

I Кореляційний аналіз

Кореляційний аналіз займається з'ясуванням питання є істотним (суттєвим, достатньо тісним) зв'язок між змінними, які досліджуються.

Основні етапи розв'язання задачі кореляційного аналізу:

- вибір числової характеристики статистичного зв'язку;
- визначення оцінки цієї характеристики статистичного зв'язку;
- на основі отриманого значення оцінки характеристики статистичного зв'язку прийняття рішення, чи є істотним зв'язок між змінними, які аналізуються.

І тільки після стверджувальної відповіді про істотність статистичного зв'язку між змінними, що розглядаються, має сенс переходити до наступного етапу – пошуку математичної моделі цього зв'язку засобами інших розділів аналізу та обробки інформації (даних).

Вимоги до характеристик статистичного зв'язку $K_{\eta\bar{\xi}}$ при їх конструюванні:

- якщо $K_{\eta\bar{\xi}} = 0$, то це відповідає ситуації, що зв'язок між η та $\bar{\xi}$ відсутній,
- зі збільшенням відхилення $K_{\eta\bar{\xi}}$ від нуля зростає суттєвість зв'язку між η та $\bar{\xi}$,

-
- нехай $\max |K_{\eta\bar{\xi}}| = K_{\max}$, ($K_{\max} < \infty$); тоді якщо $|K_{\eta\bar{\xi}}| = K_{\max}$, то це відповідає ситуації, що зв'язок між η та $\bar{\xi}$ функціональний.

Нехай $\hat{K}_{\eta\bar{\xi}}$ є оцінкою характеристики $K_{\eta\bar{\xi}}$, тоді у загальному випадку спільна процедура її використання буде мати вигляд:

- якщо $\hat{K}_{\eta\bar{\xi}} = 0$, то зв'язок між η та $\bar{\xi}$ відсутній,
- якщо $|\hat{K}_{\eta\bar{\xi}}| = K_{\max}$, ($K_{\max} < \infty$), то зв'язок між η та $\bar{\xi}$ функціональний,
- якщо $|\hat{K}_{\eta\bar{\xi}}| \in (0, K_{\max})$, то потрібно перевірити гіпотезу про те, чи значимо відхиляється від нуля коефіцієнт $K_{\eta\bar{\xi}}$, тобто

Нехай аналізується зв'язок між *залежною* скалярною змінною η та вектором *незалежних* змінних $\bar{\xi}$ розмірності q деякого типу. Припустимо, що була обрана характеристика (коефіцієнт) статистичного зв'язку $K_{\eta\bar{\xi}}$.

Якщо $q=1$, то її називають *парною характеристикою статистичного зв'язку*, інакше – *множинною характеристикою статистичного зв'язку*.

- якщо $\left| \hat{K}_{\eta\bar{\xi}} \right| \in (0, K_{\max})$, то потрібно перевірити гіпотезу про те, чи значимо відхиляється від нуля коефіцієнт $K_{\eta\bar{\xi}}$, тобто здійснити *перевірку $K_{\eta\bar{\xi}}$ на значимість*, а саме перевірити гіпотезу:

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

з деяким рівнем значущості $\alpha > 0$.

Приклади задач кореляційного аналізу:

1. чи суттєво впливає рівень безробіття на рівень злочинності у країні (+1% росту рівня безробіття дає +5% росту рівня злочинності),
2. вплив рівня алкоголю у крові водія на час його реакції,
3. вплив рівня допінгу у крові спортсмена на його спортивне досягнення,
4. вплив ряду основних економічних показників країни на рівень життя населення,
5. чи є суттєвим вплив на об'єми продажів товару/послуги має обсяг фінансування рекламної компанії товару/послуги,
6. чи істотно впливають на врожайність певної сільськогосподарської культури кількість опадів та сонячних годин протягом сезону,
7. чи суттєво залежить остаточний результат студента під час сесії від результатів його роботи протягом семестру.

Структура кореляційного аналізу:

- кореляційний аналіз кількісних змінних,
- кореляційний аналіз ординальних змінних,
- кореляційний аналіз номінальних змінних.

Кореляційний аналіз кількісних змінних

Ключове поняття у цьому розділі це функція регресії η щодо $\bar{\xi}$.

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

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

Тут $M(\eta / \bar{\xi} = \bar{x})$ - *умовне математичне сподівання випадкової величини η відносно події $\{\bar{\xi} = \bar{x}\}$* . Детальніше, відомо що:

$$M(\eta) = \int_{-\infty}^{\infty} y dF_{\eta}(y), \text{ де } F_{\eta}(y) = P\{\eta < y\}, \text{ тоді}$$

$$M(\eta / \bar{\xi} = \bar{x}) = \int_{-\infty}^{\infty} y dF_{\eta}(y / \bar{\xi} = \bar{x}), \text{ де } F_{\eta}(y / \bar{\xi} = \bar{x}) = P\{\eta < y / \bar{\xi} = \bar{x}\}.$$

Якщо $M\eta^2 < \infty$, тоді по аналогії буде існувати *умовна дисперсія випадкової величини η відносно події $\{\bar{\xi} = \bar{x}\}$* :

$$g(\bar{x}) = D(\eta / \bar{\xi} = \bar{x}).$$

Для $f(\bar{\xi})$ та $g(\bar{\xi})$ будемо використовувати такі нотації:

$$\begin{aligned} M(\eta / \bar{\xi}) &= f(\bar{\xi}), \\ D(\eta / \bar{\xi}) &= g(\bar{\xi}), \end{aligned}$$

та називати відповідно *умовним математичним сподіванням випадкової величини η відносно випадкового вектора $\bar{\xi}$* та *умовною дисперсією випадкової величини η відносно випадкового вектора $\bar{\xi}$* .

Самостійна робота №4. З навчального посібника
 «Слабоспицький О.С. Основи кореляційного аналізу даних, 2006».
 Пропрацювати матеріал наведений у Додатку 1:
 Умовні ймовірності та математичні сподівання. Основні властивості.

Згадаємо деякі властивості для них:

1. $M\{M(\eta / \bar{\xi})\} = M\eta,$
2. $M(\varphi(\bar{\xi})\eta / \bar{\xi}) = \varphi(\bar{\xi})M(\eta / \bar{\xi}), \quad \varphi(\cdot) \in \mathfrak{B}_q,$
 де \mathfrak{B}_q – множина борелівських функцій на \mathbb{R}^q .

Наслідок.

1. $Mf(\bar{\xi}) = M\eta, \quad (*)$
2. $M\left[(\eta - f(\bar{\xi})) / \bar{\xi}\right] = 0. \quad (**)$
3. $M(\varphi(\bar{\xi})[\eta - f(\bar{\xi})]) = 0, \quad \forall \varphi(\cdot) \in \mathfrak{B}_q. \quad (***)$

Доведення наслідку. 1. Очевидно з огляду на першу властивість.

2. Легко перевіряється, бо згідно другої властивості:

$$M\left[(\eta - f(\bar{\xi})) / \bar{\xi}\right] = M(\eta / \bar{\xi}) - f(\bar{\xi}) = 0.$$

3. Скористаємося послідовно спочатку властивістю 1 справа наліво, а потім наслідком 2, в результаті отримаємо:

$$\begin{aligned} M(\varphi(\bar{\xi})[\eta - f(\bar{\xi})]) &= M\{M(\varphi(\bar{\xi})[\eta - f(\bar{\xi})] / \bar{\xi})\} = \\ &= M\{\varphi(\bar{\xi})M([\eta - f(\bar{\xi})] / \bar{\xi})\} = 0. \quad \blacksquare \end{aligned}$$

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

$$D\eta = {}^I Df(\bar{\xi}) + Mg(\bar{\xi}),$$

або розгорнуто

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

Доведення. Використаємо вищенаведені властивості:

$$\begin{aligned}
 D\eta &= M(\eta - M\eta)^2 = M\left[\left(\eta - f(\bar{\xi})\right) + \left(f(\bar{\xi}) - M\eta\right)\right]^2 = \\
 &= M\left(\eta - f(\bar{\xi})\right)^2 + 2M\left[\left(\eta - f(\bar{\xi})\right)\left(f(\bar{\xi}) - M\eta\right)\right] + M\left(f(\bar{\xi}) - M\eta\right)^2 \stackrel{(***)}{=} \\
 &= M\left\{M\left[\left(\eta - f(\bar{\xi})\right)^2 / \bar{\xi}\right]\right\} + Df(\bar{\xi}) = \\
 &= M\left\{M\left[\left(\eta - M(\eta / \bar{\xi})\right)^2 / \bar{\xi}\right]\right\} + Df(\bar{\xi}) = M\left\{D(\eta / \bar{\xi})\right\} + Df(\bar{\xi}) = \\
 &= Mg(\bar{\xi}) + Df(\bar{\xi}). \quad \blacksquare
 \end{aligned}$$

Теорема (про фундаментальну властивість функції регресії). Нехай η і $\bar{\xi}$ – випадкові величина та вектор розмірності q , відповідно, причому $M\eta^2 < \infty$, \mathfrak{B}_q – множина борелівських функцій на \mathbb{R}^q , тоді:

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

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

$$f(\cdot) = \arg \min_{\varphi(\cdot) \in \mathfrak{B}_q} M\left(\eta - \varphi(\bar{\xi})\right)^2$$

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

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

$$\begin{aligned}
M[\eta - \varphi(\bar{\xi})]^2 &= M\left[\left(\eta - f(\bar{\xi})\right) + \left(f(\bar{\xi}) - \varphi(\bar{\xi})\right)\right]^2 = \\
&= M\left[\eta - f(\bar{\xi})\right]^2 + 2M\left[\left(\eta - f(\bar{\xi})\right)\left(f(\bar{\xi}) - \varphi(\bar{\xi})\right)\right] + \\
&\quad + M\left[f(\bar{\xi}) - \varphi(\bar{\xi})\right]^2 \stackrel{(\varphi \neq f)}{\geq} M\left[\eta - f(\bar{\xi})\right]^2.
\end{aligned}$$

У результаті, маємо

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

тобто отримано нижню межу для нашого функціоналу $M\left[\eta - \varphi(\bar{\xi})\right]^2$, яка досягається на функції регресії $f(\bar{x}) = M(\eta / \bar{\xi} = \bar{x})$. ■

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

модель можна взяти таке співвідношення, яке будемо називати *регресійною моделлю* η щодо $\bar{\xi}$:

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

де ε – залишкова похибка апроксимації.

Для цієї моделі мають місце такі властивості:

- | |
|---|
| <ol style="list-style-type: none"> 1) $Mf(\bar{\xi}) = M\eta$, $M\varepsilon = 0$; 2) $f(\bar{\xi})$ та ε – 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 η відносно

випадкового вектора $\vec{\xi}$ та доведемо послідовно кожну з властивостей:

$$1) \quad Mf(\vec{\xi})^{(*)} = M\eta.$$

А це у свою чергу, дозволяє стверджувати, що

$$M\varepsilon = M\{\eta - f(\vec{\xi})\} = M\eta - Mf(\vec{\xi}) = 0;$$

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

$$M[f(\vec{\xi}) - Mf(\vec{\xi})]\varepsilon = M[f(\vec{\xi}) - M\eta][\eta - f(\vec{\xi})]^{(***)} = 0.$$

3) з доведеної у другій властивості некорельованості $f(\vec{\xi})$ та ε випливає: $D\eta = D(f(\vec{\xi}) + \varepsilon) = Df(\vec{\xi}) + D\varepsilon$. ■

Коефіцієнт детермінації та його властивості. Індекс кореляції.

Введемо універсальну характеристику статистичного зв'язку для змінних η та $\vec{\xi}$, ($\vec{\xi} \in \mathbb{R}^q$). Будемо вважати, що $M|\eta| < \infty$, тоді існуватиме функція регресії η щодо $\vec{\xi}$, а саме: $f(\vec{x}) = M(\eta / \vec{\xi} = \vec{x})$. А для випадкової величини η можна використовувати таку математичну модель:

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

де ε – залишкова похибка апроксимації.

Причому $D\eta = Df(\vec{\xi}) + D\varepsilon$.

Після цих міркувань, враховуючи унікальну властивість функції регресії, доведену в теоремі, цілком природним здається використання, як характеристики статистичного зв'язку для кількісних змінних η та $\vec{\xi}$, нижчевизначеного коефіцієнта.

Означення. Нехай η і $\vec{\xi}$ – випадкові величина та вектор розмірності q , відповідно, причому $0 < D\eta < \infty$. Тоді **коефіцієнтом детермінації η щодо $\vec{\xi}$** називається величина

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

Зауваження 1. Коефіцієнт детермінації $I_{\eta\vec{\xi}}^2$ приваблює ще тим, що має прозору інтерпретацію, а саме: він вказує, яка частина дисперсії змінної η визначається варіацією (дисперсією) функції регресії $f(\vec{\xi})$.

Зауваження 2. Коефіцієнт детермінації $I_{\eta\vec{\xi}}^2$, де $\vec{\xi} \in \mathbb{R}^q$, у випадку $q = 1$ ще називають **парним коефіцієнтом детермінації η щодо $\vec{\xi}$** , а коли $q > 1$ його ще називають **множинним коефіцієнтом детермінації η щодо $\vec{\xi}$** .

Властивості коефіцієнта детермінації:

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

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

Доведення. Скористаємось відомими властивостями умовного математичного сподівання та умовної дисперсії:

1) так як $D\varepsilon \geq 0$, то отримуємо:

$$0 \leq I_{\eta\bar{\xi}}^2 = 1 - \frac{D\varepsilon}{D\eta} \leq 1;$$

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

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

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

3) нехай $I_{\eta\bar{\xi}}^2 = 1$, тоді

$$1 = I_{\eta\bar{\xi}}^2 = 1 - \frac{D\varepsilon}{D\eta} \Rightarrow D\varepsilon = 0.$$

Згідно першої властивості регресійної моделі $M\varepsilon = 0$, тоді

$$D\varepsilon = 0 \Rightarrow \varepsilon = M\varepsilon = 0.$$

Але $\varepsilon = \eta - f(\bar{\xi})$, а це, у свою чергу, дозволяє стверджувати, що з ймовірністю 1 має місце така функціональна залежність між η та $\bar{\xi}$:

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

