

GBDT算法甄别某电商的高质量用户

实验目标

通过本案例的学习和课后作业的练习：

1. 通过代码的实现，帮助大家快速理解机器学习整个流程；
2. 能够使用SKlearn实现GBDT算法。

你也可以将本案例相关的 ipynb 学习笔记分享到 [AI Gallery Notebook \(https://marketplace.huaweicloud.com/markets/aihub/notebook/list/\)](https://marketplace.huaweicloud.com/markets/aihub/notebook/list/) 版块获得成长值 (https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=9b8d7e7a-a150-449e-ac17-2dcf76d8b492)，分享方法请查看[此文档 \(https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=8afec58a-b797-4bf9-acca-76ed512a3acb\)](https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=8afec58a-b797-4bf9-acca-76ed512a3acb)。

案例内容介绍

梯度提升决策树（Gradient Boosting Decision Tree，GBDT）算法属于集成算法的一种，基分类器是回归树（分类问题也是回归树，最后再用sigmoid或者softmax函数计算类别），是一种boosting算法，即逐步拟合逼近真实值，是一个串行的算法，可以减少bias（误差）却不能减少variance（偏差），因为每次基本都是全样本参与训练，不能消除偶然性的影响，但每次都逐步逼近真实值，可以减少误差。

GBDT包括三种基本用法，一是回归，二是二分类，三是多分类。技术细节有点差异，但是整体思路都是一样的。

该算法是近年来被提及的比较多，这主要得益于其算法的性能，以及该算法在各类数据挖掘以及机器学习比赛中的卓越表现，有很多人对GBDT算法进行了开源代码的开发，比较火的是陈天奇的XGBoost和微软的LightGBM。

本案例推荐的理论学习视频：

- [《AI技术领域课程--机器学习》 GBDT \(https://education.huaweicloud.com/courses/course-v1:HuaweiX+CBUCNXE086+Self-paced/courseware/f56ce3bf3bc84b4e80dce98cd5dbf2ca/d765dbae0f5e4c69936fe8fe0474ebd9/\)](https://education.huaweicloud.com/courses/course-v1:HuaweiX+CBUCNXE086+Self-paced/courseware/f56ce3bf3bc84b4e80dce98cd5dbf2ca/d765dbae0f5e4c69936fe8fe0474ebd9/)

注意事项

1. 如果您是第一次使用 JupyterLab，请查看 [《ModelArts JupyterLab使用指导》](https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=03676d0a-0630-4a3f-b62c-07fba43d2857) (https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=03676d0a-0630-4a3f-b62c-07fba43d2857) 了解使用方法；
2. 如果您在使用 JupyterLab 过程中碰到报错，请参考 [《ModelArts JupyterLab常见问题解决办法》](https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=9ad8ce7d-06f7-4394-80ef-4dbf6cfb4be1) (https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=9ad8ce7d-06f7-4394-80ef-4dbf6cfb4be1) 尝试解决问题。

实验步骤

1、导入相关库

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn import datasets
%matplotlib inline
```

2、读取用户数据

```
In [2]: # 拷贝数据
import moxing as mox
import os

if not os.path.exists('data_carrier_svm.csv'):
    mox.file.copy('obs://modelarts-labs-bj4-v2/course/hwc_edu/machine_learning/datasets/GBDT/data_carrier_svm.csv',
                  'data_carrier_svm.csv')
```

INFO:root:Using MoXing-v1.17.3-

INFO:root:Using OBS-Python-SDK-3.20.7

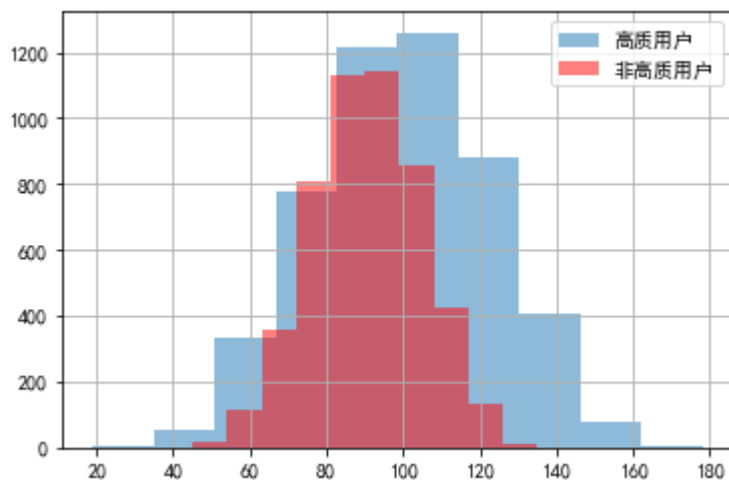
```
In [3]: data = pd.read_csv(r'data_carrier_svm.csv', encoding='gbk')
data.head()
```

Out[3]:

	用户标识符	网络类型	类目1消费金额	类目2消费金额	类目3消费金额	类目4消费金额	类目5消费金额	类目6消费金额	类目7消费金额	类目7消费金额.1	总消费次数	账户余额	是否高质用户
0	66069	3G	70	97	395	13	64	168	59	465	7	36	0
1	64410	3G	94	79	366	35	59	182	70	542	13	66	0
2	60110	3G	92	99	390	44	134	219	8	548	8	110	1
3	69600	4G	131	87	391	0	128	180	63	498	4	30	1
4	64683	4G	74	104	397	35	112	258	68	614	15	18	1

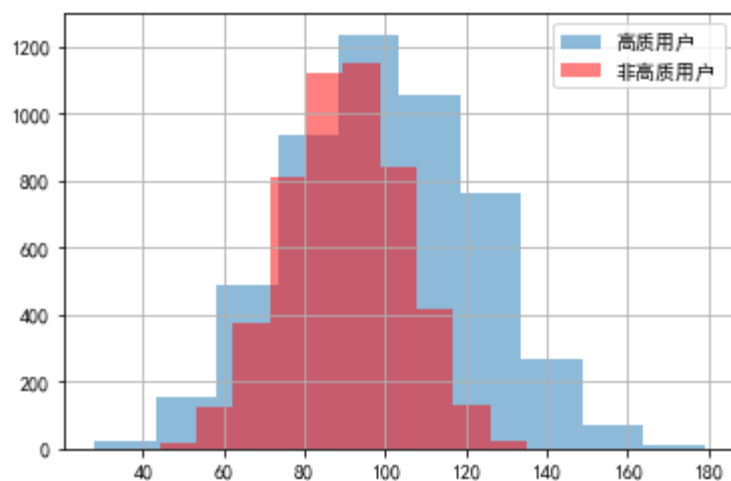
```
In [4]: # 不同用户的类目1消费金额分布情况对比
cond = data['是否高质用户'] == 1
data[cond]['类目1消费金额'].hist(alpha=0.5, label='高质用户')
data[~cond]['类目1消费金额'].hist(color='r', alpha=0.5, label='非高质用户')
plt.legend()
```

Out[4]: <matplotlib.legend.Legend at 0x7f9dfaa906a0>



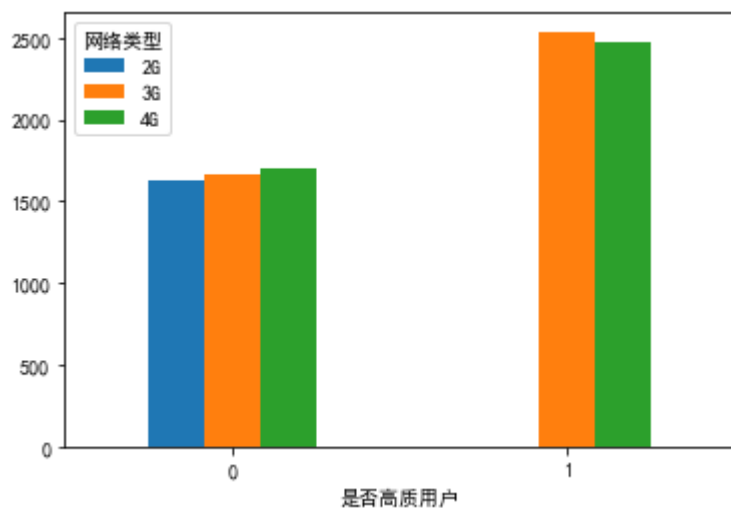
```
In [5]: # 不同用户的类目2消费金额分布情况对比
cond = data['是否高质用户'] == 1
data[cond]['类目2消费金额'].hist(alpha=0.5, label='高质用户')
data[~cond]['类目2消费金额'].hist(color='r', alpha=0.5, label='非高质用户')
plt.legend()
```

Out[5]: <matplotlib.legend.Legend at 0x7f9dfaa04e80>



```
In [6]: # 不同用户的网络类型情况对比
grouped = data.groupby(['是否高质用户', '网络类型'])['用户标识符'].count().unstack()
grouped.plot(kind='bar', alpha=1.0, rot=0)
```

Out[6]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9df8302630>

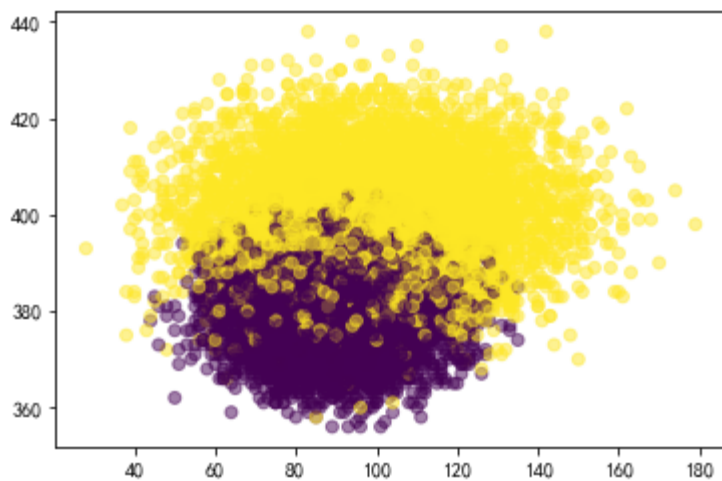


```
In [7]: data['是否高质用户'].value_counts()
```

```
Out[7]: 1    5003
        0    4997
        Name: 是否高质用户, dtype: int64
```

```
In [8]: # 生成数据可视化
y = data.loc[:, '是否高质用户']
plt.scatter(data.loc[:, '类目2消费金额'], data.loc[:, '类目3消费金额'], c=
y, alpha=0.5)
```

Out[8]: <matplotlib.collections.PathCollection at 0x7f9df5afb400>



3、数据预处理

```
In [9]: # 类别特征编码
X = data.loc[:, '网络类型':'账户余额']
y = data.loc[:, '是否高质用户']
print('The shape of X is {0}'.format(X.shape))
print('The shape of y is {0}'.format(y.shape))
```

The shape of X is (10000, 11)

The shape of y is (10000,)

```
In [10]: X.head()
```

Out[10]:

	网络 类型	类目1消 费金额	类目2消 费金额	类目3消 费金额	类目4 消费金 额	类目5消 费金额	类目6消 费金额	类目7 消费金 额	类目7消 费金额.1	总消 费次 数	账户 余额
0	3G	70	97	395	13	64	168	59	465	7	36
1	3G	94	79	366	35	59	182	70	542	13	66
2	3G	92	99	390	44	134	219	8	548	8	110
3	4G	131	87	391	0	128	180	63	498	4	30
4	4G	74	104	397	35	112	258	68	614	15	18

```
In [11]: from sklearn.preprocessing import OneHotEncoder

def service_mapping(cell):
    if cell == '2G':
        return 2
    elif cell == '3G':
        return 3
    elif cell == '4G':
        return 4

# 将网路类型的string型值映射为整数型
service_map = X['网络类型'].map(service_mapping)
service = pd.DataFrame(service_map)
# service_df
# 使用OneHotEncoder转化类型特征为0/1编码的多维特征
enc = OneHotEncoder()
service_enc = enc.fit_transform(service).toarray()
service_enc

# 0/1编码的多维特征的名称
service_names = enc.active_features_.tolist()
service_newname = [str(x) + 'G' for x in service_names]

service_df = pd.DataFrame(service_enc, columns=service_newname)
service_df.head()
X_enc = pd.concat([X, service_df], axis=1).drop('网络类型', axis=1)
X_enc.head()
```

```
/home/ma-user/anaconda3/envs/XGBoost-Sklearn/lib/python3.6/site-packages/sklearn/preprocessing/_encoders.py:363: FutureWarning: The handling of integer data will change in version 0.22. Currently, the categories are determined based on the range [0, max(values)], while in the future they will be determined based on the unique values.
```

If you want the future behaviour and silence this warning, you can specify "categories='auto'".

In case you used a LabelEncoder before this OneHotEncoder to convert the categories to integers, then you can now use the OneHotEncoder directly.

```
warnings.warn(msg, FutureWarning)
```

```
/home/ma-user/anaconda3/envs/XGBoost-Sklearn/lib/python3.6/site-packages/sklearn/utils/deprecation.py:77: DeprecationWarning: Function active_features_ is deprecated; The ``active_features_`` attribute was deprecated in version 0.20 and will be removed 0.22.
```

```
warnings.warn(msg, category=DeprecationWarning)
```

Out[11]:

	类目1 消费金额	类目2 消费金额	类目3 消费金额	类目4 消费金额	类目5 消费金额	类目6 消费金额	类目7 消费金额	类目7 消费金额.1	总消费次数	账户余额	2G	3G	4G
0	70	97	395	13	64	168	59	465	7	36	0.0	1.0	0.0
1	94	79	366	35	59	182	70	542	13	66	0.0	1.0	0.0
2	92	99	390	44	134	219	8	548	8	110	0.0	1.0	0.0
3	131	87	391	0	128	180	63	498	4	30	0.0	0.0	1.0
4	74	104	397	35	112	258	68	614	15	18	0.0	0.0	1.0

4、训练准备过程

```

In [12]: from sklearn.model_selection import cross_val_score
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.ensemble import GradientBoostingClassifier
from xgboost import XGBClassifier
from sklearn import metrics

# 分割训练集和测试集
X_train, X_test, y_train, y_test = train_test_split(X_enc, y, test_size=0.2, random_state=112)
print('The shape of X_train is {0}'.format(X_train.shape))
print('The shape of X_test is {0}'.format(X_test.shape))
print('The shape of y_train is {0}'.format(y_train.shape))
print('The shape of y_test is {0}'.format(y_test.shape))

# 生成数据可视化
plt.scatter(X_train.iloc[:, 0], X_train.iloc[:, 1], c=y_train, alpha=0.5)

# 交叉验证
def modelfit(alg, X_train, y_train, performCV=True, printFeatureImportance=True, cv_folds=5):
    alg.fit(X_train, y_train)
    # Predict training set:
    train_predictions = alg.predict(X_train)
    train_predprob = alg.predict_proba(X_train)[:, 1]
    # Perform cross-validation: here the author calculate cross-validated AUC
    if performCV:
        cv_score = cross_val_score(alg, X_train, y_train, cv=cv_folds, scoring='roc_auc')
        # Print model report:
        print("\nModel Report")
        print("Accuracy (Train): %3.4f" % metrics.accuracy_score(y_train.values, train_predictions))
        ## IMPORTANT: first argument is true values, second argument is predicted probabilities
        print("AUC Score (Train): %f" % metrics.roc_auc_score(y_train, train_predprob))

    if performCV:
        print("CV Score : Mean - %.7g | Std - %.7g | Min - %.7g | Max - %.7g" \
              % (np.mean(cv_score), np.std(cv_score), np.min(cv_score), np.max(cv_score)))

    # Print Feature Importance:
    if printFeatureImportance:
        feat_imp = pd.Series(alg.feature_importances_, X_train.columns.tolist()).sort_values(ascending=True)
        feat_imp.plot(kind='barh', title='Feature Importances')
        plt.ylabel('Feature Importance Score')
        _ = plt.xlabel('Relative importance')

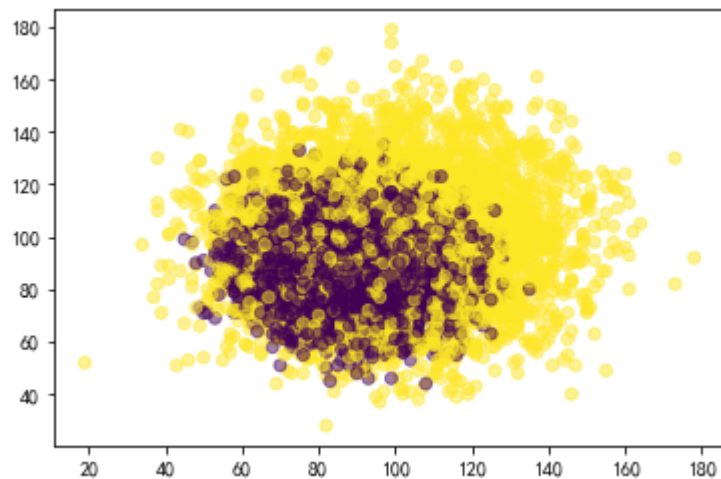
```


The shape of X_train is (8000, 13)

The shape of X_test is (2000, 13)

The shape of y_train is (8000,)

The shape of y_test is (2000,)



5、GBDT模型训练及预测

```
In [13]: from pylab import mpl

mpl.rcParams['font.sans-serif'] = ['SimHei'] # 指定默认字体
mpl.rcParams['axes.unicode_minus'] = False # 解决保存图像是负号 '-' 显示为方块的问题

# 模型实例化
clf0 = GradientBoostingClassifier(random_state=110)
# 在训练集上训练模型
clf0.fit(X_train, y_train)

# 在测试集上预测
y_pred = clf0.predict(X_test)

# 计算准备：率
score = metrics.accuracy_score(y_test, y_pred)
print('The accuracy score of the model is: {0}'.format(score))

# 查看混淆举证
metrics.confusion_matrix(y_test, y_pred)

# #模型实例化
clf0 = GradientBoostingClassifier(random_state=110)

modelfit(clf0, X_train, y_train)
```

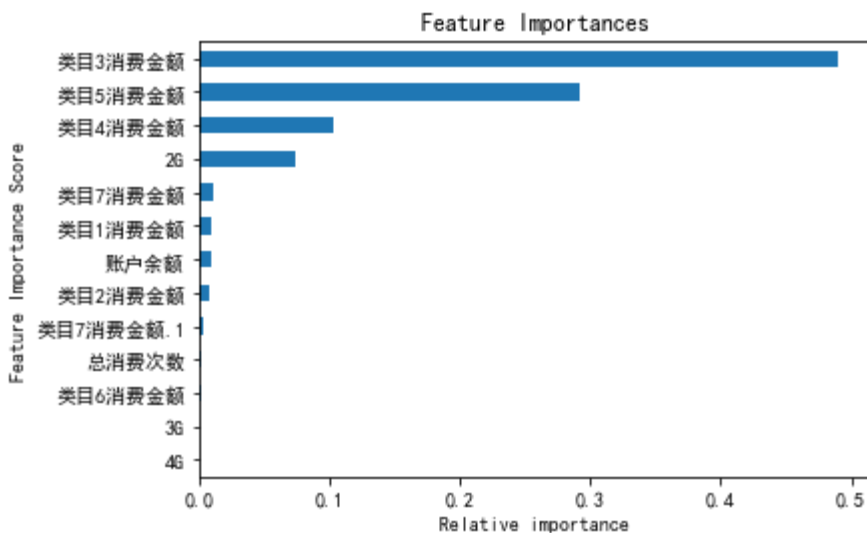
The accuracy score of the model is: 0.9825

Model Report

Accuracy (Train): 0.9862

AUC Score (Train): 0.999099

CV Score : Mean - 0.9975603 | Std - 0.0004027038 | Min - 0.9969856 |
Max - 0.9979951



小结

本实验基于脱敏后的某电商用户消费数据，使用GBDT算法根据用户的行为特征帮助甄别是否高质量用户，可以看到模型的分类准确率可以达到0.9825，可以根据实际业务指标来判断模型是否达到了需要，可以通过后续的多种策略（超参数选择，模型剪枝等）进行模型性能的提升。

以上是 GBDT 的实现方法，受限于篇幅原因，本案例未完全覆盖 GBDT 的全部操作，欢迎你将更全面的 GBDT 学习笔记分享到 [AI Gallery Notebook \(https://marketplace.huaweicloud.com/markets/aihub/notebook/list/\)](https://marketplace.huaweicloud.com/markets/aihub/notebook/list/) 版块获得成长值 (https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=9b8d7e7a-a150-449e-ac17-2dcf76d8b492)，分享方法请查看[此文档 \(https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=8afec58a-b797-4bf9-acca-76ed512a3acb\)](https://marketplace.huaweicloud.com/markets/aihub/article/detail/?content_id=8afec58a-b797-4bf9-acca-76ed512a3acb)。

作业

请你利用本实验中学到的知识点，完成以下编程题：

1. 请你尝试修改 `GradientBoostingClassifier()` 函数的 `n_estimators`（基学习器个数）参数的不同取值，看看该参数的修改对模型会有怎样的影响。 [\(https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=4ce69a2b-e194-46d1-8ea5-4cf92746d912\)](https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=4ce69a2b-e194-46d1-8ea5-4cf92746d912)
2. 请你尝试修改 `GradientBoostingClassifier()` 函数的 `max_depth`（最大深度）参数的不同取值，看看该参数的修改对模型会有怎样的影响。 [\(https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=284864aa-9262-4ce9-9316-6988f67392f8\)](https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=284864aa-9262-4ce9-9316-6988f67392f8)
3. 请你尝试修改 `GradientBoostingClassifier()` 函数的所有可调参数的不同取值，看看不同参数的不同取值组合，对模型会有怎样的影响。 [\(https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=004a94cf-d628-42e7-b745-da5a47082eb5\)](https://marketplace.huaweicloud.com/markets/aihub/notebook/detail/?id=004a94cf-d628-42e7-b745-da5a47082eb5)