

Проверка гипотез однородности

Н.П. Алексеева,

СПбГУ, математико-механический факультет

2020 г.

Пусть имеется выборка наблюдений, которая разбивается на r групп, Каждая группа содержит n_i величин $\mathcal{N}(\mu_i, \sigma)$ с одинаковой дисперсией, $n_1 + \dots + n_r = n$, x_{ij} — j -я величина ($j = 1, \dots, n_i$) в i -й группе ($i = 1, \dots, r$). Модель

$$x_{ij} = \mu_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim \mathcal{N}(0, \sigma)$$

Нулевая гипотеза о равенстве внутригрупповых средних

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_r = \mu$$

Оценки параметров модели

- $\hat{\mu}_i = \bar{x}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}$ — среднее в i -й группе;
- $\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}$ — общее среднее.

Поскольку разбиение на группы обусловлено одним фактором (методика обучения), анализ называется однофакторным.

Пример данных для однофакторного дисперсионного анализа

Три группы водителей обучались по различным методикам. После окончания обучения был проведен контрольный тест. Получены следующие результаты:

Номер группы i	Число ошибок x_{ij}	Сумма	Число водителей n_i	Среднее в группе \bar{x}_i
1	1 3 2 1 0 2 1	10	7	$10/7 = 1.43$
2	2 3 2 1 4	12	5	$12/5 = 2.4$
3	4 5 3	12	3	$12/3 = 4$
		34	$n = 15$	

Общее среднее $\bar{x} = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij} = \frac{34}{15} = 2.27$.

Значимо ли различие между результатами разных методик?

Разложение суммы квадратов отклонений от выборочного среднего

- $Q = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2$ общий источник вариации
- $Q_1 = \sum_{i=1}^r \sum_{j=1}^{n_i} (\bar{x}_i - \bar{x})^2 = \sum_{i=1}^r n_i (\bar{x}_i - \bar{x})^2$ источник вариации, обусловленный различием между группами
- $Q_2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$ источник вариации, обусловленный различием внутри групп

$$\begin{aligned} Q &= \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i + \bar{x}_i - \bar{x})^2 = \\ &= \underbrace{\sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}_{Q_2} + \underbrace{\sum_{i=1}^r \sum_{j=1}^{n_i} (\bar{x}_i - \bar{x})^2}_{Q_1}, \text{ так как} \end{aligned}$$

$$\sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(\bar{x}_i - \bar{x}) = \sum_{i=1}^r (\bar{x}_i - \bar{x}) \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i) = 0.$$

Проверка гипотез в однофакторном дисперсионном анализе

При справедливости гипотезы $H_0 : \mu_1 = \dots = \mu_r = \mu$

- $Q \sim \sigma^2 \chi^2(n - 1)$
- $Q_1 \sim \sigma^2 \chi^2(r - 1)$
- $Q_2 \sim \sigma^2 \chi^2(n - r)$
- Q_1 и Q_2 независимы.

Для проверки гипотезы H_0 используется статистика

$$F = \frac{\frac{1}{r-1} Q_1}{\frac{1}{n-r} Q_2} \sim F(r - 1, n - r).$$

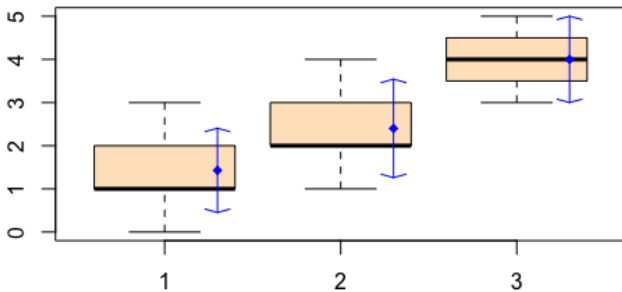
Пример вычислений

$$r = 3, n = 15$$

- $$Q = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij}^2 - 2x_{ij}\bar{x} + \bar{x}^2) =$$
$$= \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}^2 - 2\bar{x} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij} + n\bar{x}^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}^2 - n\bar{x}^2 =$$
$$= 104 - 15 \left(\frac{34}{15}\right)^2 \approx 26.93$$
- $$Q_1 = \sum_{i=1}^r n_i (\bar{x}_i - \bar{x})^2 = \sum_{i=1}^r n_i (\bar{x}_i^2 - 2\bar{x}_i\bar{x} + \bar{x}^2) = \sum_{i=1}^r n_i \bar{x}_i^2 - n\bar{x}^2 =$$
$$= 91.086 - 15 \left(\frac{34}{15}\right)^2 \approx 14.02$$
- $Q_2 = Q - Q_1 = 12.91$
- $F = \frac{Q_1/(r-1)}{Q_2/(n-r)} = \frac{14.02/2}{12.912/12} \approx 6.513$
- $\alpha_* = P\{F > 6.513\} = \text{FDIST}(6.513, 2, 12) = 0.012 < 0.05$

Методики дают значимо различные результаты тестового контроля.

Comparing boxplot(s) and non-robust mean +/- SD



Для проверки гипотез $H_0 : \mu_i = \mu_k$ значимости отклонений внутригрупповых средних используем статистики

$$\tau = \frac{(\bar{x}_i - \bar{x}_k) - (\mu_i - \mu_k)}{\sqrt{Q_2} \cdot \sqrt{\frac{1}{n_i} + \frac{1}{n_k}}} \sqrt{n-r} \sim T(n-r).$$

$$\xi_0 = \frac{(\bar{x}_i - \bar{x}_k) - (\mu_i - \mu_k)}{\sigma \sqrt{\frac{1}{n_i} + \frac{1}{n_k}}} \sim \mathcal{N}(0, 1), \quad \eta = \frac{Q_2}{\sigma^2} \sim \chi^2(n-r)$$

$$\tau = \frac{\xi_0}{\sqrt{\frac{1}{n-r}\eta}} = \frac{(\bar{x}_i - \bar{x}_k) - (\mu_i - \mu_k)}{\sigma \sqrt{\frac{1}{n_i} + \frac{1}{n_k}} \sqrt{\frac{Q_2}{\sigma^2}}} \sqrt{n-r} \sim T(n-r)$$

$$\begin{aligned}
 H_0 : \mu_1 = \mu_2, \quad \tau_{1-2} &= \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{Q_2} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sqrt{n - r} = \\
 &= \frac{(1.43 - 2.4)}{\sqrt{12.91} \cdot \sqrt{\frac{1}{7} + \frac{1}{5}}} \sqrt{15 - 3} = -1.6, \\
 p_{1-2} &= P\{\tau < -1.6\} = 0.068
 \end{aligned}$$

$$H_0 : \mu_1 = \mu_3, \quad \tau_{1-3} = -3.59, \quad p_{1-3} = P\{\tau < -3.59\} = 0.0018$$

$$H_0 : \mu_2 = \mu_3, \quad \tau_{2-3} = -2.11, \quad p_{2-3} = P\{\tau < -2.11\} = 0.028$$

Вычисление доверительных уровней вероятности

$$P\{\tau < -\tau_*\} = \text{TDIST}(\tau_*, \text{df}, 1),$$

$$P\{\tau > \tau_*\} = \text{TDIST}(\tau_*, \text{df}, 1), \quad \text{df} = n - r = 15 - 3 = 12$$

Поправки Бонферони, частный случай

Мы проверили три гипотезы о $\mu_1 = \mu_2$, $\mu_1 = \mu_3$, $\mu_2 = \mu_3$.

Число проверяемых гипотез обозначим через $m = 3$.

- Пусть α_0 вероятность ошибки первого рода (отвергнуть правильную гипотезу) в каждой из этих гипотез.
- Обозначим через $\tilde{\alpha}$ вероятность ошибиться хотя бы в одной из них, $\tilde{\alpha} = 1 - (1 - \alpha_0)^m$.
- Задача найти такое α_0 , что $\tilde{\alpha} = 1 - (1 - \alpha_0)^m < \alpha$.

$$\begin{aligned}(1 - \alpha_0)^m > 1 - \alpha &\iff 1 - \alpha_0 > (1 - \alpha)^{\frac{1}{m}} \iff \\ &\iff \alpha_0 < 1 - (1 - \alpha)^{\frac{1}{m}} = \alpha_*\end{aligned}$$

Например, при $m = 3$, $\alpha = 0.05$ имеем $\alpha_0 < \alpha_* = 0.0169$. Сравниваем $p_{1-2} = 0.068$, $p_{1-3} = 0.0018$, $p_{2-3} = 0.028$. В результате поправок получаем значимое различие только между 1 и 3 группами.

Поправки Бонферони, упрощение граничного значения

На практике предлагается вместо неравенства $\alpha_0 < 1 - (1 - \alpha)^{\frac{1}{m}}$ рассматривать $\alpha_0 < \frac{\alpha}{m}$. Например, $\frac{\alpha}{m} = 0.01667 < \alpha_* = 0.0169$. Поэтому если рассматривать значения $\alpha_0 < \frac{\alpha}{m}$, то из $\frac{\alpha}{m} < \alpha_* = 1 - (1 - \alpha)^{\frac{1}{m}}$ следует, не увеличится ошибка $\tilde{\alpha} = 1 - (1 - \alpha_0)^m < \alpha$. Осталось убедиться

$$\frac{\alpha}{m} < 1 - (1 - \alpha)^{\frac{1}{m}}$$

$1 - \frac{\alpha}{m} > (1 - \alpha)^{\frac{1}{m}} \iff (1 - \frac{\alpha}{m})^m > 1 - \alpha$, так как

$C_m^2 \frac{\alpha^2}{m^2} - C_m^3 \frac{\alpha^3}{m^3} + C_m^4 \frac{\alpha^4}{m^4} + \dots > 0$, например, обозначим $\alpha' = \frac{\alpha}{m}$ и при $m = 3$ имеем $3\alpha'^2 - \alpha'^3 = \alpha'^2(3 - \alpha') > 0$.

Пусть имеются нулевые гипотезы $\mathcal{H}_1, \dots, \mathcal{H}_m$ с соответствующими p -значениями: p_1, \dots, p_m . Например, $m = C_r^2$, если проверяются гипотезы о равенстве всех попарных средних. Определим $\tilde{\alpha}$ как вероятность совершить хотя бы одну ошибку в выводе из m выводов. Если осуществляется m независимых испытаний с вероятностью ошибки α_0 в каждом, то $\tilde{\alpha} = 1 - (1 - \alpha_0)^m$. В случае зависимых испытаний $\tilde{\alpha} \leq 1 - (1 - \alpha)^m$. Заметим, что по неравенству Буля

$$\tilde{\alpha} = 1 - (1 - m\alpha + C_m^2\alpha^2 - \dots) \leq m \cdot \alpha.$$

Таким образом, если мы хотим, чтобы вероятность наличия хотя бы одного неверного вывода из m была равна 0.05, то достаточно установить вероятность неверного отклонения нулевой гипотезы равной $0.05/m$ для каждого вывода.

Неравенство Буля: $P(A_1 + \dots + A_n) \leq P(A_1) + \dots + P(A_n)$ следует из определения вероятности суммы событий
 $P(A + B) = P(A) + P(B) - P(AB)$.

Задача на множественные сравнения

Проверим однородность данных о тяжести социальной зависимости $asi2$ emp у наркоманов в зависимости от фактора $educat$ образования: 1 - неполное среднее образование, 2 - среднее, 3 - неполное высшее, 4 - высшее. Нужно вычислить средние с ошибками среднего в каждой группе и заполнить таблицу соответствующими значимостями. Принять $\alpha = 0.1$. $Q_2 = 19.299$.

p	1	2	3	4
1				
2				
3				
4				

Задача на множественные сравнения, проверка

Проверим однородность данных о тяжести социальной зависимости asi_2 emr у наркоманов в зависимости от фактора educat образования: 1 - неполное среднее образование, 2 - среднее, 3 - неполное высшее, 4 - высшее. Заполнить таблицу соответствующими значимостями.
 $Q_2 = 19.299$, $n = 278$, $r = 4$.

$$\tau_{1-2} = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{Q_2} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sqrt{n - r}$$

p	1	2	3	4
1 $0.760 \pm 0.073(21)$		0.21	0.1	0.11
2 $0.807 \pm 0.017(221)$			0.0038	0.028
3 $0.662 \pm 0.062(27)$				0.39
4 $0.634 \pm 0.108(9)$				

$$\alpha_0 < \frac{0.05}{6} = 0.008.$$