§ 5 **Метод найменших квадратів.**

Передбачається, що залежність ознаки $Y$ від ознаки $X$ має такий вигляд

$$y = f(x, a_1, a_2, \ldots, a_m)$$

де — значення ознаки $X$; $y$ — значення ознаки $Y$; $a_1, a_2, \ldots, a_m$ — параметри, які належить визначити, та, що за результатами експерименту отримані такі емпіричні дані:

(*)

| Значення ознаки $X$ | $x_1$ | $x_2$ | $\ldots$ | $x_i$ | $\ldots$ | $x_n$ |
|---|---|---|---|---|---|---|
| Значення ознаки $Y$ | $y_1$ | $y_2$ | $\ldots$ | $y_i$ | $\ldots$ | $y_n$ |

Метод найменших квадратів стверджує, що найімовірніше значення параметрів $a_1, a_2, \ldots, a_m$ дає мінімум функції

$$S = \sum_{i=1}^{n} \left[ y_i - f(x_i, a_1, a_2, \ldots, a_m) \right]^2.$$

Коли $f(x, a_1, a_2, \ldots, a_m)$ має неперервні частинні похідні за усіма своїми параметрами, то необхідна умова мінімуму функції $S$ складає систему $m$ рівнянь з $m$ невідомими:

(*,*) $\quad \dfrac{\partial S}{\partial a_1} = 0; \quad \dfrac{\partial S}{\partial a_2} = 0; \quad \ldots; \quad \dfrac{\partial S}{\partial a_m} = 0.$

Визначення функціональної залежності між ознаками $Y$ та $X$ на основі експериментальних даних (*) називають <u>вирівнюванням емпіричних даних вздовж кривої</u> $y = f(x, a_1, a_2, \ldots, a_m)$.

Якщо $f(x, \alpha_1, \alpha_2, ..., \alpha_m) = \alpha_1 x + \alpha_2$ то зазначеного кривою буде пряма лінія $y = \alpha_1 x + \alpha_2$. У цьому випадку система $(*, *)$ може бути перетворена у так звану нормальну систему методу найменших квадратів за умови вирівнювання за прямою:

$$\begin{cases} \alpha_1 \sum_{i=1}^{n} x_i^2 + \alpha_2 \sum_{i=1}^{n} x_i = \sum_{i=1}^{n} x_i y_i; \\ \alpha_1 \sum_{i=1}^{n} x_i + \alpha_2 n = \sum_{i=1}^{n} y_i. \end{cases}$$

Система $(*, *)$ коли вирівнювання здійснюється за параболою $y = \alpha_1 x^2 + \alpha_2 x + \alpha_3$ може бути перетворена до такого вигляду:

$$\begin{cases} \alpha_1 \sum_{i=1}^{n} x_i^4 + \alpha_2 \sum_{i=1}^{n} x_i^3 + \alpha_3 \sum_{i=1}^{n} x_i^2 = \sum_{i=1}^{n} x_i^2 y_i; \\ \alpha_1 \sum_{i=1}^{n} x_i^3 + \alpha_2 \sum_{i=1}^{n} x_i^2 + \alpha_3 \sum_{i=1}^{n} x_i = \sum_{i=1}^{n} x_i y_i; \\ \alpha_1 \sum_{i=1}^{n} x_i^2 + \alpha_2 \sum_{i=1}^{n} x_i + \alpha_3 n = \sum_{i=1}^{n} y_i. \end{cases}$$

## Лінійна кореляційна залежність

Стверджується, що дві ознаки $X$ та $Y$ знаходяться в кореляційній залежності, якщо кожному значенню одного з них відповідає певний розподіл іншого. Кореляційна залежність між ознаками $X$ та $Y$ задається за допомогою кореляційної таблиці.

| X \ Y | $y_1$ | $y_2$ | ... | $y_j$ | ... | $y_m$ | $m_x$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $m_{11}$ | $m_{12}$ | ... | $m_{1j}$ | ... | $m_{1n}$ | $m_{x1}$ |
| $x_2$ | $m_{21}$ | $m_{22}$ | ... | $m_{2j}$ | ... | $m_{2n}$ | $m_{x2}$ |
| ... | | | | | | | |
| $x_i$ | $m_{i1}$ | $m_{i2}$ | ... | $m_{ij}$ | ... | $m_{in}$ | $m_{xi}$ |
| ... | | | | | | | |
| $x_k$ | $m_{k1}$ | $m_{k2}$ | ... | $m_{kj}$ | ... | $m_{kn}$ | $m_{xk}$ |
| $m_y$ | $m_{y1}$ | $m_{y2}$ | ... | $m_{yj}$ | ... | $m_{yn}$ | $N$ |

У цій таблиці $x_1, x_2, ..., x_i, ..., x_k$; $y_1, y_2, ..., y_j, ..., y_n$ — середини інтервалів або значення ознак X та Y, а $m_{x_1}, m_{x_2}, ..., m_{x_i}, ..., m_{x_k}$; $m_{y_1}, m_{y_2}, ..., m_{y_j}, ..., m_{y_n}$ — відповідні частоти; $m_{ij}$ — частота, з якого з'являється пара $(x_i; y_j)$; $N = \sum\limits_{i=1}^{k} \sum\limits_{j=1}^{n} m_{ij}$.

Кореляційна залежність між ознаками Y та X може бути замінена функціональною, якщо кожному значенню ознаки X поставити у відповідність умовне середнє ознаки Y, тобто для $X = x_i$ поставити у відповідність величину

$$\overline{Y}_{x_i} = \frac{\sum\limits_{j=1}^{n} m_{ij} y_j}{m_{x_i}},$$

Якщо потім точки $(x_i, \overline{Y}_{x_i})$ вирівняти за методом найменших квадратів вздовж кривої $y = f(x, \alpha_1, \alpha_2, ..., \alpha_m)$, то останнє

називається *лінією регресії* у на x, а її рівняння - рівнянням регресії у на x. Аналогічно визначається лінія регресії x на y. Найпростішими і найважливішими випадками кривих регресій є прямі лінії.

Кутовий коефіцієнт прямої регресії у на x (x на y) називають *коефіцієнтом регресії* у на x (x на y) і позначають так $\rho_{y/x}$ ($\rho_{x/y}$).

Коефіцієнти регресії можуть бути розраховані за формулами:

$$\rho_{y/x} = \frac{\overline{XY} - \overline{X}\,\overline{Y}}{\sigma_x^2}; \quad \rho_{x/y} = \frac{\overline{XY} - \overline{X}\,\overline{Y}}{\sigma_y^2},$$

де $\overline{XY} = \dfrac{\sum\limits_{i=1}^{k}\sum\limits_{j=1}^{n} m_{ij}x_i y_j}{N}$ — середнє значення добутку ознак $X$ та $Y$; $\overline{X}$ і $\overline{Y}$ — їх середні значення, а $\sigma_x^2$ та $\sigma_y^2$ — їх дисперсії.

Рівняння прямих регресій мають ~~вигляд~~ вигляд:

$$y - \overline{Y} = \rho_{y/x}(x - \overline{X});$$
$$x - \overline{X} = \rho_{x/y}(y - \overline{Y}).$$

Коефіцієнтом лінійної кореляції ознак $X$ та $Y$ називається величина

$$r = r(X, Y) = \frac{\overline{XY} - \overline{X}\,\overline{Y}}{\sigma_x \sigma_y} = \pm\sqrt{\rho_{y/x}\,\rho_{x/y}}.$$

Коефіцієнт лінійної кореляції г

характеризується такими властивостями:

1) $-1 \leq r \leq 1$;

2) $r[\alpha(X-x_0), \beta(Y-y_0)] = r(X,Y)$ $(\alpha > 0, \beta > 0)$;

3) якщо $r(X,Y) = \pm 1$, то між ознаками X і Y існує лінійна функціональна залежність (коли $r = 1$ — пряма залежність, а коли $r = -1$ — зворотна залежність);

4) якщо $r(X,Y) = 0$, то між ознаками X та Y відсутня лінійна кореляційна залежність;

5) $\rho_{y/x} = r \dfrac{\sigma_y}{\sigma_x}$ і $\rho_{x/y} = r \dfrac{\sigma_x}{\sigma_y}$.

Квадрат коефіцієнта лінійної кореляції дає коефіцієнт детермінації, який вимірює долю варіації Y, яка пояснюється впливом ознаки X, і навпаки.

На практиці про розподіл ознак X та Y в генеральній сукупності визначають за даними вибірки. За цими даними можливо знайти вибірковий коефіцієнт лінійної кореляції $r_B$, який є випадковою величиною. За умови достатньо великого об'єму вибірки $r_B \approx r$. Якщо розподіл ознак X та Y досить близький до нормального, то можливо наближено вважати $r_B$ також нормальною випадковою величиною, середнє квадратичне відхилення якої дорівнює

$$\dfrac{1-r^2}{\sqrt{N}},$$ де N — об'єм вибірки.

Приклад розв'язання задачі.

Задача 1.

Компанія-перевізник провела статистичне дослідження на основних маршрутах і одержала залежність між вартістю перевезення $\xi$ (в умовних одиницях за 1 км) і довжиною маршруту $\eta$ (тис. км). Результати дослідження наведені в таблиці:
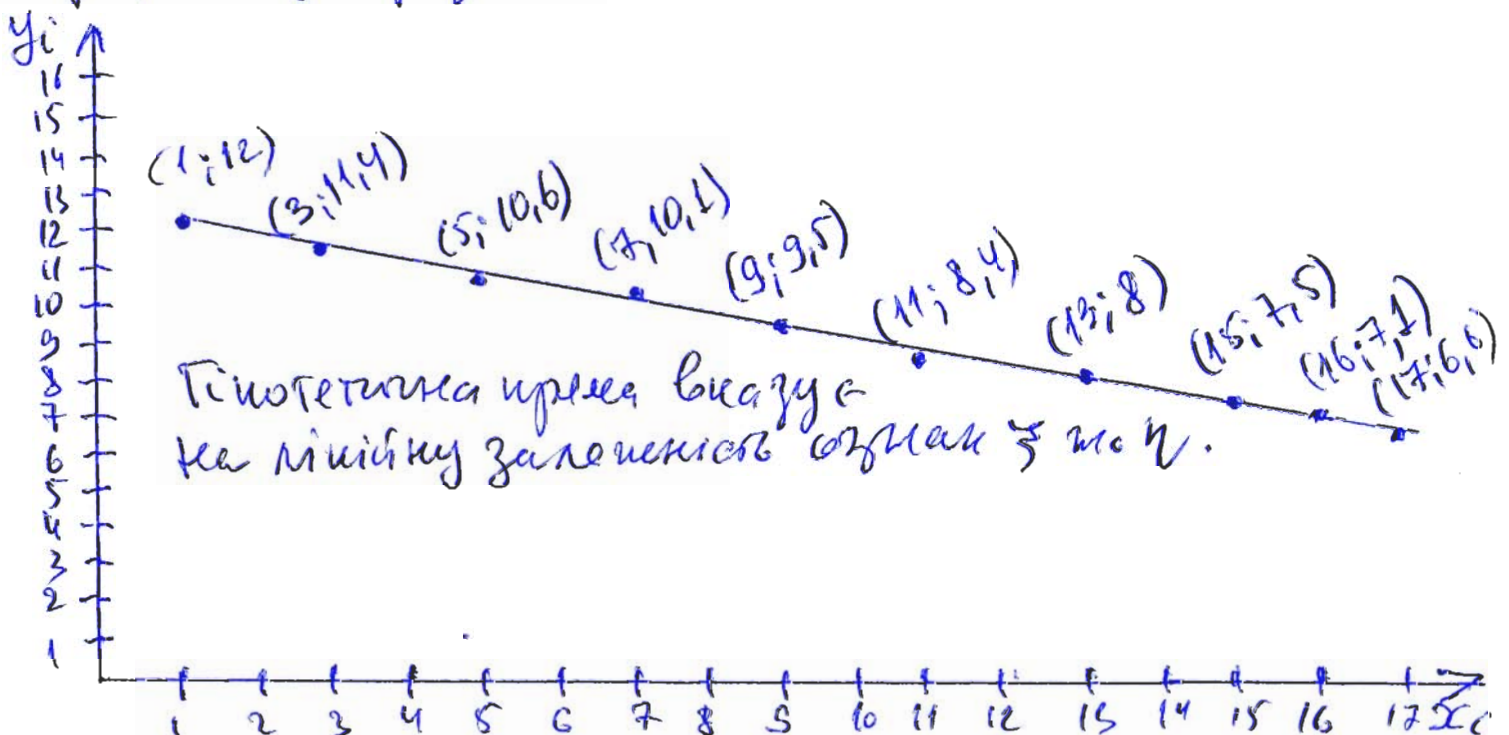
| Значення вартості $x_i$ | 1 | 3 | 5 | 7 | 9 | 11 | 13 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|
| Довжина маршруту $y_i$ | 12 | 11,4 | 10,6 | 10,1 | 9,5 | 8,4 | 8 | 7,5 | 7,1 | 6,6 |

В задачі потрібно

1) встановити форму залежності між $\xi$ і $\eta$;

2) знайти рівняння лінійної регресії $\eta$ на $\xi$ та $\xi$ на $\eta$;

3) обчислити коефіцієнт кореляції вибірки $r$ та оцінити силу лінійного зв'язку між $\xi$ та $\eta$.

Розв'язок.

1) Для обґрунтування залежності ознак $\xi$ та $\eta$ графічно зобразимо точки $(x_i, y_i)$.

Гіпотетична пряма вказує на лінійну залежність ознак $\xi$ та $\eta$.

2) Зв'язок між ознаками $\xi$ та $\eta$, який вивчається може бути вираженний рівнянням прямої лінії регресії $\eta$ на $\xi$: $\overline{y}_x = ax + b$.

Для обчислення параметрів $a$ i $b$ та коефіцієнта кореляції складемо розрахункові таблиці:

(*)

| $x_i$ | 1 | 3 | 5 | 7 | 9 | 11 | 13 | 15 | 16 | 17 | сума $\sum$ результ 97 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $y_i$ | 12 | 11,4 | 10,6 | 10,1 | 9,5 | 8,4 | 8 | 7,5 | 7,1 | 6,6 | 91,2 |
| $x_i^2$ | 1 | 9 | 25 | 49 | 81 | 121 | 169 | 225 | 256 | 289 | 1225 |
| $y_i^2$ | 144 | 129,96 | 112,36 | 102,01 | 90,25 | 70,56 | 64 | 37,5 | 50,41 | 43,56 | 844,3 |
| $x_i y_i$ | 12 | 34,2 | 53 | 70,7 | 85,5 | 92,4 | 104 | 112,5 | 113,6 | 105,6 | 783,5 |

Знайдемо рівняння лінійної регресії $\eta$ на $\xi$.

$$\overline{y}_x = ax + b$$

Застосуємо метод найменших квадратів і складемо систему рівнень для визначення параметрів $a$ і $b$.

$$\begin{cases} a\sum_{i=1}^{n} x_i^2 + b\sum_{i=1}^{n} x_i = \sum_{i=1}^{n} x_i y_i , \\ a\sum_{i=1}^{n} x_i + bn = \sum_{i=1}^{n} y_i \end{cases}$$

$n = 10$

Підставляємо значення із таблиці (*) і отримуємо систему рівнень:

$$\begin{cases} 1225 \cdot a + 97 \cdot b = 783,5 \\ 97 \cdot a + 10 \cdot b = 91,2 \end{cases}$$

$a = \dfrac{783,5 - 97b}{1225}$, $\quad 97\dfrac{783,5 - 97b}{1225} + 10b = 91,2 \Rightarrow$

$\Rightarrow 97(0,64 - 0,08b) + 10b = 91,2 \Rightarrow b = \dfrac{29,12}{2,24} = 13;$

$a = \dfrac{783,5 - 97 \cdot 13}{1225} = \dfrac{783,5 - 1261}{1225} = \dfrac{-477,5}{1225} = -0,4;$

$\underline{\overline{y}_x = -0,4x + 13.}$

Знайдемо рівняння лінійної регресії ξ на η ⑧
$$\overline{x}_y = ay + b.$$

Для визначення параметрів $a$ та $b$ за методом найменших квадратів запишемо відповідну систему рівнянь.

$$\begin{cases} a\sum\limits_{i=1}^{u} y_i^2 + b\sum\limits_{i=1}^{u} y_i = \sum\limits_{i=1}^{u} x_i y_i \\ a\sum\limits_{i=1}^{u} y_i + bu = \sum\limits_{i=1}^{u} x_i \end{cases}$$

Із таблиці (*) підставляємо розраховані дані. Отримуємо таку систему рівнянь:

$$\begin{cases} 844{,}3 \cdot a + 91{,}2\, b = 783{,}5 \\ 91{,}2\, a + 10\, b = 97 \end{cases}$$

Розв'яжемо систему рівнянь.

$$a = \frac{783{,}5 - 91{,}2\, b}{844{,}3} \Rightarrow 91{,}2(0{,}93 - 0{,}1\,b) + 10b = 97;$$

$$\Rightarrow 84{,}8 - 9{,}12\, b + 10b = 97; \Rightarrow 0{,}88\, b = 12{,}2 \Rightarrow b = \frac{12{,}2}{0{,}88} =$$

$$= 13{,}9 \approx 14$$

$$a = \frac{783{,}5 - 91{,}2 \cdot 13{,}9}{844{,}3} = \frac{783{,}5 - 1267{,}68}{844{,}3} = \frac{-484{,}18}{844{,}3} =$$

$$= -0{,}57; \Rightarrow a = -0{,}57, \; b = 14 \quad \underline{\overline{x}_y = -0{,}57x + 14}.$$

3) Обчислимо коефіцієнт кореляції $r$ вибірки. Коефіцієнт кореляції визначається за формулою

$$r(\xi, \eta) = \frac{\overline{\xi\eta} - \overline{\xi} \cdot \overline{\eta}}{\sigma_\xi \cdot \sigma_\eta},$$

Необхідно знайти вибіркові середнє добутку ознак $\xi$ та $\eta$, $\overline{\xi\eta} = \frac{1}{n}\sum\limits_{i=1}^{u} x_i y_i$, середні вибіркові

$$\bar{\xi} = \frac{1}{u}\sum_{i=1}^{u} x_i; \quad \bar{\eta} = \frac{1}{u}\sum_{i=1}^{u} y_i;$$

та вибіркові середньо квадратичні відхилення

$$\sigma_{\xi} = \sqrt{\frac{1}{n-1}\sum_{i=1}^{u}(x_i - \bar{x}_B)^2}; \quad \sigma_{\eta} = \sqrt{\frac{1}{n-1}\sum_{i=1}^{u}(y_i - \bar{y}_B)^2}$$

Застосуємо обчислені значення із таблиці (☆) і знайдемо вибіркові середні:

$$\overline{\xi\eta} = \frac{1}{u}\sum_{i=1}^{u} x_i y_i = \frac{1}{10}\cdot 783,5 = 78,35;$$

$$\bar{\xi} = \frac{1}{u}\sum_{i=1}^{u} x_i = \frac{1}{10}\cdot 97 = 9,7; \quad \bar{\eta} = \sum_{i=1}^{u} y_i = \frac{91,2}{10} = 9,12;$$

Для визначення $\sigma_{\xi}$ і $\sigma_{\eta}$ складемо таблицю для $(x_i - \bar{x}_B)^2$ та $(y_i - \bar{y}_B)^2$

(☆,☆)

| $(x_i - \bar{x}_B)^2$ | 75,69 | 44,89 | 22,09 | 7,29 | 0,49 | 1,69 | 10,89 | 28,09 | 39,69 | 53,29 |
| $(y_i - \bar{y}_B)^2$ | 8,3 | 5,2 | 2,2 | 0,77 | 0,14 | 0,52 | 1,25 | 2,62 | 4,08 | 6,35 |

Застосовуючи значення із таблиці (☆,☆) знайдемо $\sigma_{\xi}$ та $\sigma_{\eta}$

$$\sigma_{\xi} = \left[\frac{1}{9}(75,69 + 44,89 + 22,09 + 7,29 + 0,49 + 1,69 + 10,89 + 28,09 + 39,69 + 53,29)\right]^{1/2} =$$

$$= \sqrt{\frac{284,1}{9}} = 5,61; \quad \sigma_{\xi} = 5,61$$

$$\sigma_{\eta} = \left[\frac{1}{9}(8,3 + 5,2 + 2,2 + 0,77 + 0,14 + 0,52 + 1,25 + 2,62 + 4,08 + 6,35)\right]^{1/2} = \sqrt{\frac{31,43}{9}} = 1,87$$

$$\sigma_{\eta} = 1,87.$$

$$r(\xi, \eta) = \frac{78,35 - 9,7 \cdot 9,12}{5,61 \cdot 1,87} = \frac{78,35 - 88,46}{10,49} =$$

$$= \frac{-10,11}{10,49} = -0,96$$

Ознаки $\xi$ та $\eta$ від'ємно корельовані. Це визначається тим, що коли ознака $\xi$ зростає ознака $\eta$ спадає. Значення $r$ вказує на те, що кореляція ознак $\xi$ та $\eta$ висока.