

III. **Множинний t -метод.** Нехай M – кількість *a priori* обраних для аналізу контрастів. (Зазвичай контрасти вибирають для дослідження після експерименту.) Тоді наближена довірча область для множини цих вибраних контрастів $\left\{ \sum_{i=1}^I c_i^{(j)} a_i \right\}_{j=1}^M$ з рівнем довіри, не менше ніж $(1-\gamma)$, $\gamma > 0$, задається системою довірчих інтервалів для кожного контрасту $\sum_{i=1}^I c_i^{(j)} a_i$ з рівнем довіри $\left(1 - \frac{\gamma}{M}\right)$, $\gamma > 0$ і визначається таким чином:

$$\left| \sum_{i=1}^I c_i^{(j)} a_i - \sum_{i=1}^I c_i^{(j)} \bar{y}_i \right| \leq \sqrt{\bar{S}_e \left(\sum_{i=1}^I \frac{(c_i^{(j)})^2}{N_i} \right)} t_{\frac{\gamma}{2M}} (N-I) = \Delta_3^j, j = \overline{1, M}, \quad (18)$$

де $t_\gamma(v)$ – 100γ відсоткова точка t -розподілу Стюдента з v ступенями свободи.

рівнем довіри $(1-\gamma)$, $\gamma > 0$, в результаті чого отримали (для методів Шеффе та Тьюкі) відповідний довірчий інтервал виду $[\hat{k} - \Delta_i, \hat{k} + \Delta_i]$, $\Delta_i > 0, i = \overline{1, 2}$, де $\hat{k} = \sum_{i=1}^I c_i \bar{y}_i$.

Це дозволяє на другому кроці області прийняття гіпотези (13) для цих методів записати в такому вигляді:

$$(\hat{k} - \Delta_i)(\hat{k} + \Delta_i) \leq 0 \Leftrightarrow \hat{k}^2 \leq \Delta_i^2 \Leftrightarrow |\hat{k}| \leq \Delta_i, i = 1, 2,$$

тобто справедливості останньої нерівності означає, що відповідний контраст слід вважати таким, що незначимо відхиляється від нуля із

статистичної точки зору з рівнем значущості $\gamma > 0$, у протилежному випадку його треба вважати таким, що істотно відхиляється від нуля.

Остаточно область прийняття відповідної гіпотези матиме такий вигляд:

I. для методу Шеффе з урахуванням (15) одержуємо

$$\left| \sum_{i=1}^I c_i \bar{y}_i \right| \leq \sqrt{\bar{S}_e \left(\sum_{i=1}^I \frac{c_i^2}{N_i} \right) (I-1) F_\gamma (I-1, N-I)} = \Delta_1,$$

II. для методу Тьюкі, якщо взяти до уваги (17), отримуємо

$$\left| \sum_{i=1}^I c_i \bar{y}_i \right| \leq \frac{1}{2} \sum_{i=1}^I |c_i| \sqrt{\frac{\bar{S}_e}{N_0}} q_\gamma (I, N-I) = \Delta_2,$$

III. для множинного t -методу врахування (19) дозволяє записати

$$\left| \sum_{i=1}^I c_i^{(j)} \bar{y}_i \right| \leq \sqrt{\bar{S}_e \left(\sum_{i=1}^I \frac{(c_i^{(j)})^2}{N_i} \right)} t_{\frac{\gamma}{2M}} (N-I) = \Delta_3^j, j = \overline{1, M}.$$

η - залежна кількісна скалярна змінна,

а незалежні якісні скалярні змінні:

ζ_1 - фактор A , який набуває своїх значень з I_1 градацій,

ζ_2^I - фактор B , який набуває своїх значень з I_2 градацій.

Необхідно за спостереженнями над залежною змінною η при активних різних сполученнях градацій незалежних змінних ζ_1 та ζ_2 побудувати математичну модель залежності змінної η від змінних ζ_1 та ζ_2 .

Фон: Приклад 2: η – врожайність зернової культури,

ζ_1 – сорт зернової культури, всього I_1 сортів,

ζ_2 – вид добрива, всього I_2 видів добрива.

Нехай при активному сполученні i -ї градації змінної ζ_1 та j -ї градації змінної ζ_2 доступно N_{ij} спостережень y_{ijk} над η , $i = \overline{1, I_1}, j = \overline{1, I_2}, k = \overline{1, N_{ij}}$ ($N_{ij} \geq 1$). Тоді математичну модель двофакторного дисперсійного аналізу будемо шукати в такому вигляді:

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}, \quad i = \overline{1, I_1}, j = \overline{1, I_2}, k = \overline{1, N_{ij}} \quad (N_{ij} \geq 1), \quad (20)$$

де

y_{ijk} – k -те спостереження над η при активному сполученні i -ї градації змінної ζ_1 та j -ї градації змінної ζ_2 ,

μ – загальне середнє всіх спостережень у деякому розумінні,

α_i – кількісний вираз відносного впливу i -ї градації змінної ζ_1 на η відносно μ (або, іншими словами, головний ефект i -го рівня фактора A),

β_j – кількісний вираз відносного впливу j -ї градації змінної ζ_2 на η відносно μ (або, іншими словами, головний ефект j -го рівня фактора B),

γ_{ij} – кількісний вираз відносного впливу взаємодії i -ї градації змінної ζ_1 та j -ї градації змінної ζ_2 на η відносно μ (або, іншими словами, взаємодія i -го рівня фактора A та j -го рівня фактора B),

e_{ijk} – похибка моделі k -го спостереження над η при активному сполученні i -ї градації змінної ζ_1 та j -ї градації змінної ζ_2 .

Причому всього спостережень доступно в кількості

$$N = \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} N_{ij}.$$

У свою чергу, кількісний вираз абсолютного впливу i -ї градації змінної ζ_1 на η дорівнює

$$a_i = \mu + \alpha_i, \quad i = \overline{1, I_1},$$

кількісний вираз абсолютного впливу j -ї градації змінної ζ_2 на η має вигляд

$$b_j = \mu + \beta_j, \quad j = \overline{1, I_2},$$

а кількісний вираз абсолютного впливу взаємодії i -ї градації змінної ζ_1 та j -ї градації змінної ζ_2 на η визначається як

$$c_{ij} = \mu + \gamma_{ij}, \quad i = \overline{1, I_1}, j = \overline{1, I_2}.$$

Припустимо, що похибки e_{ijk} моделі (20) є:

- $e_{ijk} \sim \mathcal{N}(0, \sigma^2), \sigma^2 > 0, \quad \forall i, j, k;$
- $\{e_{ijk}\}$ - незалежні.

Необхідно за доступними скалярними спостереженнями $\{y_{ijk}, i = \overline{1, I_1}, j = \overline{1, I_2}, k = \overline{1, N_{ij}} (N_{ij} \geq 1)\}$ знайти оцінки невідомих параметрів:

$$\begin{aligned} &\mu, \\ &\alpha_1, \alpha_2, \dots, \alpha_{I_1}, \\ &\beta_1, \beta_2, \dots, \beta_{I_2}, \\ &\gamma_{11}, \gamma_{12}, \dots, \gamma_{1I_2}, \\ &\gamma_{21}, \gamma_{22}, \dots, \gamma_{2I_2}, \\ &\dots \\ &\gamma_{I_11}, \gamma_{I_12}, \dots, \gamma_{I_1I_2} \end{aligned} \tag{21}$$

математичної моделі (20). Шукаємо за допомогою МНК. Загальна кількість невідомих параметрів буде дорівнювати

$$1 + I_1 + I_2 + I_1 I_2 = (1 + I_1)(1 + I_2).$$

Далі схема розв'язання задачі двофакторного дисперсійного аналізу повністю аналогічна процедурі розв'язання задачі однофакторного дисперсійного аналізу.

Спочатку модель (20) переписується, як і раніше, а саме:

$$y = X\alpha + e,$$

причому

y – вектор-стовпчик з усіх спостережень y_{ijk} ,

α – вектор-стовпчик з усіх невідомих параметрів (21),

X – матриця відповідної розмірності, елементи кожного рядка якої всі дорівнюють нулю, окрім першого та трьох інших, що відповідають місцезнаходженню відповідних параметрів головних ефектів та попарної взаємодії у векторі α ,

e – вектор-стовпчик з усіх похибок моделі

$$\{e_{ijk}, i = \overline{1, I_1}, j = \overline{1, I_2}, k = \overline{1, N_{ij}} (N_{ij} \geq 1)\}.$$

У цьому випадку матриця X буде мати неповний ранг. Тому, щоб скористатися МНК при визначенні оцінок вектора α , необхідно врахувати додаткові лінійні обмеження, які справедливі для нього.

Дійсно, враховуючи зміст невідомих параметрів (21), можна стверджувати, що:

$$\exists \{v_i\}_{i=1}^{I_1}, \{w_j\}_{j=1}^{I_2} : \forall i \ v_i > 0, \forall j \ w_j > 0,$$

$$\begin{cases} \sum_{i=1}^{I_1} v_i \alpha_i = 0, \\ \sum_{j=1}^{I_2} w_j \beta_j = 0, \\ \sum_{i=1}^{I_1} v_i \gamma_{ij} = 0, \quad j = \overline{1, I_2}, \\ \sum_{j=1}^{I_2} w_j \gamma_{ij} = 0, \quad i = \overline{1, I_1}. \end{cases} \quad (22)$$

Визначення вагових коефіцієнтів $\{v_i\}_{i=1}^{I_1}$ та $\{w_j\}_{j=1}^{I_2}$ здійснюється відповідно до змісту конкретної постановки задачі.

Врахування лінійних обмежень (22) дозволяє однозначно визначити оцінку $\hat{\alpha}$ методом найменших квадратів у математичній моделі (20) за спостереженнями $\{y_{ijk}, i = \overline{1, I_1}, j = \overline{1, I_2}, k = \overline{1, N_{ij}} (N_{ij} \geq 1)\}$.

Окрім цього, цікавою є перевірка на значимість параметрів моделі двофакторного дисперсійного аналізу, і насамперед – перевірка з деяким рівнем значущості $\gamma > 0$ таких гіпотез:

$$H_0^A : \alpha_1 = \alpha_2 = \dots = \alpha_{I_1} = 0, \quad (23)$$

$$H_0^B : \beta_1 = \beta_2 = \dots = \beta_{I_2} = 0, \quad (24)$$

$$H_0^{AB} : \gamma_{ij} = 0, \quad i = \overline{1, I_1}, j = \overline{1, I_2}. \quad (25)$$

Для розв'язання цих задач достатньо використати той самий математичний апарат, що й при розв'язанні відповідних задач в однофакторному дисперсійному аналізі. Проте формули стануть більш громіздкими. Для їх спрощення наведемо розв'язок для випадку, коли справедливо

$$N_{ij} = N_0 (N_0 \geq 1), \quad i = \overline{1, I_1}, j = \overline{1, I_2}, \quad (26)$$

а відповідні вагові коефіцієнти $\{v_i\}_{i=1}^{I_1}, \{w_j\}_{j=1}^{I_2}$ усі однакові

$$v_i = \frac{1}{I_1}, \quad i = \overline{1, I_1}; \quad w_j = \frac{1}{I_2}, \quad j = \overline{1, I_2}. \quad (27)$$

Тепер загальна кількість спостережень визначатиметься таким чином:

$$N = I_1 I_2 N_0.$$

У результаті використання методу найменших квадратів для визначення невідомих параметрів (21) у математичній моделі (20) за наявності лінійних обмежень (22) при справедливості припущень (26), (27) отримуємо такі їх оцінки:

$$\hat{\mu} = \bar{y},$$

$$\hat{\alpha}_i = \bar{y}_{i..} - \bar{y}, \quad i = \overline{1, I_1},$$

$$\hat{\beta}_j = \bar{y}_{.j.} - \bar{y}, \quad j = \overline{1, I_2},$$

$$\hat{\gamma}_{ij} = \bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}, \quad i = \overline{1, I_1}, j = \overline{1, I_2},$$

де

$$\bar{y} = \frac{1}{N} \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} y_{ijk},$$

$$\bar{y}_{ij.} = \frac{1}{N_0} \sum_{k=1}^{N_0} y_{ijk}, \quad i = \overline{1, I_1}, j = \overline{1, I_2},$$

$$\bar{y}_{i..} = \frac{1}{I_2 N_0} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} y_{ijk}, \quad i = \overline{1, I_1},$$

$$\bar{y}_{.j.} = \frac{1}{I_1 N_0} \sum_{i=1}^{I_1} \sum_{k=1}^{N_0} y_{ijk}, \quad j = \overline{1, I_2}.$$

Таблиця двофакторного дисперсійного аналізу

Проведемо аналіз повної суми квадратів відхилень спостережень y_{ijk} від загального середнього \bar{y} . Дійсно,

$$\begin{aligned} \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} (y_{ijk} - \bar{y})^2 &= \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} \left[(y_{ijk} - \bar{y}_{ij.}) + \right. \\ &\quad \left. + (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}) + (\bar{y}_{i..} - \bar{y}) + (\bar{y}_{.j.} - \bar{y}) \right]^2 = \\ &= \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} (y_{ijk} - \bar{y}_{ij.})^2 + \\ &\quad + N_0 \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y})^2 + \quad \text{(на с/р)} \\ &\quad + I_2 N_0 \sum_{i=1}^{I_1} (\bar{y}_{i..} - \bar{y})^2 + I_1 N_0 \sum_{j=1}^{I_2} (\bar{y}_{.j.} - \bar{y})^2. \end{aligned}$$

Останнє перетворення справедливо в силу того, що в усіх подвійних добутках квадратні дужки дорівнюють нулеві. Отриманий результат скорочено можна записати таким чином:

$$S = S_e + S_A + S_B + S_{AB}, \quad (28)$$

де

$$S = \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} (y_{ijk} - \bar{y})^2, S_e = \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} \sum_{k=1}^{N_0} (y_{ijk} - \bar{y}_{ij.})^2,$$

$$S_A = I_2 N_0 \sum_{i=1}^{I_1} (\bar{y}_{i..} - \bar{y})^2, S_B = I_1 N_0 \sum_{j=1}^{I_2} (\bar{y}_{.j.} - \bar{y})^2,$$

$$S_{AB} = N_0 \sum_{i=1}^{I_1} \sum_{j=1}^{I_2} (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y})^2.$$

Отже, для S , повної суми квадратів відхилень спостережень y_{ijk} від загального середнього \bar{y} , у випадку двофакторної моделі дисперсійного аналізу отримали розклад (28), аналогічний розкладу (12), отриманому для однофакторної моделі.

Результати двофакторного дисперсійного аналізу також заносять у відповідну таблицю 2.1 двофакторного дисперсійного аналізу.

Таблиця 2.1 двофакторного дисперсійного аналізу

Джерело варіації	Сума квадратів	Кількість ступенів свободи	Середня сума квадратів	F- статистика	γ_{\max}
головні ефекти фактора A	S_A	$I_1 - 1$	$\bar{S}_A = \frac{S_A}{I_1 - 1}$	$F_A = \frac{\bar{S}_A}{\bar{S}_e}$	γ_A
головні ефекти фактора B	S_B	$I_2 - 1$	$\bar{S}_B = \frac{S_B}{I_2 - 1}$	$F_B = \frac{\bar{S}_B}{\bar{S}_e}$	γ_B
взаємодії факторів A та B	S_{AB}	$(I_1 - 1)(I_2 - 1)$	$\bar{S}_{AB} = \frac{S_{AB}}{(I_1 - 1)(I_2 - 1)}$	$F_{AB} = \frac{\bar{S}_{AB}}{\bar{S}_e}$	γ_{AB}
ПОМИЛКИ	S_e	$I_1 I_2 (N_0 - 1)$	$\bar{S}_e = \frac{S_e}{I_1 I_2 (N_0 - 1)}$		
	S	$N - 1$			

В останньому рядку табл. 2.1, таблиці двофакторного дисперсійного аналізу, наведено суми по другому та третьому стовпчикам відповідно. Зауважимо, що результат у другому стовпчику збігається з уже отриманим результатом (28). Підраховані в таблиці значення $\gamma_A, \gamma_B, \gamma_{AB}$ – значення максимальних рівнів значущості γ , за яких гіпотези (23), (24), (25) будуть справедливі, а використання статистик F_A, F_B, F_{AB} дозволяє записати області прийняття гіпотез H_0^A, H_0^B, H_0^{AB} , відповідно:

для гіпотези H_0^A :

$$F_A < F_\gamma(I_1 - 1, I_1 I_2 (N_0 - 1)),$$

для гіпотези H_0^B :

$$F_B < F_\gamma(I_2 - 1, I_1 I_2 (N_0 - 1)),$$

для гіпотези H_0^{AB} :

$$F_{AB} < F_\gamma((I_1 - 1)(I_2 - 1), I_1 I_2 (N_0 - 1)),$$

де $F_\gamma(v_1, v_2)$ – 100 γ відсоткова точка F -розподілу з параметрами v_1 та v_2 .

Багатофакторний дисперсійний аналіз

...

Самостійна робота №6. З навчального посібника
«Слабоспицький О.С. Дисперсійний аналіз даних, 2013»
пропрацювати матеріал наведений у Розділі 3.:
«Багатофакторний дисперсійний аналіз».

I

(пропустити 5 стор.)

Коваріаційний аналіз

Приклад 1. Залежна змінна η – **кількісний** показник ризику зараження та важкість перебігу COVID-19,
незалежна **якісна** змінна ζ_1 – група крові пацієнта,
незалежна **якісна** змінна ζ_2 – раса пацієнта,
незалежна **якісна** змінна ζ_3 – стать пацієнта,
незалежна **кількісна** змінна ξ_1 – вік,
незалежна **кількісна** змінна ξ_2 – рівень цукру у крові,
незалежна **кількісна** змінна ξ_3 – середня кількість викурених сигарет за добу,

Файл | C:\Users\Александр\Рабочий стол\Коваріаційний аналіз (завантаження).ppt

Коваріаційний аналіз

Коваріаційний аналіз – розділ аналізу даних, який займається побудовою математичних моделей поточних зв'язків між залежною кількісною змінною η та вектором незалежних якісних змінних ζ^T та вектором незалежних кількісних змінних ξ^T . Цей розділ аналізу даних в об'єднанні регресійного та дисперсійного аналізу даних.

Товарова група

Нехай \vec{x}_g має розмірність q , а \vec{x} має розмірність p , k -те спостереження над y є $y(k)$, тоді поточною моделю коваріаційного аналізу можна наступним чином:

$$y(k) = x_g^T(k) \alpha_g + x^T(k) \alpha + e(k), k = \overline{1, N} \quad (1)$$

Зауваження. α має розмірність p , що дорівнює кількості регресорів, а α_g ("амбда-де") має розмірність q , що дорівнює кількості невідомих параметрів дисперсійної моделі.
 $\alpha, e(k)$ можна зал. від $x_g(k)$

$$y = X_g \alpha_g + X \alpha + e \quad (2)$$

$$y = \begin{pmatrix} y(1) \\ \vdots \\ y(N) \end{pmatrix}, X_g = \begin{pmatrix} x_g^T(1) \\ \vdots \\ x_g^T(N) \end{pmatrix}, X = \begin{pmatrix} x^T(1) \\ \vdots \\ x^T(N) \end{pmatrix}, e = \begin{pmatrix} e(1) \\ \vdots \\ e(N) \end{pmatrix}$$

Класичний коваріаційний аналіз

Припущення класичного коваріаційного аналізу:

$$1. e \sim N(0, \sigma^2 E_N), \sigma^2 > 0$$

$$2. \text{rang}(X_g) = q \text{ (враховані лінійні обмеження)}$$

$$\text{rang}(X) = p$$

$$3. \alpha_g \in \mathbb{R}^q, \alpha \in \mathbb{R}^p$$

$$4. \text{Матриця } X \text{ не залежить від матриці } X_g$$

(3)

Двухэтапный метод наименьших квадратов

I шаг:

$$y - X\alpha = Xg\alpha_g + e$$

Итак $\alpha = \theta_p$

$$\hat{\alpha}_g(\theta) = (X_g^T X_g)^{-1} X_g^T y$$

S_e - суммарная сумма квадратов

$$S_e(\theta) = \|y - X_g \hat{\alpha}_g\|^2 = \|y - X_g (X_g^T X_g)^{-1} X_g^T y\|^2$$

$$Q = E - X_g (X_g^T X_g)^{-1} X_g^T$$

Q - симметрична, идемпотентна (на \mathbb{C}/\mathbb{R}).

$$\begin{aligned} \Rightarrow \| (E - X_g (X_g^T X_g)^{-1} X_g^T) y \|^2 &= y^T Q^T Q y = y^T Q^2 y = \\ &= y^T Q y = \|y\|_Q^2 \end{aligned}$$

II шаг:

$$y \rightarrow y - X\alpha$$

$$S_e(\alpha) = (y - X\alpha)^T Q (y - X\alpha)$$

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} S_e(\alpha) = (X^T Q X)^{-1} X^T Q y \quad (4)$$

$$\hat{\alpha}_g(\hat{\alpha}) = (X_g^T X_g)^{-1} X_g (y - X\hat{\alpha}) \quad (5)$$

$$S_e(\hat{\alpha}) = (y - X\hat{\alpha})^T Q (y - X\hat{\alpha})$$