

汽车价格分析与预测

案例背景

A公司是一家中国汽车生产公司。近年来，随着A公司不断的发展与壮大，管理层希望拓展海外市场，在B国建立汽车生产基地，并在当地生产与销售汽车。

为了能够在竞争中处于有利位置，A公司需要熟悉影响B国汽车定价的主要因素。这是因为，不同国家的汽车定价原则很可能是不同的，如果能够提前知道不同因素对汽车定价的影响，这对A公司是非常有利的。

任务说明

我们的任务是：

- 建立模型，确定影响汽车定价的主要因素。
- 在不同变量组合下，评估（预测）汽车的价格。

这样一来，就能够产生如下的价值：

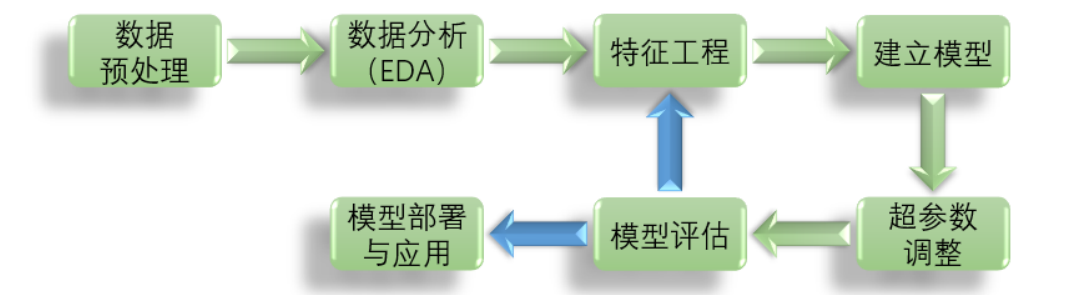
- 管理层借助于模型，熟悉海外新市场的定价动态。
- 根据特定的价格水平，进行汽车的设计，制定商业策略。

数据集描述

为了能够顺利完成任务，A公司联系了B国当地的一家汽车咨询公司（C公司）。根据市场调查，C公司收集了整个B国市场上各种类型汽车与价格的数据集。

属性	说明
car_ID	汽车编号。
symboling	保险风险等级。
CarName	汽车名称。
fueltype	燃料类型。
aspiration	增压机类型。
doornumber	车门数量。
carbody	车身类型。
drivewheel	车轮驱动。
enginelocation	发动机位置。
wheelbase	汽车的轴距。
carlength	汽车长度。
carwidth	汽车宽度。
carheight	汽车高度。
curbweight	汽车重量。
enginetype	发动机类型。
cylindernumber	气缸数量。
enginesize	发动机尺寸。
fuelsystem	燃料系统。
boreratio	口径比率。
stroke	冲程。
compressionratio	压缩比。
horsepower	马力。
peakrpm	最大每分钟转数。
citympg	城市路每加仑燃料所行英里数。
highwaympg	高速路每加仑燃料所行英里数。
price	汽车价格。

数据挖掘流程



数据预处理

加载数据集

导入所需要的库，并进行一些基本设置。

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 sns.set(style="darkgrid", font_scale=1.2)
7 plt.rcParams["font.family"] = "SimHei"
8 plt.rcParams["axes.unicode_minus"] = False

1 data = pd.read_csv("car.csv", header=0)
2 print(data.shape)
3 # 当列数较多时，会显示不完整，可以设置最大显示列数。
4 pd.set_option("display.max_columns", 100)
5 data.head()

1 | (205, 26)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wh
0	1	3	alfa-romero giulia	gas	std	two	convertible	rwd	front	88.0
1	2	3	alfa-romero stelvio	gas	std	two	convertible	rwd	front	88.0
2	3	1	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd	front	94.0
3	4	2	audi 100 ls	gas	std	four	sedan	fwd	front	99.0
4	5	2	audi 100ls	gas	std	four	sedan	4wd	front	99.0

数据清洗

缺失值

我们可以调用DataFrame对象的info方法来显示数据的缺失情况，数据类型，占用空间等信息。

```
1 data.info()

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 205 entries, 0 to 204
3 Data columns (total 26 columns):
4 #   Column              Non-Null Count  Dtype
5 ---  ---
6 0   car_ID              205 non-null    int64
7 1   symboling           205 non-null    int64
```

```
8 | 2 CarName          205 non-null object
9 | 3 fueltype         205 non-null object
10 | 4 aspiration        205 non-null object
11 | 5 doornumber        205 non-null object
12 | 6 carbody           205 non-null object
13 | 7 drivewheel        205 non-null object
14 | 8 enginelocation    205 non-null object
15 | 9 wheelbase         205 non-null float64
16 | 10 carlength         205 non-null float64
17 | 11 carwidth          205 non-null float64
18 | 12 carheight         205 non-null float64
19 | 13 curbweight        205 non-null int64
20 | 14 enginetype        205 non-null object
21 | 15 cylindernumber    205 non-null object
22 | 16 enginesize         205 non-null int64
23 | 17 fuelsystem        205 non-null object
24 | 18 boreratio         205 non-null float64
25 | 19 stroke            205 non-null float64
26 | 20 compressionratio  205 non-null float64
27 | 21 horsepower        205 non-null int64
28 | 22 peakrpm           205 non-null int64
29 | 23 citympg           205 non-null int64
30 | 24 highwaympg        205 non-null int64
31 | 25 price             205 non-null float64
32 | dtypes: float64(8), int64(8), object(10)
33 | memory usage: 41.8+ KB
```

重复值

```
1 | data.duplicated().sum()
```

```
1 | 0
```

数据转换

CarName

对于CarName列，只需要保留前部分的厂商名称即可，对于厂商后面的内容，没有分析价值，可以去掉。

```
1 | t = data["CarName"].str.split(" ", expand=True)
2 | display(t.head())
3 | data["CarName"] = t[0]
4 | data["CarName"].value_counts()
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	0	1	2	3	4
0	alfa-romero	giulia	None	None	None
1	alfa-romero	stelvio	None	None	None
2	alfa-romero	Quadrifoglio	None	None	None
3	audi	100	ls	None	None
4	audi	100ls	None	None	None

```
1 | toyota      31
2 | nissan      17
3 | mazda       15
4 | honda       13
5 | mitsubishi  13
6 | subaru      12
7 | peugeot     11
8 | volvo       11
9 | volkswagen   9
10 | dodge        9
11 | buick         8
```

```

12 | bmw            8
13 | audi           7
14 | plymouth       7
15 | saab           6
16 | isuzu          4
17 | porsche        4
18 | alfa-romero    3
19 | chevrolet      3
20 | jaguar         3
21 | vw             2
22 | maxda          2
23 | renault        2
24 | toyouta        1
25 | vokswagen      1
26 | Nissan         1
27 | mercury        1
28 | porcshce       1
29 | Name: CarName, dtype: int64

```

我们成功将厂商名称提取出来，不过，名称中存在一些问题：

- 个别厂商名称存在一些拼写错误。
- 同一个厂商的大小写不统一。
- 同一个厂商，有的用全称，有的用简称。

因此，我们需要进行如下的转换：

- maxda => mazda
- Nissan => nissan
- porcshce => porsche
- toyouta => toyota
- vokswagen => volkswagen
- vw => volkswagen

```

1 | data["CarName"] = data["CarName"].replace({"maxda": "mazda", "Nissan": "nissan", "porcshce": "porsche",
2 |      "toyouta": "toyota", "vokswagen": "volkswagen", "vw": "volkswagen"})
3 | data["CarName"].value_counts()

```

```

1 | toyota        32
2 | nissan        18
3 | mazda        17
4 | mitsubishi    13
5 | honda         13
6 | volkswagen    12
7 | subaru        12
8 | peugeot       11
9 | volvo         11
10 | dodge         9
11 | buick          8
12 | bmw           8
13 | audi           7
14 | plymouth       7
15 | saab           6
16 | porsche        5
17 | isuzu          4
18 | jaguar         3
19 | chevrolet      3
20 | alfa-romero    3
21 | renault        2
22 | mercury        1
23 | Name: CarName, dtype: int64

```

探索性数据分析

探索性数据分析（EDA——Exploratory Data Analysis），是一种分析数据集以总结其主要特征的方法，在分析过程中，通常会通过可视化来进行辅助。探索性数据分析的作用为：

- 探索发现数据的一些特性。
- 探索发现变量与变量之间的关系。
- 探索发现变量与目标值之间的关系。
- 方便于特征工程的构建。

相关性分析（两个连续变量）

变量浏览

首先，我们可以输出所有的数值类型变量。

```

1 # 方案一。
2 # 定义列表，用来存放数值类型变量的名称。
3 number_var = []
4 for k, v in data.dtypes.items():
5     if np.issubdtype(v, np.number):
6         number_var.append(k)
7 print("数值类型变量个数: ", len(number_var))
8 print(number_var)

```

```

1 数值类型变量个数: 16
2 ['car_ID', 'symboling', 'wheelbase', 'carlength', 'carwidth', 'carheight', 'curbweight', 'enginesize', 'bore_ratio', 'stroke',
  'compressionratio', 'horsepower', 'peakrpm', 'citympg', 'highwaympg', 'price']

```

```

1 # 方案二
2 # 也可以使用DataFrame对象的select_dtypes方法，该方法会返回DataFrame
3 # 对象，包含所有符合参数指定的类型列。
4 print(data.select_dtypes(np.number).columns)

```

```

1 Index(['car_ID', 'symboling', 'wheelbase', 'carlength', 'carwidth',
2        'carheight', 'curbweight', 'enginesize', 'bore_ratio', 'stroke',
3        'compressionratio', 'horsepower', 'peakrpm', 'citympg', 'highwaympg',
4        'price'],
5        dtype='object')

```



课堂练习



如果只是为了获取所有数值类型列的名称，优先使用哪种方案？

- A 方案一。
- B 方案二。
- C 都可以。
- D 根据具体情况而定。



散点图矩阵

对于两个连续类型的变量，可以绘制两个变量的散点图，来观察二者是否具有一定的关系。

```

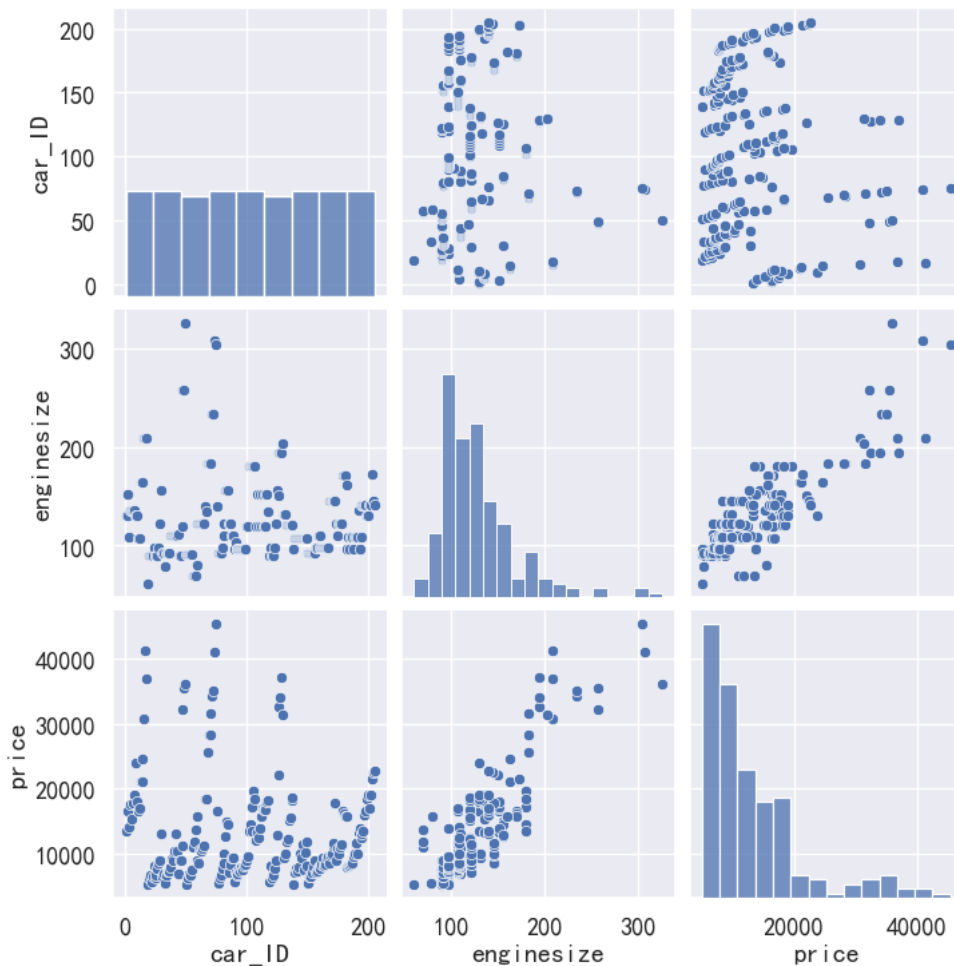
1 # 绘制变量之间的散点图矩阵。
2 # vars: 指定使用哪些变量绘制，默认使用所有变量。
3 # 此处变量较多，故只选择了三个变量。
4 sns.pairplot(data, vars=["car_ID", "enginesize", "price"])

```

```

1 <seaborn.axisgrid.PairGrid at 0x26c60d8f070>

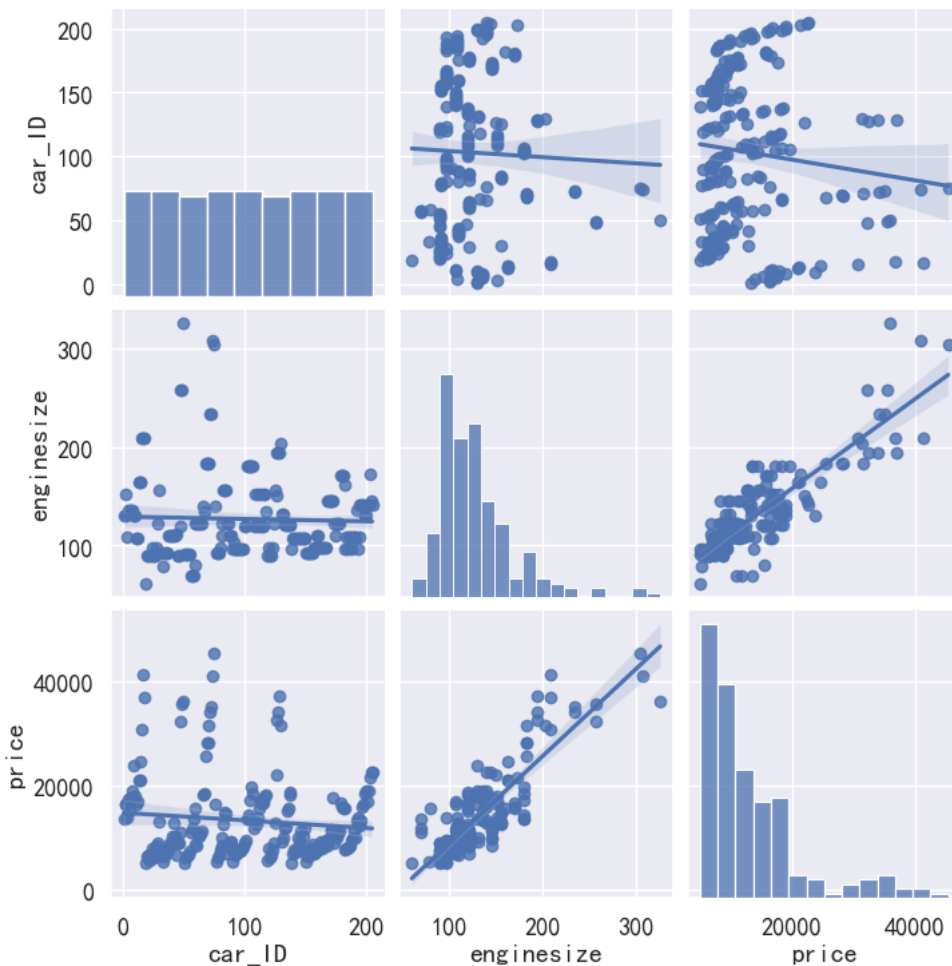
```



通过散点图矩阵，我们发现，enginesize与price存在较强的线性关系。另外，在绘制散点图矩阵的同时，也可以绘制回归线。

```
1 # kind: 取值为scatter（默认）与reg。
2 # scatter: 绘制散点图。
3 # reg: 除了散点图外，同时绘制回归线与置信区间。
4 sns.pairplot(data, vars=["car_ID", "enginesize", "price"], kind="reg")
```

```
1 <seaborn.axisgrid.PairGrid at 0x26c61c69cf0>
```



协方差

协方差，体现的是两个变量之间的分散程度以及两个变量变化步调是否一致。设变量：

$$X = (x_1, x_2, \dots, x_n) \quad Y = (y_1, y_2, \dots, y_n)$$

协方差的定义为：

$$Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$$

在样本中，计算方式如下：

$$Cov(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

- \bar{x} ：变量X的均值。

相关系数

相关系数，可以用来体现两个连续变量之间的相关性，是协方差的标准化版本，最为常用的为皮尔逊相关系数，衡量两个连续变量的线性相关性。两个变量的相关系数定义为：

$$r(X, Y) = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}$$

在样本中，计算方式如下：

$$r(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

- σ_X ：变量X的标准差。

相关系数的取值范围为[-1, 1]，我们可以根据相关系数的取值（绝对值）来衡量两个变量的相关性：

- 0.8-1.0：极强相关

- 0.6-0.8：强相关
- 0.4-0.6：中等程度相关
- 0.2-0.4：弱相关
- 0.0-0.2：极弱相关或无相关

相关系数计算

我们以汽车长度（carlength）与汽车宽度（carwidth）为例，来求解两个变量的相关系数。

```
1 | x = data["carlength"]
2 | y = data["carwidth"]
3 | # 计算carwidth与carheight的协方差。
4 | a = (x - x.mean()) * (y - y.mean())
5 | cov = np.sum(a) / (len(a) - 1)
6 | print("协方差: ", cov)
7 | # 计算carwidth与carheight的相关系数。
8 | corr = cov / (x.std() * y.std())
9 | print("相关系数: ", corr)
```

```
1 | 协方差: 22.261035150645633
2 | 相关系数: 0.8411182684818456
```

此外，Series对象也提供了cov与corr方法，来求解协方差与相关系数。

```
1 | print("协方差: ", x.cov(y))
2 | print("相关系数: ", x.corr(y))
```

```
1 | 协方差: 22.261035150645633
2 | 相关系数: 0.841118268481846
```

同样，DataFrame对象提供了计算了corr方法，可以计算任意两个数值变量的相关系数。

说明：

- 非数值类型不会在相关系数的结果中体现。

```
1 | data.corr(numeric_only=True)
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

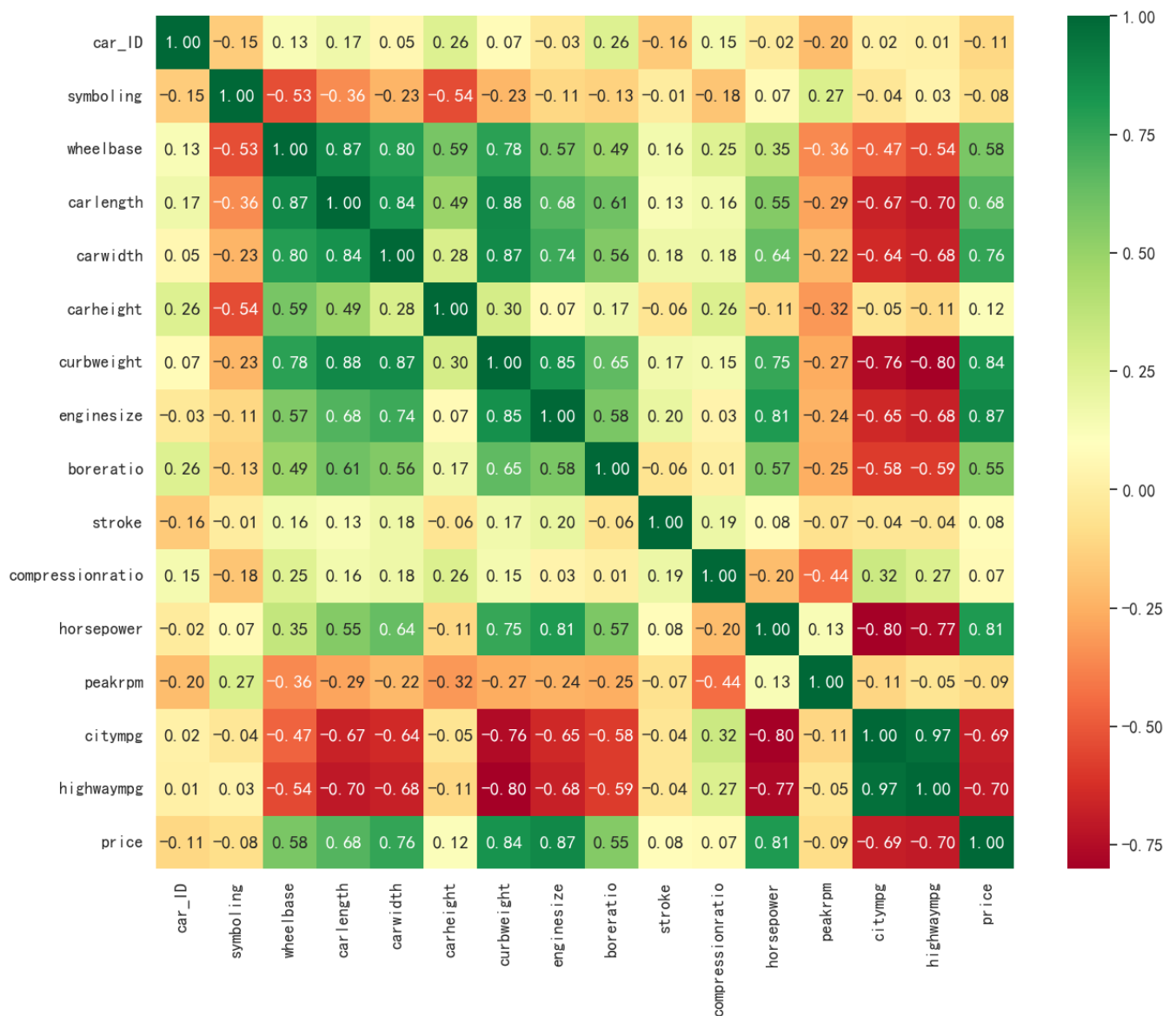
	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbweight	enginesize	boreratio	stroke
car_ID	1.000000	-0.151621	0.129729	0.170636	0.052387	0.255960	0.071962	-0.033930	0.260064	-0.160824
symboling	-0.151621	1.000000	-0.531954	-0.357612	-0.232919	-0.541038	-0.227691	-0.105790	-0.130051	-0.008735
wheelbase	0.129729	-0.531954	1.000000	0.874587	0.795144	0.589435	0.776386	0.569329	0.488750	0.160959
carlength	0.170636	-0.357612	0.874587	1.000000	0.841118	0.491029	0.877728	0.683360	0.606454	0.129533
carwidth	0.052387	-0.232919	0.795144	0.841118	1.000000	0.279210	0.867032	0.735433	0.559150	0.182942
carheight	0.255960	-0.541038	0.589435	0.491029	0.279210	1.000000	0.295572	0.067149	0.171071	-0.055307
curbweight	0.071962	-0.227691	0.776386	0.877728	0.867032	0.295572	1.000000	0.850594	0.648480	0.168790
enginesize	-0.033930	-0.105790	0.569329	0.683360	0.735433	0.067149	0.850594	1.000000	0.583774	0.203129
boreratio	0.260064	-0.130051	0.488750	0.606454	0.559150	0.171071	0.648480	0.583774	1.000000	-0.055909
stroke	-0.160824	-0.008735	0.160959	0.129533	0.182942	-0.055307	0.168790	0.203129	-0.055909	1.000000
compressionratio	0.150276	-0.178515	0.249786	0.158414	0.181129	0.261214	0.151362	0.028971	0.005197	0.186110
horsepower	-0.015006	0.070873	0.353294	0.552623	0.640732	-0.108802	0.750739	0.809769	0.573677	0.080940
peakrpm	-0.203789	0.273606	-0.360469	-0.287242	-0.220012	-0.320411	-0.266243	-0.244660	-0.254976	-0.067964
citympg	0.015940	-0.035823	-0.470414	-0.670909	-0.642704	-0.048640	-0.757414	-0.653658	-0.584532	-0.042145
highwaympg	0.011255	0.034606	-0.544082	-0.704662	-0.677218	-0.107358	-0.797465	-0.677470	-0.587012	-0.043931
price	-0.109093	-0.079978	0.577816	0.682920	0.759325	0.119336	0.835305	0.874145	0.553173	0.079443

为了能够更清晰的呈现相关系数值，我们可以使用热图来展示相关系数。


```

1 plt.figure(figsize=(15, 12))
2 # cmap: 颜色图, 控制显式的颜色风格。
3 # annot: 是否显示数值, 默认为False。
4 # fmt: 指定数值格式化方式。
5 ax = sns.heatmap(data.corr(numeric_only=True), cmap=plt.cm.RdYlGn, annot=True, fmt=".2f")

```



结果统计

从相关系数的结果, 我们可以得出以下结论:

- car_ID, symboling, carheight, stroke, compressionratio, peakrpm与price相关度较低。
- wheelbase, carlength, carwidth, boreratio, citympg, highwaympg与price中度相关。
 - 正相关: wheelbase, carlength, carwidth, boreratio。
 - 负相关: citympg, highwaympg。
- curbweight, enginesize, horsepower与price高度相关。
 - 正相关: curbweight, enginesize, horsepower。
 - 负相关: 无。
- symboling虽然与price不相关, 但该变量也可以当做类别变量看待。

相关性分析 (类别变量与连续变量)

变量浏览

我们输出所有的类别类型变量。

```

1 category_var = []
2 for k, v in data.dtypes.items():
3     if not np.issubdtype(v, np.number):
4         category_var.append(k)
5 print("类别变量个数: ", len(category_var))
6 print(category_var)

```

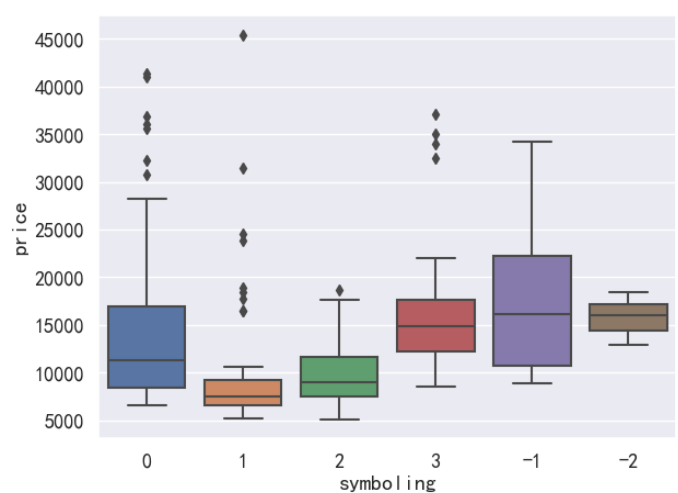
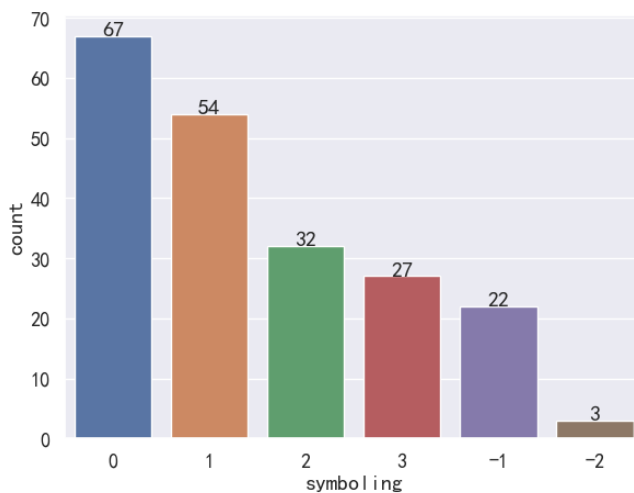
```
1 类别变量个数: 10
2 ['CarName', 'fueltype', 'aspiration', 'doornumber', 'carbody', 'drivewheel', 'enginelocation', 'enginetype', 'cylindernumber',
  'fuelsystem']
```

symboling

symboling (保险风险等级) 在数据集中呈数值类型, 但是, 我们也可以将其视为类别类型, 也许更加合理。

```
1 from scipy import stats
2
3 def category_analysis(col_name, target_name="price"):
4     """对类别变量进行分析。
5
6     分析包括绘制柱形图与箱线图。并对类别变量进行假设检验。
7
8     Parameters
9     -----
10    col_name : str
11        分类变量的名称。
12    target_name : str, 可选。
13        目标变量的名称。
14
15    """
16
17    v = data[col_name].value_counts()
18    # 如果分类数量超过10, 则显示2行, 否则显示1行。
19    if len(v) > 10:
20        row, col = 2, 1
21    else:
22        row, col = 1, 2
23    fig, ax = plt.subplots(row, col)
24    # 根据行数不同, 设置不同的高度。
25    fig.set_size_inches(15, row * 5)
26    sns.countplot(x=col_name, data=data, order=v.index, ax=ax[0])
27    # 在图像上绘制数值。
28    for x, y in enumerate(v):
29        text = ax[0].text(x, y, y)
30        # 文本数值居中对齐。
31        text.set_ha("center")
32    sns.boxplot(x=col_name, y=target_name, order=v.index, data=data, ax=ax[1])
33    if len(v) > 10:
34        for a in ax:
35            a.set_xticklabels(a.get_xticklabels(), rotation=30)
36        plt.subplots_adjust(hspace=0.3)
37    # 根据col_name进行分组。
38    g = data.groupby(col_name)[target_name]
39    group_data = []
40    g.apply(lambda s: group_data.append(s.values))
41    # Kruskal-wallis H-test, 用来检验各个分组的中位数是否相等。
42    # 原假设: 每个分组的中位数相等。
43    # 备择假设: 至少存在两个分组的中位数不等。
44    print(stats.kruskal(*group_data))
45
46
47 category_analysis("symboling")
```

```
1 KruskalResult(statistic=55.05140268551899, pvalue=1.2739605541044965e-10)
```



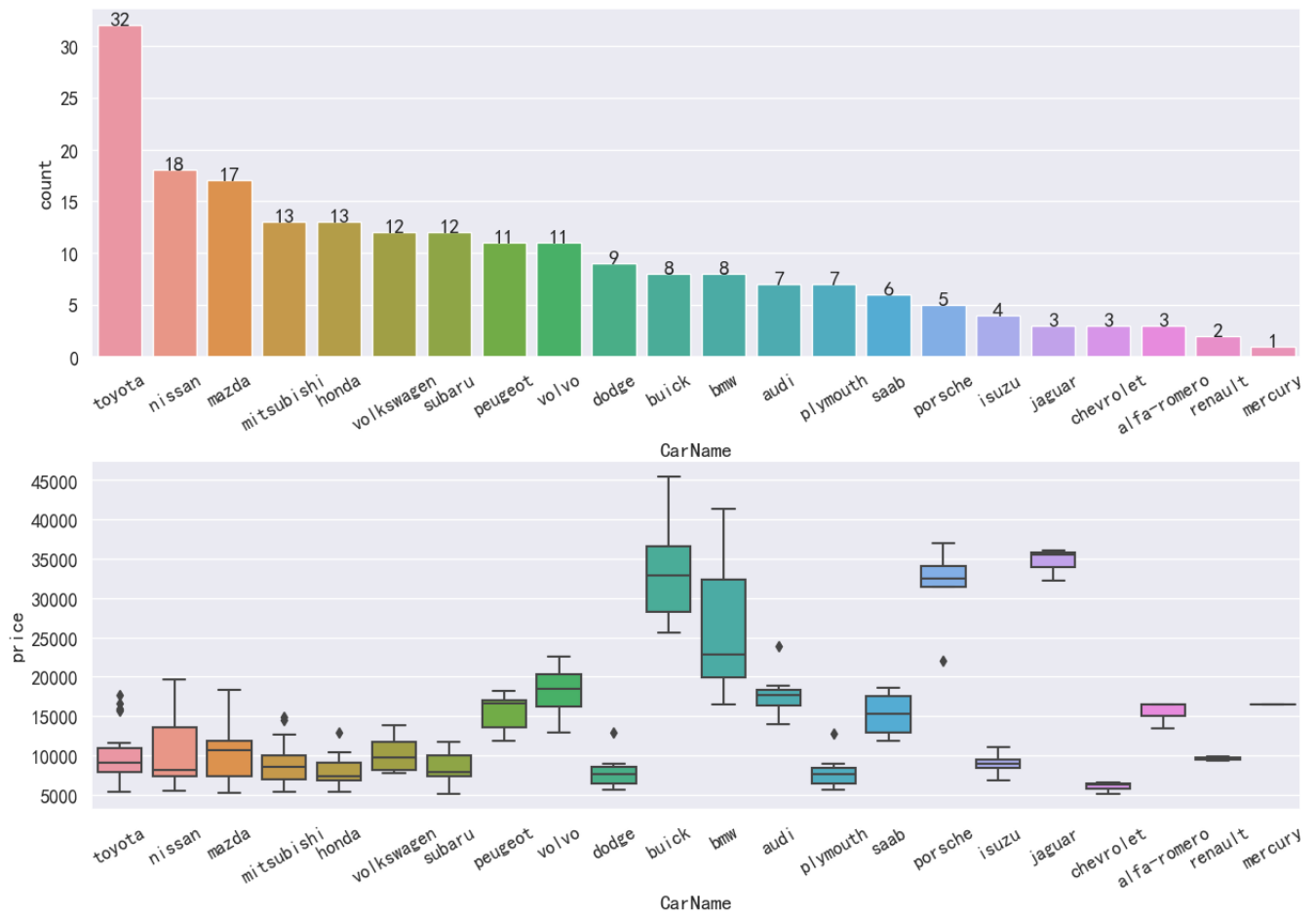
- 保险等级主要以0级与1级为主。
- Kruskal-Wallis H-test结果非常显著, 说明不同风险等级, 确实与汽车价格会有一定的关联。

- -2等级的数量很少。

CarName

```
1 | category_analysis("CarName")
```

```
1 | KruskalResult(statistic=128.0293000475285, pvalue=2.3580012078648455e-17)
```

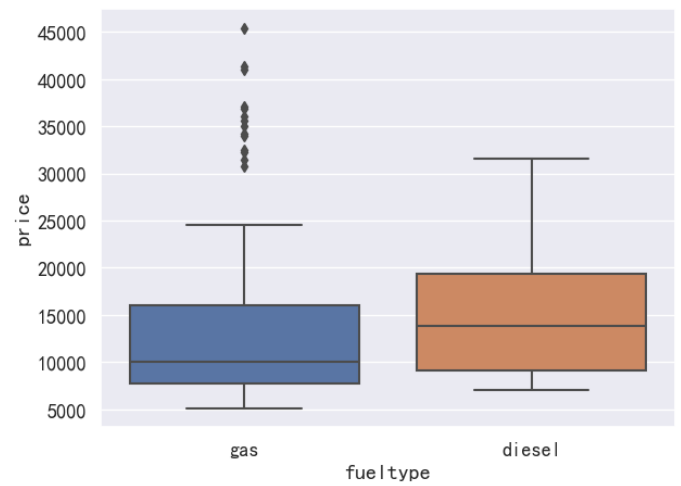
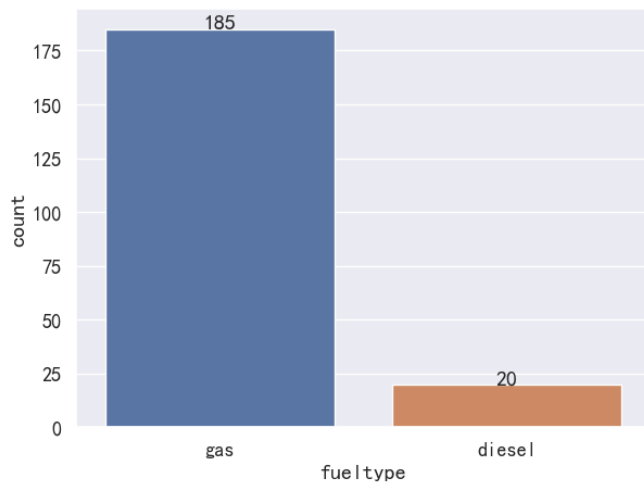


- 不同品牌的汽车，车型数量差异较大，车型数量以toyota为主。
- 不同品牌的汽车，在价格上有着很大的差别。
 - bmw, buick, porsche与jaguar的价格较高。
 - bmw与buick的价格浮动区间较大。

fueltype

```
1 | category_analysis("fueltype")
```

```
1 | KruskalResult(statistic=3.9833925913400314, pvalue=0.0459509243643051)
```

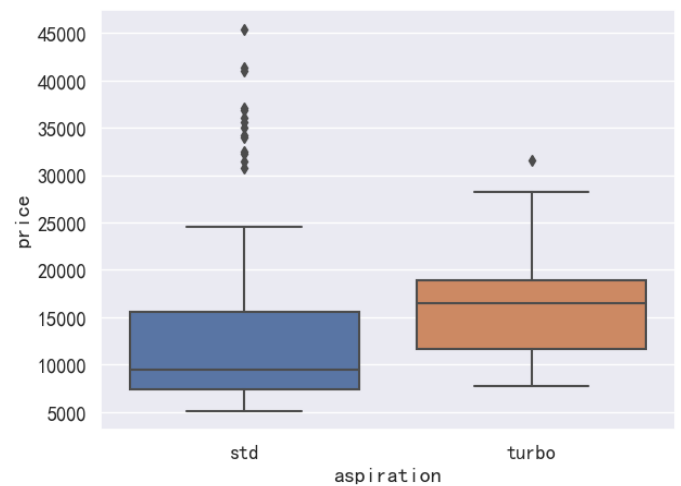
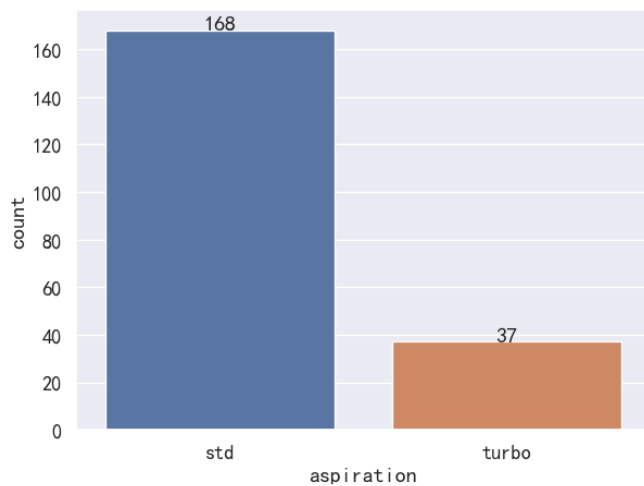


- 汽车主要以汽油为主，柴油的汽车数量较少。
- 汽油与柴油两种类型在价格上有所差异（但并不是非常显著）。
 - 柴油汽车价格普遍高于汽油汽车价格。

aspiration

```
1 | category_analysis("aspiration")
```

```
1 | KruskalResult(statistic=19.47319231451774, pvalue=1.020215018682171e-05)
```

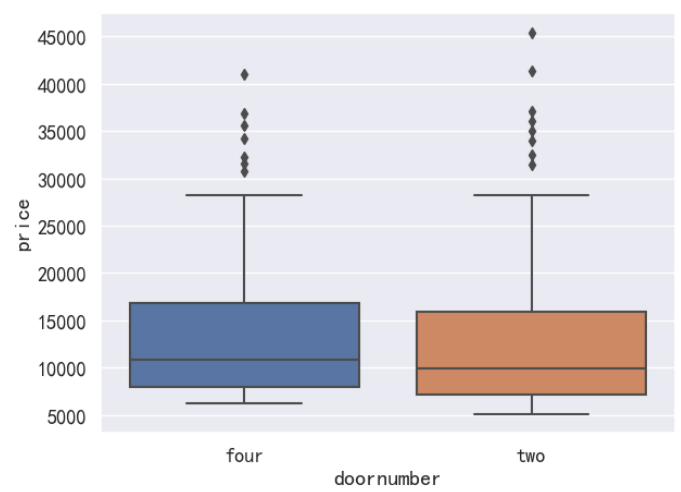
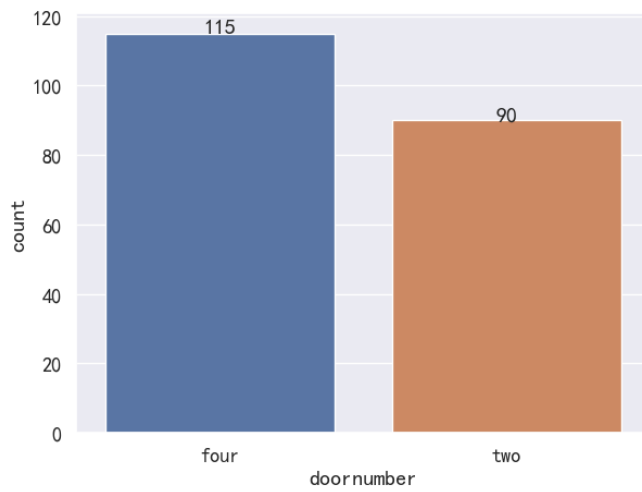


- 增压机类型以标准类型（std）为主。涡轮增压（turbo）的数量较少。
- 两种类型价格差异显著。
 - 标准类型价格低于涡轮增压类型价格。

doornumber

```
1 | category_analysis("doornumber")
```

```
1 | KruskalResult(statistic=2.8612616578123022, pvalue=0.09073630460630733)
```

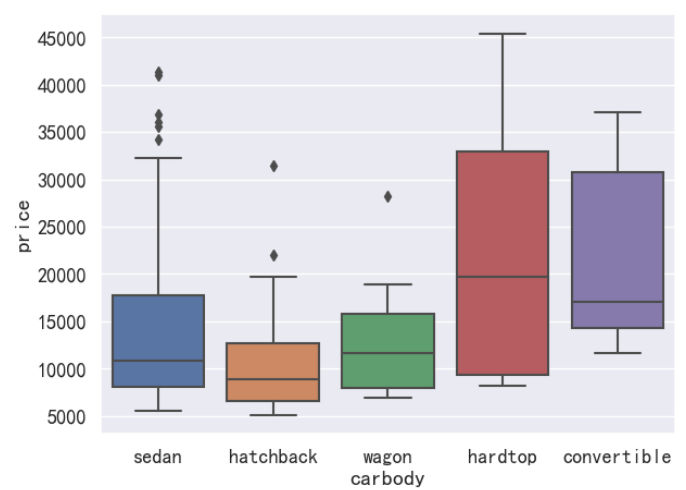
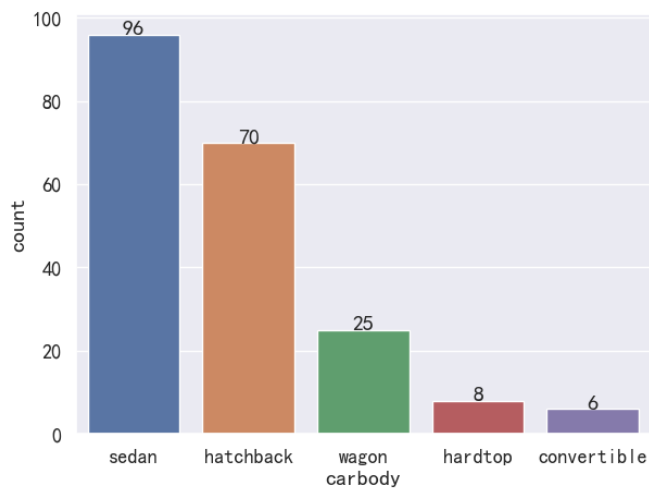


- 双门与四门汽车类型接近。
- 两个类型对价格影响不显著。

carbody

```
1 | category_analysis("carbody")
```

```
1 | KruskalResult(statistic=22.183437941390444, pvalue=0.0001842534808688567)
```

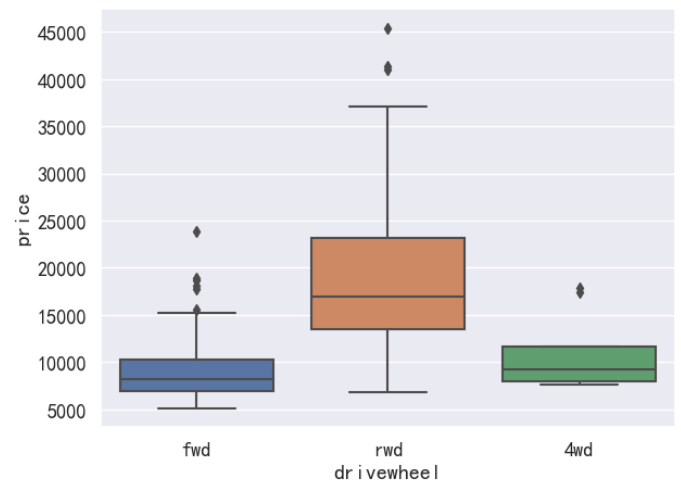
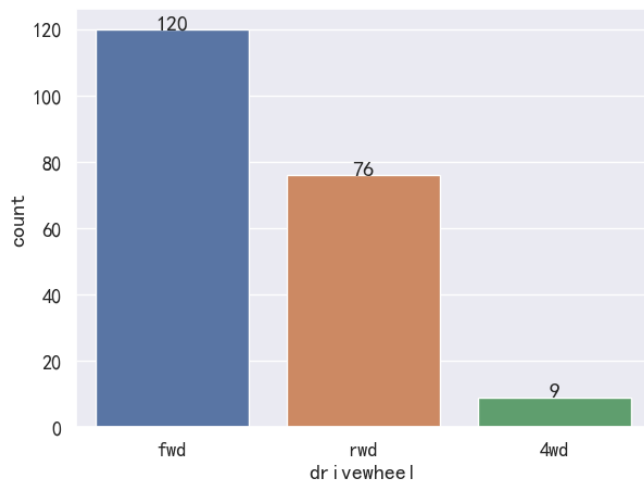


- 汽车类型以轿车（sedan）与掀背车（hatchback）为主。
- 汽车类型对价格影响显著。
- 其中轿车，硬顶轿车（hardtop）的价格浮动区间较大。掀背车与货车（wagon carbody）价格浮动区间较小。

drivewheel

```
1 | category_analysis("drivewheel")
```

```
1 | KruskalResult(statistic=95.02387066880598, pvalue=2.3218205980565562e-21)
```

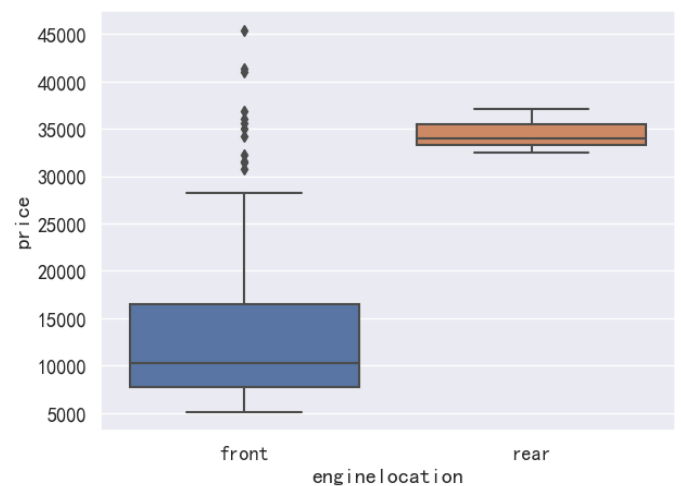
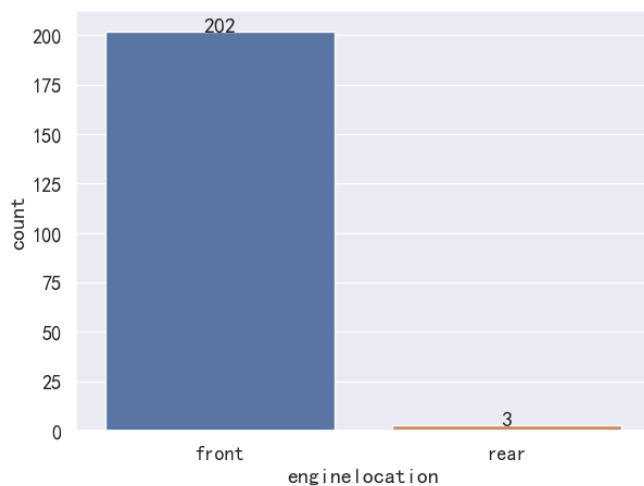


- 4轮驱动车型 (4wd) 很少，前轮驱动 (fwd) 最多，后轮驱动 (rwd) 其次。
- 不同类型对价格影响显著。
- 后轮驱动的价格最高。

engine location

```
1 | category_analysis("engine location")
```

```
1 | KruskalResult(statistic=7.753234528191141, pvalue=0.005361642811118273)
```

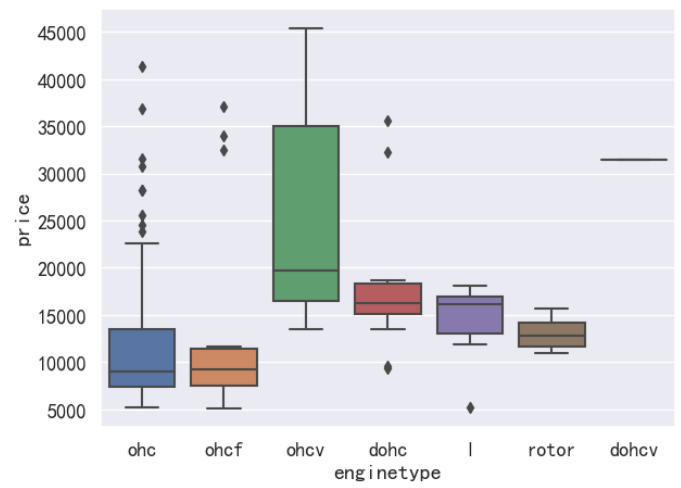
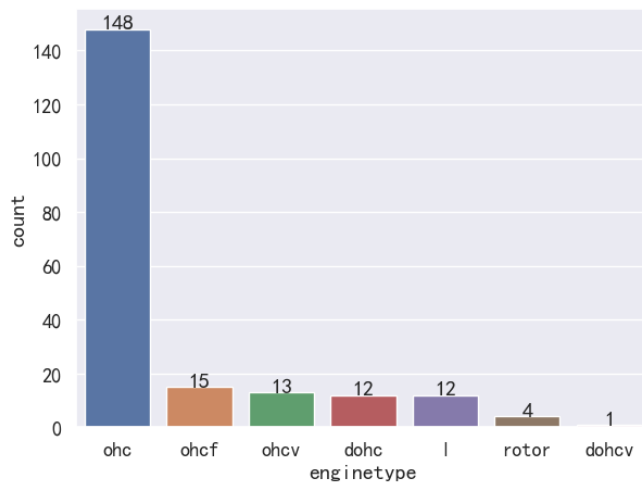


- 发动机位置几乎都是前置。
- 虽然发动机位置对价格影响显著，但是后置发动机样本数量过少，可能不足以说明问题。
- 前置与后置发动机的价格区间差距很大。

enginetype

```
1 | category_analysis("enginetype")
```

```
1 | KruskalResult(statistic=40.133921438590676, pvalue=4.287270740376635e-07)
```

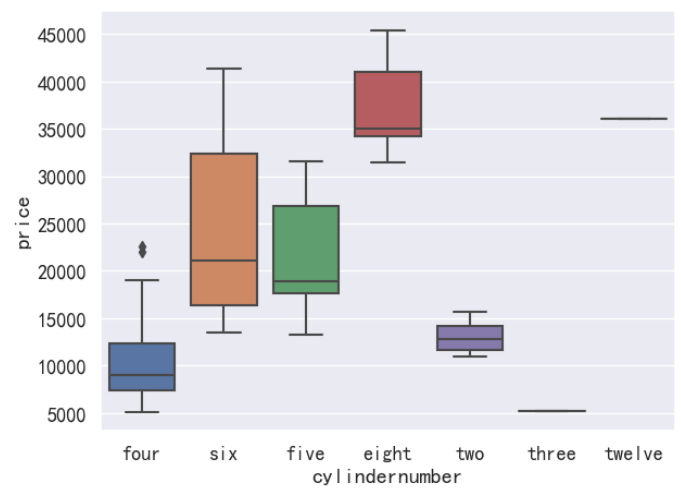
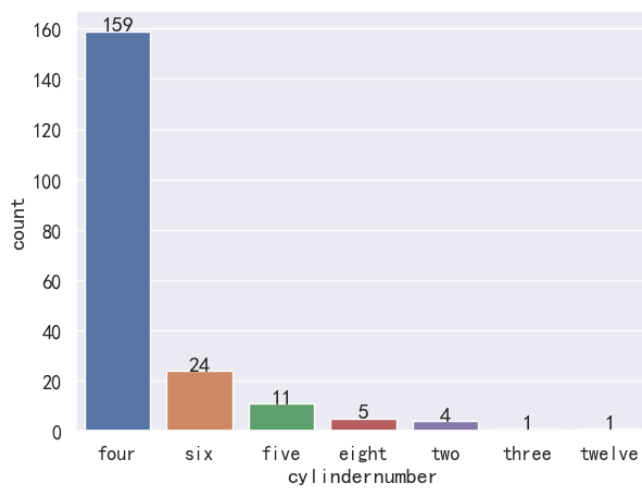


- 发动机几乎都是ohc类型，其余6个类型很少。
- 不同类型发动机，价格差异显著。

cylindernumber

```
1 | category_analysis("cylindernumber")
```

```
1 | KruskalResult(statistic=84.76297575868743, pvalue=3.696266502144558e-16)
```

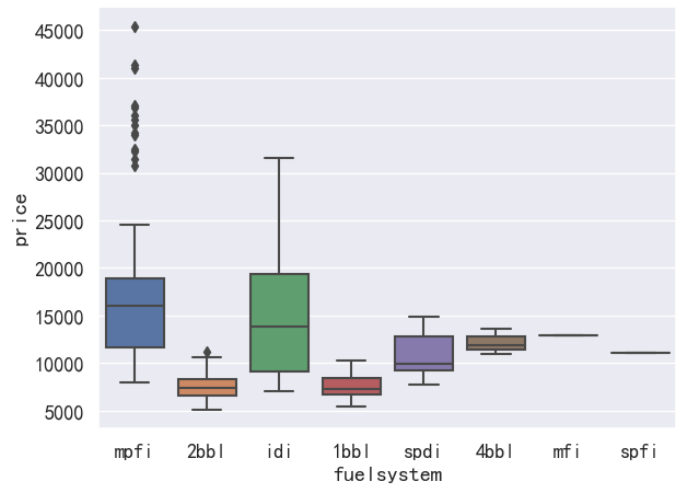
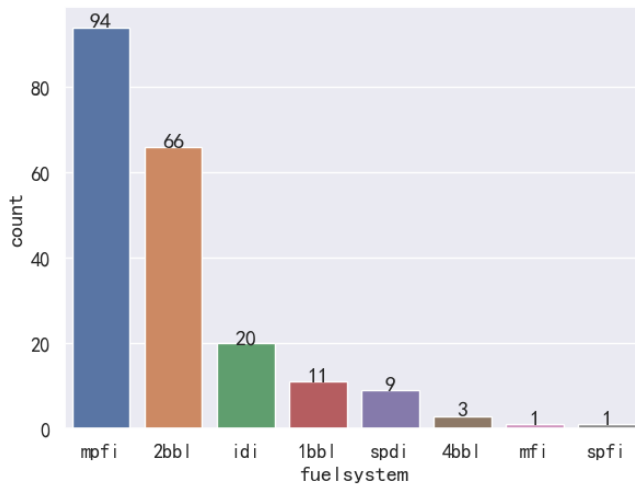


- 气缸基本都是4气缸，其余气缸车型很少。
- 三气缸车型价格较低，十二气缸与八气缸类型价格较高。
- 气缸数量对汽车价格影响显著。

fuelsystem

```
1 | category_analysis("fuelsystem")
```

```
1 | KruskalResult(statistic=125.0821070976894, pvalue=6.683871002354369e-24)
```



- 燃料系统以mpfi与2bbl为主。
- 不同燃料系统，差异显著。

特征工程

特征工程是数据挖掘 / 机器学习中非常重要的一个环节。所谓特征工程，就是从原始数据中提取有用的信息，以便于更好的建立模型，提高模型的质量。

删除无用特征

如果变量与目标（都是连续变量）之间的相关系数很低，则通常情况下，这些变量对预测因变量不会起到太大作用，我们可以考虑将相关系数很低的变量删除。

- 很明显，car_ID仅是一个汽车编号，与汽车的价格毫无关联，我们可以将其删除。
- 对于其他相关性很低的变量，我们暂时保留，待定。

```
1 data.drop("car_ID", inplace=True, axis=1)
```

建立模型

基模型（baseline）

在进行特征工程之前，我们可以事先建立一个相对简单的基模型，作为我们的基准，并且在后面的操作中，进行改进。这里，我们仅使用enginesize一个特征来建立基模型。

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3
4 SEED = 0
5 # 将Series转换为DataFrame。
6 X = data["engine_size"].to_frame()
7 y = data["price"]
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=SEED)
9 lr = LinearRegression()
10 lr.fit(X_train, y_train)
11 print(lr.score(X_train, y_train))
12 print(lr.score(X_test, y_test))
```

```
1 0.7557897278380871
2 0.7797940080893313
```

使用所有数值变量

然后，我们使用所有的数值变量，来进行拟合。

```
1 # 选择所有数值（连续）类型的变量。
2 X = data.select_dtypes(include=np.number)
3 y = X.pop("price")
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=SEED)
5 lr = LinearRegression()
6 lr.fit(X_train, y_train)
7 print(lr.score(X_train, y_train))
8 print(lr.score(X_test, y_test))
```

```
1 0.858252048258165
2 0.8214693378253439
```

类别变量处理

在使用所有数值类型变量后，我们也应该要考虑类别类型的变量。只不过，类别类型变量不能直接放到模型中使用，我们需要对其进行处理。这里，我们对每个类别变量进行分组，然后使用该组价格的均值（或中位数）来作为该变量的值。


```

1 # symboling当成类别变量看待。
2 for name in category_var + ["symboling"]:
3     # 对每一个类别变量，根据其取值进行分组，使用每组价格的均值，
4     # 来作为新的特征值。
5     v = data.groupby(name)["price"].mean()
6     data[name] = data[name].map(v)
7
8 data.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase
0	17221.296296	15498.333333	12999.7982	12611.270833	12989.924078	21890.500000	19910.809211	12961.097361	88.6
1	17221.296296	15498.333333	12999.7982	12611.270833	12989.924078	21890.500000	19910.809211	12961.097361	88.6
2	10037.907407	15498.333333	12999.7982	12611.270833	12989.924078	10376.652386	19910.809211	12961.097361	94.5
3	10109.281250	17859.166714	12999.7982	12611.270833	13501.152174	14344.270833	9239.308333	12961.097361	99.8
4	10109.281250	17859.166714	12999.7982	12611.270833	13501.152174	14344.270833	11087.463000	12961.097361	99.4

```

1 y = data["price"]
2 x = data.iloc[:, :-1]
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=SEED)
4 lr = LinearRegression()
5 lr.fit(X_train, y_train)
6 print(lr.score(X_train, y_train))
7 print(lr.score(X_test, y_test))

```

```

1 0.9497131369526582
2 0.886804982188443

```

不显著的变量

在之前类别变量分析中，我们通过假设检验（Kruskal-Wallis H-test）发现如下的信息：

- doornumber与price之间不显著。
- fueltype与price之间显著性不太高，处于临界状态。
- 其他变量与price之间显著。

这里，我们可以删除doornumber变量，而对于fueltype，暂时保留。不过，为了能够在后面证明这点，我们将两个变量全部保留。

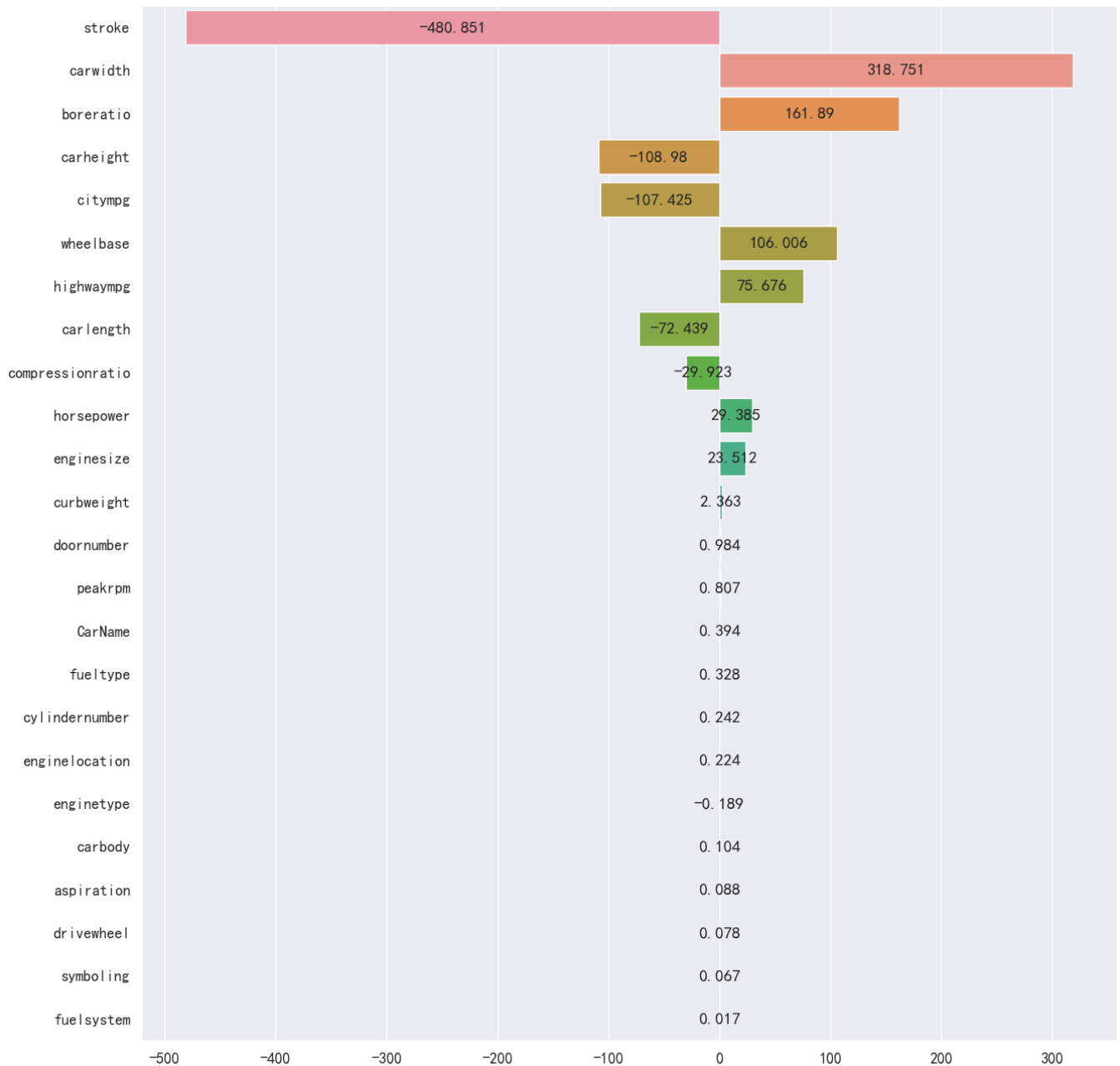
权重可视化

我们可以绘制模型的权重。

```

1 def plot_weight(model, names, figsize=None):
2     """绘制模型的权重。
3
4     会根据权重的绝对值进行降序排列。
5
6     Parameters
7     -----
8     model : object
9         模型对象。
10    names : array-like
11        权重对应的名称。
12    figsize : tuple
13        图像尺寸。
14
15    """
16    s = pd.Series(index=names, data=model.coef_)
17    s.sort_values(ascending=False, inplace=True, key=lambda a: a.abs())
18    ax = sns.barplot(y=s.index, x=s.values)
19    if not figsize:
20        figsize = (15, len(model.coef_) // 1.5)
21    ax.get_figure().set_size_inches(figsize)
22    for y, x in enumerate(s):
23        ax.text(x / 2, y, round(x, 3), va="center", ha="center")
24
25
26 plot_weight(lr, X_train.columns)

```



课堂练习



以上权重值能否客观反映对价格的解释（影响）力度？

- A 能。
- B 不能。
- C 暂时还不确定。



数据标准化

很多算法对特征的数量级都是敏感的，因此，在使用算法之前，我们最好将数据集中的特征转换成相同的量纲，从而消除不同量纲对算法造成的负面影响，我们将这个过程称为**数据标准化**。实际上，即使特征量纲相同，标准化也不会产生负面影响。

在scikit-learn中，常用的标准化方式为：

- 均值标准差标准化（StandardScaler）。
- 最小最大值标准化（MinMaxScaler）。



课堂练习



标准化与否，对线性回归会造成怎样的影响？

- A 影响 w 与 R^2 。
- B 影响 w ，不影响 R^2 。
- C 不影响 w ，影响 R^2 。
- D 不会影响 w 与 R^2 。



执行标准化

```
1 from sklearn.preprocessing import MinMaxScaler
2
3 x_scaler = MinMaxScaler()
4 y_scaler = MinMaxScaler()
5
6 X_train_scale = x_scaler.fit_transform(X_train)
7 y_train_scale = y_scaler.fit_transform(y_train.values.reshape(-1, 1)).ravel()
8 X_test_scale = x_scaler.transform(X_test)
9 y_test_scale = y_scaler.transform(y_test.values.reshape(-1, 1)).ravel()
10
11 # 重新将ndarray数组对象包装成DataFrame对象，方便查看。
12 X_train_scale = pd.DataFrame(X_train_scale, columns=X_train.columns)
13 X_test_scale = pd.DataFrame(X_test_scale, columns=X_test.columns)
14 display(X_train_scale.head(3))
15 display(X_test_scale.head(3))
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	carlength	carwi
0	0.593609	0.076162	0.0	0.0	1.0	0.335334	0.0	0.0	0.288630	0.485039	0.070
1	0.593609	0.162483	0.0	0.0	1.0	0.335334	0.0	0.0	0.355685	0.522835	0.474
2	0.985001	0.162483	0.0	0.0	0.0	0.000000	1.0	0.0	0.253644	0.384252	0.393

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	carlength	carwi
0	0.000000	0.162483	0.0	0.0	0.0	0.000000	0.0	0.0	0.189504	0.228346	0.242
1	1.000000	0.135656	0.0	0.0	1.0	0.168639	1.0	0.0	0.521866	0.680315	0.474
2	0.009787	0.414513	0.0	0.0	0.0	0.335334	0.0	0.0	0.384840	0.514961	0.454

拟合优度对比

对线性回归来讲，标准化不会影响到模型的拟合优度。

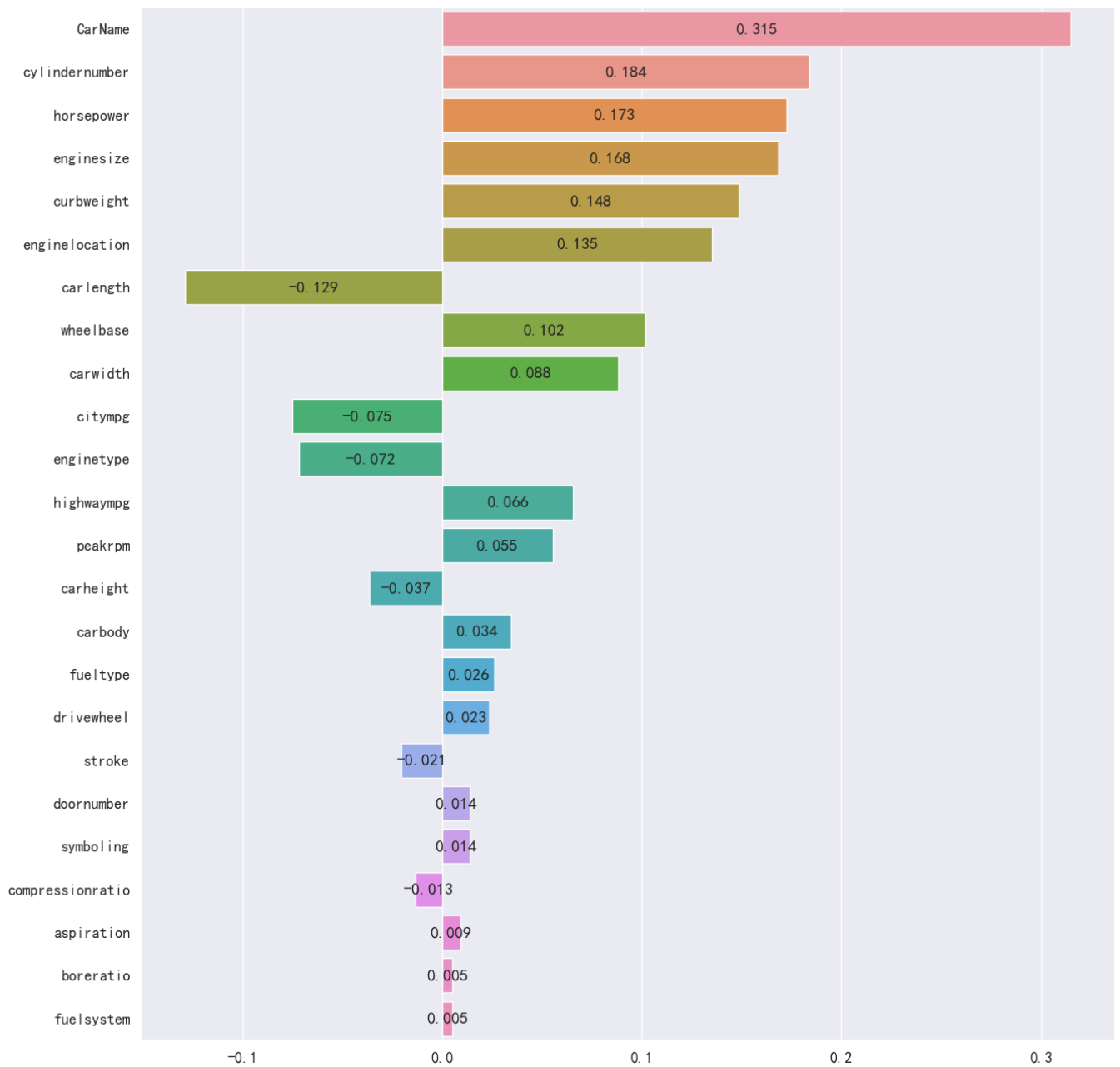
```
1 lr.fit(X_train_scale, y_train_scale)
2 print(lr.score(X_train_scale, y_train_scale))
3 print(lr.score(X_test_scale, y_test_scale))
```

```
1 0.9497131369526582
2 0.8868049821884444
```

权重系数对比

不过，标准化会对权重系数造成影响，使得权重更加客观。

```
1 | plot_weight(lr, x_train.columns)
```



多重共线性分析

权重的疑惑

从刚才的运行结果中，我们发现了一个非常严重的问题，citympg与highwaympg对目标price的影响起到了完全相反的结果！这是怎么回事呢？



概念

在多元线性回归中，自变量与因变量应该要存在线性关系。但是，如果自变量之间存在线性关系，我们称这种现象为**多重共线性**。多重共线性会造成权重的不确定性，甚至可能会出现拟合后的权重与相关系数呈现相反的结果。

不过，多重共线性并不会影响拟合优度（模型的效果），如果我们仅追求模型的拟合优度，而不考虑权重系数的大小，则可以不用处理多重共线性。

方差膨胀系数

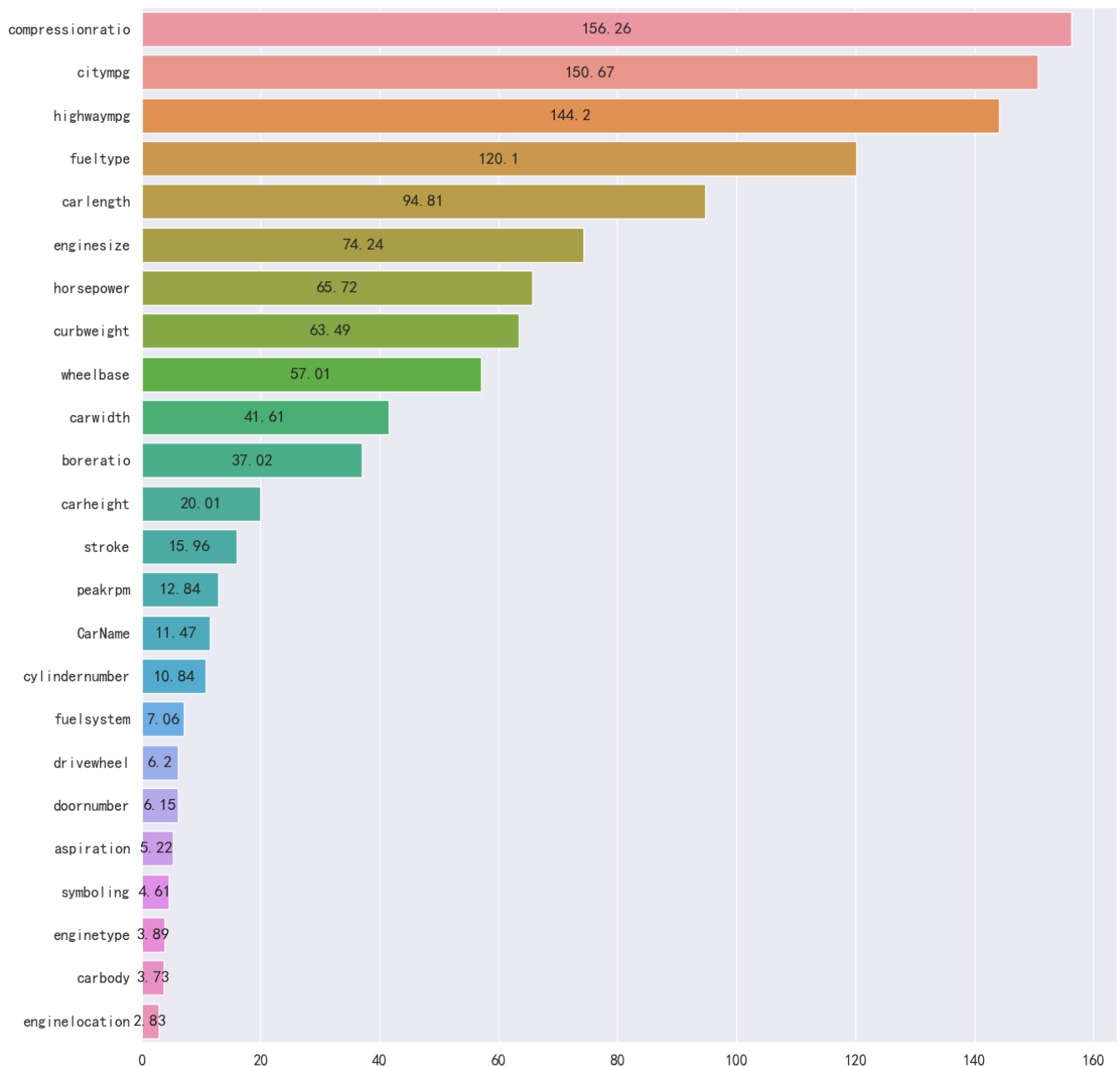
我们常使用方差膨胀系数（VIF——Variance Inflation Factor）来检测多元线性回归中的多重共线性。方差膨胀系数的计算公式如下：

$$VIF_i = \frac{1}{1 - R_i^2}$$

- VIF_i ：第*i*个变量的方差膨胀系数。
- R_i^2 ：第*i*个变量作为因变量，其他变量作为自变量进行回归时的 R^2 。

由公式可知，如果变量之间存在多重共线性，则得到的VIF值就会较高。当VIF值大于10时，说明存在严重的多重共线性。

```
1  from statsmodels.stats.outliers_influence import variance_inflation_factor
2
3  def plot_VIF(x, figsize=None):
4      """绘制特征变量的方差膨胀系数（VIF）。
5
6      将每个特征变量以条形图的方式绘制，并进行降序排列。
7
8      Parameters
9      -----
10     x : DataFrame, 形状为（样本数量， 特征数量）。
11         需要计算VIF的所有变量数据。
12     figsize : tuple, 形状为（图像宽度，图像高度），可选。
13         设置图像的宽度与高度。
14
15     """
16
17     # 对每一个特征变量，求解VIF值。
18     vif = [variance_inflation_factor(x.values, i) for i in range(x.shape[1])]
19     # 组成Series对象，方便对值排序时，可以与对应的列名相对应。
20     s = pd.Series(vif, index=x.columns)
21     s.sort_values(ascending=False, inplace=True)
22     s = s.round(2)
23     ax = sns.barplot(y=s.index, x=s.values)
24     ax.get_figure().set_size_inches(15, x.shape[1] // 1.5)
25     for y, x in enumerate(s):
26         ax.text(x / 2, y, x, va="center", ha="center")
27
28
29 plot_VIF(X_train_scale)
```



课堂练习



我们发现了一些VIF较大的特征，应该如何处理？

- A 将所有大于阈值的特征删除。
- B 将VIF最大的特征删除。
- C 将VIF最小的特征删除。
- D 不删除任何特征。



```

1 def feature_eliminate_by_VIF(X, max_vif=10, verbose=1):
2     """根据特征变量的方差膨胀系数（VIF），删除特征。
3
4     该方法每次会删除VIF最大的特征，直到所有的特征都不超过
5     指定的阈值为止。
6
7     Parameters
8     -----
9     X : DataFrame, 形状为（样本数量， 特征数量）。
10        需要删除特征变量的数据。
11     max_vif : float, 可选。
12        指定特征的最大VIF值，超过该值的特征将会删除。
13     verbose : int, 可选（1或0）。

```

```

14     是否显示详细的删除信息，默认为1。
15
16     Returns
17     -----
18     eli_col_names : list
19         所有因为VIF大于指定阈值而被删除的特征名。
20
21     """
22
23     # 对参数复制，避免修改影响到参数。
24     x = x.copy()
25     # 用来保存所有因为VIF过大而被删除的特征名。
26     eli_col_names = []
27     while x.shape[1] > 0:
28         vif = pd.DataFrame()
29         vif["features"] = x.columns
30         vif["vif"] = [variance_inflation_factor(x.values, i) for i in range(x.shape[1])]
31         vif = vif.sort_values("vif", ascending=False)
32         # 获取vif最大的行。
33         s = vif.iloc[0]
34         if s["vif"] > max_vif:
35             x.drop(s["features"], axis=1, inplace=True)
36             eli_col_names.append(s["features"])
37             if verbose:
38                 print(f'{s["features"]}的vif值过大（{s["vif"]:.2f}），删除！')
39         else:
40             # 如果已经没有VIF过大的特征，则显示最终的结果，退出循环。
41             display(vif)
42             break
43     return eli_col_names
44
45
46 drop_columns = feature_eliminate_by_VIF(X_train_scale, max_vif=10)
47 print(drop_columns)

```

```

1  compressionratio的vif值过大（156.26），删除！
2  highwaympg的vif值过大（143.99），删除！
3  carlength的vif值过大（89.89），删除！
4  enginesize的vif值过大（67.84），删除！
5  curbweight的vif值过大（53.34），删除！
6  wheelbase的vif值过大（43.58），删除！
7  horsepower的vif值过大（35.28），删除！
8  carwidth的vif值过大（24.59），删除！
9  boreratio的vif值过大（15.21），删除！
10 carheight的vif值过大（10.67），删除！

```

```

1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }

```

	features	vif
1	CarName	7.579864
13	citympg	6.661770
12	peakrpm	6.626759
11	stroke	6.533504
10	fuelsystem	5.156966
9	cylindernumber	5.105943
6	drivewheel	4.017209
0	symboling	3.434676
4	doornumber	3.270645
5	carbody	3.032210
2	fueltype	2.889652
8	enginetype	2.808392
3	aspiration	2.168334
7	enginelocation	1.744288

```
1 ['compressionratio', 'highwaympg', 'carlength', 'enginesize', 'curbweight', 'wheelbase', 'horsepower', 'carwidth', 'bore', 'stroke', 'carheight']
```

删除共线性特征

结果与之前的表现也是一致的，在数据集中，citympg与highwaympg的相关性很高，而curbweight与carlength，carwidth，enginesize等特征的相关性很高。这里，我们将删除多重共线性的特征。

```
1 X_train_scale.drop(drop_columns, axis=1, inplace=True)
2 X_test_scale.drop(drop_columns, axis=1, inplace=True)
3 lr.fit(X_train_scale, y_train_scale)
4 print(lr.score(X_train_scale, y_train_scale))
5 print(lr.score(X_test_scale, y_test_scale))
```

```
1 0.9339202978016744
2 0.852217003937153
```

```
1 plot_weight(lr, X_train_scale.columns)
```



特征选择

在建立模型时，特征并非越多越好，有些特征可能对模型质量并没有什么改善，我们可以进行删除。

RFECV

特征选择的方式有很多，这里，我们使用RFECV方法来实现特征选择。RFECV分成两个部分，分别是：

- RFE (Recursive feature elimination)：递归特征消除，用来对特征进行重要性评级。
- CV (Cross Validation)：交叉验证，在特征评级后，通过交叉验证，选择最佳数量的特征。

具体过程如下：

- RFE阶段。
 1. 初始的特征集为所有可用的特征。
 2. 使用当前特征集进行建模，然后计算每个特征的重要性。
 3. 删除最不重要的一个（或多个）特征，更新特征集。
 4. 跳转到步骤2，直到完成所有特征的重要性评级。
- CV阶段。
 1. 根据RFE阶段确定的特征重要性，依次选择不同数量的特征。
 2. 对选定的特征集进行交叉验证。
 3. 确定平均分最高的特征数量，完成特征选择。

```
1 from sklearn.feature_selection import RFECV
2
3 # estimator: 要操作的模型。
```



```

4 # step: 每次删除的变量数。
5 # cv: 使用的交叉验证折数。
6 # n_jobs: 并发的数量。
7 # scoring: 评估的方式。
8 rfecv = RFECV(estimator=lr, step=1, cv=5, n_jobs=-1, scoring="r2")
9 rfecv.fit(X_train_scale, y_train_scale)
10 # 返回经过选择之后, 剩余的特征数量。
11 print(rfecv.n_features_)
12 # 返回经过特征选择后, 使用缩减特征训练后的模型。
13 print(rfecv.estimator_)
14 # 返回每个特征的等级, 数值越小, 特征越重要。
15 print(rfecv.ranking_)
16 # 返回布尔数组, 用来表示特征是否被选择。
17 print(rfecv.support_)
18 # 返回对应数量特征时, 模型交叉验证的评分。
19 print(rfecv.cv_results_["mean_test_score"])

```

```

1 9
2 LinearRegression()
3 [1 1 3 1 4 1 1 1 5 1 2 1 6 1]
4 [ True  True False  True False  True  True  True False  True False  True
5  False  True]
6 [0.80375872 0.85457418 0.88819254 0.88920672 0.89296809 0.89665574
7  0.89781439 0.89971934 0.90554417 0.90312207 0.90114276 0.90220023
8  0.90311225 0.90348766]

```

我们也可以绘制下, 在特征选择过程中, 使用交叉验证获取的 R^2 值。

```

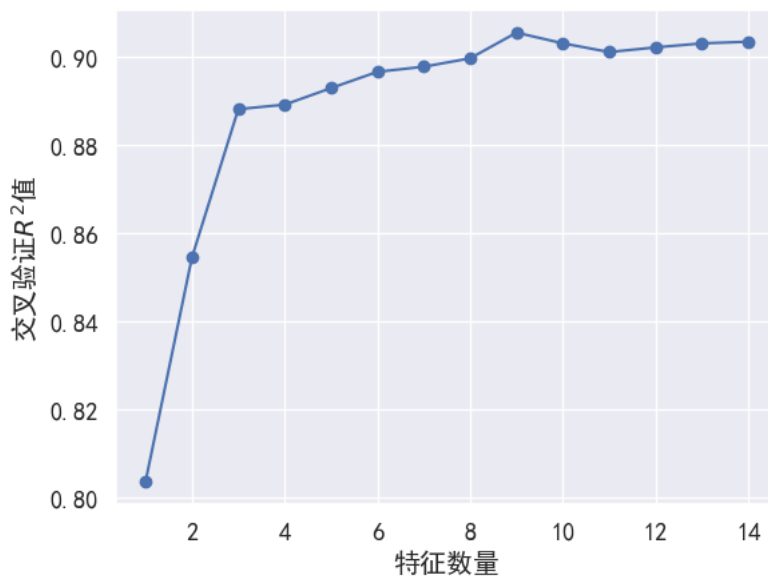
1 score = rfecv.cv_results_["mean_test_score"]
2 plt.plot(range(1, len(score) + 1), score, marker="o")
3 plt.xlabel("特征数量")
4 plt.ylabel("交叉验证$R^2$值")

```

```

1 Text(0, 0.5, '交叉验证$R^2$值')

```



特征选择后, 我们可以对现有的数据集进行转换。

```

1 print("剔除的变量: ", X_train_scale.columns[~rfecv.support_])
2 X_train_eli = rfecv.transform(X_train_scale)
3 X_test_eli = rfecv.transform(X_test_scale)
4 print(X_train_scale.shape, X_train_eli.shape)
5 # 说明: 经过transform转换之后, 返回ndarray数组类型。
6 # 这里X_train_scale是DataFrame类型, 如果依然希望返回DataFrame类型,
7 # 需要自己去调整。
8 col = X_train_scale.columns[rfecv.support_]
9 # X_train_eli = pd.DataFrame(X_train_eli, columns=col)
10 # X_test_eli = pd.DataFrame(X_test_eli, columns=col)
11 print(rfecv.estimator_.score(X_train_eli, y_train_scale))
12 print(rfecv.estimator_.score(X_test_eli, y_test_scale))

```

```
1 剔除的变量:  Index(['fueltype', 'doornumber', 'enginetype', 'fuelsystem', 'peakrpm'], dtype='object')
2  (143, 14) (143, 9)
3  0.9318950398900173
4  0.8521250777462107
```

删除变量分析

- doornumber, fueltype为不显著变量。
- peakrpm为相关性很低的变量。