

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ
ПЕТРА ВЕЛИКОГО

ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ И МЕХАНИКИ
КАФЕДРА "ПРИКЛАДНАЯ МАТЕМАТИКА"

ОТЧЁТ
ЛАБОРАТОРНАЯ РАБОТА №6
ПО ДИСЦИПЛИНЕ
"МАТЕМАТИЧЕСКАЯ СТАТИСТИКА"

ВЫПОЛНИЛ СТУДЕНТ:
САЛИХОВ С.Р.
ГРУППА: 3630102/70401

ПРОВЕРИЛ:
К.Ф-М.Н., ДОЦЕНТ
БАЖЕНОВ АЛЕКСАНДР НИКОЛАВИЧ

САНКТ-ПЕТЕРБУРГ
2020 г.

Содержание

| | Стр. |
|---|----------|
| 1. Постановка задачи | 4 |
| 2. Теория | 4 |
| 2.1. Простая линейная регрессия | 4 |
| 2.1.1 Модель простой линейной регрессии | 4 |
| 2.2. Метод наименьших квадратов | 4 |
| 2.2.1 Расчётные формулы для МНК-оценок | 4 |
| 2.3. Метод наименьших модулей | 5 |
| 3. Реализация | 5 |
| 4. Результаты | 6 |
| 4.1. Выборка без возмущений | 6 |
| 4.2. Выборка с возмущениями | 6 |
| 5. Обсуждение | 7 |
| 6. Литература | 7 |
| 7. Приложения | 7 |

Список иллюстраций

| | | |
|---|--|---|
| 1 | Графики линейной регрессии при выборке с возмущением и без | 6 |
|---|--|---|

1 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.

2 Теория

2.1 Простая линейная регрессия

2.1.1 Модель простой линейной регрессии

Регрессионную модель описания данных называют простой линейной регрессией, если

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

, $i = 1, \dots, n$, где x_1, \dots, x_n — заданные числа (значения фактора); y_1, \dots, y_n — наблюдаемые значения отклика; $\epsilon_1, \dots, \epsilon_n$ — независимые, нормально распределённые $N(0, \delta)$ с нулевым математическим ожиданием и одинаковой (неизвестной) дисперсией случайные величины (ненаблюдаемые); β_0, β_1 — неизвестные параметры, подлежащие оцениванию.

В модели 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.

2.2 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):

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

Задача минимизации квадратичного критерия (10) носит название задачи метода наименьших квадратов (МНК), а оценки $\hat{\beta}_0, \hat{\beta}_1$ параметров β_0, β_1 , реализующие минимум критерия, называют МНК-оценками.

2.2.1 Расчётные формулы для МНК-оценок

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

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

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

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

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

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

Тогда:

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

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

$$\hat{\beta}_1 = \frac{\overline{xy} - \bar{x} * \bar{y}}{\overline{x^2} - \bar{x}^2}$$

а МНК-оценку $\hat{\beta}_0$ определяем непосредственно из первого уравнения системы:

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

2.3 Метод наименьших модулей

Критерий наименьших модулей – заключается в минимизации следующей функции [?]:

$$M(a, b) = \sum_{i=1}^n |y_i - ax_i - b| \rightarrow \min \quad (1)$$

3 Реализация

Для генерации выборки был использован *Python 3.7* и модуль *numpy*. Для отрисовки графиков использовался модуль *matplotlib*. *scipy.stats* для обработки функций распределений.

4 Результаты

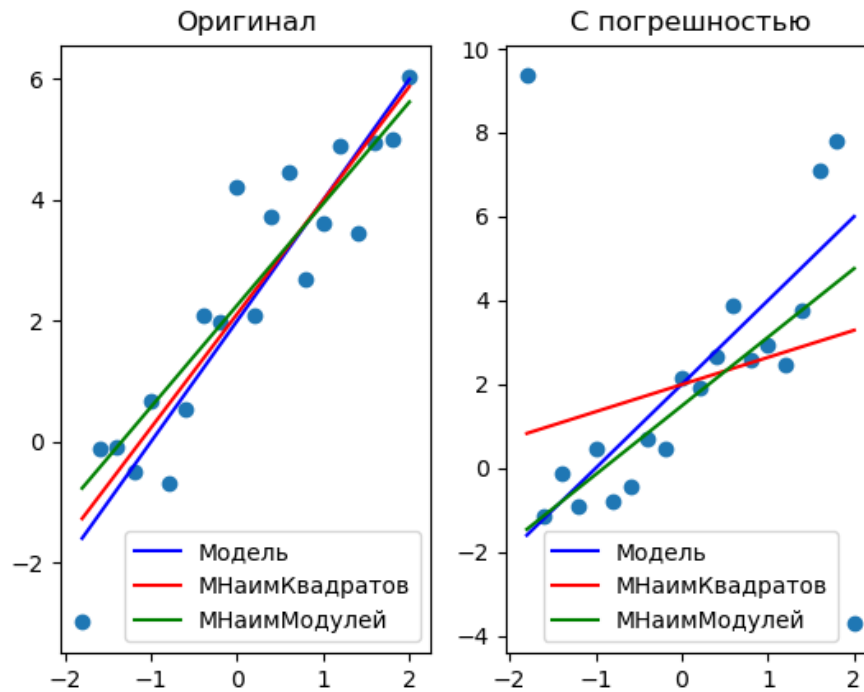


Рис. 1: Графики линейной регрессии при выборке с возмущением и без

4.1 Выборка без возмущений

Критерий наименьших квадратов:

$$\hat{a} \approx 1.93, \hat{b} \approx 2.19$$

Критерий наименьших модулей:

$$\hat{a} \approx 2.24, \hat{b} \approx 1.77$$

4.2 Выборка с возмущениями

Критерий наименьших квадратов:

$$\hat{a} \approx 0.48, \hat{b} \approx 1.76$$

Критерий наименьших модулей:

$$\hat{a} \approx 1.85, \hat{b} \approx 1.39$$

5 Обсуждение

1)МНК оценивает коэффициенты линейной регрессии точнее, на выборке без возмущений.

Для доказательства этого введём метрику суммы квадратов разностей значений по оси y между МНК и модели и МНМ и модели. $\rho_1 = \sum (y_{MNK} - y_{etl})^2$, $\rho_2 = \sum (y_{MNM} - y_{etl})^2$ и увидим, что ρ_1 всегда меньше ρ_2 .

Пример:

$y_{etl} = [-1.60, -1.20, -0.80, -0.40, 0.00, 0.40, 0.80, 1.20, 1.60, 2.00, 2.40, 2.80, 3.20, 3.60, 4.00, 4.40, 4.80, 5.20, 5.60, 6.0]$
 $y_{MNK} \approx [-1.27, -0.89, -0.52, -0.14, 0.23, 0.61, 0.98, 1.36, 1.74, 2.11, 2.491, 2.86, 3.24, 3.62, 3.99, 4.37, 4.75, 5.12, 5.50, 5.85]$
 $y_{MNM} \approx [-0.77, -0.43, -0.09, 0.23, 0.57, 0.911, 1.24, 1.58, 1.92, 2.25, 2.59, 2.93, 3.26, 3.60, 3.94, 4.27, 4.61, 4.95, 5.28, 5.62]$
МНК : $\hat{a} = 1.88, \hat{b} = 2.11$
МНМ : $\hat{a} = 2.26, \hat{b} = 1.68$
 $\rho_1 = 16.59$
 $\rho_2 = 17.96$
Таким образом, $\rho_1 < \rho_2$.

2)На выборке с возмущениями эффективнее использовать МНМ. Таким образом, метод наименьших модулей устойчив к редким выбросам, в свою очередь МНМ обладает большей сложностью вычислений, чем МНК(т.к. ,в коде, МНК - оценки вычисляются из расчётных формул, а МНМ - оценки вычисляются, через решение задачи минимизации).

6 Литература

Модуль numpy

Модуль matplotlib

Модуль scipy

Метод наименьших модулей

Метод наименьших квадратов

7 Приложения

Код лаборатрной

Код отчёта