

УДК 519.2(075.8)

ББК 22.17я73

С47

Рецензенти:

д-р фіз.-мат. наук, проф. О. Г. Наконечний,  
канд. фіз.-мат. наук, доц. В. А. Заславський

Рекомендовано до друку вченою радою факультету кібернетики  
(протокол № 2 від 16 жовтня 2006 року)

Слабоспицький О. С.

С47 Основи кореляційного аналізу даних: Навчальний посібник. – К.:  
ВПЦ "Київський університет", 2006. – 59 с.  
ISBN 966-594-869-5

Висвітлено базові методи кореляційного аналізу, які дозволяють з'ясувати наявність статистичного зв'язу між змінними, що досліджуються. Розглянуто випадки аналізу кількісних, ординальних і номінальних змінних.

Для студентів, які навчаються за спеціальністю "Аналіз даних".

УДК 519.2(075.8)

ББК 22.17я73

ISBN 966-594-869-5

© О. С. Слабоспицький, 2006  
© Київський національний університет імені Тараса Шевченка,  
ВПЦ "Київський університет", 2006

## ВСТУП

Сучасні технології кожному надають доступ до неосяжних об'ємів різноманітної інформації. Але це не означає, що в подальшому отримані дані будуть кваліфіковано оброблені, а потім з результатів цієї обробки будуть зроблені правильні висновки.

Для розв'язання такого роду проблем можна скористатися математичними методами аналізу даних. Їх використання прискореними темпами на практиці сприяла поява широкого кола відповідних програмних продуктів для сучасної обчислювальної техніки, які дозволяють обробляти інформацію незалежно від її природи. Це надало можливість ефективно розв'язувати різноманітні задачі в різних галузях, особливо в науці, бізнесі, медицині, фінансах, біології тощо.

Серед розділів аналізу даних, які найчастіше використовуються, можна назвати такі як: попередня обробка даних, кореляційний аналіз, дисперсійний аналіз, регресійний аналіз, коваріаційний аналіз, дискримінантний аналіз, кластерний аналіз, аналіз часових рядів.

Серед них чільне місце посідає кореляційний аналіз. Саме він дозволяє з'ясувати питання наявності істотного зв'язу між процесами (явищами, об'єктами), які досліджуються. Для цього, враховуючи тип даних, які обробляються, як правило, здійснюють такі кроки:

- обирають конкретну числову характеристику для парного чи множинного статистичного зв'язу;
- визначають оцінку значення характеристики, яка використовується;
- на основі оцінки значення характеристики парного чи множинного статистичного зв'язу приймають рішення, чи є істотним зв'язок між процесами (явищами, об'єктами), які досліджуються.

І тільки після стверджуваної відповіді про істотність статистичного зв'язу між змінними, що спостерігаються, має сенс переходити до наступного етапу – пошуку математичної моделі цього зв'язу засобами інших розділів аналізу даних.

У роботі висвітлено можливості, які надає кореляційний аналіз даних, при дослідженні наявності парних і множинних статистичних зв'язків. Окремо розглянуто випадки аналізу кількісних, ординальних і номінальних змінних.

Автор щиро вдячний студентам факультету кібернетики у сприянні покращенню даного посібника.

Усі зауваження та побажання щодо даної роботи будуть з вдячністю сприйняті автором. Їх можна надіслати електронною поштою ([sl@univ.kiev.ua](mailto:sl@univ.kiev.ua)).

# 1. АНАЛІЗ НАЯВНОСТІ СТАТИСТИЧНОГО ЗВ'ЯЗКУ МІЖ КІЛЬКІСНИМИ ЗМІННИМИ

У цьому розділі увага буде зосереджена на дослідженні наявності статистичного зв'язку між кількісними змінними. Спочатку аналізуються скалярні залежна та незалежна змінні, потім випадок скалярної залежної змінної та вектора незалежних змінних. Із цією метою використовують відповідні парні та множинні характеристики статистичного зв'язку, які базуються на фундаментальній властивості функції регресії. Крім цього описується частинний коефіцієнт кореляції. У всіх випадках наводяться методики їх використання.

Як правило, характеристики статистичного зв'язку вводяться таким чином, щоб у разі рівності нулю, зв'язок був відсутній, а суттєвість зв'язку між змінними, що досліджуються, зростала зі збільшенням його відхилення від нуля. Таким чином, щоб з'ясувати істотність зв'язку за допомогою деякої характеристики статистичного зв'язку  $K$ , яка прийняла ненульове значення, достатньо перевірити гіпотезу про те, чи значимо відхиляється від нуля коефіцієнт  $K$ , тобто здійснити перевірку його на значимість. Це зводиться до перевірки гіпотези:

$$H_0: K = 0,$$

з деяким рівнем значущості  $\alpha > 0$ .

Приклади задач кореляційного аналізу кількісних змінних:

- чи суттєво впливає рівень безробіття населення на рівень злочинності в країні?
- наскільки істотним є вплив на рівень життя населення основних показників економіки країни?
- чи залежить від оцінок, які отримав студент під час роботи протягом семестру, підсумковий результат у конкретного викладача з його предмету?
- чи істотним є зв'язок між температурою зовнішнього повітря та кількістю вжитих безалкогольних напоїв населенням?

## 1.1. Середньоквадратична апроксимація випадкової величини. Функція регресії та її властивості

Нехай на ймовірнісному просторі  $(\Omega, F, P)$  розглядаються скалярна випадкова величина  $\eta$  і випадковий вектор  $\bar{\xi} \in R^q$ . Будемо інтерпретувати  $\eta$ , як залежну змінну, а  $\bar{\xi}$  – як вектор незалежних змінних.

Припустимо, що  $M\eta^2 < \infty$ , тоді будуть існувати умовне математичне сподівання та умовна дисперсія випадкової величини  $\eta$  відносно події  $\{\bar{\xi} = \bar{x}\}$ , для яких у подальшому будуть використовуватися такі позначення:

$$f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x}), \quad g(\bar{x}) = D(\eta / \bar{\xi} = \bar{x}).$$

Надалі буде корисне таке твердження.

**Лема.** Якщо  $\eta$  і  $\bar{\xi}$  – випадкові величина та вектор, відповідно, а  $M\eta^2 < \infty$ , тоді для них справедливо:

$$D\eta = M \left\{ \left( M(\eta / \bar{\xi}) - M\eta \right)^2 \right\} + M \left\{ M \left\{ \left( \eta - M(\eta / \bar{\xi}) \right)^2 / \bar{\xi} \right\} \right\},$$

або скорочено  $D\eta = Df(\bar{\xi}) + Mg(\bar{\xi})$ .

З доведенням леми можна ознайомитися у додатку 2.

**Означення.** Нехай  $\eta$  і  $\bar{\xi}$  – випадкові величина та вектор, відповідно, причому  $M|\eta| < \infty$ . Тоді функцією регресії  $\eta$  на  $\bar{\xi}$  (або функцією регресії  $\eta$  щодо  $\bar{\xi}$ ) називається функція

$$f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x}).$$

Функція регресії  $f(\bar{x})$  має унікальну властивість, що підтверджує таке важливе твердження.

**Теорема.** Нехай  $M\eta^2 < \infty$ ,  $\Phi$  – множина борелівських функцій на  $R^q$ , тоді

$$f(\cdot) = \arg \min_{\varphi(\cdot) \in \Phi} M \left[ \eta - \varphi(\bar{\xi}) \right]^2,$$

де  $f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x})$ .

Доведення цього фундаментального твердження викладено в додатку 2.

Наведена теорема дозволяє стверджувати, що найкращою в середньоквадратичному розумінні апроксимацією  $\eta$  на класі борелівських функцій від  $\bar{\xi}$  є функція  $f(\bar{\xi})$ , тобто за математичну модель можна взяти співвідношення

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

де  $\varepsilon$  – залишкова помилка апроксимації.

Для цієї моделі мають місце такі властивості:

- 1)  $Mf(\bar{\xi}) = M\eta$ ,  $M\varepsilon = 0$ ;
- 2)  $f(\bar{\xi})$  та  $\varepsilon$  – некорельовані;

$$3) D\eta = Df(\bar{\xi}) + D\varepsilon.$$

Зауваження. Перша властивість дозволяє запропонувати для останнього співвідношення ще одне представлення:

$$D\eta = Df(\bar{\xi}) + M\varepsilon^2.$$

Доведення. Скористаємося властивостями умовного математичного сподівання випадкової величини та доведемо послідовно кожен з властивостей:

1)  $Mf(\bar{\xi}) = M\{M(\eta/\bar{\xi})\} = M\eta$ . А це, у свою чергу, дозволяє стверджувати, що  $M\varepsilon = M\{\eta - f(\bar{\xi})\} = M\eta - Mf(\bar{\xi}) = 0$ ;

2) оскільки, згідно з попередньою властивістю  $M\varepsilon = 0$ , то для доведення некорельованості достатньо впевнитися, що  $M[f(\bar{\xi}) - Mf(\bar{\xi})]\varepsilon = 0$ . Дійсно,

$$\begin{aligned} M[f(\bar{\xi}) - Mf(\bar{\xi})]\varepsilon &= M[f(\bar{\xi}) - M\eta][\eta - f(\bar{\xi})] = \\ &= M\{M[(f(\bar{\xi}) - M\eta)(\eta - f(\bar{\xi}))/\bar{\xi}]\} = \\ &= M\{(f(\bar{\xi}) - M\eta)M[(\eta - f(\bar{\xi}))/\bar{\xi}]\} = 0; \end{aligned}$$

3) із доведеної в другій властивості некорельованості  $f(\bar{\xi})$  та  $\varepsilon$  випливає:  $D\eta = D(f(\bar{\xi}) + \varepsilon) = Df(\bar{\xi}) + D\varepsilon$ .

Властивості доведено.

## 1.2. Індекс кореляції та його властивості. Коефіцієнт детермінації

Припустимо, що на ймовірнісному просторі  $(\Omega, F, P)$  задано скалярну залежну змінну  $\eta$  і випадковий вектор незалежних змінних  $\bar{\xi} \in R^q$ . Введемо універсальну характеристику статистичного зв'язку для змінних  $\eta$  та  $\bar{\xi}$ . Будемо вважати, що  $M\eta^2 < \infty$ , тоді існуватиме функція регресії  $\eta$  на  $\bar{\xi}$ , а саме:  $f(\bar{x}) = M(\eta/\bar{\xi} = \bar{x})$ . А для випадкової величини  $\eta$  можна використовувати таку математичну модель:  $\eta = f(\bar{\xi}) + \varepsilon$ , де  $\varepsilon$  – залишкова помилка апроксимації. Причому результати з попереднього розділу дозволяють стверджувати, що мають місце такі властивості:

$$1) Mf(\bar{\xi}) = M\eta, M\varepsilon = 0;$$

$$2) f(\bar{\xi}) \text{ та } \varepsilon - \text{некорельовані};$$

$$3) D\eta = Df(\bar{\xi}) + D\varepsilon.$$

Після цих міркувань, враховуючи унікальну властивість функції регресії, доведену в теоремі (див. п. 1.1), цілком природним здається використання, як характеристики статистичного зв'язку для кількісних змінних  $\eta$  та  $\bar{\xi}$ , нижчевизначеного індексу кореляції.

Означення. Нехай  $\eta$  і  $\bar{\xi}$  ( $\bar{\xi} \in R^q$ ) – випадкові величина та вектор, відповідно, причому  $0 < M\eta^2 < \infty$ . Тоді (множинним) індексом кореляції  $\eta$  щодо  $\bar{\xi}$  називається величина

$$I_{\eta\bar{\xi}} = \sqrt{\frac{Df(\bar{\xi})}{D\eta}}.$$

Зауваження. Оскільки для дисперсії функції регресії справедливо  $Df(\bar{\xi}) = D\eta - D\varepsilon$ , звідси випливає, що для індексу кореляції  $I_{\eta\bar{\xi}}$  можна отримати ще одне представлення:

$$I_{\eta\bar{\xi}} = \sqrt{1 - \frac{D\varepsilon}{D\eta}},$$

де  $\varepsilon = \eta - f(\bar{\xi})$ .

Властивості індексу кореляції:

$$1) 0 \leq I_{\eta\bar{\xi}} \leq 1;$$

2) якщо  $I_{\eta\bar{\xi}} = 0$ , то відсутній вплив  $\bar{\xi}$  на  $\eta$ ;

3) якщо  $I_{\eta\bar{\xi}} = 1$ , то існує функціональний зв'язок між  $\eta$  та  $\bar{\xi}$ , а саме, з ймовірністю 1 справедливо  $\eta = f(\bar{\xi})$ .

(Тут і далі для випадкових величин/векторів рівності/нерівності вважаються справедливими з ймовірністю 1, якщо не наголошується на іншому.)

Доведення. Скористаємось властивостями умовного математичного сподівання та умовної дисперсії випадкової величини  $\eta$  відносно випадкового вектора  $\bar{\xi}$  і послідовно доведемо потрібні властивості індексу кореляції:

1) згідно з вищенаведеними властивостями та враховуючи, що  $D\varepsilon \geq 0$ , отримуємо:

$$0 \leq I_{\eta\bar{\xi}} = \sqrt{\frac{Df(\bar{\xi})}{D\eta}} = \sqrt{\frac{D\eta - D\varepsilon}{D\eta}} \leq 1;$$

2) припустимо, що  $I_{\eta\bar{\xi}} = 0$ . Тоді

$$Df(\bar{\xi}) = 0 \Rightarrow f(\bar{\xi}) = Mf(\bar{\xi}) \Rightarrow f(\bar{\xi}) = M\eta = \text{const}.$$

Звідси випливає, що функція регресії, за допомогою якої апроксимується  $\eta$ , не залежить від значень свого аргументу –  $\bar{\xi}$ ;

3) нехай  $I_{\eta\bar{\xi}} = 1$ , тоді

$$1 = I_{\eta\bar{\xi}} = \sqrt{1 - \frac{D\varepsilon}{D\eta}} \Rightarrow D\varepsilon = 0.$$

Остання умова еквівалентна наступній:  $0 = D\varepsilon = M\{\eta - f(\bar{\xi})\}^2$ . А це, у свою чергу, дозволяє стверджувати, що з ймовірністю 1 має місце така функціональна залежність між  $\eta$  та  $\bar{\xi}$ :  $\eta = f(\bar{\xi})$ .

Справедливість властивостей доведено.

Поряд з індексом кореляції  $I_{\eta\bar{\xi}}$ , використовується також характеристика, яку визначає таке визначення.

**Означення.** Нехай  $\eta$  і  $\bar{\xi}$  ( $\bar{\xi} \in R^q$ ) – випадкові величина та вектор, відповідно, причому  $0 < M\eta^2 < \infty$ . Тоді (множинним) коефіцієнтом детермінації  $\eta$  щодо  $\bar{\xi}$  називається величина

$$I_{\eta\bar{\xi}}^2 = \frac{Df(\bar{\xi})}{D\eta}.$$

**Зауваження 1.** Коефіцієнт детермінації  $I_{\eta\bar{\xi}}^2$  є більш привабливим порівняно з індексом кореляції, тому що має більш прозору інтерпретацію, а саме: він вказує, яка частина дисперсії змінної  $\eta$  визначається варіацією функції регресії  $f(\bar{\xi})$ .

**Зауваження 2.** Аналогічно, для коефіцієнта детермінації також має інше представлення:

$$I_{\eta\bar{\xi}}^2 = 1 - \frac{D\varepsilon}{D\eta}.$$

У підсумку, для кількісних змінних, випадкової величини  $\eta$  ( $0 < M\eta^2 < \infty$ ) і випадкового вектора  $\bar{\xi}$  ( $\bar{\xi} \in R^q$ ), введено теоретичну характеристику статистичного зв'язку – індекс кореляції  $I_{\eta\bar{\xi}}$  (коефіцієнт детермінації  $I_{\eta\bar{\xi}}^2$ ).

### 1.3. Аналіз парних статистичних зв'язків кількісних змінних

При дослідженні зв'язку між двома скалярними змінними маємо справу з так званим *парним статистичним зв'язком*. Нехай на ймовірнісному просторі  $(\Omega, F, P)$  розглядаються дві скалярні випадкові величини:  $\eta$  – залежна змінна,  $\xi$  – незалежна змінна. Потрібно з'ясувати, чи є істотним зв'язок з статистичної точки зору, між цими змінними з рівнем значущості  $\alpha > 0$ .

Якщо припустити, що  $0 < M\eta^2 < \infty$ , то можна скористатися в загальному випадку для аналізу статистичного зв'язку між змінними  $\eta$  та  $\xi$  універсальною характеристикою – (парним) індексом кореляції  $\eta$  щодо  $\xi$ :

$$I_{\eta\xi} = \sqrt{\frac{Df(\xi)}{D\eta}},$$

де  $f(x) = M(\eta/\xi = x)$ .

Для залежної змінної  $\eta$  можна використовувати представлення  $\eta = f(\xi) + \varepsilon$ , де  $\varepsilon$  – залишкова помилка апроксимації, для якої справедливі:

- 1)  $Mf(\xi) = M\eta$ ,  $M\varepsilon = 0$ ;
- 2)  $f(\xi)$  та  $\varepsilon$  – некорельовані;
- 3)  $D\eta = Df(\xi) + D\varepsilon$ .

Остання властивість дозволяє використовувати також таке представлення для індексу кореляції  $\eta$  щодо  $\xi$ :

$$I_{\eta\xi} = \sqrt{1 - \frac{D\varepsilon}{D\eta}}.$$

Усі властивості індексу кореляції (див. розд. 1.2) залишаються в силі для  $I_{\eta\xi}$ :

- 1)  $0 \leq I_{\eta\xi} \leq 1$ ;
- 2) якщо  $I_{\eta\xi} = 0$ , то відсутній вплив  $\xi$  на  $\eta$ ;
- 3) якщо  $I_{\eta\xi} = 1$ , то існує функціональний зв'язок між  $\eta$  та  $\xi$ , а саме: з ймовірністю 1 справедливо  $\eta = f(\xi)$ .

Відповідно, (парний) коефіцієнт детермінації  $\eta$  щодо  $\xi$  обчислюється згідно з виразами

$$I_{\eta\xi}^2 = \frac{Df(\xi)}{D\eta} = 1 - \frac{D\varepsilon}{D\eta}.$$

Таким чином, маємо для скалярних кількісних змінних  $\eta$  ( $0 < M\eta^2 < \infty$ ) і  $\xi$  теоретичну характеристику парного статистичного зв'язку – індекс кореляції  $I_{\eta\xi}$  (коефіцієнт детермінації  $I_{\eta\xi}^2$ ).

### 1.3.1. Вибіркові значення характеристики парного статистичного зв'язку в загальному випадку

Для практичного використання теоретичних показників кореляційного зв'язку потрібно вміти обчислювати їх вибіркові (емпіричні) значення, на основі отриманих спостережень над кількісними змінними, які досліджуються. Процедура побудови відповідної оцінки для індексу кореляції, у подальшому, буде ґрунтуватися на його теоретичному визначенні у вигляді

$$I_{\eta\xi} = \sqrt{\frac{Df(\xi)}{D\eta}} = \sqrt{1 - \frac{D(\eta - f(\xi))}{D\eta}}.$$

Обчислення конкретних емпіричних значень індексу кореляції (коефіцієнта детермінації) буде залежати в кожному випадку від умов проведення експерименту. Перейдемо до розгляду цих ситуацій.

#### 1.3.1.1. Кореляційне відношення

Розглянемо спочатку випадок, коли проведено групування  $n$  спостережень над незалежною змінною  $\xi$ . Нехай отримано  $s$  інтервалів групування, причому виміри над залежною змінною  $\eta$ , які відповідають спостереженням над незалежною змінною з  $i$ -го інтервалу групування позначатимемо як  $y_{i1}, y_{i2}, \dots, y_{iv_i}$ ,  $i = \overline{1, s}$ . Тут  $v_i$  – кількість вимірів, які

потрапили в  $i$ -й інтервал групування  $\left( n = \sum_{i=1}^s v_i \right)$ .

**Зауваження.** У подальшому скрізь, як правило, в оцінках різних характеристик будемо опускати аргумент  $n$ , який вказує на об'єм вибірки, якщо кількість спостережень незмінна.

Скористаємося позначенням для загального середнього:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^s v_i \bar{y}_{i\bullet},$$

де  $\bar{y}_{i\bullet}$  – середнє на  $i$ -му інтервалі групування, яке обчислюється таким чином:

$$\bar{y}_{i\bullet} = \frac{1}{v_i} \sum_{j=1}^{v_i} y_{ij}, \quad i = \overline{1, s}.$$

Зауважимо, що крапка в  $\bar{y}_{i\bullet}$  є позначенням сумування за тим індексом, замість якого вона фігурує. Зважаючи на те, що  $D\eta$  можна оцінити згідно з

$$s_{\eta}^2 = \frac{1}{n} \sum_{i=1}^s \sum_{j=1}^{v_i} (y_{ij} - \bar{y})^2,$$

а  $Df(\xi)$  – відповідно, за формулою  $s_f^2 = \frac{1}{n} \sum_{i=1}^s v_i (\bar{y}_{i\bullet} - \bar{y})^2$ , отримаємо таку оцінку для індексу кореляції  $\eta$  щодо  $\xi$ :

$$\hat{r}_{\eta\xi} = \sqrt{\frac{s_f^2}{s_{\eta}^2}}.$$

Оцінку  $\hat{r}_{\eta\xi}$  називають парним кореляційним відношенням  $\eta$  щодо  $\xi$ .

Методика використання  $\hat{r}_{\eta\xi}$  є традиційною:

- 1) якщо  $\hat{r}_{\eta\xi} = 0$ , то вважаємо, що зв'язок між  $\eta$  і  $\xi$  відсутній;
- 2) якщо  $\hat{r}_{\eta\xi} = 1$ , то вважаємо, що зв'язок між  $\eta$  і  $\xi$  функціональний;
- 3) якщо  $\begin{cases} \hat{r}_{\eta\xi} \neq 0 \\ \hat{r}_{\eta\xi} \neq 1 \end{cases}$ , то істотність зв'язку з'ясовується шляхом перевірки на значимість індексу кореляції  $\eta$  щодо  $\xi$ :

$$H_0 : I_{\eta\xi} = 0,$$

з деяким рівнем значущості  $\alpha > 0$ . Якщо гіпотеза не справедлива, то статистичний зв'язок між  $\eta$  і  $\xi$  вважається істотним, інакше зв'язок вважається не істотним.

Для перевірки останньої гіпотези скористаємося тим, що при справедливості гіпотези  $H_0$ , можна вважати, що розподіл статистики

$$F = \frac{\hat{r}_{\eta\xi}^2}{1 - \hat{r}_{\eta\xi}^2} \cdot \frac{n-s}{s-1}$$

можна наблизити  $F$ -розподілом з  $(s-1)$  та  $(n-s)$  ступенями свободи. Тоді, з огляду на структуру статистики  $F$ , як критичну область для гіпотези  $H_0$  потрібно взяти область великих значень, а відповідна область прийняття для нашої гіпотези матиме вигляд

$$F < F_{\alpha}(s-1, n-s),$$

де  $F_{\alpha}(m, n)$  –  $100\alpha\%$ -на точка  $F$ -розподілу з  $m$  та  $n$  ступенями свободи.

### 1.3.1.2. Оцінка індексу кореляції

Розглянемо тепер ситуацію, коли спостереження відповідних змінних  $\eta: y_1, y_2, \dots, y_n$ ,  $\xi: x_1, x_2, \dots, x_n$ , є доступними, а також є можливість припустити, що функція регресії  $\eta$  на  $\xi$  може бути наближена на деякому класі параметричних функцій  $f(x, \lambda)$ ,  $\lambda \in R^p$ .

Якщо отримана оцінка  $\hat{\lambda}$  для вектора невідомих параметрів  $\lambda$  деяким відомим методом, то вибіркове значення  $D_{\eta\xi}$  можна обчислити таким чином:

$$\frac{1}{n-p} \sum_{i=1}^n (y_i - f(x_i, \hat{\lambda}))^2,$$

а емпіричне (вибіркове) значення індексу кореляції набуде вигляду

$$\hat{I}_{\eta\xi} = \sqrt{1 - \frac{\frac{1}{n-p} \sum_{i=1}^n (y_i - f(x_i, \hat{\lambda}))^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}},$$

$$\text{де } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Для практичного використання оцінки  $\hat{I}_{\eta\xi}$  не вистачає тільки процедури перевірки індексу кореляції на значимість, тобто гіпотези

$$H_0: I_{\eta\xi} = 0$$

з деяким рівнем значущості  $\alpha > 0$ . Для цього скористаємося тим фактом, що при справедливості гіпотези  $H_0$  розподіл статистики

$$F = \frac{\hat{I}_{\eta\xi}^2}{1 - \hat{I}_{\eta\xi}^2} \cdot \frac{n-p}{p-1}$$

можна апроксимувати  $F$ -розподілом з  $(p-1)$  та  $(n-p)$  ступенями свободи. Тоді область прийняття для гіпотези  $H_0$  матиме вигляд:

$$F < F_{\alpha}(p-1, n-p),$$

де  $F_{\alpha}(m, n)$  –  $100\alpha\%$ -на точка  $F$ -розподілу з  $m$  та  $n$  ступенями свободи.

У розглянутому випадку для індексу кореляції отримано формулу емпіричного значення та процедури перевірки на значиме відхилення від нуля.

### 1.3.2. Коефіцієнт кореляції – характеристика парного статистичного зв'язку у нормальному випадку

Припустимо, що  $\eta$  і  $\xi$  – нормально розподілені випадкові величини, задані на ймовірнісному просторі  $(\Omega, F, P)$ , а саме:

$$\eta \sim N(m_{\eta}, \sigma_{\eta}^2), \xi \sim N(m_{\xi}, \sigma_{\xi}^2), \sigma_{\eta}^2, \sigma_{\xi}^2 > 0.$$

Перш ніж проводити обчислення емпіричного значення характеристики парного статистичного зв'язку для нормальних змінних, проведемо аналіз теоретичної характеристики. Виявляється, що з'ясування питання про істотність зв'язку між змінними  $\eta$  та  $\xi$  з рівнем значущості  $\alpha > 0$  спрощується. Дійсно, за цих припущень індекс кореляції  $I_{\eta\xi}$  легко виражається через (парний) коефіцієнт кореляції  $r_{\eta\xi}$  випадкових величин  $\eta$  та  $\xi$ , який обчислюється як:

$$r_{\eta\xi} = \frac{\text{cov}(\eta, \xi)}{\sigma_{\eta}\sigma_{\xi}},$$

де  $\text{cov}(\eta, \xi) = M(\eta - m_{\eta})(\xi - m_{\xi})$  – коефіцієнт коваріації випадкових величин  $\eta$  та  $\xi$ .

У нормальному випадку для функцій  $f(x) = M(\eta/\xi = x)$ ,  $g(x) = D(\eta/\xi = x)$  справедливе твердження.

**Т е о р е м а .** Нехай пара  $(\eta, \xi)$  має гауссівський розподіл і  $\sigma_{\xi}^2 > 0$ , тоді справедливо

$$f(x) = m_{\eta} + \frac{\text{cov}(\eta, \xi)}{\sigma_{\xi}^2}(x - m_{\xi}), \quad g(x) = \sigma_{\eta}^2 - \frac{\text{cov}^2(\eta, \xi)}{\sigma_{\xi}^2}.$$

Якщо ж припустити, що  $\sigma_{\eta}^2 > 0$ , то має місце

$$f(x) = m_{\eta} + r_{\eta\xi} \frac{\sigma_{\eta}}{\sigma_{\xi}}(x - m_{\xi}), \quad g(x) = \sigma_{\eta}^2(1 - r_{\eta\xi}^2).$$

Остання теорема дозволяє перекоонатися, що має місце таке твердження.

**Т е о р е м а .** Нехай  $\eta$  і  $\xi$  – нормально розподілені випадкові величини:  $\eta \sim N(m_{\eta}, \sigma_{\eta}^2)$ ,  $\xi \sim N(m_{\xi}, \sigma_{\xi}^2)$ ,  $\sigma_{\eta}^2, \sigma_{\xi}^2 > 0$ . Тоді індекс кореляції  $I_{\eta\xi} = |r_{\eta\xi}|$ , а відповідно, коефіцієнт детермінації  $I_{\eta\xi}^2 = r_{\eta\xi}^2$ .

**Доведення.** Згідно з попередньою теоремою індекс кореляції  $I_{\eta\xi}$  можна обчислити як:

$$I_{\eta\xi} = \sqrt{\frac{Df(\xi)}{D\eta}} = \sqrt{\frac{D\left(m_{\eta} + r_{\eta\xi} \frac{\sigma_{\eta}}{\sigma_{\xi}} (\xi - m_{\xi})\right)}{\sigma_{\eta}^2}} = \sqrt{\frac{r_{\eta\xi}^2 \frac{\sigma_{\eta}^2}{\sigma_{\xi}^2} D(\xi - m_{\xi})}{\sigma_{\eta}^2}} = |r_{\eta\xi}|.$$

Тоді представлення для коефіцієнта детермінації у вигляді  $I_{\eta\xi}^2 = r_{\eta\xi}^2$  стає очевидним. Доведення завершено.

У результаті, у нормальному випадку як характеристику парного статистичного зв'язку можна використовувати не індекс кореляції  $I_{\eta\xi}$ , а звичайний коефіцієнт кореляції  $r_{\eta\xi}$ . Його властивості загальновідомі:

- 1)  $-1 \leq r_{\eta\xi} \leq 1$ ;
- 2) якщо  $r_{\eta\xi} = 0$ , то  $\eta$  не залежить від  $\xi$ ;
- 3) якщо  $|r_{\eta\xi}| = 1$ , то з ймовірністю 1 існує лінійний зв'язок між  $\eta$  та  $\xi$ , а саме, справедливо

$$\eta = m_{\eta} + r_{\eta\xi} \frac{\sigma_{\eta}}{\sigma_{\xi}} (\xi - m_{\xi});$$

- 4)  $r_{\eta\xi} = r_{\xi\eta}$ .

**Доведення.** Для перевірки їх справедливості скористаємося властивостями індексу кореляції:

- 1) дійсно, з першої властивості  $I_{\eta\xi}$  випливає

$$1 \geq I_{\eta\xi} = |r_{\eta\xi}| \Rightarrow -1 \leq r_{\eta\xi} \leq 1;$$

- 2) оскільки  $r_{\eta\xi} = 0$ , то  $I_{\eta\xi} = 0$ . Згідно з другою властивістю індексу кореляції це означає, що  $\eta$  не залежить від  $\xi$ ;

- 3) так як  $|r_{\eta\xi}| = 1$ , то  $I_{\eta\xi} = 1$ . Тоді третя властивість  $I_{\eta\xi}$  дозволяє стверджувати, що з ймовірністю 1 справедливо  $\eta = f(\xi)$ . Але в норма-

льному випадку  $f(x) = m_{\eta} + r_{\eta\xi} \frac{\sigma_{\eta}}{\sigma_{\xi}} (x - m_{\xi})$ . Властивість доведено;

- 4) очевидно.

Доведення завершено.

**Зауваження.** Вид функції регресії  $f(x)$ , який наведено в теоремі вказує на те, що зв'язок між  $\eta$  і  $\xi$  має монотонний характер, а знак коефіцієнта кореляції  $r_{\eta\xi}$  конкретизує його, а саме: якщо  $r_{\eta\xi} > 0$  ( $r_{\eta\xi} < 0$ ), то функція залежності буде зростаючою (спадною).

Перейдемо до практичного використання коефіцієнта кореляції  $r_{\eta\xi}$  як характеристики парного статистичного зв'язку для скалярних кількісних змінних. Нехай спостереження над випадковими величинами

$$\eta: y_1, y_2, \dots, y_n, \quad \xi: x_1, x_2, \dots, x_n$$

є доступними. Тоді, на основі цих вимірів *вибіркове (емпіричне) значення*  $\hat{r}_{\eta\xi}$  (парного) коефіцієнта кореляції  $r_{\eta\xi}$  обчислюється згідно з

$$\hat{r}_{\eta\xi} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (x_i - \bar{x})^2}},$$

$$\text{де } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Методика аналізу статистичного зв'язку між змінними  $\eta$  та  $\xi$  така:

- 1) якщо  $\hat{r}_{\eta\xi} = 0$ , то вважаємо, що зв'язок між  $\eta$  та  $\xi$  відсутній;
- 2) якщо  $|\hat{r}_{\eta\xi}| = 1$ , то зв'язок між  $\eta$  та  $\xi$  функціональний (точніше лінійний) і має вигляд

$$\eta = m_{\eta} + r_{\eta\xi} \frac{\sigma_{\eta}}{\sigma_{\xi}} (\xi - m_{\xi});$$

- 3) якщо  $\left\{ \begin{array}{l} \hat{r}_{\eta\xi} \neq 0 \\ |\hat{r}_{\eta\xi}| \neq 1 \end{array} \right.$ , то потрібно перевірити гіпотезу про те, чи суттєво

відхиляється від нуля коефіцієнт кореляції  $r_{\eta\xi}$  із статистичної точки зору, тобто здійснити *перевірку* його на *значимість*. Для цього потрібно перевірити гіпотезу:

$$H_0: r_{\eta\xi} = 0,$$

з деяким рівнем значущості  $\alpha > 0$ . І якщо гіпотеза виявиться справедливою, то вважається, що статистичний зв'язок між  $\eta$  та  $\xi$  є не істотним, у протилежному випадку – істотним.

Конкретизуємо процедуру перевірки останньої гіпотези. Скористаємося тим фактом, що при справедливості гіпотези  $H_0$ , статистика

$$t(n) = \frac{\sqrt{n-2} \hat{r}_{\eta\bar{\xi}}}{\sqrt{1-\hat{r}_{\eta\bar{\xi}}^2}}$$

має  $t$ -розподіл Стюдента з  $(n-2)$  ступенями свободи. Тоді логічно за критичну область цієї гіпотези взяти області набуття статистикою своїх екстремальних значень. А область прийняття гіпотези  $H_0$  матиме представлення:

$$|t(n)| < t_{\alpha/2}(n-2),$$

де  $t_{\alpha/2}(n)$  – 100 $\alpha$ %-на точка  $t$ -розподілу Стюдента з  $n$  ступенями свободи.

**Зауваження.** Скрізь у подальшому, як правило, будемо вважати справедливим припущення про нормальність відповідних спостережень у разі необхідності перевірки гіпотез або побудови довірчих інтервалів, якщо не обумовлено інше.

Таким чином, у нормальному випадку за характеристику парного статистичного зв'язку можна використовувати звичайний коефіцієнт кореляції  $r_{\eta\bar{\xi}}$ , який має прозору інтерпретацію й дозволяє стверджувати про наявність лінійного зв'язку, коли він набуває значення  $\pm 1$ . У загальному випадку коефіцієнт кореляції вже не має такої яскравої інтерпретації і тому виникає потреба у зверненні до універсальної характеристики парного статистичного зв'язку – індексу кореляції  $I_{\eta\bar{\xi}}$ .

#### 1.4. Аналіз множинних статистичних зв'язків кількісних змінних

Будемо стверджувати, що маємо справу з *множинним статистичним зв'язком*, якщо досліджується зв'язок між більше, ніж двома скалярними змінними. До цього розглядалися зв'язки тільки між двома скалярними змінними, тобто аналізувалися парні статистичні зв'язки.

Нехай на ймовірнісному просторі  $(\Omega, F, P)$  задані скалярна залежна змінна  $\eta$  і вектор незалежних змінних  $\bar{\xi} \in R^q$  ( $\bar{\xi}^T = (\xi_1, \xi_2, \dots, \xi_q)$ ). Проведемо дослідження питання про істотність зв'язку між змінними  $\eta$  та  $\bar{\xi}$ , з статистичної точки зору, з рівнем значущості  $\alpha > 0$ .

Припустимо, що  $0 < M\eta^2 < \infty$ . Тоді, якщо відсутня додаткова апіорна інформація, можна скористатися множинним індексом кореляції  $\eta$  щодо  $\bar{\xi}$ , який було введено раніше:

$$I_{\eta\bar{\xi}} = \sqrt{\frac{Df(\bar{\xi})}{D\eta}} = \sqrt{1 - \frac{D(\eta - f(\bar{\xi}))}{D\eta}},$$

де  $f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x})$ . Зауважимо, що залежна змінна  $\eta$  може бути представлена у вигляді  $\eta = f(\bar{\xi}) + \varepsilon$ , де  $\varepsilon$  – помилка моделі. Причому, як було доведено, мають місце властивості:

- 1)  $Mf(\bar{\xi}) = M\eta$ ,  $M\varepsilon = 0$ ;
- 2)  $f(\bar{\xi})$  та  $\varepsilon$  – некорельовані;
- 3)  $D\eta = Df(\bar{\xi}) + D\varepsilon$ .

А для індексу кореляції  $I_{\eta\bar{\xi}}$ , у свою чергу, справедливо:

- 1)  $0 \leq I_{\eta\bar{\xi}} \leq 1$ ;
- 2) якщо  $I_{\eta\bar{\xi}} = 0$ , то відсутній вплив  $\bar{\xi}$  на  $\eta$ ;
- 3) якщо  $I_{\eta\bar{\xi}} = 1$ , то існує функціональний зв'язок між  $\eta$  та  $\bar{\xi}$  у вигляді  $\eta = f(\bar{\xi})$ .

У результаті, у загальному випадку можна використовувати для аналізу статистичного зв'язку між скалярною залежною змінною  $\eta$  і вектором незалежних змінних  $\bar{\xi}$  теоретичну характеристику множинного статистичного зв'язку – індекс кореляції  $I_{\eta\bar{\xi}}$ .

##### 1.4.1. Емпіричні значення характеристики множинного статистичного зв'язку в загальному випадку

Як і при аналізі парних зв'язків, обчислення вибіркового значення індексу кореляції  $I_{\eta\bar{\xi}}$  залежатиме від конкретної ситуації під час обробки даних. Окремо розглянемо випадки обробки згрупованих спостережень і наявності можливості апроксимації функції регресії на деякому класі параметричних функцій.

###### 1.4.1.1. Обробка згрупованих даних

Нехай отримано  $n$  спостережень над скалярною випадковою величиною  $\eta$  і випадковим вектором  $\bar{\xi}$  ( $\bar{\xi}^T = (\xi_1, \xi_2, \dots, \xi_q)$ ). Проведемо групування за вектором незалежних змінних  $\bar{\xi}$ . Якщо раніше у випадку однієї



незалежної змінної під час групування її область значень розбивалася на інтервали групування, то у векторному випадку область значень  $\bar{\xi}$  розбивається на їх аналоги – гіперпаралелепіеди групування. Якщо область значень змінної  $\xi_i$  розбилася на  $s_i$  ( $i = \overline{1, q}$ ) інтервалів, то отримуємо

$s = \prod_{i=1}^q s_i$  гіперпаралелепіедів групування за вектором незалежних змінних

$\bar{\xi}$ . Тоді спостереження над залежною змінною  $\eta$ , які відповідають вимірам над вектором незалежних змінних  $\bar{\xi}$  з  $i$ -го гіперпаралелепіеда групування будемо позначати як  $y_{i1}, y_{i2}, \dots, y_{iv_i}$ ,  $i = \overline{1, s}$ , де  $v_i$  – кількість вимірів, які потрапили в  $i$ -й гіперпаралелепіед групування  $\left( n = \sum_{i=1}^s v_i \right)$ .

У подальшому процедура визначення оцінки індексу кореляції буде подібною до випадку скалярної незалежної змінної. Введемо позначення для загального середнього

$$\bar{y} = \frac{1}{n} \sum_{i=1}^s v_i \bar{y}_{i\bullet}, \text{ де } \bar{y}_{i\bullet} = \frac{1}{v_i} \sum_{j=1}^{v_i} y_{ij}, i = \overline{1, s}.$$

Тоді для підрахування оцінок  $D\eta$  та  $Df(\bar{\xi})$  можна скористатися такими виразами відповідно:

$$s_{\eta}^2 = \frac{1}{n} \sum_{i=1}^s \sum_{j=1}^{v_i} (y_{ij} - \bar{y})^2, s_f^2 = \frac{1}{n} \sum_{i=1}^s v_i (\bar{y}_{i\bullet} - \bar{y})^2.$$

Це дає змогу обчислювати вибіркове значення індексу кореляції  $I_{\eta\bar{\xi}}$  таким чином:

$$\hat{\rho}_{\eta\bar{\xi}} = \sqrt{\frac{s_f^2}{s_{\eta}^2}}.$$

Аналогічно до вже розглянутих парних кореляційних відношень, оцінку  $\hat{\rho}_{\eta\bar{\xi}}$  будемо називати *множинним кореляційним відношенням*  $\eta$  щодо  $\bar{\xi}$ .

Процедура використання статистики  $\hat{\rho}_{\eta\bar{\xi}}$  є ідентичною до методики використання парного кореляційного відношення.

#### 1.4.1.2. Використання наближення функції регресії

Припустимо, що функцію регресії  $\eta$  на  $\bar{\xi}$  можна апроксимувати на класі параметричних функцій  $f(\bar{x}, \lambda)$ ,  $\bar{x} \in R^q$ ,  $\lambda \in R^p$ . Нехай під час

аналізу можна використовувати такі спостереження над залежною змінною  $\eta$  і вектором незалежних змінних  $\bar{\xi}$ :

$$\eta: y_1, y_2, \dots, y_n, \quad \bar{\xi}: \bar{x}_1, \bar{x}_2, \dots, \bar{x}_n.$$

Наявність оцінки  $\hat{\lambda}$  для вектора параметрів  $\lambda$  дозволяє підрахувати емпіричне (вибіркове) значення індексу кореляції  $I_{\eta\bar{\xi}}$  таким чином:

$$\hat{I}_{\eta\bar{\xi}} = \sqrt{1 - \frac{\frac{1}{n-p} \sum_{i=1}^n (y_i - f(\bar{x}_i, \hat{\lambda}))^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}},$$

$$\text{де } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Методика використання вибіркового значення  $\hat{I}_{\eta\bar{\xi}}$  у множинному випадку не чим не відрізняється від відповідної процедури в парному випадку. *Перевірка індексу кореляції на значимість*

$$H_0: I_{\eta\bar{\xi}} = 0,$$

з деяким рівнем значущості  $\alpha > 0$ , здійснюється за допомогою статистики

$$F = \frac{\hat{I}_{\eta\bar{\xi}}^2}{1 - \hat{I}_{\eta\bar{\xi}}^2} \cdot \frac{n-p}{p-1},$$

розподіл якої можна наблизити  $F$ -розподілом з  $(p-1)$  та  $(n-p)$  ступенями свободи. Це дозволяє область прийняття для гіпотези  $H_0$  записати у вигляді

$$F < F_{\alpha}(p-1, n-p),$$

де  $F_{\alpha}(m, n) - 100\alpha\%$  -на точка  $F$ -розподілу з  $m$  та  $n$  ступенями свободи.

#### 1.4.2. Обробка спостережень у нормальному випадку

Розглянемо на ймовірнісному просторі  $(\Omega, F, P)$  нормально розподілені  $\eta$  та  $\bar{\xi}^T = (\xi_1, \xi_2, \dots, \xi_q)$  з параметрами:

$$M\eta = m_{\eta}, D\eta = \sigma_{\eta}^2; M\bar{\xi} = \bar{m}_{\bar{\xi}}, M(\bar{\xi} - \bar{m}_{\bar{\xi}})(\bar{\xi} - \bar{m}_{\bar{\xi}})^T = \Sigma_{\bar{\xi}\bar{\xi}},$$

або скорочено  $\eta \sim N(m_{\eta}, \sigma_{\eta}^2)$ ,  $\bar{\xi} \sim N(\bar{m}_{\bar{\xi}}, \Sigma_{\bar{\xi}\bar{\xi}})$ ,  $\sigma_{\eta}^2 > 0, \Sigma_{\bar{\xi}\bar{\xi}} > 0$ .

Для уніфікації позначень будемо у подальшому використовувати позначення  $\xi_0$  для змінної  $\eta$  у разі потреби, тобто  $\xi_0 \equiv \eta$ . Аналіз наявності статистичного зв'язку між гауссівськими змінними  $\eta$  та  $\bar{\xi}$  має свою специфіку. Дійсно, індекс кореляції  $I_{\eta\bar{\xi}}$  за цих припущень виражається через множинний коефіцієнт кореляції  $r_{\eta\bar{\xi}}$  випадкових змінних  $\eta$  та  $\bar{\xi}$ . А з'ясування наявності залежності для пари змінних  $(\eta, \xi_i)$ , призводить до необхідності врахування впливу на них інших не врахованих змінних  $\xi_j$ ,  $j \neq i$ . Останнє дозволяють зробити частинні коефіцієнти кореляції, які вводяться далі.

#### 1.4.2.1. Частинний коефіцієнт кореляції та його властивості

Практика дослідження парного статистичного зв'язку між скалярними змінними  $\xi_i$  та  $\xi_j$  вказує, що звичайний коефіцієнт кореляції  $r_{\xi_i\xi_j}$  (або скорочено  $r_{ij}$ ) може неадекватно віддзеркалювати наявний зв'язок, коли існують інші змінні, наприклад  $\xi_l$ ,  $l \in \{0, 1, 2, \dots, q\} \setminus \{i, j\}$ , які мають суттєвий вплив як на  $\xi_i$ , так і  $\xi_j$ . Тому виникає бажання удосконалити звичайний коефіцієнт кореляції  $r_{\xi_i\xi_j}$  таким чином, щоб модифікована характеристика парного статистичного зв'язку була не чутлива до впливу сторонніх змінних.

Для цього запропоновано підраховувати коефіцієнт кореляції для  $\xi_i$  та  $\xi_j$  не за їх сумісним розподілом, як це робилося раніше для парного коефіцієнта кореляції  $r_{\xi_i\xi_j}$ , а за їх умовним сумісним розподілом, де в умові сторонні змінні, наприклад  $\xi_l$ ,  $l \in \{0, 1, 2, \dots, q\} \setminus \{i, j\}$ , фіксуються на певному рівні.

**Означення.** Частинним коефіцієнтом кореляції випадкових величин  $\xi_i$  та  $\xi_j$  при сторонніх змінних  $\xi_l$ ,  $l \in I$  називається числова характеристика, яка підрахована за формулою парного коефіцієнта кореляції для  $\xi_i$  та  $\xi_j$ , але за умовним сумісним розподілом випадкових величин  $\xi_i$  та  $\xi_j$  за умови, що змінні  $\xi_l$ ,  $l \in I$  набувають фіксованих значень, де  $I$  – множина індексів сторонніх змінних. Позначення –  $r_{ij(I)}$ .

**Зауваження.** Множина  $I$  не містить у собі індекси  $i$  та  $j$ .

У загальному випадку, така модифікована характеристика парного статистичного зв'язку для  $\xi_i$  та  $\xi_j$  буде залежати від постійних значень,

яких набувають сторонні змінні, і тому породжує певні труднощі в її використанні. Але з'ясувалося, що у нормальному випадку частинний коефіцієнт кореляції для випадкових величин  $\xi_i$  та  $\xi_j$  не залежить від фіксованих значень, яких набувають інші змінні  $\xi_l$ ,  $l \in I$ , і тому знайшов своє широке застосування.

Якщо випадкові величини  $\xi_0, \xi_1, \xi_2, \dots, \xi_q$  мають невироджений сумісний гауссівський розподіл ( $\xi_0 \equiv \eta$ ), а множина індексів сторонніх змінних  $I = \{0, 1, 2, \dots, q\} \setminus \{i, j\}$ , то частинний коефіцієнт кореляції  $r_{ij(I)}$  можна обчислити згідно з

$$r_{ij(I)} = -\frac{R_{ij}}{\sqrt{R_{ii}R_{jj}}},$$

де  $R_{kl}$  – алгебраїчне доповнення до елемента  $r_{kl}$  у матриці парних коефіцієнтів кореляції для змінних  $\xi_0, \xi_1, \xi_2, \dots, \xi_q$ , а саме:

$$R = \begin{pmatrix} 1 & r_{01} & r_{02} & \dots & r_{0q} \\ r_{10} & 1 & r_{12} & \dots & r_{1q} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{q0} & r_{q1} & r_{q2} & \dots & 1 \end{pmatrix}.$$

**Зауваження.** У позначенні  $r_{ij(I)} = r_{ij(\{0, 1, 2, \dots, q\} \setminus \{i, j\})}$  у дужках вказується множина індексів усіх сторонніх змінних, зафіксованих на певних рівнях.

Якщо підрахувати частинний коефіцієнт кореляції для  $\xi_0$  та  $\xi_1$ , зафіксувавши в умові лише значення змінної  $\xi_2$ , то матимемо

$$r_{01(2)} = \frac{r_{01} - r_{02}r_{12}}{\sqrt{(1 - r_{02}^2)(1 - r_{12}^2)}}.$$

Множину сторонніх змінних можна розширювати, послідовно додаючи їх по одній. Це дозволяє для обчислення частинного коефіцієнта кореляції  $r_{01(2, 3, \dots, q)}$  запропонувати рекурентну процедуру

$$r_{01(2, 3, \dots, k)} = \frac{r_{01(2, 3, \dots, k-1)} - r_{0k(2, 3, \dots, k-1)}r_{1k(2, 3, \dots, k-1)}}{\sqrt{(1 - r_{0k(2, 3, \dots, k-1)}^2)(1 - r_{1k(2, 3, \dots, k-1)}^2)}}, \quad k = \overline{3, q}.$$

**Зауваження.** З огляду на означення  $r_{ij(I)}$ , частинного коефіцієнта кореляції випадкових величин  $\xi_i$  та  $\xi_j$  при сторонніх змінних  $\xi_l$ ,  $l \in I$ , можна

стверджувати, що його властивості збігаються з властивостями парного коефіцієнта кореляції для випадкових величин  $\xi_i$  та  $\xi_j$ , тобто  $r_{ij}$ .

Для обчислення емпіричного значення частинного коефіцієнта кореляції  $\hat{r}_{ij(I)}$ ,  $I = \{0, 1, 2, \dots, q\} \setminus \{i, j\}$  можна скористатися виразом

$$\hat{r}_{ij(I)} = -\frac{\hat{R}_{ij}}{\sqrt{\hat{R}_{ii}\hat{R}_{jj}}},$$

де  $\hat{R}_{kl}$  — алгебраїчне доповнення до елемента  $\hat{r}_{kl}$  у матриці вибірових парних коефіцієнтів кореляції для змінних  $\xi_0, \xi_1, \xi_2, \dots, \xi_q$ , а саме:

$$\hat{R} = \begin{pmatrix} 1 & \hat{r}_{01} & \hat{r}_{02} & \dots & \hat{r}_{0q} \\ \hat{r}_{10} & 1 & \hat{r}_{12} & \dots & \hat{r}_{1q} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \hat{r}_{q0} & \hat{r}_{q1} & \hat{r}_{q2} & \dots & 1 \end{pmatrix}.$$

Рекурентний алгоритм для обчислення емпіричного значення частинного коефіцієнта кореляції  $\hat{r}_{01(2,3,\dots,q)}$  набуває вигляду

$$\hat{r}_{01(2,3,\dots,k)} = \frac{\hat{r}_{01(2,3,\dots,k-1)} - \hat{r}_{0k(2,3,\dots,k-1)}\hat{r}_{1k(2,3,\dots,k-1)}}{\sqrt{(1-\hat{r}_{0k(2,3,\dots,k-1)}^2)(1-\hat{r}_{1k(2,3,\dots,k-1)}^2)}}, \quad k = \overline{3, q},$$

з початковим значенням

$$\hat{r}_{01(2)} = \frac{\hat{r}_{01} - \hat{r}_{02}\hat{r}_{12}}{\sqrt{(1-\hat{r}_{02}^2)(1-\hat{r}_{12}^2)}}.$$

Нехай  $I$  — деяка множина з  $m$  індексів сторонніх змінних, які зафіксовані на деякому рівні при обчисленні частинного коефіцієнта кореляції  $r_{ij(I)}$  для змінних  $\xi_i$  та  $\xi_j$ . (Очевидно, що множина  $I$  не містить у собі індекси  $i$  та  $j$ ).

Методика аналізу статистичного зв'язку між змінними  $\xi_i$  та  $\xi_j$  за допомогою вибірового значення частинного коефіцієнта кореляції  $\hat{r}_{ij(I)}$  збігається з відповідною процедурою для значення парного коефіцієнта кореляції  $\hat{r}_{ij}$  в усьому, крім етапу перевірки цієї характеристики на значиме відхилення від нуля, у який потрібно внести незначні корекції.

Дійсно, якщо  $\begin{cases} \hat{r}_{ij(I)} \neq 0 \\ |\hat{r}_{ij(I)}| \neq 1 \end{cases}$ , то для з'ясування суттєвості статистичного

зв'язку між змінними  $\xi_i$  та  $\xi_j$  потрібно здійснити перевірку на значимість частинного коефіцієнта кореляції  $r_{ij(I)}$ , а саме, перевірити гіпотезу:

$$H_0 : r_{ij(I)} = 0$$

з деяким рівнем значущості  $\alpha > 0$ .

Виявляється, що при справедливості гіпотези  $H_0$ , статистика

$$t(n, m) = \frac{\sqrt{n-m-2} \hat{r}_{ij(I)}}{\sqrt{1-\hat{r}_{ij(I)}^2}}$$

має  $t$ -розподіл Стюдента з  $(n-m-2)$  ступенями свободи. Тоді область прийняття гіпотези  $H_0$  набуває вигляду:

$$|t(n, m)| < \frac{t_{\alpha}(n-m-2)}{2},$$

де  $t_{\alpha}(n)$  —  $100\alpha$  % -на точка  $t$ -розподілу Стюдента з  $n$  ступенями свободи.

Якщо гіпотеза  $H_0$  відхиляється, то вважається, що статистичний зв'язок між змінними  $\xi_i$  та  $\xi_j$  є істотним, у протилежному випадку — не істотним.

#### 1.4.2.2. Множинний коефіцієнт кореляції та його властивості

Припустимо, що на ймовірнісному просторі  $(\Omega, F, P)$  задані випадкова величина  $\eta$  і випадковий вектор  $\bar{\xi}^T = (\xi_1, \xi_2, \dots, \xi_q)$ , які мають сумісний нормальний розподіл:

$$\eta \sim N(m_{\eta}, \sigma_{\eta}^2), \quad \bar{\xi} \sim N(\bar{m}_{\bar{\xi}}, \Sigma_{\bar{\xi}\bar{\xi}}), \quad \sigma_{\eta}^2 > 0, \Sigma_{\bar{\xi}\bar{\xi}} > 0;$$

$$M(\eta - m_{\eta})(\bar{\xi} - \bar{m}_{\bar{\xi}})^T = \Sigma_{\eta\bar{\xi}}.$$

**Означення.** Нехай  $\eta$  і  $\bar{\xi}$  випадкові величина та вектор відповідно. Тоді *множинним коефіцієнтом кореляції*  $\eta$  і  $\bar{\xi}$  називається парний коефіцієнт кореляції випадкових величин  $\eta$  і лінійної функції від  $\bar{\xi}$ , яка є найкращою в середньоквадратичному розумінні апроксимацією для  $\eta$ . Позначення  $r_{\eta\bar{\xi}}$ .

**Зауваження.** У нормальному випадку такою лінійною функцією є функція регресії  $f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x})$ , тому  $r_{\eta\bar{\xi}} = r_{\eta f(\bar{\xi})}$ .

Крім цього для функцій  $f(\bar{x})$  і  $g(\bar{x}) = D(\eta / \bar{\xi} = \bar{x})$  має місце твердження.

**Т е о р е м а .** Нехай пара  $(\eta, \bar{\xi})$  має сумісний гауссівський розподіл, причому

$$\eta \sim N(m_\eta, \sigma_\eta^2), \quad \bar{\xi} \sim N(\bar{m}_\xi, \Sigma_{\bar{\xi}\bar{\xi}}), \quad \sigma_\eta^2 > 0, \Sigma_{\bar{\xi}\bar{\xi}} > 0;$$

$$M(\eta - m_\eta)(\bar{\xi} - \bar{m}_\xi)^T = \Sigma_{\eta\bar{\xi}}.$$

Тоді справедливо

$$f(\bar{x}) = m_\eta + \Sigma_{\eta\bar{\xi}} \Sigma_{\bar{\xi}\bar{\xi}}^{-1} (\bar{x} - \bar{m}_\xi), \quad g(\bar{x}) = \sigma_\eta^2 (1 - r_{\eta\bar{\xi}}^2).$$

**Н а с л і д о к .** Якщо додатково припустити, що  $\xi_1, \xi_2, \dots, \xi_q$  неко-рельовані, то

$$f(\bar{x}) = m_\eta + \sum_{i=1}^q r_{\eta\xi_i} \frac{\sigma_\eta}{\sigma_{\xi_i}} (x_i - m_{\xi_i}),$$

де  $\bar{x}^T = (x_1, x_2, \dots, x_q)$ ,  $\sigma_{\xi_i}^2 = D\xi_i$ ,  $i = \overline{1, q}$ .

Наступне твердження дозволяє встановити зв'язок між множинним коефіцієнтом кореляції  $r_{\eta\bar{\xi}}$  та індексом кореляції  $I_{\eta\bar{\xi}}$ .

**Т е о р е м а .** Припустимо, що  $\eta$  та  $\bar{\xi}$  нормально розподілені

$$\eta \sim N(m_\eta, \sigma_\eta^2), \quad \bar{\xi} \sim N(\bar{m}_\xi, \Sigma_{\bar{\xi}\bar{\xi}}), \quad \sigma_\eta^2 > 0, \Sigma_{\bar{\xi}\bar{\xi}} > 0;$$

$$M(\eta - m_\eta)(\bar{\xi} - \bar{m}_\xi)^T = \Sigma_{\eta\bar{\xi}}.$$

Тоді індекс кореляції  $I_{\eta\bar{\xi}} = |r_{\eta\bar{\xi}}|$ , а коефіцієнт детермінації  $I_{\eta\bar{\xi}}^2 = r_{\eta\bar{\xi}}^2$ .

**Зауваження.** У гауссівському випадку для парного/множинного статистичного зв'язку індекс кореляції збігається з модулем, а коефіцієнт детермінації з квадратом парного/множинного коефіцієнта кореляції.

**Доведення.** Скористаємося попередньою теоремою для обчислення індексу кореляції  $I_{\eta\bar{\xi}}$ . Оскільки  $D\eta = Df(\bar{\xi}) + Mg(\bar{\xi})$ , можна стверджувати, що

$$I_{\eta\bar{\xi}} = \sqrt{\frac{Df(\bar{\xi})}{D\eta}} = \sqrt{\frac{\sigma_\eta^2 - Mg(\bar{\xi})}{\sigma_\eta^2}} = \sqrt{\frac{\sigma_\eta^2 - \sigma_\eta^2 (1 - r_{\eta\bar{\xi}}^2)}{\sigma_\eta^2}} = |r_{\eta\bar{\xi}}|.$$

Що й треба було довести.

Таким чином, у нормальному випадку як характеристику множинного статистичного зв'язку можна використовувати замість індексу кореляції  $I_{\eta\bar{\xi}}$  множинний коефіцієнт кореляції  $r_{\eta\bar{\xi}}$ .

Остання теорема дозволяє легко перекоонатись у справедливості властивостей  $r_{\eta\bar{\xi}}$ :

- 1)  $-1 \leq r_{\eta\bar{\xi}} \leq 1$ ;
- 2) якщо  $r_{\eta\bar{\xi}} = 0$ , то  $\eta$  не залежить від  $\bar{\xi}$ ;
- 3) якщо  $|r_{\eta\bar{\xi}}| = 1$ , то з ймовірністю 1 існує лінійний зв'язок між  $\eta$  та  $\bar{\xi}$ , а саме, справедливо  $\eta = m_\eta + \Sigma_{\eta\bar{\xi}} \Sigma_{\bar{\xi}\bar{\xi}}^{-1} (\bar{\xi} - \bar{m}_\xi)$ .

Для визначення коефіцієнта детермінації можна скористатися одним з таких підходів:

- 1) нехай випадкова величина  $\eta$  і випадковий вектор  $\bar{\xi}^T = (\xi_1, \xi_2, \dots, \xi_q)$  мають невироджений сумісний гауссівський розподіл. Для єдиності позначень покладемо  $\xi_0 \equiv \eta$ . Тоді

$$r_{\eta\bar{\xi}}^2 = 1 - \frac{\det(R)}{R_{00}},$$

де  $R_{kl}$  – алгебраїчне доповнення до елемента  $r_{kl}$  у матриці звичайних парних коефіцієнтів кореляції для змінних  $\xi_0, \xi_1, \xi_2, \dots, \xi_q$ , а саме:

$$R = \begin{pmatrix} 1 & r_{01} & r_{02} & \dots & r_{0q} \\ r_{10} & 1 & r_{12} & \dots & r_{1q} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{q0} & r_{q1} & r_{q2} & \dots & 1 \end{pmatrix};$$

- 2) коефіцієнт детермінації також можна визначити через частинні коефіцієнти кореляції таким чином:

$$r_{\eta\bar{\xi}}^2 = 1 - (1 - r_{01}^2) (1 - r_{02(\{1\})}^2) (1 - r_{03(\{1,2\})}^2) \dots (1 - r_{0q(\{1,2,\dots,q-1\})}^2).$$

Для використання на практиці введених характеристик множинного зв'язку потрібно вміти знаходити їх емпіричні значення. Припустимо, що доступні спостереження:

$$\eta: y_1, y_2, \dots, y_n, \quad (\xi_0 \equiv \eta), \quad \bar{\xi}: \bar{x}_1, \bar{x}_2, \dots, \bar{x}_n.$$

Тоді, враховуючи останні результати, для обчислення вибіркового значення коефіцієнта детермінації  $\hat{r}_{\eta\bar{\xi}}^2$  можна використати вираз

$$\hat{r}_{\eta\bar{\xi}}^2 = 1 - \frac{\det(\hat{R})}{\hat{R}_{00}},$$

де  $\hat{R}_{kl}$  – алгебраїчне доповнення до елемента  $\hat{r}_{kl}$  у матриці вибірових парних коефіцієнтів кореляції для змінних  $\xi_0, \xi_1, \xi_2, \dots, \xi_q$ , а саме:

$$\hat{R} = \begin{pmatrix} 1 & \hat{r}_{01} & \hat{r}_{02} & \dots & \hat{r}_{0q} \\ \hat{r}_{10} & 1 & \hat{r}_{12} & \dots & \hat{r}_{1q} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \hat{r}_{q0} & \hat{r}_{q1} & \hat{r}_{q2} & \dots & 1 \end{pmatrix}.$$

З іншого боку, якщо доступні емпіричні значення для потрібних частинних коефіцієнтів кореляції, множинний коефіцієнт детермінації  $\hat{r}_{\eta\bar{\xi}}^2$  можна підрахувати таким чином:

$$\hat{r}_{\eta\bar{\xi}}^2 = 1 - (1 - \hat{r}_{01}^2) \left( 1 - \hat{r}_{02(\{1\})}^2 \right) \left( 1 - \hat{r}_{03(\{1,2\})}^2 \right) \dots \left( 1 - \hat{r}_{0q(\{1,2,\dots,q-1\})}^2 \right).$$

Методика аналізу множинного статистичного зв'язку між випадковою змінною  $\eta$  і випадковим вектором  $\bar{\xi}$  за допомогою  $\hat{r}_{\eta\bar{\xi}}$  така:

- 1) якщо  $\hat{r}_{\eta\bar{\xi}} = 0$ , то вважаємо, що зв'язок між  $\eta$  та  $\bar{\xi}$  відсутній;
- 2) якщо  $|\hat{r}_{\eta\bar{\xi}}| = 1$ , то зв'язок між  $\eta$  та  $\bar{\xi}$  є лінійний і має представлення:  $\eta = m_{\eta} + \Sigma_{\eta\bar{\xi}} \Sigma_{\bar{\xi}\bar{\xi}}^{-1} (\bar{\xi} - m_{\bar{\xi}})$ ;

- 3) якщо  $\begin{cases} \hat{r}_{\eta\bar{\xi}} \neq 0 \\ |\hat{r}_{\eta\bar{\xi}}| \neq 1 \end{cases}$ , то потрібно перевірити гіпотезу про те, чи суттєво

відхиляється від нуля множинний коефіцієнт кореляції  $r_{\eta\bar{\xi}}$  зі статистичної точки зору, а саме, зробити перевірку його на значимість, тобто потрібно перевірити гіпотезу:

$$H_0 : r_{\eta\bar{\xi}} = 0,$$

з деяким рівнем значущості  $\alpha > 0$ . І якщо гіпотеза виявиться справедливою, то вважається, що статистичний зв'язок між  $\eta$  та

$\bar{\xi}$  ( $\bar{\xi} \in R^q$ ) є не істотним, у протилежному випадку – істотним.

Для перевірки останньої гіпотези скористаємося тим фактом, що при справедливості гіпотези  $H_0$  розподіл статистики

$$F = \frac{\hat{r}_{\eta\bar{\xi}}^2}{1 - \hat{r}_{\eta\bar{\xi}}^2} \cdot \frac{n - q - 1}{q}$$

можна наблизити  $F$ -розподілом з  $q$  і  $(n - q - 1)$  ступенями свободи. А для обчислення самої статистики  $F$  потрібно вміти обраховувати вибірове значення множинного коефіцієнта детермінації  $\hat{r}_{\eta\bar{\xi}}^2$ , але, згідно з вищенаведеними результатами, це не викликає труднощів. Останнє дозволяє область прийняття для гіпотези  $H_0$  записати у вигляді

$$F < F_{\alpha}(q, n - q - 1),$$

де  $F_{\alpha}(m, n)$  –  $100\alpha$  %-на точка  $F$ -розподілу з  $m$  та  $n$  ступенями свободи.

У підсумку, у гауссівському випадку як характеристику множинного статистичного зв'язку можна використовувати множинний коефіцієнт кореляції  $r_{\eta\bar{\xi}}$ , застосування якого на практиці є більш логічним, оскільки при набутті ним значень  $\pm 1$  можна зробити висновок про наявність лінійного зв'язку між  $\eta$  та  $\bar{\xi}$ . У загальному випадку множинний коефіцієнт кореляції не надає можливості в такій інтерпретації, тому буде необхідно звернутися до використання індексу кореляції  $I_{\eta\bar{\xi}}$ .

## 2. Кореляційний аналіз ординальних змінних

Розглянемо проблему аналізу наявності статистичного зв'язку між ординальними (порядковими) змінними. Введемо такі поняття, як ранг, ранжировка, таблиця рангів. Наведемо методику для дослідження парного зв'язку між ординальними змінними за допомогою рангових коефіцієнтів кореляції Спірмена та Кендела. Для аналізу множинного зв'язку використаємо коефіцієнт конкордації. Для всіх запропонованих характеристик статистичного зв'язку наведено основні властивості та процедури їх перевірки на значиме відхилення від нуля.

### 2.1. Ранги та таблиці рангів

Нехай необхідно провести аналіз ординальних змінних  $\eta$ ,  $\xi_1$ ,  $\xi_2$ , ...,  $\xi_q$ . У подальшому для змінної  $\eta$  (для єдиності позначень) будемо використовувати позначення  $\xi_0$ , тобто  $\xi_0 \equiv \eta$ . Їх дослідження про-

ведемо на основі доступних спостережень проявів цих змінних для деяких  $n$  об'єктів.

Спочатку проаналізуємо випадок, коли за кожною ординальною змінною (властивістю)  $\xi_i$  її прояви в різних об'єктів не однакові ( $i = \overline{0, q}$ ). Розглянемо прояви конкретної змінної  $\xi_i$  для  $n$  об'єктів. Оскільки в різних об'єктів прояв властивості  $\xi_i$  інший, то всі об'єкти можна впорядкувати в міру спадання прояву змінної  $\xi_i$ . Це дозволяє кожному  $k$ -му об'єкту зіставляти натуральне число  $x_k^{(i)}$  – номер місця, яке  $k$ -й об'єкт зайняв у цій упорядкованій послідовності.

**Означення.** Рангом  $k$ -го об'єкта за ординальною змінною  $\xi_i$  називається натуральне число  $x_k^{(i)}$ , яке вказує номер місця, яке зайняв  $k$ -й об'єкт у послідовності всіх досліджуваних об'єктів, упорядкованій у міру спадання прояву властивості  $\xi_i$ ,  $k = \overline{1, n}$ ,  $i = \overline{0, q}$ .

Зрозуміло, що для якісного визначення рангів об'єктів потрібно залучати експертів відповідного профілю.

**Означення.** Ранжировкою за ординальною змінною  $\xi_i$  називається вектор-стовпчик  $x^{(i)}$ ,  $k$ -та компонента якого дорівнює рангу  $k$ -го об'єкта за ординальною змінною  $\xi_i$ , тобто  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$ .

У розглянутому випадку, ранжировки для змінних, які потрібно дослідити, будуть являти собою вектор-стовпчики, утворені внаслідок деяких перестановок компонент у вектор-стовпчику  $(1, 2, \dots, n)^T$ .

**Означення.** Ранжировка  $y = (y_1, y_2, \dots, y_n)^T$  називається протилежною ранжировкою до ранжировки  $x = (x_1, x_2, \dots, x_n)^T$ , якщо справедливо  $y_i = n + 1 - x_i$ ,  $i = \overline{1, n}$ .

Наприклад: протилежною ранжировкою до ранжировки  $(5, 1, 4, 2, 3)^T$  є ранжировка  $(1, 5, 2, 4, 3)^T$ .

Розглянемо тепер більш загальний випадок, коли існує ординальна змінна для якої існує група об'єктів, яка складається з не менше, ніж двох членів, для яких прояв цієї ординальної змінної однаковий. Тоді маємо справу з наявністю груп нерозрізних об'єктів за змінною, яка аналізується. Узагальнимо для цього випадку поняття рангу.

**Означення.** Зв'язаним рангом об'єкта з групи нерозрізних об'єктів за ординальною змінною, називається середнє арифметичне номерів місць, які припали на цю групу об'єктів.

Наприклад, якщо на групу нерозрізних об'єктів за ординальною змінною, припали місця з номерами 3, 4, 5, 6, то кожному з цих об'єктів присвоюється зв'язаний ранг, який дорівнює

$$\frac{3 + 4 + 5 + 6}{4} = 4,5.$$

Аналогічно обчислюється зв'язаний ранг для об'єктів з інших груп нерозрізних об'єктів. Приходимо до того, що в загальному випадку ранг, що віддзеркалює степінь прояву властивості, не завжди буде натуральним числом.

На основі рангів, які присвоєні об'єктам, за всіма ординальними змінними формується таблиця рангів, у якій стовпчик з номером  $i$  являє собою не що інше, як ранжировку  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$  за ординальною змінною  $\xi_i$ ,  $i = \overline{0, q}$ , а саме, таблицю такого вигляду:

№ об'єкта \ № змінної	0	1	2	...	$i$	...	$q$
1	$x_1^{(0)}$	$x_1^{(1)}$	$x_1^{(2)}$	...	$x_1^{(i)}$	...	$x_1^{(q)}$
2	$x_2^{(0)}$	$x_2^{(1)}$	$x_2^{(2)}$	...	$x_2^{(i)}$	...	$x_2^{(q)}$
...	...	...	...	...	...	...	...
$k$	$x_k^{(0)}$	$x_k^{(1)}$	$x_k^{(2)}$	...	$x_k^{(i)}$	...	$x_k^{(q)}$
...	...	...	...	...	...	...	...
$n$	$x_n^{(0)}$	$x_n^{(1)}$	$x_n^{(2)}$	...	$x_n^{(i)}$	...	$x_n^{(q)}$

Заповнена таблиця є основою при проведенні кореляційного аналізу ординальних змінних. Поява в ній значень зв'язаних рангів буде вимагати модифікації відповідних процедур для дослідження парних і множинних зв'язків.

З огляду на введені поняття, стає зрозумілим, чому розділ математики з досліджень істотності статистичного зв'язку між ординальними змінними називають також аналізом рангових кореляцій.

## 2.2. Аналіз парних рангових кореляцій

Розглянемо дві ординальні змінні. Нехай для них щодо кожного об'єкта відомий його ранг, тобто доступні відповідні ранжировки. Необхідно, спираючись на цю інформацію, з'ясувати істотність статистичного зв'язку між цими порядковими змінними. Пропонується використовувати рангові коефіцієнти кореляції. Розглянемо найбільш вживані з них.

### 2.2.1. Ранговий коефіцієнт кореляції Спірмена

Для аналізу статистичного зв'язку між двома ординальними змінними  $\xi_i$  та  $\xi_j$  з відомими для них ранжировками  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$  та  $x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)})^T$ , відповідно, К. Спірмен запропонував оригінальну характеристику парного зв'язку.

Спочатку розглянемо ситуацію, коли групи об'єктів з однаковим проявом властивостей  $\xi_i$  та  $\xi_j$  відсутні.

**Означення.** Ранговим коефіцієнтом кореляції Спірмена ординальних змінних  $\xi_i$  та  $\xi_j$  з ранжировками  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$  та  $x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)})^T$ , відповідно, називається характеристика, яка обчислюється таким чином:

$$\hat{\tau}_{ij}^{(S)} = 1 - \frac{\|x^{(i)} - x^{(j)}\|^2}{\frac{n^3 - n}{6}},$$

$$\text{де } \|x^{(i)} - x^{(j)}\|^2 = \sum_{k=1}^n (x_k^{(i)} - x_k^{(j)})^2.$$

Цей показник має такі властивості:

- 1)  $-1 \leq \hat{\tau}_{ij}^{(S)} \leq 1$ ;
- 2) якщо  $\hat{\tau}_{ij}^{(S)} = 0$ , то зв'язок між  $i$ -ою та  $j$ -ою змінними відсутній;
- 3) якщо  $\hat{\tau}_{ij}^{(S)} = 1$ , то ранжировки  $x^{(i)}$  та  $x^{(j)}$  рівні;

4) якщо  $\hat{\tau}_{ij}^{(S)} = -1$ , то ранжировки  $x^{(i)}$  та  $x^{(j)}$  протилежні;

5)  $\hat{\tau}_{ij}^{(S)} = \hat{\tau}_{ji}^{(S)}$ .

У загальному випадку, коли наявні групи нерозрізних об'єктів принаймні за однією зі змінних  $\xi_i$  чи  $\xi_j$ , використовується *модифікований ранговий коефіцієнт кореляції Спірмена*:

$$\hat{\tau}_{ij}^{(S)} = \frac{\frac{n^3 - n}{6} - \|x^{(i)} - x^{(j)}\|^2 - T^{(i)} - T^{(j)}}{\sqrt{\left[\frac{n^3 - n}{6} - 2T^{(i)}\right] \left[\frac{n^3 - n}{6} - 2T^{(j)}\right]}},$$

де  $T^{(k)} = \frac{1}{12} \sum_{l=1}^{m^{(k)}} \left[ \left( n_l^{(k)} \right)^3 - n_l^{(k)} \right]$ ;  $m^{(k)}$  – кількість груп нерозрізних об'єктів за змінною  $\xi_k$ ;  $n_l^{(k)}$  – кількість об'єктів, які ввійшли в  $l$ -ту групу нерозрізних об'єктів за змінною  $\xi_k$ ;  $k = i, j$ .

**Зауваження.** Коли відсутні групи нерозрізних об'єктів за змінними  $\xi_i$  та  $\xi_j$ , справедливим є  $T^{(i)} = T^{(j)} = 0$ , оскільки  $m^{(i)} = m^{(j)} = n$ , а  $n_l^{(i)} = n_l^{(j)} = 1$ ,  $l = \overline{1, n}$ . Звідси випливає, що  $\hat{\tau}_{ij}^{(S)}$  буде дорівнювати  $\hat{\tau}_{ij}^{(S)}$ .

Методика аналізу статистичного зв'язку між ординальними змінними  $\xi_i$  та  $\xi_j$  за допомогою рангового коефіцієнта кореляції Спірмена у випадках, коли він набуває одне зі значень  $-1, 0, 1$  очевидна і впливає з його властивостей. А коли

$$\left\{ \begin{array}{l} \hat{\tau}_{ij}^{(S)} \neq 0 \\ \left| \hat{\tau}_{ij}^{(S)} \right| \neq 1 \end{array} \right\},$$

потрібно з'ясувати, чи суттєво відхиляється від нуля ранговий коефіцієнт кореляції Спірмена із статистичної точки зору, тобто здійснити *перевірку його на значимість*, а саме перевірити гіпотезу

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

з деяким рівнем значущості  $\alpha > 0$ . Якщо гіпотеза буде відхилена, то вважається, що зв'язок між ординальними змінними  $\xi_i$  та  $\xi_j$  зі статистичної точки зору є істотним, у протилежному випадку – не суттєвим.

При  $n = 4, 10$  таку перевірку можна здійснити за допомогою спеціальних таблиць. А при  $n \geq 11$  область прийняття гіпотези набуває вигляду

$$\left| \frac{\sqrt{n-2} \hat{\tau}_{ij}^{(S)}}{\sqrt{1 - (\hat{\tau}_{ij}^{(S)})^2}} \right| < t_{\alpha}(n-2),$$

де  $t_{\alpha}(n)$  – 100 $\alpha$  %-на точка  $t$ -розподілу Стюдента з  $n$  ступенями свободи.

### 2.2.2. Ранговий коефіцієнт кореляції Кендела

Розглянемо ще одну характеристику, яку у свій час запропонував М. Кендел для аналізу статистичного зв'язку між двома ординальними змінними  $\xi_i$  та  $\xi_j$  з доступними для них ранжировками

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T \text{ та } x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)})^T.$$

Почнемо з випадку, коли відсутні групи об'єктів з однаковим проявом властивостей  $\xi_i$  та  $\xi_j$ .

**Означення.** Ранговим коефіцієнтом кореляції Кендела ординальних змінних  $\xi_i$  та  $\xi_j$  з ранжировками  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$  і

$x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)})^T$ , відповідно, називається характеристика, що визначається таким чином:

$$\hat{\tau}_{ij}^{(K)} = 1 - \frac{v(x^{(i)}, x^{(j)})}{\frac{n(n-1)}{4}},$$

де  $v(x^{(i)}, x^{(j)})$  – мінімальна кількість обмінів сусідніх компонент у ранжировці  $x^{(i)}$ , яка приводить її до ранжировки  $x^{(j)}$ .

Властивості рангового коефіцієнта кореляції Кендела збігаються з властивостями коефіцієнта Спірмена і мають вигляд:

$$1) \quad -1 \leq \hat{\tau}_{ij}^{(K)} \leq 1;$$

$$2) \quad \text{якщо } \hat{\tau}_{ij}^{(K)} = 0, \text{ то зв'язок між } i\text{-ою та } j\text{-ою змінними відсутній};$$

$$3) \quad \text{якщо } \hat{\tau}_{ij}^{(K)} = 1, \text{ то ранжировки } x^{(i)} \text{ та } x^{(j)} \text{ рівні};$$

$$4) \quad \text{якщо } \hat{\tau}_{ij}^{(K)} = -1, \text{ то ранжировки } x^{(i)} \text{ та } x^{(j)} \text{ протилежні};$$

$$5) \quad \hat{\tau}_{ij}^{(K)} = \hat{\tau}_{ji}^{(K)}.$$

За наявності груп нерозрізних об'єктів принаймні за однією зі змінних  $\xi_i$  чи  $\xi_j$ , ранговий коефіцієнт кореляції Кендела корегується й використовується у вигляді модифікованого рангового коефіцієнта кореляції Кендела:

$$\hat{\tau}_{ij}^{(K)} = \frac{\hat{\tau}_{ij}^{(K)} - \frac{2(U^{(i)} + U^{(j)})}{n(n-1)}}{\sqrt{\left(1 - \frac{2U^{(i)}}{n(n-1)}\right) \left(1 - \frac{2U^{(j)}}{n(n-1)}\right)}},$$

де  $U^{(k)} = \frac{1}{2} \sum_{l=1}^{m^{(k)}} n_l^{(k)} (n_l^{(k)} - 1)$ ;  $m^{(k)}$  – кількість груп нерозрізних об'єктів

за змінною  $\xi_k$ ;  $n_l^{(k)}$  – кількість об'єктів, які ввійшли в  $l$ -ту групу нерозрізних об'єктів за змінною  $\xi_k$ ;  $k = i, j$ .

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

Аналіз статистичного зв'язку між ординальними змінними  $\xi_i$  та  $\xi_j$  у випадках, коли ранговий коефіцієнт кореляції Кендела набуває одного із значень  $-1, 0, 1$ , впливає з вищенаведених його властивостей. А для прийняття рішення, коли

$$\left\{ \begin{array}{l} \hat{\tau}_{ij}^{(K)} \neq 0 \\ \left| \hat{\tau}_{ij}^{(K)} \right| \neq 1 \end{array} \right\},$$

потрібно звернутися до перевірки на значимість рангового коефіцієнта кореляції Кендела, тобто перевірити гіпотезу

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



з деяким рівнем значущості  $\alpha > 0$ .

При  $n = 4, 10$  така перевірка здійснюється за допомогою спеціальних таблиць. А при  $n \geq 11$  аналіз показав, що область прийняття гіпотези набуває такого представлення:

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

де  $u_{\beta}$  – 100 $\beta$  %-на точка нормального розподілу з параметрами 0 і 1.

Зв'язок між ординальними змінними  $\xi_i$  та  $\xi_j$  із статистичної точки зору вважається суттєвим при відхиленні гіпотези  $H_0$ , а в протилежному разі його слід сприймати як не істотний.

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

### 2.3. Коефіцієнт конкордації – характеристика множинних статистичних зв'язків для ординальних змінних

Перейдемо до аналізу наявності статистичного зв'язку між кількома ординальними змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ , причому нехай для цих змінних

доступні ранжировки  $x^{(i_j)} = (x_1^{(i_j)}, x_2^{(i_j)}, \dots, x_n^{(i_j)})^T$ ,  $j = \overline{1, m}$ ,  $2 \leq m \leq q+1$ ,

де  $n$  – кількість об'єктів, що досліджуються. Для розв'язання цієї проблеми М. Кендел запропонував використовувати спеціальний показник.

Розглянемо спочатку випадок, коли відсутні групи об'єктів з однако-вим проявом у змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ ,  $2 \leq m \leq q+1$ . Введемо позначення

$$\bar{\xi} = (\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m})^T.$$

**Означення.** Коефіцієнтом конкордації ординальних змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  з ранжировками  $x^{(i_1)}, x^{(i_2)}, \dots, x^{(i_m)}$  називається характеристика, що задається згідно з:

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

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

Для коефіцієнта конкордації мають місце такі властивості:

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

Додаткової уваги потребує випадок наявності груп об'єктів з однаковим проявом у змінних  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$ ,  $2 \leq m \leq q+1$ . Після відповідної корекції, коефіцієнт конкордації можна використовувати у вигляді такого модифікованого коефіцієнта конкордації:

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

де  $T^{(i_j)} = \frac{1}{12} \sum_{l=1}^{m^{(i_j)}} \left[ \left( n_l^{(i_j)} \right)^3 - n_l^{(i_j)} \right]$ ;  $m^{(i_j)}$  – кількість груп нерозрізних

об'єктів за змінною  $\xi_{i_j}$ ;  $n_l^{(i_j)}$  – кількість об'єктів, які ввійшли в  $l$ -ту групу нерозрізних об'єктів за змінною  $\xi_{i_j}$ ;  $j = \overline{1, m}$ ;  $2 \leq m \leq q+1$ .

Використання на практиці коефіцієнта конкордації, коли він набуває екстремальних значень 0 чи 1, з огляду на його властивості, не викликає труднощів. Зупинимось на ситуації, коли  $0 < \hat{W}_{\bar{\xi}} < 1$ . Знову звертаємось до перевірки на значиме відхилення від нуля коефіцієнта, який використовується, тобто до перевірки гіпотези

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

з деяким рівнем значущості  $\alpha > 0$ .

При  $n = 3, 7$ ,  $m = 3, 20$  потрібна перевірка здійснюється за допомогою спеціальних таблиць. А при  $n \geq 8$  можна скористатися тим фактом, що

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

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

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

Зв'язок між ординальними змінними  $\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_m}$  із статистичної точки зору будемо вважати суттєвим при відхиленні гіпотези  $H_0$ , а в протилежному випадку вважаємо не істотним.

### 3. Дослідження наявності зв'язку між номінальними змінними

Кореляційний аналіз номінальних (класифікаційних) змінних має свою специфіку, оскільки значення, яких вони набувають, не є впорядкованими. При роботі з двома змінними, які мають скінченні множини значень, на базі спостережень формується таблиця спряженості, яка є основою для побудови цілої низки характеристик статистичного зв'язку.

Деякі з них базуються на статистиці квадратичної спряженості, що й буде продемонстровано далі. З іншого боку, при побудові потрібного показника можна скористатися поняттям ентропії і на його основі побудувати інформаційну міру зв'язку. Аналіз цих характеристик дозволяє зробити висновок відносно істотності статистичного зв'язку між номінальними змінними, які досліджуються.

#### 3.1. Таблиця спряженості

Розглянемо номінальні змінні  $\eta$  та  $\xi$ . Нехай множина значень у змінної  $\eta$  складається з  $r_1$  градацій, а у змінної  $\xi$  — з  $r_2$  градацій відповідно. Припустимо, що проведено всього  $n$  спостережень над деякими об'єктами. Вважатимемо, що у  $n_{ij}$  випадках  $\eta$  набула значень своєї  $i$ -ої градації, а  $\xi$  —  $j$ -ої градації. Тоді результати всіх спостережень можна занести до таблиці, у якій у комірку на перетині  $i$ -го рядка та  $j$ -го стовпчика заноситься значення  $n_{ij}$ ,  $i = \overline{1, r_1}$ ,  $j = \overline{1, r_2}$ . Побудована таким чином таблиця називається *таблицею спряженості* змінних  $\eta$  та  $\xi$  і матиме такий вигляд:

$\eta \backslash \xi$	1	2	...	$j$	...	$r_2$	
1	$n_{11}$	$n_{12}$	...	$n_{1j}$	...	$n_{1r_2}$	$n_{1\bullet}$
2	$n_{21}$	$n_{22}$	...	$n_{2j}$	...	$n_{2r_2}$	$n_{2\bullet}$
$\vdots$	$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$	$\vdots$
$i$	$n_{i1}$	$n_{i2}$	...	$n_{ij}$	...	$n_{ir_2}$	$n_{i\bullet}$
$\vdots$	$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$	$\vdots$
$r_1$	$n_{r_1 1}$	$n_{r_1 2}$	...	$n_{r_1 j}$	...	$n_{r_1 r_2}$	$n_{r_1 \bullet}$
	$n_{\bullet 1}$	$n_{\bullet 2}$	...	$n_{\bullet j}$	...	$n_{\bullet r_2}$	$n$

де  $n_{i\bullet} = \sum_{j=1}^{r_2} n_{ij}$ ,  $n_{\bullet j} = \sum_{i=1}^{r_1} n_{ij}$ ,  $n = \sum_{i=1}^{r_1} n_{i\bullet} = \sum_{j=1}^{r_2} n_{\bullet j}$ . Причому, в останніх сто-

впчику та рядку таблиці підраховані суми значень по рядках і стовпчиках таблиці відповідно.

Таблиці спряженості можуть бути також використані при аналізі й інших класів змінних, які набувають скінченні множини значень, а саме: порядкових, дискретних кількісних зі скінченню множиною значень, неперервних кількісних для яких провели процедуру групування.

Використовуючи таблицю спряженості, можна конструювати показники, потрібні для дослідження наявності статистичного зв'язку між номінальними змінними.

#### 3.2. Квадратична спряженість і характеристики парного зв'язку на її основі

Для подальшої побудови характеристик парного статистичного зв'язку введемо статистику *квадратичної спряженості* для змінних  $\eta$  та  $\xi$ :

$$\hat{\chi}_{\eta\xi}^2 = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \frac{\left( n_{ij} - \frac{n_{i\bullet} n_{\bullet j}}{n} \right)^2}{\frac{n_{i\bullet} n_{\bullet j}}{n}} = n \left[ \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \left( \frac{n_{ij}^2}{n_{i\bullet} n_{\bullet j}} \right) - 1 \right].$$

Очевидно, що ця статистика набуває значень нуль тоді й тільки тоді, коли має місце

$$\frac{n_{ij}}{n} = \frac{n_{i\bullet}}{n} \frac{n_{\bullet j}}{n}, \quad i = \overline{1, r_1}, j = \overline{1, r_2}.$$

Але ця умова є не що інше, як умова незалежності змінних  $\eta$  та  $\xi$ , але записана у вибіркового представленні.

Останнє було використано різними авторами при побудові низки характеристик парного статистичного зв'язку на базі статистики квадратичної спряженості для змінних  $\eta$  та  $\xi$ , а саме:

1) коефіцієнт спряженості – квадратний корінь з середньої квадратичної спряженості  $\hat{\phi}_{\eta\xi} = \sqrt{\frac{\hat{\chi}_{\eta\xi}^2}{n}}$ ;

2) коефіцієнт спряженості Пірсона  $\hat{P}_{\eta\xi} = \sqrt{\frac{\hat{\chi}_{\eta\xi}^2}{n + \hat{\chi}_{\eta\xi}^2}}$ ;

3) коефіцієнт спряженості Чупрова  $\hat{T}_{\eta\xi} = \sqrt{\frac{\hat{\chi}_{\eta\xi}^2}{n\sqrt{(r_1 - 1)(r_2 - 1)}}}$ ;

4) коефіцієнт спряженості Крамера  $\hat{C}_{\eta\xi} = \sqrt{\frac{\hat{\chi}_{\eta\xi}^2}{n \cdot \min(r_1 - 1, r_2 - 1)}}$ .

Для цих характеристик справедливі такі властивості:

1)  $\hat{\phi}_{\eta\xi}, \hat{P}_{\eta\xi}, \hat{T}_{\eta\xi}, \hat{C}_{\eta\xi} \geq 0$ ;

2)  $\hat{P}_{\eta\xi} < 1, \hat{C}_{\eta\xi} \leq 1$ ;

3) якщо  $\begin{cases} \hat{\phi}_{\eta\xi} = 0 \\ \hat{P}_{\eta\xi} = 0 \\ \hat{T}_{\eta\xi} = 0 \\ \hat{C}_{\eta\xi} = 0 \end{cases}$ , то  $\eta$  не залежить від  $\xi$ .

При практичному використанні цих коефіцієнтів, коли вони набувають нульового значення, робимо висновок, що  $\eta$  не залежить від  $\xi$ , а якщо вони набувають значення відмінне від нуля, то потрібно звернутися до перевірки їх на значимість, але з огляду на їх структуру та залежність від квадратичної спряженості достатньо перевірити гіпотезу

$$H_0: \chi_{\eta\xi}^2 = 0$$

з деяким рівнем значущості  $\alpha > 0$ .

Враховуючи той факт, що розподіл статистики  $\hat{\chi}_{\eta\xi}^2$  можна наблизити  $\chi^2$ -розподілом з  $(r_1 - 1)(r_2 - 1)$  ступенями свободи за умови справедливості  $H_0$ , то область прийняття гіпотези набуває вигляду:

$$\hat{\chi}_{\eta\xi}^2 < \chi_{\alpha}^2((r_1 - 1)(r_2 - 1)).$$

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

Якщо остання нерівність буде справедлива, то будемо стверджувати, що  $\eta$  не залежить від  $\xi$ , у протилежному випадку йдеться про істотність зв'язку між  $\eta$  та  $\xi$  із статистичної точки зору.

### 3.3. Інформаційна міра парного статистичного зв'язку

Інший підхід до введення характеристики статистичного зв'язку базується на поняттях з теорії інформації. Нехай змінна  $\eta$  набуває своїх значень  $y_i$  з ймовірностями  $p_{\eta}(y_i)$ , змінна  $\xi$  – значень  $x_j$  з ймовірностями  $p_{\xi}(x_j)$ , а пара  $(\eta, \xi)$  набуває значень  $(y_i, x_j)$  з ймовірностями  $p_{\eta\xi}(y_i, x_j)$ ,  $i = \overline{1, r_1}, j = \overline{1, r_2}$ .

**Означення.** Ентропією  $\eta$  називається величина, яка обчислюється як:

$$H_{\eta} = - \sum_{i=1}^{r_1} p_{\eta}(y_i) \ln(p_{\eta}(y_i)) = -M \ln(p_{\eta}(\eta)).$$

У свою чергу, ентропію  $\xi$  можна підрахувати згідно з:

$$H_{\xi} = - \sum_{j=1}^{r_2} p_{\xi}(x_j) \ln(p_{\xi}(x_j)) = -M \ln(p_{\xi}(\xi)).$$

**Означення.** Ентропією  $\eta$  та  $\xi$  називається значення, яке задається виразом

$$H_{\eta\xi} = - \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} p_{\eta\xi}(y_i, x_j) \ln(p_{\eta\xi}(y_i, x_j)) = -M \ln(p_{\eta\xi}(\eta, \xi)).$$

**Означення.** Інформаційна міра зв'язку  $\eta$  та  $\xi$  визначається згідно з

$$N_{\eta\xi} = H_{\eta} + H_{\xi} - H_{\eta\xi}.$$

Основні властивості інформаційної міри зв'язку  $\eta$  та  $\xi$ :

1)  $N_{\eta\xi} \geq 0$ ;

2) якщо  $N_{\eta\xi} = 0$ , то  $\eta$  не залежить від  $\xi$ ;

3)  $N_{\eta\eta} = H_{\eta}$ .

Спираючись на доступні значення з таблиці спряженості, підрахуємо вибіркове значення для інформаційної міри зв'язку  $\eta$  та  $\xi$ . Оскільки

$$\hat{H}_\eta = -\sum_{i=1}^{r_1} \frac{n_{i\bullet}}{n} \ln\left(\frac{n_{i\bullet}}{n}\right), \quad \hat{H}_\xi = -\sum_{j=1}^{r_2} \frac{n_{\bullet j}}{n} \ln\left(\frac{n_{\bullet j}}{n}\right), \quad \hat{H}_{\eta\xi} = -\sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \frac{n_{ij}}{n} \ln\left(\frac{n_{ij}}{n}\right),$$

ТО

$$\begin{aligned} \hat{N}_{\eta\xi} &= -\sum_{i=1}^{r_1} \frac{n_{i\bullet}}{n} \ln\left(\frac{n_{i\bullet}}{n}\right) - \sum_{j=1}^{r_2} \frac{n_{\bullet j}}{n} \ln\left(\frac{n_{\bullet j}}{n}\right) + \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \frac{n_{ij}}{n} \ln\left(\frac{n_{ij}}{n}\right) = \\ &= \frac{1}{n} \left\{ \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} n_{ij} \ln(n_{ij}) - \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} n_{ij} \ln(n) - \sum_{i=1}^{r_1} n_{i\bullet} \ln(n_{i\bullet}) + \sum_{i=1}^{r_1} n_{i\bullet} \ln(n) - \right. \\ &\quad \left. - \sum_{j=1}^{r_2} n_{\bullet j} \ln(n_{\bullet j}) + \sum_{j=1}^{r_2} n_{\bullet j} \ln(n) \right\} = \\ &= \frac{1}{n} \left\{ \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} n_{ij} \ln(n_{ij}) - \sum_{i=1}^{r_1} n_{i\bullet} \ln(n_{i\bullet}) - \sum_{j=1}^{r_2} n_{\bullet j} \ln(n_{\bullet j}) + n \ln(n) \right\}. \end{aligned}$$

Коли  $\hat{N}_{\eta\xi} = 0$ , будемо приймати рішення про відсутність статистичного зв'язку між  $\eta$  та  $\xi$ , а у випадку  $\hat{N}_{\eta\xi} > 0$ , потрібно звернутися до перевірки інформаційної міри зв'язку  $\eta$  та  $\xi$  на значимість, тобто перевірити гіпотезу

$$H_0 : N_{\eta\xi} = 0$$

з деяким рівнем значущості  $\alpha > 0$ . Для цього розглянемо статистику

$$\begin{aligned} \hat{N}_{\eta\xi} &= 2n\hat{N}_{\eta\xi} - n_0 = \\ &= 2 \left[ \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} n_{ij} \ln(n_{ij}) - \sum_{i=1}^{r_1} n_{i\bullet} \ln(n_{i\bullet}) - \sum_{j=1}^{r_2} n_{\bullet j} \ln(n_{\bullet j}) + n \ln(n) \right] - n_0, \end{aligned}$$

де  $n_0$  — кількість нульових значень у таблиці спряженості.

Виявляється, що при справедливості цієї гіпотези розподіл статистики  $\hat{N}_{\eta\xi}$  можна наблизити  $\chi^2$ -розподілом з  $(r_1 - 1)(r_2 - 1)$  ступенями свободи. Тоді область прийняття гіпотези  $H_0$  набуде вигляду

$$\hat{N}_{\eta\xi} < \chi_\alpha^2((r_1 - 1)(r_2 - 1)).$$

де  $\chi_\alpha^2(n)$  — 100 $\alpha$  %-на точка  $\chi^2$ -розподілу з  $n$  ступенями свободи.

Якщо остання нерівність не справедлива, то будемо говорити про істотність статистичного зв'язку між  $\eta$  та  $\xi$ , інакше вважаємо, що зв'язок між ними не є суттєвим.

## Умовні ймовірності та математичні сподівання.

### Основні властивості

Додаток містить довідкову інформацію необхідну при використанні умовних ймовірностей і математичних сподівань. Згадаємо основні необхідні поняття.

**Означення.** Нехай  $d_R$  є сукупністю усіх інтервалів виду  $(a, b]$ ,  $a, b \in R$ . Борелівською  $\sigma$ -алгеброю  $\sigma_R$  в  $R$  називається найменша  $\sigma$ -алгебра, яка містить у собі систему  $d_R$ .

**Зауваження 1.** З означення борелівської  $\sigma$ -алгебри  $\sigma_R$  в  $R$  випливає, що в неї, крім інтервалів виду  $(a, b]$ ,  $a, b \in R$ , також входять інтервали:

$$(a, b), [a, b], [a, b], (-\infty, b), (-\infty, b], (a, \infty), [a, \infty), a, b \in R,$$

порожня множина  $(a, a)$  і множини, які складаються з однієї точки  $[a, a]$ .

**Зауваження 2.** В означенні борелівської  $\sigma$ -алгебри  $\sigma_R$  в  $R$  замість інтервалів виду  $(a, b]$  можна використовувати інтервали одного з видів:

$$(a, b), [a, b], [a, b], (-\infty, b), (-\infty, b], (a, \infty), [a, \infty), a, b \in R.$$

**Означення.** Борелівською множиною в  $R$  (множиною вимірною за Борелем в  $R$ ) називається довільна множина з борелівської  $\sigma$ -алгебри  $\sigma_R$  в  $R$ .

**Означення.** Борелівською функцією на  $R$  (функцією вимірною за Борелем на  $R$ ) називається дійсна функція задана на  $R$ , у якій прообраз довільної борелівської множини в  $R$  є борелівською множиною в  $R$ .

Аналогічно вводиться борелівська функція з областю визначення в просторі  $R^q$ ,  $q \in N$ .

**Означення.** Нехай  $d_{R^q}$  є сукупністю всіх множин виду  $\{(x_1, x_2, \dots, x_q) \in R^q : x_i \in (a_i, b_i]; a_i, b_i \in R, i = \overline{1, q}\}$ . Борелівською  $\sigma$ -алгеброю  $\sigma_{R^q}$  в  $R^q$  називається найменша  $\sigma$ -алгебра, яка містить у собі систему  $d_{R^q}$ .

**Зауваження 1.** З означення борелівської  $\sigma$ -алгебри  $\sigma_{R^q}$  в  $R^q$  випливає, що в неї також входять множини виду

$$\{(x_1, x_2, \dots, x_q) \in R^q : x_i \in B_i, i = \overline{1, q}\},$$

де  $B_1, B_2, \dots, B_q$  – борелівські множини в  $R$ .

**Зауваження 2.** В означенні борелівської  $\sigma$ -алгебри  $\sigma_{R^q}$  в  $R^q$  замість множин

$$\{(x_1, x_2, \dots, x_q) \in R^q : x_i \in (a_i, b_i]; a_i, b_i \in R, i = \overline{1, q}\}$$

можна використовувати множини виду

$$\{(x_1, x_2, \dots, x_q) \in R^q : x_i \in B_i, i = \overline{1, q}\},$$

де  $B_1, B_2, \dots, B_q$  – борелівські множини в  $R$ .

**Означення.** Борелівською множиною в  $R^q$  (множиною вимірною за Борелем в  $R^q$ ) називається довільна множина з борелівської  $\sigma$ -алгебри  $\sigma_{R^q}$  в  $R^q$ .

**Означення.** Борелівською функцією на  $R^q$  (функцією вимірною за Борелем на  $R^q$ ) називається дійсна функція, задана на  $R^q$ , у якій прообраз довільної борелівської множини в  $R$  є борелівською множиною в  $R^q$ .

Нехай на ймовірнісному просторі  $(\Omega, F, P)$  задані дійсні випадкові величини  $\eta = \eta(\omega)$  та  $\xi = \xi(\omega)$ , а  $G$  – деяка  $\sigma$ -алгебра ( $G \subset F$ ).

**Означення.** Випадкова величина  $\xi$  є вимірною відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ), якщо прообраз довільної борелівської множини в  $R$  належить  $\sigma$ -алгебрі  $G$ , тобто  $\forall A \in \sigma_R \quad \xi^{-1}(A) \in G$ .

**Означення.** Умовним математичним сподіванням невід'ємної випадкової величини  $\xi$  відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ) називається невід'ємна випадкова величина, яка позначається  $M(\xi/G)$  або  $M(\xi/G)(\omega)$  така, що

- $M(\xi/G)$  є вимірною відносно  $\sigma$ -алгебри  $G$ ,
- $\forall A \in G$  справедливо  $\int_A \xi dP = \int_A M(\xi/G) dP$ .

Нескладно з'ясувати питання існування та єдиності  $M(\xi/G)$ . Дійсно, якщо випадкова величина  $\xi$  є інтегрованою і  $\xi \geq 0$  майже всюди, тоді на вимірному просторі  $(\Omega, G)$  функція множин

$$\lambda(A) = \int_A \xi dP, \quad \forall A \in G$$

є мірою, яка абсолютно неперервна відносно ймовірнісної міри  $P$  (розглядається на вимірному просторі  $(\Omega, G)$ ). Очевидно, що міра  $\lambda$  буде скінченною та  $\sigma$  – адитивною.

Далі скористаємося відомим твердженням.

**Т е о р е м а (Радона – Никодима).** Нехай  $(\Omega, \sigma_\Omega)$  – вимірний простір, на  $\sigma$ -алгебрі  $\sigma_\Omega$  задані скінченні  $\sigma$ -адитивні міри  $\lambda$  та  $\mu$ . Якщо міра  $\lambda$  є абсолютно неперервною відносно міри  $\mu$ , то існує  $\sigma_\Omega$ -вимірна єдина (з точністю до  $\mu$ -еквівалентності) така функція  $\varphi(\cdot)$ , задана на  $\Omega$ , яка задовольняє вимоги

- $\varphi \in L_1(\Omega, \sigma_\Omega, \mu)$ , (множина функцій інтегровних на  $\Omega$  за мірою  $\mu$ );
- $\varphi(\omega) \geq 0$  майже всюди на  $\Omega$  відносно міри  $\mu$ , що справедливо

$$\lambda(A) = \int_A \varphi d\mu, \quad \forall A \in \sigma_\Omega.$$

З теореми Радона – Никодима випливає, що існує невід'ємна інтегровна  $G$ -вимірна єдина (з точністю до  $P$ -еквівалентності) функція  $M(\xi/G)$  задана на  $\Omega$  така, що

$$\lambda(A) = \int_A \xi dP = \int_A M(\xi/G) dP \quad \forall A \in G.$$

Таким чином, для невід'ємної випадкової величини  $\xi$  існує єдине (з точністю до  $P$ -еквівалентності) умовне математичне сподівання відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ):  $M(\xi/G)$ .

Якщо  $\xi$  випадкова величина, то введемо позначення для таких випадкових величин  $\xi^+ = \max(\xi, 0)$ ,  $\xi^- = -\min(\xi, 0)$ . Очевидно, що  $\xi = \xi^+ - \xi^-$ .

**Означення.** Умовним математичним сподіванням довільної випадкової величини  $\xi$  відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ) називається випадкова величина  $M(\xi/G)$  (або  $M(\xi/G)(\omega)$ ), яка визначається як:

$$M(\xi/G) = M(\xi^+/G) - M(\xi^-/G).$$



Зауваження. З означення випливає, що з ймовірністю 1 справедливо

$$\begin{cases} M(\xi^+ / G) \geq 0, \\ M(\xi^- / G) \geq 0. \end{cases}$$

**Означення.**  $M(\xi / G)$  існує (або визначено), якщо з ймовірністю 1 справедливо  $\min(M(\xi^+ / G), M(\xi^- / G)) < \infty$ .

**Означення.** Умовною ймовірністю події  $B$  ( $B \in F$ ) відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ), яка позначається  $P(B / G)$  (або  $P(B / G)(\omega)$ ), називається невід'ємна випадкова величина, що визначається згідно з

$$P(B / G) = M(I_B / G),$$

де  $I_B$  — індикатор множини  $B$ , тобто

$$I_B(\omega) = \begin{cases} 1, & \omega \in B, \\ 0, & \omega \notin B. \end{cases}$$

Можна використовувати й інше означення.

**Означення.** Умовною ймовірністю події  $B$  ( $B \in F$ ) відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ) називається невід'ємна випадкова величина, яка позначається  $P(B / G)$  (або  $P(B / G)(\omega)$ ) така, що:

- $P(B / G)$  є вимірною відносно  $\sigma$ -алгебри  $G$ ,
- $\forall A \in G$  справедливо  $P(A \cap B) = \int_A P(B / G) dP$ .

**Означення.**  $\sigma$ -алгебри  $G_1, G_2$  ( $G_1, G_2 \subset F$ ) незалежні, якщо

$$P(A_1 \cap A_2) = P(A_1)P(A_2), \quad \forall A_1 \in G_1, A_2 \in G_2.$$

**Означення.** Випадкові величини  $\eta$  та  $\xi$  незалежні, якщо незалежні  $\sigma$ -алгебри, які породжені випадковими величинами  $\eta$  та  $\xi$ , а саме,  $\sigma$ -алгебри  $\sigma_\eta$  та  $\sigma_\xi$ .

**Означення.**  $\eta$  не залежить від  $\sigma$ -алгебри  $G$  ( $G \subset F$ ), якщо не залежні  $\sigma$ -алгебра  $G$  і  $\sigma$ -алгебра, яка породжена випадковою величиною  $\eta$ , тобто  $\sigma_\eta$ .

Основні властивості умовних математичних сподівань випадкових величин відносно  $\sigma$ -алгебр мають такий вигляд:

- 1) нехай  $M|\eta| < \infty$ ,  $G$  —  $\sigma$ -алгебра ( $G \subset F$ ). Тоді існує  $M(\eta / G)$ ;

- 2) нехай існує  $M(\eta)$  і  $Q = \{\emptyset, \Omega\}$  — тривіальна  $\sigma$ -алгебра. Тоді

$$M(\eta / Q) = M(\eta);$$

- 3) нехай існує  $M(\eta)$  та  $\eta$  не залежить від  $\sigma$ -алгебри  $G$  ( $G \subset F$ ). Тоді  $M(\eta / G) = M(\eta)$ ;

- 4) нехай  $\eta$  — вимірна відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ), тоді  $M(\eta / G) = \eta$ . Тим паче  $M(\eta / F) = \eta$ ;

- 5) нехай  $M|\eta_1| < \infty$ ,  $M|\eta_2| < \infty$ ,  $G$  —  $\sigma$ -алгебра ( $G \subset F$ ),  $a, b \in R$ , тоді  $M(a\eta_1 + b\eta_2 / G) = aM(\eta_1 / G) + bM(\eta_2 / G)$ ;

- 6) нехай  $M|\eta| < \infty$ ,  $M|\xi\eta| < \infty$ ,  $\xi$  — вимірна відносно  $\sigma$ -алгебри  $G$  ( $G \subset F$ ), тоді  $M\{\xi\eta / G\} = \xi M(\eta / G)$ ;

- 7) нехай  $c \in R$ ,  $G$  —  $\sigma$ -алгебра ( $G \subset F$ ), тоді  $M(c / G) = c$ ;

- 8) нехай  $M|\eta| < \infty$ ,  $G_1, G_2$  —  $\sigma$ -алгебри ( $G_1 \subset G_2 \subset F$ ), тоді

$$M\{M(\eta / G_2) / G_1\} = M(\eta / G_1),$$

$$M\{M(\eta / G_1) / G_2\} = M(\eta / G_1).$$

**Означення.** Нехай  $\xi$  — випадкова величина. Розглянемо сукупність множин  $\{\omega : \xi(\omega) \in B\}$ ,  $B \in \sigma_R$ . Найменша  $\sigma$ -алгебра, що містить у собі ці множини, називається  $\sigma$ -алгеброю, яка породжена випадковою величиною  $\xi$ . Позначення —  $\sigma_\xi$ .

Аналогічне поняття вводиться для випадкового вектора.

**Означення.** Нехай випадковий вектор  $\bar{\xi} \in R^q$ . Розглянемо сукупність множин  $\{\omega : \bar{\xi}(\omega) \in B\}$ ,  $B \in \sigma_{R^q}$ . Найменша  $\sigma$ -алгебра, що містить у собі ці множини, називається  $\sigma$ -алгеброю, яка породжена випадковим вектором  $\bar{\xi}$ . Позначення —  $\sigma_{\bar{\xi}}$ .

Умовні математичне сподівання та дисперсія випадкової величини  $\eta$  зручно вводити одразу відносно випадкового вектора  $\bar{\xi}$  ( $\bar{\xi} \in R^q$ ,  $q \in N$ ), а відносно випадкової величини отримуємо, як частинний випадок.

**Означення.** Умовним математичним сподіванням випадкової величини  $\eta$  відносно випадкового вектора  $\bar{\xi}$  називається умовне математичне сподівання випадкової величини  $\eta$  відносно  $\sigma$ -алгебри, яка

породжена випадковим вектором  $\bar{\xi}$ , тобто відносно  $\sigma_{\bar{\xi}}$ . Позначення –  $M(\eta/\bar{\xi})$  (або  $M(\eta/\bar{\xi})(\omega)$ ).

**Означення.** Умовною дисперсією випадкової величини  $\eta$  відносно випадкового вектора  $\bar{\xi}$  називається випадкова величина, яка позначається  $D(\eta/\bar{\xi})$  (або  $D(\eta/\bar{\xi})(\omega)$ ) і визначається як:

$$D(\eta/\bar{\xi}) = M \left\{ \left[ \eta - M(\eta/\bar{\xi}) \right]^2 / \bar{\xi} \right\}.$$

**Наслідок.** Для умовної дисперсії випадкової величини  $\eta$  відносно випадкового вектора  $\bar{\xi}$  справедливо

$$D(\eta/\bar{\xi}) = M \left\{ \left[ \eta - M(\eta/\bar{\xi}) \right]^2 / \bar{\xi} \right\} = M(\eta^2/\bar{\xi}) - \left[ M(\eta/\bar{\xi}) \right]^2.$$

Подібним чином вводиться відповідне поняття для умовної ймовірності події.

**Означення.** Умовною ймовірністю події  $B$  ( $B \in F$ ) відносно випадкового вектора  $\bar{\xi}$  називається умовна ймовірність події  $B$  відносно  $\sigma$ -алгебри, яка породжена випадковим вектором  $\bar{\xi}$ . Позначення –  $P(B/\bar{\xi})$  (або  $P(B/\bar{\xi})(\omega)$ ).

У свою чергу, відповідні властивості умовних математичних сподівань випадкової величини відносно випадкового вектора  $\bar{\xi}$  ( $\bar{\xi} \in R^q$ ) мають такий вигляд:

1) нехай  $M|\eta| < \infty$ , тоді існує  $M(\eta/\bar{\xi})$ ;

2) нехай  $M|\eta_1| < \infty$ ,  $M|\eta_2| < \infty$ ,  $a, b \in R$ , тоді

$$M(a\eta_1 + b\eta_2/\bar{\xi}) = aM(\eta_1/\bar{\xi}) + bM(\eta_2/\bar{\xi}).$$

3) нехай  $M|\eta| < \infty$ ,  $M|\varphi(\bar{\xi})\eta| < \infty$ ,  $\varphi(\cdot)$  – борелівська функція на  $R^q$ , тоді  $M\{\varphi(\bar{\xi})\eta/\bar{\xi}\} = \varphi(\bar{\xi})M(\eta/\bar{\xi})$ ;

4) нехай  $c \in R$ , тоді  $M(c/\bar{\xi}) = c$ ;

5) нехай  $M|\eta| < \infty$ , тоді  $M\{M(\eta/\bar{\xi})\} = M\eta$ .

**Теорема.** Нехай на ймовірнісному просторі  $(\Omega, F, P)$  задані дійсна випадкова величина  $\eta$  і випадковий вектор  $\bar{\xi} \in R^q$ , причому  $\eta$  є вимірною відносно  $\sigma$ -алгебри, яка породжена випадковим вектором  $\bar{\xi}$ , тобто відносно  $\sigma_{\bar{\xi}}$ . Тоді існує борелівська функція  $\varphi(\cdot)$  на  $R^q$  така, що  $\eta = \varphi(\bar{\xi})$ .

**Зауваження.** Зворотний результат очевидний. Нехай на ймовірнісному просторі  $(\Omega, F, P)$  задані дійсна випадкова величина  $\eta$  і випадковий вектор  $\bar{\xi} \in R^q$ . Припустимо, що існує борелівська функція  $\varphi(\cdot)$  на  $R^q$  така, що  $\eta = \varphi(\bar{\xi})$ , тоді  $\eta$  є вимірною відносно  $\sigma$ -алгебри, яка породжена випадковим вектором  $\bar{\xi}$ , тобто відносно  $\sigma_{\bar{\xi}}$ .

Скористаємося останньою теоремою, аналізуючи  $M(\eta/\bar{\xi})$ . Оскільки згідно з означенням  $M(\eta/\bar{\xi})$  є вимірною відносно  $\sigma$ -алгебри  $\sigma_{\bar{\xi}}$ , то існує борелівська функція  $f(\cdot)$  на  $R^q$  така, що  $M(\eta/\bar{\xi}) = f(\bar{\xi})$ . Останнє дозволяє сформулювати таке означення.

**Означення.** Для  $M(\eta/\bar{\xi})$  існує борелівська функція  $f(\cdot)$  на  $R^q$  така, що  $M(\eta/\bar{\xi}) = f(\bar{\xi})$ . Умовне математичне сподівання випадкової величини  $\eta$  відносно події  $\{\bar{\xi} = \bar{x}\}$  (або умовне математичне сподівання випадкової величини  $\eta$  за умови, що  $\bar{\xi} = \bar{x}$ ) позначається  $M(\eta/\bar{\xi} = \bar{x})$  і визначається таким чином:  $M(\eta/\bar{\xi} = \bar{x}) = f(\bar{x})$ ,  $\bar{x} \in R^q$ .

**Зауваження 1.** Функцію  $f(\bar{x}) = M(\eta/\bar{\xi} = \bar{x})$  називають функцією регресії  $\eta$  на  $\bar{\xi}$  (або функцією регресії  $\eta$  щодо  $\bar{\xi}$ ).

**Зауваження 2.** Легко бачити, що  $Mf(\bar{\xi}) = M\{M(\eta/\bar{\xi})\} = M\eta$ .

**Означення.** Умовна дисперсія випадкової величини  $\eta$  відносно події  $\{\bar{\xi} = \bar{x}\}$  (або умовна дисперсія випадкової величини  $\eta$  за умови, що  $\bar{\xi} = \bar{x}$ ) позначається  $D(\eta/\bar{\xi} = \bar{x})$  і визначається таким чином:

$$D(\eta/\bar{\xi} = \bar{x}) = M \left\{ \left[ \eta - M(\eta/\bar{\xi} = \bar{x}) \right]^2 / \bar{\xi} = \bar{x} \right\}.$$

**Н а с л і д о к .** Для умовної дисперсії випадкової величини  $\eta$  відносно події  $\{\bar{\xi} = \bar{x}\}$  справедливо

$$D(\eta/\bar{\xi} = \bar{x}) = M \left\{ \left[ \eta - M(\eta/\bar{\xi} = \bar{x}) \right]^2 / \bar{\xi} = \bar{x} \right\} = \\ = M(\eta^2/\bar{\xi} = \bar{x}) - \left[ M(\eta/\bar{\xi} = \bar{x}) \right]^2.$$

**Зауваження 1.** Для нотації  $D(\eta/\bar{\xi} = \bar{x})$  також буде використовуватися таке позначення:  $g(\bar{x}) = D(\eta/\bar{\xi} = \bar{x})$ .

**Зауваження 2.** Очевидно, що

$$g(\bar{\xi}) = M \left\{ \left[ \eta - M(\eta/\bar{\xi}) \right]^2 / \bar{\xi} \right\}.$$

## ДОДАТОК 2

### Основні властивості функції регресії

**Л е м а .** Якщо  $\eta$  та  $\bar{\xi}$  – випадкові величина й вектор, відповідно, а  $M\eta^2 < \infty$ , тоді для них справедливо:

$$D\eta = M \left\{ \left( M(\eta/\bar{\xi}) - M\eta \right)^2 \right\} + M \left\{ M \left\{ \left( \eta - M(\eta/\bar{\xi}) \right)^2 / \bar{\xi} \right\} \right\},$$

або скорочено  $D\eta = Df(\bar{\xi}) + Mg(\bar{\xi})$ .

**Доведення.** Використовуючи властивості умовного математичного сподівання випадкових величин відносно випадкового вектора  $\bar{\xi}$ , дисперсію  $\eta$  можна розписати таким чином:

$$D\eta = M(\eta - M\eta)^2 = M \left\{ M \left\{ (\eta - M\eta)^2 / \bar{\xi} \right\} \right\} = \\ = M \left\{ M \left\{ \left[ (\eta - M(\eta/\bar{\xi})) + (M(\eta/\bar{\xi}) - M\eta) \right]^2 / \bar{\xi} \right\} \right\} = \\ = M \left\{ M \left\{ \left[ (\eta - M(\eta/\bar{\xi}))^2 + 2(\eta - M(\eta/\bar{\xi}))(M(\eta/\bar{\xi}) - M\eta) + \right. \right. \right. \\ \left. \left. \left. + (M(\eta/\bar{\xi}) - M\eta)^2 \right] / \bar{\xi} \right\} \right\} = M \left\{ M \left\{ (\eta - M(\eta/\bar{\xi}))^2 / \bar{\xi} \right\} \right\} +$$

$$+ 2M \left\{ \left( M(\eta/\bar{\xi}) - M\eta \right) M \left\{ (\eta - M(\eta/\bar{\xi})) / \bar{\xi} \right\} \right\} + \\ + M \left\{ M \left\{ \left( M(\eta/\bar{\xi}) - M\eta \right)^2 / \bar{\xi} \right\} \right\} = \\ = M \left\{ M \left\{ (\eta - M(\eta/\bar{\xi}))^2 / \bar{\xi} \right\} \right\} + M \left\{ M \left\{ \left( M(\eta/\bar{\xi}) - M\eta \right)^2 / \bar{\xi} \right\} \right\} = \\ = M \left\{ M \left\{ (\eta - M(\eta/\bar{\xi}))^2 / \bar{\xi} \right\} \right\} + M \left\{ \left( M(\eta/\bar{\xi}) - M\eta \right)^2 \right\} = \\ = Mg(\bar{\xi}) + Df(\bar{\xi}).$$

Що й треба було довести.

**Т е о р е м а .** Нехай  $M\eta^2 < \infty$ ,  $\Phi$  – множина борелівських функцій на  $R^q$ , тоді

$$f(\cdot) = \arg \min_{\varphi(\cdot) \in \Phi} M \left[ \eta - \varphi(\bar{\xi}) \right]^2,$$

де  $f(x) = M(\eta/\bar{\xi} = \bar{x})$ .

**Доведення.** Без втрати загальності будемо розглядати ті борелівські функції на  $R^q$ , для яких  $M\varphi^2(\bar{\xi}) < \infty$ . Скориставшись властивостями умовного математичного сподівання випадкової величини  $\eta$  відносно випадкового вектора  $\bar{\xi}$  і врахувавши, що  $M[f(\bar{\xi}) - \varphi(\bar{\xi})]^2 \geq 0$ , легко бачити, що має місце ланцюжок перетворень:

$$M[\eta - \varphi(\bar{\xi})]^2 = M \left[ (\eta - f(\bar{\xi})) + (f(\bar{\xi}) - \varphi(\bar{\xi})) \right]^2 = \\ = M[\eta - f(\bar{\xi})]^2 + 2M \left[ (\eta - f(\bar{\xi}))(f(\bar{\xi}) - \varphi(\bar{\xi})) \right] + \\ + M[f(\bar{\xi}) - \varphi(\bar{\xi})]^2 \geq M[\eta - f(\bar{\xi})]^2 + \\ + 2M \left\{ M \left[ (\eta - f(\bar{\xi}))(f(\bar{\xi}) - \varphi(\bar{\xi})) / \bar{\xi} \right] \right\} = M[\eta - f(\bar{\xi})]^2 + \\ + 2M \left\{ (f(\bar{\xi}) - \varphi(\bar{\xi})) M \left[ (\eta - f(\bar{\xi})) / \bar{\xi} \right] \right\} = M[\eta - f(\bar{\xi})]^2.$$

В останньому переході скористалися тим, що  $M[(\eta - f(\bar{\xi})) / \bar{\xi}] = 0$ .

У результаті, маємо



$$M[\eta - \varphi(\bar{\xi})]^2 \geq M[\eta - f(\bar{\xi})]^2,$$

тобто отримано нижню межу для нашого функціоналу  $M[\eta - \varphi(\bar{\xi})]^2$ , яка

досягається на функції регресії  $f(x) = M(\eta / \bar{\xi} = \bar{x})$ .

Доведення завершено.

## СПИСОК ЛІТЕРАТУРИ

1. Айвазян С.А., Енюков И.С., Мешалкин Л.Д. Прикладная статистика. Исследование зависимостей. — М.: Финансы и статистика, 1985.
2. Афифи А., Эйзен С. Статистический анализ. Подход с использованием ЭВМ. — М.: Мир, 1982.
3. Большев Л.Н., Смирнов Н.В. Таблицы математической статистики. — М.: Наука, 1983.
4. Кендалл М., Стьюарт А. Статистические выводы и связи. — М.: Наука, Гл. редакция физ.-мат. лит-ры, 1973.
5. Кендэл М. Ранговые корреляции. — М.: Статистика, 1975.
6. Крамер Г. Математические методы статистики. — М.: Мир, 1975.
7. Кульбак С. Теория информации и статистика. — М.: Наука, 1967.
8. Рао С.Р. Линейные статистические методы и их применение. — М.: Наука, 1968.
9. Ширяев А.Н. Вероятность. — М.: Наука, Гл. редакция физ.-мат. лит-ры, 1980.

## Предметний покажчик

### Аналіз

- кореляційний змінних
- — кількісних 4
- — класифікаційних 36
- — номінальних 36
- — ординальних 27
- — порядкових 27
- — якісних 27, 36
- рангових кореляцій 29

### Борелівська

- множина
- — в  $R$  41
- — в  $R^q$  42
- функція
- — на  $R$  41
- — на  $R^q$  42
- $\sigma$ -алгебра
- — в  $R$  41
- — в  $R^q$  42

### Вибіркове значення

- індексу кореляції 12, 19
- коефіцієнта кореляції
- — парного 15
- — множинного 26

### Випадкова величина

- вимірна відносно  $\sigma$ -алгебри 42
- незалежна від  $\sigma$ -алгебри 44

### Випадкові величини незалежні 44

### Відношення кореляційне

- парне 11
- множинне 18

### Властивості

- індексу кореляції 7
- коефіцієнта кореляції
- — парного 14
- — множинного 25
- умовного математичного сподівання відносно
- — випадкового вектора 46
- —  $\sigma$ -алгебри 44–45

**Детермінації коефіцієнт 8**

- для нормального випадку 13, 25
- — вибіркове значення 26
- — емпіричне значення 26

**Емпіричне значення**

- індексу кореляції 12, 19
- коефіцієнта кореляції
- — парного 15
- — множинного 26

**Ентропія 39**

**Зв'язаний ранг 29**

**Зв'язок статистичний**

- парний 9
- множинний 16

**Змінна**

- залежна 4
- незалежна 4

**Індекс кореляції 7**

- вибіркове значення 12, 19
- властивості 7
- для нормального випадку 13, 24
- емпіричне значення 12, 19
- перевірка на значимість 12, 19

**Індикатор множини 44**

**Інформаційна міра зв'язку 39**

- вибіркове значення 39–40
- властивості 39
- перевірка на значимість 40

**Квадратична спряженість 37**

**Коваріації коефіцієнт 13**

**Коефіцієнт**

- детермінації 8
- — для нормального випадку 13, 25
- — — вибіркове значення 26
- — — емпіричне значення 26
- коваріації 13
- конкордації 34
- — властивості 35
- — модифікований 35
- — перевірка на значимість 35–36
- кореляції 13
- — множинний 23
- — — властивості 25

— / — перевірка на значимість 26–27

— — парний 13

— — — властивості 14

— — — вибіркове значення 15

— — — емпіричне значення 15

— — — перевірка на значимість 15–16

— — ранговий

— — — Кендела 32

— — — — властивості 32–33

— — — — модифікований 33

— — — — перевірка на значимість 33–34

— — — Спірмена 30

— — — — властивості 30–31

— — — — модифікований 31

— — — — перевірка на значимість 31–32

— — частинний 20–21

— — — вибіркове значення 22

— — — емпіричне значення 22

— — — перевірка на значимість 23

— — — рекурентна процедура підрахунку 21

— спряженості

— — властивості 38

— — квадратний корінь із середньої квадратичної спряженості 38

— — Крамера 38

— — перевірка на значимість 38

— — Пірсона 38

— — Чупрова 38

**Конкордації**

— коефіцієнт 34

— — властивості 35

— — модифікований 35

— — перевірка на значимість 35–36

**Кореляції**

— індекс 7

— — вибіркове значення 12, 19

— — властивості 7

— — для нормального випадку 13, 24

— — емпіричне значення 12, 19

— — перевірка на значимість 12, 19

— коефіцієнт 13

— — множинний 23

— — — властивості 25

— — — перевірка на значимість 26–27

— — парний 13

- — властивості 14
- — вибіркове значення 15
- — емпіричне значення 15
- — перевірка на значимість 15–16
- — ранговий
- — Кендела 32
- — властивості 32–33
- — модифікований 33
- — перевірка на значимість 33–34
- — Спірмена 30
- — властивості 30–31
- — модифікований 31
- — перевірка на значимість 31–32
- — частинний 20–21
- — вибіркове значення 22
- — емпіричне значення 22
- — перевірка на значимість 23
- — рекурентна процедура підрахунку 21
- Кореляційне відношення
  - парне 11
  - множинне 18
- Кореляційний аналіз змінних
  - кількісних 4
  - класифікаційних 36
  - номінальних 36
  - порядкових 27
  - ординальних 27
  - якісних 27, 36
- Множина борелівська
  - в  $R$  41
  - в  $R^q$  42
- Множини індикатор 44
- Множинне кореляційне відношення 18
- Множинний
  - коефіцієнт
  - детермінації 8, 16
  - кореляції 23
  - властивості 25
  - перевірка на значимість 26–27
  - статистичний зв'язок 16
- Незалежні
  - випадкові величини 44
  - $\sigma$ -алгебри 44

- Парне кореляційне відношення 11
- Парний
  - коефіцієнт кореляції 13
  - статистичний зв'язок 9
- Перевірка на значимість 4
  - індексу кореляції 12, 19
  - коефіцієнта кореляції
  - парного 15–16
  - множинного 26–27
  - частинного 23
- Протилежна ранжировка 28
- Радона – Никодима теорема 43
- Ранг об'єкта 28
  - зв'язаний 29
  - з групи нерозрізних об'єктів 29
- Ранговий коефіцієнт кореляції
  - Кендела 32
  - властивості 32–33
  - модифікований 33
  - перевірка на значимість 33–34
  - Спірмена 30
  - властивості 30–31
  - модифікований 31
  - перевірка на значимість 31–32
- Ранжировка 28
  - протилежна 28
- Регресії функція 5
  - для нормального випадку 13, 24
- Рівень значущості 4
- Спряженості
  - коефіцієнт
  - властивості 38
  - квадратний корінь із середньої квадратичної спряженості 38
  - Крамера 38
  - перевірка на значимість 38
  - Пірсона 38
  - Чупрова 38
  - таблиця 37
- Спряженість квадратична 37
- Статистичний зв'язок
  - парний 9
  - множинний 16

Таблиця рангів 29

Таблиця спряженості 37

Теорема Радона – Никодима 43

Умовна дисперсія випадкової величини відносно

— випадкового вектора 46

— події  $\{\bar{\xi} = \bar{x}\}$  47

Умовна ймовірність події відносно

— випадкового вектора 46

—  $\sigma$ -алгебри 44

Умовне математичне сподівання випадкової величини відносно

— випадкового вектора 45

— властивості 46

— події  $\{\bar{\xi} = \bar{x}\}$  47

—  $\sigma$ -алгебри 43

— властивості 44–45

Функція

— борелівська

— на  $R$  41

— на  $R^q$  42

— регресії 5

— для нормального випадку 13, 24

Частинний коефіцієнт кореляції 20–21

— вибіркове значення 22

— емпіричне значення 22

— перевірка на значимість 23

— рекурентна процедура підрахунку 21

$\sigma$ -алгебра

— борелівська

— в  $R$  41

— в  $R^q$  42

— породжена

— випадковим вектором 45

— випадковою величиною 45

$\sigma$ -алгебри незалежні 44

## Показчик позначень

$\hat{C}_{\eta\bar{\xi}}$ 38	$M(\eta/\bar{\xi})$ 45–46
$d_R$ 41	$M(\eta/\bar{\xi} = \bar{x})$ 47
$d_{R^q}$ 41	$n_{ij}$ 36
$D(\eta/\bar{\xi})$ 46	$n_{i\bullet}$ 37
$D(\eta/\bar{\xi} = \bar{x})$ 47	$n_{\bullet j}$ 37
$f(\bar{x})$ 5	$n_i^{(k)}$ 31, 33, 35
$F_\alpha(m, n)$ 11	$N_{\eta\bar{\xi}}$ 39
$g(\bar{x})$ 5	$\hat{N}_{\eta\bar{\xi}}$ 40
$H_{\bar{\xi}}$ 39	$\hat{\hat{N}}_{\eta\bar{\xi}}$ 40
$\hat{H}_{\bar{\xi}}$ 40	$P(B/G)$ 44
$H_{\eta\bar{\xi}}$ 39	$P(B/\bar{\xi})$ 46
$\hat{H}_{\eta\bar{\xi}}$ 40	$\hat{P}_{\eta\bar{\xi}}$ 38
$I_B$ 44	$r_{ij}$ 20
$I_{\eta\bar{\xi}}$ 9	$r_{ij(I)}$ 20
$\hat{I}_{\eta\bar{\xi}}$ 12	$\hat{\hat{r}}_{ij(I)}$ 22
$I_{\eta\bar{\xi}}$ 7	$r_{\eta\bar{\xi}}$ 13
$\hat{I}_{\eta\bar{\xi}}$ 19	$\hat{r}_{\eta\bar{\xi}}$ 15
$m^{(k)}$ 31, 33, 35	$r_{\eta\bar{\xi}}$ 23
$M(\eta/G)$ 43	$\hat{r}_{\eta\bar{\xi}}^2$ 26

$t_{\alpha}(n)$	16	$\hat{\rho}_{\eta\xi}$	18
$T^{(k)}$	31, 35	$\sigma_R$	41
$\hat{T}_{\eta\xi}$	38	$\sigma_{R^q}$	41
$U^{(k)}$	33	$\sigma_{\xi}$	45
$\hat{W}_{\xi}$	34	$\sigma_{\bar{\xi}}$	45
$\hat{W}_{\xi}$	35	$\hat{\tau}_{ij}^{(K)}$	32
$x^{(i)}$	28	$\hat{\tau}_{ij}^{(K)}$	33
$x_k^{(i)}$	28	$\hat{\tau}_{ij}^{(S)}$	30
$\bar{x}, \bar{y}$	15	$\hat{\tau}_{ij}^{(S)}$	31
$\xi^+$	43	$\hat{\phi}_{\eta\xi}$	38
$\xi^-$	43	$\chi_{\alpha}^2(n)$	36
$\hat{\rho}_{\eta\xi}$	11	$\hat{\chi}_{\eta\xi}^2$	37

## З М І С Т

ВСТУП	3
1. АНАЛІЗ НАЯВНОСТІ СТАТИСТИЧНОГО ЗВ'ЯЗКУ МІЖ КІЛЬКІСНИМИ ЗМІННИМИ	4
1.1. Середньоквадратична апроксимація випадкової величини. Функція регресії та її властивості	4
1.2. Індекс кореляції та його властивості. Коефіцієнт детермінації	6
1.3. Аналіз парних статистичних зв'язків кількісних змінних	9
1.3.1. Вибіркові значення характеристики парного статистичного зв'язку в загальному випадку	10
1.3.1.1. Кореляційне відношення	10
1.3.1.2. Оцінка індексу кореляції	12
1.3.2. Коефіцієнт кореляції – характеристика парного статистичного зв'язку у нормальному випадку	13
1.4. Аналіз множинних статистичних зв'язків кількісних змінних	16
1.4.1. Емпіричні значення характеристики множинного статистичного зв'язку в загальному випадку	17
1.4.1.1. Обробка згрупованих даних	17
1.4.1.2. Використання наближення функції регресії	18
1.4.2. Обробка спостережень у нормальному випадку	19
1.4.2.1. Частинний коефіцієнт кореляції та його властивості	20
1.4.2.2. Множинний коефіцієнт кореляції та його властивості	23
2. КОРЕЛЯЦІЙНИЙ АНАЛІЗ ОРДИНАЛЬНИХ ЗМІННИХ	27
2.1. Ранги та таблиці рангів	27
2.2. Аналіз парних рангових кореляцій	30
2.2.1. Ранговий коефіцієнт кореляції Спірмена	30
2.2.2. Ранговий коефіцієнт кореляції Кендела	32
2.3. Коефіцієнт конкордації – характеристика множинних статистичних зв'язків для ординальних змінних	34
3. ДОСЛІДЖЕННЯ НАЯВНОСТІ ЗВ'ЯЗКУ МІЖ НОМІНАЛЬНИМИ ЗМІННИМИ	36
3.1. Таблиця спряженості	36
3.2. Квадратична спряженість і характеристики парного зв'язку на її основі	37
3.3. Інформаційна міра парного статистичного зв'язку	39
ДОДАТОК 1. Умовні ймовірності та математичні сподівання. Основні властивості	41
ДОДАТОК 2. Основні властивості функції регресії	48
СПИСОК ЛІТЕРАТУРИ	50
ПРЕДМЕТНИЙ ПОКАЖЧИК	51
ПОКАЖЧИК ПОЗНАЧЕНЬ	57