

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ ИМЕНИ
ПЕТРА ВЕЛИКОГО

ВЫСШАЯ ШКОЛА ПРИКЛАДНОЙ МАТЕМАТИКИ И ФИЗИКИ

Математическая статистика
Отчёт по лабораторной работе №9

Выполнил:

Студент: Теплов Андрей
Сергеевич Группа: 5030102/10201

Принял:

к. ф.-м. н., доцент
Баженов Александр Николаевич

Санкт-Петербург
2024г..

Содержание

1. Постановка задачи	2
2. Теория	3
2.1. Представление данных	3
2.2. Линейная регрессия	3
2.2.1. Описание модели	3
2.2.2. Метод наименьших модулей	3
2.3. Предварительная обработка данных	4
2.4. Коэффициент Жаккара	4
2.5. Процедура оптимизации	5
3. Реализация	5
4. Результаты	5
5. Обсуждение	8
6. Ссылка на репозиторий	8
Список литературы	8

Список иллюстраций

1. Схема установки для исследования фотоэлектрических характеристик. .	2
2. Исходные данные из экспериментов	5
3. Интервальное представление исходных данных	6
4. Линейная модель дрейфа данных	6
5. Гистограммы значений множителей коррекции w	6
6. Скорректированные модели данных	7
7. Гистограммы скорректированных данных	7
8. Значение коэффициента Жаккара от калибровочного множителя от R_{21} .	7
9. Гистограмма объединённых данных при оптимальном значении R_{21} . .	8

1. Постановка задачи

Исследование из области солнечной энергетики [1]. На рис 1 показана схема установки для исследования фотоэлектрических характеристик.

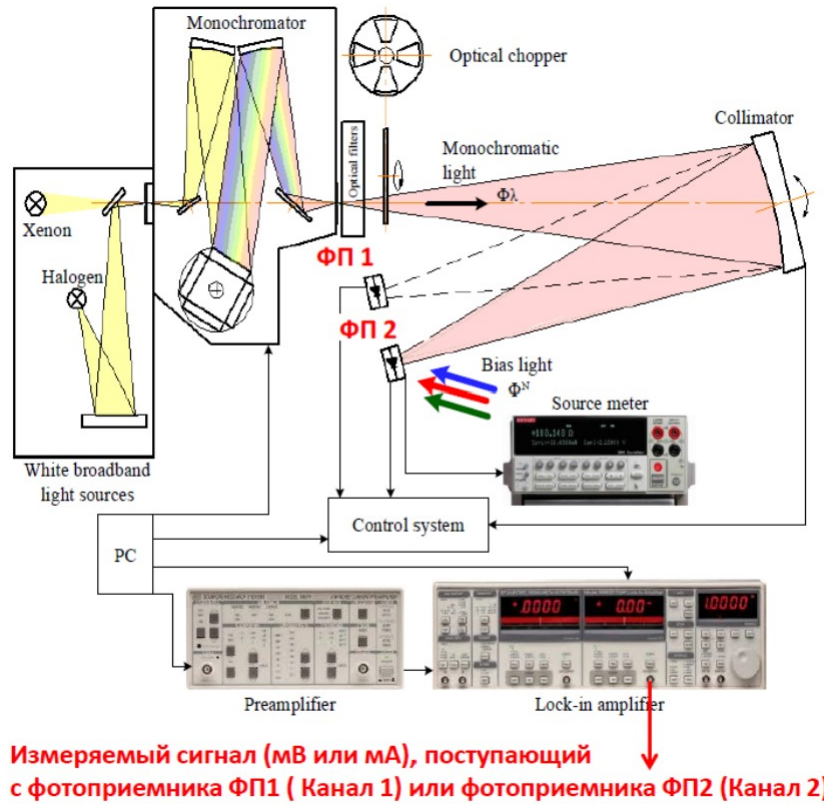


Рис. 1. Схема установки для исследования фотоэлектрических характеристик.

Калибровка датчика ФП1 производится по эталону ФП2. Зависимость между квантовыми эффективностями датчиков предполагается одинаковой для каждой пары измерений

$$QE_2 = \frac{I_2}{I_1} * QE_1 \quad (1)$$

QE_2, QE_1 - эталонная эффективность эталонного и исследуемого датчика, а I_2, I_1 - измеренные точки. Данные с датчиков находятся в файлах Ch2_800nm_2m.csv и Ch1_800nm_2m.csv.

Требуется определить коэффициент калибровки

$$R_{21} = \frac{I_2}{I_1} \quad (2)$$

при помощи линейной регрессии на множестве интервальных данных и коэффициента Жаккара.

2. Теория

2.1. Представление данных

В первую очередь представим данные таким образом, чтобы применить понятия статистики данных с интервальной неопределенностью.

Один из распространённых способов получения интервальных результатов в первичных измерениях - это "обинтерваливание" точечных значений, когда к точечному базовому значению \dot{x} , которое считывается по показаниям измерительного прибора, прибавляется *интервал погрешности* ϵ :

$$\mathbf{x} = \dot{x} + \epsilon \quad (3)$$

Интервал погрешности зададим как

$$\epsilon = [-\epsilon; \epsilon]$$

В конкретных измерениях примем $\epsilon = 10^{-4}$ мВ.

Согласно терминологии интервального анализа, рассматриваемая выборка - это вектор интервалов, или интервальный вектор $x = (x_1, x_2, \dots, x_n)$.

2.2. Линейная регрессия

2.2.1. Описание модели

Линейная регрессия - регрессионная модель зависимости одной переменной от другой с линейной функцией зависимости:

$$y_i = X_i \beta_i + \epsilon_i$$

где X - заданные значения, y - параметры отклика, ϵ - случайная ошибка модели. В случае, если у нас y_i зависит от одного параметра x_i , то модель выглядит следующим образом:

$$y_i = \beta_0 + \beta_1 * x_i + \epsilon_i \quad (4)$$

В данной модели мы пренебрегаем погрешностью и считаем, что она получается при измерении y_i .

2.2.2. Метод наименьших модулей

Для наиболее точного приближения входных с фотоприемников данных y_i линейной регрессией $f(x_i)$ используется метод наименьших квадратов. Этот метод основывается на минимизации нормы разности последовательности:

$$\|f(x_i) - y_i\|_{l^1} \rightarrow \min \quad (5)$$

В данном случае ставится задача линейного программирования, решение которой дает нам коэффициенты β_0 и β_1 , а также вектор множителей коррекции данных w . По итогу получается следующая задача линейного программирования

$$\sum_{i=1}^n |w_i| \rightarrow \min \quad (6)$$

$$\beta_0 + \beta_1 * x_i - w_i * \epsilon \leq y_i, i = 1..n \quad (7)$$

$$\beta_0 + \beta_1 * x_i + w_i * \epsilon \leq y_i, i = 1..n \quad (8)$$

$$1 \leq w_i, i = 1..n \quad (9)$$

2.3. Предварительная обработка данных

Для оценки постоянной, как можно будет увидеть далее, необходима предварительная обработка данных. Займемся линейной моделью дрейфа.

$$Lin_i(n) = A_i + B_i * n, n = 1, 2, \dots N \quad (10)$$

Поставив и решив задачу линейного программирования, найдем коэффициенты A_i , B_i и вектор w_i множителей коррекции данных (где $i = 1$ соответствует данным с ФП1, а $i = 2$ соответственно ФП2). В последствии множитель коррекции данных необходимо применить к погрешностям выборки, чтобы получить данные, которые согласовывались с линейной моделью дрейфа:

$$I_i^f(n) = \dot{x}(n) + \epsilon * w_i(n), n = 1, 2, \dots N \quad (11)$$

По итоге необходимо построить "спрямленные" данные выборки: получить их можно путем вычитания из исходных данных линейную компоненту:

$$I_i^c(n) = I_i^f(n) - B_i * n, n = 1, 2, \dots N \quad (12)$$

2.4. Коэффициент Жаккара

Коэффициент Жаккара - мера сходства множеств. В интервальных данных рассматривается некоторая модификация этого коэффициента: в качестве меры множества (в данном случае интервала) рассматривается его длина, а в качестве пересечения и объединения - взятие минимума и максимума в интервальной арифметике Каухера соответственно. Можно заметить, что в силу возможности минимума по включению быть неправильным интервалом, коэффициент Жаккара может достигать значения только в интервале $[-1; 1]$.

$$JK(x) = \frac{width(\wedge x_i)}{width(\vee x_i)} \quad (13)$$

2.5. Процедура оптимизации

Чтоб найти оптимальный параметр калибровки R_{21} необходимо поставить и решить задачу максимизации коэффициента Жаккара, зависящего от параметра калибровки:

$$JK(I_1^c(n) * R \cup I_2^c(n)) \Rightarrow \max \quad (14)$$

где I_1^c и I_2^c - полученные спрямленные выборки, а R - параметр калибровки. Найденный таким образом R и будет искомым оптимальным R_{21} в силу наибольшего совпадения, оцененного коэффициентом Жаккара.

3. Реализация

Лабораторная работа выполнена на языке программирования Python(3.7) с использованием следующих библиотек: Numpy, Scipy, Tabulate, Statsmodels, Matplotlib.

Отчет написан в онлайн редакторе LaTeX - Overleaf.

4. Результаты

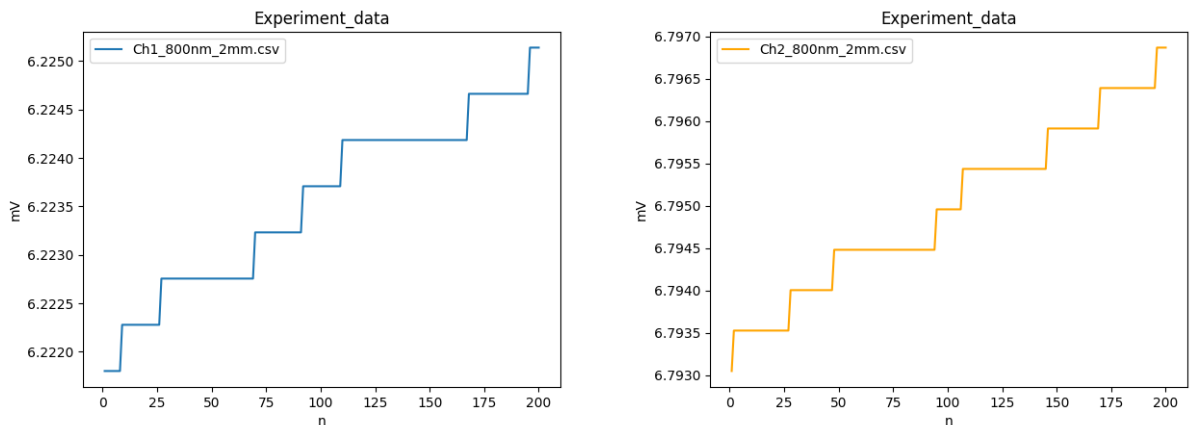


Рис. 2. Исходные данные из экспериментов

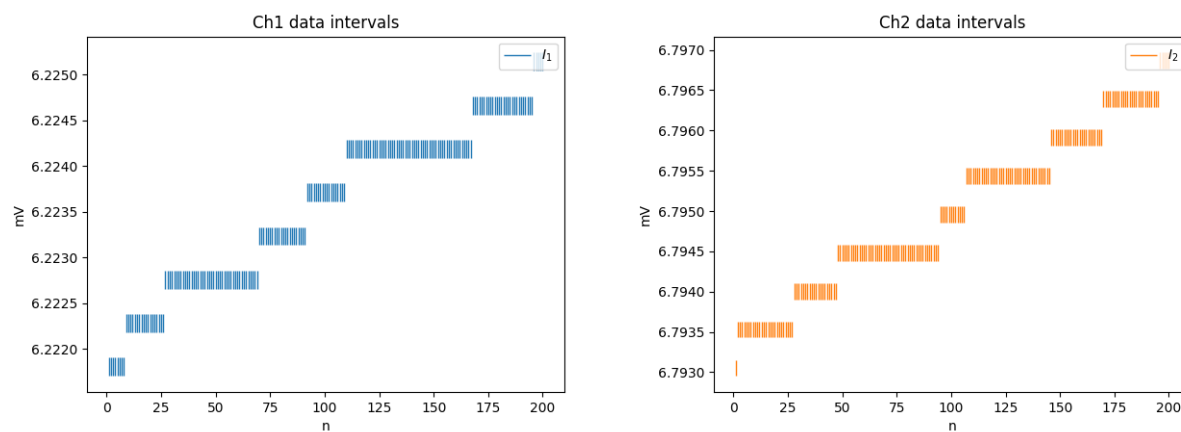


Рис. 3. Интервальное представление исходных данных

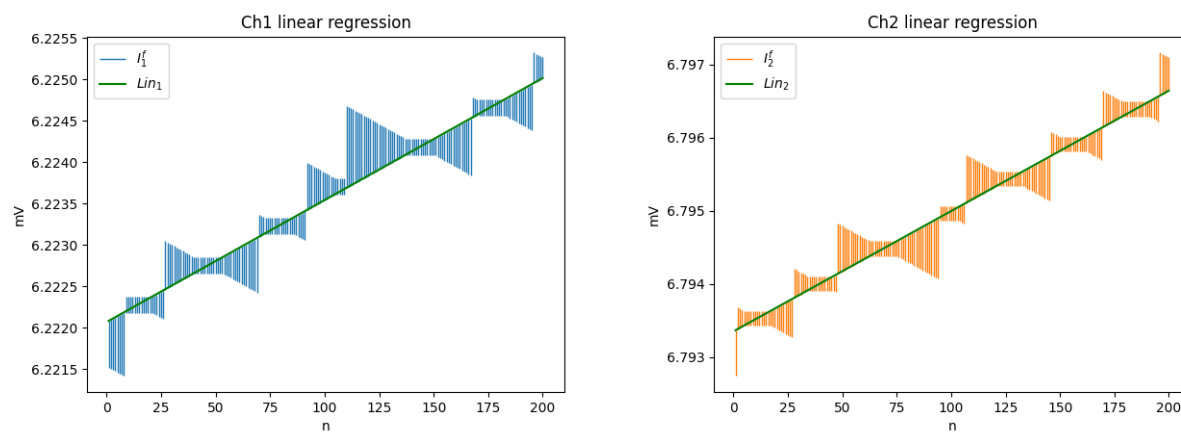


Рис. 4. Линейная модель дрейфа данных

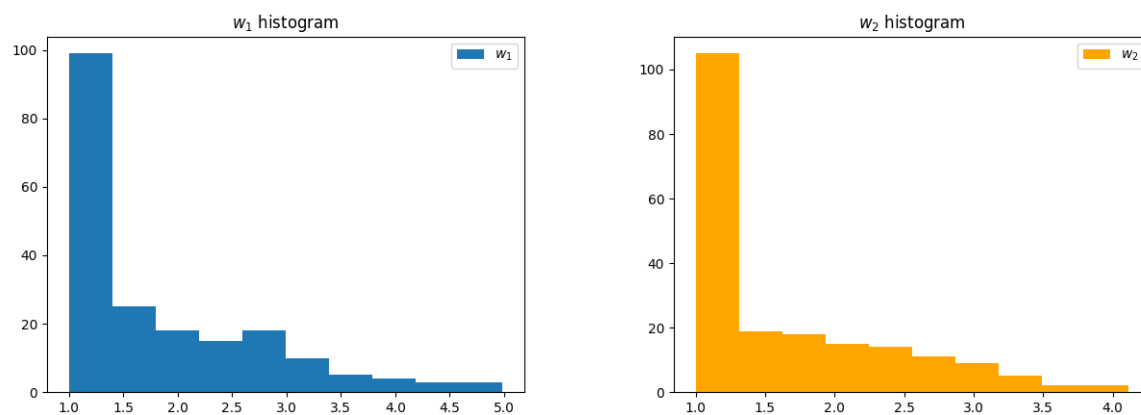


Рис. 5. Гистограммы значений множителей коррекции w

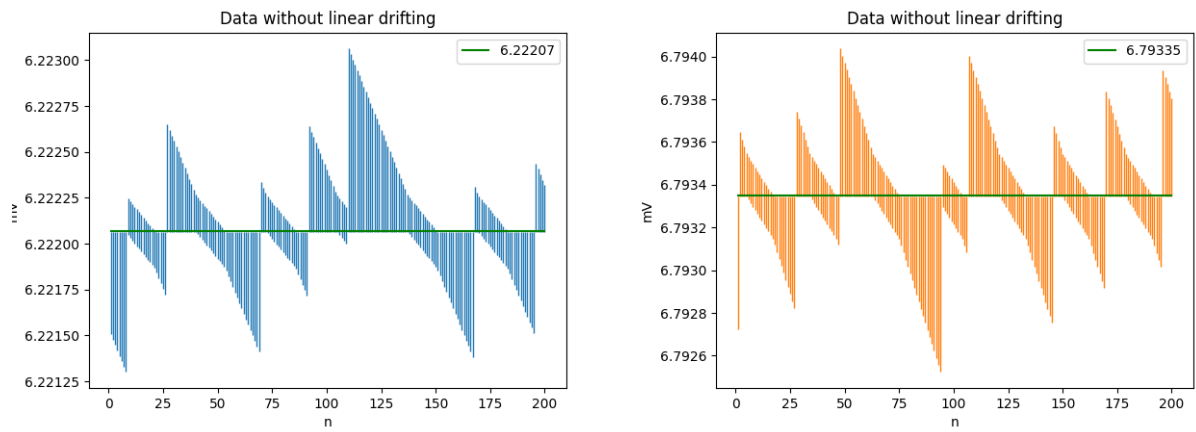


Рис. 6. Скорректированные модели данных

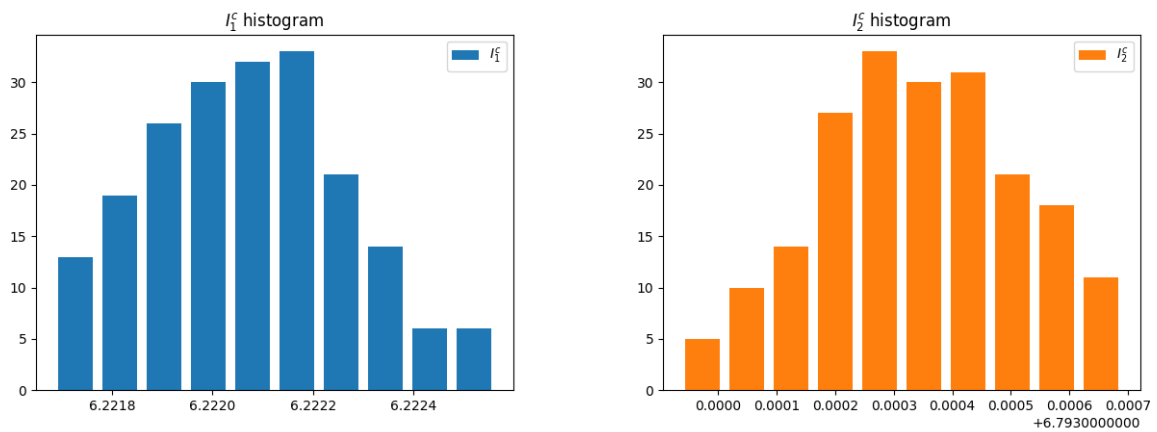
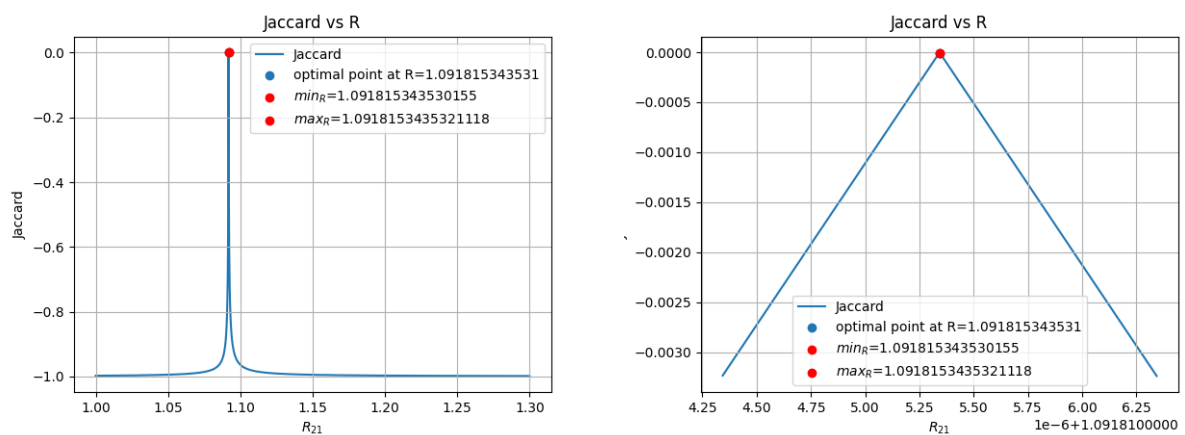


Рис. 7. Гистограммы скорректированных данных

Рис. 8. Значение коэффициента Жаккара от калибровочного множителя от R_{21}

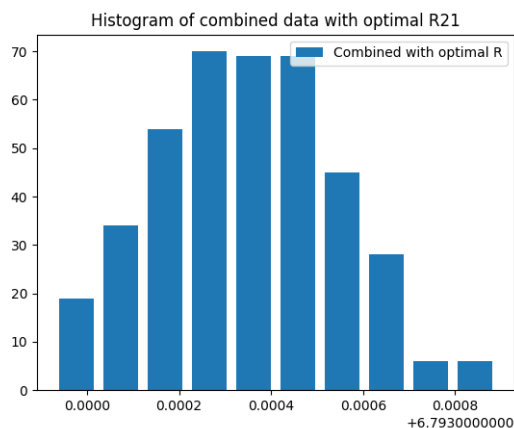


Рис. 9. Гистограмма объединённых данных при оптимальном значении R_{21}

5. Обсуждение

Множители коррекции w . На гистограммах значений множителей коррекции (Рис.5), видно, что половина (для эталлоного фотопередатчика даже больше) не требует коррекции. Этот факт свидетельствует о том, что линейная модель дрейфа данных является разумным приближением.

Коэффициент Жаккара На рис.8 видно, что оптимальным множителем R_{21} является число, равное 1.091815343531. Однако видно, что коэффициент Жаккара при оптимальном значении едва-едва превышает 0, а интервал, при котором $JK \geq 0$, соизмерим с точкой (длина интервала оценивается $10^{-9} - 10^{-10}$). Это показывает на то, что исходные данные имеют ряд неточностей, которые сложно устранить.

Гистограмма объединённых данных при оптимальном значении R . Рассматривая гистограммы скорректированных данных (Рис. 7), можно заметить, что выборки I_1^f и I_2^f имеют характерные "пики" около центра, что имеет отражение на гистограмме объединённых данных (Рис. 9).

6. Ссылка на репозиторий

<https://github.com/PopovIV/MathematicalStatistics>

Список литературы

- [1] А.Н.Баженов Введение в анализ данных с интервальной неопределенностью. 2022.