

## 19. Свойства коэффициентов уравнения парной линейной регрессии.

Модель парной линейной регрессии имеет вид:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

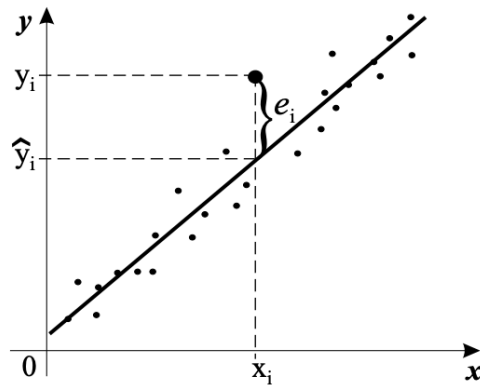
где:

- $y$  — зависимая переменная,
- $x$  — независимая переменная,
- $\beta_0, \beta_1$  — неизвестные параметры (коэффициенты),
- $\varepsilon$  — случайная ошибка

### Метод наименьших квадратов (МНК)

$$y = \beta_0 + \beta_1 x + \varepsilon$$

будем аппроксимировать  $\hat{y} = b_0 + b_1 x$



Величина  $\hat{y}_i$  — это расчетное значение переменной  $y_i$ , соответствующее  $x_i$ . Наблюдаемые значения  $y_i$  не лежат в точности на линии регрессии, т.е. не равны  $\hat{y}_i$ .

Определим остаток  $e_i$  в  $i$ -м наблюдении как разность между фактическим и расчетным значениями зависимой переменной:

$$e_i = y_i - \hat{y}_i \quad (7.16)$$

Неизвестные значения  $(b_0, b_1)$  определим **методом наименьших квадратов** (МНК), суть которого заключается в минимизации суммы квадратов остатков:

$$Q = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - b_0 - b_1 x_i)^2 \rightarrow \min. \quad (7.17)$$

Здесь  $(x_i, y_i)$  — известные значения наблюдения (числа),  $b_0, b_1$  — неизвестные.

Запишем необходимые условия экстремума — частные производные от  $Q$  по неизвестным параметрам  $b_0$  и  $b_1$  должны быть равны нулю:

$$\begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum (y_i - b_0 - b_1 x_i) = 0; \\ \frac{\partial Q}{\partial b_1} = -2 \sum (y_i - b_0 - b_1 x_i) x_i = 0. \end{cases}$$

В результате получим систему нормальных уравнений

$$\begin{cases} b_0 n - b_1 \sum x_i = \sum y_i, \\ b_0 \sum x_i + b_1 \sum x_i^2 = \sum x_i y_i, \end{cases} \quad \text{или} \quad \begin{cases} b_0 + b_1 \frac{1}{n} \sum x_i = \frac{1}{n} \sum y_i; \\ b_0 \frac{1}{n} \sum x_i + b_1 \frac{1}{n} \sum x_i^2 = \frac{1}{n} \sum x_i y_i. \end{cases}$$

Заметим, что  $\frac{1}{n} \sum x_i = \bar{x}$ ,  $\frac{1}{n} \sum y_i = \bar{y}$ ,  $\frac{1}{n} \sum x_i^2 = \overline{(x^2)}$  и  $\frac{1}{n} \sum x_i y_i = \overline{xy}$ . Тогда система примет простой вид:

$$\begin{cases} b_0 + b_1 \bar{x} = \bar{y}, \\ b_0 \bar{x} + b_1 \overline{(x^2)} = \overline{xy}. \end{cases}$$

Решая эту систему относительно  $b_0$  и  $b_1$ , получим:

$$b_1 = \frac{\overline{xy} - \bar{x} \bar{y}}{\overline{x^2} - (\bar{x})^2}, \quad b_0 = \bar{y} - b_1 \bar{x}. \quad (7.18)$$

Проверим достаточное условие существования минимума функции  $Q(b_0, b_1)$  при найденных значениях  $b_0$   $b_1$  (??). Напомним, что для функции двух переменных в некоторой точке наблюдается минимум, если выполнены условия:

$$\Delta = \begin{vmatrix} A & B \\ B & C \end{vmatrix} > 0 \quad \text{и} \quad A > 0$$

где  $A = \frac{\partial^2 Q}{\partial b_0^2}$ ,  $B = \frac{\partial^2 Q}{\partial b_0 \partial b_1}$  и  $C = \frac{\partial^2 Q}{\partial b_1^2}$ .

Вычислим вторые частные производные:

$$\begin{aligned} A &= \frac{\partial^2 Q}{\partial b_0^2} = \frac{\partial}{\partial b_0} [-2 \sum (y_i - b_0 - b_1 x_i)] = 2n, \\ B &= \frac{\partial^2 Q}{\partial b_0 \partial b_1} = \frac{\partial}{\partial b_1} [-2 \sum (y_i - b_0 - b_1 x_i)] = 2 \sum x_i, \\ C &= \frac{\partial^2 Q}{\partial b_1^2} = \frac{\partial}{\partial b_1} [-2 \sum (y_i - b_0 - b_1 x_i) x_i] = 2 \sum x_i^2. \end{aligned}$$

Составим требуемый определитель и проверим выполнение достаточных условий существования минимума:

$$\Delta = \begin{vmatrix} 2n & 2 \sum x_i \\ 2 \sum x_i & 2 \sum x_i^2 \end{vmatrix} = 2n \begin{vmatrix} 1 & \bar{x} \\ \bar{x} & \overline{x^2} \end{vmatrix} = 2n (\overline{x^2} - (\bar{x})^2) = 2n \hat{D}(x) > 0, \quad A = 2n > 0.$$

Значит при значениях  $b_0$   $b_1$ , определяемых по формулам (??), функция  $Q(b_0, b_1)$  достигает своего минимума.

## Свойства коэффициентов уравнения парной линейной регрессии

Рассматриваемая нами модель регрессии (7.14)

$$y_i = \beta_0 + \beta_1 \cdot x_i + \varepsilon_i. \quad (1)$$

Её оценкой по имеющимся данным является выборочное уравнение регрессии (7.15)

$$\hat{y}_i = b_0 + b_1 \cdot x_i. \quad (2)$$

МНК-оценки коэффициентов  $\beta_0$  и  $\beta_1$ , найденные нами ранее:

$$b_0 = \bar{y} - b_1 \bar{x}, \quad (3.a)$$

$$b_1 = \frac{\text{cov}(x, y)}{D(x)}. \quad (3.b)$$

Рассмотрим выражение (3.b). Подставим в него вместо  $y$  его выражение из (1):

$$\begin{aligned} b_1 &= \frac{\text{cov}(x, y)}{D(x)} = \frac{\text{cov}(x, \beta_0 + \beta_1 \cdot x + \varepsilon)}{D(x)} = \frac{\text{cov}(x, \beta_0) + \beta_1 \text{cov}(x, x) + \text{cov}(x, \varepsilon)}{D(x)} = \\ &= \frac{0 + \beta_1 D(x) + \text{cov}(x, \varepsilon)}{D(x)} = \beta_1 + \frac{\text{cov}(x, \varepsilon)}{D(x)}. \end{aligned}$$

Здесь учтено, что  $\text{cov}(x, \beta_0) = 0$ , т.к.  $\beta_0 = \text{const}$ , и то, что  $\text{cov}(x, x) = D(x)$ .

Рассмотрим выражение (3.b), подставив в него вместо  $b_1$  полученное только что выражение, а вместо  $\bar{y}$  — выражение  $\frac{1}{n} \sum y_i$ , где  $y_i$  определяется в (1):

$$\begin{aligned} b_0 &= \bar{y} - b_1 \bar{x} = \frac{1}{n} \sum (\beta_0 + \beta_1 \cdot x_i + \varepsilon_i) - \bar{x} \cdot \left[ \beta_1 + \frac{\text{cov}(x, \varepsilon)}{D(x)} \right] = \\ &= \frac{1}{n} \sum \beta_0 + \beta_1 \frac{1}{n} \sum x_i + \frac{1}{n} \sum \varepsilon_i - \bar{x} \beta_1 - \frac{\bar{x} \cdot \text{cov}(x, \varepsilon)}{D(x)} = \\ &= \frac{1}{n} \cdot n \cdot \beta_0 + \beta_1 \bar{x} + \frac{1}{n} \sum \varepsilon_i - \bar{x} \beta_1 - \frac{\bar{x} \cdot \text{cov}(x, \varepsilon)}{D(x)} = \beta_0 + \left[ \frac{1}{n} \sum \varepsilon_i - \frac{\bar{x} \cdot \text{cov}(x, \varepsilon)}{D(x)} \right]. \end{aligned}$$

Получили, что коэффициенты  $b_0$  и  $b_1$  можно выразить через «истинные» коэффициенты  $\beta_0$ ,  $\beta_1$  и значения независимой переменной  $x$ :

$$b_1 = \beta_1 + \frac{\text{cov}(x, \varepsilon)}{D(x)}, \quad (4)$$

$$b_0 = \beta_0 + \left[ \frac{1}{n} \sum \varepsilon_i - \frac{\bar{x} \cdot \text{cov}(x, \varepsilon)}{D(x)} \right]. \quad (5)$$

## Несмещённость и дисперсии оценок

Покажем, что оценки  $b_0$  и  $b_1$  коэффициентов  $\beta_0$  и  $\beta_1$  являются несмещенными и дисперсии их равны следующим величинам:

$$D(b_0) = \sigma_{b_0}^2 = \frac{\overline{x^2} \cdot \sigma_\varepsilon^2}{n \cdot D(x)}; \quad D(b_1) = \sigma_{b_1}^2 = \frac{\sigma_\varepsilon^2}{n \cdot D(x)}$$

Предварительно оценим (учитывая, что  $\bar{\varepsilon} = 0$ )

$$\text{cov}(x, \varepsilon) = \frac{1}{n} \sum_i (x_i - \bar{x})(\varepsilon_i - \bar{\varepsilon}) = \frac{1}{n} \sum_i (x_i - \bar{x})\varepsilon_i - \frac{\bar{\varepsilon}}{n} \sum_i (x_i - \bar{x}) = \frac{1}{n} \sum_i (x_i - \bar{x})\varepsilon_i$$

Принимая предположение теоремы Гаусса-Маркова о том, что  $x$  - неслучайная (детерминированная) величина, оценим математическое ожидание и дисперсию коэффициентов, используя их представление (4) и (5). Из (4) оценим:

$$M[\text{cov}(x, \varepsilon)] = M \left[ \frac{1}{n} \sum_i (x_i - \bar{x})\varepsilon_i \right] = \frac{1}{n} \sum_i (x_i - \bar{x})M[\varepsilon_i] = 0, \quad (6)$$

$$D[\text{cov}(x, \varepsilon)] = D \left[ \frac{1}{n} \sum_i (x_i - \bar{x})\varepsilon_i \right] = \frac{1}{n^2} \sum_i (x_i - \bar{x})^2 D(\varepsilon_i) = \frac{D(x)}{n} \sigma_\varepsilon^2. \quad (7)$$

Из (5) оценим:

$$M \left[ \frac{1}{n} \sum_i \varepsilon_i \right] = \frac{\sum M(\varepsilon_i)}{n} = 0; \quad (8)$$

$$D \left[ \frac{1}{n} \sum_i \varepsilon_i \right] = \frac{\sum D(\varepsilon_i)}{n^2} = \frac{\sigma_\varepsilon^2 \cdot n}{n^2} = \frac{\sigma_\varepsilon^2}{n}; \quad (9)$$

$$M \left[ \frac{\bar{x} \text{ cov}(x, \varepsilon)}{D(x)} \right] = \frac{\bar{x}}{D(x)} M[\text{cov}(x, \varepsilon)]_{(6)} = 0; \quad (10)$$

$$D \left[ \frac{\bar{x} \text{ cov}(x, \varepsilon)}{D(x)} \right] = \frac{(\bar{x})^2}{D^2(x)} D[\text{cov}(x, \varepsilon)]_{(7)} = \frac{(\bar{x})^2}{D^2(x)} \frac{D(x)}{n} \sigma_\varepsilon^2 = \frac{(\bar{x})^2 \sigma_\varepsilon^2}{n D(x)}. \quad (11)$$

Используя полученные оценки, найдем математическое ожидание и дисперсию коэффициентов  $b_0$  и  $b_1$ .

$$M(b_1) = M \left[ \beta_1 + \frac{\text{cov}(x, \varepsilon)}{D(x)} \right] = \beta_1 + \frac{M[\text{cov}(x, \varepsilon)]}{D(x)}_{(6)} = \beta_1$$

$$D(b_1) = D \left[ \beta_1 + \frac{\text{cov}(x, \varepsilon)}{D(x)} \right] = \frac{D[\text{cov}(x, \varepsilon)]}{D^2(x)}_{(7)} = \frac{D(x)}{n} \sigma_\varepsilon^2 \frac{1}{D^2(x)} = \frac{\sigma_\varepsilon^2}{n \cdot D(x)}$$

$$M(b_0) = M[\beta_0] + M \left[ \frac{1}{n} \sum_i \varepsilon_i \right]_{(8)} - M \left[ \frac{\bar{x} \text{ cov}(x, \varepsilon)}{D(x)} \right]_{(10)} = \beta_0;$$

$$\begin{aligned} D(b_0) &= D \left[ \frac{1}{n} \sum_i \varepsilon_i \right]_{(9)} + D \left[ \frac{\bar{x} \text{ cov}(x, \varepsilon)}{D(x)} \right]_{(11)} = \frac{\sigma_\varepsilon^2}{n} + \frac{(\bar{x})^2 \sigma_\varepsilon^2}{n D(x)} = \frac{\sigma_\varepsilon^2}{n} \left[ 1 + \frac{(\bar{x})^2}{D(x)} \right] = \\ &= \frac{\sigma_\varepsilon^2 \overline{x^2}}{n D(x)} \quad \text{так как} \quad 1 + \frac{(\bar{x})^2}{D(x)} = \frac{D(x) + (\bar{x})^2}{D(x)} = \frac{(\overline{x^2} - (\bar{x})^2) + (\bar{x})^2}{D(x)} = \frac{\overline{x^2}}{D(x)} \end{aligned}$$

Оценкой дисперсии ошибки  $\sigma_\varepsilon^2$  по выборке служит величина

$$S_e^2 = \frac{Q_e}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}.$$

Тогда для коэффициентов  $b_0$  и  $b_1$  выборочного уравнения регрессии имеем:

$$M(b_0) = \beta_0, \quad M(b_1) = \beta_1.$$

Это означает, что найденные оценки являются несмещенными.

Дисперсии и средние квадратические отклонения коэффициентов  $b_0$  и  $b_1$ :

$$\begin{aligned} D(b_0) &= \frac{S_e^2 \overline{x^2}}{nD(x)}, & S_{b_0} &= \frac{S_e \sqrt{\overline{x^2}}}{\sqrt{n}\sigma_x}, \\ D(b_1) &= \frac{S_e^2}{nD(x)}, & S_{b_1} &= \frac{S_e}{\sqrt{n}\sigma_x}. \end{aligned}$$