

Министерство науки и высшего образования Российской Федерации
Новосибирский государственный технический университет
Кафедра теоретической и прикладной информатики

Статистические методы анализа данных
Лабораторная работа №5

Факультет:	ФПМИ
Группа:	ПМ-63
Студенты:	Кожекин М.В. Майер В.А. Назарова Т.А. Утюганов Д.С.
Вариант:	1

Новосибирск

2019

1. Цель работы

Разработать программу для определения мультиколлинеарности модели.

2. Решение

1. В соответствии с вариантом задания сгенерировать экспериментальные данные, в которых в явном виде присутствует эффект мультиколлинеарности.

Согласно заданию регрессия на 6 факторах. Эффект мультиколлинеарности создают 3 фактора. Имеется разброс в масштабах факторов.

$$\begin{aligned}\eta &= f^T(x) = x_1 + x_2 + x_3 + x_4 + x_5 + x_6 \\ x_1 &\in [-1, 1] \\ x_2 &\in [-2, 2] \\ x_3 &\in [-100, 100] \\ x_4 &\in [-10, 10] \\ x_5 &\in [-1, 1] \\ x_6 &= 2x_4 - 3x_5 + N(0, 0.01) \\ \rho &= 0.15, n = 100, m = 6\end{aligned}$$

2. Рассчитать ряд показателей, характеризующих эффект мультиколлинеарности. Определить факторы, ответственные за возникновение эффекта мультиколлинеарности.

В качестве мер измерения эффекта мультиколлинеарности рассмотрим следующее:

1. определитель информационной матрицы

$$|X^T X| = \prod_{i=1}^m \lambda_i = 59987206$$

2. минимальное собственное число матрицы

$$\lambda_1 = \lambda_{\min}(X^T X) = 2.1923124189023676e - 06$$

3. мера обусловленности матрицы по Нейману-Голдстейну

$$\frac{\lambda_{\max}(X^T X)}{\lambda_{\min}(X^T X)} = 704048310351$$

4. максимальная парная сопряжённость

$$\max_{i,j} x|r_{ij}|, i \neq j$$

$R = (r_{ij})$ – матрица сопряжённости

$$r_{ij} = \cos(\underline{x}_i, \underline{x}_j) = \sum \frac{X_{it}X_{jt}}{\sqrt{(\sum X_{it}^2 \sum X_{jt}^2)}}$$

$$\max_{i,j} x|r_{ij}| = 0.082419$$

5. максимальная сопряжённость

В качестве меры мультиколлинеарности возьмём

$$\max_i |R_i|, \text{ где:}$$

$$R_i = \sqrt{1 - \frac{1}{R_{ii}^{-1}}}, i = 1, m$$

$$\max_i |R_i| = 0.926$$

3. Построить ридж-оценки параметров при различных значениях параметров регуляризации. Выбрать оптимальное значение параметра регуляризации. Построить графики изменения квадрата евклидовой нормы оценок параметров и остаточной суммы квадратов от параметра регуляризации.

Один из способов оценивания параметров в условиях мультиколлинеарности состоит в управлении масштабом полученных оценок. Однако это смещение значительно меньше, чем у обычных МНК оценок.

С целью управления масштабом оценок введём в рассмотрение функцию стоимости

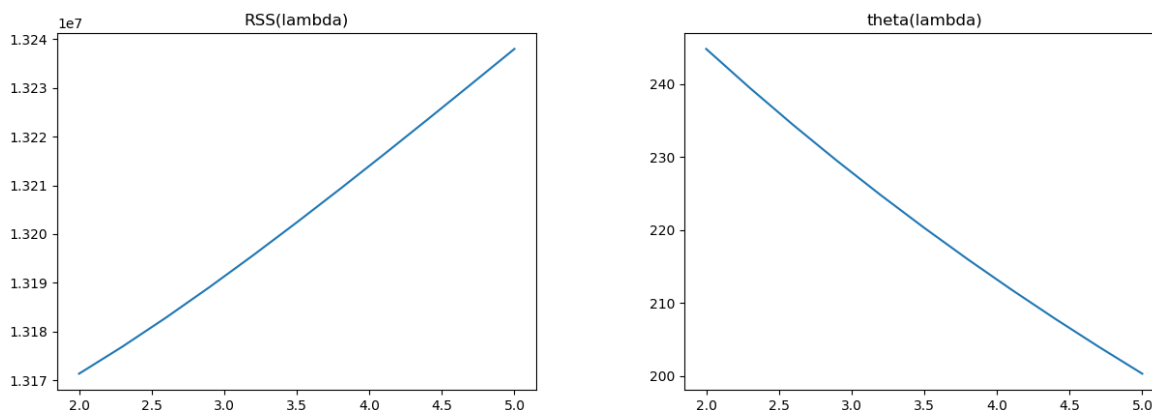
$$C = \sum_{i=1}^n (y_i - f(X_i)\theta)^2 + \sum_{j=1}^m \lambda_j \theta_j^2$$

где второе слагаемое рассматривается как штраф при условии, что $\lambda_j \geq 0$

Минимизируя C , получим $\hat{\theta} = (X^T X + \Lambda)^{-1} X^T y$, которые известны как ридж-оценки. Часто матрицу Λ задают диагональной в виде $\Lambda_{ii} = \lambda (X^T X)_{ii}$, $\lambda \geq 0$.

Оптимальное значение $\lambda = m\hat{\sigma}^2 / \|\hat{\beta}\|^2$

Графики:



4. Провести оценивание модели регрессии по методу главных компонент. Перейти к описанию в исходном пространстве факторов. Сравнить решение с ридж-оцениванием по смещению оценок и точности предсказания отклика.

Для начала мы центрируем нашу исходную модель

$$y_t = \Theta_0 + \Theta_1 x_{1t} + \dots + \Theta_k x_{kt} + E_t$$

Тогда модель наблюдения будет иметь вид

$$y_t^* = \beta_0 + \beta_1 x_{1t}^* + \dots + \beta_k x_{kt}^* + E_t$$

Т.к. меру изменчивости можно измерить при помощи собственных значений, то мы можем получить матрицу главных компонент через матрицу

$$V = (v_1, \dots, v_k) \\ Z = X^* V$$

После отбора главных компонент можно оценить регрессию $y^* : y_t^* = b_1 z_{1t} + \dots + b_l z_{lt} + E_t$

theta = -55.3567551900075, 41.0306349638696, 3.23326603674220, 15.9659268561517, 297.857791838169, 7.93397529930516

3. Текст программы

Файл run_lab5.py

```
1 import os
2
3 #os.system('python 1/lab1.py 0.15 6 False')
4 os.system('python 5/lab5.py 0.15')
```

Файл model_6parameters.py

```
1 # файл параметров модели
2 import numpy as np
3
4
5 # параметры уравнения
6 #  $th_0 \cdot x_1 + th_1 \cdot x_2 + th_2 \cdot x_3 + th_3 \cdot x_4 + th_4 \cdot x_5 + th_5 \cdot x_6$ 
7 theta_true = [1, 1, 1, 1, 1, 1]
8 n, m = 30, 6 # число точек, параметров и осей
9 # a, b = -1, 1
10 a = [-10, -2, -100, -10, -1]
11 b = [ 10, 2, 100, 10, 1]
12
13 #  $x_6 = 2 \cdot x_2 + x_3 + 4 \cdot x_4$ 
14 rho = 0.15 # шум в диапазоне [0.05, 0.15] или [0.50, 0.70]
15
16
17 def sample_x():
18     x = np.ndarray((m,n))
19     x[0] = np.random.uniform(low=a[0], high=b[0], size=n)
20     x[1] = np.random.uniform(low=a[1], high=b[1], size=n)
21     x[2] = np.random.uniform(low=a[2], high=b[2], size=n)
22     x[3] = np.random.uniform(low=a[3], high=b[3], size=n)
23     x[4] = np.random.uniform(low=a[4], high=b[4], size=n)
24     x[5] = 2*x[0] - 3*x[3] + np.random.uniform(0, 0.01)
25     # print(x[5])
26     return x.transpose()
27
28
29 def f(th, x_vector):
30     return th[0] * x_vector[0] + \
31            th[1] * x_vector[1] + \
32            th[2] * x_vector[2] + \
33            th[3] * x_vector[3] + \
34            th[4] * x_vector[4] + \
35            th[5] * x_vector[5]
36
37
38 def f_vector(x_vector):
39     return [x_vector[0],
40            x_vector[1],
41            x_vector[2],
42            x_vector[3],
43            x_vector[4],
44            x_vector[5]]
45
46
47 def f_vector(th, x_vector):
48     return [th[0] * x_vector[0],
49            th[1] * x_vector[1],
```

```

50         th[2] * x_vector[2],
51         th[3] * x_vector[3],
52         th[4] * x_vector[4],
53         th[5] * x_vector[5] ]

```

Файл lab1.py

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from mpl_toolkits.mplot3d import Axes3D
4 from scipy import stats
5 from sympy import *
6 from sys import argv, path
7
8
9 #
10 #
11 #
12 class lab1():
13     def __init__(self):
14         self.x = sample_x()
15         self.u = np.ndarray(n)
16         for i in range(n):
17             self.u[i] = f(theta_true, self.x[i])
18         u_mean = np.full((n), np.mean(self.u))
19         w_squared = np.dot((self.u - u_mean), (self.u - u_mean)) / (n-1)
20         self.dist = rho * w_squared
21         self.y = np.copy(self.u)
22         for i in range(n):
23             self.y[i] += np.random.normal(0, self.dist)
24
25     def save_table_1(self):
26         with open('1/report/table_1_u_y_' + str(int(rho*100)) + '.txt', 'w') as
file:
27             file.write('i\t')
28             for i in range(1, m+1):
29                 file.write('x%d\t' % i)
30             file.write('u\ty\n')
31             for i in range(n):
32                 file.write('{:d}\t'.format(i))
33                 for j in range(m):
34                     file.write('{:.17f}\t'.format(self.x[i][j]))
35
36                 file.write('{:.17f}\t{:.17f}\n'.format(self.u[i], self.y[i]))
37
38     def draw(self, doDrawWithNoize, title, name, elevation, azimuth):
39         path_to_save = '1/pics/' + name + '_' + \
40             str(int(rho*100)) + '_' + str(elevation) + '_' + str(azimuth) + '.
png'
41         fig = plt.figure('1')
42         ax = fig.gca(projection='3d')
43         tmp_range = np.linspace(a, b, sqrt(n))
44         x1, x2 = np.meshgrid(tmp_range, tmp_range)
45         u = f_2_parameters(theta_true, x1, x2)
46         fig.subplots_adjust(bottom=-0.05, top=1, left=-0.05, right=1.05)
47         ax.view_init(elevation, azimuth)
48         ax.plot_surface(x1, x2, u, zorder=2, alpha=0.2)
49         if doDrawWithNoize:
50             title += ' ' + str(int(rho * 100)) + '%'
51             ax.scatter(self.x1, self.x2, self.y, c='black', zorder=1)
52         else:

```

```

53         ax.scatter(self.x1, self.x2, self.u, c='black', zorder=1)
54     plt.title(title, fontsize=19)
55     plt.xlabel('x1')
56     plt.ylabel('x2')
57     plt.grid(alpha=0.5)
58     plt.savefig(path_to_save)
59     plt.clf()
60
61
62 #
63 #
64 #
65 if __name__ == "__main__":
66     path.insert(1, '')
67     global rho
68     rho = float(argv[1])
69
70     params_count = int(argv[2])
71     if params_count == 3:
72         from model_3parameters import *
73     elif params_count == 4:
74         from model_4parameters import *
75     elif params_count == 6:
76         from model_6parameters import *
77
78     doDraw = bool(argv[3])
79
80     print('Запущен код 1 лабораторной работы: {:d} параметров, шум {:d}%'.format(int(
81         argv[2]), int(rho*100)))
82
83     l1 = lab1()
84     l1.save_table_1()
85     doDraw = False
86     if doDraw:
87         l1.draw(False, 'исходная модель', 'before', None, None)
88         l1.draw(False, 'исходная модель', 'before', 0, 0)
89         l1.draw(False, 'исходная модель', 'before', 0, 90)
90         l1.draw(True, 'помеха', 'after', None, None)
91         l1.draw(True, 'помеха', 'after', 0, 0)
92         l1.draw(True, 'помеха', 'after', 0, 90)

```

Файл lab5.py

```

1 import matplotlib.pyplot as plt
2 #import numpy as np
3 from mpl_toolkits.mplot3d import Axes3D
4 from scipy import stats
5 from sympy import *
6 from sys import argv, path
7 import math
8
9
10 def dotproduct(v1, v2):
11     return sum((a*b) for a, b in zip(v1, v2))
12
13
14 def length(v):
15     return math.sqrt(abs(dotproduct(v, v)))
16
17
18 def angle2(v1, v2):

```

```

19     return math.acos(abs(dotproduct(v1, v2)) / (length(v1) * length(v2)))
20
21
22 def angle(v1, v2):
23     return np.dot(a, b) / (np.linalg.norm(a)*np.linalg.norm(b))
24
25 #

```

```

26 #

```

```

27 #

```

```

28
29
30 class lab5():
31     def __init__(self):
32         self.x = np.ndarray((n, m))
33         self.X = np.ndarray((n, m))
34         self.u = np.ndarray(n)
35         self.y = np.ndarray(n)
36
37     with open('1/report/table_1_u_y_' + str(int(rho*100)) + '.txt', 'r') as
file:
38         line = file.readline().rstrip()
39         for i in range(n):
40             line = file.readline().rstrip().rsplit()
41             for j in range(m):
42                 self.x[i][j] = float(line[j + 1])
43                 self.u[i] = float(line[m + 1])
44                 self.y[i] = float(line[m + 2])
45
46         for i in range(n):
47             self.X[i] = f_vector_(self.x[i])
48
49     def calc_det_of_inf_matrix(self):
50         X = Matrix(self.X)
51         self.inf_matr = np.array(X.T * X).astype(np.float64)
52         d = np.linalg.det(self.inf_matr)
53         print('det = ', d)
54
55     # def calc_det_of_inf_matrix(self):
56     #     X = Matrix(self.X)
57     #     self.inf_matr_sp = X.T * X
58     #     self.inf_matr_sp = self.inf_matr_sp / Trace(self.inf_matr_sp)
59     #     print(self.inf_matr_sp)
60     #     self.inf_matr = np.array(self.inf_matr_sp).astype(np.float64)
61     #     d = np.linalg.det(self.inf_matr)
62     #     print('det = ', d)
63
64     def calc_min_eig_val(self):
65         self.eig_vals = np.linalg.eigvals(self.inf_matr)
66         print('lambda_min = ', np.min(self.eig_vals))
67
68     def calc_cond_number(self):
69         self.cond_number = np.max(self.eig_vals) / np.min(self.eig_vals)
70         print('cond_number = ', self.cond_number)
71

```

```

72 def calc_pair_sopr(self):
73     X = self.X.transpose()
74     self.R = np.ndarray((m, m))
75     self.pair_sopr = -9000.0
76     for i in range(m):
77         for j in range(m):
78             self.R[i][j] = 0.0
79             for t in range(n):
80                 numerator = X[i][t] * X[j][t]
81                 denominator1 = 0.0
82                 for t1 in range(n):
83                     denominator1 += X[i][t1]**2
84                 denominator2 = 0.0
85                 for t2 in range(n):
86                     denominator2 += X[i][t2]**2
87                 self.R[i][j] += numerator / np.sqrt(denominator1*denominator2)
88             self.R[i][i] = 1.0
89
90     #print(max(abs(self.R.max()), abs(self.R.min()))))
91
92     for i in range(m):
93         for j in range(m):
94             if i != j and (abs(self.R[i][j]) > self.pair_sopr):
95                 self.pair_sopr = abs(self.R[i][j])
96
97     # print(self.R)
98     print('pair_sopr = %f' % self.pair_sopr)
99
100 # def calc_pair_sopr(self):
101 #     X = self.x.transpose()
102 #     # print(X)
103 #     self.R = np.ndarray((m, m))
104 #     self.pair_sopr = -9000
105 #     for i in range(m):
106 #         for j in range(m):
107 #             self.R[i][j] = angle(X[i], X[j])
108 #             self.R[i][i] = 1.0
109
110 #     for i in range(m):
111 #         for j in range(m):
112 #             if i != j and abs(self.R[i][j]) > self.pair_sopr:
113 #                 self.pair_sopr = abs(self.R[i][j])
114
115 #     # print(self.R)
116 #     print('pair_sopr = %d' % self.pair_sopr)
117
118 def calc_max_sopr(self):
119     # R_inv = Matrix(self.R)**-1
120     R_inv = np.linalg.inv(self.R)
121     self.R_list = np.ndarray(m)
122     # print(R_inv)
123     self.max_sopr = -9000.0
124     for i in range(m):
125         R_i = sqrt(1.0 - (1.0 / R_inv[i][i]))
126         if R_i > self.max_sopr and R_i != float('inf'):
127             self.max_sopr = R_i
128     #print(max(abs(self.R.max()), abs(self.R.min()))))
129
130     print('max_sopr = %f' % self.max_sopr)
131

```



```

132 def calc_theta_ridge(self):
133     X = Matrix(self.X)
134     y = Matrix(self.y)
135     _theta_true = Matrix(theta_true)
136     points = 11
137     RSS = np.ndarray(points)
138     theta_norm = np.ndarray(points)
139     theta_calc_ridge = np.ndarray((points, m))
140     lambda_params = np.linspace(2, 5, points)
141     i = 0
142     for lambda_param in lambda_params:
143         lambdas = Matrix(np.diag(np.full(m, lambda_param)))
144         _th = ((X.T*X + lambdas) ** -1) * X.T * y
145         _RSS = (y - X*_th).T * (y - X*_th)
146         _theta_norm = (_th - _theta_true).norm()
147
148         RSS[i] = np.array(_RSS).astype(np.float64)
149         theta_norm[i] = np.array(_theta_norm).astype(np.float64)
150         theta_calc_ridge[i] = np.array(_th.T).astype(np.float64)
151         i += 1
152
153     path_to_save = '5/pics/RSS_from_lambda.png'
154     plt.title('RSS(lambda)')
155     plt.plot(lambda_params, RSS)
156     plt.savefig(path_to_save)
157     plt.clf()
158
159     path_to_save = '5/pics/theta_from_lambda.png'
160     plt.title('theta(lambda)')
161     plt.plot(lambda_params, theta_norm)
162     plt.savefig(path_to_save)
163     plt.clf()
164
165 def calc_theta_main_components(self):
166     X = self.X.T
167     for i in range(m):
168         X[i] -= np.mean(X[i])
169     X_centered = Matrix(X.T)
170     X_inf_matr_centered = X_centered.T * X_centered
171     X_centered_np = np.array(X_inf_matr_centered).astype(np.float64)
172
173     V = Matrix(np.diag(np.linalg.eigvals(X_centered_np)))
174     y = Matrix(self.y)
175     Z = X_centered * V
176     b = ((Z.T*Z) ** -1) * Z.T * y
177     self.theta_calc_mc = V.dot(b)
178     print(self.theta_calc_mc)
179
180
181 #


---


182 #


---


183 #


---


184 if __name__ == "__main__":
185     path.insert(1, '')

```

```

186     from model_6parameters import *
187
188     global rho
189     rho = float(argv[1])
190     print('Запущен код 5 лабораторной работы: {:d} параметров, шум {:d}%'.format(
191         m, int(rho*100)))
192     l5 = lab5()
193     l5.calc_det_of_inf_matrix()
194     l5.calc_min_eig_val()
195     l5.calc_cond_number()
196     l5.calc_pair_sopr()
197     l5.calc_max_sopr()
198     l5.calc_theta_ridge()
199     l5.calc_theta_main_components()

```