

Эконометрика

Модель линейной регрессии

Шишкин Владимир Андреевич

Пермский государственный национальный исследовательский
университет

Вероятностью $P(A)$ события A называется численная мера степени объективной возможности появления этого события.

Случайная величина X — функция, заданная на множестве элементарных событий.

Закон распределения случайной величины X — всякое соотношение, устанавливающее связь между возможными значениями случайной величины и соответствующими им вероятностями.

Случайные величины

Характеристики случайных величин

- Математическое ожидание:

$$E(X) = \sum_{i=1}^n x_i p_i \quad E(X) = \int_{-\infty}^{\infty} xp(x) dx$$

- Дисперсия и среднее квадратическое отклонение (стандартное отклонение):

$$D(x) = E(X - E(X))^2 \quad \sigma_X = \sqrt{D(X)}$$

- x_q — квантиль уровня q , если

$$F(x_q) = P(x \leq x_q) = q$$

- Начальные ν_k и центральные μ_k моменты k -го порядка:

$$\nu_k = E(X^k) \quad \mu_k(X) = E(X - E(X))^k$$

Характеристики многомерных случайных величин

- Ковариация (ковариационный момент):

$$\text{Cov}(X, Y) = E\left(\left(X - E(X)\right)\left(Y - E(Y)\right)\right)$$

- Коэффициент корреляции:

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Оценка $\tilde{\theta}_n$ параметра θ

Всякая функция результатов наблюдений над случайной величиной X , с помощью которой судят о величине θ .

- *Несмещённая* оценка: $E(\tilde{\theta}_n) = \theta$.
- *Состоятельная* оценка: для любого $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\tilde{\theta}_n - \theta| \leq \varepsilon)$$

- *Эффективная* оценка — несмещённая оценка, имеющая наименьшую дисперсию среди всех несмещённых оценок, вычисленных по выборкам одного и того же объёма.

Оценки параметров нормального распределения

Точечные оценки

- Математическое ожидание (среднее арифметическое):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Дисперсия (исправленная выборочная дисперсия):

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- Стандартное отклонение:

$$s = \sqrt{s^2}$$

Оценки параметров нормального распределения

Интервальные оценки

Интервальная оценка параметра θ

Числовой интервал $(\tilde{\theta}_n^{(1)}, \tilde{\theta}_n^{(2)})$, который с заданной вероятностью γ накрывает неизвестное значение параметра θ .

- Доверительный интервал для математического ожидания на уровне значимости α :

$$\left(\bar{x} + t_{\frac{\alpha}{2}, n-1} \frac{s}{\sqrt{n-1}}, \bar{x} + t_{1-\frac{\alpha}{2}, n-1} \frac{s}{\sqrt{n-1}} \right)$$

- Доверительный интервал для генеральной дисперсии на уровне значимости α :

$$\left(\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right)$$

Статистическая гипотеза

Любое предположение о виде или параметре неизвестного закона распределения.

Гипотеза H_0	Принимается	Отвергается
Верна	Правильно	Ошибка первого рода
Не верна	Ошибка второго рода	Правильно

Уровень значимости критерия — вероятность α допустить ошибку первого рода.

Мощность критерия — вероятность $(1 - \beta)$ не допустить ошибку второго рода.

Проверка статистических гипотез для величин, имеющих нормальное распределение

$$H_0: \bar{x} = \mu:$$

t -критерий Стьюдента

$$\frac{\bar{x} - \mu}{s} \sqrt{n} \sim t_{n-1}$$

T^2 -критерий Хотеллинга

$$x = (x_1, \dots, x_m)$$

$$n(\bar{x} - \mu)^\top S^{-1}(\bar{x} - \mu) \sim \frac{m(n-1)}{n-m} F_{m, n-m}$$

Проверка статистических гипотез для величин, имеющих нормальное распределение

$$H_0: \bar{x}_1 = \bar{x}_2:$$

t-критерий Стьюдента

$$\frac{(\bar{x}_1 - \bar{x}_2)\sqrt{n_1 n_2}}{s_* \sqrt{n_1 + n_2}} \sim t_{n_1 + n_2 - 2} \quad s_* = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

T^2 -критерий Хотеллинга

$$n_1 n_2 (\bar{x}_1 - \bar{x}_2)^\top S_*^{-1} (\bar{x}_1 - \bar{x}_2) \sim \frac{(n_1 + n_2 - 2)m}{n_1 + n_2 - m - 1} F_{m, n_1 + n_2 - m - 1}$$

$$S_* = \frac{1}{n_1 + n_2 - 2} (K_1^\top K_1 + K_2^\top K_2) \quad K_{ij} = x_{ij} - \bar{x}_j$$

Парная линейная регрессионная модель

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

x — объясняющая переменная (регрессор)

y — объясняемая переменная

ε — ошибка модели

β_1 — величина изменения y при изменении x на единицу

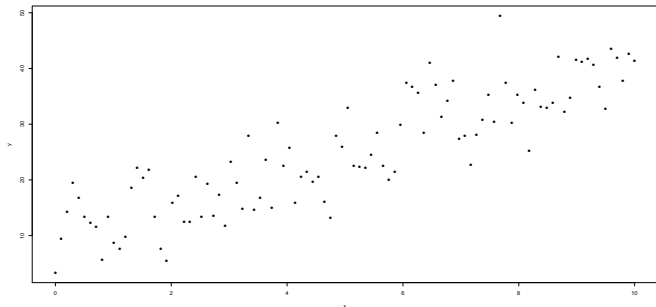
β_0 — величина y при нулевом значении x

Исходные данные

$(x_i, y_i), i = 1, 2, \dots, n.$

Случайные
величины

Парная
линейная
регрессия



Задача подгонки кривой

Подобрать функцию $f(x)$ из параметрического семейства $f(x | \beta)$ вида

$$f(x | \beta) = \beta_0 + \beta_1 x$$

наилучшим образом описывающую зависимость y от x :

$$y_i \approx f(x_i), \quad i = 1, \dots, n.$$

Оценка близости

- сумма квадратов отклонений $\sum_{i=1}^n (y_i - f(x_i | \beta))^2$;
- сумма модулей $\sum_{i=1}^n |y_i - f(x_i | \beta)|$;
- ...

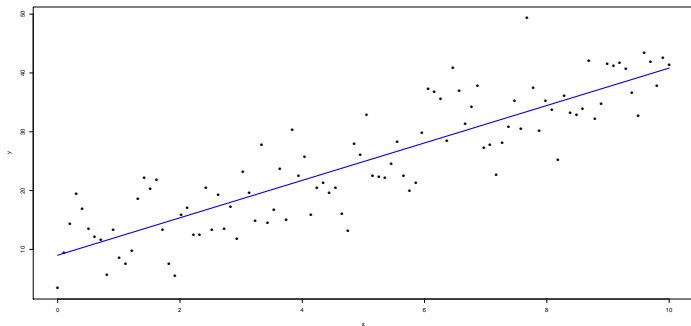
Случайные
величиныПарная
линейная
регрессия

Парная линейная регрессия

Линейная
регрессия

Случайные
величины

Парная
линейная
регрессия



Парная линейная регрессия

Оценка коэффициентов регрессии

Случайные
величины

Парная
линейная
регрессия

Метод наименьших квадратов

$$ESS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \min$$

ESS — сумма квадратов ошибок (Errors' Sum of Square)

Парная линейная регрессия

Оценка коэффициентов регрессии

Случайные
величины

Парная
линейная
регрессия

Метод наименьших квадратов

$$ESS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \min$$

Необходимые условия существования экстремума:

$$\frac{\partial ESS}{\partial \beta_0} = 0, \quad \frac{\partial ESS}{\partial \beta_1} = 0$$

Парная линейная регрессия

Оценка коэффициентов регрессии

Метод наименьших квадратов

$$ESS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \min$$

Необходимые условия существования экстремума:

$$\frac{\partial ESS}{\partial \beta_0} = 0, \quad \frac{\partial ESS}{\partial \beta_1} = 0$$

Стандартная система нормальных уравнений:

$$b_0 n + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i, \quad b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i$$

Случайные
величины

Парная
линейная
регрессия

Парная линейная регрессия

Оценка коэффициентов регрессии

Случайные
величины

Парная
линейная
регрессия

Метод наименьших квадратов

$$ESS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \min$$

Необходимые условия существования экстремума:

$$\frac{\partial ESS}{\partial \beta_0} = 0, \quad \frac{\partial ESS}{\partial \beta_1} = 0$$

Стандартная система нормальных уравнений:

$$b_0 + b_1 \bar{x} = \bar{y}, \quad b_0 \bar{x} + b_1 \overline{x^2} = \overline{yx}$$

Парная линейная регрессия

Оценка коэффициентов регрессии

Случайные
величиныПарная
линейная
регрессия

$$b_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \quad b_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2}$$

$$b_0 = \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n x_i b_1 \quad b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\text{cov}(x, y)}{s_x^2}$$

Парная линейная регрессия

Оценка коэффициентов регрессии (матричная форма записи)

Систему уравнений

$$y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_2 + \varepsilon_1$$

...

$$y_n = \beta_0 + \beta_1 x_n + \varepsilon_1$$

можно записать в виде

$$y = X\beta + \varepsilon$$

где

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Случайные
величины

Парная
линейная
регрессия

Парная линейная регрессия

Оценка коэффициентов регрессии (матричная форма записи)

Случайные
величины

Парная
линейная
регрессия

$$y = X\beta + \varepsilon$$

$$X^T y = X^T X \beta + X^T \varepsilon$$

$$X^T y = X^T X b$$

$$b = (X^T X)^{-1} X^T y$$

Парная линейная регрессионная модель

Основные гипотезы

- 1 $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i = 1, 2, \dots, n$ — модель правильно специфицирована
- 2 x — детерминированная величина, вектор (x_1, \dots, x_n) не коллинеарен вектору $(1, 1, \dots, 1)^T$
- 3 $E \varepsilon_i = 0, E (\varepsilon_i^2) = D (\varepsilon_i) = \sigma^2$
- 4 $E (\varepsilon_i \varepsilon_j) = 0$ при $i \neq j$
- 5 $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Теорема Гаусса–Маркова

В предположениях модели 1–4 оценки b_0, b_1 , полученные по методу наименьших квадратов, имеют наименьшую дисперсию в классе всех линейных несмещённых оценок.

Парная линейная регрессионная модель

Дисперсия ошибок

Случайные
величиныПарная
линейная
регрессия

Прогноз значения y_i в точке x_i (модельное значение):

$$\hat{y}_i = b_0 + b_1 x_i, \quad i = 1, 2, \dots, n$$

Остатки модели:

$$e_i = y_i - \hat{y}_i, \quad i = 1, 2, \dots, n$$

Несмещённая оценка дисперсии ошибок σ_ε^2 :

$$s_e^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2$$

Парная линейная регрессионная модель

Оценки дисперсии параметров модели

Случайные
величиныПарная
линейная
регрессия

$$s_{b_1}^2 = \frac{s_e^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$
$$s_{b_0}^2 = \frac{s_e^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}$$
$$\text{cov}(b_0, b_1) = -\frac{s_e^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Парная линейная регрессионная модель

Проверка гипотезы $H_0: \beta = \beta^*$

$$t_{i,n} = \frac{b_i - \beta_i^*}{sb_i} \sim t_{n-2}$$

Случайные
величины

Парная
линейная
регрессия

t -критерий Стьюдента

Если

$$\frac{|b_i|}{sb_i} > t_{1-\frac{\alpha}{2}, n-2},$$

то коэффициент b_i линейной модели значим при уровне значимости α (отклоняется нулевая гипотеза $\beta_i = 0$).

Доверительный интервал для β_i :

$$b_i + t_{\frac{\alpha}{2}, n-2} \cdot sb_i \leq \beta_i \leq b_i + t_{1-\frac{\alpha}{2}, n-2} \cdot sb_i$$

Парная линейная регрессионная модель

Анализ вариации зависимой переменной в регрессии

Случайные
величины

Парная
линейная
регрессия

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{TSS}} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{ESS}} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{RSS}}$$

TSS — вся дисперсия y (total sum of squares)

ESS — не объяснённая дисперсия (error sum of squares)

RSS — объяснённая моделью часть дисперсии y (regression sum of squares)

Коэффициент детерминации

$$R^2 = 1 - \frac{\text{ESS}}{\text{TSS}} = \frac{\text{RSS}}{\text{TSS}}$$

Парная линейная регрессионная модель

Проверка значимости модели

Случайные
величиныПарная
линейная
регрессия

F -критерий Фишера

Если выполняется условие

$$F_H = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\frac{1}{n-2} \sum_{i=1}^n e_i^2} > F_{1-\alpha, 1, n-2},$$

то модель значима на уровне значимости α .

$$F_H = (n - 2) \frac{R^2}{1 - R^2}$$