

ПРИЛОЖЕНИЕ 7

*В 3/9 царстве, в 3/10 государстве все были помешаны на дробях
Результат наблюдения*

Статистические оценки для интервального ранжированного частотного ряда

В некоторых исследованиях исходные данные для статистической оценки задаются интервальным ранжированным частотным рядом. Например, в виде

| Возраст особи | | частота встречаемости |
|---------------|--------|--------------------------|
| лет от | лет до | |
| 18 | 21 | 1 |
| 21 | 24 | 3 |
| 24 | 27 | 6 |
| 27 | 30 | 10 |
| 30 | 33 | 5 |
| 33 | 36 | 3 |
| 36 | 39 | 2 |



Требуется вычислить статистические оценки данных (в списке стандартных функций электронных таблиц необходимые средства для данной формы входных данных отсутствуют). А именно, требуется определить среднее, медиану, моду и параметры вариации.

Медиана (Me) – это величина, которая соответствует варианту, находящемуся в середине ранжированного ряда.

Для ранжированного ряда с нечетным числом индивидуальных величин (например, 1, 2, 4, 4, 6, 7, 8, 8, 10) медианой будет величина, которая расположена в центре ряда, т.е. пятая величина.

Для ранжированного ряда с четным числом индивидуальных величин (например, 1, 5, 7, 10, 11, 14) медианой будет средняя арифметическая величина, которая рассчитывается из двух смежных величин. В данном случае медиана равна $(7+10) : 2 = 8,5$.

Таким образом, для нахождения медианы сначала необходимо определить ее порядковый номер (ее положение в ранжированном ряду) по формуле

$$N_{Me} = \frac{n + 1}{2},$$

где n – объем совокупности.

Численное значение медианы определяют по накопленным частотам дискретного вариационного ряда, для чего требуется указать интервал нахождения медианы в интервальном ряду распределения. Медианным называют первый интервал, где сумма накопленных частот превышает половину наблюдений от общего числа всех наблюдений.

Численное значение медианы определяются по формуле

$$Me = x_{Me} + \Delta \frac{\frac{n + 1}{2} - S_{-1}}{f_{Me}},$$

где x_{Me} – нижняя граница медианного интервала;
 Δ – величина интервала;
 S_{-1} – накопленная частота интервала, которая предшествует медианному;
 f_{Me} – частота медианного интервала.

Модой (M_o) называют значение признака, которое встречается наиболее часто у единиц совокупности. Для дискретного ряда модой будет являться вариант с наибольшей частотой. Для определения моды интервального ряда сначала определяют модальный интервал (интервал, имеющий наибольшую частоту). Затем в пределах этого интервала находят то значение признака, которое может являться модой.

Чтобы найти конкретное значение моды, необходимо использовать формулу

$$M_o = x_{M_o} + \Delta \frac{f_{M_o} - f_{M_o-1}}{(f_{M_o} - f_{M_o-1}) + (f_{M_o} - f_{M_o+1})}$$

где x_{M_o} – нижняя граница модального интервала;
 f_{M_o} – частота модального интервала;
 f_{M_o-1} – частота интервала, предшествующего модальному;
 f_{M_o+1} – частота интервала, следующего за модальным.

Размах вариации R . Это самый доступный по простоте расчета абсолютный показатель, который определяется как разность между самым большим и самым малым значениями признака у единиц данной совокупности:

Вычисление размаха количественной вариации (выборки) производится по формуле:

$$R = x_{max} - x_{min} \quad \text{где } x_{max} \text{ – значение максимальной варианты,} \\ x_{min} \text{ – значение минимальной варианты выборки.}$$

Функции MS Excel: МИН(данные); МАКС(данные).
Функция Calc (Open Office): =MIN(данные) и =MAX(данные).

Размах вариации (размах разброса данных) – важный показатель, но только крайних отклонений. Для более точной характеристики рассеяния вариации признака используются другие показатели.

Среднее отклонение (выборочная оценка среднего отклонения), подобно стандартному отклонению, характеризует разброс эмпирической выборки относительно среднего значения и вычисляется по формуле

$$\bar{d} = \frac{1}{n} \sum_i |x_i - \bar{x}|, \quad n - \text{численность выборки,}$$

$x_i - \text{значения вариант выборки.}$

Выборочное среднее значение \bar{x}
определяется очевидным соотношением

$$\bar{x} = \frac{\sum x_i f_i}{\sum f_i}.$$

Среднее отклонение отражает так называемый модульный подход к вычислению меры отклонения между величинами в противоположность тому, что стандартное отклонение отражает квадратический подход.

Функция MS Excel: СРОТКЛ(данные);
Функция Calc (Open Office): AVEDEV(данные).

Частотно-взвешенная оценка среднего отклонения вычисляется соотношением

$$\bar{d} = \frac{\sum (|x_i - \bar{x}| f_i)}{\sum f_i}.$$

При использовании показателя среднего линейного отклонения возникают определенные неудобства, связанные с расчетом модуля разности. Широкое распространение получили обобщающие показатели, найденные с использованием вторых степеней отклонений. В предположении принадлежности выборки генеральной совокупности к таким показателям относятся несмещенная оценка дисперсии D^2

$$D^2 = \frac{\sum(x_i - \bar{x})^2}{n}$$

и среднее квадратическое отклонение D

$$D = \sqrt{D^2} = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n}}.$$

Функция MS Excel: ДИСПР(данные);

Функция Calc (Open Office): VARP(данные).

Частотно-взвешенная оценка дисперсии вычисляется соотношением

$$D^2 = \frac{\sum[(x_i - \bar{x})^2 f_i]}{\sum f_i} = \frac{\sum x_i^2 f_i}{\sum f_i} - \bar{x}^2$$

Кроме показателей вариации, выраженных абсолютными величинами, в статистических исследованиях используются относительные показатели вариации V_* , в частности для сравнения разброса признаков нескольких совокупностей.

Данные показатели рассчитываются как отношение размаха вариации к средней величине признака (коэффициент осцилляции), отношение среднего линейного отклонения к средней величине признака (линейный коэффициент вариации), отношение среднего квадратического отклонения к средней величине признака (коэффициент вариации) и, как правило, выражаются в процентах.

Формулы расчета относительных показателей вариации:

| | |
|---------------------------------------|---|
| – коэффициент осцилляции V_R | $V_R = \frac{R}{\bar{x}} \cdot 100\%$, |
| – линейный коэффициент вариации V_A | $V_A = \frac{\bar{d}}{\bar{x}} \cdot 100\%$, |
| – коэффициент вариации V_D | $V_D = \frac{D}{\bar{x}} \cdot 100\%$, |

Отметим следующее – из приведенных формул видно, что чем больше коэффициент V_* ближе к нулю, тем меньше вариация значений признака.

В статистической практике наиболее часто применяется коэффициент вариации V_D . Он используется не только для сравнительной оценки вариации, но и для характеристики однородности совокупности. Для распределений, близких к нормальному, совокупность считается однородной, если коэффициент вариации не превышает 33%.

Пример П7. Определить среднее, медиану, моду и параметры вариации, заданной таблицей справа.

Расчет описанных выше статистических параметров для интервального ранжированного частотного ряда можно выполнить по следующей схеме.

| Возраст особи | | частота встречаемости |
|---------------|--------|-----------------------|
| лет от | лет до | |
| 18 | 21 | 1 |
| 21 | 24 | 3 |
| 24 | 27 | 6 |
| 27 | 30 | 10 |
| 30 | 33 | 5 |
| 33 | 36 | 3 |
| 36 | 39 | 2 |

1. В соответствующие ячейки вносятся текстовые поясняющие данные (см. рис. П7.1).
2. Заносятся исходные данные в диапазоны A5:B11, D5:D11 (на рис. П7.1 выделены голубым фоном).
3. В ячейке C5 формулой $=(A5+B5)/2$ определяется средний по диапазону возраст. Далее автозаполнение на диапазон C6:C11.
4. Формулой $=СУММ(D5:D11)$ в ячейке D12 находится сумма частот.
5. Столбец накопленных частот строится внесением в E1 формулы $=D5$, в E2 формулы $=E5+D6$ и тиражированием последней на диапазон E7:E11.

| | A | B | C | D | E | F | G | H | |
|----|---------------|--------|--------------------------------------|-------------|---------|------|----------------------|---|--|
| 1 | | | | | | | $= (A5+B5)/2$ | | |
| 2 | | | | | | | $= D5$ | | |
| 3 | возраст особи | лет | лет | частота | частота | | | | |
| 4 | лет от | лет до | среднее | f_i | S_i | d | $= E5+D6$ | | |
| 5 | 18 | 21 | 19,5 | 1 | 1 | 9,2 | $= ABS(C5-B\$14)*D5$ | | |
| 6 | 21 | 24 | 22,5 | 3 | 4 | 18,6 | | | |
| 7 | 24 | 27 | 25,5 | 6 | 10 | 19,2 | | | |
| 8 | 27 | 30 | 28,5 | 10 | 20 | 2 | | | |
| 9 | 30 | 33 | 31,5 | 5 | 25 | 14 | | | |
| 10 | 33 | 36 | 34,5 | 3 | 28 | 17,4 | $= СУММ(F5:F11)$ | | |
| 11 | 36 | 39 | 37,5 | 2 | 30 | 17,6 | | | |
| 12 | | | | 30 | | 98 | $= СУММ(D5:D11)$ | | |
| 13 | | | | | | | | | |
| 14 | $\bar{x} =$ | 28,70 | $= СУММПРОИЗВ(C5:C11;D5:D11)/D12$ | | | | | | |
| 15 | $N =$ | 15,50 | $= (D12+1)/2$ | $= B5-A5$ | | | | | |
| 16 | $\Delta =$ | 3,00 | | | | | | | |
| 17 | $Me =$ | 28,65 | $= A8+B16*(B15-E7)/D8$ | | | | | | |
| 18 | $Mo =$ | 28,33 | $= A8+B16*(D8-D7)/(D8-D7+D8-D9)$ | | | | | | |
| 19 | | | | | | | | | |
| 20 | $R =$ | 18,00 | $= C11-C5$ | $= F12/D12$ | | | | | |
| 21 | $d =$ | 3,27 | $= СУММПРОИЗВ(C5:C11;C5:C11;D5:D11)$ | | | | | | |
| 22 | $D^2 =$ | 18,56 | $/D12-B14*B14$ | | | | | | |
| 23 | $D =$ | 4,31 | $= КОРЕНЬ(B22)$ | | | | | | |
| 24 | | | | | | | | | |
| 25 | $V_D =$ | 0,15 | $= B23/B14$ | | | | | | |
| 26 | $V_R =$ | 0,63 | $= B20/B14$ | | | | | | |
| 27 | $V_A =$ | 0,11 | $= B21/B14$ | | | | | | |
| 28 | | | | | | | | | |

Рис. П7.1 Скриншот схемы вычислений

6. В столбце F5:F11 подсчитываются частичные суммы среднего отклонения: в F5 заносится формула =ABS(C5-B\$14)*D5 и тиражируется автозаполнением на диапазон F6:F11. В ячейке F12 определяется сумма этих величин =СУММ(F5:F11).

7. В соответствии с нижеприведенной таблицей рассчитываются все требуемые параметры распределения.

| адрес ячейки | Формула в ячейке | Пояснение |
|--------------|--------------------------------|---|
| A14 | =СУММПРОИЗВ(C5:C11;D5:D11)/D12 | Подсчитывается среднее по формуле $\bar{x} = \frac{\sum x_i f_i}{\sum f_i}$. |
| A15 | =(D12+1)/2 | Определяется порядковый номер медианы $N_{Me} = \frac{n+1}{2}$ где $n = \sum x_i$ |
| A16 | =B5-A5 | Определяется величина интервала Δ |
| A17 | =A8+B16*(B15-E7)/D8 | Рассчитывается значение медианы. По величине N_{Me} определяется номер строки, где значение накопленной частоты "содержит" N_{Me} . В данном примере это строка 8, соответствующее накопленная частота отмечена красной рамкой. $Me = x_{Me} + \Delta \frac{\frac{n+1}{2} - S_{-1}}{f_{Me}}$ |

| | | |
|-----|---|--|
| A18 | =A8+B16*(D8-D7)/(D8-D7+D8-D9) | <p>Рассчитывается значение моды. По максимальному значению частоты в столбце D определяется номер строки, по положению которой в ряду распределения рассчитывается мода:</p> $M_o = x_{M_o} + \Delta \frac{f_{M_o} - f_{M_o-1}}{2f_{M_o} - f_{M_o-1} - f_{M_o+1}}$ <p>Максимальное значение частоты отмечено красной рамкой.</p> |
| A20 | =C11-C5 | <p>Рассчитывается размах вариации. Для ранжированного ряда это весьма просто.</p> $R = x_{max} - x_{min}$ |
| A21 | =F12/D12 | <p>Вычисляется частотно-взвешенная оценка среднего отклонения</p> $\bar{d} = \frac{\sum(x_i - \bar{x} f_i)}{\sum f_i}$ |
| A22 | =СУММПРОИЗВ(C5:C11;C5:C11;D5:D11)/D12-B14*B14 | <p>Считается дисперсия выборки:</p> $D^2 = \frac{\sum x_i^2 f_i}{\sum f_i} - \bar{x}^2$ |
| A23 | =КОРЕНЬ(B22) | <p>Среднеквадратическое отклонение:</p> $D = \sqrt{D^2}$ |

| | | |
|-----|----------|--|
| A25 | =B23/B14 | Коэффициент вариации $V_D = \frac{D}{\bar{x}} \cdot 100\%$ |
| A26 | =B20/B14 | Коэффициент осцилляции $V_R = \frac{R}{\bar{x}} \cdot 100\%$ |
| A27 | =B21/B14 | Линейный коэффициент вариации $V_A = \frac{\bar{d}}{\bar{x}} \cdot 100\%$ |

7. Окончательно:

| | |
|---------------------------------|-------|
| Средневыборочное | 28,70 |
| Медиана | 28,65 |
| Мода | 28,33 |
| Размах вариации | 18,00 |
| Среднее отклонение | 3,27 |
| Дисперсия выборки | 18,56 |
| Среднеквадратическое отклонение | 4,31 |
| Коэффициент вариации | 0,15 |
| Линейный коэффициент вариации | 0,11 |
| Коэффициент осцилляции | 0,63 |