

保险行业数据可视化分析

制作人： 吕艾（CDA持证人、会员）

联系方式： IVYLUE@163.COM

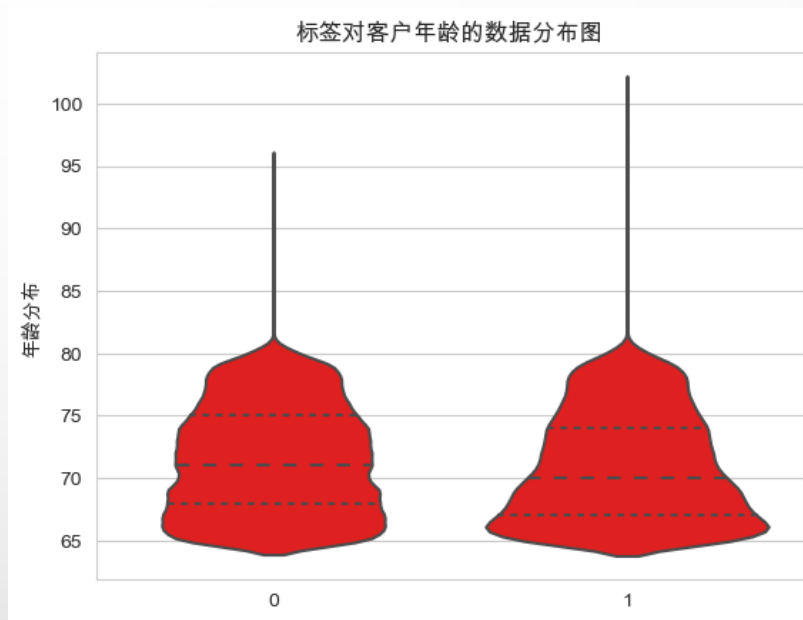
WECHAT ID: IVYLUE

数据建模与业务分析概览

- 第一阶段：数据集的准备与处理
- 第二阶段：特征处理
- 第三阶段：建模、模型评估
- 第四阶段：可视化分析
- 总结：面板可视化总结

第一阶段：数据集的准备与处理

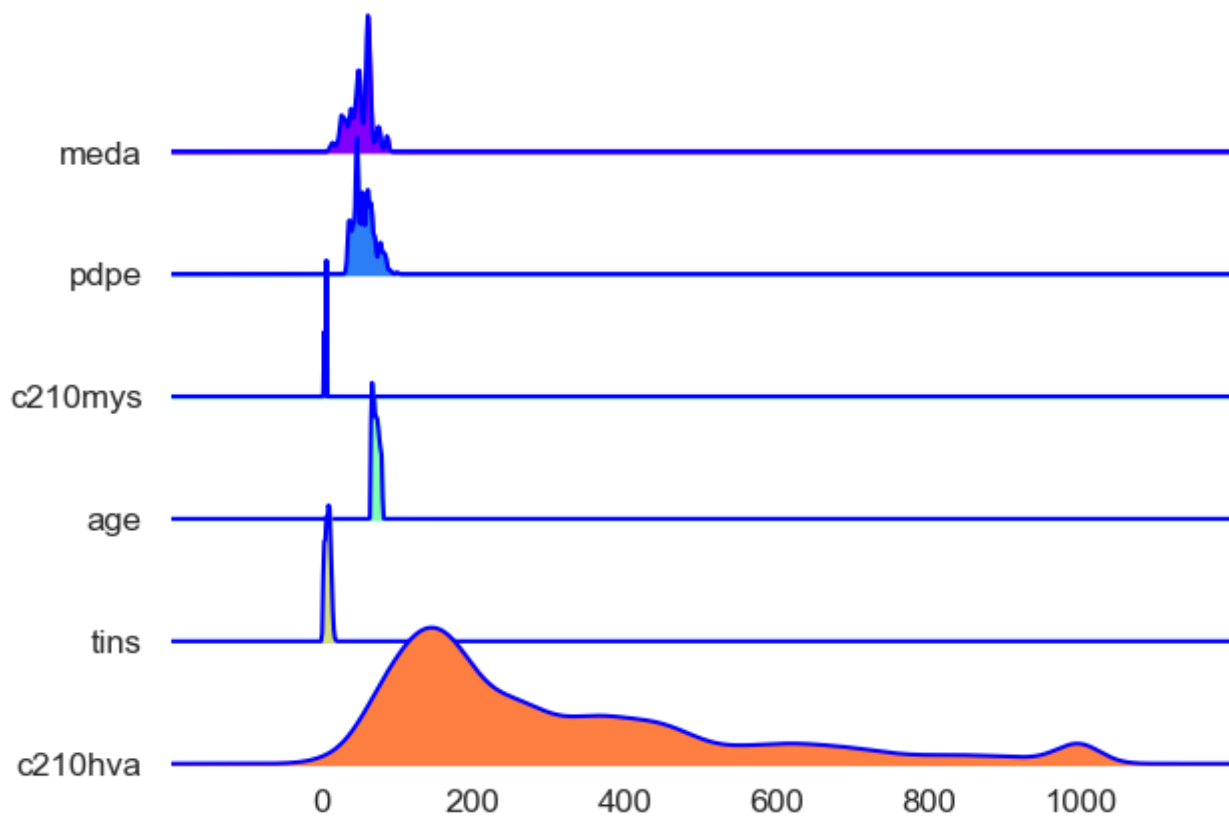
- 处理缺失值
- 查看数据类型
- 处理重复记录
- 查看标签数据的分布



第二阶段：特征处理

- 分类型变量做编码处理
- 数值型变量做标准化处理
- 分割测试集和训练集数据

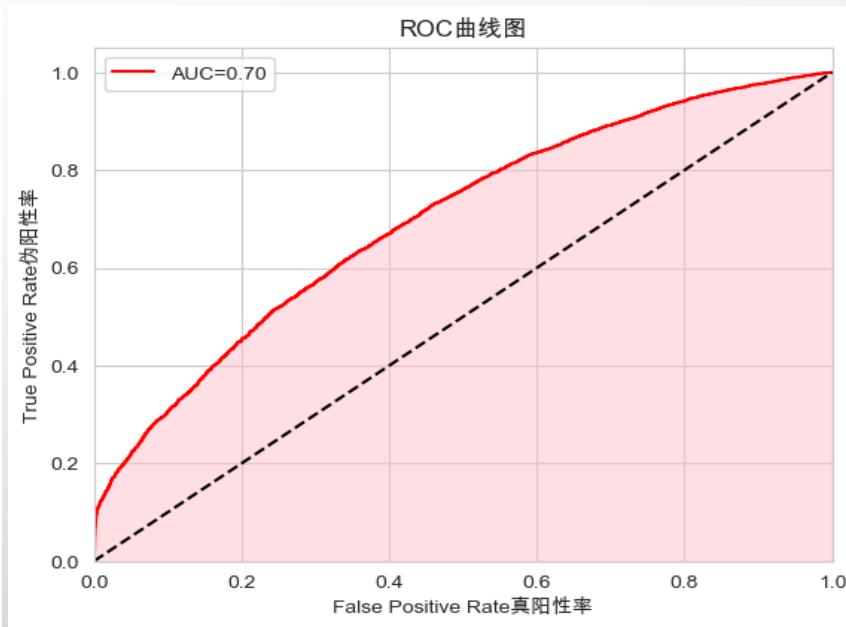
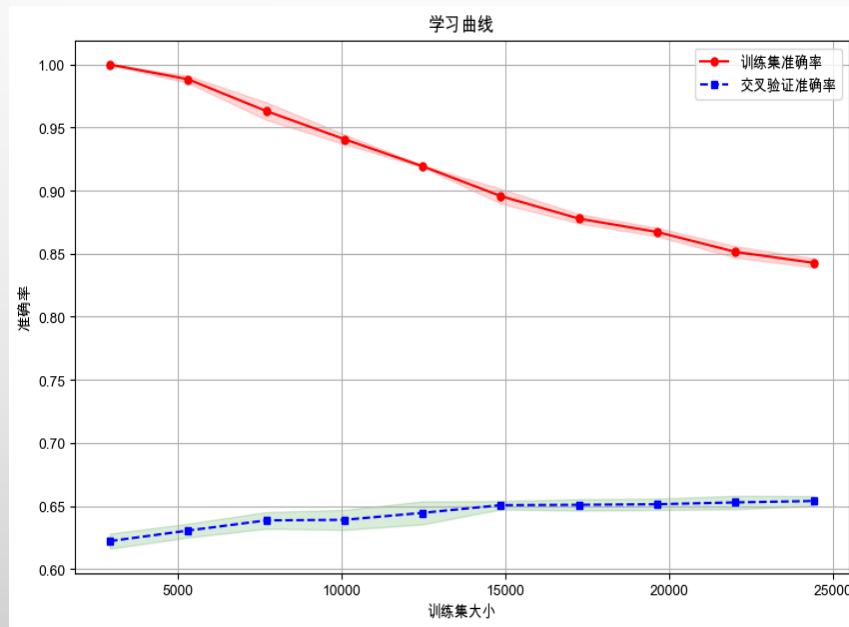
重要特征的数据分布



第三阶段：建模、模型评估

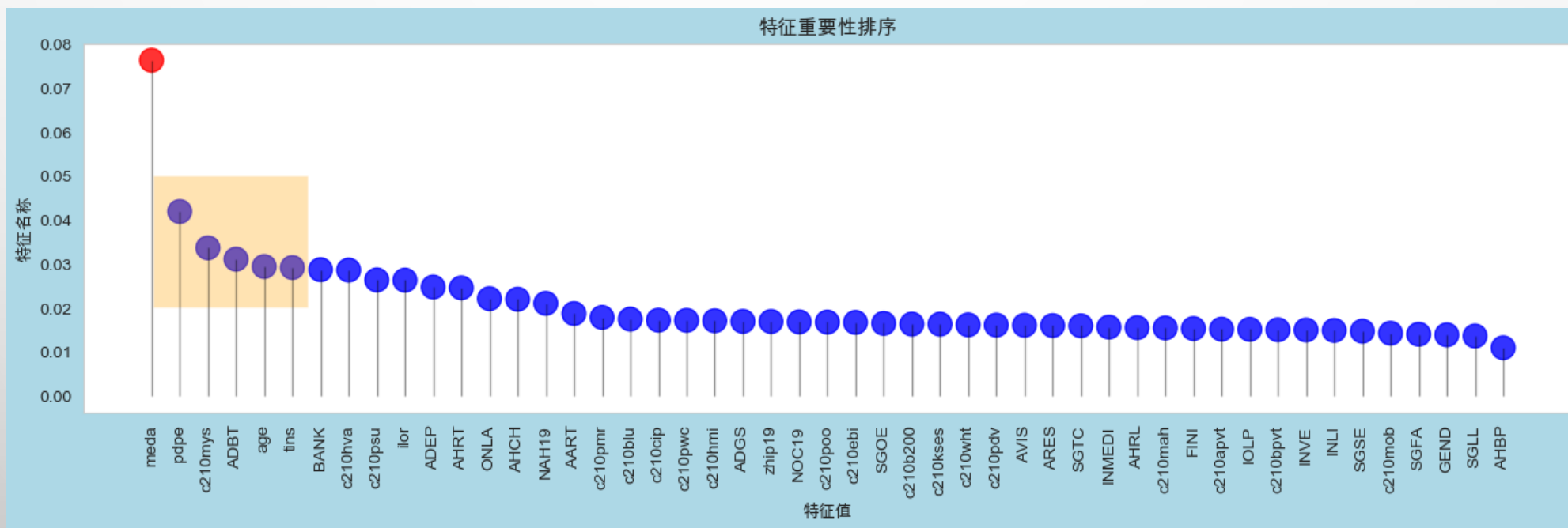
- XGBOOST模型
- ROC曲线图
- 学习曲线

	precision	recall	f1-score	support
0	0.68	0.81	0.74	7832
1	0.61	0.45	0.51	5247
accuracy				0.66 13079
macro avg	0.65	0.63	0.63	13079
weighted avg	0.65	0.66	0.65	13079



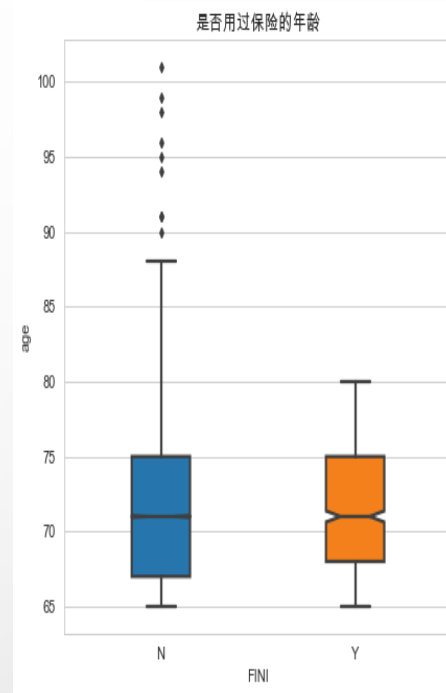
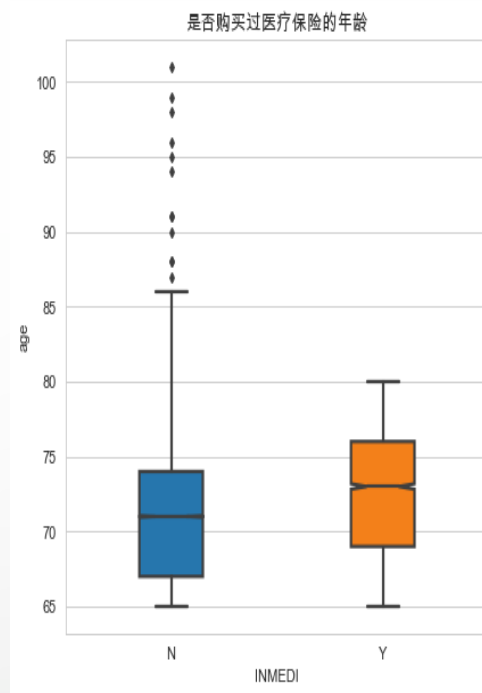
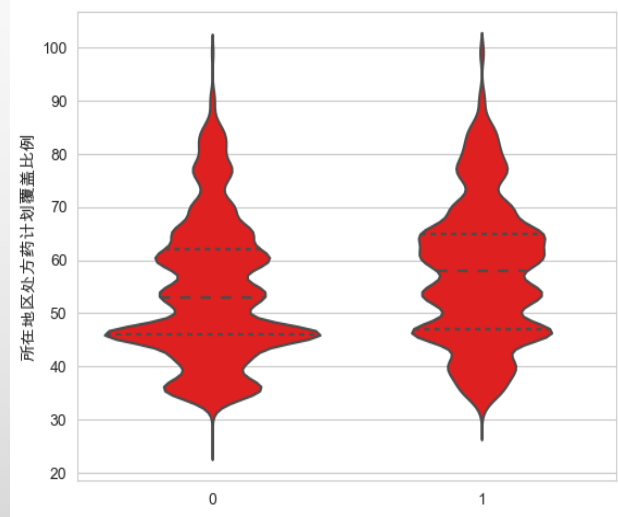
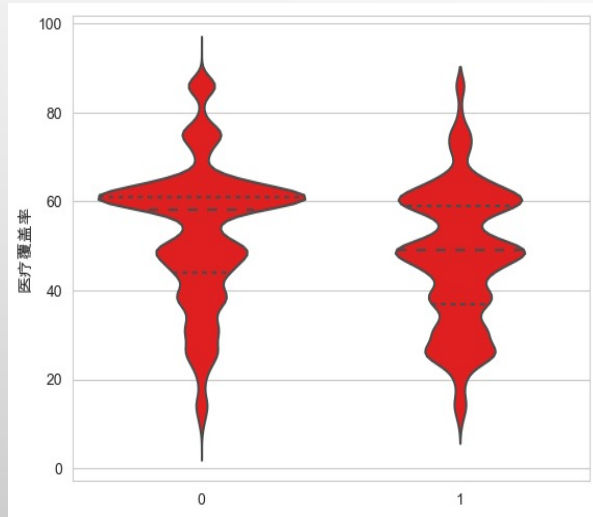
第四阶段：可视化分析

- 特征重要性的排序和选择



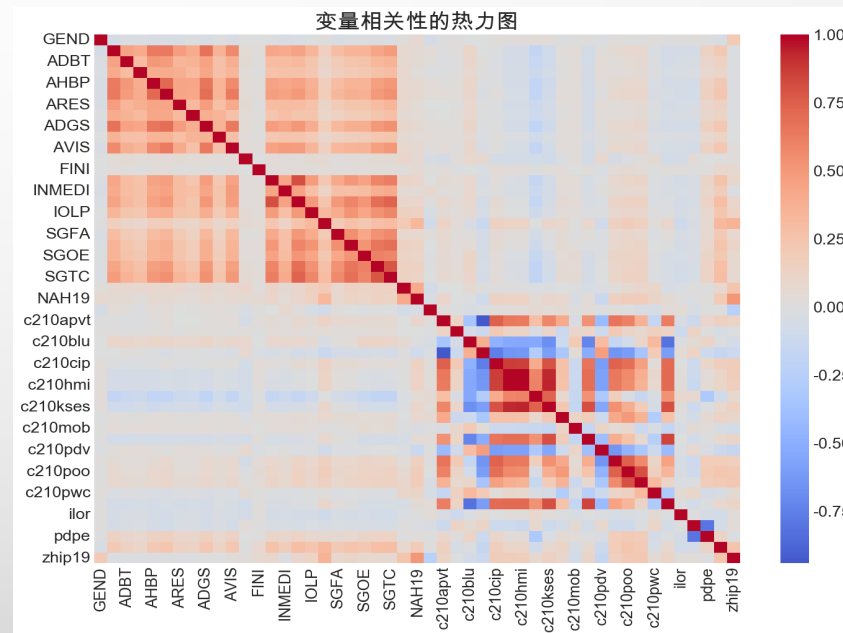
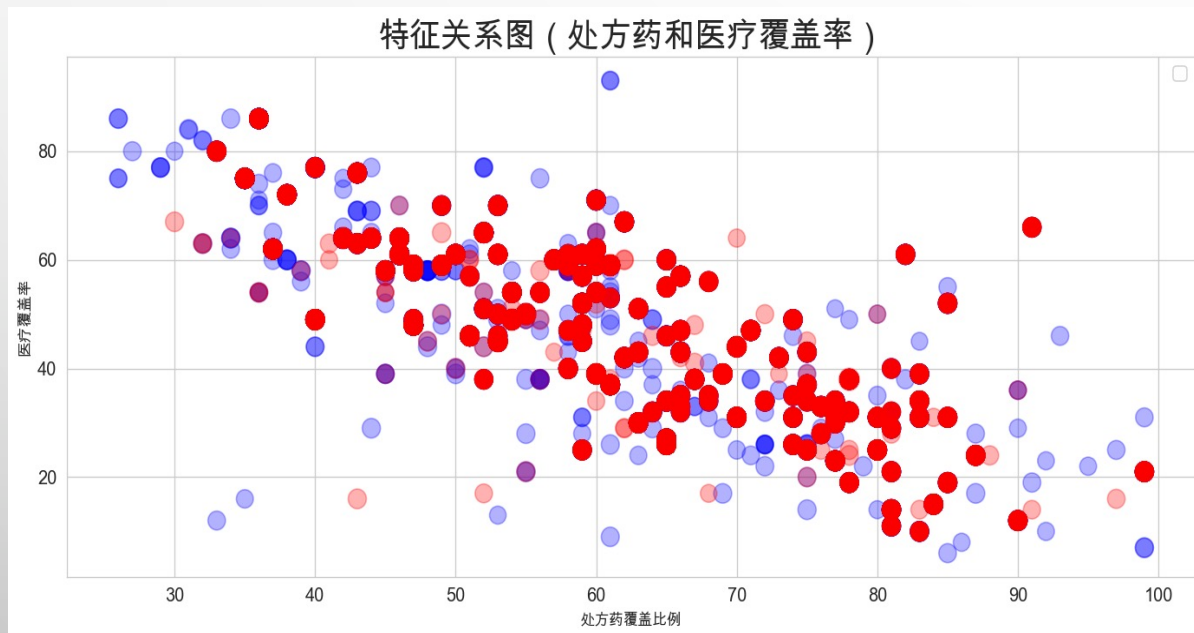
第四阶段：可视化分析

- 特征的业业务分析



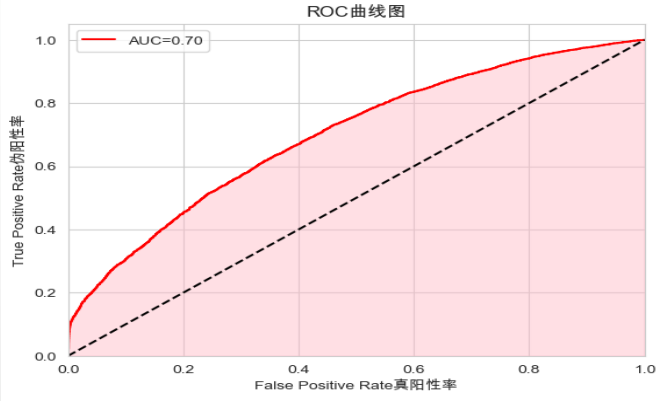
第四阶段：可视化分析

- 特征关系的分析

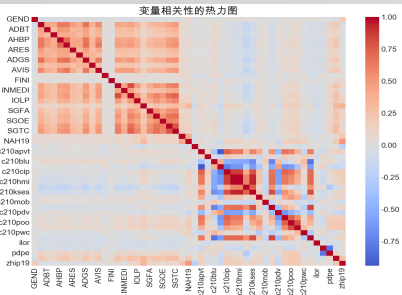
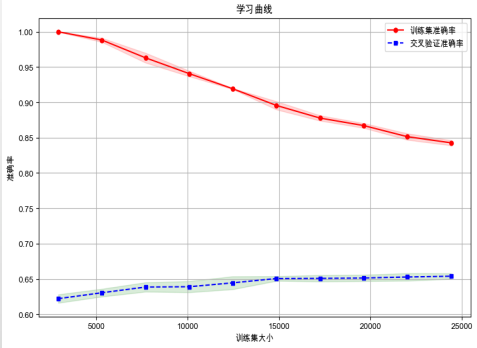


保险行业数据分析

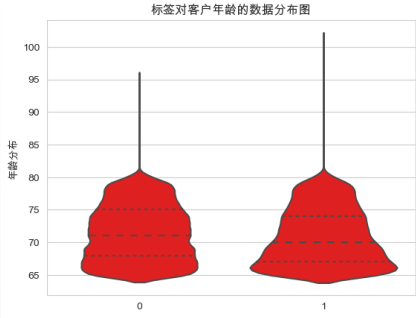
收集了四万多条记录包含个特征一列二分类的标签特征通过数据分析和可视化展示了报销售需要重点关注医疗覆盖率处方药覆盖率不同程度的地区客户学历和年龄情况另外糖尿病和破产几率也是重要的影响因素。



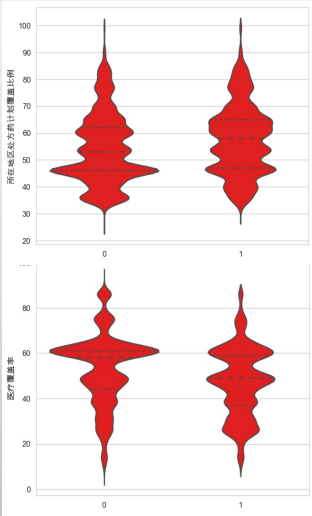
模型评估指标
AUC=0.7，因为数据经过脱敏处理。学习曲线能够看出模型的学习能力和泛化能力在不断收敛和趋于平衡，但受制于数据的质量程度有限。



	precision	recall	f1-score	support
0	0.68	0.81	0.74	7832
1	0.61	0.45	0.51	5247
accuracy	0.66			13079
macro avg	0.65	0.63	0.63	13079
weighted avg	0.65	0.66	0.65	13079

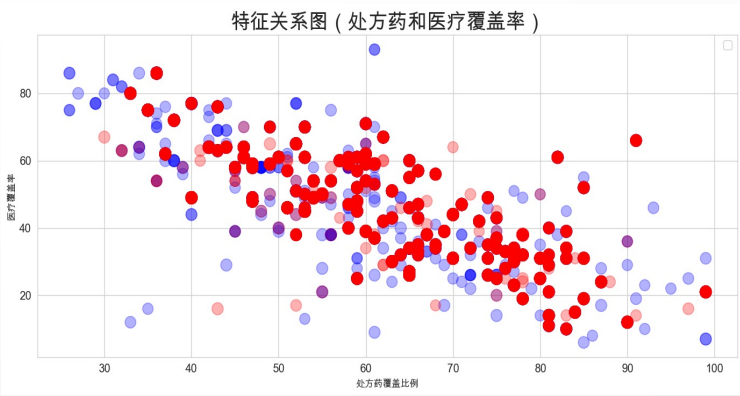
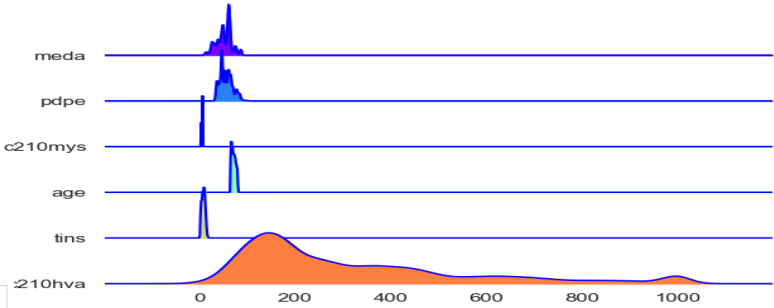


标签二分类数据分布均匀，对保险产品有反馈和没反馈的人群年龄类似。

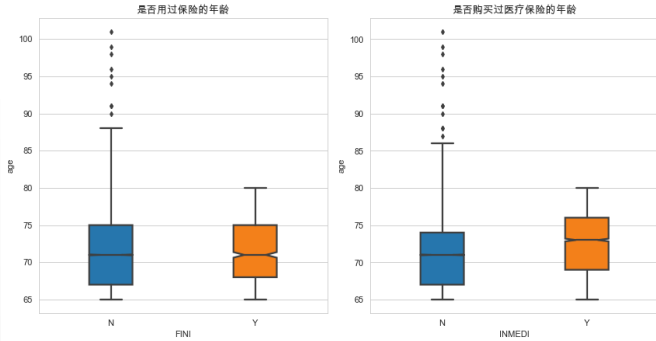


医疗覆盖率在 60% 左右对保险产品的反馈并不积极，处方药覆盖率 45% 的时候相对减少了很多对保险产品响应率。

重要特征的数据分布



购买过医疗保险的人群年龄在 68-77 岁左右，用过保险的人群年龄主要在 67-75 岁范围内。



对保险反馈最重要的特征是医疗覆盖率遥遥领先其他因素，但是对整个业务来说“处方药的覆盖率”、“学历”、“是否有糖尿病”、“是否有过破产记录”等等都是值得业务人员关注的因素。

特征重要性排序

