

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.

$$\eta = f^T(x) = x_1 + x_2 + x_3 + x_4 + x_5 + x_6$$

$$x_1 \in [-1, 1]$$

$$x_2 \in [-1, 1]$$

$$x_3 \in [-1, 1]$$

$$x_4 \in [-1, 1]$$

$$x_5 \in [-1, 1]$$

$$x_6 = 2x_4 - 3x_5 + N(0, 0.01)$$

$$\rho = 0.15, n = 100, m = 6$$

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

$$|X^T X| = \prod_{i=1}^m \lambda_i = 7.22$$

2. минимальное собственное число матрицы

$$\lambda_1 = \lambda_{\min}(X^T X) = 4.77e - 06$$

3. мера обусловленности матрицы по Нейману-Голдстейну

$$\frac{\lambda_{\max}(X^T X)}{\lambda_{\min}(X^T X)} = 34699460.47$$

4. максимальная парная сопряжённость

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

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

$$r_{ij} = \cos(\underline{x}_i, \underline{x}_j) = \frac{(x_i, x_j)}{|x_i||x_j|}$$

$$\max_{i,j} |r_{ij}| = 0.84$$

5. максимальная сопряжённость

В качестве меры мультиколлинеарности возьмём

$$\max_i |R_i|, \text{ где :}$$

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

$$\max_i |R_i| = 0.99$$

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

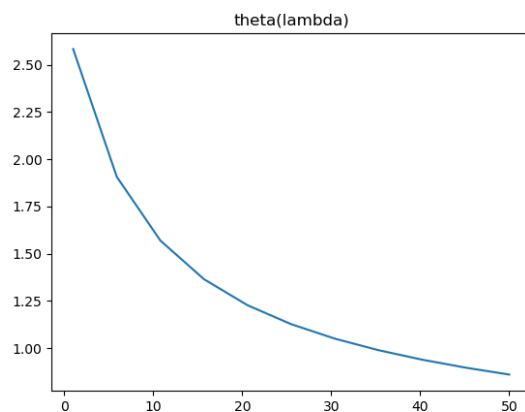
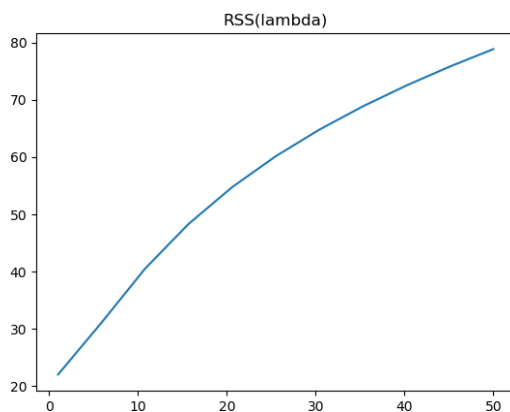
$$C = \sum_{i=1}^n (y_i - f(X_i)\theta)^2 + \sum_{j=1}^m \lambda_j \theta_j^2$$

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^* на главные компоненты:

$$y_t^* = b_1 z_{1t} + \dots + b_l z_{lt} + E_t \\ \hat{b} = (Z^T Z)^{-1} Z^T y^*$$

$$\theta = 1, 1, 1, 1, 1, 1$$

$$\hat{\theta} = 1.78, 3.58, 0.56, 1.30, 1.13, 0.63$$

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 # th0*x1 + th1*x2 + th2*x3 + th3*x4 + th4*x5 + th5*x6
7 # x6 = 2*x2 + x3 + 4*x4
8 theta_true = [1, 1, 1, 1, 1, 1]
9 n, m = 30, 6 # число точек, параметров и осей
10 rho = 0.15 # шум в диапазоне [0.05, 0.15] или [0.50, 0.70]
11 a = [-1, -1, -1, -1, -1]
12 b = [ 1, 1, 1, 1, 1]
13
14
15 def sample_x():
16     x = np.ndarray((m,n))
17     x[0] = np.random.uniform(low=a[0], high=b[0], size=n)
18     x[1] = np.random.uniform(low=a[1], high=b[1], size=n)
19     x[2] = np.random.uniform(low=a[2], high=b[2], size=n)
20     x[3] = np.random.uniform(low=a[3], high=b[3], size=n)
```

```

21     x[4] = np.random.uniform(low=a[4], high=b[4], size=n)
22     x[5] = 2*x[0] - 3*x[3] + np.random.uniform(0, 0.01)
23     return x.transpose()
24
25
26 def f(theta, x):
27     return np.dot(theta, x)
28
29 def f_vector_T(x):
30     return np.array(x)

```

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, doDrawWithNoise, 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 numpy.linalg import inv, det, eigvals, norm
4 from sys import argv, path
5
6
7 def cos(a, b):
8     ''' косинус между углами (a,b) = |a|*|b|*cos(a,b) '''
9     return np.dot(a, b) / (norm(a)*norm(b))
10
11 #
12 #

```

```

13 #
14 class lab5():
15     def __init__(self):
16         ''' Выделение памяти под массивы '''
17         self.x = np.ndarray((n, m))
18         self.X = np.ndarray((n, m))
19         self.u = np.ndarray(n)
20         self.y = np.ndarray(n)
21         self.R = np.ndarray((m, m))
22
23     def read_from_file(self, filename):
24         ''' Считываем данные из файла '''
25         data = np.loadtxt(filename, float, skiprows=1).transpose()
26         self.x = data[1:m+1].transpose()
27         self.u = data[-2].transpose()
28         self.y = data[-1].transpose()
29         for i in range(n):
30             self.X[i] = f_vector_T(self.x[i])
31
32     def calc_det_of_inf_matrix(self):
33         ''' Вычисляем определитель информационной матрицы '''
34         X = self.X
35         self.inf_matr = X.T @ X
36         d = det(self.inf_matr)
37         print('определитель информационной матрицы: ', d)
38
39     def calc_min_eig_val(self):
40         ''' Поиск минимального собственного значения '''
41         self.eig_vals = eigvals(self.inf_matr)
42         self.min_eig_val = np.min(self.eig_vals)
43         self.max_eig_val = np.max(self.eig_vals)
44         print('минимальное собственное значение: ', self.min_eig_val)
45
46     def calc_cond_number(self):
47         ''' Мера обусловленности по НеймануГолдстейну- '''
48         self.cond_number = self.max_eig_val / self.min_eig_val
49         print('число обусловленности: ', self.cond_number)
50
51     def calc_max_pair_sopr(self):
52         '''
53         Определяем максимальную парную сопряжённость.
54         3 вектора могут быть коллинеарны, а попарно нет
55         Максимальный коэффициент сопряжённости не несёт на себе эффекта масштаба
56         '''
57         X = self.X.transpose()
58         self.max_pair_sopr = -9000.0
59         # строим матрицу сопряжённости
60         for i in range(m):
61             for j in range(m):
62                 self.R[i][j] = cos(X[i], X[j])
63                 self.R[i][i] = 1.0
64
65         for i in range(m):
66             for j in range(m):
67                 if i != j and (abs(self.R[i][j]) > self.max_pair_sopr):
68                     self.max_pair_sopr = abs(self.R[i][j])
69         print('максимальная парная сопряжённость: ', self.max_pair_sopr)
70
71     def calc_max_sopr(self):
72         ''' Максимальная сопряжённость '''

```

```

73     R_inv = inv(self.R)
74     self.max_sopr = -9000.0
75     for i in range(m):
76         R_i = np.sqrt(1.0 - (1.0 / R_inv[i][i]))
77         if R_i > self.max_sopr and R_i != float('inf'):
78             self.max_sopr = R_i
79     print('максимальная сопряжённость: ', self.max_sopr)
80
81 def calc_theta_ridge(self):
82     '''
83     Ридж-оценка—
84     Мультиколлинеарные( данные выглядят как хребты на карте)
85     '''
86     X = self.X
87     y = self.y
88     points = 11
89     RSS = np.ndarray(points)
90     theta_norm = np.ndarray(points)
91     theta_calc_ridge = np.ndarray((points, m))
92     lambda_params = np.linspace(1, 50, points)
93     i = 0
94     for lambda_param in lambda_params:
95         lambdas = np.diag(np.full(m, lambda_param))
96         theta = inv(X.T @ X + lambdas) @ X.T @ y
97         RSS[i] = (y - X @ theta).T @ (y - X @ theta)
98         theta_norm[i] = norm(theta)
99         theta_calc_ridge[i] = theta.T
100        i += 1
101
102    # print(theta_calc_ridge)
103
104    path_to_save = '5/pics/RSS_from_lambda.png'
105    plt.title('RSS(lambda)')
106    plt.plot(lambda_params, RSS)
107    plt.savefig(path_to_save)
108    plt.clf()
109
110    path_to_save = '5/pics/theta_from_lambda.png'
111    plt.title('theta(lambda)')
112    plt.plot(lambda_params, theta_norm)
113    plt.savefig(path_to_save)
114    plt.clf()
115
116 def calc_theta_main_components(self):
117     '''
118     Метод главных компонент
119     '''
120     X = self.X.T
121     y = self.y
122     for i in range(m):
123         X[i] -= np.mean(X[i])
124     X_centered = X.T
125     X_inf_matr_centered = X_centered.T @ X_centered
126
127     V = np.diag(eigvals(X_inf_matr_centered))
128     Z = X_centered @ V
129     b = inv(Z.T @ Z) @ Z.T @ y
130     self.theta_calc_mc = V.dot(b)
131     print(self.theta_calc_mc)
132

```



```
133
134 #
135 #
136 #
137 if __name__ == "__main__":
138     path.insert(1, '')
139     from model_6parameters import *
140
141     global rho
142     rho = float(argv[1])
143     print('5 Лабораторная работа:', m, 'параметров, шум', int(rho*100), '\n')
144
145     l5 = lab5()
146     l5.read_from_file('1/report/table_1_u_y_' + str(int(rho*100)) + '.txt')
147     l5.calc_det_of_inf_matrix()
148     l5.calc_min_eig_val()
149     l5.calc_cond_number()
150     l5.calc_max_pair_sopr()
151     l5.calc_max_sopr()
152     l5.calc_theta_ridge()
153     l5.calc_theta_main_components()
```