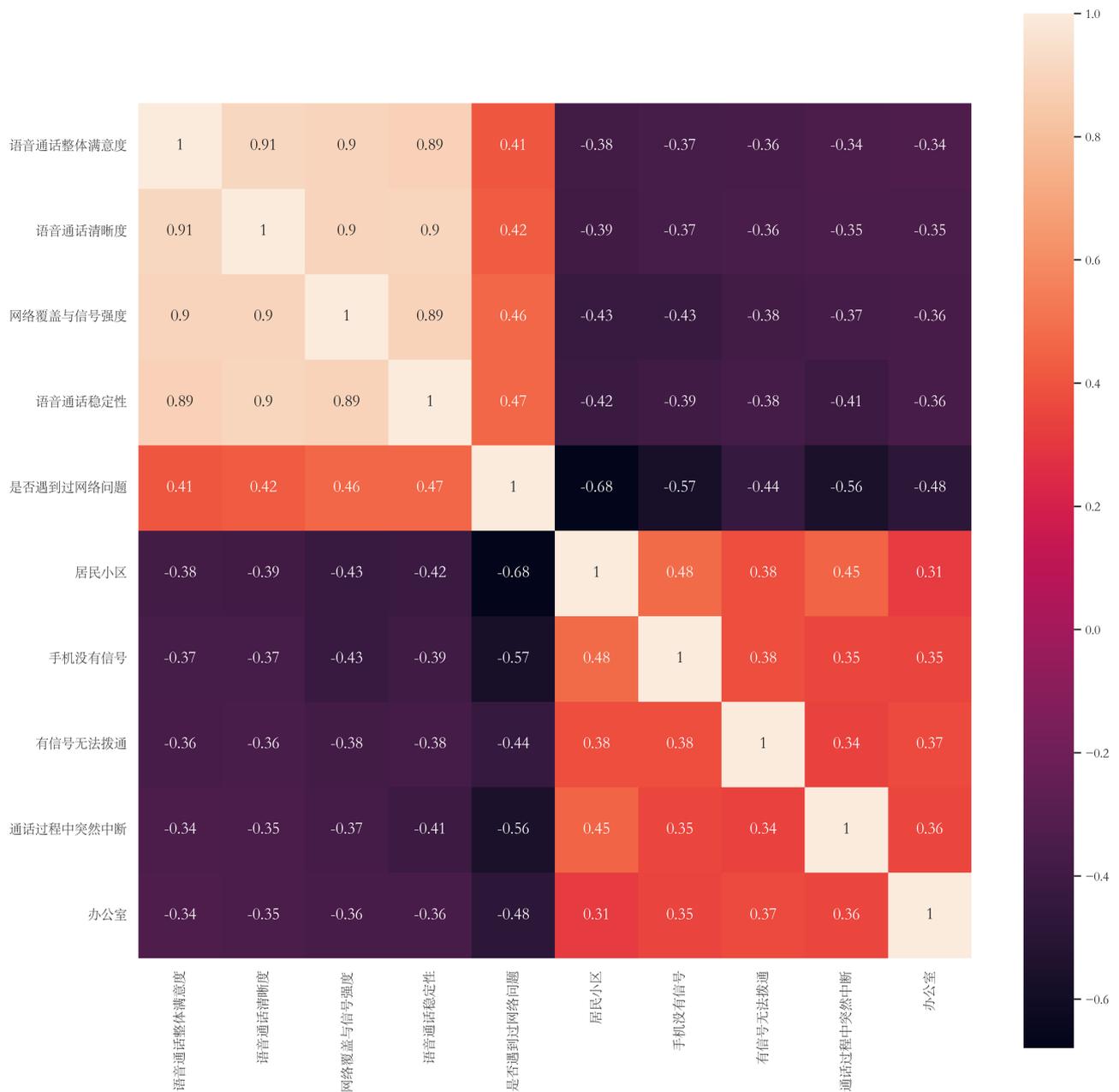


Figure 1. Heatmap of correlation of factors influencing voice service satisfaction

图 1. 语音业务满意度影响因素相关性热图

预测性能, 融合后, 才有可能达到更好的效果。各个单模型的调参方法不尽相同, 但都主要以准确率为目标进行调参。

随机森林调参。随机森林模型调参一般经过两个步骤, 一是确定最佳决策树数, 二是确定最佳最大深度。首先从 0 到 200 决策树数遍历搜索, 绘制模型准确率 acc 曲线, 初步确定稳定的最优决策树数范围在 25 到 100 较好。继续在 50 到 80 之间搜索, 最终确定最佳决策树数去两个极大值 70 和 71。树深度经过网格搜索[12]深度的最佳数为 1。

逻辑回归参数调优。逻辑回归属于凸优化问题, 实际要调的参数不多。首先使用网格搜索方式,

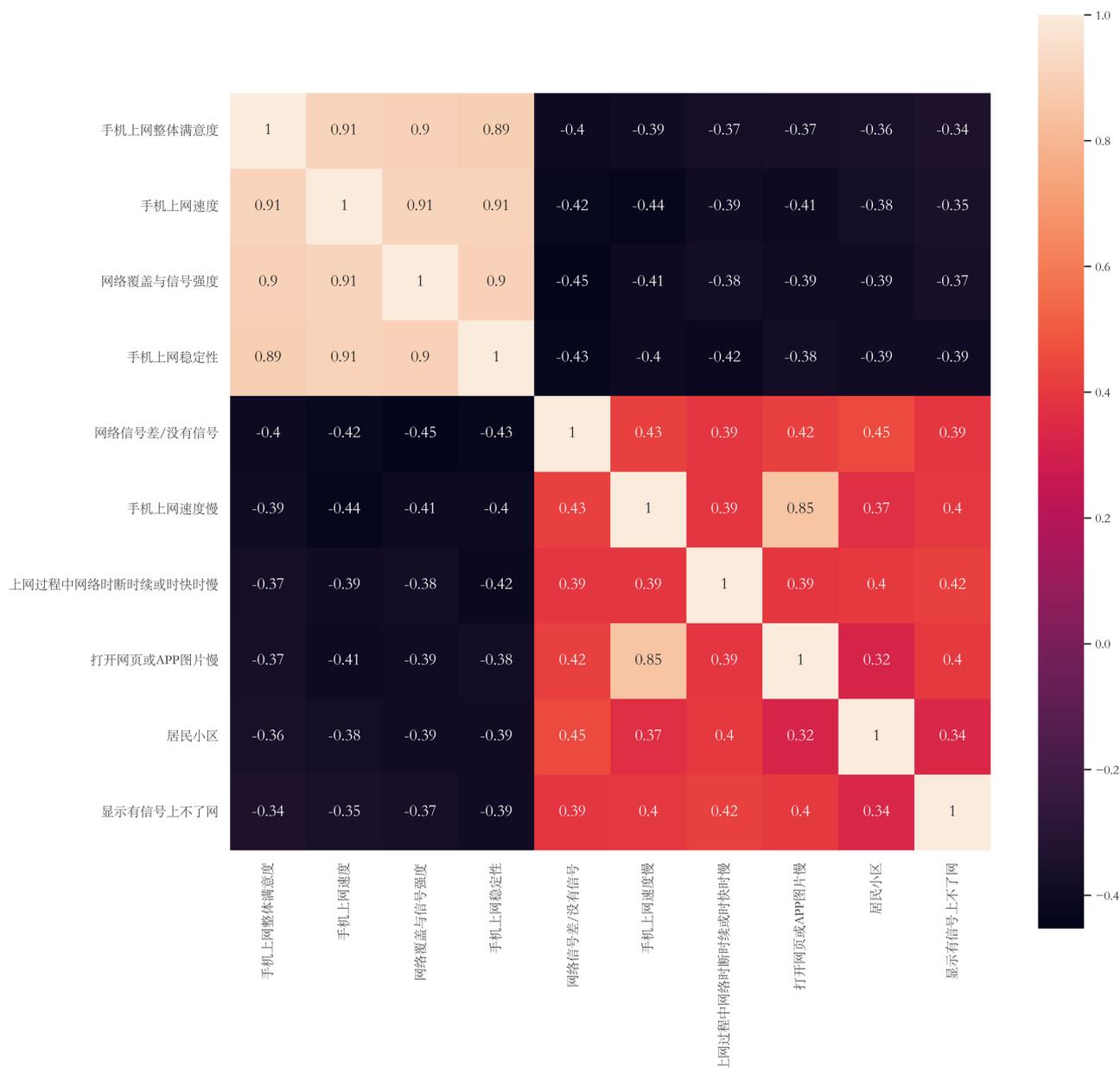


Figure 2. Heatmap of correlation of factors influencing satisfaction with Internet service

图 2. 上网业务满意度影响因素相关性热图

搜索 C (正则化系数 λ 的倒数)和 solver (优化算法选择参数), 最终确定最佳 C 值为 0.05, 最佳 solver 为 liblinear。

K 近邻参数调优。KNN 参数调优主要是寻找最佳 K 值和距离的权重, KNN 的 K 值选取较为重要。通过网络搜索, 在 1 到 11 搜索最佳 K 值, 距离权重在 1 到 6 搜索, 最终得出最佳 K 值为 10。

Adaboost 参数调优。在集成学习中, 参数调优一般是先选择框架的参数进行调整, 再选择基学习器的参数。在框架参数调优方面, 首先对弱学习器个数进行择优, 同样是设定一个大的范围数, 然后遍历尝试, 若范围不好确定, 可以再次缩小范围搜索。在基学习器参数调优方面, 迭代尝试树的最大深度 max_depth 和节点再划分所需最少样本数 min_samples_split。最后还应当用模型准确率对学习率与基学习

器个数进行迭代测试, 选出最佳参数。最终确定基学习器数为 10, 最大深度 `max_depth` 为 1, 节点再划分所需最少样本数 `min_samples_split` 为 18, 弱学习器的权重缩减系数为 0.01。

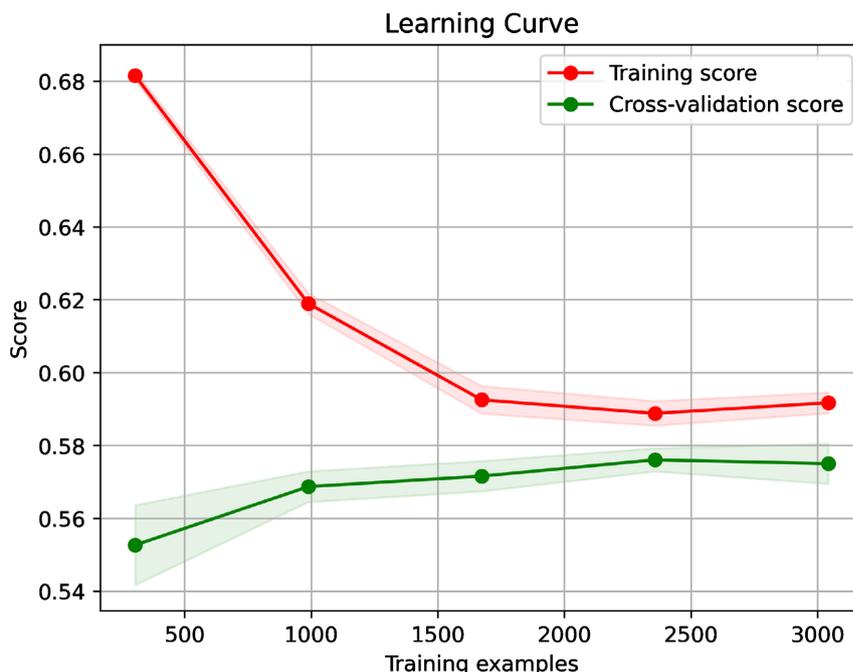


Figure 3. XGboost learning curve
图 3. XGboost 学习曲线

经过测试, 对随机森林、逻辑回归、K 近邻、Adaboost、XGBoost 共 5 个模型分别训练和调参, 5 种算法在语音业务数据中预测准确率分别为 0.582、0.575、0.583、0.584、0.576, 在网上业务数据中预测准确率分别为 0.573、0.581、0.582、0.577、0.581。

4.5. Stacking 融合预测实验

本文使用 Stacking 算法, 设置 5 折交叉验证, 将随机森林、逻辑回归、K 近邻、Adaboost、XGBoost 共 5 中算法融合预测, 其流程图如图 4 所示。K 折交叉验证是一种常见的验证模型性能的统计分析方法, 其基本思想是将输入数据分为 K 组, 其 K - 1 组作为训练集, 剩下的一组作为测试集, K 折的预测结果会用平均值计算代替。

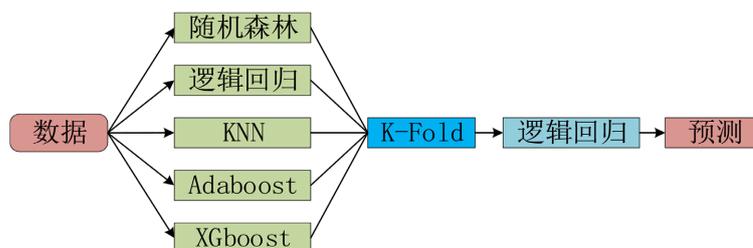


Figure 4. Stacking fusion prediction chart
图 4. Stacking 融合预测图

具体实现过程: 设置输入各个经过调优的单模型, 保证在参数设置上一致; 第一层指定为各个调优

的单模型, 第二层指定为逻辑输出; 从第一层分类器中依次加载单模型, 并设置 5 折交叉验证; 训练 Stacking 并进行测试准确率。

经过实验, 在语音业务数据中, 融合 5 种模型后的预测准确率为 0.613。在上网业务数据中, 融合 5 种模型后的预测准确率为 0.606。

4.6. 结果分析

相较于各个单模型预测, Stacking 融合预测显著提升了模型的预测准确性。在语音业务数据中, 单模型准确率最高为 Adaboost, 0.584, Stacking 融合预测准确率为 0.623, 提升了 6.7%。在上网业务数据中, 单模型准确率最高为 K 近邻, 0.582, Stacking 融合预测准确率为 0.616, 提升了 5.8%。满意度评分因素分析和预测分析通常较为复杂, 本文的预测也一定程度上受限于样本数据较小的问题, 准确难以达到理想的指标。

5. 总结

对于海量复杂的电信用户满意度数据, 常规预测方法往往难以胜任预测的任务。本文实验结果证明了 Stacking 融合预测方法的有效性, 增强了模型的泛化能力和鲁棒性, 有效提升了满意度的预测准确性, 可以助力运营商用户在业务问题中的早期预警和使用体验研究中的决策赋能, 为电信用户满意度评分预测和分析研究提供了一定的理论参考。

参考文献

- [1] 我国明确四项举措推进网络提速降费[J]. 新疆农垦科技, 2021, 44(2): 68.
- [2] 刘凡. 电信用户满意度预测研究[D]: [硕士学位论文]. 杭州: 浙江工商大学, 2020.
- [3] 陈燕, 卢小静. 新疆电信宽带用户满意度影响因素研究[J]. 山东纺织经济, 2010(9): 88-90.
- [4] 周志华. 机器学习[M]. 北京: 清华大学出版社, 2016.
- [5] 张碧清, 韩啸, 贺瑞军, 等. 基于灰色关联度分析的不锈钢离子渗氮层对磨损性能的影响[J]. 金属热处理, 2024, 49(3): 174-181.
- [6] 许振腾, 王琪. 基于随机森林算法的航班延误时间预测模型研究[J]. 滨州学院学报, 2024, 40(2): 28-35.
- [7] 赖红清. 基于逻辑回归的企业二次创业金融数据分类方法研究[J]. 重庆工商大学学报(自然科学版), 2021, 38(5): 114-119.
- [8] 刘晴晴. 基于 K 近邻的变权组合预测模型及应用[J]. 科学技术创新, 2021(14): 28-29.
- [9] Cao, Y., Miao, Q., Liu, J. and Gao, L. (2013) Advance and Prospects of Adaboost Algorithm. *Acta Automatica Sinica*, 39, 745-758. [https://doi.org/10.1016/s1874-1029\(13\)60052-x](https://doi.org/10.1016/s1874-1029(13)60052-x)
- [10] Chen, T. and Guestrin, C. (2016) XGBoost. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, 13-17 August 2016, 785-794. <https://doi.org/10.1145/2939672.2939785>
- [11] 王琳, 周捷, 林海飞, 等. 基于 Stacking 集成模型的煤层瓦斯含量预测研究[J]. 煤炭工程, 2024, 56(4): 125-132.
- [12] 王昭栋, 王自法, 李兆焱, 等. 基于机器学习-网格搜索优化的砂土液化预测[J]. 振动与冲击, 2024, 43(5): 82-93.