

系统工程导论作业四——黑箱建模 2

彭程 2020011075

1. 试说明：病态线性回归问题中，显著性检验是否需要？

结论：

病态线性回归中需要显著性检验，且在自变量降维去线性之后进行检验。

需要进行检验：

显著性检验就是事先对总体（随机变量）的参数或总体分布形式做出一个假设，然后利用样本信息来判断这个假设（备择假设）是否合理，即判断总体的真实情况与原假设是否有显著性差异。在线性回归中体现为验证线性回归模型的有效性。病态线性回归问题也是利用样本信息来判断总体的参数，所以也需要显著性检验来确定总体分布形式假设是否合理，判断总体的真实情况与原假设是否有显著性差异，因此病态线性回归问题中也需要进行显著性检验。

降维之后检验：

病态回归问题中，由于自变量成线性相关或近似线性相关关系，所以自变量降维去线性之前，参数估计的误差可能会被严重放大，线性相关的变量的系数是不稳定的，因此即便是估计出来参数，意义也并不大，更没有必要进行显著性检验。而在自变量降维去线性之后，线性回归误差较小，因此应该在此时进行显著性检验。

2. 编程实现多元线性回归。

2.1 算法思路

1. 样本数据规范化，消除单位影响。

$$\bar{x}_i(t) = \frac{x_i(t) - e(x_i)}{\sqrt{\delta^2(x_i)}} \quad \forall i, t$$

其中， $e(x_i) = \frac{1}{N} \sum_{t=1}^N x_i(t)$ 为样本均值， $\delta^2(x_i) = \frac{1}{N-1} \sum_{t=1}^N (x_i(t) - e(x_i))^2$ 为样本方差

$$\bar{y}(t) = \frac{y(t) - e(y)}{\sqrt{\delta^2(y)}} \quad \forall t$$

2. 根据阈值确定最优维度 m

对归一化的 XX^T 进行特征值分解，得到 $XX^T = Q\Lambda Q^T$ 。按照特征值从大到小选取特征向量，直到相对逼近误差（未选取的特征值占特征值的比例）小于选取的阈值。

3. 求出降维后的矩阵。

将特征向量组成矩阵 $Q_m = [q(1)q(2)\cdots q(m)]$ ，从而降维后的矩阵为 $Z = Q_m^T X$ 。

4. 计算降维后的回归系数 d，从而确定降维前的回归系数 c，恢复得到归一化前的回归系数

$$\begin{aligned} \hat{d} &= (ZZ^T)^{-1} ZY^T \\ \hat{c} &= Q_m \hat{d} \end{aligned}$$

5. 进行显著性检验

对计算得到的回归系数进行显著性检验，注意其中 n 为降维之后的维度 m ：

$$F = \frac{(N - n - 1)ESS}{nRSS}$$

7. 求取置信区间

给定显著性水平 α ，对某一 x_0 ，相应的 y_0 将以 $1 - \alpha$ 的概率落在置信区间：

$$(\hat{y}_0 - Z_{\alpha/2}S_\delta, \hat{y}_0 + Z_{\alpha/2}S_\delta)$$

其中 $Z_{\alpha/2}$ 是标准正态分布上 $\alpha/2$ 百分位点的值，剩余均方差 $S_\delta = \sqrt{\frac{RSS}{N-n-1}}$ 。

2.2 实验结果

```
特征值从大到小依次为： [21.06874612 11.2273649  7.68243347  0.02145552]
降维后的维度为： 3
降维后的系数为： [ 0.66742279 -0.00376926 -0.25422719]
规范化后的系数为： [0.48232919 0.21752131 0.00671087 0.47968412]
原始系数为： [0.0729661  0.59856205 0.00187169 0.10548166]
原始偏移为： -9.1515
F检验结果为： F = 195.1165 > F_alpha = 4.3468 说明x与y存在线性关系
置信区间为： (y - 1.1557, y + 1.1557)
回归方程为： y = -9.1515 + 0.0730x1 + 0.5986x2 + 0.0019x3 + 0.1055x4
```

回归方程: $y = -9.1515 + 0.0730x_1 + 0.5986x_2 + 0.0019x_3 + 0.1055x_4$

F 检验: $F = 195.1165 > 4.3468$ ，因此 x, y 存在线性关系。

置信区间: $[y - 1.1557, y + 1.1557]$

2.3 代码

```
1 import numpy as np
2 from scipy import stats
3
4 def linear_regression1(Y, X, alpha=0.05, error=0.01):
5     """
6     输入: N×N的矩阵X, 1×N的矩阵Y
7     输出: N×1的回归系数THETA
8     功能: 实现Y=C^T*X+B的多元线性回归, 自适应病态线性回归
9     """
10    n, N = X.shape
11
12    # 数据规范化
13    x_mean = np.mean(X, axis=1, keepdims=True)
14    delta_x = np.sqrt(
15        np.sum(np.multiply(X - np.mean(X, axis=1, keepdims=True), X - np.mean(X, axis=1, keepdims=True)),
16              axis=1,
17              keepdims=True) / (N - 1))
18    X_normal = np.divide(X - x_mean, delta_x)
```

```

18 y_mean = np.mean(Y)
19 delta_y = np.sqrt(np.sum(np.multiply(Y - y_mean, Y - y_mean)) / (N - 1))
20 Y_normal = (Y - y_mean) / delta_y
21
22 # 根据阈值确定最优维度m
23 eigenvalue, eigenvector = np.linalg.eig(np.dot(X_normal, X_normal.T))
24 eigen_index = np.argsort(-eigenvalue) # 返回排序的下标
25 eigen_sum = np.sum(eigenvalue)
26 for m in range(0, n):
27     eigen_error = 0
28     for i in range(0, m):
29         eigen_error += eigenvalue[eigen_index[n - 1 - i]]
30     eigen_error = eigen_error / eigen_sum
31     if eigen_error > error:
32         break
33 m = n - m + 1 # 则前m项是要保留的
34
35 # 计算降维后的回归系数D, 从而确定降维前的回归系数c_n, 恢复得到归一化前的回归系数c
36 Q = eigenvector[:, eigen_index[0: m]]
37 Z = Q.T.dot(X_normal)
38 d = np.linalg.inv(Z.dot(Z.T)).dot(Z.dot(Y_normal.T))
39 c_n = Q.dot(d)
40 c = (c_n.T / np.squeeze(delta_x) * delta_y).T
41 b = y_mean - np.sum(c * np.squeeze(x_mean))
42 print("特征值从大到小依次为: ", eigenvalue[eigen_index[:]])
43 print("降维后的维度为: ", m)
44 print("降维后的系数为: ", d)
45 print("规范化后的系数为: ", c_n)
46 print("原始系数为: ", c)
47 print("原始偏移为: {:.4f}".format(b))
48
49 # F检验操作
50 Y_estimate = np.dot(c.T, X) + b
51 ESS = np.dot((Y_estimate - y_mean), (Y_estimate - y_mean).T)
52 RSS = np.dot((Y - Y_estimate), (Y - Y_estimate).T)
53 F = ((N - m - 1) * ESS) / (m * RSS) # ESS自由度为 (N-m-1), RSS自由度为m
54 F_alpha = stats.f.isf(alpha, m, N - m - 1)
55 if F > F_alpha:
56     print("F检验结果为: F = {:.4f} > F_alpha = {:.4f}".format(F, F_alpha), " 说明x与y存在线性关系")
57 elif F <= F_alpha:
58     print("F检验结果为: ")
59     print("F-value = {} <= F_alpha = {}".format(F, F_alpha))
60     print("x与y不存在线性关系")
61     return 0
62
63 # 置信区间
64 S_delta = np.sqrt(RSS / (N - m - 1))
65 Z_alpha_div2 = stats.norm.isf(alpha / 2, 0, 1)
66 interval = Z_alpha_div2 * S_delta
67 print("置信区间为: (y - {:.4f}, y + {:.4f})".format(interval, interval))
68
69 # 打印回归方程
70 equation = "y = {:.4f}" % b
71 for i in range(n):
72     equation += " + {:.4f}X%d" % (c[i], i + 1)
73 print("回归方程为: ", equation)

```

```
74
75     return
76
77
78 if __name__ == '__main__':
79     f = open('data.txt', 'r', encoding='utf-8')
80     # RAW_DATA = NP.ARRAY([LIST(MAP(FLOAT, I[:-1].SPLIT(' '))) FOR I IN F.READLINES()[1:]]
81     raw_data = np.array([i[:-1].split(' ') for i in f.readlines()[1:]], dtype=float)
82     x = raw_data[:, 1: 5]
83     y = raw_data[:, -1]
84     linear_regression1(y.T, x.T, 0.05, 0.01)
```