

МСПС ANOVA

7

Будем рассматривать теперь проверку гипотезы однородности (ну или более слабее гипотезы: отсутствие дива, $EX=0$ и т.д.) для нескольких (больше 2) выборок.

Сначала рассмотрим случай нормального выборок, и проверим сначала гипотезу о равенстве их дисперсий (как мы это делали в t-критерии Стьюдента)

Пусть $X_{ij} \sim N(\mu_j, \sigma_j^2)$ ($i=1 \dots n_j, j=1 \dots k$) - независ., μ_j и σ_j^2 неизвестны
Обозначим $\bar{X}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} X_{ij}$ - несмещ. оц. для μ_j ; $s_j^2 = \frac{1}{n_j-1} \sum_{i=1}^{n_j} (X_{ij} - \bar{X}_j)^2$ - несмещ. оц. для σ_j^2 ; $\bar{X} = \frac{1}{N} \sum_{j=1}^k \sum_{i=1}^{n_j} X_{ij}$.
 $N = \sum_{j=1}^k n_j$

Для проверки гипотезы $H': \sigma_1^2 = \dots = \sigma_k^2$ используют крит. Барнетта;

Статистика критерия $B = \left(\frac{1}{N-k} \sum_{j=1}^k (n_j - 1) s_j^2 \right) / \sqrt{\prod_{j=1}^k (s_j^2)^{n_j-1}}$

Теорема Пусть H' верна и $n_j \geq 3 \forall j$. Тогда

$$B^2 = \chi^2_{k-1} \xrightarrow{d} \chi^2_{k-1}, \text{ где } \chi = 1 + \frac{1}{3(k-1)} \left[\left(\sum_{j=1}^k \frac{1}{n_j} \right) - \frac{1}{N} \right].$$

Замечание Крит. Барнетта крайне чувств. к отклонениям от нормального распредел. Например, если при $k=10$ заменить $N(0,1)$ на t_7 (распр. Стьюдента), то факт. уровень значимости критерия $\{B^2 > u_{0.05}\}$ возрастет с 0.05 до 0.49.

Теперь проверим гипотезу $H'': \mu_1 = \mu_2 = \dots = \mu_k$ (при условии, что H' не отвергнута)

Для этого используем критерий Фишера (однофактор. дисперс. анализа)

$$\text{Введём } R = \frac{\frac{1}{k-1} \sum_{j=1}^k n_j (\bar{X}_j - \bar{X})^2}{\frac{1}{N-k} \sum_{j=1}^k \sum_{i=1}^{n_j} (X_{ij} - \bar{X}_j)^2}.$$

Если $H_0: \{\text{все выборки о.р.}\}$ верна, то $R \sim F_{k-1, N-k}$, если $\forall n_j > 1$.

Замечание 1) Окажется, критерий работает не только для нормального выборок.

Если $N-k-1 \geq 20$ и отношение $\frac{\max_j s_j^2}{\min_j s_j^2} < 10$, то при верной H_0

$R \xrightarrow{d} F_{k-1, N-k}$ (т.е. можно приближенно пользоваться)

2) При $n_1 = \dots = n_k$ метод устойчив к нарушению H_0 в этих 2 предположениях, если же объёмы выборок неравны, то вер. ошибки I рода может возрасти при невыполнении условия о равенстве дисперсий.

3) Не устойчив к выбросам.

Примечание

Ещё критерии проверки однородности дисперсий: тест Левина (Levene's test), тест Брауна-Форсайта (Brown-Forsythe test)

ММРС | ANOVA

(2)

Фактор., мы рассматриваем сейчас т.н. ~~однофакторную~~ модель

$X_{ij} = \mu + \beta_j + \varepsilon_{ij}$, где μ - общее сред., β_j - эффект воздействия, ε_{ij} - сл. ошибка. Пусть теперь ε_{ij} независ. и одинаково распредел.

Проверим гипотезу равенства $H_0: \beta_1 = \beta_2 = \dots = \beta_k$ с помощью крит. Краскала - Уоллиса:

Пусть $R_{ij} = R(X_{ij})$ - ранг наблюд. X_{ij} в общей совокупности.

$$\bar{R}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} R_{ij}, \quad \bar{R} = \frac{\sum_{j=1}^k \sum_{i=1}^{n_j} R_{ij}}{N} = \frac{N+1}{2}$$

$$\text{Обозначим } W := (N-1) \frac{\sum_{j=1}^k n_j (\bar{R}_j - \bar{R})^2}{\sum_{j=1}^k \sum_{i=1}^{n_j} (R_{ij} - \bar{R})^2} = \frac{12}{N(N+1)} \sum_{j=1}^k n_j R_j^2 - 3(N+1).$$

(Оказывается, знаменатель равен $(N^2 - N)/12$)

При $n_j > 5$ можно пользоваться аппроксимацией $W \sim \chi^2_{k-1}$, но при малых n_j типа 4-6 лучше брать поправку Шварца

$$\tilde{W} = \frac{1}{2} W \left[\frac{N-k}{N-1-W} + 1 \right], \quad \text{приблиз. крит. значения } \tilde{W}_{1-\alpha} = \frac{1}{2} \left[\chi^2_{k-1, 1-\alpha} + F_{k-1, N-k, 1-\alpha} \right]$$

Гипотезу равенства проверяет и критерий Феллхера (против $H_A: \beta_1 \leq \beta_2 \leq \dots \leq \beta_k$)

$$S = \sum_{j=1}^k \sum_{i=1}^{n_j} a_{ij}, \quad \text{где } a_{ij} - \text{число наблюд. у первых } j-1 \text{ выборок, } < X_{ij}.$$

Симмет. таблица распредел., при $n_j > 10$ можно пользоваться приближением: $S \sim N(\mu, \sigma^2)$, где $\mu = \frac{1}{4} (N^2 - \sum_{j=1}^k n_j^2)$, $\sigma^2 = \frac{1}{72} (N^2(2N+3) - \sum_{j=1}^k n_j^2(2n_j+3))$.

Оба этих теста устойчивы к выбросам. Фелхер точнее крит. Краскала-Уол. на альтернативе H_A .

Если у нас гипотеза однородности отверглась, то появляется естеств. желание узнать, какие пары (группы) выборок неоднородны. Это делается (в основном) с помощью двухвыбороч. критериев, пройденных на пред. лекции. А также см. сф. 3.

Двухфактор. дисперсионный анализ

Данные представлены в виде $X_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$, $i=1..n$, $j=1..k$.

Эта пар-зация удобна и не требует оржон. восст. μ, α и β .

Пример

Блоки	Обработки			
	1	2	...	k
1	X_{11}	X_{12}		X_{1k}
2	X_{21}	X_{22}		X_{2k}
...				
n	X_{n1}	X_{n2}		X_{nk}

α_i - эффект блока

β_j - эффект обработки (он нас и интересует)

ε_{ij} - н.о.р. с непер. распредел., $\sum \varepsilon_{ij} = 0$.

МСПС | ANOVA

(3)

Проверка гипотезы
Критерий Фридмана

$H_0: \beta_1 = \dots = \beta_k$ против альтерн. H_1 : не все равны.

Его статистика $F = \frac{12}{nk(k+1)} \sum_{j=1}^k T_j^2 - 3n(k+1) = \frac{12n}{k(k+1)} \sum_{j=1}^k (\bar{R}_j - \bar{R})^2$
где R_{ij} = ранг X_{ij} в строке i , $\bar{R}_j = \frac{1}{n} \sum_{i=1}^n R_{ij}$; $\bar{R} = \frac{k+1}{2}$ - сред. ранг по матрице
 $T_j = n \bar{R}_j$.

- 1) При $n > 15, k > 4$ $F \sim \chi^2_{k-1}$.
- 2) Более точная: $\frac{(n-1)F}{n(k-1) - F} \sim F(n-1, (n-1)(k-1))$.

Следует ещё обратить внимание на крит. Неймана, он проверяет гипотезу $H_0: \beta_1 = \dots = \beta_k$ против альтерн. $H_1: \beta_1 \leq \dots \leq \beta_k$ (т.е. для такой альтернативы он является состоятельным)

Статистика: $L = \sum_{j=1}^k j R_j$.

Аппроксимация: $L^* = \frac{L - nk(k+1)^2/4}{\sqrt{\frac{1}{144(k-1)} n(k^3 - k)^2}} \sim N(0, 1)$, приближением можно пользоваться при $n > 15, k > 10$.

Замечание Критерии для проверки неоднородности 2 факторов в заданном факторе дисперсионного анализа: критерий Шеффе (для норм. факторов), критерий Даннета, LSD Фишера, MSD Тьюки, Кемпбелл (см. презентацию Рубенко).

Наконец, мы ещё не говорили о ситуации, когда факторы у нас авт. сведены. Моделируется эта ситуация так:
 $X_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$, где μ - "глобальное" среднее прироста X ,
 α_i - отклонение от μ , вызванное влиянием фактора i ,
 β_j - отклонение от μ и α_i , вызван. влиянием фактора j ,
 ε_{ij} - н.о.р. с. ошибки.

Фактор	1	j	k
Объект			
1	X_{11}		X_{1k}
2			
\vdots			
i			
\vdots			
n	X_{n1}		X_{nk}

Хотим понять, влияют ли факторы f_1, \dots, f_k на среднее объектов X_{ij} .

Как легко видеть, это очень похоже на двухфактор. дисперсионный анализ. Считать, что факторы сведены как-то по-другому, довольно сложно. Т.е. можно пользоваться крит. Фридмана и Неймана для проверки гипотезы $H_0: \beta_1 = \dots = \beta_k$, а также критерием Фишера: (против альтерн. H_1 : H_0 неверна)

$$F = \frac{n \cdot \sum_{j=1}^k (\bar{X}_j - \bar{X})^2 \cdot (n-1)}{\sum_{i=1}^n \sum_{j=1}^k (X_{ij} - \bar{X})^2 - k \sum_{i=1}^n (\bar{X}_i - \bar{X})^2}, \text{ где } \bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}; \bar{X}_i = \frac{1}{k} \sum_{j=1}^k X_{ij}; \bar{X} = \frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k X_{ij}$$

При верной H_0 $F \sim F(k-1, (n-1)(k-1))$

(лучше всего работает в предположении о нормальности данных и "сферичности" или "сферичности симметрии": $\forall j, X_{ij} - X_{kj}$ имеют одинаков. дисперсию, $\forall i \neq k$).