

СТАТИСТИЧЕСКИЙ АНАЛИЗ НАБЛЮДЕНИЙ, ИМЕЮЩИХ ИНТЕРВАЛЬНОЕ ПРЕДСТАВЛЕНИЕ

Б.Ю. ЛЕМЕШКО*, С.Н. ПОСТОВАЛОВ*

В статье рассмотрены вопросы проверки гипотез о согласии и вопросы построения точечных оценок по исходным данным, имеющим интервальное представление. Предложена модификация критериев согласия χ^2 Пирсона, отношения правдоподобия, Колмогорова, Смирнова, ω^2 и Ω^2 Мизеса и метода максимального правдоподобия вычисления точечных оценок по интервальной выборке. Приведены результаты программной реализации.

1. ВВЕДЕНИЕ

Наиболее общим представлением исходных данных в классической статистике будет частично группированная выборка, частными случаями которой являются негруппированные, группированные и цензурированные данные. Дальнейшим обобщением частично группированной выборки является *интервальная* выборка, в которой каждое наблюдение представлено интервалом $[a_i, b_i]$, которому принадлежит неизвестное точно значение x_i . Классификация одномерных выборок показана на рис.1. Группированная выборка задана непересекающимися интервалами, негруппированная - вырожденными интервалами, у которых $a_i = b_i$. Интервальное представление наблюдения можно интерпретировать как неточное измерение случайной величины, связанное либо с заведомо известной погрешностью измерительного прибора, либо с особенностями измеряемой величины. Например, уровень воды в реке невозможно измерить точно, так как он непрерывно колеблется, но можно представить в виде интервала задающего его верхнюю и нижнюю границы [1]. Интервальное представление можно получить также в результате обработки точных наблюдений, сопровождающихся усечением части данных (например, при *группировании* или *цензуре*).

Подробная библиография работ, связанных с вопросами интервальной статистики, приведена в [2].

Для статистического анализа интервальных наблюдений большинство классических методов непосредственно неприменимо, требуется их модификация, либо разработка новых. В настоящей статье рассмотрена модификация непараметрических критериев согласия Колмогорова, Смирнова,

* Доцент кафедры прикладной математики, канд. техн. наук

♥ Аспирант кафедры прикладной математики

нова, ω^2 и Ω^2 Мизеса и предложен способ построения точечных оценок, минимизирующих статистики соответствующих критериев согласия.

По сравнению с частным случаем поставленной задачи, рассмотренным в [3], нам удалось уточнить и упростить полученные выражения.



Рис. 1 Классификация выборочных наблюдений

1. ЭМПИРИЧЕСКАЯ ФУНКЦИЯ РАСПРЕДЕЛЕНИЯ

Пусть задана интервальная выборка:

$$\mathbf{x}_k = \left\{ \langle a_i, b_i, n_i \rangle \mid a_i \leq x_{ij} \leq b_i, j=1, \dots, n_i; a_i \in R, b_i \in R, i=1, \dots, k \right\},$$

где k -число интервалов; $n = \sum_{i=1}^k n_i$ - объём выборки; n_i - число наблюдений

в i -м интервале; a_i и b_i - границы i -го интервала; x_{ij} - точные значения наблюдений. Упорядочим граничные точки интервалов:

$$a_{(1)} < a_{(2)} < \dots < a_{(k)};$$

$$b_{(1)} < b_{(2)} < \dots < b_{(k)}.$$

Обозначим через $n_{a_{(i)}}$ - число наблюдений, принадлежащих интервалу с левой границей $a_{(i)}$, и $n_{b_{(i)}}$ - число наблюдений, принадлежащих интервалу с правой границей $b_{(i)}$. Тогда

наб:
 $F_n(:$

Так

На
для

а

и овалов

ценок,

исмотр-

я.

$$N_{a(i)} = \sum_{j=1}^i n_{a(j)}; \quad N_{b(i)} = \sum_{j=1}^i n_{b(j)}, \quad i=1, \dots, k.$$

Эмпирическая функция распределения $F_n(x)$, построенная по точным наблюдениям x_{ij} , будет ограничена снизу и сверху функциями $\underline{F}_n(x)$ и $\overline{F}_n(x)$, имеющими следующий вид:

$$\underline{F}_n(x) = \begin{cases} 0, & x < a(1) \\ \frac{N_{a(i)}}{n}, & a(i) \leq x < a(i+1), \quad i=1, \dots, k-1, \\ 1, & x \geq a(k) \end{cases}$$

$$\overline{F}_n(x) = \begin{cases} 0, & x < b(1) \\ \frac{N_{b(i)}}{n}, & b(i) \leq x < b(i+1), \quad i=1, \dots, k-1. \\ 1, & x \geq b(k) \end{cases}$$

Таким образом,

$$\underline{F}_n(x) \leq F_n(x) \leq \overline{F}_n(x), \quad \forall x \in R. \quad (1)$$

На рис. 2 показана эмпирическая функция распределения ($\underline{F}_n(x)$ и $\overline{F}_n(x)$) для различных типов исходных данных.

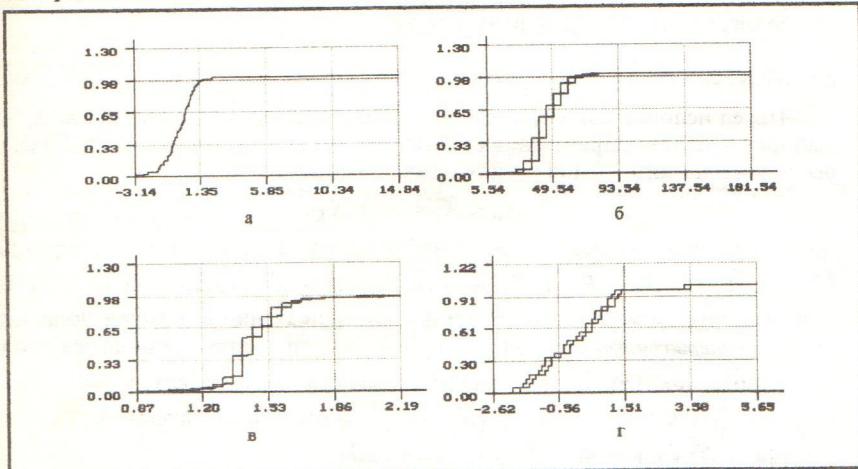


Рис. 2 Эмпирическая функция распределения для различных типов выборочных данных:
а - негруппированные; б - группированные; в - частично группированные; г - интервальные

2. ПРОВЕРКА ГИПОТЕЗ О СОГЛАСИИ

При проверке гипотез о согласии для найденного значения соответствующей статистики S^* вычисляется вероятность

$$p = P\{S > S^*\} = \int_{S^*}^{\infty} g(s)ds,$$

где $g(s)$ - плотность распределения статистики при условии истинности нулевой гипотезы. При заданном уровне значимости α гипотеза о согласии не отвергается, если $p > \alpha$. Когда выборка задана неточно, то статистика принадлежит интервалу $[S^*, \bar{S}^*]$, где на основании (1) границы определяются следующим неравенством:

$$\underline{S}^* = \inf_{F_n \leq F_n \leq \bar{F}_n} S^*(F_n, F) \leq S^*(F_n, F) \leq \sup_{F_n \leq F_n \leq \bar{F}_n} S^*(F_n, F) = \bar{S}^*$$

Вероятность $P\{S > S^*\}$ будет принадлежать интервалу $[p_{\min}, p_{\max}]$, где

$$p_{\min} = \int_{S^*}^{\infty} g(s)ds, \quad p_{\max} = \int_{\bar{S}^*}^{\infty} g(s)ds. \text{ Тогда, при заданном уровне значимости}$$

α , гипотезу о согласии следует отклонить, если $p_{\max} \leq \alpha$; гипотезу о согласии не следует отвергать, если $p_{\min} > \alpha$.

2.1. КРИТЕРИЙ χ^2 -ПИРСОНА

Перед использованием критерия необходимо сгруппировать исходную выборку. Область определения случайной величины разбивается на K непересекающихся интервалов граничными точками

$$X_0 < X_1 < \dots < X_K,$$

после чего подсчитывается число наблюдений m_i , попавших в интервалы $(X_i, X_{i+1}]$, $i = 0, 1, \dots, K - 1$. Процедура подсчета не представляет трудности, если выборка задана точно. В случае же, если граничная точка попадает внутрь интервального наблюдения $a_i < X_j < b_i$ процедура становится неоднозначной, так как точное значение можно отнести как к интервалу $(X_{j-1}, X_j]$, так и к интервалу $(X_j, X_{j+1}]$. Возможные значения m_i в соответствии с (1) удовлетворяют ограничениям:

толов

ответ-

ти ну-
ии не
стисти-
зделя-

], где

мости

зогла-

одную
непе-звалы
юсти,
адает
неод-
рвалу
соот-

$$n \cdot \underline{F}_n(X_i) \leq \sum_{j=0}^i m_j \leq n \cdot \overline{F}_n(X_i), \quad m_i \geq 0, \quad i = 0, 1, \dots, K-1; \quad (2)$$

$$\sum_{i=0}^{K-1} m_i = n. \quad (3)$$

Статистика критерия имеет вид

$$\chi^2 = \sum_{i=0}^{K-1} \frac{(m_i - np_i)^2}{p_i}.$$

Если m_i определяются неоднозначно, то можно найти максимум и минимум статистики χ^2 на области, заданной формулами (2) и (3):

$$\underline{\chi^2} = \min_{m_i} \chi^2; \quad \overline{\chi^2} = \max_{m_i} \chi^2.$$

Аналогично ограничивается статистика критерия отношения правдоподобия.

2.2. КРИТЕРИЙ СОГЛАСИЯ КОЛМОГОРОВА

Статистика критерия имеет вид

$$D_n = \sup_x |F_n(x) - F(x)|,$$

где $F_n(x)$ - эмпирическая функция распределения, $F(x)$ - теоретическая, согласие с которой проверяется, n - объем выборки. Преобразуем неравенство (1) к виду:

$$\underline{F}_n(x) - F(x) \leq F_n(x) - F(x) \leq \overline{F}_n(x) - F(x);$$

$$F(x) - \overline{F}_n(x) \leq F(x) - F_n(x) \leq F(x) - \underline{F}_n(x).$$

Эти неравенства выполняются для всех x , поэтому они сохраняются при взятии супремума:

$$\sup_x (\underline{F}_n(x) - F(x)) \leq \sup_x (F_n(x) - F(x)) \leq \sup_x (\overline{F}_n(x) - F(x)); \quad (4)$$

$$\sup_x (F(x) - \overline{F}_n(x)) \leq \sup_x (F(x) - F_n(x)) \leq \sup_x (F(x) - \underline{F}_n(x)). \quad (5)$$

Объединим эти неравенства в одно:

$$\underline{D}_n = \max \left\{ \sup_x (\underline{F}_n(x) - F(x)), \sup_x (F(x) - \overline{F}_n(x)) \right\} \leq$$

$$\leq D_n = \max \left\{ \sup_x (F_n(x) - F(x)), \sup_x (F(x) - F_n(x)) \right\} \leq \\ \leq \overline{D_n} = \max \left\{ \sup_x (\overline{F}_n(x) - F(x)), \sup_x (F(x) - \underline{F}_n(x)) \right\}. \quad (6)$$

2.3. КРИТЕРИЙ СОГЛАСИЯ СМИРНОВА

Статистика критерия имеет вид $D_n^+ = \sup_x (F_n(x) - F(x))$. Из неравенства (2) следует:

$$\underline{D_n^+} = \sup_x (\underline{F}_n(x) - F(x)), \quad \overline{D_n^+} = \sup_x (\overline{F}_n(x) - F(x)).$$

2.4. КРИТЕРИЙ СОГЛАСИЯ ω^2 МИЗЕСА

Построим вариационный ряд для точных значений наблюдений:

$$x_{(1)} < x_{(2)} < \dots < x_{(n)}.$$

Если наблюдения заданы неточно, то каждый член вариационного ряда известен с точностью до интервала

$$\underline{x_{(i)}} < x_{(i)} < \overline{x_{(i)}}, \quad (7)$$

где $\underline{x_{(i)}}$ и $\overline{x_{(i)}}$ можно определить из неравенства (1), так как между вариационным рядом и эмпирической функцией распределения имеется взаимно-однозначное соответствие (см. рис. 3) и $F_n(x_{(i)}) = i/n$:

$$\underline{x_{(i)}} = \overline{F_n^{-1}}(i/n), \quad \overline{x_{(i)}} = \underline{F_n^{-1}}(i/n).$$

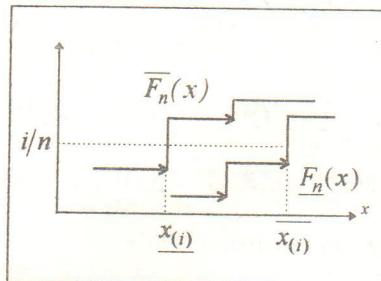


Рис. 3 Определение границ i -го члена вариационного ряда

Статистика критерия имеет вид

$$(6) \quad n\omega_n^2 = n \int_{-\infty}^{+\infty} (F_n(x) - F(x))^2 dF(x) = \frac{1}{12n} + \sum_{i=1}^n \left[F(x_{(i)}) - \frac{2i-1}{2n} \right]^2.$$

Пусть $s_i = \left[F(x_{(i)}) - \frac{2i-1}{2n} \right]^2$. Тогда из монотонности функции распределения $F(x)$ и условия (7) следует, что $\underline{s}_i \leq s_i \leq \bar{s}_i$, где \underline{s}_i и \bar{s}_i имеют вид:

$$\underline{s}_i = \begin{cases} \left[F(\overline{x_{(i)}}) - \frac{2i-1}{2n} \right]^2, & \text{при } F(\underline{x_{(i)}}) \leq F(\overline{x_{(i)}}) \leq \frac{2i-1}{2n}, \\ 0, & \text{при } F(\underline{x_{(i)}}) \leq \frac{2i-1}{2n} \leq F(\overline{x_{(i)}}), \\ \left[F(\underline{x_{(i)}}) - \frac{2i-1}{2n} \right]^2, & \text{при } \frac{2i-1}{2n} \leq F(\underline{x_{(i)}}) \leq F(\overline{x_{(i)}}), \end{cases}$$

$$\bar{s}_i = \max \left\{ \left[F(\overline{x_{(i)}}) - \frac{2i-1}{2n} \right]^2, \left[F(\underline{x_{(i)}}) - \frac{2i-1}{2n} \right]^2 \right\}.$$

Тогда

$$n\omega_n^2 = \frac{1}{12n} + \sum_{i=1}^n \underline{s}_i, \quad \bar{n}\omega_n^2 = \frac{1}{12n} + \sum_{i=1}^n \bar{s}_i.$$

2.5. КРИТЕРИЙ СОГЛАСИЯ Ω^2 МИЗЕСА

Статистика имеет вид

$$n\Omega_n^2 = -n - 2 \sum_{i=1}^n \left\{ \frac{2i-1}{2n} \ln F(x_{(i)}) + \left(1 - \frac{2i-1}{2n}\right) \ln(1 - F(x_{(i)})) \right\}.$$

Из неравенства (7) и монотонности функции распределения следует:

$$\left(F(x_{(i)}) \right)^{\frac{2i-1}{2n}} \left(1 - F(x_{(i)}) \right)^{1 - \frac{2i-1}{2n}} \leq \left(F(\overline{x_{(i)}}) \right)^{\frac{2i-1}{2n}} \left(1 - F(\underline{x_{(i)}}) \right)^{1 - \frac{2i-1}{2n}},$$

$$\left(F(x_{(i)}) \right)^{\frac{2i-1}{2n}} \left(1 - F(x_{(i)}) \right)^{1 - \frac{2i-1}{2n}} \geq \left(F(\underline{x_{(i)}}) \right)^{\frac{2i-1}{2n}} \left(1 - F(\overline{x_{(i)}}) \right)^{1 - \frac{2i-1}{2n}}.$$

Отсюда

$$\Omega_n^2 = -n - 2 \ln \prod_{i=1}^n \left(F(x_{(i)}) \right)^{\frac{2i-1}{2n}} \left(1 - F(x_{(i)}) \right)^{1-\frac{2i-1}{2n}};$$

$$\overline{\Omega_n^2} = -n - 2 \ln \prod_{i=1}^n \left(F(\bar{x}_{(i)}) \right)^{\frac{2i-1}{2n}} \left(1 - F(\bar{x}_{(i)}) \right)^{1-\frac{2i-1}{2n}},$$

2.6. ЧИСЛЕННЫЙ ПРИМЕР

Рассмотрим следующий пример. Была смоделирована негруппированная выборка объемом 100 наблюдений по нормальному закону с параметрами $(\bar{\mu}, \bar{\sigma})$. Оценки максимального правдоподобия параметров сдвига и масштаба по негруппированной выборке равны соответственно $\hat{\mu} = 0.048574$ и $\hat{\sigma} = 1.07144$ (рис. 4). Далее, предположили, что исходные наблюдения получены с некоторой ошибкой, не превосходящей по модулю $\Delta = 0.05$. Тогда каждое i -е наблюдение является интервалом вида $[x_i - \Delta, x_i + \Delta]$. Зафиксировав значение параметра сдвига равным оценке максимального правдоподобия по негруппированной выборке, было проверено согласие по всем рассмотренным критериям в зависимости от параметра масштаба (рис. 5). На рисунке изображены функции $p_{\max}(\sigma)$ и $p_{\min}(\sigma)$.

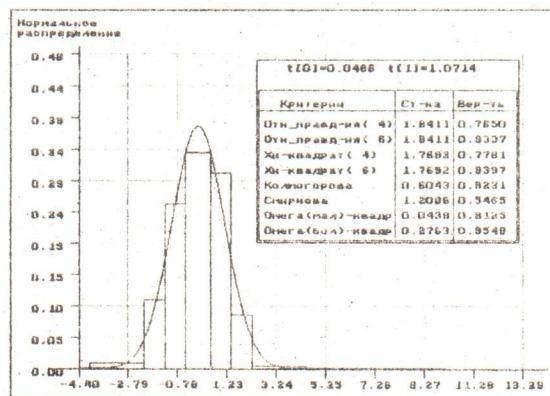


Рис. 4. Оценивание параметров нормального распределения по негруппированной выборке

На основании той же негруппированной выборки, была исследована зависимость длины интервала $p_{\max} - p_{\min}$ от Δ . Проверялось согласие с нормальным распределением с параметрами, оцененными по методу максимального правдоподобия. Очевидно, что с ростом неопределенности ис-

ходных данных (Δ), происходит увеличение неопределенности вероятности согласия ($\Delta p = p_{max} - p_{min}$). При этом критерии можно упорядочить по возрастанию Δp следующим образом: Колмогорова, Смирнова, ω^2 Мизеса, отношения правдоподобия, χ^2 Пирсона, Ω^2 Мизеса. Поэтому, при интервальном задании исходной выборки наиболее предпочтительным оказывается использование критерии *Колмогорова, Смирнова, ω^2 Мизеса.*

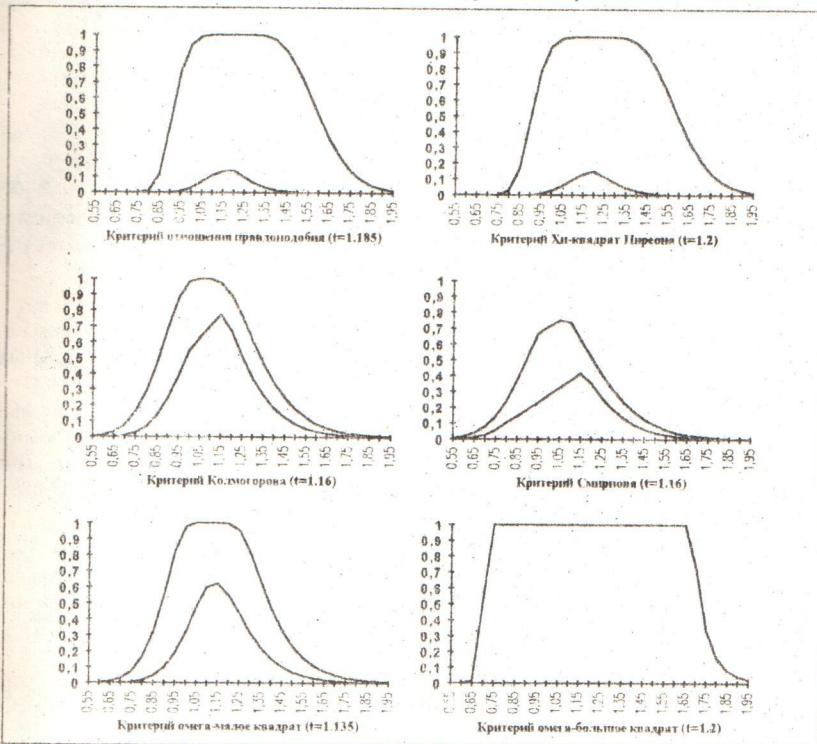


Рис. 5 Согласие интервальной выборки с нормальным распределением с параметрами (μ, σ) , где $\mu = 0.048574$, а σ изменяется от 0.55 до 1.95. Верхняя линия соответствует p_{max} , нижняя — p_{min}

3. ОЦЕНИВАНИЕ ПАРАМЕТРОВ РАСПРЕДЕЛЕНИЙ

Для получения оценок параметров распределения ни один из известных методов не применим в явном виде. Однако, используя полученные выше выражения оценок для статистик критериев согласия, не представляет труда построить оценки параметров распределений, аналогичные оценкам

типа минимума χ^2 . Очевидно, что чем меньше статистика критерия согласия, тем выше будет вероятность не отклонения гипотезы. Поскольку точное значение статистики неизвестно, то можно минимизировать ее верхнюю границу:

$$\hat{\theta} = \arg \min_{\theta} S^*(\underline{F_n}, \bar{F_n}, F(\theta)).$$

Оценку такого типа рассматриваем в терминах теории принятия решений в условиях неопределенности: как оценку, максимизирующую вероятность согласия при наихудшем расположении точных значений. На рис. 5 показаны оценки параметра σ в виде значений t , полученные таким способом. По всем критериям оценка масштабного параметра оказалась больше, чем оценка максимального правдоподобия по исходной точной выборке.

Более точные оценки удалось получить, используя модификацию *метода максимального правдоподобия*, когда при достаточно малых интервалах $[a_i, b_i]$ использовали разностную аппроксимацию производной функции распределения. Оценка масштабного параметра нормального распределения, полученная таким образом, оказалась равной $\sigma = 1,0657$, совсем не значительно отличаясь от точного значения.

4. ЗАКЛЮЧЕНИЕ

Получены общие выражения оценок снизу и сверху для статистик непараметрических критериев согласия в случае частично группированных данных или интервального (нечеткого) их представления.

Показано, что в случае интервального представления исходной выборки достаточно определенные выводы о согласии можно делать на основании полученных оценок для статистик критериев согласия. При этом наиболее эффективными оказываются непараметрические критерии Колмогорова, ω^2 Мизеса и Смирнова.

При оценивании параметров распределения было отмечено, что при малых длинах интервалов наилучшим является использование предложенной модификации метода максимального правдоподобия с разностной аппроксимацией производной функции распределения.

[1] FRUHWIRTH-SCHNATTER S. *On statistical inference for fuzzy data with applications to descriptive statistics* // *Fuzzy Sets and Systems*, 1992. - № 50. - С. 143 - 165.

[2] ОРЛОВ А.И. *Интервальная статистика* // *Интервальные вычисления*. - 1992. - № 1. - С. 44 - 52.

[3] ЛЕМЕШКО Б.Ю., ПОСТОВАЛОВ С.Н. *К использованию непараметрических критериев по частично группированным данным* // Сб. науч. тр. НГТУ. - Новосибирск, 1995. - Вып. 2. - С. 21 - 30.