

Санкт-Петербургский политехнический университет  
Высшая школа прикладной математики и вычислительной физики, ФизМех

Направление подготовки  
«01.03.02 Прикладная математика и информатика»

Отчет по лабораторной работе № 4(10)  
дисциплина "Математическая статистика"

Выполнил студент гр. 5030102/00201  
Преподаватель:

Соболев Д.В.  
Баженов А.Н.

Санкт-Петербург

**2023**

# Содержание

<b>1</b>	<b>Постановка задачи</b>	<b>4</b>
<b>2</b>	<b>Теория</b>	<b>5</b>
<b>3</b>	<b>Результаты</b>	<b>15</b>
3.1	Диаграмма рассеяния . . . . .	15
3.2	Варьирование неопределенности измерений . . . . .	16
3.3	Варьирование неопределённости измерений с расширением и сужением интервалов . . . . .	17
3.4	Анализ регрессионных остатков . . . . .	18
3.5	Информационное множество задачи . . . . .	19
3.6	Коридор совместных зависимостей . . . . .	20
3.7	Построение прогноза внутри и вне области данных . . . . .	20
<b>4</b>	<b>Обсуждение</b>	<b>21</b>
4.1	Варьирование неопределенности измерений . . . . .	21
4.2	Варьирование неопределенности измерений с расширением и сужением интервалов . . . . .	21
4.3	Анализ регрессионных остатков . . . . .	21
4.4	Информационное множество задачи . . . . .	21
4.5	Коридор совместных зависимостей . . . . .	21
4.6	Построение прогноза внутри и вне области данных . . . . .	21
<b>5</b>	<b>Реализация</b>	<b>22</b>

## Список иллюстраций

1	Диаграмма рассеяния выборки $\mathbf{X}_1$ с уравновешенным интервалом погрешности. . . . .	5
2	Диаграмма рассеяния выборки $\mathbf{X}_1$ и регрессионная прямая по модели (4) и (5). . . . .	6
3	Диаграмма рассеяния выборки $\mathbf{X}_1$ и регрессионная прямая по модели (10) и (11) . . . . .	7
4	Векторы $\omega_1$ и $\omega_0$ . . . . .	8
5	Диаграмма рассеяния по модели (4) и (5). . . . .	9
6	Диаграмма рассеяния регрессионных остатков выборки $\mathbf{X}_1$ по (10) и (11). . . . .	10
7	Частоты элементарных подинтервалов регрессионных остатков выборки $\mathbf{X}_1$ по модели (4) и (5) — красный график, и (10) и (11) — синий график. . . . .	10
8	Информационное множество по модели (10) и (11), интервальная оболочка — красный брус. . . . .	12
9	Коридор совместных зависимостей (23). . . . .	13
10	Коридор совместных зависимостей (23). Построение прогноза. . . .	14
11	Кусочно-линейная регрессионная зависимость. . . . .	15
12	Диаграмма рассеяния выборки $\mathbf{X}_1$ с уравновешенным интервалом погрешности . . . . .	16
13	Диаграмма рассеяния выборки $\mathbf{X}_1$ и регрессионная прямая по модели (2.35) и (2.36) . . . . .	16
14	Диаграмма рассеяния выборки $\mathbf{X}_1$ и регрессионная прямая по модели (2.41) и (2.42) . . . . .	17
15	Векторы $\omega_1$ и $\omega_2$ . . . . .	17
16	Диаграмма рассеяния по модели (2.35) и (2.36) . . . . .	18
17	Диаграмма рассеяния регрессионных остатков выборки $\mathbf{X}_1$ по (2.41) и (2.42) . . . . .	18
18	Частоты элементарных подинтервалов регрессионных остатков выборки $\mathbf{X}_1$ по модели (2.35) и (2.36) — красный график, и (2.41) и (2.42) — синий график . . . . .	19
19	Информационное множество по модели (2.41) и (2.42), интервальная оболочка — красный брус . . . . .	19
20	Коридор совместных зависимостей (2.54) . . . . .	20
21	Коридор совместных зависимостей (2.54). Построение прогноза . .	20

# 1 Постановка задачи

Дадим общую формулировку задачи восстановления функциональной зависимости. Пусть некоторая величина  $y$  является функцией от независимых переменных  $x_1, x_2, \dots, x_m$ :

$$y = f(\beta, x) \quad (1)$$

где  $x = (x_1, x_2, \dots, x_m)$  является вектором независимых переменных,  $\beta = (\beta_1, \beta_2, \dots, \beta_p)$  — вектор параметров функции. Заметим, что переменные  $x_1, x_2, \dots, x_m$  также называются входными, а переменные  $y_1$  — выходной.

Задача восстановления функциональной зависимости заключается в том, чтобы, располагая набором значений  $x$  и  $y$ , найти такие  $\beta_1, \beta_2, \dots, \beta_p$  в выражении (1), которые соответствуют конкретной функции  $f$  из параметрического семейства.

Если функция  $f$  является линейной, то можно записать

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m \quad (2)$$

В общем случае результаты измерений величин  $x_1, x_2, \dots, x_m$  и  $y$  являются интервальнозначными

$$x_1^{(k)}, x_2^{(k)}, \dots, x_m^{(k)}, y^{(k)}.$$

Индекс  $k$  пробегает значения от 1 до  $n$ , равного полному числу измерений.

**Определение 2.2.1** Брусом неопределенности  $k$ -го измерения функциональной зависимости будем называть интервальный вектор-брус, образованный интервальными результатами измерений с одинаковыми значениями индекса  $k$  [1]:

$$(x_{k1}, x_{k2}, \dots, x_{km}, y_k) \subset \mathbb{R}^{m+1}, k = 1, 2, \dots, n. \quad (3)$$

Брус неопределенности измерения является прямым декартовым произведением интервалов неопределенности независимых переменных и зависимой переменной.

## 2 Теория

**Данные выборки.** Имеется выборка данных  $X_1$  с интервальной неопределённостью. Число отсчётов в выборке равно 200.

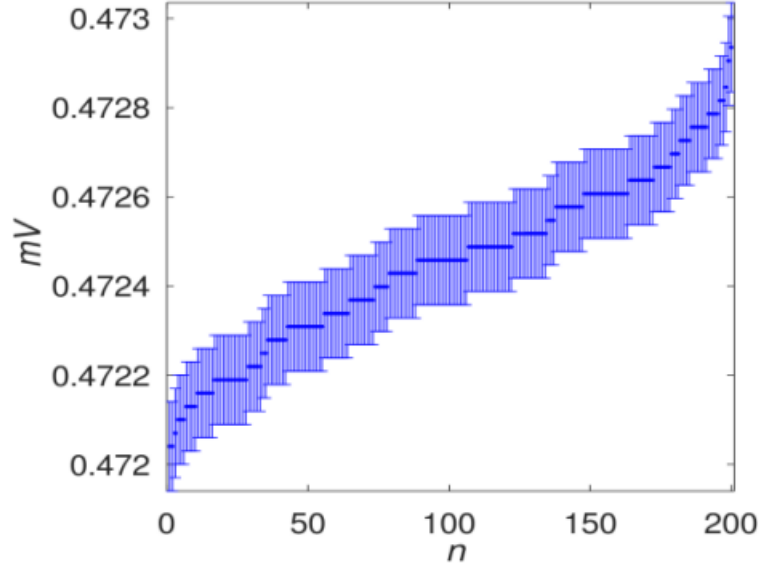


Рис. 1: Диаграмма рассеяния выборки  $X_1$  с уравновешенным интервалом погрешности.

На Рис. 1 представлены данные с прибора [23] с учётом погрешности измерительного прибора.

Построим линейную модель данных и посмотрим, насколько удачно она описывает линейный тренд.

**Варьирование неопределённости измерений.** Если величину коррекции каждого интервального наблюдения выборки выражать коэффициентом его уширения  $\omega_i \geq 1$ , а общее изменение выборки характеризовать суммой этих коэффициентов, то минимальная коррекция выборки в виде вектора коэффициентов  $\omega = (\omega_1, \dots, \omega_n)$ , необходимая для совместности задачи построения зависимости  $x = \beta_0 + \beta_1 * i$  может быть найдена решением задачи условной оптимизации

$$\min_{\omega, \beta} \sum_{i=1}^n \omega_i \quad (4)$$

при ограничениях

$$\begin{cases} mid x_i - \omega_i \epsilon_i \leq \beta_0 + \beta_1 * i \leq mid x_i + \omega_i \epsilon_i, \\ \omega_i \geq 1, \end{cases} \quad i = 1, \dots, n. \quad (5)$$

Результирующие значения коэффициентов  $\omega_i$ , строго превосходящие единицу, указывают на наблюдения, которые требуют уширения интервалов неопределённости для обеспечения совместности данных и модели.

Проведём вычисление параметров линейной регрессии по данным интервальной выборки  $\mathbf{X}_1$  с использованием программ С.И.Жилина [8] и оформленных применительно к задаче на [23]. Синтаксис вызова программы

$$[tau, w, yint] = DataLinearModel(input1, epsilon0) \quad (6)$$

В (6) входами программы служат значения  $mid \mathbf{X}_1$  и величин неопределённости  $\epsilon$ , а выходами  $tau$  — значения параметров регрессии  $\beta_0, \beta_1$   $w$  — вектор весов расширения интервалов.

На Рис. 2 красным цветом приведена регрессионная прямая.

Вычисления с использованием программы (6) дают следующие результаты для регрессионных коэффициентов

$$\beta_0 = tau(1) = 4.7203e - 01, \quad (7)$$

$$\beta_1 = tau(2) = 4.0915e - 06. \quad (8)$$

Все компоненты вектора  $\omega$  оказались равны 1, то есть, расширения интервалов измерений не понадобилось. Таким образом, величина (4) равна числу элементов выборки.

$$\min_{\omega, \beta} \sum_{i=1}^n \omega_i = 200 \quad (9)$$

Недостатком полученного решения с единичными значениями  $\omega_i$  является неучёт расстояний точек регрессионной зависимости до данных интервальной выборки. Таким образом, прямая с параметрами

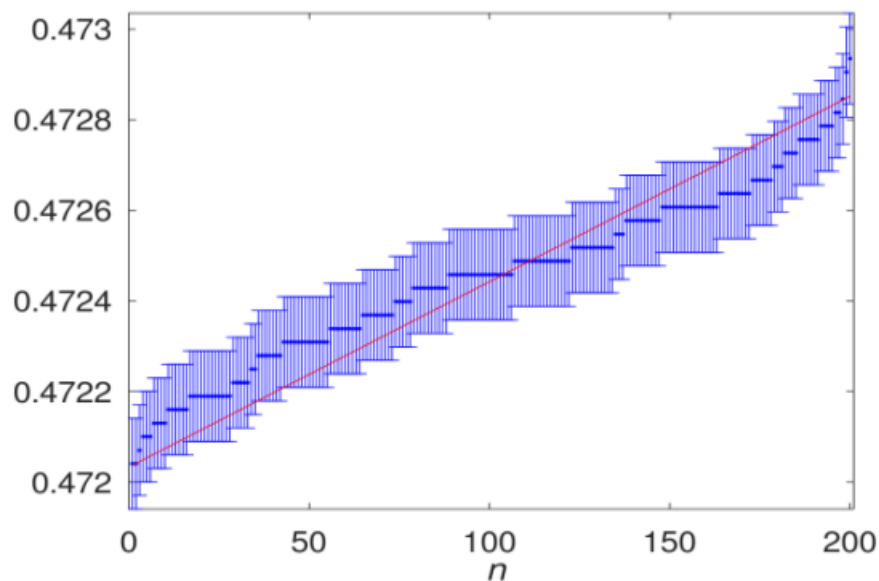


Рис. 2: Диаграмма рассеяния выборки  $\mathbf{X}_1$  и регрессионная прямая по модели (4) и (5).

(7) и (8) «не чувствует» отклонений измерений от прямой на концах выборки — неопределённости измерений достаточно велики, чтобы покрыть этот эффект.

**Варьирование неопределённости измерений с расширением и сужением интервалов.** Выясним, что даёт решение задачи оптимизации другим способом, с расширением и сужением интервалов.

Поставим задачу условной оптимизации следующим образом:

$$\min_{\omega, \beta} \sum_{i=1}^n \omega_i \quad (10)$$

при ограничениях

$$\begin{cases} mid x_i - \omega_i \epsilon_i \leq \beta_0 + \beta_1 * i \leq mid x_i + \omega_i \epsilon_i, \\ \omega_i \geq 0, \end{cases} \quad i = 1, \dots, n. \quad (11)$$

Отличие постановки от (4) и (5) состоит в том, что интервалы измерений могут как расширяться в случае  $\omega_i \geq 1$ , так и сужаться при  $0 \leq \omega_i \leq 1$ . Вычисление параметров линейной регрессии по данным интервальной выборки  $\mathbf{X}_1$  производится как и в случае (6) с использованием программ С.И.Жилина [8] и оформленных применительно к задаче на [23]. Синтаксис вызова программы

$$[tau, w, yint] = DataLinearModelZ(input1, epsilon0) \quad (12)$$

Входы и выходы функции DataLinearModelZ такие же, как и для DataLinearModelZ (6).

На Рис. 3 красным цветом приведена регрессионная прямая.

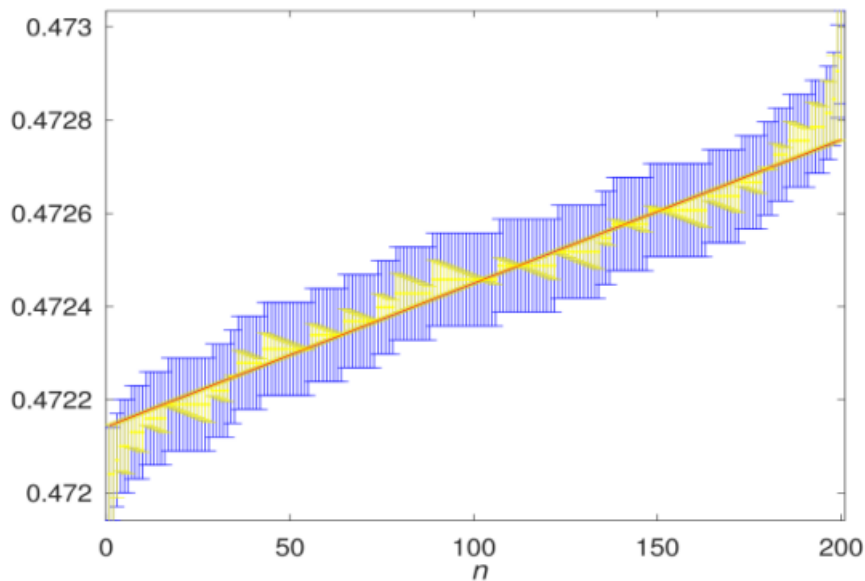


Рис. 3: Диаграмма рассеяния выборки  $\mathbf{X}_1$  и регрессионная прямая по модели (10) и (11)

Жёлтым цветом на Рис. 3 показаны скорректированные интервалы выборки  $\mathbf{X}_1$ . Небольшая часть интервалов на границах области расширилась, а большинство интервалов в диапазоне замеров примерно от 20 до 180 — сузилось.

Величина меры (4) уменьшилась более, чем в 4 раза.

$$\min_{\omega, \beta} \sum_{i=1}^n \omega_i = 45.7 \leq 200 \quad (13)$$

Таким образом, постановка задачи с возможностью одновременного увеличения и уменьшения радиусов неопределённости измерений позволяет более гибко подходить к задаче оптимизации.

На Рис. 4 приведены графики векторов  $\omega_0$  и  $\omega_1$ , полученных при использовании двух рассмотренных подходов.

В конкретном случае график вектора  $\omega_0$  для постановки задачи оптимизации (10) и (11) содержит большое количество информации.

Например, задавшись каким-то порогом  $\alpha$ :  $0 < \alpha \leq 1$ , можно выделить области входного аргумента  $\Psi$ , в которых регрессионная зависимость хуже соответствует исходным данным. Например:

$$\Psi = \arg_i \omega_i \geq \alpha \quad (14)$$

Для конкретного примера имеем две области  $\Psi$  в начале и конце области данных.

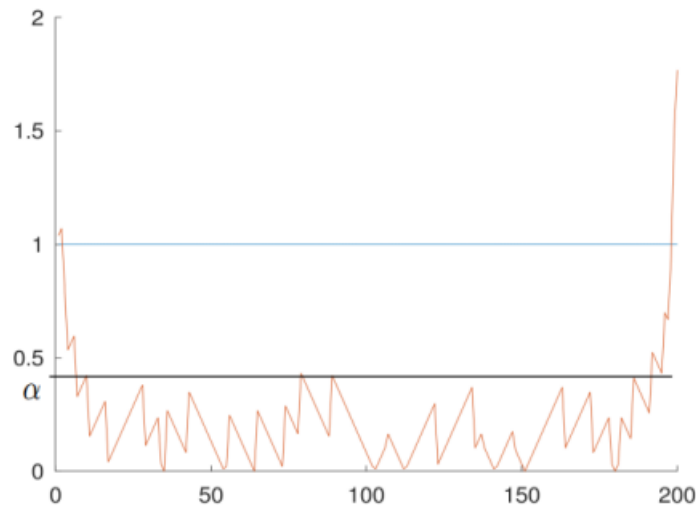


Рис. 4: Векторы  $\omega_1$  и  $\omega_0$ .

Для объективного использования этого приёма параметр  $\alpha$  можно брать, например, из анализа гистограммы распределения вектора  $\omega$ .

Использование выделения «подозрительных» областей даёт основу для других приёмов. Например, для построения кусочно-линейной регрессионной зависимости.



**Анализ регрессионных остатков.** В теоретико-вероятностной математической статистике анализ регрессионных остатков — один из приёмов оценки качества регрессии.

Приведём пример пояснения этого приёма. «Если выбранная регрессионная модель хорошо описывает истинную зависимость, то остатки должны быть независимыми, нормально распределёнными случайными величинами с нулевым средним, и в их значениях должен отсутствовать тренд. Анализ регрессионных остатков — это процесс проверки выполнения этих условий.» <https://wiki.loginom.ru/articles/discrep>

В случае интервальных выборок мы не задаёмся вопросом о виде распределения остатков, а будем использовать те возможности которые появляются при описании объектов и результатов вычислений в виде интервалов.

На Рис. 5 приведена диаграмма рассеяния регрессионных остатков выборки  $X_1$  по модели (4) и (5).

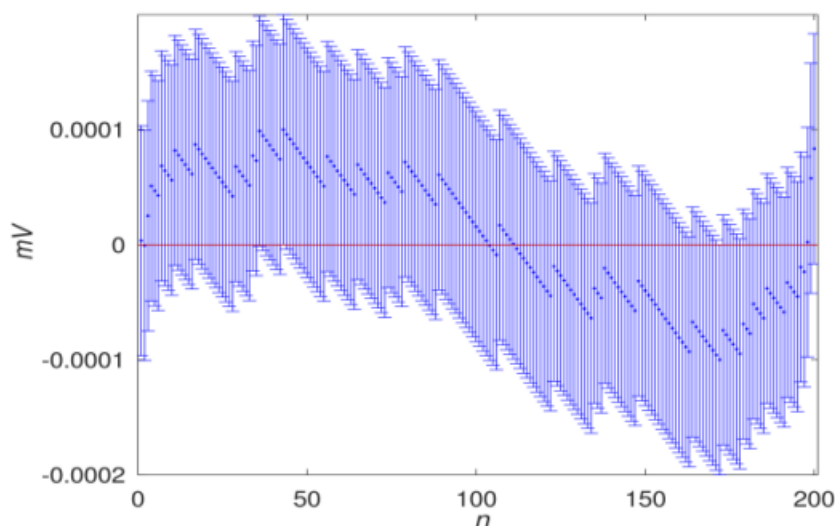


Рис. 5: Диаграмма рассеяния по модели (4) и (5).

На Рис. 6 приведена диаграмма рассеяния регрессионных остатков выборки  $X_1$  по модели (10) и (11).

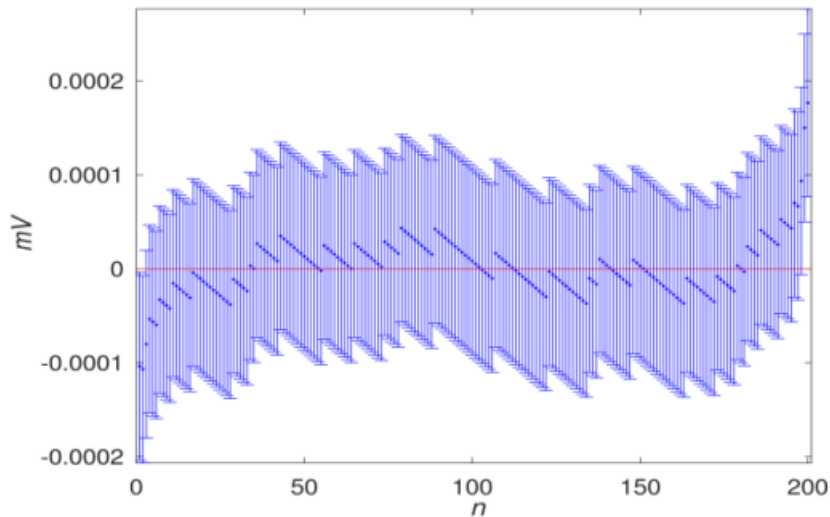


Рис. 6: Диаграмма рассеяния регрессионных остатков выборки  $X_1$  по (10) и (11).

Из сравнения Рис. 5 и На Рис. 6 видно, что интервальные выборки остатков получились с весьма разными свойствами. Формально диаграмма рассеяния на первом рисунке ‘уже, то есть внешняя оценка более компактная. В то же время вторая диаграмма рассеяния выглядит более естественно.

На Рис. 7 приведены графики частот элементарных подинтервалов при вычислении интервальной моды для двух моделей.

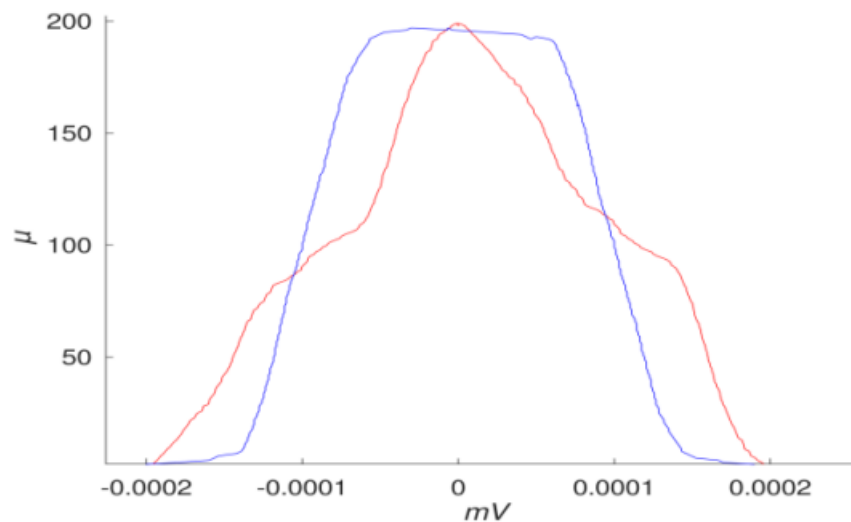


Рис. 7: Частоты элементарных подинтервалов регрессионных остатков выборки  $X_1$  по модели (4) и (5) — красный график, и (10) и (11) — синий график.

Как и в случае анализа диаграмм рассеяния, второй график выглядит более естественно. Его внутренняя оценка существенно шире, что соответствует большей устойчивостью к возмущениям данных.

К остаткам можно применить и другие меры совместности оценки постоянной величины, описанные ранее.

$$mode \mathbf{X}^1 = \dots \quad (15)$$

$$Ji(\mathbf{X})^1 = \dots \quad (16)$$

$$\vdots \quad (17)$$

$$mode \mathbf{X}^2 = \dots \quad (18)$$

$$Ji(\mathbf{X})^2 = \dots \quad (19)$$

$$\vdots \quad (20)$$

здесь  $\mathbf{X}^{1,2}$  — регрессионные остатки выборки  $\mathbf{X}_1$ , вычисленные с использованием разных условий оптимизации.

**Информационное множество задачи.** Интервальные оценки параметров.

Один из главных вопросов при построении регрессии — оценивание её параметров. В зависимости от прикладных целей характер и назначение искомых оценок могут существенно разниться.

Внешняя интервальная оценка параметра определяется минимальным и максимальным значениями, которых может достигать значение параметра в информационном множестве.

В совокупности интервальные оценки параметров задают брус, описанный вокруг информационного множества и именуемый внешней интервальной оболочкой информационного множества:

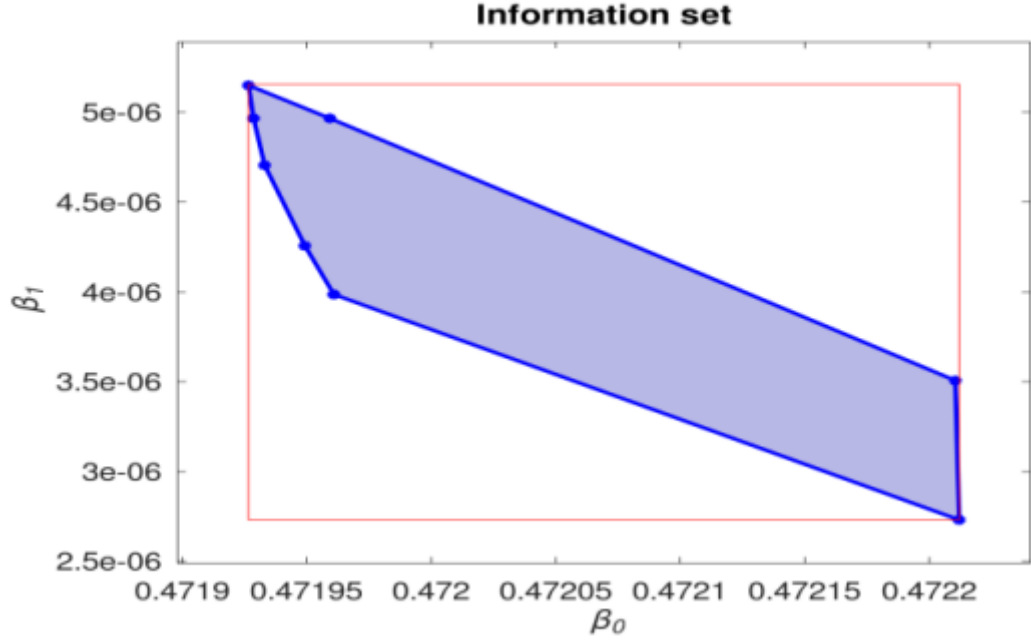


Рис. 8: Информационное множество по модели (10) и (11), интервальная оболочка — красный брус.

Проведём вычисление параметров линейной регрессии по данным интервальной выборки  $\mathbf{X}_1$  с использованием программ С.И.Жилина [8].

Синтаксис вызова программ:

Решение задачи линейного программирования

$$SS = ir\_problem(A, x, max(w0) * epsilon, lb);$$

Вершины информационного множества задачи построения интервальной регрессии

$$vertices = ir\_beta2poly(SS);$$

Внешние интервальные оценки параметров модели  $y = \beta_1 + \beta_2 * x$

$$b_{int} = ir\_outer(SS).$$

Входами программы служат значения  $mid \mathbf{X}_1$  и величин неопределённости  $\epsilon$ , умноженные на расчётное уширение по модели (10) и (11), матрица  $A$ , составленная из нулевой и первой степеней номеров замеров, параметры условной оптимизации. Структура  $SS$  содержит значения параметров регрессии.

**Коридор совместных зависимостей.** Информационное множество задачи определяется в пространстве параметров. Каждая его точка задаёт зависимость в пространстве переменных. Множество всех таких моделей именуется коридором совместных зависимостей.

Выше мы нашли внешние интервальные оценки параметров модели

$$\text{mid } \beta_0 = [4.7193e - 01, 4.7221e - 01], \quad (21)$$

$$\text{mid } \beta_1 = [2.7304e - 06, 5.1571e - 06]. \quad (22)$$

Подставляя значения (21) и (22) в уравнение регрессии, получаем

$$x(k) = \text{mid } \beta_0 + \text{mid } \beta_1 * k, \quad (23)$$

где  $k$  — номер измерения.

На Рис. 9 приведён коридор совместных зависимостей для модели (23). Визуально видно, что внутри коридор совместных зависимостей можно провести множество прямых.

**Построение прогноза внутри и вне области данных.** Одним из способов использования регрессионной модели является предсказание значений выходной переменной для заданных значений входной. С помощью построенной выше модели (23) можно получить прогнозные значения выходной переменной в точках эксперимента.

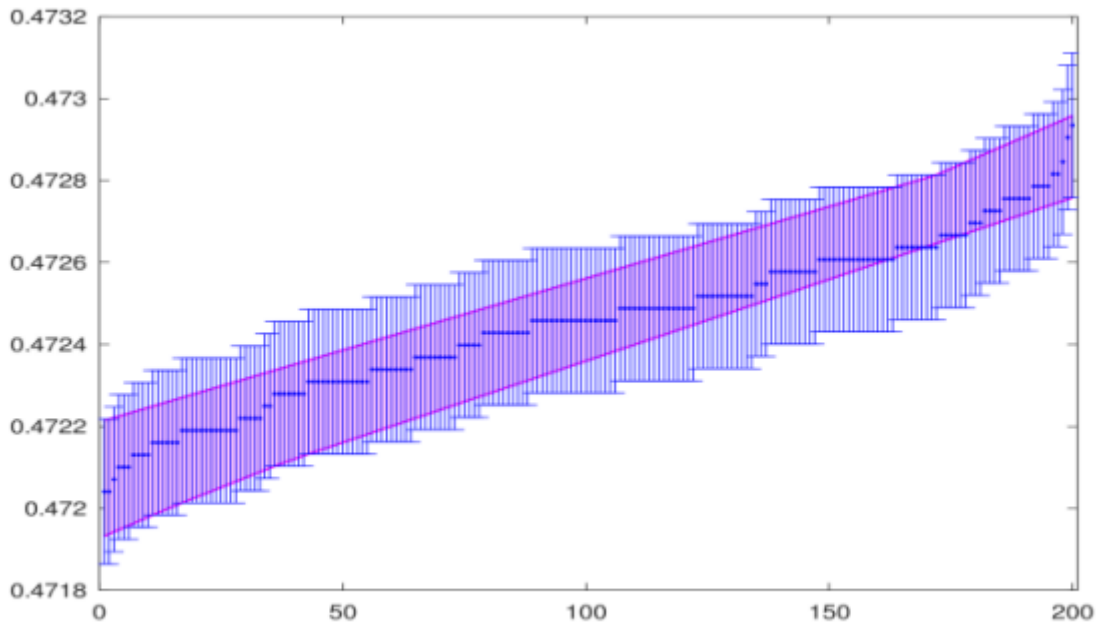


Рис. 9: Коридор совместных зависимостей (23).

Ценность модели также заключается в возможности её употребления для предсказания выходной переменной в точках, где измерения не производились.

Расширив область определения аргумента для модели (23), можно получить оценки для значений выходной переменной (экстраполяция). На Рис. 10 сплошной

заливкой дан прогноз в том числе за пределами данных интервальной выборки  $X_1$ .

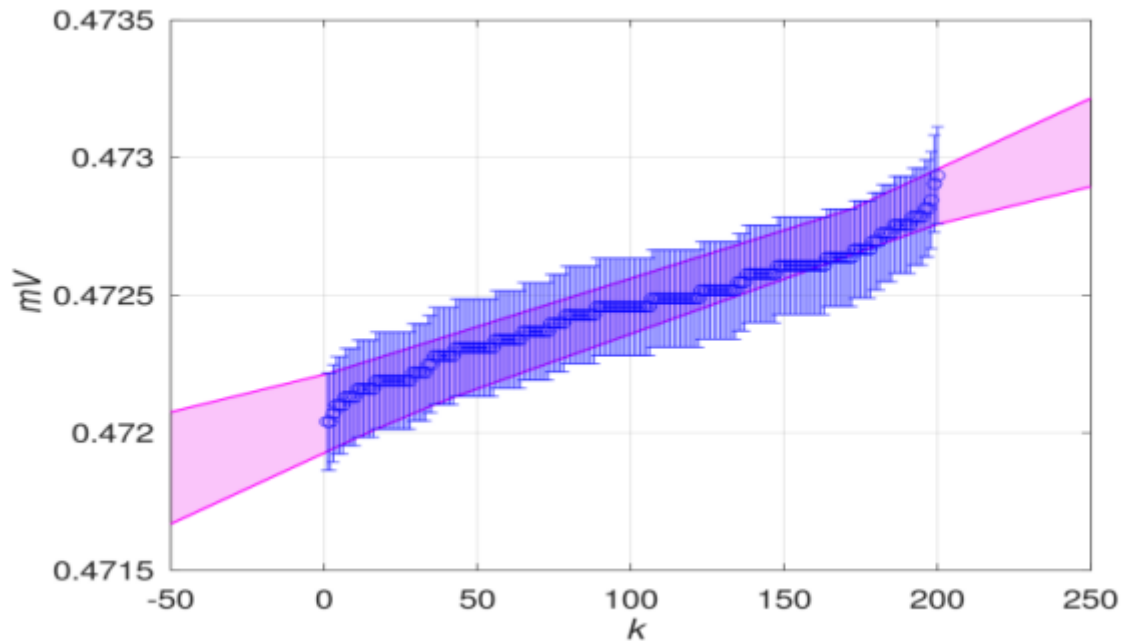


Рис. 10: Коридор совместных зависимостей (23). Построение прогноза.

Следует обратить внимание, что величина неопределённости прогнозов растёт по мере удаления от области, в которой производились исходные измерения. Это обусловлено видом коридора зависимостей, расширяющимся за пределами области измерений, и согласуется со здравым смыслом.

**Уточнение структуры модели. Кусочно-линейная регрессионная зависимость.** Рис. 5 и Рис. 6 регрессионных остатков свидетельствуют о том, что линейные регрессионные модели не вполне точно отражают характер зависимости для интервальной выборки  $X_1$ . Наиболее простым способом учёта этого факта является использование кусочно-линейная регрессионной зависимости.

В разделе «Варьирование неопределённости измерений» были вычислены векторы весов  $\omega$  расширения неопределённости измерений для достижения совместности — см. Рис. 4. Резкое возрастание весов  $\omega$  на границах области определения свидетельствует о несоответствии данных и модели. Эти точки и можно взять как «угловые» для определения линейных участков.

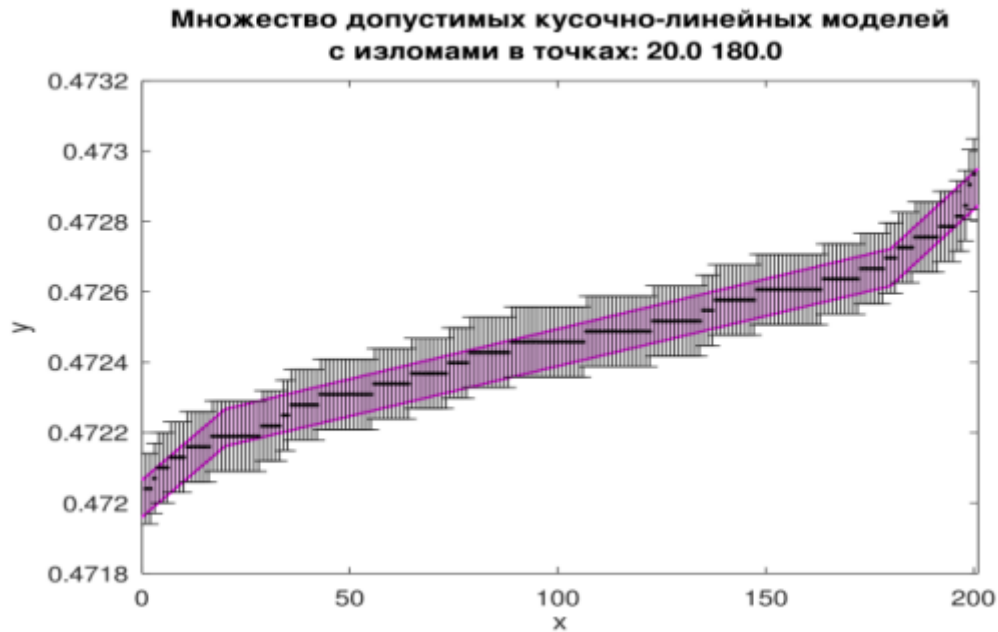


Рис. 11: Кусочно-линейная регрессионная зависимость.

На Рис. 11 показан пример построения кусочно-линейной регрессионной зависимости и коридора совместных зависимостей. После вычитания модели, можно переходить к анализу отстатков регрессии и другим приёмам анализа.

В более общей постановке ставится задача автоматического определения точек излома [29], [30]. Имеется программное обеспечение С.И.Жилина, реализующее идеи этого подхода.

## 3 Результаты

### 3.1 Диаграмма рассеяния

Данные для выборки взяты из файла

octave/Channel\_1\_500nm\_0\_23mm.csv, погрешность прибора  $\epsilon = 10^{-4}$ .

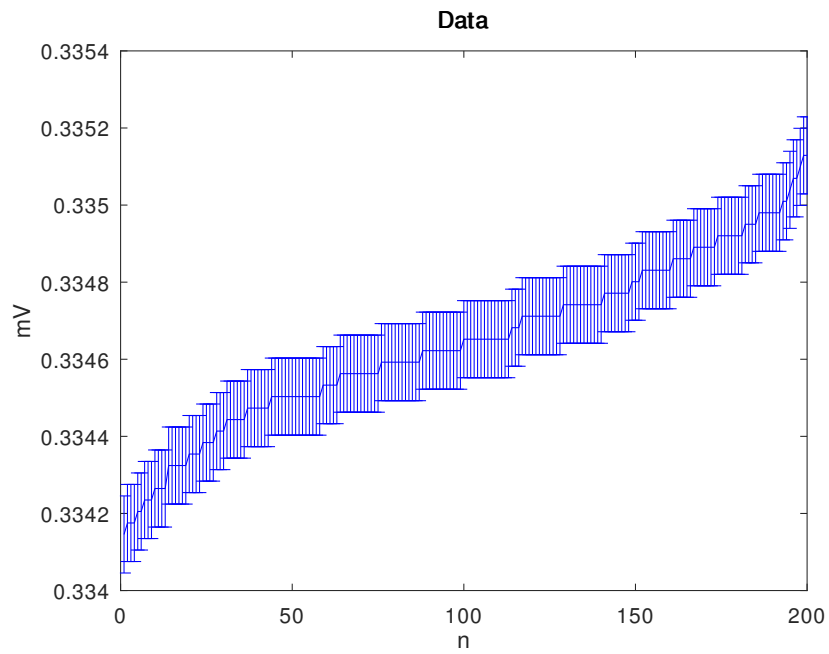


Рис. 12: Диаграмма рассеяния выборки  $\mathbf{X}_1$  с уравновешенным интервалом погрешности

### 3.2 Варьирование неопределенности измерений

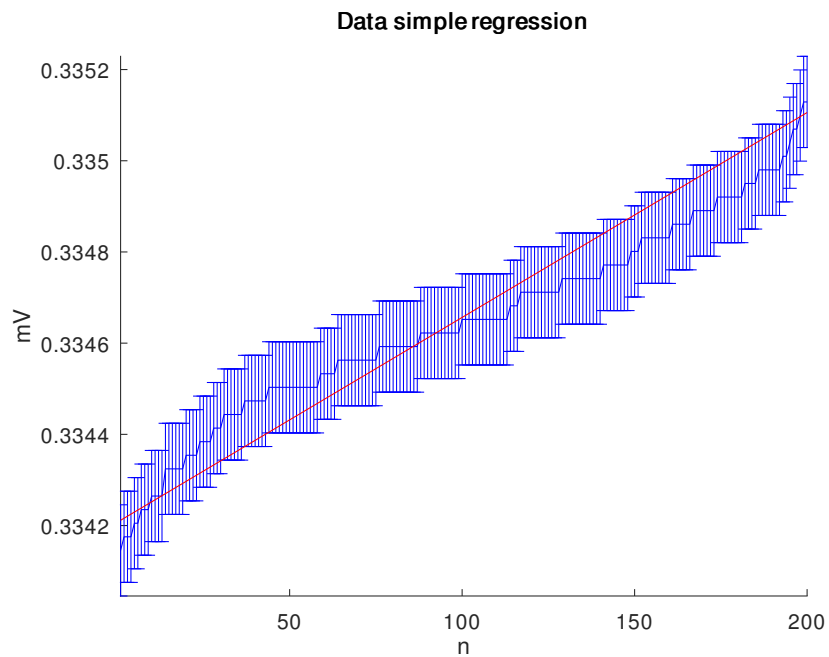


Рис. 13: Диаграмма рассеяния выборки  $\mathbf{X}_1$  и регрессионная прямая по модели (2.35) и (2.36)

$$\sum_{i=1}^n \omega_i = 200, \beta_0 = 0.334207, \beta_1 = 4.5 \cdot 10^{-6}$$



### 3.3 Варьирование неопределённости измерений с расширением и сужением интервалов

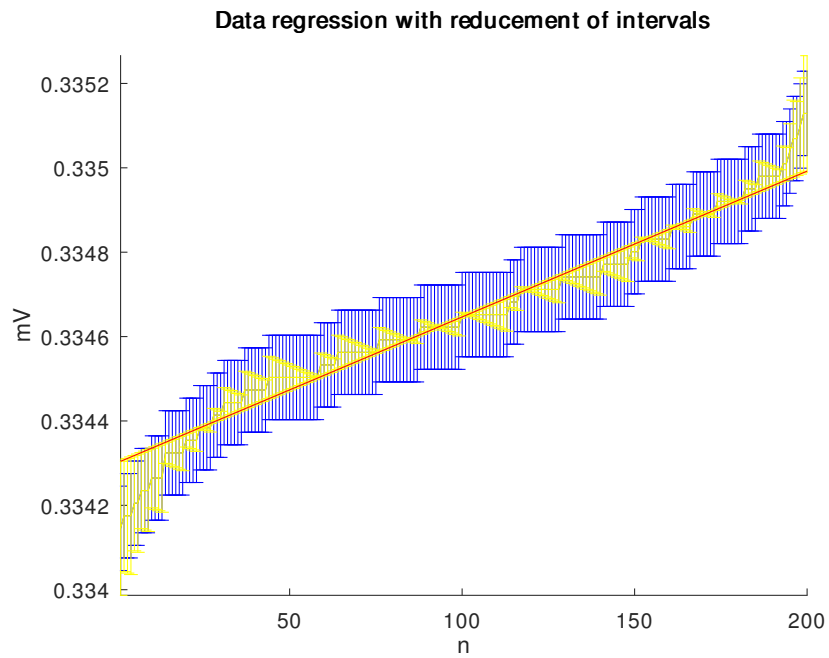


Рис. 14: Диаграмма рассеяния выборки  $\mathbf{X}_1$  и регрессионная прямая по модели (2.41) и (2.42)

$$\sum_{i=1}^n \omega_i = 54.021, \beta_0 = 0.334301, \beta_1 = 3.5 \cdot 10^{-6}$$

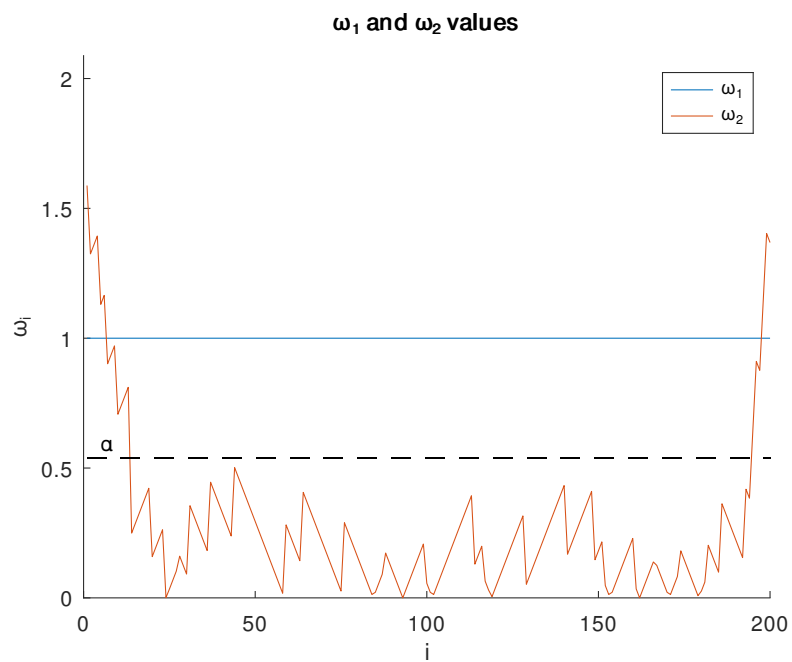


Рис. 15: Векторы  $\omega_1$  и  $\omega_2$

### 3.4 Анализ регрессионных остатков

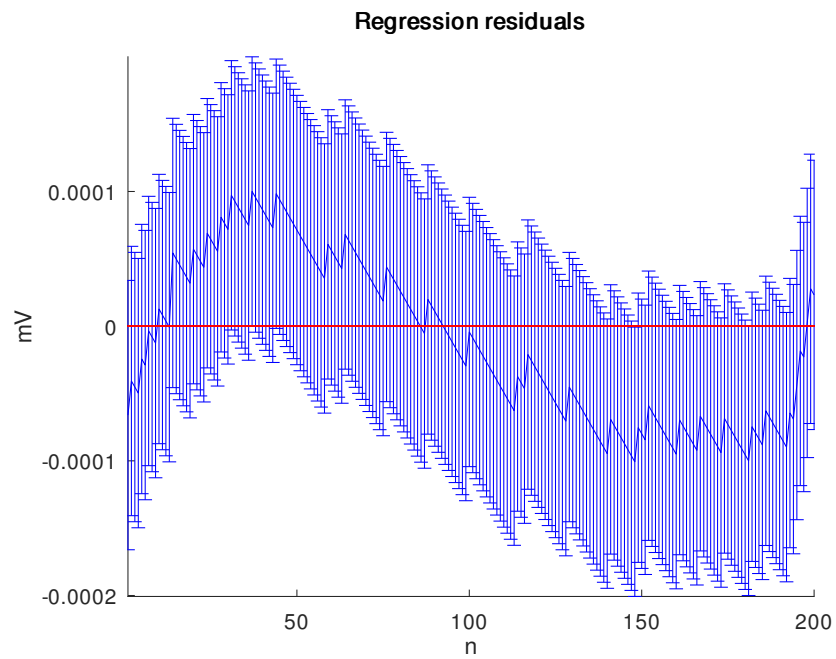


Рис. 16: Диаграмма рассеяния по модели (2.35) и (2.36)

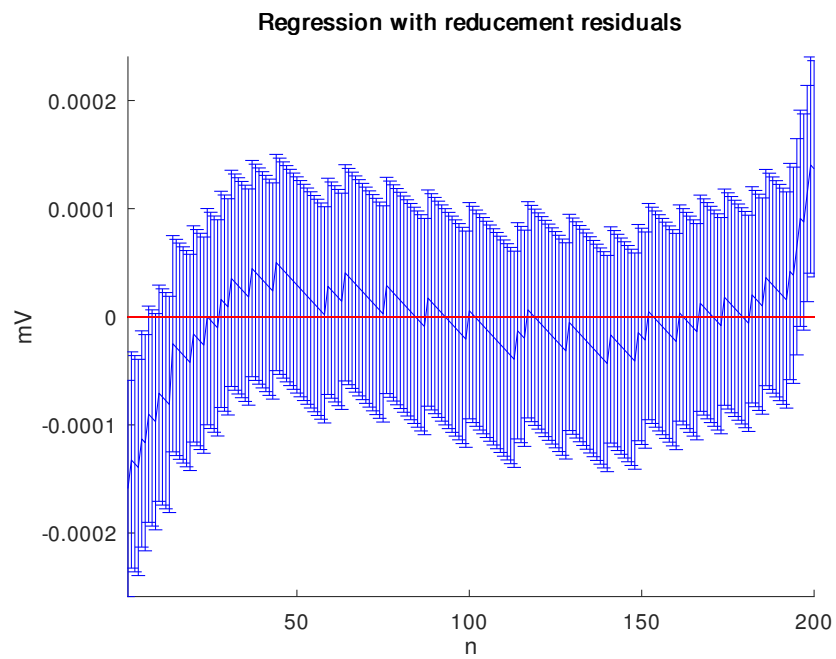


Рис. 17: Диаграмма рассеяния регрессионных остатков выборки  $\mathbf{X}_1$  по (2.41) и (2.42)

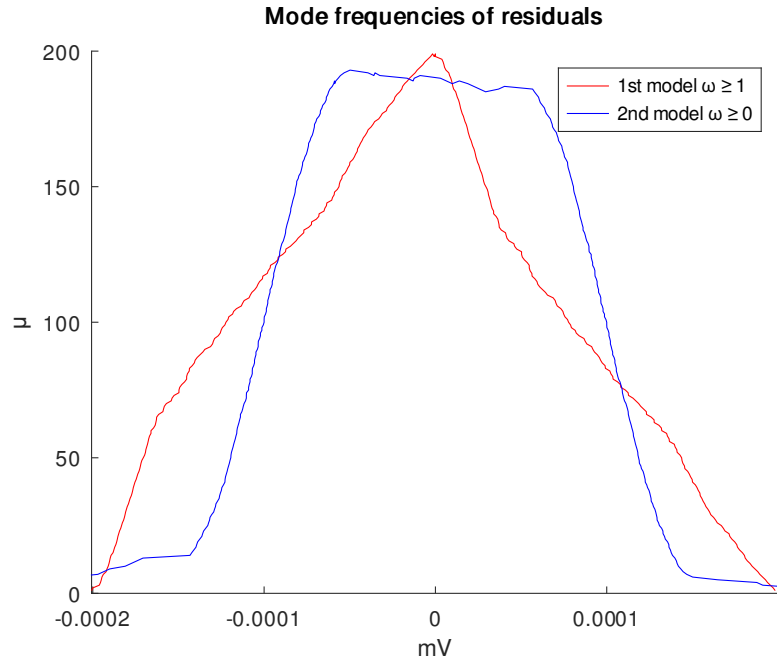


Рис. 18: Частоты элементарных подинтервалов регрессионных остатков выборки  $\mathbf{X}_1$  по модели (2.35) и (2.36) — красный график, и (2.41) и (2.42) — синий график

### 3.5 Информационное множество задачи

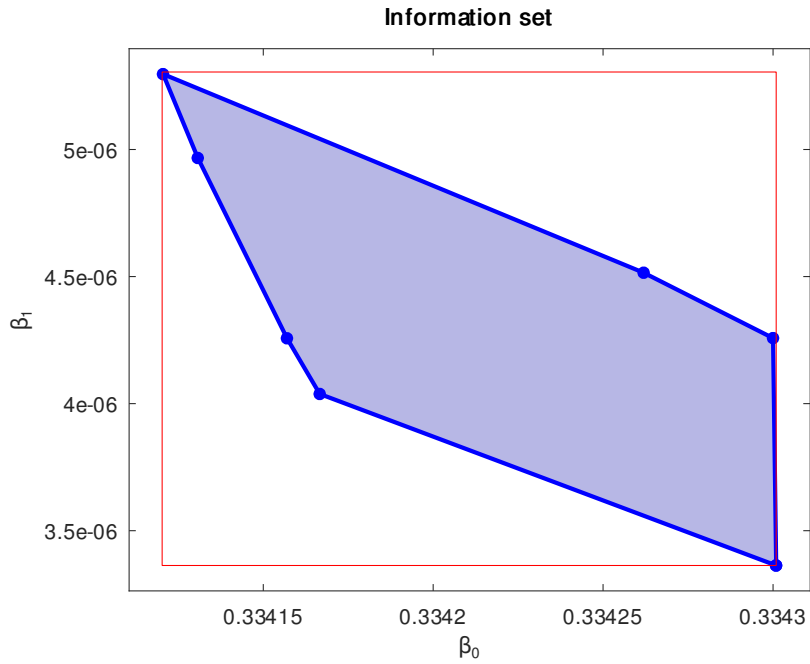


Рис. 19: Информационное множество по модели (2.41) и (2.42), интервальная оболочка — красный брус

$$\text{mid } \beta_0 = [0.334120, 0.334301]$$

$$\text{mid } \beta_1 = [3.4 \cdot 10^{-6}, 5.3 \cdot 10^{-6}]$$

### 3.6 Коридор совместных зависимостей

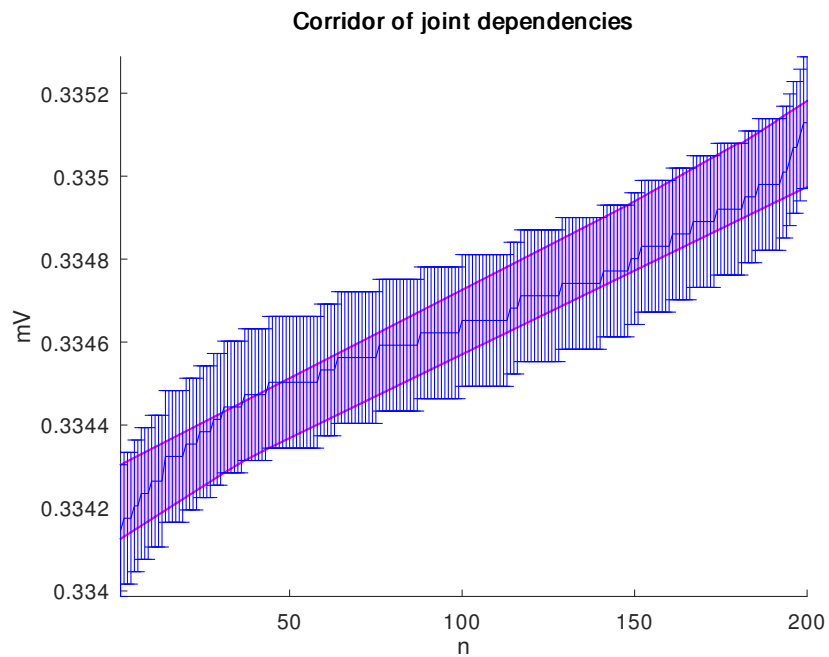


Рис. 20: Коридор совместных зависимостей (2.54)

### 3.7 Построение прогноза внутри и вне области данных

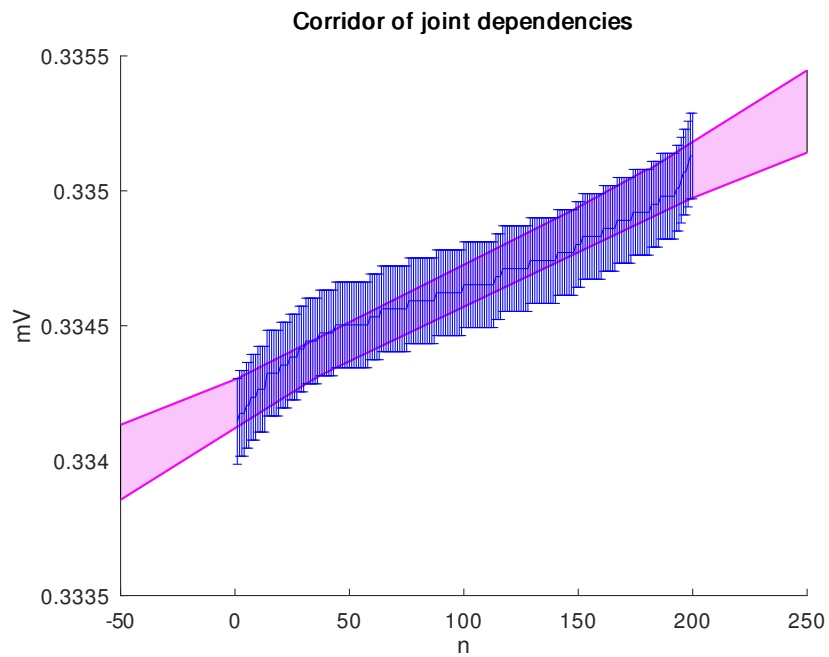


Рис. 21: Коридор совместных зависимостей (2.54). Построение прогноза

## 4 Обсуждение

### 4.1 Варьирование неопределенности измерений

Для модели регрессии с  $\omega_i \geq 1$  видим, что все  $\omega_i = 1$ , а регрессионная прямая действительно пересекает каждый отрезок без необходимости увеличения какого-либо из них.

### 4.2 Варьирование неопределенности измерений с расширением и сужением интервалов

Для модели регрессии с  $\omega_i \geq 0$  видим, что для большинства интервалов  $\omega_i < 1$ , однако в начале и конце имеются выбросы  $\omega_i \approx 1.5$ . Также из рисунка видно, что регрессионная прямая пересекает уже не все интервалы. Это объясняется тем, что некоторые из них были увеличены, и регрессионная прямая пересекает измененные интервалы (желтые), притом пересекая увеличенный интервал она вовсе не обязана пересечь исходный.

### 4.3 Анализ регрессионных остатков

По результатам вычислений для регрессионных остатков можно сделать вывод, что мода регрессионных остатков по модели с  $\omega_i \geq 0$  представляет собой более широкую окрестность нуля. Это означает, что регрессия по этой модели качественнее, нежели по модели  $\omega_i \geq 1$ .

### 4.4 Информационное множество задачи

Из графика для информационного множества задачи видим, что решение классических численных задач в интервальных постановках является вовсе не интервалом, а многогранным множеством. Притом независимые оценки для компонент решения можно дать, построив интервальную оболочку.

### 4.5 Коридор совместных зависимостей

По результатам построения коридора совместных зависимостей получено, как нетрудно видеть, множество, любая прямая, лежащая в котором, будет являться совместной регрессионной зависимостью для данной интервальной выборки.

### 4.6 Построение прогноза внутри и вне области данных

На уменьшенном рисунке для коридора совместных зависимостей видим, что он сужается ближе к центру выборки по  $n$  и расширяется при отдалении от цен-

тра.

## 5 Реализация

Лабораторная работа выполнена с помощью пакета GNU Octave 8.2.0.

Ссылка на исходный код лабораторной: <https://github.com/Parampaika/MatStat/tree/main/Lab10>