# WQD7005 DATA MINING

**ALTERNATIVE ASSESSMENT 1**

**S2121005 XU HUANDI**

**Case Study: E-Commerce Customer Behaviour Analysis**
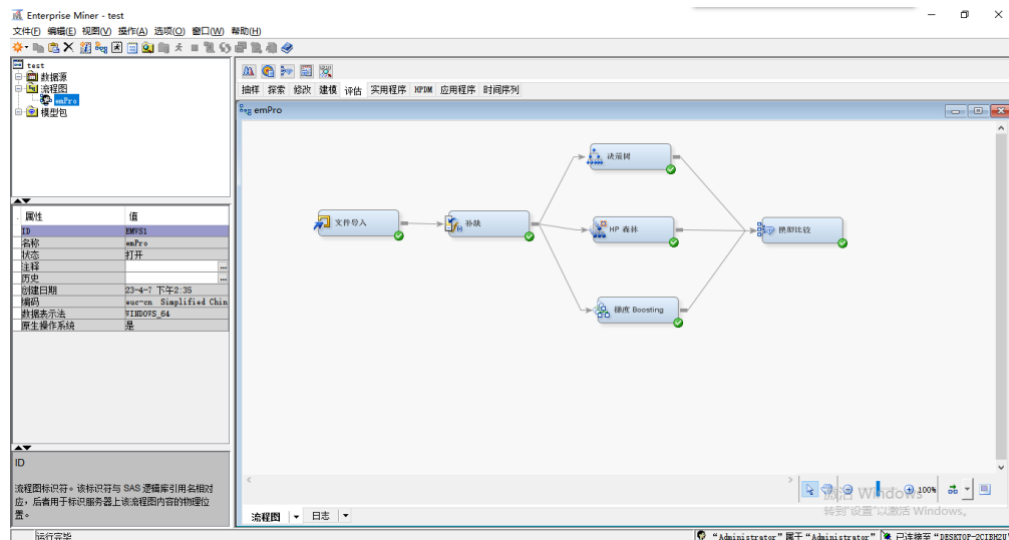**Data sources：https://www.kaggle.com**
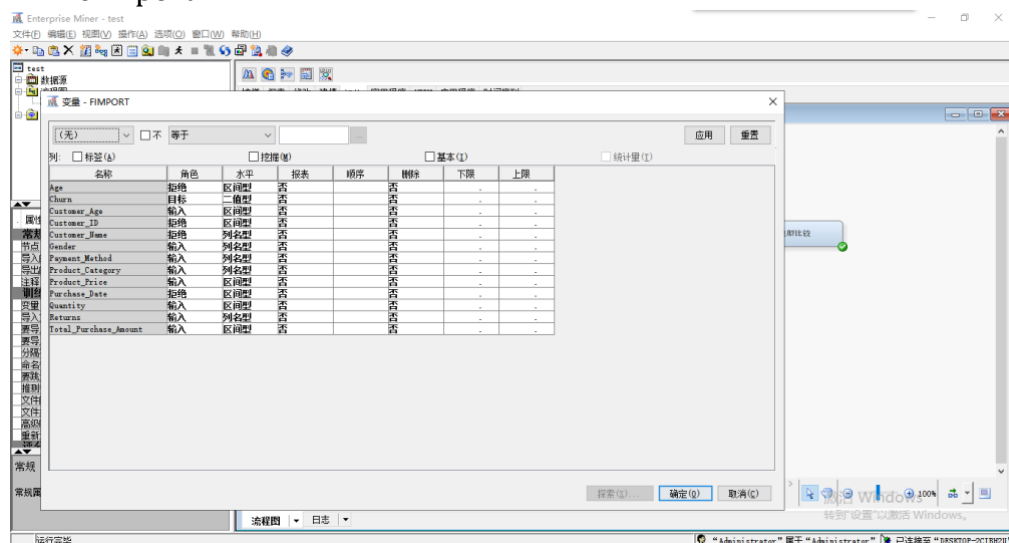**Dataset link：**
**https://drive.google.com/drive/folders/17oIfGGlBI_cSjsItHpvvMQCT6VQ7**
**eNmm?usp=sharing**

1. Data Import and Preprocessing: Import your dataset into SAS Enterprise Miner, handle missing values, and specify variable roles.
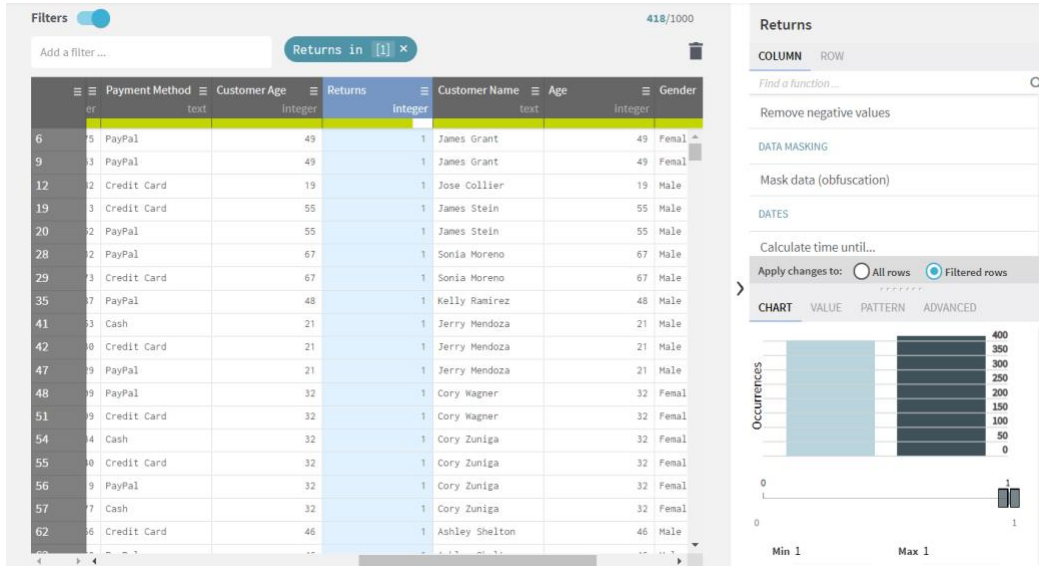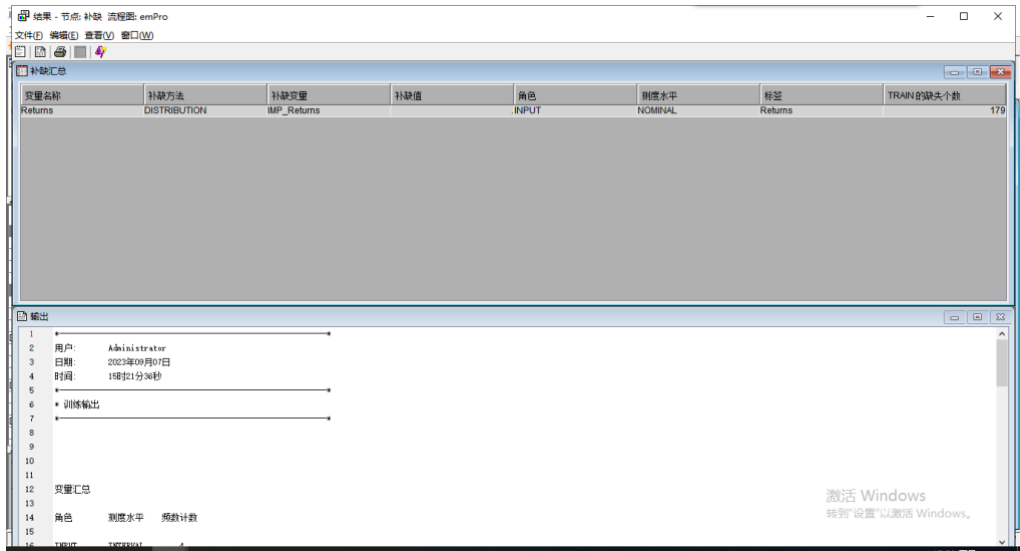   —Overall Flow Chart



   —File import

—First use talend (DP) to handle missing values



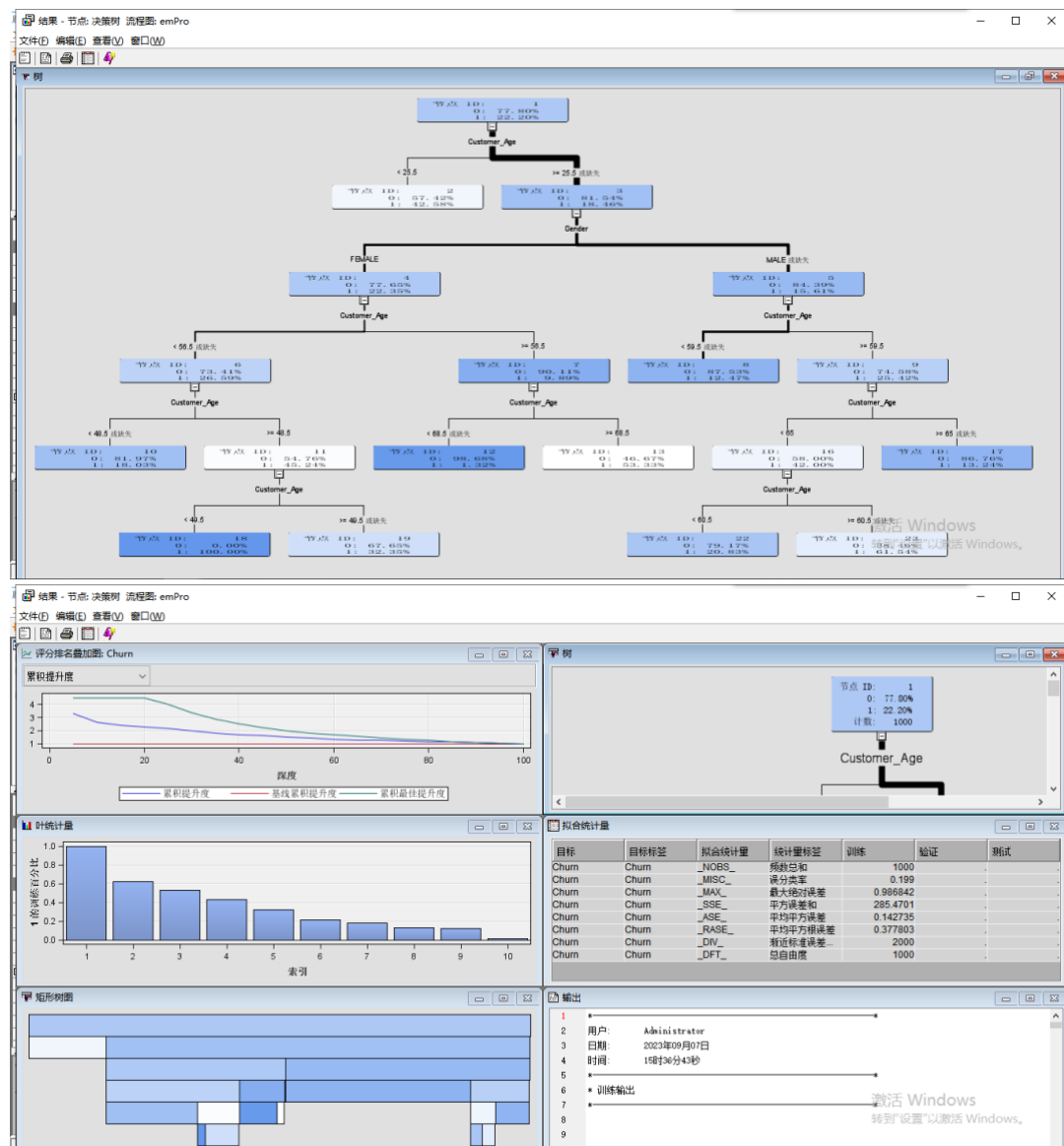—Noise and missing values are removed with SAS (EM).



In the dataset, the dependent variable (target variable) should be "Churn" because this is the outcome I want to predict. All other variables can be used as independent variables (explanatory variables) to predict whether a customer will churn.

Because there are two ages in the data. So I set up a rejection for one of them, because the name, ID, and date cannot be used for modeling, so I also set up a rejection.

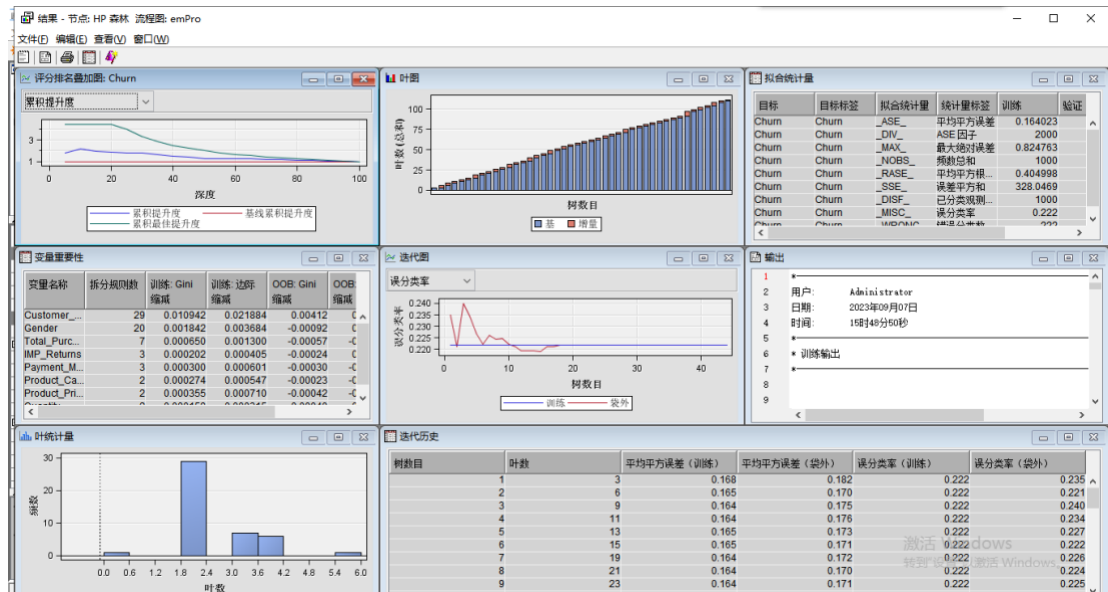| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Customer ID | Purchase Date | Product Cat| Product Pric| Quantity | Total Purch| Payment M| Customer A| Returns | Customer Name| Age | Gender | Churn |
| 2 | 46251 | 2020/9/8 09:38 | Electronics | 12 | 3 | 740 | Credit Card | 37 | 0 | Christine Hernar | 37 | Male | 0 |
| 3 | 46251 | 2022/3/5 12:56 | Home | 468 | 4 | 2739 | PayPal | 37 | 0 | Christine Hernar | 37 | Male | 0 |
| 4 | 46251 | 2022/5/23 18:18 | Home | 288 | 2 | 3196 | PayPal | 37 | 0 | Christine Hernar | 37 | Male | 0 |
| 5 | 46251 | 2020/11/12 13:13 | Clothing | 196 | 1 | 3509 | PayPal | 37 | 0 | Christine Hernar | 37 | Male | 0 |

2. Decision Tree Analysis: Create a decision tree model in SAS Enterprise Miner to analyse customer behaviour
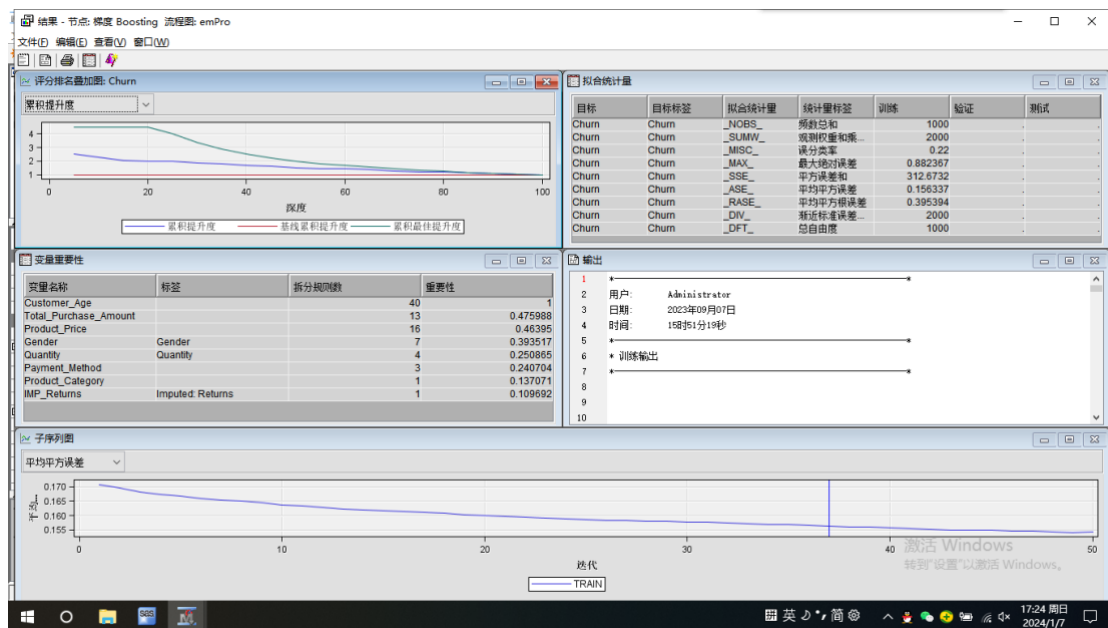   —Create a decision



The problem is that an overfitted model relies too much on specific features in the training data and therefore may perform poorly on new, unknown data because those features may not apply to the new data. On the contrary, underfitting models perform poorly on both training and test data because they do not fully learn the characteristics of the training data and do not understand the data well enough. To solve these problems, we can use techniques such as limiting the depth of the decision tree or the minimum number of samples at a node to prevent overfitting, or consider using ensemble methods such as random forest or gradient boosting, which work by combining multiple models. Improve generalization capabilities to achieve better performance on new data. These methods can effectively improve the performance of the model.

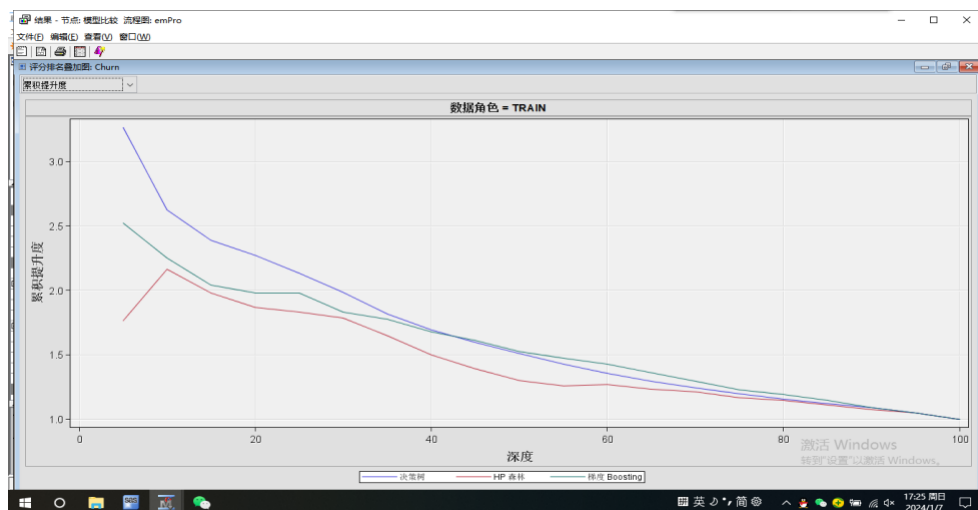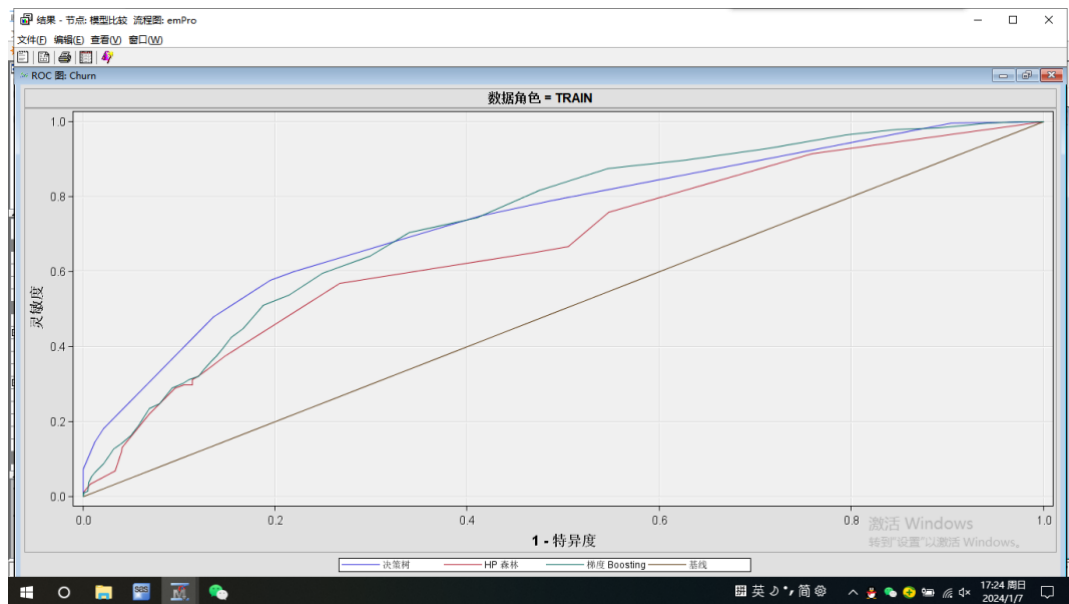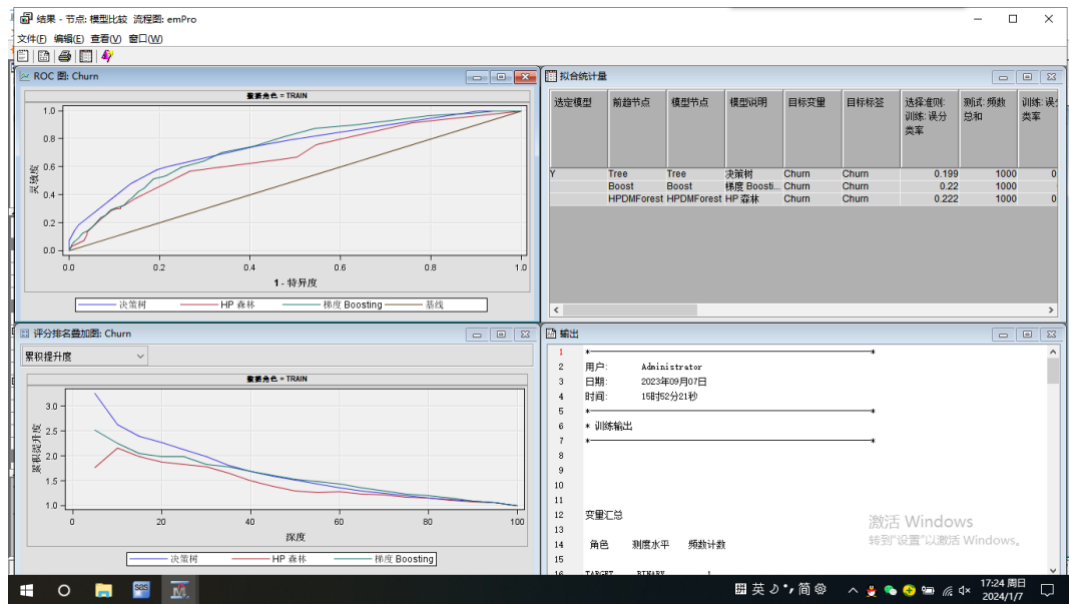3. Ensemble Methods: Apply Bagging and Boosting, using the Random Forest algorithm as a Bagging example.

—Random Forest



—Boosting

## 4. Model comparison

| 前趋节点 | 模型节点 | 模型说明 | 目标变量 | 目标标签 | 选择准则 训练: 误分类率 | 测试: 频数总和 | 训练: 误分类率 | 训练: 最大绝对误差 | 训练: 平方误差和 | 训练: 平均平方误差 | 训练: 平均平方根误差 | 训练: 渐近标准误差的预数 | 训练: 总自由度 | 训练: 已分类观测的频数 | 训练: 错误分类数 | 训练: 观测权重和频以频数 | 训练: RO索引 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Tree | Tree | 决策树 | Churn | Churn | 0.199 | 1000 | 0.199 | 0.986842 | 285.4701 | 0.142735 | 0.377803 | 2000 | 1000 | | | | 0 |
| Boost | Boost | 梯度 Boosti... | Churn | Churn | 0.22 | 1000 | 0.22 | 0.882367 | 312.6732 | 0.156337 | 0.395394 | 2000 | 1000 | | | 2000 | 0 |
| HPDMForest | HPDMForest | HP 森林 | Churn | Churn | 0.222 | 1000 | 0.222 | 0.824763 | 328.0469 | 0.164023 | 0.404998 | 2000 | | 1000 | 222 | | 0 |

SAS helped me figure out which factors were most important in predicting whether a customer would leave us. The age of our customers, how much they spend, and the price of our products are particularly critical. It's interesting that older customers seem to have different departure patterns than younger customers. And it seems that customers who spend more are more likely to stay. Gender seems to have little impact on whether a customer will churn.

The models we used, especially ensemble methods like random forests and gradient boosting, worked well, especially in terms of accuracy and generalization, compared to decision trees alone.

With this information, I think we can do this:

1. Create some holiday or specific discounts: For those customers who may be about to leave, especially young people or those who don't spend much overall, we can create some discounts or loyalty programs to make them feel more valued and increase their trust in the brand. of loyalty.

2. Refining customer groups: We can use decision tree analysis to divide customers into different groups , and customize different marketing strategies based on their shopping habits and basic information.

3. Re-examine pricing: Since product price is so important, we may need to rethink our pricing strategy to ensure that the price is both competitive and that customers feel they are worth their money.

4. More feedback and engagement: We need to understand why customers choose to leave and improve based on their feedback. This way we can better retain customers and increase their satisfaction.

Dr, I'm sorry about the Chinese in the software, because the Mac I use can only be downloaded using a virtual machine. I don't know why it turned into Chinese, but I don't have much time to download it again. I apologize to you again.