

## 19. Свойства коэффициентов уравнения парной линейной регрессии.

*Определение 7.13.* Регрессионный анализ — статистический метод, предназначенный для определения **вида** (аналитического выражения) **связи** зависимой случайной величины  $y$  (называемой результативным признаком) от независимых случайных величин  $X_1, X_2, \dots, X_p$  (называемых факторами или объясняющими переменными).

Форма связи результативного признака  $Y$  с факторами  $X = (X_1, \dots, X_p)^T$  получила название **уравнения регрессии**:

$$y = f(X, \varepsilon),$$

где случайная составляющая  $\varepsilon$  характеризует, с одной стороны, влияние на  $y$  не входящих в  $X$  факторов, а с другой — погрешность в измерении показателя  $y$ .

**Модель парной линейной регрессии имеет вид:**

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Индексы  $i$  показывают, что модель записана для **каждого наблюдения** в выборке.

где:

- $y_i$  — **зависимая переменная** (отклик) для  $i$ -го наблюдения,
- $x_i$  — **независимая переменная** (фактор, предиктор) для  $i$ -го наблюдения,
- $\varepsilon_i$  — **случайная ошибка** для  $i$ -го наблюдения.
- $\beta_0$  — **свободный член** (константа, intercept),
- $\beta_1$  — **коэффициент наклона** (slope),

**Смысл коэффициентов:**

- $\beta_0$  (intercept) — это **ожидаемое значение**  $y$ , когда  $x = 0$  (если такое значение имеет смысл в контексте задачи).
- $\beta_1$  (slope) — показывает, **насколько в среднем изменяется**  $y$  при увеличении  $x$  на одну единицу.

Так как мы оцениваем  $\beta_0 = b_0$  и  $\beta_1 = b_1$  стоит пояснить что за коэффициенты  $b_0$  и  $b_1$ :

Коэффициент  $b_1$  — угловой коэффициент регрессии, он показывает, на сколько единиц в среднем меняется переменная  $y$  при увеличении независимой переменной  $x$  на единицу.

Постоянная  $b_0$  дает прогнозируемое значение зависимой переменной при  $x = 0$ . Это может иметь смысл в зависимости от того, как далеко находится  $x = 0$  от выборочных значений  $x$ .

**Пример:**

Если  $y$  — зарплата (в тыс. руб.),  $x$  — стаж (в годах), и

$$y = 30 + 5x + \varepsilon,$$

то:

- $\beta_0 = 30$  — ожидаемая зарплата при нулевом стаже,
- $\beta_1 = 5$  — при увеличении стажа на 1 год зарплата в среднем растёт на 5 тыс. руб.

## 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):

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

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

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

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

## Несмешённость и дисперсии оценок

Покажем, что оценки  $b_0$  и  $b_1$  коэффициентов  $\beta_0$  и  $\beta_1$  являются несмешенными 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:

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

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

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

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

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

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

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

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

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

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 несмешенными.

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

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

### Статистические свойства (Что мы рассмотрели выше)

При использовании метода наименьших квадратов (МНК) для оценки коэффициентов, они обладают набором оптимальных свойств.:

- **Несмешенность:** Математическое ожидание оценок  $b_0$  и  $b_1$  совпадает с истинными значениями параметров в генеральной совокупности, что гарантирует **отсутствие систематической ошибки** при оценивании.
- **Эффективность:** Оценки МНК имеют **минимально возможную дисперсию** (разброс) среди всех линейных несмешенных оценок, что делает их наиболее надежными.
- **Состоительность:** С увеличением объема выборки значения коэффициентов стремятся (по вероятности) к своим теоретическим значениям.

**Важное замечание:** Эти свойства выполняются только при соблюдении предпосылок регрессионного анализа, таких как **гомоскедастичность** (равенство дисперсий ошибок) и отсутствие их автокорреляции.