

МИНОБРНАУКИ РОССИИ
САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ
ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ «ЛЭТИ»
ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)
Кафедра алгоритмической математики

ОТЧЁТ

по практической работе №5

по дисциплине «Статистический анализ»

**Тема: Элементы регрессионного анализа. Выборочные прямые
среднеквадратической регрессии. Корреляционное отношение**

Студент гр. 9372

Иванов Р. С.

Преподаватель

Сучков А. И.

Санкт-Петербург
2021

Цель работы

Ознакомление с основными положениями метода наименьших квадратов (МНК), со статистическими свойствами МНК оценок, с понятием функции регрессии и роли МНК в регрессионном анализе, с корреляционным отношением, как мерой тесноты произвольной (в том числе и линейной) корреляционной связи.

Основные теоретические положения

Алгоритм поиска линейной зависимости методом наименьших квадратов

Даны пары чисел (x_i, y_i) , которые расположены вдоль прямой линии $y = kx + b$. Угловым коэффициентом $\rho_{YX} = k$ называется *выборочным коэффициентом регрессии*. Найдем k и b с помощью метода наименьших квадратов. Введем функцию F :

$$\begin{aligned} F(k, b) &= \sum_{i=1}^k (y_i - f(x_i, k, b))^2 \rightarrow \min \Leftrightarrow \\ &\Leftrightarrow \begin{cases} \frac{\partial F}{\partial k} = \sum_{i=1}^k 2x_i(kx_i + b - y_i) = 0 \\ \frac{\partial F}{\partial b} = \sum_{i=1}^k 2(kx_i + b - y_i) = 0 \end{cases} \Leftrightarrow \\ &\Leftrightarrow \begin{cases} k\overline{x^2} + b\overline{x} - \overline{xy} = 0 \\ k\overline{x} + b - \overline{y} = 0 \end{cases} \Rightarrow \boxed{\rho_{YX} = k = \frac{\overline{xy} - \overline{x} \cdot \overline{y}}{\overline{x^2} - (\overline{x})^2} = \frac{\overline{xy} - \overline{x} \cdot \overline{y}}{\sigma_X^2}} \end{aligned}$$

Таким образом, уравнение прямой линии регрессии Y на X выглядит

$$y_x - \overline{y} = \rho_{YX} \cdot (x - \overline{x})$$

X на Y :

$$x_y - \overline{x} = \rho_{XY} \cdot (y - \overline{y})$$

Постановка задачи

Для заданной двумерной выборки (X, Y) построить уравнения выборочных прямых среднеквадратической регрессии. Полученные линейные функ-

ции регрессии отобразить графически. Найти выборочное корреляционное отношение. Полученные результаты содержательно проинтерпретировать.

Выполнение работы

1. Построим график выборки.

Рисунок 1 – График выборки

2. Построим уравнения среднеквадратичной регрессии.

$$\bar{x} = 17111.87$$

$$\bar{y} = 50225.33$$

$$\overline{xy} = 629742536.30$$

$$\sigma_X^2 = 11878.18$$

$$\sigma_Y^2 = 39568.61$$

$$\rho_{YX} = -1.6281$$

$$\rho_{XY} = -0.1467$$

- Y на X :

$$y_x = -1.6281 \cdot (x - 17111.87) + 50225.33$$

- X на Y :

$$x_y = -0.1467 \cdot (y - 50225.33) + 17111.87$$

Напомним, что наша выборка содержит данные о стоимости машин на вторичном рынке и их пробеге. Из графика видно, что можно разделить точки на 2 категории: те, кто лежат непосредственно около прямой, и те, кто лежат от нее далеко. Можно сделать вывод, что для большинства машин стоимость линейно зависит от пробега. Для второй категории людей можно предположить, что они относятся к специфичной категории, возможно это какие-то редкие модели или что-то в этом роде.

3. Составим корреляционную таблицу.

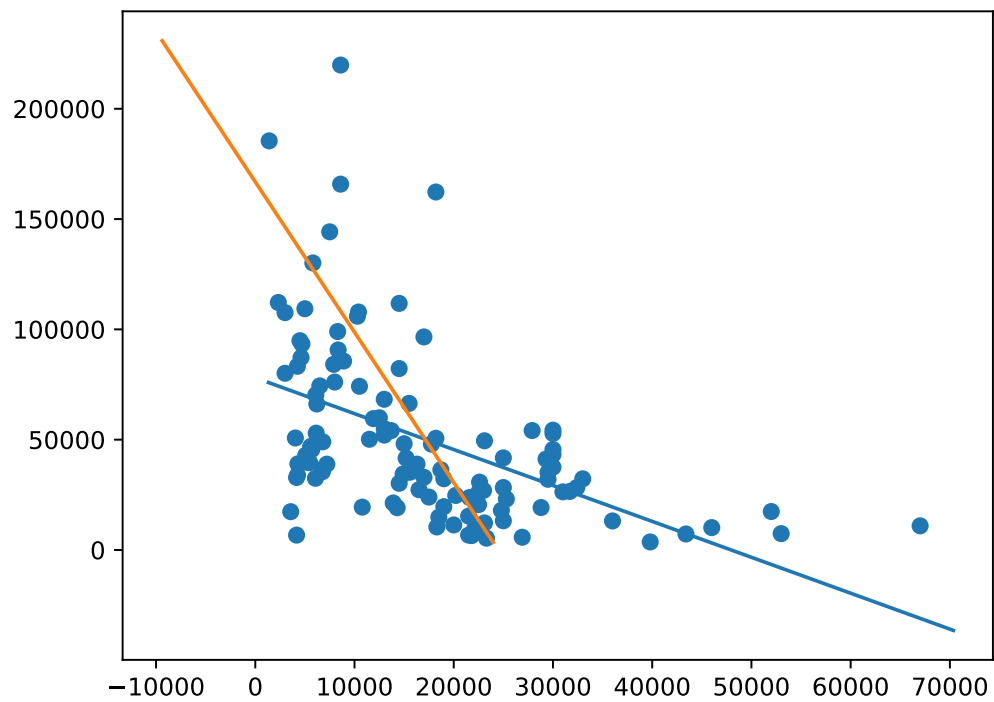


Рисунок 2 – Прямые среднеквадратичной регрессии

X	Y									m_{x_i}
	y_1	y_2	y_3	y_4	y_5	y_6	y_7	y_8	y_9	
x_1	2	3	16	7	2	2	2	0	1	35
x_2	16	9	8	10	1	0	0	0	0	44
x_3	5	7	0	2	0	0	0	0	0	14
x_4	6	6	1	0	0	0	0	0	0	13
x_5	3	1	0	0	0	0	0	0	0	4
x_6	1	0	0	0	0	0	0	0	0	1
x_7	1	1	1	0	0	0	0	0	0	3
x_8	0	0	0	0	0	0	0	0	0	0
x_9	0	1	0	0	0	0	0	0	0	1
m_{y_j}	34	28	26	19	3	2	2	0	1	115

Таблица 1 – Корреляционная таблица

Найдем межгрупповую дисперсию

$$\delta_X^2 = 6684$$

$$\delta_Y^2 = 21402$$

корреляционное отношение и выборочный коэффициент корреляции:

$$\eta_{XY} = 0.563$$

$$\eta_{YX} = 0.541$$

$$r_{XY} = -0.489$$

Как можно видеть, соотношение $\eta \geq |r_{XY}|$ выполняется.

4. Найдем уравнение параболической зависимости:

$$y = 3.014 \cdot x^2 - 3.118 \cdot x + 90492$$

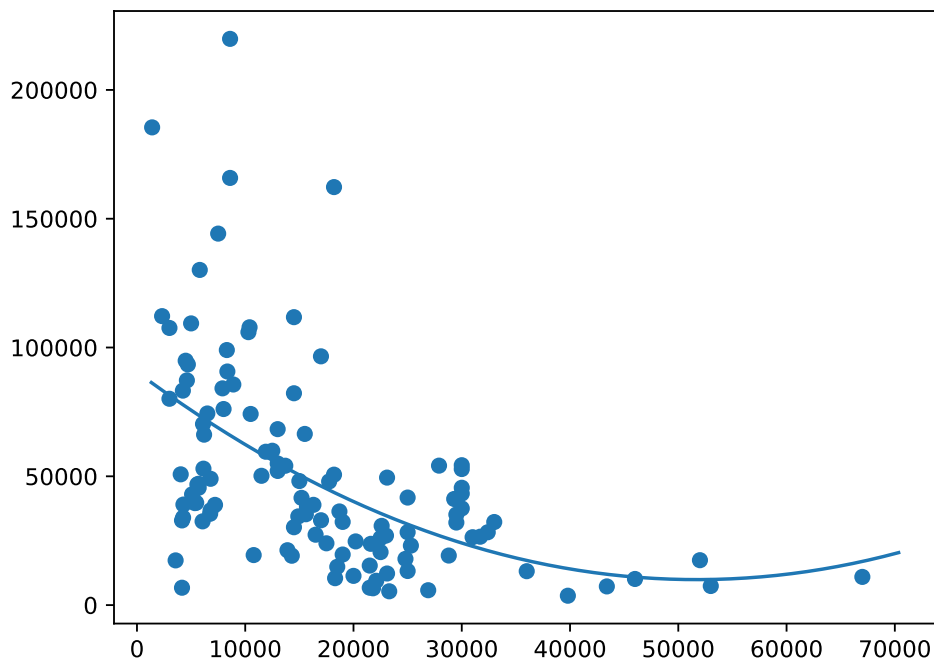


Рисунок 3 – Корреляционная кривая параболического вида

5. Построим линию корреляции вида $y = a \cdot e^{bx}$

$$y = 73866 \cdot e^{-4.04 \cdot x}$$

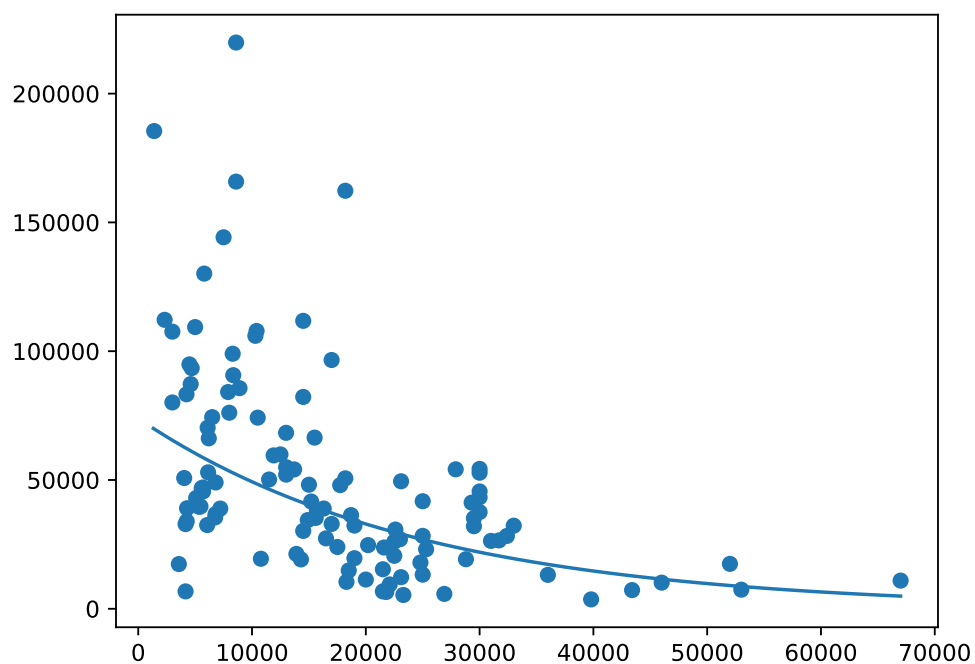


Рисунок 4 – Корреляционная кривая дробно-линейного вида

6. Вычислим показатели качества регрессии для каждой кривой:

- $y = kx + b$:

- Теоретический коэффициент детерминации:

$$R^2 = 0.2389$$

- Средняя квадратическая ошибка:

$$S_{\varepsilon} = 10454$$

- Средняя ошибка аппроксимации:

$$A = 0.829$$

- $x = ky + b$:

- Теоретический коэффициент детерминации:

$$R^2 = 0.2389$$

– Средняя квадратическая ошибка:

$$S_{\varepsilon} = 34825$$

– Средняя ошибка аппроксимации:

$$A = 0.949$$

• $y = ax^2 + bx + c$:

– Теоретический коэффициент детерминации:

$$R^2 = 0.2698$$

– Средняя квадратическая ошибка:

$$S_{\varepsilon} = 34261$$

– Средняя ошибка аппроксимации:

$$A = 0.841$$

• $y = \frac{1}{ax + b}$:

– Теоретический коэффициент детерминации:

$$R^2 = 0.215$$

– Средняя квадратическая ошибка:

$$S_{\varepsilon} = 35783$$

– Средняя ошибка аппроксимации:

$$A = 0.691$$

Выводы

Мы ознакомились с основными положениями метода наименьших квадратов, со статистическими свойствами МНК оценок, с понятием функции регрессии и роли МНК в регрессионном анализе, с корреляционным отношением, как мерой тесноты произвольной корреляционной связи. Научились строить кривые корреляции для произвольных функций, а также корреляционную таблицу.

ПРИЛОЖЕНИЕ А

```
1 import csv
2 import math
3 from math import sqrt, floor, log2
4 from itertools import islice
5 from collections import Counter
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 def read_dataset():
10     with open('Price-Mileage.csv', 'r') as csv_file:
11         reader = csv.reader(csv_file, delimiter=';')
12         return [(int(row[0]), int(row[1])) for row in
13                 islice(reader, 1, None)]
14
15 # 1d sorted dataset
16 def make_interval_dataset(dataset):
17     k = 1 + floor(log2(len(dataset)))
18     h = (dataset[-1] - dataset[0]) / k
19
20     def x(i):
21         x0 = dataset[0] - h / 2
22         return x0 + h * i
23
24     j = 0
25     counts = []
26     xs = [(x(i), x(i + 1)) for i in range(k + 1)]
27     for i in range(k + 1):
28         count_i = 0
29         while j < len(dataset) and dataset[j] <= x(i + 1):
30             count_i += 1
31             j += 1
```

```

30     counts.append(count_i)
31     frequencies = [c / len(dataset) for c in counts]
32     return list(zip(xs, counts, frequencies))
33
34 def paragraph(i):
35     print()
36     print('-' * 15 + f' {i} ' + '-' * 15)
37
38 dataset = read_dataset()
39 #dataset = list(filter(lambda p: p[1] <= 18000,
40     read_dataset()))
41 xs = list(map(lambda p: p[0], dataset))
42 ys = list(map(lambda p: p[1], dataset))
43
44 xs_int = make_interval_dataset(sorted(xs))
45 ys_int = make_interval_dataset(sorted(ys))
46
47 # 1
48 paragraph(1)
49
50 fig, ax = plt.subplots()
51 ax.scatter(xs, ys)
52 plt.savefig('scatter.pdf')
53
54 # 2
55 paragraph(2)
56
57 print('y = kx + b; x = ky + b')
58 x_avg = sum(xs) / len(xs)
59 y_avg = sum(ys) / len(ys)
60 xys = [xs[i] * ys[i] for i in range(len(xs))]
61 xy_avg = sum(xys) / len(xys)

```

```

62 print(f'averages:\n\t x = {x_avg} y = {y_avg} xy = {
    xy_avg}')
63
64
65 def calc_sdev(a, avg):
66     return sqrt(sum(map(lambda x: (x - avg) ** 2, a))
    / len(a))
67
68 x_sdev = calc_sdev(xs, x_avg)
69 y_sdev = calc_sdev(ys, y_avg)
70
71 print(f'sdevs:\n\t x_sdev = {x_sdev} y_sdev = {y_sdev}
    ')
72
73 rho_yx = (xy_avg - x_avg * y_avg) / x_sdev ** 2
74 rho_xy = (xy_avg - x_avg * y_avg) / y_sdev ** 2
75
76 print(f'pho:\n\t rho_yx = {rho_yx} rho_xy = {rho_xy}')
77
78 def y(x):
79     return rho_yx * (x - x_avg) + y_avg
80
81 def x(y):
82     return rho_xy * (y - y_avg) + x_avg
83
84 max_x = max(xs)
85 max_y = max(ys)
86 min_x = min(xs)
87 min_y = min(ys)
88
89 def make_y_plot():
90     xs = np.linspace(0.95 * min_x, max_x * 1.05, 20)
91     ys = np.linspace(0.95 * min_y, max_y * 1.05, 20)
92     ax.plot(xs, y(xs))

```

```

93     ax.plot(x(ys), ys)
94     plt.savefig('regression_lines.pdf')
95
96 make_y_plot()
97
98
99 # parameters
100 paragraph('parameters')
101
102 R2 = rho_xy * rho_yx
103 print(f'R^2 = {R2}')
104
105 R2_y = sum(map(lambda i: (y(xs[i]) - y_avg) ** 2,
106                range(len(xs)))) / (y_sdev ** 2
107                * len(ys))
108
109 R2_x = sum(map(lambda i: (x(ys[i]) - x_avg) ** 2,
110                range(len(xs)))) / (x_sdev ** 2
111                * len(ys))
112
113 print(f'R^2_y = {R2_y} R^2_x = {R2_x}')
114
115 eps_y = list(map(lambda i: ys[i] - y(xs[i]), range(len
116                (xs))))
117
118 eps_x = list(map(lambda i: xs[i] - x(ys[i]), range(len
119                (xs))))
120
121 S_eps_y = sqrt(sum(map(lambda i: i ** 2, eps_y)) / (
122                len(xs) - 2))
123
124 S_eps_x = sqrt(sum(map(lambda i: i ** 2, eps_x)) / (
125                len(xs) - 2))
126
127 print(f'S_eps_y = {S_eps_y}, S_eps_x = {S_eps_x}')
128
129 Ay = sum(map(lambda i: abs(eps_y[i] / ys[i]), range(
130                len(xs)))) / len(xs)
131
132 Ax = sum(map(lambda i: abs(eps_x[i] / xs[i]), range(
133                len(xs)))) / len(xs)
134
135 print(f'Ax = {Ax} Ay = {Ay}')
136
137

```

```

119 # ----- 3 -----
120 paragraph(3)
121
122 def get_x_interval(x):
123     for i in range(len(xs_int)):
124         if xs_int[i][0][0] <= x <= xs_int[i][0][1]:
125             return i
126         raise f'incorrect x = {x}'
127
128 def get_y_interval(y):
129     for i in range(len(ys_int)):
130         if ys_int[i][0][0] <= y <= ys_int[i][0][1]:
131             return i
132         raise f'incorrect y = {y}'
133
134 count = Counter([(get_x_interval(x), get_y_interval(y)
135                   ) for x, y in zip(xs, ys)])
136
137 corr_table = [[count[(i, j)] for j in range(len(ys_int))]]
138               for i in range(len(xs_int))]
139
140 m_y = [sum(map(lambda m: m[j], corr_table)) for j in
141         range(len(ys_int))]
142
143 assert sum(m_y) == len(dataset)
144 print(f'm_y = {m_y}')
145
146
147 m_x = [sum(line) for line in corr_table]
148
149 assert sum(m_x) == len(dataset)
150 print(f'm_x = {m_x}')
151
152 def get_x(i):
153     return (xs_int[i][0][1] + xs_int[i][0][0]) / 2

```

```

150 def get_y(j):
151     return (ys_int[j][0][1] + ys_int[j][0][0]) / 2
152
153 def x_y_avg(j):
154     s1 = 0
155     for i in range(len(xs_int)):
156         s1 += get_x(i) * corr_table[i][j]
157     s2 = 0
158     for i in range(len(xs_int)):
159         s2 += corr_table[i][j]
160     return s1 / s2
161
162 def y_x_avg(i):
163     s1 = 0
164     for j in range(len(ys_int)):
165         s1 += get_y(j) * corr_table[i][j]
166     s2 = 0
167     for j in range(len(ys_int)):
168         s2 += corr_table[i][j]
169     return s1 / s2
170
171
172 new_x_avg = 0
173 for i in range(len(xs_int)):
174     new_x_avg += get_x(i) * m_x[i]
175 new_x_avg /= len(dataset)
176
177 print(new_x_avg, x_avg)
178
179 delta_x = 0
180 for j in range(len(ys_int)):
181     delta_x += (x_y_avg(j) - x_avg) ** 2 * m_y[j]
182 delta_x = sqrt(delta_x / len(dataset))
183

```

```

184 delta_y = 0
185 for i in range(len(xs_int)):
186     delta_y += (y_x_avg(i) - y_avg) ** 2 * m_x[i]
187 delta_y = sqrt(delta_y / len(dataset))
188
189
190 print(f'delta_x = {delta_x} delta_y = {delta_y}')
191 nu_xy = delta_x / x_sdev
192 nu_yx = delta_y / y_sdev
193
194 print(f'nu_xy = {nu_xy} nu_yx = {nu_yx}')
195
196
197 r = (xy_avg - x_avg * y_avg) / (y_sdev * x_sdev)
198 print(f'r = {r}')
199
200 # ----- 4 -----
201 paragraph(4)
202
203 def get_degree_sum(i):
204     s = 0
205     for x in xs:
206         s += x ** i
207     return float(s)
208
209 A = np.array([
210     [get_degree_sum(4), get_degree_sum(3),
211      get_degree_sum(2)],
212     [get_degree_sum(3), get_degree_sum(2),
213      get_degree_sum(1)],
214     [get_degree_sum(2), get_degree_sum(1),
215      get_degree_sum(0)]
216 ])

```



```

215 x2y_sum = 0
216 for i in range(len(xs)):
217     x2y_sum += xs[i] ** 2 * ys[i]
218
219 xy_sum = sum(xys)
220 y_sum = sum(ys)
221
222 B = np.array([
223     float(x2y_sum),
224     float(xy_sum),
225     float(y_sum)
226 ])
227
228 a, b, c = np.linalg.solve(A, B)
229 print('y = ax^2 + bx + c\n\t', f'a = {a}, b = {b}, c =
      {c}', sep='')
230
231 def y(x):
232     return a * x ** 2 + b * x + c
233
234 def make_parab_plot():
235     x = np.linspace(0.95 * min_x, max_x * 1.05, 50)
236     fig, ax = plt.subplots()
237     ax.scatter(xs, ys)
238     ax.plot(x, y(x))
239     plt.savefig('parab.pdf')
240
241 make_parab_plot()
242
243 #parameters
244 paragraph('parameters')
245
246 R2 = sum(map(lambda i: (y(xs[i]) - y_avg) ** 2, range(
      len(xs)))) / (y_sdev ** 2

```

```

247         * len(ys))
248 print(f'R^2 = {R2}')
249 eps_y = list(map(lambda i: ys[i] - y(xs[i]), range(len
    (xs))))
250 S_eps_y = sqrt(sum(map(lambda i: i ** 2, eps_y)) / (
    len(xs) - 3))
251 print(f'S_eps_y = {S_eps_y}')
252 Ay = sum(map(lambda i: abs(eps_y[i] / ys[i]), range(
    len(xs)))) / len(xs)
253 print(f'Ay = {Ay}')
254
255 # ----- 5 -----
256 # var 3
257 # y = 1 / (ax + b)
258 paragraph(5)
259
260 xs_stroke = xs
261 ys_stroke = list(map(lambda y: math.log(y), ys))
262
263 x_stroke_avg = sum(xs_stroke) / len(xs_stroke)
264 y_stroke_avg = sum(ys_stroke) / len(ys_stroke)
265 xys_stroke = [xs_stroke[i] * ys_stroke[i] for i in
    range(len(xs_stroke))]
266 xy_stroke_avg = sum(xys_stroke) / len(xys_stroke)
267
268 print(f'averages:\n\t x = {x_avg} y = {y_avg} xy = {
    xy_avg}')
269
270
271 def calc_sdev(a, avg):
272     return sqrt(sum(map(lambda x: (x - avg) ** 2, a))
    / len(a))
273
274 x_stroke_sdev = calc_sdev(xs_stroke, x_stroke_avg)

```

```

275 y_stroke_sdev = calc_sdev(ys_stroke, y_stroke_avg)
276
277 print(f'sdevs:\n\t x_sdev = {x_sdev} y_sdev = {y_sdev}
      ')
278
279 rho_yx = (xy_stroke_avg - x_stroke_avg * y_stroke_avg)
      / x_stroke_sdev ** 2
280 b = rho_yx
281 a = math.exp(y_stroke_avg - rho_yx * x_stroke_avg)
282 print(f'y = a*exp(bx):\n\t a = {a}, b = {b}')
283
284 def y(x):
285     return a*np.exp(x * b)
286     #return 1 / ( rho_yx * (x - x_stroke_avg) +
      y_stroke_avg)
287
288
289 max_x = max(xs)
290 max_y = max(ys)
291 min_x = min(xs)
292 min_y = min(ys)
293
294 def make_y_plot():
295     fig, ax = plt.subplots()
296     ax.scatter(xs, ys)
297     # ax.set_yscale('log')
298     x = np.linspace(0.95 * min_x, max_x, 50)
299     ax.plot(x, y(x))
300     plt.savefig('regression_lines2.pdf')
301
302 make_y_plot()
303
304 #params
305 paragraph('parameters')

```

```

306
307 R2 = sum(map(lambda i: (y(xs[i]) - y_avg) ** 2, range(
    len(xs)))) / (y_sdev ** 2
308               * len(ys))
309 print(f'R^2 = {R2}')
310 eps_y = list(map(lambda i: ys[i] - y(xs[i]), range(len
    (xs))))
311 S_eps_y = sqrt(sum(map(lambda i: i ** 2, eps_y)) / (
    len(xs) - 2))
312 print(f'S_eps_y = {S_eps_y}')
313 Ay = sum(map(lambda i: abs(eps_y[i] / ys[i]), range(
    len(xs)))) / len(xs)
314 print(f'Ay = {Ay}')
315
316 # ----- 6 -----
317
318
319 plt.show()

```

Листинг 1 – Исходный код программы