

Мода интервальной выборки и алгоритм для мультимоды

А.Н. Баженов, С.И. Жилин, С.И. Кумков, С.П. Шарый

10 октября 2022 г.

0.1 Мода интервальной выборки

В традиционной статистике важной характеристикой выборки является её *мода* — значение из выборки, которое встречается наиболее часто. Если же рассматривается случайная величина с непрерывным вероятностным распределением, то её мода — это точка (или точки), в которых плотность вероятности имеет локальный максимум. Выборки и распределения с одной модой называются, как известно, *унимодальными*, а с двумя и более модами — *мультимодальными* (бимодальными т. д.).

Мода лучше, чем среднее значение (математическое ожидание) характеризует выборки с большим разбросом значений. Кроме того, мода, как характеристика «средней величины», может применяться при обработке данных, имеющих нечисловую природу.

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 содержательный смысл лишь для накрывающих выборок. Иначе, если выборка ненакрывающая, то не имеет смысла говорить о «частоте» тех или иных значений в пределах рассматриваемых интервалов этой выборки.

Следуя работе [1], введём

Определение 0.1.1 *Модой интервальной выборки назовём совокупность интервалов пересечения наибольших совместных подвыборок рассматриваемой выборки. Наибольшая длина совместных подвыборок данной выборки называется частотой моды.*

Так как сама выборка, очевидно, является своей подвыборкой, то в случае её совместности мода совпадает с пересечением всех интервалов выборки. Если же выборка несовместна, то в ней может найтись несколько совместных подвыборок максимальной длины, и их пересечения нужно рассматривать в совокупности друг с другом. Как следствие, мода может быть мультиинтервалом (см. §??). Это совершенно аналогично ситуации с обычными неинтервальными (точечными) данными, где выборка или распределение могут иметь несколько мод.

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

 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

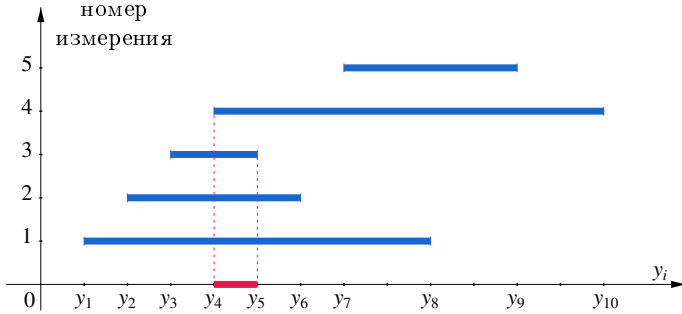


Рис. 1: Вычисление моды интервальной выборки.

Мы также будем использовать термины унимодальный, бимодальный, мультимодальный по отношению к интервальным выборкам.

Псевдокод алгоритма для нахождения моды выборки интервальных измерений и её частоты приведён в Табл. 1. Число N , не превосходящее общего числа $2n$ концов интервалов выборки, необходимо для работы алгоритма потому, что некоторые концы интервалов обрабатываемой выборки могут совпадать друг с другом. Отметим также, что мода интервальной выборки — это интервал или мультиинтервал, который не обязан совпадать с каким-либо из интервалов обрабатываемой выборки.

Пример 0.1.1. Рассмотрим пример вычисления моды интервальной выборки.

Пусть имеется интервальная выборка из 4 элементов

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

Диаграмма рассеяния выборки \mathbf{X} приведена на Рис. 2.

В соответствии с алгоритмом Табл. 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 после ключевого слова ELSE.

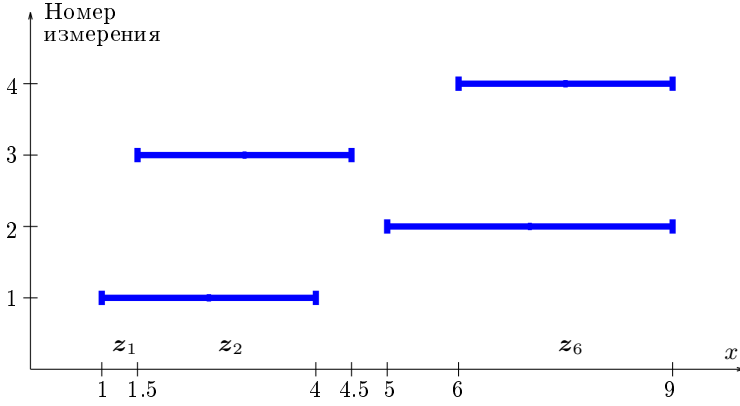


Рис. 2: Диаграмма рассеяния интервальной выборки (1) и элементы массива \mathbf{z} .

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

$$\mathbf{z} = \{ [1, 1.5], [1.5, 4], [4, 4.5], [4.5, 5], [5, 6], [6, 9], [9, 9] \}. \quad (2)$$

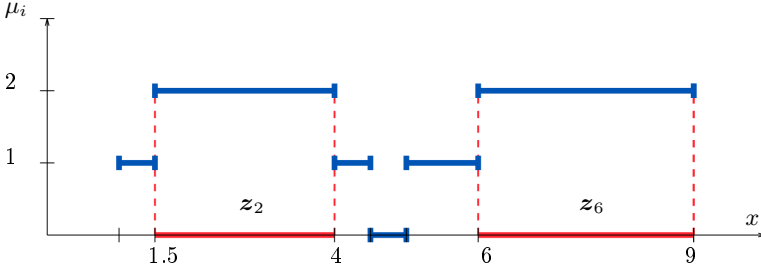


Рис. 3: Значения частот μ_i и интервальная мода $\text{mode } \mathbf{X}$ выборки (1), а также элементы \mathbf{z}_k , $k \in K$, массива \mathbf{z} .

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

$$\{1, 2, 1, 0, 1, 2, 2\}. \quad (3)$$

Максимальные μ_i , равные 2, достигаются для индексного множества

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

так что частота моды равна $\mu = 2$. Как итог, мода является мультиинтервалом

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

На Рис. 3 значения частот μ_i (3) показаны синим цветом, а интервальная мода $\text{mode } \mathbf{X}$ (4) — красным цветом. ■

0.2 Выборки унимодальные и мультимодальные

Тот факт, что выборка не является унимодальной, указывает на её неоднородность и может служить признаком сложной внутренней структуры описываемого ею явления. Получается, что из всего диапазона охватываемых выборкой значений выделяются тогда два или более изолированных друг от друга участка, одинаково доминирующих над остальными значениями по частоте. Если это доминирование велико, то исследуемая величина может, к примеру, не быть постоянной, а является «смесью» нескольких близких постоянных величин.

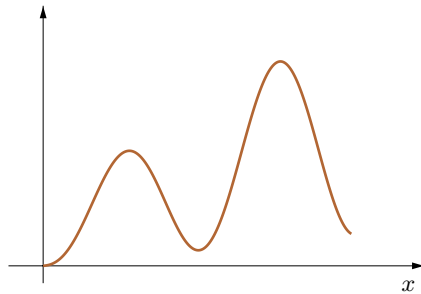


Рис. 4: Плотность вероятности бимодального распределения.

При обработке неинтервальных (точечных) данных распределения вида, показанного на Рис. 4, обоснованно считаются уже не унимодальными, так как имеют более одного явно выраженного пика, хотя и разной высоты. Но в конструкциях предшествующего параграфа возможное различие частот подвыборок, на которые распадается исходная

выборка, никак не учитывалось, и это является их существенным недостатком. Соответственно, необходимо распространить наши рассуждения на ситуации, когда интервальная выборка имеет несколько пиков частоты с разными значениями.

Интересное и содержательное обсуждение различных причин отсутствия унимодальности выборок и распределений на примере медицинских данных можно найти в работе [2].

Алгоритм вычисления мод мультимодального распределения
Представим основные идеи возможного алгоритма вычисления мод многомодального распределения.

Неформально процесс можно представить как последовательное исключение из массива частот μ_i областей пиков. Например, на Рис. 4 необходимо сначала устранить область правого пика, и перейти к левому. Границы этих областей трудно определить формально в силу их неизвестной формы и возможного наличия участков немонотонности «склонов» пиков. Мы предлагаем приём, похожий на алгоритм сегментации областей методом «водораздела» (watershed) в компьютерной графике. Отличием от поиска водоразделов является использование максимумов, а не минимумов функции распределения.

Предполагается, что все моды достаточно «компактны» и не содержат «тонкой структуры». Имеется в виду, что если текущая мода является мультиинтервалом, то его интервальная оболочка не имеет непустых пересечений с областями других мод. Такое непустое пересечение может иметь место в спектральном анализе атомов и особенно молекул, где мультиплетные распределения случаются часто [3].

Основой алгоритма является «Алгоритм для нахождения моды интервальной выборки», представленный в Табл. 1 §0.1. Этот алгоритм используется для вычисления первой («старшей») моды $\text{mode } \mathbf{X} = \bigcup_{k \in K} z_k$. Здесь множество K состоит из номеров k интервалов z_k , входящих в первую моду. Также на этом этапе вычисляется массив частот μ .

Результатом работы алгоритма является вычисление массива мод $\mathbf{MX} = \{\text{mode } \mathbf{X}^{(k)}\}_{k=1}^n$ выборки \mathbf{X} и соответствующего массива частот $M\mu(\mathbf{MX})$. Ввиду того, что массив мод полностью определяется массивом частот $M\mu(\mathbf{MX})$, дальнейшее изложение будет посвящено нахождению этого массива. Здесь n — полное число найденных мод.

Вычисления носят итерационный характер. Пусть найдена $k - 1$ мода. Для следующих («младших») мод из интервальной выборки сле-

Таблица 2: Алгоритм для нахождения частот мод
мультимодальной интервальной выборки

Вход

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

Массив частот $\{\mu\}$. Первая мода $\text{mode } \mathbf{X}^{(1)}$.

Выход

Массив областей $\mathbf{MX} = \{\mathbf{IH}\}$ выборки \mathbf{X} и массив частот $M\mu(\mathbf{MX})$.

Алгоритм

$\mathbf{MX} \leftarrow \text{mode } \mathbf{X}^{(1)}; M\mu \leftarrow \max \mu$

Установка признака выполнения задания $\mathbf{G0}=1$

WHILE $\mathbf{G0}$

 Обработка текущей моды k .

Вычисление $\text{mode } \mathbf{X}^{(k)}$ — Алгоритм Табл. 1 §0.1

 Вычисление границ текущей моды k :

 Установка границ интервала $\mathbf{IH}^{(k)} = \text{mode } \mathbf{X}^{(k)}$.

Расширение границ $\mathbf{IH}^{(k)}$ — Алгоритм Табл. 3.

 Зануление значений массива частот μ с индексами $\mathbf{IH}^{(k)}$:

$\mu(\mathbf{IH}^{(k)}) := 0$.

$\mathbf{MX} \leftarrow \mathbf{MX} \cup \mathbf{IH}^{(k)}, \quad M\mu \leftarrow M\mu \cup \mu^{(k)}$.

 Вычисление распределения частот μ .

 IF $\max \mu \leq 1$ THEN

 Окончание работы $\mathbf{G0}=0$

 ELSE

 Переход к следующей моде

 ENDIF

END WHILE

дует удалить последовательно области $\mathbf{IH}^{(k)} \supseteq \text{mode } \mathbf{X}^{(k)}$, в которых содержатся уже найденные моды $\text{mode } \mathbf{X}^{(k)}$ с номерами $1, 2, \dots, n-1$. Обозначение \mathbf{IH} указывает на *interval hall* — «расширенную» интервальную оболочку моды. Области $\mathbf{IH}^{(k)}$ расширяют интервалы мод $\text{mode } \mathbf{X}^{(k)}$.

Критерием останова расширения (окончанием сегментации в алгоритме водораздела) является тот факт, что значения массива частот μ на границах $\underline{\mathbf{IH}}^{(k)}$ и $\overline{\mathbf{IH}}^{(k)}$ становятся меньше, чем значение частоты следующей моды

$$\max\{ \mu(\underline{\mathbf{IH}}^{(k)}), \mu(\overline{\mathbf{IH}}^{(k)}) \} < \max \mu. \quad (5)$$

При выполнении условия (5) полагаем, что текущая мода распределения грубо оконтурена.

Исключаем область расширенной области текущей моды $\mathbf{IH}^{(k)}$ из общей выборки частот, следующая мода становится текущей и для неё вычисляется расширенная область.

Множество интервалов

$$\mathbf{IH} = \{ \mathbf{IH}^{(1)}, \mathbf{IH}^{(2)}, \dots \} \quad (6)$$

грубо описывает структуру интервальной выборки \mathbf{X} , а множество интервалов

$$\bigcup_{1 \leq k \leq n} \text{mode } \mathbf{X}^{(k)} = \{ \text{mode } \mathbf{X}^{(1)}, \text{mode } \mathbf{X}^{(2)}, \dots \} \quad (7)$$

даёт набор соответствующих интервальных мод со значениями частот $\mu_n = \max_{k \in K^{(n)}} \mu_k$. Множества $K^{(k)}$ — суть множества номеров индексов $\mathbf{IH}^{(k)}$.

Наиболее проблемным является этап алгоритма Табл. 2 «Расширение границ $\mathbf{IH}^{(k)} \supseteq \text{mode } \mathbf{X}^{(k)}$ ». Он, как и основной алгоритм, носит итерационный характер. В графике частот могут содержаться участки немонотонности «склонов» графика μ или очень пологие склоны. В большинстве случаев для устойчивой работы алгоритма необходимы дополнительные параметры настройки. Пример реализации алгоритма приведён в Табл. 3. С примером реализации алгоритма на языке Octave можно ознакомиться на [4].

В случае унимодальной выборки конструкции (6) и (7) соответствуют определению интервальной моды 0.1.1.

Пример 0.2.1. Рассмотрим численный пример вычисления частот мультимоды по алгоритму Табл. 2 .

Пусть уже рассчитан массив частот с помощью «Алгоритма для нахождения моды интервальной выборки», представленного в Табл. 1 §0.1.

$$\mu = \{1, 2, 3, 4, 4, 3, 2, 1, 3, 4, 5, 5, 4, 3, 0, 1, 3, 3, 1, 1, 0\}.$$

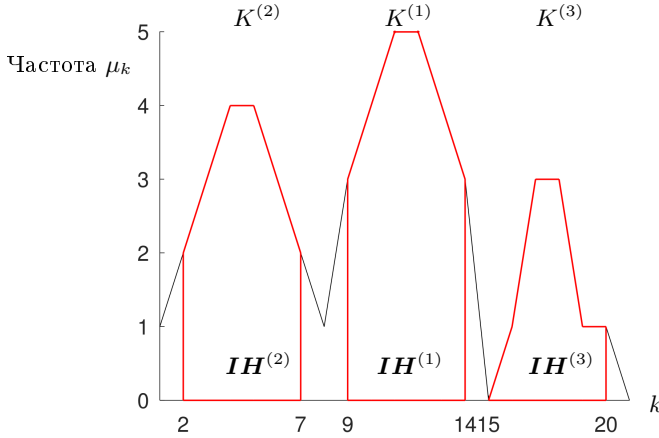


Рис. 5: Пример вычисления частот мультимоды.

На Рис. 5 представлен результат вычисления трёх областей \mathbf{IH} . Основной график представляет массив частот μ . Красным цветом даны области \mathbf{IH} для трёх мод.

Начальные значения индексов мод и частот

$$K = \{ [11, 12], [4, 5], [17, 18] \}, \quad M\mu = \{4, 5, 3\}.$$

Вычисления дают значения расширенных областей мод

$$\mathbf{IH} = \{[9, 14], [2, 7], [15, 20]\}.$$

Значения интервальных мод находятся как

$$\text{mode } \mathbf{X} = \bigcup_{1 \leq k \leq n} X(\mathbf{IH}^{(k)}).$$

■

Таблица 3: Расширение границ $\underline{IH}^{(k)} \supseteq \text{mode } X^{(k)}$

Вход
Мода k интервальной выборки $X = \text{mode } X^{(k)}$
Выход
Интервал $\underline{IH}^{(k)}$
Алгоритм
$\underline{IH}^{(k)} \leftarrow [\min \text{mode } X^{(k)}, \max \text{mode } X^{(k)}]$
<u>Вычисление начальных значений параметров алгоритма:</u>
$step$ — шаг по индексу, $noise$ — уровень шума
$addstep$ — увеличение шага по индексу
Проверяем монотонность μ на границах $\underline{IH}^{(k)}, \overline{IH}^{(k)}$
Проверяем условие COND валидности расширения - формула (10)
WHILE COND
IF COND
$\underline{IH}^{(k)} \leftarrow \underline{IH}^{(k)} - step,$
$\overline{IH}^{(k)} \leftarrow \overline{IH}^{(k)} + step$
ELSE
$step \leftarrow step + addstep$
END IF
Проверяем условие COND (10)
END WHILE
$\underline{IH}^{(k)} \leftarrow [\underline{IH}^{(k)}, \overline{IH}^{(k)}].$

В алгоритме Табл. 3 выполняется проверка монотонности μ на границах $\underline{IH}^{(k)}$, $\overline{IH}^{(k)}$. В обозначениях алгоритма это означает проверку условий

$$\mu(\underline{IH}^{(k)}) \geq \mu(\underline{IH}^{(k)} - step), \quad (8)$$

$$\mu(\overline{IH}^{(k)}) \geq \mu(\overline{IH}^{(k)} + step). \quad (9)$$

Условия (8) и (9) используются независимо. Например, на Рис. 6 для первой моды левый и правый склоны пика не симметричны.

Начальный этап Вычисление начальных значений параметров алгоритма требует отдельного рассмотрения и довольно существенно зависит от характера зависимости μ и задач ее интерпретации. В разделе §0.3 представлен практический пример.

0.3 ПРИЛОЖЕНИЕ - не для книги

Приведём практический пример для иллюстрации работы алгоритма Табл. 3 Вычисление начальных значений параметров алгоритма и проблем, возникающих при работе с реальными данными.

Пример относится к спектроскопии электромагнитных излучений при высокой энергии. Речь идёт о детектировании гамма-квантов с энергией около 1 МэВ, что является типичной величиной для ядерных реакций, идущих в промышленных и исследовательских ядерных реакторах.

На Рис. 6 приведён типичный спектр регистрации такого кванта в сцинтилляционном детекторе. По оси абсцисс представлены энергии вторичных частиц (интервалы энергии), возникающих при регистрации. По оси ординат — абсолютное количество соответствующих событий. Спектр получен при помощи хорошо апробированной программы МСС3d [5]. Пик с наибольшей интенсивностью соответствует наиболее полному поглощению энергии гамма-кванта в детекторе. Между тем, в детекторе при попадании в него гамма-кванта возможны различные физические процессы, которые идут с различными вероятностями и приводят к сложному виду спектра, простирающегося влево вплоть до нулевых энергий, что соответствует прохождению гамма-кванта через детектор практически без взаимодействия.

Для того, чтобы придать процессу расширения границ объективный характер, необходимо учесть типичную разность между соседними ве-

Таблица 4: Вычисление значения параметра *noise* алгоритма Табл. 3

<p>Вход</p> <p>Массив частот $\{\mu\}$, значение уровня представительности $\beta < 1$</p> <p>Выход</p> <p><i>noise</i> — уровень шума</p> <p>Алгоритм</p> <p>$\{\Delta\mu\}_{k=1}^{n-1}$ — численное дифференцирование массива $\{\mu\}_{k=1}^n$</p> <p>Зададим число бинов гистограммы M</p> <p>Установим условие NOK=TRUE</p> <p>WHILE NOK</p> <p> Построение гистограммы $[z, N_z] = \text{HIST}(\Delta\mu)$,</p> <p> z_k — массив интервалов, N_z — сумма $\Delta\mu$ в бине k</p> <p> Отбрасываем выбросы гистограммы $\Delta\mu - (1)$</p> <p> Находим интервал $z_i = \arg \max_i N_z$</p> <p> IF $N_z > \beta \cdot \sum_{i=1}^n \mu_k$</p> <p> $noise \leftarrow \text{mid } z$, NOK = FALSE</p> <p> ELSE</p> <p> <u>изменить значение M</u> — (2)</p> <p> END IF</p> <p>END WHILE</p> <p>Возвращаем <i>noise</i></p> <p>(1) — выбросы гистограммы возникают в результате зануления элементов массива $\{\mu\}$ при последовательном нахождении мод</p> <p>(2) — отдельный алгоритм</p>
--

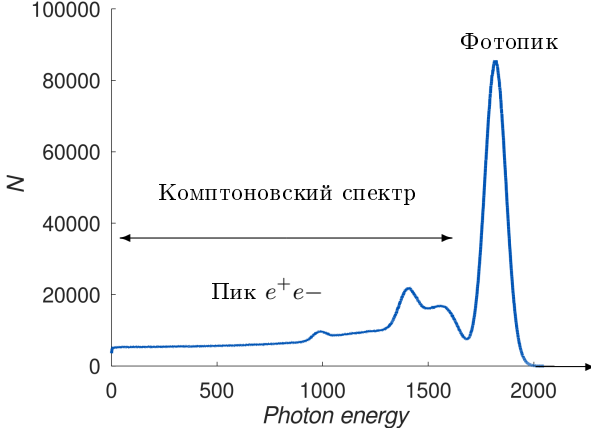


Рис. 6: Массив частот μ при регистрации гамма-квантов

личинами массива частот μ . Характерное значение шума опоставимо с наиболее вероятным значением гистограммы Рис. 7.

Способ вычисления параметра *noise* представлен в Табл. 4.

Выполним численное дифференцирование этого массива и построим гистограмму для массива N_z результата вычислений $\text{HIST}(\Delta\mu)$ с числом бинов гистограммы $M = 10$.

Интервал гистограммы $i = 4$ имеет наибольшее значение $N_4 \simeq 900$. Этот интервал, как видно из Рис.7, имеет середину $\text{mid } z_4 = -172$. Величина N_z немного меньше половины $\sum_{i=1}^n \mu_k$, так что при параметре $\beta < 0.4$ работа алгоритма 2 даёт рабочее значение параметра *noise* $\simeq 200$ алгоритма Табл. 4.

Уточним условие останова алгоритма (5):

$$\max \mu > \max \{ \mu(\underline{IH}^{(k)}), \mu(\overline{IH}^{(k)}) \} + \text{noise}. \quad (10)$$

По мере выявления мод, из распределения μ удаляются соответствующие участки распределения, так что и параметр *noise* в (10) тоже изменяется.

Параметр *step*, который используются для расширения границ $\underline{IH}^{(n)}$ задает степень подробности определения мод исходной выборки. Их величина не фиксирована чётко, и служит для исследований. Например, при обзорном рассмотрении распределений можно оконтурить

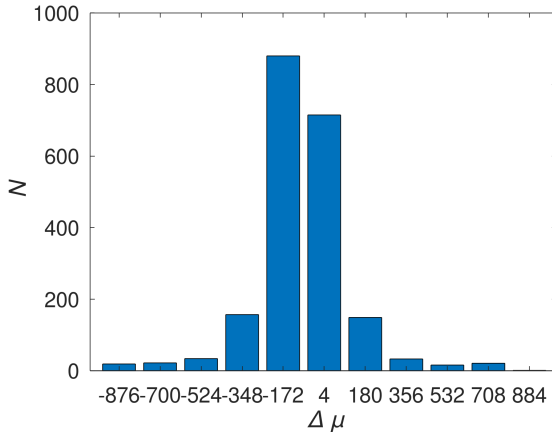


Рис. 7: Гистограмма разностей $\Delta\mu$.

мультиплетные области и далее разрешить их структуру с использованием меньшего шага.

Параметр *addstep* необходим для работы алгоритма в областях распределения, где шумы приводят к немонотонности.

Конкретно, на Рис. 8 показаны результаты работы алгоритма при величинах шага *step* = 100 и 50. Параметр *addstep* = 50.

Приведем численные результаты. Примем значение шага, равное 100. Получаем массив расширенных областей мод

$$\mathbf{IH} = \{ [1618, 2018], [1156, 1617], [744, 1155], \\ [481, 743], [217, 480], [1, 216] \}.$$

Более грубый шаг, равный 100, даёт и более грубое выявление структуры. Конкретно, область $\mathbf{IH}^{(2)} = [1156, 1617]$, на Рис. 8 а содержит два локальных максимума. Кроме того, область $\mathbf{IH}^{(1)}$ чрезмерно широка и содержит часть соседней области слева.

В случае шага, равного 50, удается разрешить структуру $\mathbf{IH}^{(2)}$, и выделить две моды на краю комптоновского распределения, [1305, 1506] и [1506, 1702], что видно из Рис. 8 б.

Область фотопика [1702, 1918] на Рис. 8 б также определена более корректно. Обратной стороной более тонкого вычисления является по-

Частота μ_k Шаг алгоритма $step = 100$

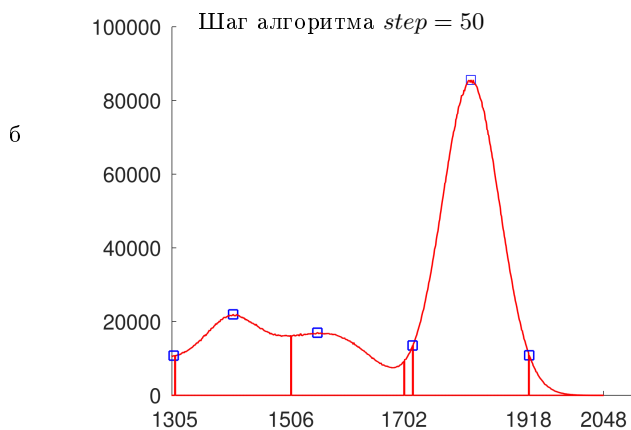
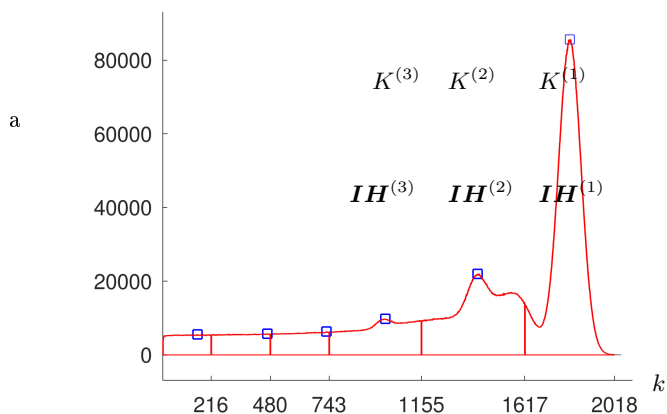


Рис. 8: Пример работы алгоритма Табл. 2
а — грубое распознавание структуры, б — уточнённое распознавание
структуры распределения мод мультимодальной выборки

явление артефактов малой ширины, например [1918, 2048]. Для их отбрасывания нужны дополнительные алгоритмы.

Литература

- [1] HU C., HU Z.H. On statistics, probability, and entropy of interval-valued datasets // Lesot M.J. et al. (eds) Information Processing and Management of Uncertainty in Knowledge-Based Systems. IPMU 2020. Communications in Computer and Information Science, vol 1239. – Cham: Springer, 2020.
- [2] MURPHY E.A. One cause? Many causes? The argument from the bimodal distribution // Journal of Chronic Diseases. – 1964. – Vol. 17. – P. 301–324.
- [3] Ельяшевич М.А. Атомная и молекулярная спектроскопия. — М.: Физматгиз, 1962; М.: Эдиториал УРСС, 2001.
- [4] Жилин С.И. Примеры и программы анализа интервальных данных в системе компьютерной математики Octave — <https://github.com/szhilin/octave-interval-examples>
- [5] БАГАЕВ К.А., КОЗЛОВСКИЙ С.С. Применение компьютерного моделирования для калибровки детекторов в водной среде // Журнал «Научнотехнические ведомости СПбГПУ. Физико-математические науки» 2011. № 2 (122). с. 106 - 111.