

МЕДИАНА ДЛЯ ИНТЕРВАЛЬНЫХ ДАННЫХ

Александр Вячеславович Пролубников

Омский государственный университет

a.v.prolubnikov@mail.ru

23.01.23

Определения медианы для точечных данных

Вариационный ряд: варианты — x_i . $x_1 \leq x_2 \leq \dots \leq x_N$.

Для вариационного ряда $\{x_i\}_{i=1}^N$, $x_i \in \mathbb{R}$:

- (1) Me — значение варианты*, для которой половина вариант с учётом их частот лежит слева, а половина — справа.
- (2) Me — значение варианты*, для которой минимальна сумма расстояний от неё до других вариант с учётом их частот: $\sum |x_i - Me| \rightarrow \min$.

Me — одна из характеристик

- выборки (вариационного ряда);
- вероятностного распределения.

Не чувствительна к выбросам.

(используется при оценивании распределения данных и нормализации (центрировании) данных)

Определения медианы интервального вар. ряда

Определение медианы Me интервального вариационного ряда требует задания:

- (1) *линейного* порядка на \mathbb{IR} ,
- (2) расстояния ρ между интервалами.

ПРИНЦИП СООТВЕТСТВИЯ:

- точечное значение — это *точное значение*;
- интервал — *приближённая оценка* точного значения.

— при уточнении измерений — *при узких интервалах* — значения медианы для интервальных данных рассчитываются эквивалентно расчёту медианы для точечных данных.

Способы задания отношения порядка на \mathbb{IR}

Возможные способы задания отношения порядка на \mathbb{IR} определены стандартом 1788 IEEE:

Говорят, что неравенство $\mathbf{a} \leq \mathbf{b}$ выполняется

- а) *в сильном смысле*, если $(\forall a \in \mathbf{a}) (\forall b \in \mathbf{b}) (a \leq b)$; $(\bar{\mathbf{a}} \leq \underline{\mathbf{b}})$;
- б) *в слабом смысле*, если $(\exists a \in \mathbf{a}) (\exists b \in \mathbf{b}) (a \leq b)$; $(\underline{\mathbf{a}} \leq \bar{\mathbf{b}})$;
- в) *в $\forall\exists$ -смысле*, если $(\forall a \in \mathbf{a}) (\exists b \in \mathbf{b}) (a \leq b)$; $(\bar{\mathbf{a}} \leq \bar{\mathbf{b}})$;
- г) *в $\exists\forall$ -смысле*, если $(\exists a \in \mathbf{a}) (\forall b \in \mathbf{b}) (a \leq b)$; $(\underline{\mathbf{a}} \leq \underline{\mathbf{b}})$;
- д) *(используем далее) в центральном смысле*, если $(\bar{\mathbf{a}} + \underline{\mathbf{a}})/2 \leq (\bar{\mathbf{b}} + \underline{\mathbf{b}})/2$.

Выбор отношения на \mathbb{IR} , задающего линейный порядок, определяет получаемое значение медианы интервальной выборки

Медиана интервальной выборки

Дана выборка: $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$, $x_i \in \mathbb{IR}$.

- 1) Как упорядочить \mathbf{X} ?
- 2) Как задать частоту интервальных данных?
 - возможно, все интервалы x_i окажутся разными, хотя и окажутся близки друг к другу.

От исходных интервалов x_i переходим к работе с их подинтервалами — *интервалами разбиения* с порядком в центральном смысле на них.

Минимальное разбиение \mathbf{R} , дающее множество линейно упорядоченных интервалов:

$\mathbf{R} = \{r_1, \dots, r_N\}$ — множество упорядоченных элементарных подинтервалов таких, что

$$(\forall x_i \in \mathbf{X}) (\exists j_1, \dots, j_t) (x_i = \cup_{i=1}^t r_{j_i}).$$

Частота элементарного подинтервала r_i — кол-во $x_i \in \mathbf{X}$ т.ч. $r_i \in x_i$.

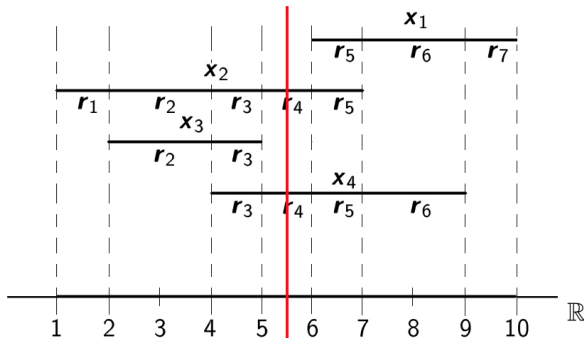
Разбиение на элементарные подинтервалы

Разбиение интервалов на подинтервалы концами интервалов:

Выборка: $\mathbf{X} = \{x_1, x_2, x_3\} \Rightarrow$ упорядоченные подинтервалы:

$$R = \{r_1, r_2, r_3, r_4, r_5, r_6, r_7\}.$$

R — вариационный ряд \Rightarrow может быть определена медиана Me .



Хаусдорфово расстояние

Пусть

- A и B — компактные множества в \mathbb{R}^n ;
- $\rho : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ — расстояние на \mathbb{R}^n .

$$\rho(A, B) = \max\left\{\sup_{x \in A} \inf_{y \in B} \rho(x, y), \sup_{y \in B} \inf_{x \in A} \rho(x, y)\right\}.$$

$\sup_{x \in A} \inf_{y \in B} \rho(x, y)$ — максимум минимального расстояния

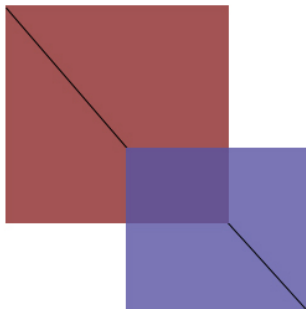
по точкам из B до точек из A .

Если $\mathbf{a}, \mathbf{b} \in \mathbb{I}\mathbb{R}$, то

$$\rho(\mathbf{a}, \mathbf{b}) = \max\left\{\sup_{x \in \mathbf{a}} \inf_{y \in \mathbf{b}} \rho(x, y), \sup_{y \in \mathbf{b}} \inf_{x \in \mathbf{a}} \rho(x, y)\right\} = \max\{|\underline{\mathbf{b}} - \underline{\mathbf{a}}|, |\bar{\mathbf{b}} - \bar{\mathbf{a}}|\}.$$

Хаусдорфово расстояние между интервалами

$\mathbf{a}, \mathbf{b} \in \mathbb{I}\mathbb{R}^2$:



Если $\mathbf{a}, \mathbf{b} \in \mathbb{I}\mathbb{R}$, то

$$\rho(\mathbf{a}, \mathbf{b}) = \max\left\{\sup_{x \in \mathbf{a}} \inf_{y \in \mathbf{b}} \rho(x, y), \sup_{y \in \mathbf{b}} \inf_{x \in \mathbf{a}} \rho(x, y)\right\} = \max\{|\underline{\mathbf{b}} - \underline{\mathbf{a}}|, |\overline{\mathbf{b}} - \overline{\mathbf{a}}|\}.$$

ТОЧЕЧНАЯ ИНТЕРВАЛЬНАЯ МЕДИАНА ИНТЕРВАЛЬНОЙ ВЫБОРКИ

Точечная медиана для сгруппированных данных

Интервалы задают разбиение всего интервала возможных значений, по элементам разбиения (интервалам) распределяются точечные значения элементов выборки.

Пример. Распределение студентов по возрастным группам в соответствии с заданным значением некоторого признака:

Возр. гр.	Кол-во студ.	Σ накопл. частот
<20	346	346
20–25	872	1218
25–30	1054	2272
30–35	781	3053
35–40	212	3265
40–45	121	3386
≥ 45	76	3462

Точечная медиана для сгруппированных данных

Вариационный ряд: *варианты* — это значения возраста студентов, точное значение которого в каждом случае не известно, а известна только его принадлежность некоторому интервалу разбиения интервала возможных значений.

⇒ вводим разбиение, чтобы упорядочить
точно не заданные значения вариант.

Частота интервала разбиения — частота попадания неточно заданных вариант в интервал.

В рассматриваемом примере:

$X = \{x_1, \dots, x_{3462}\}$, x_i — возраст i -го студента.

Точечная медиана для сгруппированных данных

Шаг 1. Определяем **медианный интервал** — интервал, содержащий варианту (либо две), которая делит вариационный ряд на две равные части: m — номер такой варианты x_m .

Шаг 2. Вычисляем **медиану** Me :

$$Me = b_0 + h \frac{\frac{1}{2} \sum_{i=1}^N f_i - S_{m-1}}{f_m},$$

b_0 — нижняя граница медианного интервала;

h — ширина этого интервала;

S_{m-1} — сумма накопленных частот для интервалов, предшествующих медианному;

f_m — частота медианного интервала — сколько неточно измеренных вариант x_i попадает в медианный интервал.

В точечном случае:

$$Me = b_0 + h \frac{\frac{1}{2} \sum_{i=1}^N f_i - S_{m-1}}{f_m}.$$

— точные значения некой величины (возраста) неизвестны, известна только их принадлежность *интервалам разбиения* всей области её значений

Me — ПРИБЛИЖЁННАЯ ОЦЕНКА МЕДИАНЫ ТОЧНЫХ ТОЧЕЧНЫХ ДАННЫХ

(которых нет — точный возраст каждого студента не известен)

— точка в медианном интервале, смещённая относительно его начала соответственно разнице суммарных частот до медианного интервала и точного значения полусуммы частот.

Точечная медиана для сгруппированных данных

Возр. гр.	Кол-во студ.	Σ накопл. частот
<20	346	346
20–25	872	1218
25–30	1054	2272
30–35	781	3053
35–40	212	3265
40–45	121	3386
≥ 45	76	3462

Интервалы —

$[0, 20], [20, 25], [25, 30], [30, 35], [35, 40], [40, 45], [45, 60]$

1. Медианный интервал — $\sum f_i = 1731 \Rightarrow$ это $[25, 30]$.

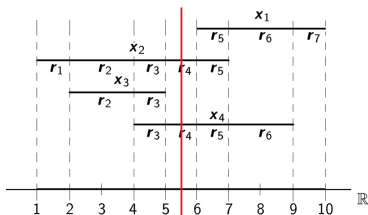
2.

$$Me = b_0 + h \frac{\sum f_i / 2 - S_{m-1}}{f_m} = 25 + 5 \frac{\frac{3462}{2} - 1218}{1054} = 27.4$$

Точечная медиана Me интервальной выборки

Интервал разбиения — элементарный подинтервал —
— это оценка значения, заданного с неопределённостью.

Частота f_i элементарного подинтервала r_i —
— это частота попадания r_i в $x_i \in X$.

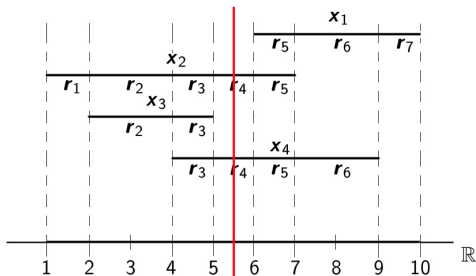


Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5	r_6	r_7
1	2	3	2	3	2	1

Точечная медиана Me интервальной выборки

Разбиение:



Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5	r_6	r_7
1	2	3	2	3	2	1

\Rightarrow Me принадлежит интервалу r_4 ,

$$Me = 5 + \frac{\frac{14}{2} - 6}{2} = 5.5$$

Точечная медиана Me интервальной выборки

Me :

$$Me = b_0 + h \frac{\frac{1}{2} \sum_{i=1}^N f_i - S_{m-1}}{f_m},$$

b_0 — нижняя граница медианного элементарного подинтервала r_m ;

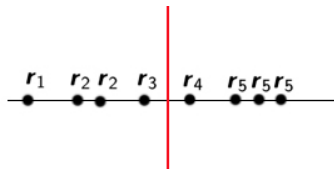
h — ширина r_m ;

S_{m-1} — сумма частот элементарных подинтервалов, предшествующих r_m ;

f_m — частота r_m — для скольких неточно измеренных значений x_i в них имеются значения, которые попадают в r_m .

Точечная медиана Me интервальной выборки

В случае, если **нет** элементарного подинтервала, который с учётом частот интервалов лежал бы точно посередине вариационного ряда, то есть середина в упорядоченной последовательности значений с повторениями лежит между r_{i-1} и r_i :



Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5
1	2	1	1	3

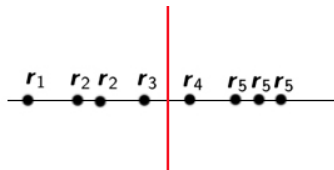
Точечная медиана Me интервальной выборки

В случае, если **нет** элементарного подинтервала, который с учётом частот интервалов, лежал бы точно посередине вариационного ряда, в формуле

$$Me = b_0 + h \frac{\frac{1}{2} \sum_{i=1}^N f_i - S_{m-1}}{f_m}$$

используем значения

- 1) $b_0 = \underline{r}_{i-1}$;
- 2) $h = \bar{r}_i - \underline{r}_{i-1}$;
- 3) $f_m = f_{i-1} + f_i$.



ИНТЕРВАЛЬНАЯ МЕДИАНА ИНТЕРВАЛЬНОЙ ВЫБОРКИ

Выборка точечных значений: $X = \{x_i\}_{i=1}^N$, $x_i \in \mathbb{R}$.

- (1) Me — значение варианты*, для которой половина членов ряда с учётом их частот лежит слева, а половина — справа.
- (2) Me — значение варианты*, для которой минимальна сумма расстояний от него до других вариантов с учётом их частот.

Медиана интервальной выборки $X = \{x_i\}_{i=1}^N$, $x_i \in \mathbb{IR}$:

- (1) $Me_f = r_m$ — накопленная сумма частот слева от r_m равна накопленной сумме частот справа от r_m :

$$\sum_{i=1}^{m-1} f_i = \sum_{i=m+1}^N f_i;$$

- (2) $Me_\rho = r_m$ — элементарный подинтервал r_m , для которого минимальна сумма хаусдорфовых расстояний от него до других элементарных интервалов с учётом их частот.

Интервальная выборка, содержащая точечные данные

Интервальная выборка \mathbf{X} содержит точечные данные —

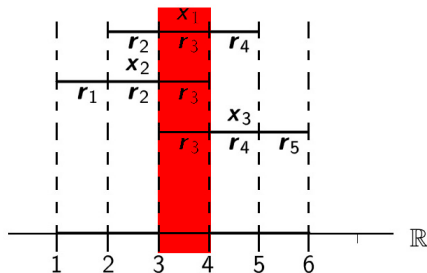
— в \mathbf{X} есть вырожденные интервалы: $\mathbf{x}_i = [v, v]$.

$\mathbf{r}_j = [v, v]$ — отдельный элементарный подинтервал.

Если есть $\mathbf{r}_k = [v, w]$, то частоты \mathbf{r}_j и \mathbf{r}_k считаются отдельно.

- Me_ρ : расстояние: $\rho(\mathbf{r}_j, \mathbf{r}_k) = |\bar{\mathbf{r}}_k - \underline{\mathbf{r}}_j| = |w - v|$.

ПРИМЕРЫ ВЫЧИСЛЕНИЯ ИНТЕРВАЛЬНОЙ МЕДИАНЫ



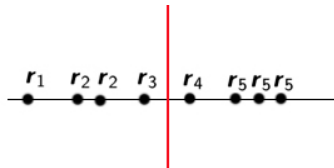
Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5
1	2	3	2	1

$$Me_f = r_3 = [3, 4]$$

Вычисление Me_f

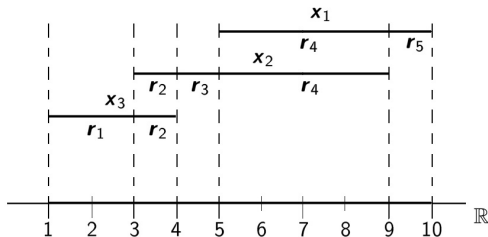
В случае, если **нет** элементарного подинтервала, который с учётом частот подинтервалов, лежал бы точно посередине вариационного ряда:



Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5
1	2	1	1	3

$$Me_f = r_3 \cup r_4$$

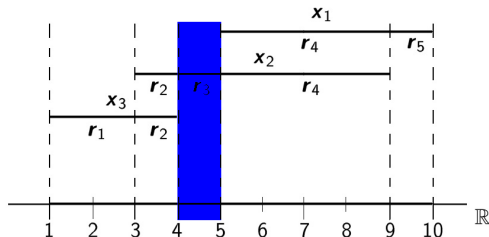


Частоты интервалов разбиения r_i :

r_1	r_2	r_3	r_4	r_5
1	2	1	2	1

$$Me_\rho = r_3 = [4, 5]$$

Вычисление Me_ρ



$$\sum_{k \neq 1} (f_k \cdot \rho(\mathbf{r}_1, \mathbf{x}_k)) = 27, \quad \sum_{k \neq 2} (f_k \cdot \rho(\mathbf{r}_2, \mathbf{x}_k)) = 19, \quad \sum_{k \neq 3} (f_k \cdot \rho(\mathbf{r}_3, \mathbf{x}_k)) = 18,$$

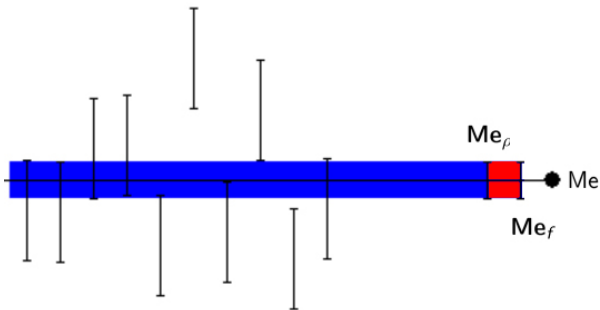
$$\sum_{k \neq 4} (f_k \cdot \rho(\mathbf{r}_4, \mathbf{x}_k)) = 24, \quad \sum_{k \neq 5} (f_k \cdot \rho(\mathbf{r}_5, \mathbf{x}_k)) = 33.$$

$$Me_\rho = \mathbf{r}_3 = [4, 5]$$

КАК СООТНОСЯТСЯ Me , Me_f И Me_ρ ?

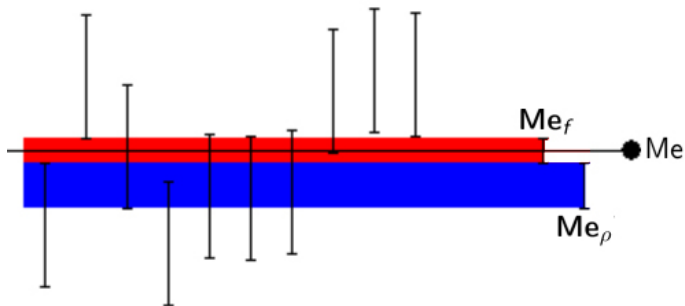
Генерирование интервалов: интервалы исходной выборки $X_0 = \{x_i\}_{i=1}^N$, $x_i = [-1, 1]$, случайно смещаются вверх и вниз.

Наиболее типичная ситуация:



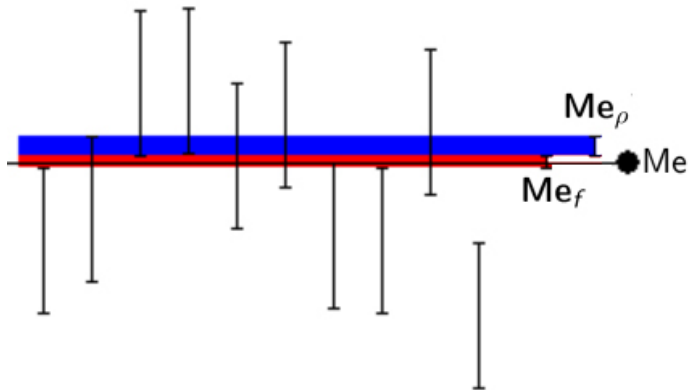
Me_f и Me_ρ — совпадающие элементарные подинтервалы

Возможно:



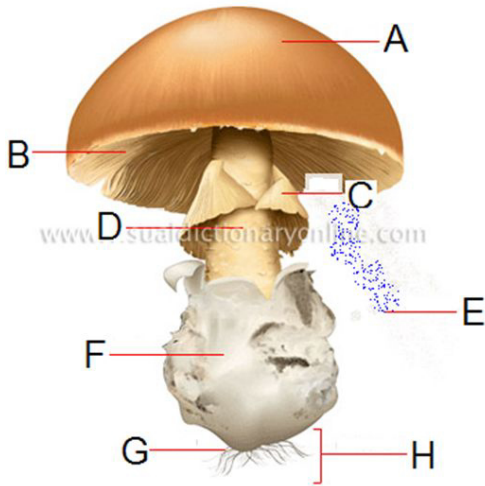
Me_f и Me_ρ — соседние элементарные подинтервалы

Возможно:



Me_f и Me_ρ — соседние элементарные подинтервалы

Mushroom structure



A: Cap(pileus)

B: Gill

C: Ring

D: Stalk,Stipe

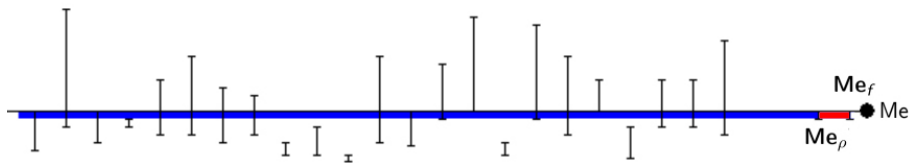
E: Spores

F: Volva

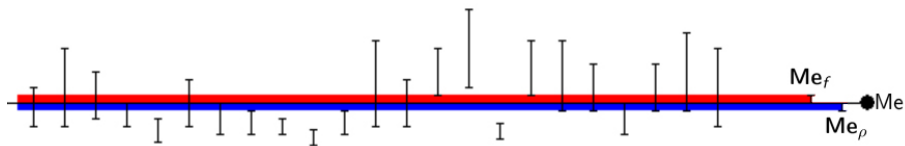
G: Hypa(e)

H: Mycelium

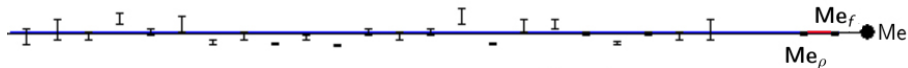
Выборка параметров грибов. Размер шляпы.



Выборка параметров грибов. Длина ножки.



Выборка параметров грибов. Ширина ножки.



Обобщение на случай многомерных интервалов

$$\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N, \mathbf{x}_i \in \mathbb{I}\mathbb{R}^n.$$

ВАРИАНТ 1.

$$\mathbf{Me} = (\mathbf{Me}_1, \dots, \mathbf{Me}_n) \in \mathbb{I}\mathbb{R}^n,$$

где \mathbf{Me}_j — медианы компонент \mathbf{x}_j .

ВАРИАНТ 2.

- 1). Находим разбиение всех интервалов выборки \mathbf{X} концами интервалов \mathbf{x}_j на элементарные подбрусы.
- 2). $\mathbf{Me} \in \mathbb{I}\mathbb{R}^n$ — брус, выбираемый из элементарных подбрусов, такой, что расстояние от него до всех остальных элементарных подбрусов минимально.