

2. Двофакторний варіансний аналіз

Нехай дані про деяку мінливу величину поділяються на m груп за ознакою A і n груп за ознакою B . Одержимо mn класифікаційних підгруп. Припустимо, що для кожної підгрупи проводиться лише одне спостереження.

Позначимо через x_{ij} - спостереження в i -й групі за ознакою A , та в j -й групі за ознакою B . Тоді всі mn спостережень можна записати в наступній таблиці

$$\begin{array}{l}
 1 \dots \dots \dots j \dots \dots \dots n \\
 x_{11} \dots \dots \dots x_{1j} \dots \dots \dots x_{1n} \\
 x_{i1} \dots \dots \dots x_{ij} \dots \dots \dots x_{in} \\
 x_{m1} \dots \dots \dots x_{mj} \dots \dots \dots x_{mn}
 \end{array}
 \quad (*)$$

Позначимо через $x_{i\cdot}$ - середнє i -ої групи за ознакою A (i - рядка)

$$x_{i\cdot} = \frac{1}{n} \sum_{j=1}^n x_{ij}, \quad (i = 1, m)$$

через $x_{\cdot j}$ - середнє j -ої групи за ознакою B

$$x_{\cdot j} = \frac{1}{m} \sum_{i=1}^m x_{ij}, \quad (j = 1, n)$$

через $x_{\cdot\cdot}$ - загальне середнє всіх спостережень

$$x_{\cdot\cdot} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n x_{ij}$$

Повна мінливість всіх спостережень виражається девіацією

$$(3) \quad \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{\cdot\cdot})^2 = \sum_{i=1}^m \sum_{j=1}^n [(x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{\cdot\cdot}) + (x_{i\cdot} - x_{\cdot\cdot}) + (x_{\cdot j} - x_{\cdot\cdot})]^2 =$$

яку запишемо у вигляді

$$= n \sum_{i=1}^m (x_{i\cdot} - x_{\cdot\cdot})^2 + m \sum_{j=1}^n (x_{\cdot j} - x_{\cdot\cdot})^2 + \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{\cdot\cdot})^2$$

(мінливість між групами ознаки A або девіація, мінливість між групами ознаки B та залишкова мінливість). Тут три суми подвійних добутків дорівнюють нулю (за властивістю середнього арифметичного).

Таким чином повна мінливість розкладається на мінливість між групою A , мінливість між групою B та залишкову.

Кожна з цих дивіацій має своє число $d.f.$ (ступенів вільності)

$$mn - 1 \quad m - 1 \quad n - 1 \quad mn - (m + n - 1) = (m - 1)(n - 1)$$

Сума чисел ступенів вільності справа = числу ступенів вільності зліва. Якщо тотожність (3) поділити на $(mn - 1)$, то одержимо, що повна варіанса є опуклою лінійною комбінацією варіанс між групами ознак A , між групами ознаки B та залишковою варіансою.

Варіанси справа позначимо через S_A^2, S_B^2, S_r^2 .
Тоді три варіанси:

$$S_A^2 = \frac{1}{m-1} n \sum_{i=1}^m (x_{i\cdot} - x_{..})^2$$

$$S_B^2 = \frac{1}{n-1} \cdot m \sum_{j=1}^n (x_{\cdot j} - x_{..})^2$$

$$S_r^2 = \frac{1}{(m-1)(n-1)} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{..})^2$$

$$\frac{1}{mn-1} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{..})^2 = \frac{m-1}{mn-1} n \cdot \frac{1}{m-1} \sum_{i=1}^m (x_{i\cdot} - x_{..})^2 + \frac{n-1}{mn-1} m \frac{1}{n-1} \sum_{j=1}^n (x_{\cdot j} - x_{..})^2$$

$$+ \frac{(m-1)(n-1)}{mn-1} \cdot \frac{1}{m-1} \cdot \frac{1}{n-1} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{..})^2$$

Якщо припустити, що спостереження однорідні і взяті з нормальної генеральної сукупності, то варіанси S_A^2, S_B^2, S_r^2 є незалежними оцінками дисперсії генеральної сукупності (тобто справа кожна з варіанс незалежна від двох інших)

Звідси слідує для перевірки H : однорідності можна вибрати статистику Фішера.

Формулюємо паралельно дві гіпотези:

H_A : коли вплив груп ознаки A не істотний ($= 0$)

Два доведення гіпотези вибираємо статистику

$$(1^*) \quad F_A = \frac{S_A^2}{S_r^2}$$

тобто гіпотезу H_A : доводимо незалежно від впливу груп ознак B , а гіпотезу H_B - незалежно від впливу груп ознаки A .

H_B : коли вплив груп ознаки B не істотний ($= 0$)

Для доведення цієї гіпотези вибираємо статистику

$$(2^*) \quad F_B = \frac{S_B^2}{S_r^2}$$

Якщо гіпотеза про однорідність даних табл. (*) з генеральної популяції вірна, то статистики $(1^*), (2^*)$ мають розподіл Фішера відповідно з $d.f. = (m-1, (m-1)(n-1))$ та $d.f. = (n-1, (m-1)(n-1))$. Це дозволяє визначити критичні значення для обидвох гіпотез.

Обчислення при двофакторному варіансному аналізі оформляємо у вигляді таблиці при одному спостереженні в кожній підгрупі

Мінливість	Девіація	$d.f.$	Варіанса
------------	----------	--------	----------

між групами A	$n \sum_{i=1}^m (x_{i\bullet} - x_{\bullet\bullet})^2$	$m - 1$	S_A^2
між групами B	$m \sum_{j=1}^n (x_{\bullet j} - x_{\bullet\bullet})^2$	$n - 1$	S_B^2
Залишкова	$\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{i\bullet} - x_{\bullet j} + x_{\bullet\bullet})^2$	$(m - 1)(n - 1)$	S_r^2
Повна	$\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - x_{\bullet\bullet})^2$	$mn - 1$	-

Останній рядок є сумою трьох попередніх, а це служить контролем правильності обчислень

Приклад. Затрати матеріалу на виготовлення деякого виробу трьома різними технологіями (A) на 4-х різних заводах (B) були такі:

B A	1	2	3	4
1	25	20	30	25
2	30	40	40	50
3	23	18	20	27

При рівні значущості $\alpha = 0,10$ перевірити гіпотезу про те, що рівень затрат матеріалу на виріб не впливає ані на вибір технології ані на вибір заводу.

H_A : вплив технології $A = 0$

H_B : вплив заводу $B = 0$

Для перевірки H проводимо варіансний аналіз

$$m = 3$$

$$n = 4$$

$$x_{1\bullet} = 25 \quad x_{1\bullet} = \frac{1}{4} \sum_{i=1}^4 x_i \quad x_{\bullet 1} = 26$$

$$x_{2\bullet} = 40 \quad x_{\bullet 2} = 26$$

$$x_{3\bullet} = 22 \quad x_{\bullet 3} = 30$$

$$x_{\bullet 4} = 34$$

$$x_{\bullet\bullet} = 29$$

Результати

Мінливість	Девіація	$d.f.$	Варіанса
між технологіями A	744	2	$S_A^2 = 372$
між заводами B	132	3	$S_B^2 = 44$
Залишкова	164	6	$S_r^2 = 27,33$
Повна	1040	11	-

$$F_{Aem} = \frac{S_A^2}{S_r^2} = \frac{372}{27,33} = 13,61 \quad F_{Aem} \square F_{Акк}$$

H_A – відкидаємо. Тип технології істотно впливає на рівень затрат матеріалу при виготовлені.

$$\alpha = 0,10 \quad \left| \quad d.f. = (2, 6) \right| \quad F_{Акк} = 4,76 \quad F_{Bem} < F_{Акк}$$

H_B – приймаємо.

Заводи не впливають істотно на рівень затрат матеріалу.