导入所需模块 In []: import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt from sklearn model_selection import cross_val_score from sklearn model_selection import train_tost_onlit
from sklearn.model_selection import train_test_split from sklearn import linear_model In []: data = pd.read_csv('Admission_Predict.csv') # 查看 data 中的数据 data.head() Out[]: Serial No. GRE Score TOEFL Score University Rating SOP LOR CGPA Research Chance of Admit
0 1 337 118 4 4.5 4.5 9.65 1 0.92 1 2 324 107 4 4.0 4.5 8.87 1 0.76 2 3 316 104 3 3.0 3.5 8.00 1 0.72 3 4 322 110 3 3.5 2.5 8.67 1 0.80 4 5 314 103 2 2.0 3.0 8.21 0 0.65
生成描述性统计数据
std 115.614301 11.473646 6.069514 1.143728 1.006869 0.898478 0.596317 0.498362 0.142609 min 1.000000 290.000000 92.000000 1.000000 1.000000 6.800000 0.000000 0.340000 25% 100.750000 308.00000 103.00000 2.500000 3.00000 8.170000 0.00000 0.640000 50% 200.500000 317.00000 107.00000 3.50000 3.50000 8.610000 1.000000 0.730000 75% 300.250000 325.00000 112.00000 4.00000 4.000000 9.062500 1.000000 0.830000
max 400.000000 340.000000 120.000000 5.00000 5.00000 9.92000 1.000000 0.970000 In []: # 查看 data 数据中是否有缺失数据 data.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 400 entries, 0 to 399</class>
Data columns (total 9 columns): # Column
6 CGPA 400 non-null float64 7 Research 400 non-null int64 8 Chance of Admit 400 non-null float64 dtypes: float64(4), int64(5) memory usage: 28.2 KB In []: # 取出相关字段, 并生成相关性的可视化图 field = ['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR', 'CGPA', 'Research', 'Chance of Admit'] sns.heatmap(data[field].corr(), annot=True, vmin=-1, vmax=1)
Out[]: - 1.00 GRE Score - 1
University Rating - 0.67
CGPA - 0.83
GRE Score - TOEFL Score - Sop - Chance of Admit - Chance of Admit -
In []: # GRE成绩与入学承认率的关系 sns.regplot(x='GRE Score', y='Chance of Admit', data=data) Out[]: AxesSubplot:xlabel='GRE Score', ylabel='Chance of Admit'>
0.9 - 0.8 - 19 0.7 -
0.7 - Upg 0.7 - 0.6 - 0.5 - 0.4 - 0.
290 300 310 320 330 340 GRE Score In []: # 托福分数与入学承认率的关系 sns.regplot(x='TOEFL Score', y='Chance of Admit', data=data)
Out[]: <axessubplot:xlabel='toefl score',="" ylabel="Chance of Admit "> 1.0</axessubplot:xlabel='toefl>
0.8 - 0.7 - 0.6 - 0.6 -
0.5 - 0.4 -
95 100 105 110 115 120 TOEFL Score En []: # 大学评分与入学承认率的关系 sns.boxplot(x='University Rating', y='Chance of Admit', data=data) Out[]: <axessubplot:xlabel='university rating',="" ylabel="Chance of Admit"></axessubplot:xlabel='university>
Chance of Admit
0.5 0.4 1 2 3 4 5 University Rating
In []: # 算目的陈泚与入学承认率的关系 sns.boxplot(x='SOP', y='Chance of Admit ', data=data) Out[]: AxesSubplot:xlabel='SOP', ylabel='Chance of Admit '> 1.0
General Control of the control of th
0.4 - 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 SOP In []: # 本科GPA与入学承认率的关系
sns.regplot(x='CGPA', y='Chance of Admit ', data=data) Out[]: <axessubplot:xlabel='cgpa', ylabel="Chance of Admit "> 1.0 -</axessubplot:xlabel='cgpa',>
0.9 - 10.8 - 40 0.7 - 90 0.6 - 0.6 -
0.5 - 0.4 -
0.3 -
0.9
Chance of Admit
0.5 - 0.4 -
通过可视化,上述 field 的字段中与 'Chance of Admit' 的都有一定的相关性,且相关性均超过0.5,即上述特征均考虑在预测模型中 In []: # 用迭代增加特征的模型优化对特征的选择 Im = linear_model.LinearRegression() features = ['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR ', 'CGPA', 'Research'] y = data['Chance of Admit ']
<pre>selected_features = [] rest_features = features[:] best_acc = -0x7f7f7f7f while len(rest_features)>0: temp_best_i = '' temp_best_acc = -0x7f7f7f7f for feature_i in rest_features: temp_features = selected_features + [feature_i] X = data[temp_features]</pre>
scores = cross_val_score(lm, X, y, cv=5, scoring='neg_mean_absolute_error') # 负均万差来作为分数评价标准 acc = np.mean(scores) if acc > temp_best_acc: temp_best_ac = acc temp_best_i = feature_i print("select", temp_best_i, "acc:", temp_best_acc) if temp_best_acc > best_acc best_acc = temp_best_acc selected_features += [temp_best_i]
rest_features.remove(temp_best_i) else: break print("best feature set: ",selected_features,"acc: ",best_acc) # 因为用交叉检验得到的 scores 是负值且结果均大于负无穷 # 所以 temp_best_acc 和 best_acc 的起始值设置为负无穷
select CGPA acc: -0.052544478165713907 select GRE Score acc: -0.05074151205914499 select LOR acc: -0.04893266934603155 select TOEFL Score acc: -0.04847770855315416 select Research acc: -0.048381681109844143 select SOP acc: -0.04842869515572236 best feature set: ['CGPA', 'GRE Score', 'LOR', 'TOEFL Score', 'Research'] acc: -0.048381681109844143 通过上述模型,我们能得到 best_feature = ['CGPA', 'GRE Score', 'LOR', 'TOEFL Score', 'Research']
n []: # 对原始数据进行测试集和训练集的划分 best_features = ['CGPA', 'GRE Score', 'LOR', 'TOEFL Score', 'Research'] # 对原始数据进行测试集和训练集的划分 char_data = data[best_features] X_train, X_test, y_train, y_test = train_test_split(char_data, data['Chance of Admit'], train_size=0.8) # 建立线性回归模型
model = lm.fit(X_test, y_test) a = list(zip(best_features, model.coef_)) b = model.intercept_ print(f'回归系数: {a}', f'截距: {b}') 回归系数: [('CGPA', 0.19679873470548737), ('GRE Score', -0.0013725299837086297), ('LOR ', 0.008505992924868566), ('TOEFL Score', 0.0015098297196899512), ('Research', 0.04404154072065522)] 截距: -0.749817501794831 n []: # 评估模型的精确度 score = model.score(X_test, y_test) print(score)
0.6 -
0.5 -
plt.plot(range(len(y_pred)),y_test,'r',label="test") plt.legend() Out[]: <matplotlib.legend.legend 0x1d13774b5b0="" at=""> 1.0 -</matplotlib.legend.legend>
0.8 - 0.7 -
0.6 - 0.5 - 0.4 — predict — test
test
from sklearn.tree import DecisionTreeClassifier, export_graphviz from sklearn.model selection import train_test_split from sklearn.metrics import accuracy_score import pandas as pd import numpy as np import graphviz import os
对承认机会进行人工分类,超过0.8的为录取,0.5-0.8之间的为待考虑,低于0.5为不录取 data['Admit_level'] = np.nan # 创建Admit_level例并初始化为NaN data['Admit_level'] = data['Admit_level'].astype(str) data.iloc[data['Chance of Admit '] >= 0.8, data.columns.get_loc('Admit_level')] = 'Admission' data.iloc[(data['Chance of Admit '] >= 0.5) & (data['Chance of Admit '] < 0.8), data.columns.get_loc('Admit_level')] = 'Pending consideration' data.iloc[data['Chance of Admit '] < 0.5, data.columns.get_loc('Admit_level')] = 'Not admitted' print(data['Admit_level'].unique()) features = ['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR ', 'CGPA', 'Research']
reatures = ['GRE Score', 'University Rating', 'SOP', 'LOR', 'CGPA', 'Research'] y = data['Admit_level'] X = data[features] # 选择正确的特征数据集 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) # 随机种子 dt_classifier = DecisionTreeClassifier(random_state=42) dt_classifier.fit(X_train, y_train) y_pred = dt_classifier.predict(X_test) accuracy = accuracy_score(y_test, y_pred) print(f'Accuracy: {accuracy}') import os
指定 Graphviz 可执行文件路径 os.environ["PATH"] += os.pathsep + 'D:\ProgramData\Anaconda3\Lib\site-packages\Graphviz-10.0.1-win64\bin' dot_data = export_graphviz(dt_classifier, out_file=None, feature_names=features,
<pre>feature_names=features, class_names=['Admission', 'Pending consideration', 'Not admitted'], filled=True, rounded=True, special_characters=True,) graph = graphviz.Source(dot_data)</pre>
<pre>engine='dot' graph.render("decision_tree") graph.view() ['Admission' 'Pending consideration' 'Not admitted'] Accuracy: 0.8 'decision_tree.pdf'</pre>