

Специальный практикум (V курс)

Применение критерия согласия Пирсона.

Тестирование бинормальной модели для выборочной бимодальной функции распределения

И. И. Никифоров

1. Критерии согласия

Критерий согласия (КС) служит для проверки согласия между выборочным и гипотетическим распределениями или между двумя выборочными распределениями. КС позволяют оценить вероятность того, что полученная выборка не противоречит сделанному предположению о виде закона распределения рассматриваемой случайной величины X . Для этого выбирается некоторая статистика (*критериальная статистика*) χ , являющаяся мерой расхождения наблюдаемого и теоретического законов распределения. Функция распределения этой статистики, $\varphi(\chi)$, принимается известной, как правило, асимптотически близкой (при $N \rightarrow \infty$) к некоторой функции распределения, рассматриваемой в теореме.

На опыте получают значение меры расхождения χ_q , по которому определяют вероятность $\alpha_q \equiv P(\chi \geq \chi_q)$, которая называется *уровнем значимости*. Далее сравнивают α_q с *критическим* уровнем значимости α — малой величиной, значение которой устанавливается в соответствии с характером задачи. Как правило, принимают $\alpha = 0.05$ (уровень значимости $\approx 2\sigma$). Т.е. рассматривают нуль-гипотезу

$$H_0: X \sim \varphi_{\text{mod}}(x). \quad (1)$$

Гипотеза H_0 принимается или отвергается по следующему правилу:

$$\begin{aligned} \alpha_q \geq \alpha &\implies \text{расхождение не противоречит } H_0, \\ \alpha_q < \alpha &\implies H_0 \text{ отвергается с } P = 1 - \alpha_q \text{ (на уровне значимости } 1 - \alpha_q). \end{aligned} \quad (2)$$

Значения α_q , весьма близкие к 1 (очень хорошее согласие), могут указывать на недоброкачественность выборки (например, из первоначальной выборки без основания выброшены элементы, дающие большие отклонения от среднего).

Наиболее популярные критерии согласия.

1. КС Карла Пирсона (Karl Pearson, 1900): $\chi = \chi^2(\varphi_{\text{obs}}(x), \varphi_{\text{mod}}(x))$. Применим в том числе и в случае, если параметры модельного закона были найдены по данной выборке.

2. КС А. Н. Колмогорова: $\varkappa = \max |\Phi_{\text{obs}}(x) - \Phi_{\text{mod}}(x)|$. Применим в случае, когда параметры теоретического закона распределения определяются не по данным исследуемой выборки.
3. КС Н. В. Смирнова – А. Н. Колмогорова: $\varkappa = \max |\Phi_{\text{obs},1}(x) - \Phi_{\text{obs},2}(x)|$. Применяется для проверки гипотезы о принадлежности двух выборок одной генеральной совокупности.

2. Применение критерия Пирсона к тестированию бинормальной модели распределения металличностей шаровых скоплений Галактики

Из наблюдений известны

$$f_i \equiv [\text{Fe}/\text{H}]_i = \left[\lg \left(\frac{\text{Fe}}{\text{H}} \right)_* - \lg \left(\frac{\text{Fe}}{\text{H}} \right)_{\odot} \right]_i, \quad i = 1, \dots, N. \quad (3)$$

Ранее была решена задача разложения наблюдаемого распределения на две нормальные составляющие:

$$\mathcal{N}_1(\bar{F}_1, \sigma_1^2), \quad \mathcal{N}_2(\bar{F}_2, \sigma_2^2). \quad (4)$$

Т.е. были найдены точечные оценки параметров $\bar{F}_1, \sigma_1, \bar{F}_2, \sigma_2, c$ модельного дифференциального закона распределения

$$\varphi_{\text{mod}}(f) = \frac{c}{\sqrt{2\pi}\sigma_1} e^{-\frac{(f - \bar{F}_1)^2}{2\sigma_1^2}} + \frac{1-c}{\sqrt{2\pi}\sigma_2} e^{-\frac{(f - \bar{F}_2)^2}{2\sigma_2^2}}. \quad (5)$$

Задача. Применить критерий Пирсона для модели (5) и для гауссианы

$$\varphi_{\text{mod}}(f) = \mathcal{N}(f; \bar{F}_0, \sigma_0^2) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(f - \bar{F}_0)^2}{2\sigma_0^2}}. \quad (6)$$

Для этих двух модельных распределений найти статистику

$$\chi_q^2 = \sum_{j=1}^J \frac{(N_j - Np_j)^2}{Np_j}, \quad (7)$$

Здесь N — объем выборки; J — число ячеек (разрядов); N_j — число объектов, попавших в ячейку j ; p_j — теоретическое значение вероятности попадания случайной величины f в j -й интервал, вычисленная для модельного закона распределения по формуле

$$p_j = \int_{\hat{f}_j}^{\hat{f}_{j+1}} \varphi_{\text{mod}}(f) df. \quad (8)$$

При $N \rightarrow \infty$ закон распределения χ_q^2 независимо от вида закона распределения случайной величины стремится к закону χ^2 -распределения с $k = J - M - 1$ степенями свободы, где M — число параметров модельного закона распределения, найденных по гистограмме — по числам N_j объектов, попавшим в соответствующие ячейки. В наших случаях можно принять $M = 0$, т.к. параметры искались методом наибольшего правдоподобия, а не по гистограммам. Единица в k вычитается из-за условия нормировки $\sum N_j = N$.

Для применения критерия Пирсона в общем случае необходимо, чтобы объем выборки N и численности разрядов N_j были достаточно велики (практически считается достаточным, чтобы было $N \geq 50 \div 60$, $N_j \geq 5 \div 8$).

Выбор границ ячеек:

$$\hat{f}_1 = -\infty, \quad \hat{f}_{J+1} = +\infty; \quad (9)$$

$$\hat{f}_{j+1} = \hat{f}_j + \Delta \hat{f}, \quad j = 2, \dots, J-1; \quad (10)$$

$$\Delta \hat{f} = 0.1 \implies \hat{f}_2 = -2.3, \quad \hat{f}_J = -0.1; \quad (11)$$

$$\Delta \hat{f} = 0.2 \implies \hat{f}_2 = -2.2, \quad \hat{f}_J = -0.2. \quad (12)$$

Пусть ячейки включают правые границы $(\dots --] --] --] \dots)$. Т.е.

$$(-\infty, \hat{f}_2]; (\hat{f}_2, \hat{f}_3]; \dots; (\hat{f}_{J-1}, \hat{f}_J]; (\hat{f}_J, +\infty).$$

Или для $\Delta \hat{f} = 0.1$

$$(-\infty, -2.30]; (-2.30, -2.20]; \dots; (-0.2, -0.1]; (-0.1, +\infty);$$

для $\Delta \hat{f} = 0.2$

$$(-\infty, -2.20]; (-2.20, -2.00]; \dots; (-0.4, -0.2]; (-0.2, +\infty).$$

Для нормального модельного распределения параметры можно найти по формулам

$$\bar{F}_0 = \frac{1}{N} \sum_{i=1}^N f_i, \quad \sigma_0^2 = \frac{1}{N} \sum_{i=1}^N (f_i - \bar{F}_0)^2. \quad (13)$$

Значения вероятностей $\alpha_q(\chi_q^2, k) = P(\chi^2 \geq \chi_q^2)$ приводятся в таблицах.

3. Вычисление p_j и $\alpha_q(\chi_q^2, k)$

П/п ndtr.f — вычисляет $P(x \leq x_q)$, $x \in \mathcal{N}(0, 1)$.

П/п cdtr.f — вычисляет $P(\chi^2 \leq \chi_q^2)$ в зависимости от χ_q^2 и k . Нужно взять дополнение до 1, чтобы найти α_q . П/п использует ndtr.f и dlgam.f.

Сверить с таблицами в книжке у Агеяна.

4. Задание и данные

Для заданного набора $f_i, i = 1, \dots, N$, и для распределений (5) и (6), значений $\Delta\hat{f} = 0.1$ и 0.2 получить/привести результаты.

1. Привести параметры $\bar{F}_1, \sigma_1, \bar{F}_2, \sigma_2, c$ или \bar{F}_0, σ_0 .
2. Привести $N, J, k, \chi_q^2, \alpha_q$.
3. Графики сравнения наблюдаемого распределения с двумя модельными (для $\Delta\hat{f} = 0.1$ и 0.2).

Данные:

1. all.dat
2. wo_p1.dat
3. wo_p12.dat
4. wo_p123.dat

5. Приложение: визуализация

Нужно сопоставить модельную функцию $\varphi_{\text{mod}}(f; \bar{F}_1, \sigma_1, \bar{F}_2, \sigma_2, c)$ с гистограммой наблюдаемого распределения f на одном рисунке.

Гистограмма наблюдаемого распределения. J ячеек длиной $\Delta f = 0.1 \text{ dex}$; в каждой ячейке N_j объектов, $j = 1, 2, \dots, J$. $\sum N_j = N$ (при нормировке столбцов гистограммы на N). Пусть ячейки гистограммы включают правые границы $(\dots --] --] --] \dots)$. Например, $(-2.40, -2.30]; (-2.30, -2.20], \dots$

Модельная функция $\varphi_{\text{mod}}(f)$. Нужно согласовать нормировки гистограммы $\varphi_{\text{hist}}(f)$ и модельной функции $\varphi_{\text{mod}}(f)$:

$$\int_{-\infty}^{+\infty} \varphi_{\text{hist}}(f) d\varphi = \sum_{j=1}^J N_j \Delta f = N \Delta f, \quad (14)$$

$$\int_{-\infty}^{+\infty} \varphi_{\text{mod}}(f) d\varphi = 1. \quad (15)$$

Т.е. совместно с гистограммой вместо модельной функции $\varphi_{\text{mod}}(f)$ нужно рисовать

$$\boxed{N \Delta f \varphi_{\text{mod}}(f)}. \quad (16)$$

Используя формулу (16), изобразить φ_{mod} и две ее нормальные составляющие [см. (5)].