

Математические методы в физике.
Курс для аспирантов ФТИ им.
А.Ф.Иоффе.
Практикум.

А.Н.Баженов, А.А.Карпова

26 апреля 2023 г.

Оглавление

1	Интервальный анализ	3
1.1	Особенность матрицы	3
1.2	Решение линейной задачи о допусках	3
1.2.1	Множества АЕ-решений интервальных систем линейных алгебраических уравнений	3
2	Анализ данных с интервальной неопределённостью	4
2.1	Оценка постоянной величины для интервальной выборки	4
2.1.1	Постановка задачи и предварительные оценки . .	4
2.1.2	Меры совместности выборок.	6
2.2	Задача восстановления линейной зависимости	15
2.2.1	Постановка задачи	16
2.2.2	Совместность линейной функциональной зависимости	17
2.2.3	Информационное множество и коридор совместных зависимостей	19

Глава 1

Интервальный анализ

1.1 Особенность матрицы

1.2 Решение линейной задачи о допусках

1.2.1 Множества АЕ-решений интервальных систем линейных алгебраических уравнений

Перейдем к рассмотрению *интервальных систем линейных алгебраических уравнений* (ИСЛАУ) вида

$$\left\{ \begin{array}{l} \mathbf{a}_{11}x_1 + \mathbf{a}_{12}x_2 + \dots + \mathbf{a}_{1n}x_n = \mathbf{b}_1, \\ \mathbf{a}_{21}x_1 + \mathbf{a}_{22}x_2 + \dots + \mathbf{a}_{2n}x_n = \mathbf{b}_2, \\ \vdots \qquad \qquad \qquad \vdots \qquad \qquad \ddots \qquad \qquad \vdots \qquad \qquad \vdots \\ \mathbf{a}_{m1}x_1 + \mathbf{a}_{m2}x_2 + \dots + \mathbf{a}_{mn}x_n = \mathbf{b}_m, \end{array} \right. \quad (1.1)$$

с интервалами \mathbf{a}_{ij} и \mathbf{b}_j , $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$, или, кратко,

$$\mathbf{A}x = \mathbf{b}, \quad (1.2)$$

где $\mathbf{A} = (\mathbf{a}_{ij})$ — интервальная матрица размеров $m \times n$, а $\mathbf{b} = (\mathbf{b}_i)$ является интервальным m -вектором.

Глава 2

Анализ данных с интервальной неопределённостью

2.1 Оценка постоянной величины для интервальной выборки

2.1.1 Постановка задачи и предварительные оценки

Предметная область относится к физике полупроводников — исследованиям фотоэлектрических характеристик испытываемого датчика, проводимым специалистами лаборатории фотоэлектрических преобразователей Физико-технического института им. А. Ф. Иоффе [23]. Развёрнутое описание задачи дано в статье [5].

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

На Рис. 2.1 представлены сырые данные с прибора [23].

Для дальнейшей работы используется модель данных с уравновешенным интервалом погрешности.

$$x = \hat{x} + \epsilon; \quad \epsilon = [-\epsilon, \epsilon] \quad \text{для некоторого } \epsilon > 0, \quad (2.1)$$

Здесь \hat{x} —данные прибора, $\epsilon = 10^{-4}$ — погрешность прибора [23].

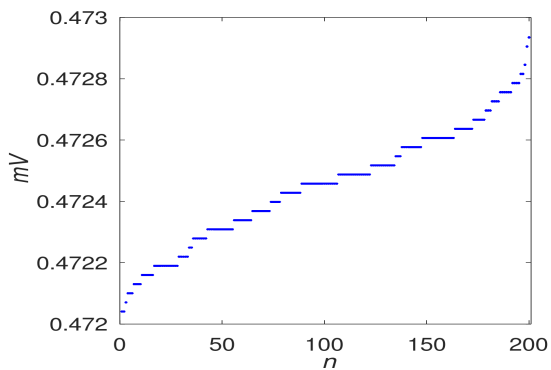


Рис. 2.1: Данные выборки X_1 [23].

Диаграмма рассеяния. Привести диаграмму рассеяния выборки с учётом погрешности прибора.

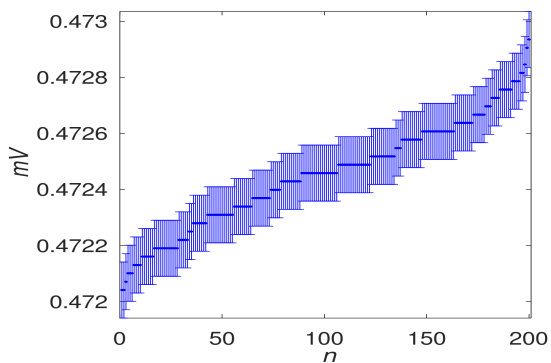


Рис. 2.2: Диаграмма рассеяния выборки X_1 с уравновешенным интервалом погрешности (2.1).

Оценки исходной выборки. Вычислим базовые оценки исходной выборки. Эти выборки несовместны и их информационные множества

пусты. Внешние оценки найдем как

$$\underline{J} = \min_{1 \leq k \leq n} \underline{x}_k, \quad \overline{J} = \max_{1 \leq k \leq n} \overline{x}_k. \quad (2.2)$$

Вычисления дают следующие результаты:

$$J_1 = 0.47194, \quad \text{wid } J_1 = 0.47304. \quad (2.3)$$

Верхние и нижние вершины оценок J_1 совпадают с границами отображения на рис. 2.2.

2.1.2 Меры совместности выборок.

Оценим выборку с помощью набора мер совместности [4].

Сначала значения мер совместности качества возьмём в исходном, ненормированном виде:

1. Размер максимальной клики $\max \mu_j$
2. Величина коэффициента вариабельности по Оскорбину k_O
3. Мера совместности Жаккара J_i

Вычисление моды выборки и максимальной клики. Имеет смысл распространить понятие моды на обработку интервальных данных, где она будет обозначать интервал тех значений, которые наиболее часты, т. е. встречаются в интервалах обрабатываемых данных наиболее часто. Фактически, это означает, что точки из моды интервальной выборки накрываются наибольшим числом интервалов этой выборки. Ясно, что по самому своему определению понятие моды имеет наибольшее значение (и наибольший смысл) лишь для накрывающих выборок. Иначе, если выборка ненакрывающая, то смысл «частоты» тех или иных значений в пределах рассматриваемых интервалов этой выборки в значительной мере теряется, хотя и не обесценивается.

Мода является пересечением интервалов максимальной совместной подвыборки, и если максимальных подвыборок имеется более одной, то мода будет объединением их пересечений, т. е. мультиинтервалом. Простой алгоритм вычисления моды интервальной выборки можно найти в [6]. Псевдокод специализированного алгоритма для нахождения моды выборки интервальных измерений и её частоты приведён в Табл. 2.1.

Ключевым в алгоритме Табл. 2.1 является формирование множества *элементарных подинтервалов измерений* из упорядоченных вершин (концов интервалов) $\underline{x}_1, \bar{x}_1, \underline{x}_2, \bar{x}_2, \dots, \underline{x}_n, \bar{x}_n$ исходной выборки \mathbf{X} .

Отметим также, что мода интервальной выборки — это интервал или мультиинтервал, который не обязан совпадать с каким-либо из интервалов обрабатываемой выборки.

Пример 2.1.1. Рассмотрим пример вычисления моды интервальной выборки из 4 элементов

$$\mathbf{X} = \{ [1, 4], [5, 9], [1.5, 4.5], [6, 9] \}. \quad (2.4)$$

В соответствии с алгоритмом Табл. 2.1, проверим совместность \mathbf{X} . Пересечение элементов выборки пусто

$$\mathbf{I} = \bigcap_{i=1}^n \mathbf{x}_i = \emptyset.$$

Таким образом, необходимо выполнить шаги алгоритма Табл. 2.1 после ключевого слова ELSE.

Сформируем массив интервалов \mathbf{z} из концов интервалов \mathbf{X}

$$\mathbf{z} = \{ [1, 1.5], [1.5, 4], [4, 4.5], [4.5, 5], [5, 6], [6, 9], [9, 9] \}. \quad (2.5)$$

Для каждого интервала \mathbf{z}_i подсчитываем число μ_i интервалов из выборки \mathbf{X} , включающих \mathbf{z}_i , получаем массив значений μ_i в виде

$$\{1, 2, 1, 0, 1, 2, 2\}. \quad (2.6)$$

Максимальные μ_i , равные $\max \mu = 2$, достигаются для индексного множества

$$K = \{2, 6, 7\},$$

Как итог, мода является мультиинтервалом

$$\text{mode } \mathbf{X} = \bigcup_{k \in K} \mathbf{z}_k = [1.5, 4] \cup [6, 9] \cup [9, 9] = [1.5, 4] \cup [6, 9]. \quad (2.7)$$

■

Перейдём к практическому примеру.

Проведём вычисление моды выборки \mathbf{X}_1 по алгоритму Табл. 2.1.

Таблица 2.1: Алгоритм для нахождения моды
интервальной выборки

Вход

Интервальная выборка $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ длины n .

Выход

Мода $\text{mode } \mathbf{X}$ выборки \mathbf{X} и её частота μ .

Алгоритм

$\mathbf{I} \leftarrow \bigcap_{i=1}^n \mathbf{x}_i$;

IF $\mathbf{I} \neq \emptyset$ THEN

$\text{mode } \mathbf{X} \leftarrow \mathbf{I}$;

$\mu \leftarrow n$

ELSE

 помещаем все концы $\underline{\mathbf{x}}_1, \overline{\mathbf{x}}_1, \underline{\mathbf{x}}_2, \overline{\mathbf{x}}_2, \dots, \underline{\mathbf{x}}_n, \overline{\mathbf{x}}_n$

 интервалов рассматриваемой выборки \mathbf{X} в один

 массив $Y = (y_1, y_2, \dots, y_{2n})$;

 упорядочиваем элементы в Y по возрастанию значений;

 порождаем интервалы $\mathbf{z}_i = [y_i, y_{i+1}]$, $i = 1, 2, \dots, 2n - 1$

 (назовём их *элементарными подинтервалами измерений*);

 для каждого \mathbf{z}_i подсчитываем число μ_i интервалов

 из выборки \mathbf{X} , включающих интервал \mathbf{z}_i ;

 вычисляем $\mu \leftarrow \max_{1 \leq i \leq 2n-1} \mu_i$;

 выбираем номера k интервалов \mathbf{z}_k , для которых μ_k

 равно максимальному, т. е. $\mu_k = \mu$, и формируем

 из таких k множество $K = \{k\} \subseteq \{1, 2, \dots, 2n - 1\}$;

$\text{mode } \mathbf{X} \leftarrow \bigcup_{k \in K} \mathbf{z}_k$

END IF

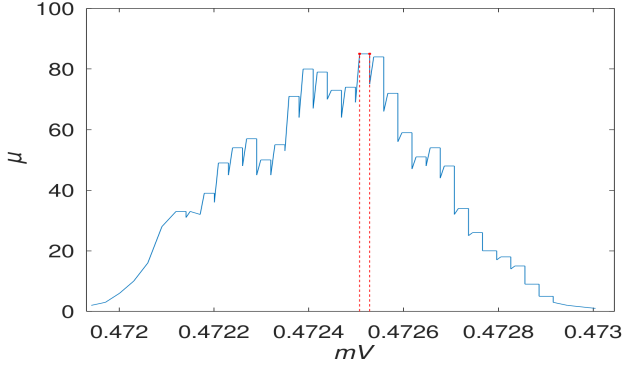


Рис. 2.3: График частот при вычислении моды выборки \mathbf{X}_1

По результатам вычисления моды выборки находим размер максимальной клики $\max \mu_j$:

$$\max \mu_j(\mathbf{X}_1) = 85. \quad (2.8)$$

Индексы таких элементов образуют множество K , а из них образуется мода

$$K = \{79, 80, \dots, 163\}, \quad (2.9)$$

$$\text{mode}(\mathbf{X}_1) = \bigcup_{k \in K} z_k = [0.47251, 0.47253]. \quad (2.10)$$

На рис. 2.4 показаны элементы выборки \mathbf{X}_1 , в которые входит мода.

Варьирование неопределённости измерений. Один из приёмов выявления достижения совместности выборки интервальных наблюдений основан на представлении о причине несовместности как недооценённой величины неопределённости [27, 28]. Закономерным шагом в этом случае становится поиск некоторой минимальной коррекции величин неопределённости интервальных наблюдений, необходимой для обеспечения совместности задачи построения зависимости. Если величину коррекции каждого интервального наблюдения $\mathbf{y}_i = [\underline{y}_i - \epsilon_i, \underline{y}_i + \epsilon_i] = (\underline{y}_i, \epsilon_i)$ выборки \mathcal{S}_n выражать коэффициентом его уширения $w_i \geq 1$, а общее изменение выборки характеризовать суммой этих коэффициентов, то минимальная коррекция выборки в виде

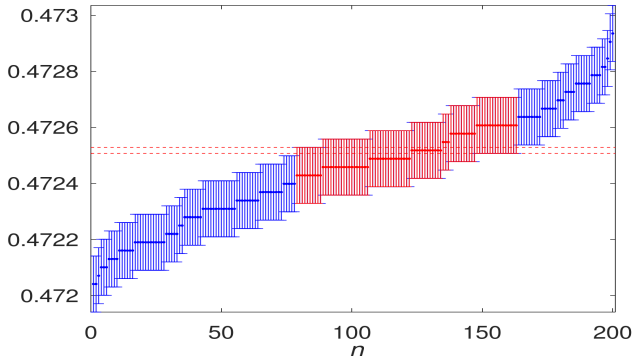


Рис. 2.4: Элементы выборки X_1 , в которые входит мода (2.3).

вектора коэффициентов $w^* = (w_1^*, \dots, w_n^*)$, необходимая для совместности задачи построения зависимости $y = f(x, \beta)$ может быть найдена решением задачи условной оптимизации

$$\text{найти} \quad \min_{w, \beta} \sum_{i=1}^n w_i \quad (2.11)$$

при ограничениях

$$\begin{cases} \hat{y}_i - w_i \epsilon_i \leq f(x_i, \beta) \leq \hat{y}_i + w_i \epsilon_i, \\ w_i \geq 1, \end{cases} \quad i = 1, \dots, n. \quad (2.12)$$

Результирующие значения коэффициентов w_i^* , строго превосходящие единицу, указывают на наблюдения, которые требуют уширения интервалов неопределённости для обеспечения совместности данных и модели. Именно такие наблюдения заслуживают внимания при анализе данных на выбросы. Значительное количество подобных наблюдений может говорить либо о неверно выбранной структуре зависимости, либо о том, что величины неопределённости измерений занижены во многих наблюдениях (например, в результате неверной оценки точности измерительного прибора).

Следует отметить значительную гибкость языка неравенств. Он даёт возможность переформулировать и расширять систему ограничений (2.12) для учёта специфики данных и задачи при поиске допустимой коррекции данных, приводящей к разрешению исходной несовместно-

сти. Например, если имеются основания считать, что величина неопределённости некоторой группы наблюдений одинакова и при коррекции должна увеличиваться синхронно, то система ограничений (2.12) может быть пополнена равенствами вида

$$w_{i_1} = w_{i_2} = \dots = w_{i_K},$$

где i_1, \dots, i_K — номера наблюдений группы. В случае, когда в надёжности каких-либо наблюдений исследователь уверен полностью, при решении задачи (2.11)–(2.12) соответствующие им величины w_i можно положить равными единице, т.е. запретить варьировать их неопределённость.

Задача поиска коэффициентов масштабирования величины неопределённости (2.11)–(2.12) сформулирована для распространённого случая уравновешенных интервалов погрешности и подразумевает синхронную подвижность верхней и нижней границ интервалов неопределённости измерений \mathbf{y}_i при сохранении базовых значений интервалов \dot{y}_i неподвижными. При необходимости постановка задачи легко обобщается. Например, если интервалы наблюдений не уравновешены относительно базовых значений (то есть $\mathbf{y}_i = [\dot{y}_i - \epsilon_i^-, \dot{y}_i + \epsilon_i^+]$ и $\epsilon^- \neq \epsilon^+$), то границы интервальных измерений можно варьировать независимо, масштабируя величины неопределённости ϵ_i^- и ϵ_i^+ с помощью отдельных коэффициентов w_i^- и w_i^+ :

$$\text{найти} \quad \min_{w^-, w^+, \beta} \quad \sum_{i=1}^n (w_i^- + w_i^+) \quad (2.13)$$

при ограничениях

$$\left\{ \begin{array}{l} \dot{y}_i - w_i^- \epsilon_i^- \leq f(x_i, \beta) \leq \dot{y}_i + w_i^+ \epsilon_i^+, \\ w_i^- \geq 1, \\ w_i^+ \geq 1, \end{array} \quad i = 1, \dots, n. \right. \quad (2.14)$$

Для линейной по параметрам β зависимости $y = f(x, \beta)$ задача (2.11)–(2.12) представляет собою задачу линейного программирования, для решения которой широко доступны хорошие и апробированные программы в составе библиотек на различных языках программирования, в виде стандартных процедур систем компьютерной математики, а также в виде интерактивных подсистем электронных таблиц.

Оптимизация по Оскорбину. Перейдём к практическому примеру выборки \mathbf{X}_1 . Поставим задачу линейного программирования (2.13) — (2.14) в простейшем виде

$$\text{найти} \quad \min_{w, \beta} w \quad (2.15)$$

при ограничениях

$$\begin{cases} \text{mid } \mathbf{x}_i - w \epsilon_i \leq \beta \leq \text{mid } \mathbf{x}_i + w \epsilon_i, \\ w \geq 1, \end{cases} \quad i = 1, \dots, n. \quad (2.16)$$

Проведём вычисление моды выборки \mathbf{X}_1 с использованием программ С.И.Жилина [8]. Синтаксис вызова программы

$$[\text{oskorbin_center_k}, w] = \text{estimate_uncertainty_center}(\mathbf{X}_1) \quad (2.17)$$

Вычисления дают следующие результаты

$$\text{oskorbin_center_k} = 0.4725, \quad (2.18)$$

$$w = 4.4705. \quad (2.19)$$

Оценка постоянной (2.18) очень близка с вычисленной ранее модой. Величина однородного расширения интервалов (2.19) достаточно велика, что соответствует весьма большой степени несовместности выборки \mathbf{X}_1 .

На рис. 2.5 приведена диаграмма рассеяния выборки \mathbf{X}_1 с увеличенным в w раз интервалом неопределённости.

Красной пунктирной линией показана оценка постоянной (2.18).

Индекс Жаккара Для описания выборок, помимо оценок их размеров, желательно иметь дополнительную информацию о мере сходства элементов выборки. В различных областях анализа данных, биологии, информатике, в науках о Земле часто используют различные меры сходства множеств (см. [22]). Меру сходства между объектами A и B можно определить как двухместную вещественнозначную функцию $S(A, B)$, которая обладает следующими свойствами:

- ограниченность: $0 \leq S(A, B) \leq 1$;
- симметричность: $S(A, B) = S(B, A)$;

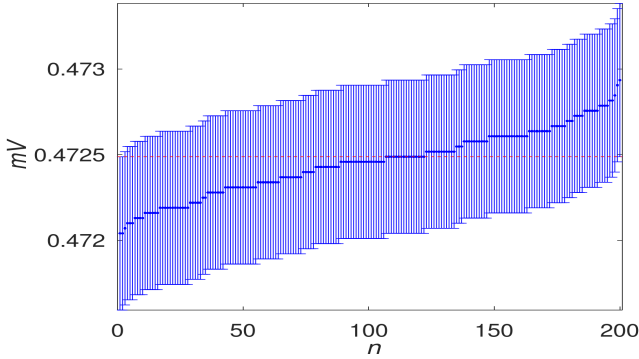


Рис. 2.5: Диаграмма рассеяния выборки X_1 с увеличенным в w раз интервалом неопределённости.

- неразличимость совпадающих элементов: $S(A, B) = 1 \Leftrightarrow A = B$;
- монотонность: $A \subseteq B \subseteq C \Rightarrow S(A, B) \geq S(A, C)$.

Таким образом, значение 1 этой меры соответствует совпадению множеств A и B , а значение 0 означает их полное несходство. Отметим, что существуют и иные системы аксиом сходства. В компьютерных приложениях (обработка изображений, машинное обучение) меру сходства множеств обозначают как *IoU* (*Intersection over Union*). В математике и её приложениях для подобных конструкций часто используется термин *индекс Жаккара*, по имени исследователя, впервые предложившего эту меру.

В процессе развития интервального анализа были введены различные определения и конструкции оценки меры совместности интервальных объектов. Вместе с тем в практике обработки данных часто необходимо оперировать относительными величинами. В частности, это нужно в связи с необходимостью сопоставления допусков и размеров деталей, погрешности измерителей и значений измеряемых величин и т. п. [?].

Представим обобщение меры Жаккара на выборки интервалов [5]. В качестве числовой характеристики степени совпадения двух интер-

валов \mathbf{x} , \mathbf{y} рассмотрим величину

$$Ji(\mathbf{x}, \mathbf{y}) := \frac{\text{wid}(\mathbf{x} \wedge \mathbf{y})}{\text{wid}(\mathbf{x} \vee \mathbf{y})}. \quad (2.20)$$

В выражении (2.20) используется ширина интервала (см. стр. ??), а вместо операций пересечения и объединения множеств — операции взятия точной нижней грани “ \wedge ” (инфимума, см. (??)) и точной верхней грани “ \vee ” (супремума, см. (??)) относительно включения для двух величин в полной интервальной арифметике Каухера. В обозначении $Ji(\mathbf{x}, \mathbf{y})$ буква J указывает на фамилию «Jassard», а i — на интервальность его применения.¹ В общем случае инфимум по включению в числителе выражения (2.20) может быть неправильным интервалом, и его ширина тогда отрицательна.

Рассмотренная мера обобщает обычное понятие меры совместности на различные типы взаимной совместности интервалов. Если пересечение интервалов \mathbf{x}, \mathbf{y} пусто, т. е. $\mathbf{x} \cap \mathbf{y} = \emptyset$, то $\mathbf{x} \wedge \mathbf{y}$ — неправильный интервал и числитель формулы (2.20) имеет отрицательное значение. В предельном случае несовпадающих вещественных вырожденных интервалов $\mathbf{x} = x$ и $\mathbf{y} = y$, $x \neq y$, имеем

$$Ji(x, y) = -1.$$

В целом получаем

$$-1 \leq Ji(\mathbf{x}, \mathbf{y}) \leq 1. \quad (2.21)$$

Таким образом, величина Ji непрерывно описывает ситуации от полной несовместности вещественных значений $x \neq y$ до полного перекрытия интервалов $\mathbf{x} = \mathbf{y}$. Следует заметить, что в отличие от случая вещественных величин, для которых индекс Жаккара может принимать только два значения, 0 и 1, формула (2.20) даёт характеризацию различных отношений сходства интервалов с помощью непрерывного ряда значений между -1 и 1 .

Мера совместности, введённая для двух интервалов в форме (2.20), допускает естественное обобщение на случай интервальной выборки $\mathbf{X} = \{\mathbf{x}_i\}$, $i = 1, 2, \dots, n$. Определим меру $Ji(\mathbf{X})$ для этой выборки как

$$Ji(\mathbf{X}) = \frac{\text{wid}(\bigwedge_i \mathbf{x}_i)}{\text{wid}(\bigvee_i \mathbf{x}_i)}. \quad (2.22)$$

¹Можно встретить также другие обозначения для этой конструкции, в частности, «JK», как это сделано в работе [5].

Видно, что выражение (2.22) переходит в случае интервальной выборки из двух элементов в выражение (2.20).

В связи несовместностью выборки будем использовать следующую меру, которая имеет место и в случае несовместных выборок.

$$\rho(\text{mode}(\mathbf{X})) = \frac{\text{wid}(\text{mode}\mathbf{X})}{\text{wid}(\bigvee_i \mathbf{x}_i)}. \quad (2.23)$$

Назовём конструкцию (2.23) *относительная ширина моды*. В отличие от минимума по включению, мода выборки всегда является правильным интервалом. В целом получаем

$$0 \leq \rho(\text{mode}(\mathbf{X})) \leq 1. \quad (2.24)$$

Вычисление меры совместности. Перейдём к практическому примеру выборки \mathbf{X}_1 . Проведём вычисление с использованием программ на ресурсе [23].

$$Ji(\mathbf{X}) = \frac{\text{wid}(\bigwedge_i \mathbf{x}_i)}{\text{wid}(\bigvee_i \mathbf{x}_i)} = -0.6344. \quad (2.25)$$

Отрицательность меры (2.25) соответствует несовместности выборки \mathbf{X}_1 , а её модуль — высокой степени этой несовместности.

Относительная ширина моды (2.23) равна

$$\rho(\text{mode}(\mathbf{X})) = \frac{\text{wid}(\text{mode}\mathbf{X})}{\text{wid}(\bigvee_i \mathbf{x}_i)} = 0.039. \quad (2.26)$$

Величина (2.26) составляет менее 4% внешней оценки выборки \mathbf{X}_1 .

2.2 Задача восстановления линейной зависимости

До настоящего момента нас интересовали результаты серии измерений некоторой единственной постоянной величины, о которой мы «накапливали» информацию с целью найти её точечную и интервальную оценки в результате обработки элементов выборки.

В данном разделе нам предстоит рассмотреть задачу восстановления функциональной зависимости по интервальным результатам измерений. Будем считать, что вид функциональной зависимости нам заранее

известен — она является линейной. При этом большая часть рассуждений, которые будут проделаны, может быть успешно перенесена на случай нелинейной зависимости.

2.2.1 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 :

$$y = f(\beta, x), \quad (2.27)$$

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 (2.27), которые соответствуют конкретной функции f из параметрического семейства.

Если функция f является линейной, то можно записать

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m. \quad (2.28)$$

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

$$\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \dots, \mathbf{x}_m^{(k)} \text{ и } \mathbf{y}^{(k)}.$$

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]:

$$(\mathbf{x}_{k1}, \mathbf{x}_{k2}, \dots, \mathbf{x}_{km}, \mathbf{y}_k) \subset \mathbb{R}^{m+1}, \quad k = 1, 2, \dots, n. \quad (2.29)$$

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

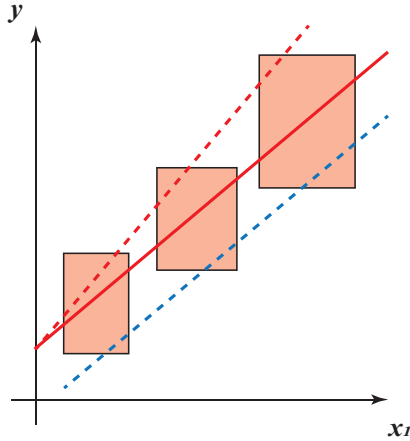


Рис. 2.6: Различные способы пересечения графика линейной зависимости $y = \beta_0 + \beta_1 x_1$ с брусами неопределенности измерений

прохождения можно описать с помощью формального языка логического исчисления предикатов:

$$(\forall x_{k1} \in \mathbf{x}_{k1}) \cdots (\forall x_{km} \in \mathbf{x}_{km}) (\exists y_k \in \mathbf{y}_k) \quad (2.31)$$

$$\beta_0 + \beta_1 x_{k1} + \beta_2 x_{k2} + \dots + \beta_m x_{km} = y_k, \quad k = 1, \dots, n.$$

Определение 2.2.2 *Функциональную зависимость назовем сильно совместной с интервальными данными, если ее график проходит через каждый брус неопределенности измерений для любого значения аргументов из интервалов неопределенности входных переменных [?].*

Значит, линейная зависимость, график которой показан на Рис. 2.6 сплошной линией красного цвета, является сильно совместной с результатами измерений. Множеством решений ИСЛАУ (2.30), выбор каждой точки (β_0, β_1) из которого гарантирует сильную совместность линейной зависимости с интервальными данными, является *допусковое множество* АЕ-решений (??).

Что означает сильная совместность функциональной зависимости с интервальными результатами измерений? Этот результат говорит о

том, что независимо от того, какими именно являются значения входных переменных, находящиеся при этом в пределах интервалов неопределенности, выходная величина будет неизменно принадлежать интервалу своей неопределенности.

Теперь перейдем к рассмотрению графиков линейных зависимостей, показанных на Рис. 2.6 пунктирными линиями красного и голубого цветов. Характер их пересечения с брусами неопределенности измерений не может быть описан так, как определено выражением (2.31). Обе эти линейные зависимости являются слабо совместными с интервальными результатами измерений.

Определение 2.2.3 *Функциональную зависимость назовем слабо совместной с интервальными данными, если ее график проходит через каждый брус неопределенности хотя бы для одного значения аргумента [1].*

В случае слабой совместности линейной зависимости с интервальными результатами измерений характер прохождения ее графика через брусы неопределенности можно описать следующим образом:

$$\begin{aligned} & (\exists x_{k1} \in \mathbf{x}_{k1}) \cdots (\exists x_{km} \in \mathbf{x}_{km}) (\exists y_k \in \mathbf{y}_k) \\ & \beta_0 + \beta_1 x_{k1} + \beta_2 x_{k2} + \dots + \beta_m x_{km} = y_k, \quad k = 1, \dots, n. \end{aligned} \quad (2.32)$$

Следовательно, выбор любой точки (β_0, β_1) из объединенного множества АЕ-решений (??) ИСЛАУ (2.30) обеспечивает слабую совместность линейной зависимости с брусами неопределенности измерений.

2.2.3 Информационное множество и коридор совместных зависимостей

Информационным множеством задачи восстановления функциональной зависимости (2.27) является множество всех значений параметров функции β , получаемое в результате обработки интервальных значений входных переменных x и выходной переменной y . Если говорить конкретно о задаче восстановления линейной зависимости (2.28), то ее информационное множество представлено вектором оценки параметров β — определенным множеством АЕ-решений ИСЛАУ (2.30). Такой подход к определению параметров восстанавливаемой зависимости принципиально отличен от того, как производится их поиск в

случае точечных результатов измерений: вместо решения задачи минимизации отклонений значений функции от данных, полученных эмпирическим путем, мы составляем систему интервальных уравнений или неравенств и находим ее решение.

Остановимся чуть подробнее на случае восстановления функциональной зависимости по точечным данным. На самом деле для одного и того же набора результатов измерений мы можем получить целое семейство зависимостей, которые «более или менее» хорошо приближают экспериментальные данные и могут быть использованы для предсказания своих значений, соответствующих интересующим нас аргументам из областей определения. Причиной этому может быть различный выбор метода регрессионного анализа (например, метод наименьших квадратов, метод наименьших модулей, чебышевская аппроксимация и проч.) либо присутствие неточностей и неопределенностей в самих обрабатываемых данных. Как результат, значения всех восстановленных функциональных зависимостей в каждой точке их определения будут задавать *коридор совместных зависимостей*. На Рис. 2.7 показаны графики трех нелинейных функций (зеленая, фиолетовая и оранжевая линии), восстановленных по одному набору точечных результатов измерений. Коридор совместных зависимостей показан заливкой розовым цветом, а его сечение в точке x^* выделено черной линией.

Вернемся к рассмотрению интервальных результатов измерений, по которым проводится восстановление функциональной зависимости. Если информационное множество задачи непусто и «не вырождено» (не точка или отрезок), то мы так же, как и в рассмотренном выше случае, получим семейство зависимостей, совместных с результатами измерений и образующих коридор.

Дадим строгое определение коридора совместных зависимостей, которое используется в интервальном анализе данных. Для этого необходимо ввести понятие *многозначного отображения*.

Определение 2.2.4 Для произвольных множеств X и Y многозначным отображением F из X в Y называется соответствие (правило), сопоставляющее каждому аргументу $x \in X$ непустое подмножество $F(x) \subset Y$, называемое значением или образом x [?].

В отличие от традиционного отображения многозначное отображение сопоставляет каждому элементу области определения X не один элемент, а целое множество из Y .

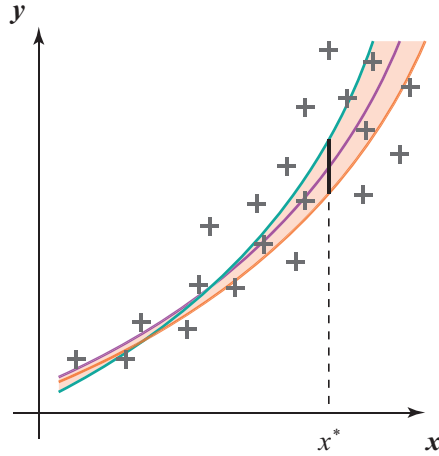


Рис. 2.7: Коридор совместных зависимостей и его сечение при некотором значении входной переменной x^*

Определение 2.2.5 Пусть в задаче восстановления зависимостей информационное множество Ω параметров зависимостей $y = f(x, \beta)$, совместных с данными, является непустым.

Коридором совместных зависимостей рассматриваемой задачи называется многозначное отображение Υ , сопоставляющее каждому значению аргумента x множество

$$\Upsilon(x) = \bigcup_{\beta \in \Omega} f(x, \beta).$$

Значение $\Upsilon(\tilde{x})$ коридора совместных зависимостей при определенном аргументе \tilde{x} (сечение коридора) — это множество

$$\bigcup_{\beta \in \Omega} f(\tilde{x}, \beta),$$

образованное всевозможными значениями, которые принимают на этом аргументе функциональные зависимости, совместные с интервальными данными измерений [?].

Внешней оценкой сечения коридора совместных зависимостей является интервал

$$\left[\min_{\beta \in \Omega} f(\tilde{x}, \beta), \max_{\beta \in \Omega} f(\tilde{x}, \beta) \right].$$

Пример 2.2.1.

Рассмотрим набор интервальных результатов измерений, представленный на Рис. 2.8. Мы должны восстановить по нему линейную зависимость $y = \beta_0 + \beta_1 x_1$. **нижний индекс 1 не нужен**

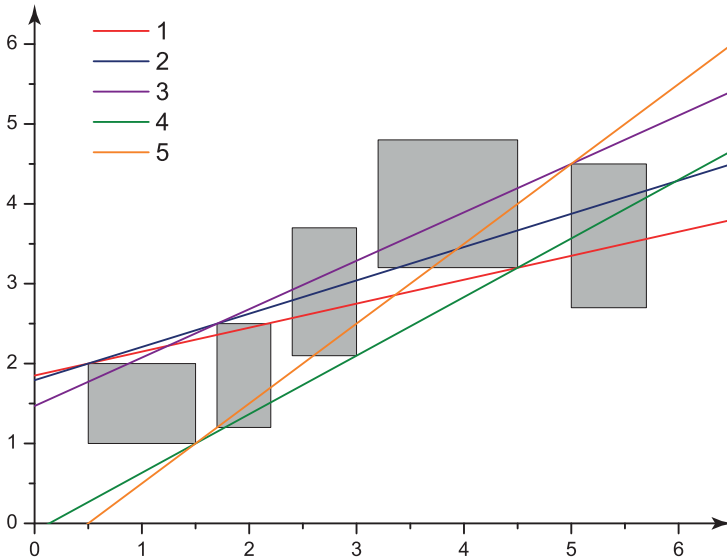


Рис. 2.8: Интервальные результаты измерений, соответствующие ИСЛАУ (2.33), а также восстановленные линейные зависимости, задающие коридор совместных зависимостей $y = 1.85 + 0.30 x_1$ (1), $y = 1.79 + 0.42 x_1$ (2), $y = 1.47 + 0.61 x_1$ (3), $y = -0.10 + 0.73 x_1$ (4), $y = -0.50 + 1.00 x_1$ (5) **подписи осей**

Для этого составим интервальную систему линейных алгебраиче-

ских уравнений:

$$\begin{pmatrix} 1 & [0.5, 1.5] \\ 1 & [1.7, 2.2] \\ 1 & [2.4, 3.0] \\ 1 & [3.2, 4.5] \\ 1 & [5.0, 5.7] \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} [1.0, 2.0] \\ [1.2, 2.5] \\ [2.1, 3.7] \\ [3.2, 4.8] \\ [2.7, 4.5] \end{pmatrix} \quad (2.33)$$

Допусковое множество АЕ-решений ИСЛАУ (2.33) пусто, а значит, по имеющимся интервальным данным невозможно восстановить линейную зависимость, которая была бы с ними сильно совместна. Как было установлено в подразделе ??, объединенное множество решений шире, чем допусковое множество решений. Непустое объединенное множество АЕ-решений ИСЛАУ (2.33) показано на Рис. 2.9.

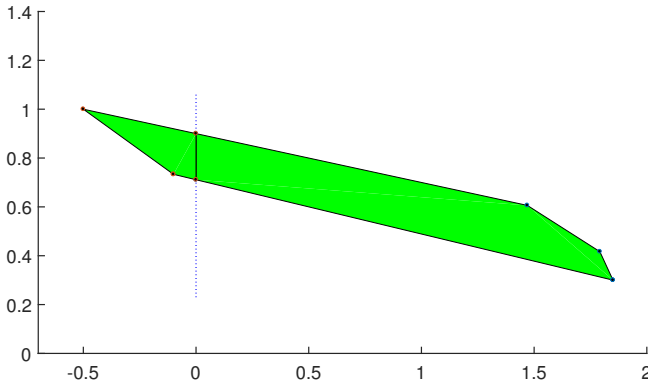


Рис. 2.9: Объединенное множество решений ИСЛАУ (2.33), построенное при использовании функции `EqnWeak2D` пакета `IntLinInc2D` для MATLAB [25] подписи осей

Замкнутое объединенное множество решений ИСЛАУ (2.33) имеет пять вершин

$$(1.85, 0.30), (1.79, 0.42), (1.47, 0.61), (-0.10, 0.73), (-0.50, 1.00),$$

причем каждая из них представляет собой такие параметры (β_0, β_1) , которые соответствуют восстановленным линейным зависимостям, задающим коридор совместных зависимостей.

Все линейные зависимости (1)–(5), показанные на Рис. 2.8, являются слабо совместными с интервальными результатами измерений. Любая из них проходит через два угла, принадлежащие двум различным интервалам, и, как можно заметить, границы коридора совместных зависимостей являются кусочно-линейными. ■

Пример 2.2.2. Предположим, что нам необходимо восстановить линейную зависимость $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. Имеющиеся интервальные результаты измерений позволяют составить следующую ИСЛАУ:

$$\begin{pmatrix} 1 & [-2.0, 0.5] & [1.0, 3.0] \\ 1 & [8.0, 9.0] & [-2.0, -0.1] \\ 1 & [-6.0, -4.0] & [10.0, 12.0] \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} [-5.0, -2.7] \\ [-0.6, 4.1] \\ [5.0, 10.0] \end{pmatrix} \quad (2.34)$$

Объединенное множество решений ИСЛАУ (2.34), имеющее замкнутую форму, показано на Рис. 2.10. Выбор каждой точки из этого множества будет гарантировать слабую совместность восстановленной линейной зависимости с интервальными данными. Допусковое множество АЕ-решений ИСЛАУ (2.34) пусто.

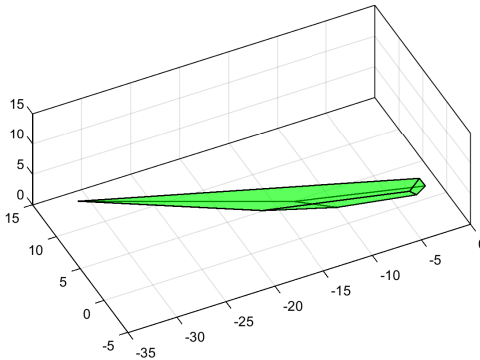


Рис. 2.10: Объединенное множество решений ИСЛАУ (2.34), построенное при использовании функции `EqnWeak3D` пакета `IntLinInc3D` для MATLAB [26]

Перечислим восемь вершин многогранника, которым является объ-

единенное множество решений ИСЛАУ (2.34):

$$\begin{aligned} &(-34.22, 6.94, 8.58), (-17.88, 1.96, 3.97), (-14.20, 2.93, 2.58), \\ &(-10.59, -1.13, 1.68), (-3.57, 0.35, 1.57), (-2.92, 0.27, 0.75), \\ &(-2.13, 1.27, 1.98), (-1.69, 0.94, 0.87). \end{aligned}$$

■

Данные выборки. Имеется выборка данных \mathbf{X}_1 с интервальной неопределённостью. Число отсчётов в выборке равно 200.

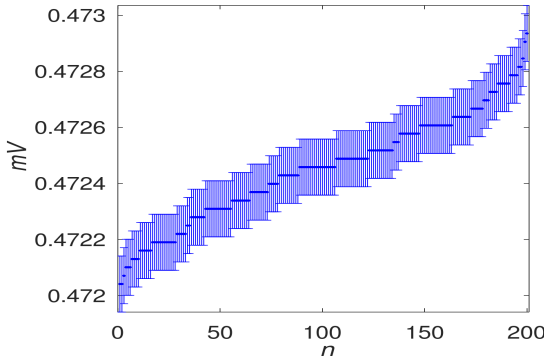


Рис. 2.11: Диаграмма рассеяния выборки \mathbf{X}_1 с уравновешенным интервалом погрешности (2.1).

На Рис. 2.11 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 $w_i \geq 1$, а общее изменение выборки характеризовать суммой этих коэффициентов, то минимальная коррекция выборки в виде вектора коэффициентов $w = (w_1, \dots, w_n)$, необходимая для совместности задачи построения зависимости $x = \beta_0 + \beta_1 \cdot i$ может

быть найдена решением задачи условной оптимизации

$$\text{найти} \quad \min_{w, \beta} \sum_{i=1}^n w_i \quad (2.35)$$

при ограничениях

$$\begin{cases} \text{mid } \mathbf{x}_i - w_i \epsilon_i \leq \beta_0 + \beta_1 \cdot i \leq \text{mid } \mathbf{x}_i + w_i \epsilon_i, \\ w_i \geq 1, \end{cases} \quad i = 1, \dots, n. \quad (2.36)$$

Результирующие значения коэффициентов w_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

$$[\text{tau}, \mathbf{w}, \text{yint}] = \text{DataLinearModel}(\text{input1}, \text{epsilon0}) \quad (2.37)$$

В (2.37) входами программы служат значения $\text{mid } \mathbf{X}_1$ и величин неопределённости ϵ , а выходами tau — значения параметров регрессии β_0, β_1 , и \mathbf{w} — вектор весов расширения интервалов.

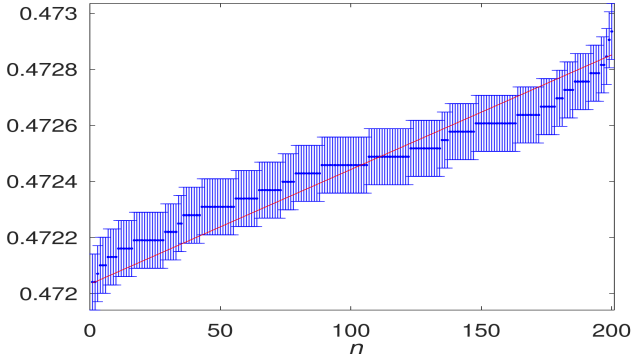


Рис. 2.12: Диаграмма рассеяния выборки \mathbf{X}_1 и регрессионная прямая по модели (2.35) и (2.36).

На Рис. 2.12 красным цветом приведена регрессионная прямая.

Вычисления с использованием программы (2.37) дают следующие результаты для регрессионных коэффициентов

$$\beta_0 = \text{tau}(1) = 4.7203e - 01, \quad (2.38)$$

$$\beta_1 = \text{tau}(2) = 4.0915e - 06. \quad (2.39)$$

Все компоненты вектора w оказались равны 1, то есть, расширения интервалов измерений не понадобилось. Таким, образом, величина (2.35) равна числу элементов выборки.

$$\min_{w, \beta} \sum_{i=1}^n w_i = 200. \quad (2.40)$$

Недостатком полученного решения с единичными значениями w_i является неучёт расстояний точек регрессионной зависимости до данных интервальной выборки. Таким образом, прямая с параметрами (2.38) и (2.39) «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:

$$\text{найти} \quad \min_{w, \beta} \sum_{i=1}^n w_i \quad (2.41)$$

при ограничениях

$$\begin{cases} \text{mid } \mathbf{x}_i - w_i \epsilon_i \leq \beta_0 + \beta_1 \cdot i \leq \text{mid } \mathbf{x}_i + w_i \epsilon_i, \\ w_i \geq 0, \end{cases} \quad i = 1, \dots, n. \quad (2.42)$$

Отличие постановки от (2.35) и (2.36) состоит в том, что интервалы измерений могут как расширяться в случае $w_i \geq 1$, так и сужаться при $0 \leq w_i < 1$.

Вчисление параметров линейной регрессии по данным интервальной выборки \mathbf{X}_1 производится как и в случае (2.37) с использованием программ С.И. Жилина [8] и оформленных применительно к задаче на [23]. Синтаксис вызова программы

$$[\text{tau}, w, \text{yint}] = \text{DataLinearModelZ}(\text{input1}, \text{epsilon0}) \quad (2.43)$$

Входы и выходы функции `DataLinearModelZ` такие же, как и для `DataLinearModelZ` (2.37).

На Рис. 2.13 красным цветом приведена регрессионная прямая.

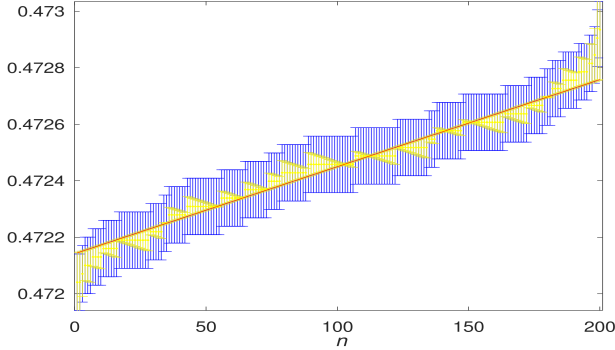


Рис. 2.13: Диаграмма рассеяния выборки \mathbf{X}_1 и регрессионная прямая по модели (2.41) и (2.42).

Жёлтым цветом на Рис. 2.13 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 (2.35) уменьшилась более, чем в 4 раза.

$$\min_{w, \beta} \sum_{i=1}^n w_i = 45.7 < 200. \quad (2.44)$$

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

На Рис. 2.14 приведены графики векторов w_1 и w_0 , полученных при использовании двух рассмотренных подходов.

В конкретном случае график вектора w_0 для постановки задачи оптимизации (2.41) и (2.42) содержит большое количество информации.

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

$$\Psi = \arg_i w_i \geq \alpha. \quad (2.45)$$

Для конкретного примера имеем две области Ψ в начале и конце области данных.

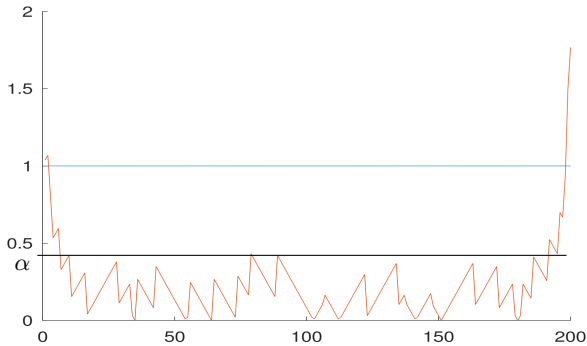


Рис. 2.14: Векторы w_1 и w_0 .

Для объективного использования этого приёма параметр α можно брать, например, из анализа гистограммы распределения вектора вектора w .

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/discrepancy.html>

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.

Построение прогноза внутри и вне области данных.

Коррекция модели данных. Кусочно-линейная регрессионная зависимость.

Литература

- [1] А.Н. Баженов, С.И. Жилин, С.И. Кумков, С.П. Шарый. «Обработка и анализ данных с интервальной неопределённостью». 2023.
- [2] Баженов А.Н. Интервальный анализ. Основы теории и учебные примеры: учебное пособие. — СПб. 2020 – 78 с. <https://elib.spbstu.ru/dl/2/s20-76.pdf/info>
- [3] Баженов А.Н.. Обобщение мер совместности для анализа данных с интервальной неопределённостью. – СПб., 2022. – 80 с. <https://elib.spbstu.ru/dl/5/tr/2022/tr22-142.pdf/info>
- [4] Баженов А.Н.. Примеры обработки измерений постоянной величины в анализе данных с интервальной неопределённостью. – СПб., 2023. <https://elib.spbstu.ru/dl/5/tr/2023/tr23-29.pdf/info>
- [5] Баженов А.Н., Тельнова А.Ю. Обобщение коэффициента Жаккара для анализа данных с интервальной неопределённостью // Измерительная техника. – 2022. № 12, С. 12-15.
- [6] Баженов А.Н. Введение в анализ данных с интервальной неопределённостью : учеб. пособие / А.Н. Баженов - СПб. : ПОЛИТЕХ-ПРЕСС, 2022. - 92 с. Табл.3. Ил.36. Библиогр.: 48 назв. ISBN 978-5-7422-7910-5 Санкт-Петербургский политехнический университет Петра Великого, 2022 doi:10.18720/SPBPU/2/id22-247
- [7] Оскорбин Н.М., Максимов А.В., Жилин С.И. Построение и анализ эмпирических зависимостей методом центра неопределенности // Известия Алтайского государственного университета. – 1998. – №1. – С. 35–38.
- [8] Примеры анализа интервальных данных в Octave <https://github.com/szhilin/octave-interval-examples>
- [9] Шарый С.П. Конечномерный интервальный анализ. – ФИЦ ИВТ: Новосибирск, 2022. Электронная книга, доступная на <http://www.nsc.ru/interval/Library/InteBooks/SharyBook.pdf>

- [10] ШРЕЙДЕР Ю.А. Равенство, сходство, порядок. – М.: Наука, 1971. 256 с.
- [11] KEARFOTT, R.B., NAKAO, M., NEUMAIER, A., RUMP, S., SHARY, S.P., VAN HENTENRYCK, P. Standardized notation in interval analysis // Вычислительные Технологии. – 2010. – Т. 15, №1. – С. 7–13.
- [12] ZHILIN, S.I. On fitting empirical data under interval error // Reliable Computing. – 2005. – Vol. 11. – P. 433–442. DOI: 10.1007/s11155-005-0050-3
- [13] E. GARDENES, A.TREPAT, J.M. JANER Approaches to simulation and to the linear problem in the SIGLA system. // Gardenes, E., Trepata, A., and Janer, J. M.: Approaches to Simulation and to the Linear Problem in the SIGLA System, Freiburger Interval-Berichte 81 (8) (1981), pp. 1-28.
- [14] NESTEROV, V.M. Interval and Twin Arithmetics. Reliable Computing 3, 369–380 (1997). <https://doi.org/10.1023/A:1009945403631>
- [15] В.М. НЕСТЕРОВ Твинные арифметики и их применение в методах и алгоритмах двустороннего интервального оценивания. дисс. д.ф.-м.н. г.Санкт-Петербург, Санкт-Петербургский институт информатики и автоматизации Российской академии наук, 1999, с. 234.
- [16] ОСКОРБИН Н. М. Построение и анализ эмпирических зависимостей методом центра неопределенности / Н.М. Оскорбин, А.В. Максимов, С.И. Жилин. — Барнаул : Изв. Алтайского гос. ун-та, 1998. — № 1. — С. 35-38.
- [17] С.И.Жилин Интервальная арифметика Каухера <https://github.com/szhilin/kinterval>
- [18] А.Андросов Библиотека `intvalpy` <https://github.com/AndrosovAS/intvalpy>,
- [19] А. ЖАВОРОНКОВА Арифметика твинов Нестерова. <https://github.com/Zhavoronkova-Alina/twin>
- [20] Т. ЯВОРУК Вычисления с изотопами `MendeleevTwin` — <https://github.com/Tatiana655/MendeleevTwin>
- [21] Tables and charts for isotope-abundance variations and atomic weights of selected elements: 2016 (Ver. 1.1, May 2018) <https://www.sciencebase.gov/catalog/item/580e719ae4b0f497e794b7d8>
- [22] Коэффициент Жаккара. URL: https://en.wikipedia.org/wiki/Jaccard_index
- [23] М.З. ШВАРЦ <https://github.com/AlexanderBazhenov/Solar-Data>

- [24] ШАРАЯ И.А. Строение допустимого множества решений интервальной линейной системы // Вычислительные технологии – 2005. – Том 10, No. 5. – P. 103-119.
- [25] ШАРАЯ И.А. Пакет IntLinInc2D для визуализации множеств решений интервальных линейных систем с двумя неизвестными. – Программное обеспечение, доступное на <http://www.nsc.ru/interval/sharaya/>. Описание <http://www.nsc.ru/interval/Programing/MCodes/IntLinInc2D.pdf>
- [26] ШАРАЯ И.А. Пакет IntLinInc3D для визуализации множеств решений интервальных линейных систем с двумя неизвестными. – Программное обеспечение, доступное на <http://www.nsc.ru/interval/sharaya/>. Описание <http://www.nsc.ru/interval/Programing/MCodes/IntLinInc3D.pdf>
- [27] Жилин С.И. Нестатистические методы и модели построения и анализа зависимостей. – Барнаул, 2004. – Диссертация на соискание учёной степени канд. физ.-мат. наук по специальности 05.13.01 «системный анализ, управление и обработка информации». Доступна на <http://www.nsc.ru/interval/Library/App1Diss/Zhilin.pdf>
- [28] ZHILIN S.I. Simple method for outlier detection in fitting experimental data under interval error // Chemometrics and Intelligent Laboratory Systems. – 2007. – Vol. 88, No. 1. – P. 60-68.
- [29] J.-R. Yu, G.-H. Tzeng, Han-Lin Li. General piecewise necessity regression analysis based on linear programming Fuzzy Sets and Systems Volume 105, Issue 3, 1 August 1999, Pages 429-436.
- [30] J.-R. Yu, G.-H. Tzeng, Han-Lin Li. General fuzzy piecewise regression analysis with automatic change-point detection. Fuzzy Sets and Systems Volume 119 Issue 2 April 16, 2001 pp 247-257 [https://doi.org/10.1016/S0165-0114\(98\)00384-4](https://doi.org/10.1016/S0165-0114(98)00384-4)
- [31] Circular polarization of γ -quanta in the $np \rightarrow d\gamma$ reactions with polarized neutrons / A.N.Bazhenov [et al.] // Physics Letters B. 3. – September 1992. Vol. 289. – No. 1-2. – P. 17-21.