

Процедура використання  $\hat{\tau}_{ij}^{(K)}$ :

- 1) якщо  $\hat{\tau}_{ij}^{(K)} = 0$ , то вважаємо, що зв'язок між змінними  $\xi_i$  та  $\xi_j$  відсутній;
- 2) якщо  $|\hat{\tau}_{ij}^{(K)}| = 1$ , то зв'язок між змінними  $\xi_i$  та  $\xi_j$  функціональний причому:
  - якщо  $\hat{\tau}_{ij}^{(K)} = 1$ , то ранжування  $x^{(i)}$  та  $x^{(j)}$  рівні;
  - якщо  $\hat{\tau}_{ij}^{(K)} = -1$ , то ранжування  $x^{(i)}$  та  $x^{(j)}$  протилежні;
- 3) якщо  $|\hat{\tau}_{ij}^{(K)}| \in (0,1)$ , то потрібно звернутися до перевірки на значимість рангового коефіцієнта кореляції Кендела, тобто перевірити гіпотезу

$$H_0 : \tau_{ij}^{(K)} = 0$$

з деяким рівнем значущості  $\alpha > 0$ . Зв'язок між ординальними змінними  $\xi_i$  та  $\xi_j$  зі статистичної точки зору будемо вважати суттєвим при відхиленні гіпотези  $H_0$ , а в протилежному разі - не істотним.

При  $n = 4, 10$ , таку перевірку можна здійснити за допомогою спеціальних таблиць.

А при  $n \geq 11$ , з'ясувалося, що розподіл статистики

$$3\hat{\tau}_{ij}^{(K)} \sqrt{\frac{n(n-1)}{2(2n+5)}}$$

може бути наближений стандартним нормальним розподілом, тоді враховуючи, що до критичної області  $H_0$  потрібно віднести області екстремальних значень останньої статистики, то область прийняття гіпотези буде мати таке представлення:

$$\left| 3\hat{\tau}_{ij}^{(K)} \sqrt{\frac{n(n-1)}{2(2n+5)}} \right| < u_{\frac{\alpha}{2}},$$

де  $u_\beta$  – 100 $\beta$  відсоткова точка нормального розподілу з параметрами 0 і 1.

**II випадок.** За наявності груп об'єктів з однаковим проявом принаймні за однією зі змінних  $\xi_i$  чи  $\xi_j$ , ранговий коефіцієнт кореляції Кендела  $\hat{\tau}_{ij}^{(K)}$  корегується й використовується у вигляді **вибіркового модифікованого рангового коефіцієнта кореляції Кендела**:

$$\hat{\tau}_{ij}^{(K)} = \frac{\hat{\tau}_{ij}^{(K)} - (\delta^{(i)} + \delta^{(j)})}{\sqrt{(1 - \delta^{(i)})(1 - \delta^{(j)})}},$$

де  $m^{(k)}$  – кількість груп об'єктів з однаковим проявом за змінною  $\xi_k$ ;  $n_l^{(k)}$  – кількість об'єктів, які увійшли в  $l$ -ту групу об'єктів з однаковим проявом за змінною  $\xi_k$ ;

$$\delta^{(k)} = \frac{\frac{1}{4} \sum_{l=1}^{m^{(k)}} n_l^{(k)} (n_l^{(k)} - 1)}{\frac{1}{4} n(n-1)}; \quad k = i, j.$$

**Зауваження 1.** Для рангового коефіцієнта кореляції Кендела, коли відсутні групи об'єктів з однаковим проявом за змінними  $\xi_i$  та  $\xi_j$ , корегуючі величини  $\delta^{(i)} = \delta^{(j)} = 0$ , оскільки  $m^{(i)} = m^{(j)} = n$ , а  $n_l^{(i)} = n_l^{(j)} = 1, l = \overline{1, n}$ . У цій ситуації модифікований коефіцієнт  $\hat{\tau}_{ij}^{(K)} = \hat{\tau}_{ij}^{(K)}$ .

У підсумку, для аналізу парного статистичного зв'язку ординальних змінних маємо ще одну характеристику – ранговий коефіцієнт кореляції Кендела.

**Зауваження 2.** Рангові коефіцієнти кореляції Спірмена та Кендела зв'язані між собою. Так, коли їх модулі не дуже близькі до одиниці та  $n$  вже велике, для них має місце таке наближення  $\hat{\tau}_{ij}^{(S)} \approx \frac{3}{2} \hat{\tau}_{ij}^{(K)}$ .

## Аналіз множинних зв'язків для ординальних змінних. Коефіцієнт конкордації.

Нехай серед усіх ординальних змінних  $\eta, \xi_1, \xi_2, \dots, \xi_q$  ( $\xi_0 \equiv \eta$ ) відібрано  $m$  змінних:  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ ,  $2 \leq m \leq q+1$ . Перейдемо до аналізу наявності статистичного зв'язку між цими вибраними ординальними змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ , причому нехай для них доступні відповідні ранжування

$$x^{(ij)} = \left( x_1^{(ij)}, x_2^{(ij)}, \dots, x_n^{(ij)} \right)^T, \quad j = \overline{1, m}, \quad 2 \leq m \leq q+1,$$

де  $n$  – кількість об'єктів, що досліджуються.

Для розв'язання цієї проблеми М. Кендел запропонував використовувати спеціальний показник.

**I випадок.** Розглянемо спочатку випадок, коли відсутні групи об'єктів з однаковим проявом у змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ ,  $2 \leq m \leq q+1$ .

Введемо позначення  $\vec{\zeta} = (\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m})^T$ .

**Означення.** Вибірковим коефіцієнтом конкордації ординальних змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  з ранжуваннями  $x^{(i_1)}, x^{(i_2)}, \dots, x^{(i_m)}$  називається числова характеристика, що задається таким чином:

$$\hat{W}_{\vec{\zeta}} = \frac{\sum_{k=1}^n \left( \sum_{j=1}^m x_k^{(i_j)} - \frac{m(n+1)}{2} \right)^2}{m^2(n^3 - n)},$$

де  $x^{(ij)} = \left( x_1^{(ij)}, x_2^{(ij)}, \dots, x_n^{(ij)} \right)^T, \quad j = \overline{1, m}, \quad 2 \leq m \leq q+1$ .



Для коефіцієнта конкордації мають місце такі властивості:

- 1)  $0 \leq \hat{W}_{\xi} \leq 1$ ;
- 2) якщо  $\hat{W}_{\xi} = 0$ , то зв'язок між змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  відсутній;
- 3) якщо  $\hat{W}_{\xi} = 1$ , то зв'язок функціональний, а саме ранжування  $x^{(i_1)}, x^{(i_2)}, \dots, x^{(i_m)}$  будуть рівні;
- 4) зміна порядку розташування змінних у послідовності  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  не змінює значення коефіцієнта конкордації.

У загальному випадку, коли кількість об'єктів  $n$  не є постійною, для вибіркового значення коефіцієнта конкордації, яке вже може залежати від  $n$ , будемо використовувати позначення  $\hat{W}_{\xi}(n)$ .

В свою чергу введемо поняття.

**Означення.** Теоретичним коефіцієнтом конкордації ординальних змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  називається числова характеристика  $W_{\xi}$ , яка визначається таким чином:

$$W_{\xi} = \lim_{n \rightarrow \infty} \hat{W}_{\xi}(\bar{n}).$$

Процедура використання  $W_{\xi}$ :

- 1) якщо  $\hat{W}_{\xi} = 0$ , то зв'язок між змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  відсутній;
- 2) якщо  $\hat{W}_{\xi} = 1$ , то зв'язок функціональний, а саме ранжування  $x^{(i_1)}, x^{(i_2)}, \dots, x^{(i_m)}$  будуть рівні;
- 3) якщо  $\hat{W}_{\xi} \in (0,1)$ , то потрібно з'ясувати, чи суттєво відхиляється від нуля коефіцієнт конкордації зі статистичної точки зору, тобто здійснити перевірку його на значимість, а саме перевірити гіпотезу

$$H_0 : W_{\xi} = 0$$

з деяким рівнем значущості  $\alpha > 0$ . Якщо гіпотеза буде відхилена, то вважається, що зв'язок між ординальними змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  зі статистичної точки зору є істотним, у протилежному випадку – не суттєвим.

При  $n = \overline{3, 7}$ ,  $m = \overline{2, 20}$ , перевірка гіпотези  $H_0$  здійснюється за допомогою спеціальних таблиць.

А при  $n \geq 8$ , можна скористатися тим фактом, що розподіл статистики  $m(n-1)\hat{W}_{\xi}$  можна наблизити  $\chi^2$ -розподілом з  $(n-1)$  ступенями свободи за умови справедливості гіпотези  $H_0$ . У результаті область прийняття гіпотези набуває вигляду:

$$m(n-1)\hat{W}_{\xi} < \chi_{\alpha}^2(n-1),$$

де  $\chi_{\alpha}^2(n)$  –  $100\alpha$  %-ва точка  $\chi^2$ -розподілу з  $n$  ступенями свободи.

**II випадок.** Додаткової уваги потребує випадок наявності груп об'єктів з однаковим проявом принаймні за однією зі змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ ,  $2 \leq m \leq q+1$ . Після відповідної корекції, коефіцієнт конкордації  $\hat{W}_{\xi}$ , можна використовувати у вигляді такого **вибіркового модифікованого коефіцієнта конкордації**:

$$\hat{W}_{\xi} = \frac{\sum_{k=1}^n \left( \sum_{j=1}^m \left( x_k^{(ij)} \right) - \frac{m(n+1)}{2} \right)^2}{\frac{m^2(n^3 - n)}{12} - \frac{m}{2} \sum_{j=1}^m \Delta^{(ij)}},$$

де  $m^{(ij)}$  – кількість груп об'єктів з однаковим проявом за змінною  $\xi_{i_j}$ ;

$n_l^{(ij)}$  – кількість об'єктів, які ввійшли в  $l$ -ту групу об'єктів з однаковим проявом за змінною  $\xi_{i_j}$ ;

$$\Delta^{(ij)} = \frac{1}{6} \sum_{l=1}^{m^{(ij)}} \left( \left( n_l^{(ij)} \right)^3 - n_l^{(ij)} \right); \quad j = \overline{1, m}; \quad 2 \leq m \leq q+1.$$

## Кореляційний аналіз номінальних змінних

...

Самостійна робота №5. З навчального посібника  
«Слабоспицький О.С. Основи кореляційного аналізу даних, 2006».  
Пропрацювати матеріал наведений у Розділі 3.:  
«Дослідження наявності зв'язку між номінальними змінними.»  
(пропустити мінімум 4 стор.)

■ №7(ЗПМ)(2020.10.08). буде перед знаходженням  
оцінки ЗМНК у Р.А.

## Регресійний аналіз

*Регресійний аналіз* один з розділів аналізу даних, який  
займається побудовою математичних моделей істотних зв'язків між  
кількісними змінними. I

Нехай аналізується зв'язок між *залежною* скалярною змінною  $\eta$  та  
вектором *незалежних* змінних  $\bar{\xi}$ :

$$\eta \in \mathbb{R}, \quad \bar{\xi} \in \mathbb{R}^q.$$

Якщо кореляційний аналіз встановив, що статистичний зв'язок між  
цими змінними є істотний, то математичну модель залежності  $\eta$  від  $\bar{\xi}$   
можна шукати у вигляді *регресійної моделі*  $\eta$  *щодо*  $\bar{\xi}$ :

$$\eta = f(\bar{\xi}) + \varepsilon,$$



де  $f(\vec{x}) = M(\eta / \vec{\xi} = \vec{x})$  – функція регресії  $\eta$  щодо  $\vec{\xi}$ ,  $\varepsilon$  – залишкова похибка апроксимації, причому  $f(\vec{\xi})$  буде найкращою у середньоквадратичному розумінні апроксимацією  $\eta$  на класі борелівських функцій на  $\mathbb{R}^q$ .

Здається, що залишається тільки знайти вигляд функції  $f(\vec{x})$ , але для цього потрібно знати відповідний розподіл для  $\eta$  та  $\vec{\xi}$ . На практиці цей розподіл фактично є невідомим. І, як правило, єдиною доступною інформацією є тільки спостереження над цими змінними:

$$\eta: y(1), y(2), \dots, y(N),$$

$$\vec{\xi}: \vec{x}'(1), \vec{x}'(2), \dots, \vec{x}'(n).$$

Тобто потрібно спираючись тільки на ці спостереження знайти апроксимацію для  $f(\vec{x})$ .

**Основні етапи розв'язання задачі регресійного аналізу:**

1. вибір класу апроксимуючих функцій  $\tilde{F}$  для  $f(\vec{x})$ , тобто для функції регресії  $\eta$  щодо  $\vec{\xi}$ :

$$\tilde{f}(\vec{x}, \alpha) \in \tilde{F}, \quad \vec{x} \in \mathbb{R}^q, \alpha \in \mathbb{R}^p,$$

де  $\alpha$  - вектор невідомих параметрів;

2. визначення оптимальної, згідно деякого критерію якості, точкової оцінки  $\hat{\alpha}$  для  $\alpha$  та її характеристики розсіювання  $M(\hat{\alpha} - \alpha)(\hat{\alpha} - \alpha)^T$ , або множинної оцінки для  $\alpha$ , тобто довірчої області для  $\alpha$ ;
3. якщо модель лінійна по  $\alpha$ , то здійснюємо перевірку на значимість параметри моделі, а саме перевіряємо гіпотези:

$$H_0: \alpha = \theta, \gamma > 0, \quad \text{або}$$

$$|H_0: \alpha_i = 0, \gamma > 0.$$

4. перевірка на адекватність отриманої моделі.

**Зауваження.** В якості критерію якості на другому етапі при обчисленні точкової оцінки  $\hat{\alpha}$  для  $\alpha$  найчастіше використовують такий функціонал:

$$M\left(\eta - \tilde{f}(\vec{\xi}, \alpha)\right)^2.$$

Вибіркове представлення його має такий вигляд:

$$\frac{1}{N} \sum_{k=1}^n \left( y(k) - \tilde{f}(\vec{x}'(k), \alpha) \right)^2. \quad (1.1)$$

Саме його і буде використано у подальшому.

## Класичний регресійний аналіз

Нехай потрібно побудувати математичну модель зв'язку залежної скалярної змінної  $\eta$  від вектора незалежних змінних  $\vec{\xi} = (\xi_1, \xi_2, \dots, \xi_q)^T$ :

$$\eta \in \mathbb{R}, \quad \vec{\xi} \in \mathbb{R}^q.$$

У класичному регресійному аналізі в якості класу апроксимуючих функцій для  $f(\vec{x})$  беруть клас функцій лінійних по вектору невідомих параметрів  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_p)^T$ . Тобто математичну модель шукають у вигляді:

$$\eta = \sum_{i=1}^p \alpha_i \varphi_i(\vec{\xi}) + \varepsilon_{\mathbb{P}} \quad (1.2)$$

де  $\{\varphi_i(\bullet)\}_{i=1}^p$  - обраний (відомий) набір функцій,  $\varepsilon$  - похибка моделі.



Як правило, у якості  $\{\varphi_i(\bullet)\}_{i=1}^p$  обирають набір незалежних (ще краще ортогональних) функцій.

У подальшому,  $\varphi_i(\bar{\xi})$  будемо називати  $i$ -им регресором,  $i = \overline{1, p}$ .  
А відповідно  $\xi_i$  —  $i$ -ою незалежною змінною,  $i = \overline{1, q}$ .

У  $k$ -ий момент модель (1.2) набуває вигляду:

$$y(k) = \sum_{i=1}^p \alpha_i \varphi_i(\bar{x}'(k)) + e(k), k = \overline{1, N} \quad (1.3)$$

Введемо позначення  $\text{I}$

$$x(k) = \begin{pmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_p(k) \end{pmatrix} = \begin{pmatrix} \varphi_1(\bar{x}'(k)) \\ \varphi_2(\bar{x}'(k)) \\ \vdots \\ \varphi_p(\bar{x}'(k)) \end{pmatrix}.$$

Тоді (1.3) можна переписати таким чином

$$y(k) = \sum_{i=1}^p \alpha_i x_i(k) + e(k), k = \overline{1, N}.$$

Або у такому представленні

$$y(k) = x^T(k) \alpha + e(k), k = \overline{1, N}. \quad (1.4)$$

Останню систему рівнянь можна записати у матричному вигляді:

$$y = X\alpha + e \quad (1.5)$$

де

$$y = \begin{pmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{pmatrix}, X = \begin{pmatrix} x^T(1) \\ x^T(2) \\ \vdots \\ x^T(N) \end{pmatrix} \in M_{N,p}(\mathbb{R}), e = \begin{pmatrix} e(1) \\ e(2) \\ \vdots \\ e(N) \end{pmatrix}. \quad \text{I}$$

**Припущення класичного регресійного аналізу:**

- I.  $e \sim \mathcal{N}(\theta_N, \sigma^2 E_N), \sigma^2 \in \mathbb{R}_+,$
- II.  $\text{rank}(X) = p,$
- III. немає ніяких обмежень на  $\alpha$ , тобто  $\alpha \in \mathbb{R}^p$ .

**Зауваження.** З цих припущень випливає:

1. з першого, що  $e(k) \sim \mathcal{N}(0, \sigma^2)$  та є незалежними,  $k = \overline{1, N}$ .
2. з другого, що матриця  $X$  має повний ранг по стовпчикам, а  $N \geq p$ .
3. з третього, що  $\alpha$  може набувати довільного значення з простору  $\mathbb{R}^p$ , бо на  $\alpha$  не накладено ніяких обмежень.

Переходимо до 2-го етапу розв'язання задачі регресійного аналізу, а саме до пошуку оптимальної, згідно критерію якості (1.1), точкової оцінки  $\hat{\alpha}$  для  $\alpha$  та її характеристики розсіювання  $M(\hat{\alpha} - \alpha)(\hat{\alpha} - \alpha)^T$ .

Тобто потрібну оптимальну оцінку  $\hat{\alpha}$  шукаємо таким чином:

$$\begin{aligned}\hat{\alpha} &= \arg \min_{\alpha} \left\{ \frac{1}{N} \sum_{k=1}^N \left( y(k) - \tilde{f}(\bar{x}'(k), \alpha) \right)^2 \right\} = \\ &= \arg \min_{\alpha} \left\{ \sum_{k=1}^N \left( y(k) - x^T(k) \alpha \right)^2 \right\} = \\ &= \arg \min_{\alpha} \left\{ \sum_{k=1}^N e^2(k) \right\} = \arg \min_{\alpha} \|e\|^2,\end{aligned}$$

де  $\|e\|$  - евклідова норма вектора  $e$ .

**Означення.** Оцінка  $\hat{\alpha}$  для вектора невідомих параметрів  $\alpha$  моделі (1.5), яка розв'язком задачі

$$\hat{\alpha} = \arg \min_{\alpha} \|e\|^2,$$

називається оцінкою методу найменших квадратів (МНК).

Ця оцінка була запропонована у 1795 році (тобто у 18 столітті) Карлом Фрідріхом Гауссом, коли йому було тільки 18 років.

Потрібну оцінку  $\hat{\alpha}$  знайдемо як частинний випадок більш загальної оцінки.

**Означення.** Оцінка  $\hat{\alpha}_W$  для вектора невідомих параметрів  $\alpha$  моделі (1.5), яка розв'язком задачі

$$\hat{\alpha}_W = \arg \min_{\alpha} \|e\|_W^2, \quad W > 0,$$

називається оцінкою зваженого методу найменших квадратів (ЗМНК), де  $\|e\|_W^2 = e^T W e$ .

Спочатку знайдемо оцінку  $\hat{\alpha}_W$ , а оцінку МНК  $\hat{\alpha}$  отримаємо, як частинний випадок оцінки ЗМНК, бо  $\hat{\alpha} = \hat{\alpha}_E$ .

**Зауваження.** Відомо, що

$$\text{grad}_{\alpha} (\alpha^T \beta) = \text{grad}_{\alpha} (\beta^T \alpha) = \beta,$$

$$\text{grad}_{\alpha} (\alpha^T A \alpha) = (A + A^T) \alpha,$$

де  $\alpha, \beta$  - вектори, а  $A$  - матриця відповідних розмірностей.

Переходимо до знаходження оцінки  $\hat{\alpha}_W$ . Розпишемо спочатку вираз критерію