

Санкт-Петербургский Политехнический университет
Петра Великого
Институт Прикладной Математики и Механики
Кафедра «Прикладная Математика и Информатика»

Отчет
По лабораторной работе № 6
По Дисциплине «Математическая статистика»

Выполнил:
Студент Селянкин Федор
Группа 3630102/70301
Проверил:
к.ф. – м.н., доцент
Баженов Александр Николаевич

Санкт-Петербург
2020 г.

Содержание

| | |
|---|---|
| Постановка задачи | 3 |
| Теория..... | 3 |
| Модель простой линейной регрессии..... | 3 |
| Метод наименьших квадратов..... | 3 |
| Расчётные формулы для МНК-оценок..... | 4 |
| Робастные оценки коэффициентов линейной регрессии | 5 |
| Реализация | 6 |
| Результаты | 6 |
| Выборка без возмущений..... | 6 |
| Выборка с возмущениями | 7 |
| Литература | 7 |
| Обсуждения | 8 |
| Список иллюстраций: | |
| Рисунок 1 Выборка без возмущений | 7 |
| Рисунок 2 Выборка с возмущениями..... | 7 |

Постановка задачи

Найти оценки коэффициентов линейной регрессии $y_i = a + bx_i + e_i$, используя 20 точек на отрезке $[-1.8; 2]$ с равномерным шагом равным 0.2. Ошибку e_i считать нормально распределённой с параметрами $(0, 1)$. В качестве эталонной зависимости взять $y_i = 2 + 2x_i + e_i$. При построении оценок коэффициентов использовать два критерия: критерий наименьших квадратов и критерий наименьших модулей. Прodelать то же самое для выборки, у которой в значения y_1 и y_{20} вносятся возмущения 10 и -10.

Теория

Модель простой линейной регрессии

Регрессионную модель описания данных называют простой линейной регрессией, если

(1)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i = 1, \dots, n,$$

где x_1, \dots, x_n – заданные числа (значения фактора); y_1, \dots, y_n – наблюдаемые значения отклика; $\varepsilon_1, \dots, \varepsilon_n$ – независимые, нормально распределённые $N(0, \sigma)$ с нулевым математическим ожиданием и одинаковой (неизвестной) дисперсией случайные величины (ненаблюдаемые); β_0, β_1 – неизвестные параметры, подлежащие оцениванию.

В модели (1) отклик y зависит от одного фактора x , и весь разброс экспериментальных точек объясняется только погрешностями наблюдений (результатов измерений) отклика y . Погрешности результатов измерений x в этой модели полагают существенно меньшими погрешностей результатов измерений y , так что ими можно пренебречь.

Метод наименьших квадратов

При оценивании параметров регрессионной модели используют различные методы. Один из наиболее распространённых подходов заключается в следующем: вводится мера (критерий) рассогласования отклика и регрессионной функции, и оценки параметров регрессии определяются так, чтобы сделать это рассогласование наименьшим. Достаточно простые расчётные формулы для оценок получают при выборе критерия в виде суммы квадратов отклонений значений отклика от значений регрессионной функции (сумма квадратов остатков):

(2)

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \rightarrow \min_{\beta_0, \beta_1}$$

Задача минимизации квадратичного критерия (2) носит название задачи метода наименьших квадратов (МНК), а оценки $\widehat{\beta}_0, \widehat{\beta}_1$ параметров β_0, β_1 , реализующие минимум критерия (2), называют МНК-оценками.

Расчётные формулы для МНК-оценок

МНК-оценки параметров $\widehat{\beta}_0$ и $\widehat{\beta}_1$ находятся из условия обращения функции $Q(\beta_0, \beta_1)$ в минимум.

Для нахождения МНК-оценок $\widehat{\beta}_0$ и $\widehat{\beta}_1$ выпишем необходимые условия экстремума:

(3)

$$\begin{cases} \frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0 \\ \frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0 \end{cases}$$

Далее для упрощения записи сумм будем опускать индекс суммирования. Из системы (3) получим

$$\begin{cases} n\widehat{\beta}_0 + \widehat{\beta}_1 \sum x_i = \sum y_i \\ \widehat{\beta}_0 \sum x_i + \widehat{\beta}_1 \sum x_i^2 = \sum x_i y_i \end{cases}$$

Разделим оба уравнения на n:

(4)

$$\begin{cases} \widehat{\beta}_0 + \left(\frac{1}{n} \sum x_i\right) \widehat{\beta}_1 = \frac{1}{n} \sum y_i \\ \left(\frac{1}{n} \sum x_i\right) \widehat{\beta}_0 + \left(\frac{1}{n} \sum x_i^2\right) \widehat{\beta}_1 = \frac{1}{n} \sum x_i y_i \end{cases}$$

и, используя известные статистические обозначения для выборочных первых и вторых начальных моментов

$$\bar{x} = \frac{1}{n} \sum x_i, \bar{y} = \frac{1}{n} \sum y_i, \overline{x^2} = \frac{1}{n} \sum x_i^2, \overline{xy} = \frac{1}{n} \sum x_i y_i,$$

получим

(5)

$$\begin{cases} \widehat{\beta}_0 + \bar{x} \widehat{\beta}_1 = \bar{y} \\ \bar{x} \widehat{\beta}_0 + \overline{x^2} \widehat{\beta}_1 = \overline{xy} \end{cases}$$

Откуда МНК-оценку $\widehat{\beta}_1$ наклона прямой регрессии находим по формуле Крамера

(6)

$$\widehat{\beta}_1 = \frac{\overline{xy} - \bar{x} * \bar{y}}{\overline{x^2} - (\bar{x})^2},$$

а МНК-оценку $\widehat{\beta}_0$ определяем непосредственно из первого уравнения системы (5):

(7)

$$\widehat{\beta}_0 = \bar{y} - \bar{x} \widehat{\beta}_1$$

Заметим, что определитель системы (5)

$$\overline{x^2} - (\bar{x})^2 = n^{-1} \sum (x_i - \bar{x})^2 = s_x^2 > 0,$$

Если среди значений x_1, \dots, x_n есть различные, что и будем предполагать.

Доказательство минимальности функции $Q(\beta_0, \beta_1)$ в стационарной точке проведём с помощью известного достаточного признака экстремума функции двух переменных. Имеем:

$$\frac{\partial^2 Q}{\partial \beta_0^2} = 2n, \frac{\partial^2 Q}{\partial \beta_1^2} = 2 \sum x_i^2 = 2n\overline{x^2}, \frac{\partial^2 Q}{\partial \beta_0 \partial \beta_1} = 2 \sum x_i = 2n\bar{x}.$$

$$\Delta = \frac{\partial^2 Q}{\partial \beta_0^2} * \frac{\partial^2 Q}{\partial \beta_1^2} - \left(\frac{\partial^2 Q}{\partial \beta_0 \partial \beta_1} \right)^2 = 4n^2 \overline{x^2} - 4n^2 (\bar{x})^2 = 4n^2 [\overline{x^2} - (\bar{x})^2] = 4n^2 \left[\frac{1}{n} \sum (x_i - \bar{x})^2 \right] = 4n^2 s_x^2 > 0$$

Этот результат вместе с условием $\frac{\partial^2 Q}{\partial \beta_0^2} = 2n > 0$ означает, что в стационарной точке функция Q имеет минимум.

Робастные оценки коэффициентов линейной регрессии

Робастность оценок коэффициентов линейной регрессии (т.е. их устойчивость по отношению к наличию в данных редких, но больших по величине выбросов) может быть обеспечена различными способами. Одним из них является использование метода наименьших модулей вместо метода наименьших квадратов:

(8)

$$\sum_{i=1}^n |y_i - \beta_0 - \beta_1 x_i| \rightarrow \min_{\beta_0, \beta_1}$$

Напомним, что использование метода наименьших модулей в задаче оценивания параметра сдвига распределений приводит к оценке в виде выборочной медианы, обладающей робастными свойствами. В отличие от этого случая и от задач метода наименьших квадратов, на практике задача (8) решается численно. Соответствующие процедуры представлены в некоторых современных пакетах программ по статистическому анализу.

Здесь мы рассмотрим простейшую в вычислительном отношении робастную альтернативу оценкам коэффициентов линейной регрессии по МНК. Для этого сначала запишем выражения для оценок (6) и (7) в другом виде:

(9)

$$\widehat{\beta}_1 = \frac{\overline{xy} - \bar{x} * \bar{y}}{\overline{x^2} - (\bar{x})^2} = \frac{k_{xy}}{s_x^2} = \frac{k_{xy}}{s_x s_y} * \frac{s_y}{s_x} = r_{xy} \frac{s_y}{s_x}$$

$$\widehat{\beta}_0 = \bar{y} - \bar{x} \widehat{\beta}_1$$

В формулах (17) заменим выборочные средние \bar{x} и \bar{y} соответственно на робастные выборочные медианы $\text{med } x$ и $\text{med } y$, среднеквадратические отклонения s_x и s_y на робастные нормированные интерквартильные широты q_x^* и q_y^* , выборочный коэффициент корреляции r_{xy} - на знаковый коэффициент корреляции r_Q :

(10)

$$\widehat{\beta}_{1R} = r_Q \frac{q_y^*}{q_x^*}$$

(11)

$$\widehat{\beta}_{0R} = \text{med } y - \widehat{\beta}_{1R} \text{ med } x,$$

(12)

$$r_Q = \frac{1}{n} \sum_{i=1}^n \text{sign}(x_i - \text{med } x) \text{sign}(y_i - \text{med } y)$$

(13)

$$q_y^* = \frac{y_{(j)} - y_{(l)}}{k_q(n)}, \quad q_x^* = \frac{x_{(j)} - x_{(l)}}{k_q(n)}$$

$$l = \begin{cases} \lceil n/4 \rceil + 1 & \text{при } n/4 \text{ дробном} \\ n/4 & \text{при } n/4 \text{ целом} \end{cases}$$

$$j = n - l + 1$$

$$\text{sign } z = \begin{cases} 1 & \text{при } z > 0 \\ 0 & \text{при } z = 0 \\ -1 & \text{при } z < 0 \end{cases}$$

Уравнение регрессии здесь имеет вид

(14)

$$y = \widehat{\beta}_{0R} + \widehat{\beta}_{1R} x$$

Статистики выборочной медианы и интерквартильной широты обладают робастными свойствами в силу того, что основаны на центральных порядковых статистиках, малочувствительных к большим по величине выбросам в данных. Статистика выборочного знакового коэффициента корреляции робастна, так как знаковая функция $\text{sign } z$ чувствительна не к величине аргумента, а только к его знаку. Отсюда оценка прямой регрессии (14) обладает очевидными робастными свойствами устойчивости к выбросам по координате y , но она довольно груба.

Реализация

Лабораторная работа выполнена с помощью встроенных средств языка программирования Python в среде разработки PyCharm, с использованием дополнительных библиотек для отображения и расчетов. Исходный код лабораторной выложен на веб-сервисе GitHub [1].

Результаты

Выборка без возмущений

- Критерий наименьших квадратов

$$\hat{a} \approx 2.17, \quad \hat{b} \approx 2.22$$

- Критерий наименьших модулей

$$\hat{a} \approx 2.50, \quad \hat{b} \approx 1.67$$

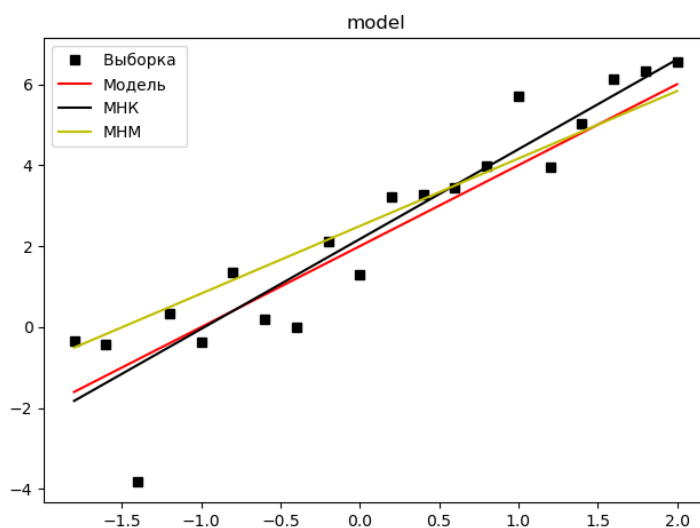


Рисунок 1 Выборка без возмущений

Выборка с возмущениями

- Критерий наименьших квадратов

$$\hat{a} \approx 2.49, \quad \hat{b} \approx 0.63$$

- Критерий наименьших модулей

$$\hat{a} \approx 2.38, \quad \hat{b} \approx 1.88$$

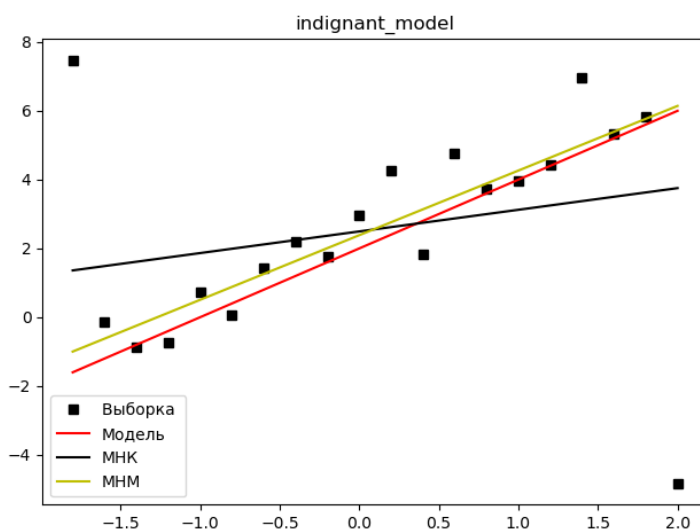


Рисунок 2 Выборка с возмущениями

Литература

1. Ссылка на репозиторий GitHub
<https://github.com/SelyankinFyodor/math-statistics>
2. Вероятностные разделы математики. Учебник для бакалавров технических направлений.
3. //Под ред. Максимова Ю.Д. – Спб.: Иван Федоров, 2001. – 592 с.Максимов Ю.Д.
Математика. Теория и практика по математической статистике. Конспект-справочник по теории вероятностей : учеб. пособие Ю.Д. Максимов; под ред. В.И. Антонова. — Спб. : Изд-во Политехн. ун-та, 2009. — 395 с. (Математика в политехническом университете).

Обсуждения

Критерий наименьших квадратов точнее оценивает коэффициенты линейной регрессии на выборке без возмущений.

Критерий наименьших модулей точнее оценивает коэффициенты линейной регрессии на выборке с возмущениями.

Критерий наименьших модулей устойчив к редким выбросам.(страница. 5)