

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
6 plt.style.use("seaborn")

In [2]: 1 csvpath='/home/ubuntu/桌面/Git/Deeplearning/Python课件/5-机器学习/J老师/others/USA_Housing.csv'
2 USAhousing = pd.read_csv(csvpath)

In [3]: 1 def address2district(address):
2     temp = address.split('\n')[1] #地址是以\n来进行分割的，所以[1]部分是我们所需要的
3     if ',' in temp: #其中一部分的地址是以','来进行分割的
4         district = temp.split(',')[1].split()[0]
5     else:
6         district=temp.split(' ')[1]
7     return district
8
9 USAhousing.loc[:, 'district'] = USAhousing['Address'].apply(address2district)
```

▼ 建立模型

```
In [4]: 1 # 特征编码
2 from sklearn.preprocessing import LabelEncoder
3
4 le = LabelEncoder()
5
6 le.fit(USAhousing['district'].values)
7 USAhousing['district_code'] = le.transform(USAhousing['district'].values)
8
9 # 特征筛选
10 X = USAhousing[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
11                'Avg. Area Number of Bedrooms', 'Area Population', 'district_code']]
12
13 y = USAhousing['Price']

In [5]: 1 # 分配训练集和测试集
2 from sklearn.model_selection import train_test_split
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=3)

In [6]: 1 # 特征缩放：数据集归一化
2 from sklearn.preprocessing import StandardScaler # 数据预处理类
3
4 # 归一化方法：标准差标准化(zero-mean) 转化函数: x = (x-mean) / std
5 # 经过处理后得到的数据集符合标准正态分布，即均值为0，标准差为1
6 # 适用于本身服从正态分布的数据
7
8 scaler = StandardScaler()
9
10 # 下面两行代码是固定用法，不能颠倒顺序
11 X_train = scaler.fit_transform(X_train) # 求得 训练集 的平均值和方差并应用在 训练集 上，同时也会保存
12 X_test = scaler.transform(X_test)      # 用保存的 训练集 的平均值和方差来应用在 测试集 上
13
14 # 数据预处理中的方法：
15
16 # - fit():
17 # 解释：简单来说，就是求得训练集x的均值啊，方差啊，最大值啊，最小值，这些训练集x固有的属性。可以理解为一个训练过程
18
19 # - transform():
20 # 解释：在Fit的基础上，进行标准化，降维，归一化等操作（看具体用的是哪个工具，如PCA，StandardScaler等）
21
22 # - fit_transform():
23 # 解释：fit_transform是fit和transform的组合，既包括了训练又包含了转换

In [7]: 1 # 模型训练
2 from sklearn.linear_model import LassoCV,ElasticNetCV
3
4 lscv = LassoCV(alphas=[1.0, 0.1, 0.01, 0.005, 0.0025, 0.001, 0.00025])
5
6 lscv.fit(X_train, y_train)
7
8 print(lscv.alpha_)
9
10 encv = ElasticNetCV(alphas=[0.1, 0.01, 0.005, 0.0025, 0.001],l1_ratio=[0.1, 0.25, 0.5, 0.75, 0.8])
11
12 encv.fit(X_train, y_train)
13
14 print((encv.alpha_, encv.l1_ratio_))

1.0
(0.001, 0.8)
```

▼ 模型评估

```
In [8]: 1 from sklearn import metrics
2
3 # 模型评估函数
4 def print_evaluate(y_test, y_predict):
5     mse = metrics.mean_squared_error(y_test, y_predict) #MSE
6     mae = metrics.mean_absolute_error(y_test, y_predict) #MAE
7     rmse = np.sqrt(mse) #RMSE
8     r2 = metrics.r2_score(y_test, y_predict) #R2 Square
9     print(f'MSE: {mse}\nRMSE: {rmse}\nMAE: {mae}\nR2: {r2}\n_____ \n')
10
11 # 模型预测 输出评估结果
12 test_pred1 = lscv.predict(X_test)
13 test_pred2 = encv.predict(X_test)
14 print("Lasso测试集计算结果: \n_____ ")
15 print_evaluate(y_test, test_pred1)
16 print("ElasticNet测试集计算结果: \n_____ ")
17 print_evaluate(y_test, test_pred2)
```

Lasso测试集计算结果：

MSE:	10776609405.682356
RMSE:	103810.44940506884
MAE:	83684.59354924301
R2:	0.9151206757696965

ElasticNet测试集计算结果：

MSE:	10777410385.78996
RMSE:	103814.30723069899
MAE:	83688.86956334059
R2:	0.9151143670460807

