

好未来AI Lab

数据挖掘

图像处理

自然语言处理

语音

项目： 业务部门数据分析与预测
数据挖掘平台建设

优势： 数据丰富
迭代速度快

学生续报预测

1 数据分析

2 续报预测

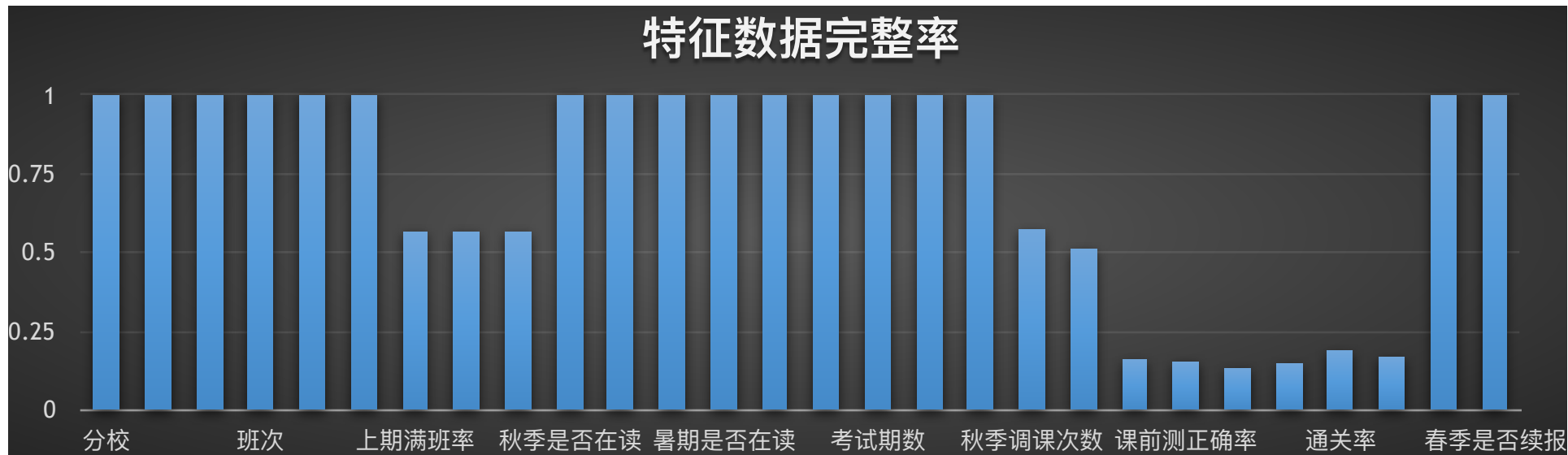
3 总结

数据分析

长期班数据： 65万个学生样本 (34维特征 、 2个待预测量： 春季和寒假是否续报)

去除不相关特征： 学员ID、班级ID、学员编码

去除数据缺失较多特征： 学员评分、课前测正确率



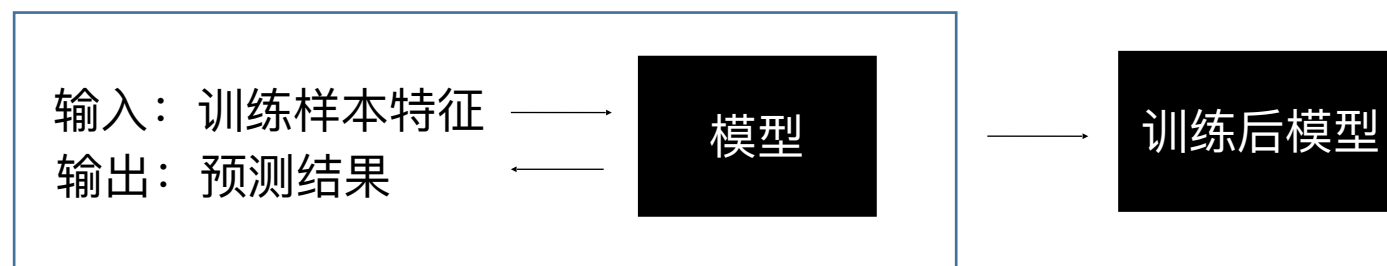
剩余20维特征：分校 年级 学科 班次 教师 是否新生 报班总课次 长期班期数....

特征预处理：异常数据消除、缺失数据填充、类别变量编码、标准化

续报预测

机器学习：

训练：



预测：



模型： 随机森林(树的数目、深度、特征数目等)

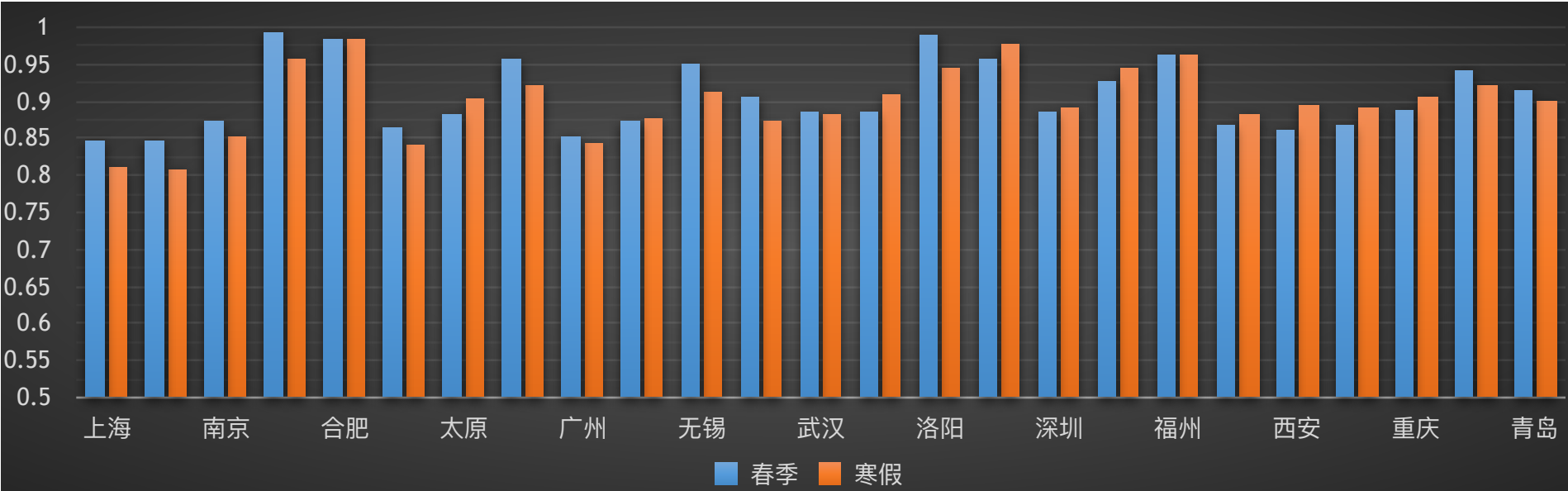
数据： 65w条学生样本（训练，验证，测试）

续报预测

春季学生续报准确率：87.7%

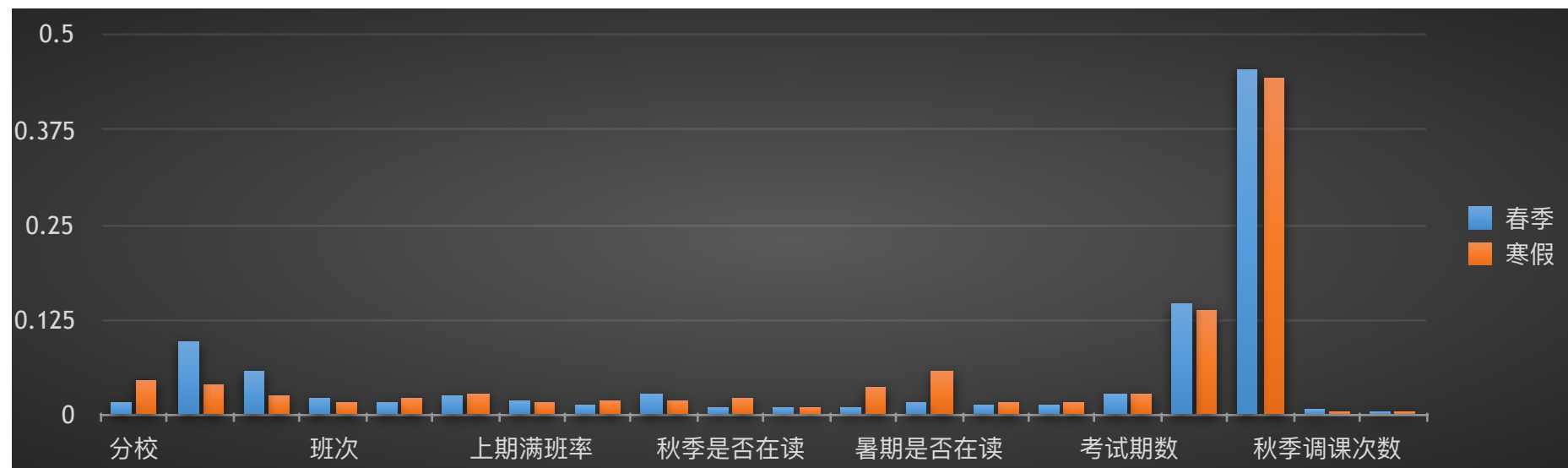
寒假学生续报准确率：86.4%

各分校学生准确度：



续报预测

特征重要性打分：



春季续报前三重要特征：长期班期数、报班总课数、年级

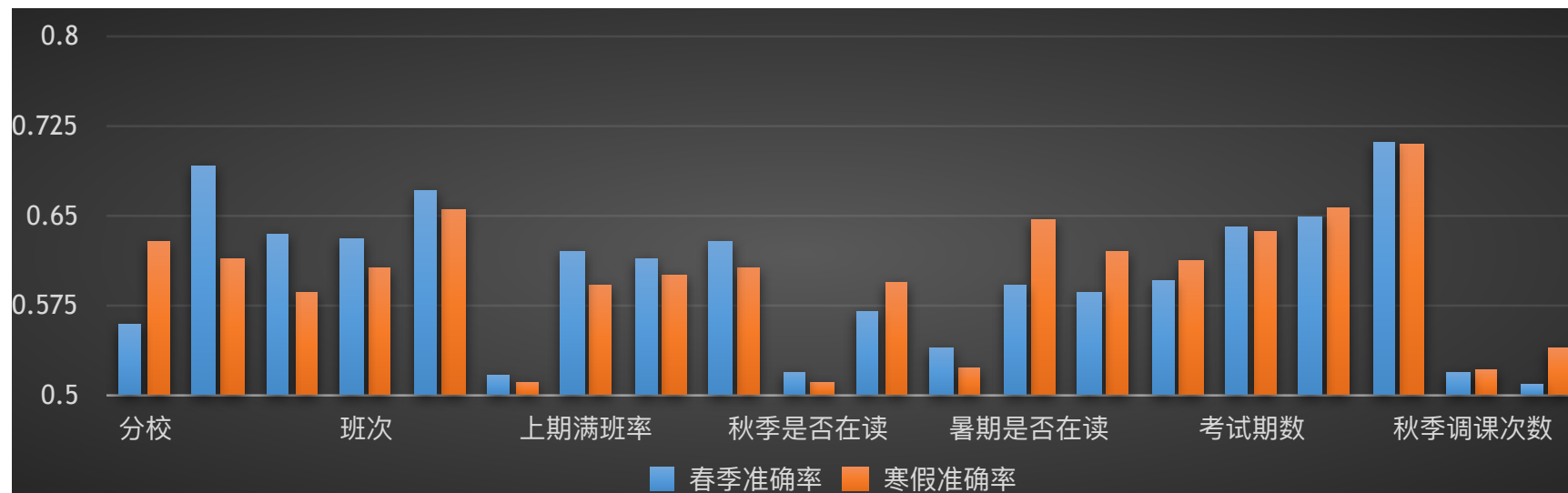
寒假续报前三重要特征：长期班期数、报班总课数、暑期是否在读

只采用前三重要特征，春季和寒假续报准确率

	春季	寒假
随机森林	85.5%	85.8%
Adaboost	85.2%	85.3%

续报预测

单个特征对应续报预测精度（每次只用一个特征）：



春季续报预测准确度前五特征：长期班期数、年级、**教师**、报班总课数、考试期数

寒假续报预测准确度前五特征：长期班期数、**教师**、暑期是否在读、报班总课数、考试期数

续报预测

每个学生续报的概率值

分校	学员ID	班级ID	学员编号	学员姓名	年份	考试期数	课程总课时	长期班期数	秋季调课	秋季转班	寒假续报预测概率	寒假是否续报	春季续报预测概率	春季是否续报
分校	0b1d0cbecc	228080815	04021118	胡奕	2016	1.0	22.0	5.0	2.0		0.866971	是	0.935164	是
分校	159ba9b68	228080815	20000621	黄二乐	2016	3.0	44.0	13.0	2.0		0.669726	是	0.908701	是
分校	17d916276	228080815	50000607	徐浩亮	2016	0.0	20.0	5.0			0.883049	是	0.944	是
分校	27a434hhr	228080815	40000608	崔晴莉	2016	0.0	2.0	4.0	6.0	4.0	0.48012	否	0.902431	是
分校	28be6d075	228080815	40000609	叶了晋	2016	1.0	74.0	11.0	4.0		0.877557	是	0.890346	是
分校	2aea207hc	228080815	50000609	倪丞佑	2016	0.0	44.0	7.0		4.0	0.828748	是	0.890512	是

不同续报(不续报)阈值对应精确率和召回率

春季			寒假	
续报阈值	精确率	召回率	精确率	召回率
P > 0.7	91.5%	45.8%	92.05%	43.7%
P > 0.9	97.4%	28.3%	98.1%	30.2%

春季			寒假	
不续报阈值	精确率	召回率	精确率	召回率
P < 0.3	92.0%	50.0%	92.8%	48.0%
P < 0.1	98.5%	31.2%	98.3%	30.8%

总结

- 1、利用机器学习模型，可以得到每个学生续报概率（准确度超过85%）。
- 2、利用机器学习模型，可以得到影响学生续报的特征的打分

谢谢