

Регрессионный анализ данных



Оценки модели линейной регрессии. Построение модели

$$(1) Y_i = \eta(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

где $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_q]'$ - q -мерная *неслучайная* векторная переменная

ε_i – случайная ошибка

$$\eta(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_q \mathbf{x}_q.$$

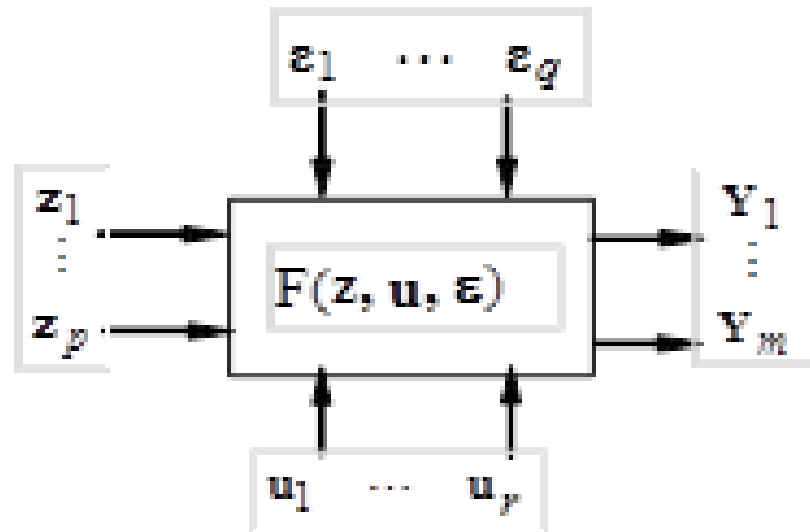
где $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_q]'$ – вектор неизвестных параметров (коэффициентов)

$$(2) \quad \mathbf{Y} = \beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_q \mathbf{x}_q + \boldsymbol{\varepsilon}$$

Если в (2.2) $\mathbf{M}\boldsymbol{\varepsilon} = 0$

$$(3) \quad \mathbf{M}[\mathbf{Y} / \mathbf{x}] = \beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_q \mathbf{x}_q$$

Пример 1



1. Контролируемые неуправляемые переменные $z = (z_1, \dots, z_p)$,
2. Контролируемые управляющие переменные $u = (u_1, \dots, u_r)$,
3. Неконтролируемые неуправляемые переменные $\epsilon = (\epsilon_1, \dots, \epsilon_q)$,
4. Контролируемые управляемые переменные $Y = (Y_1, \dots, Y_m)$,

Линейная регрессионная модель

$\eta(\mathbf{x}, \beta)$ - нелинейная относительно вектора параметров β

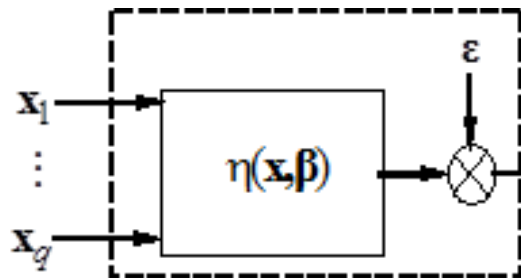


Рис. 2. Структурная схема объекта

$$Y_i = \eta(0, \dots, 0) + \sum_{j=1}^q \left(\frac{\partial \eta}{\partial x_j} \right)_0 x_{ij} + \frac{1}{2} \sum_j \sum_k \left(\frac{\partial^2 \eta}{\partial x_j \partial x_k} \right)_0 x_{ij} x_{ik} + \dots + \varepsilon_i$$

Обозначив постоянные

$$\beta_0 = (0, \dots, 0) \quad \beta_j = (\partial \eta / \partial x_j)_0 \quad \beta_{jk} = (1/2)(\partial^2 \eta / \partial x_j \partial x_k)_0$$

получим

$$Y_i = \beta_0 + \sum_{j=1}^q \beta_j x_{ij} + \sum_j \sum_k \beta_{jk} x_{ij} x_{ik} + \dots + \varepsilon_i$$

полиномиальная модель q-го порядка одной переменной

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_q x_i^q + \varepsilon_i.$$

Нелинейная модель

$$Y_i = \beta_0 + \beta_1 e^{-\beta_2 x_i} + \varepsilon_i$$

Оценивание параметров

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & \cdots & x_{1g} \\ x_{20} & x_{21} & \cdots & x_{2g} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \cdots & x_{ng} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_g \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

(4) $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

Метод наименьших квадратов (МНК) (least squares method)

Необходимо по наблюдениям $(x_{i1}, \dots, x_{ig}, Y_i)$, $i = 1, \dots, n$ найти наилучшую оценку вектора $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_g)'$ уравнения регрессии $\eta(x, \boldsymbol{\beta})$.

Метод наименьших квадратов

1. Число регрессоров q полинома (3) априори известно точно.
2. Все регрессоры измеряются без ошибок, а вычисления проводятся абсолютно точно.
3. Остаток ε является независимой нормально распределенной случайной величиной с нулевым средним $M[\varepsilon] = 0$ и неизвестной постоянной дисперсией σ_ε^2 при всех $i = 1, \dots, n$.
4. Дисперсия отклика Y_i постоянна, или является известной функцией номера наблюдения $i = 1, \dots, n$.
5. Распределение Y_i одинаково при всех $i = 1, \dots, n$.
6. Число опытов n существенно больше числа регрессоров q .

Если условия 1÷6 соблюдаются и $X'X$ обратимая матрица, то согласно фундаментальной теореме Гаусса-Маркова наилучшей оценкой вектора коэффициентов β является оценка $\hat{\beta}$, доставляющая минимум суммы квадратов остатков (невязок, ошибок, помех):

$$Q(\hat{\beta}) = \sum_{i=1}^n \left(Y_i - \sum_{k=0}^q x_{ik} \hat{\beta}_k \right)^2 \rightarrow \min$$

Система нормальных уравнений МНК

$$\partial Q / \partial \hat{\beta}_m = -2 \sum_{i=1}^n (Y_i - \sum_{k=0}^q x_{ik} \hat{\beta}_k) x_{im} = 0$$

$$(5) \quad \sum_{i=1}^n x_{im} \sum_{k=0}^q x_{ik} \hat{\beta}_k = \sum_{i=1}^n x_{im} Y_i$$
$$m = 0, 1, \dots, q$$

Свойства оценок

Обозначим $\theta = \mathbf{X}\beta$, $\hat{\theta} = \mathbf{X}\hat{\beta}$. Будем минимизировать величину $\varepsilon'\varepsilon = \|\mathbf{Y} - \theta\|^2$ по отношению к $\theta \in \Omega$, где Ω – подпространство оценок $\hat{\theta}$. Если изменять значения вектора θ в пределах Ω , то квадрат длины вектора $\|\mathbf{Y} - \theta\|^2$ достигнет минимума при значении $\theta = \hat{\theta}$, которое является проекцией вектора \mathbf{Y} на подпространство Ω . Тогда справедливо $(\mathbf{Y} - \hat{\theta}) \perp \hat{\theta}$ и, следовательно, $(\mathbf{Y} - \hat{\theta}) \perp \mathbf{X}$ (рис. 2.3). Отсюда для скалярного произведения $(\mathbf{Y} - \hat{\theta})$ и \mathbf{X} получаем

$$(\mathbf{Y} - \hat{\theta}, \mathbf{X}) = 0$$

$$(6) \quad \mathbf{X}'\hat{\theta} = \mathbf{X}'\mathbf{Y}$$

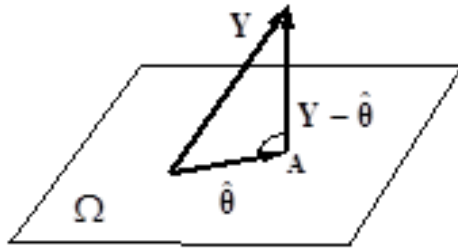


Рис. 2.3. Нахождение точки А, для которой норма минимальна

Свойства матрицы C

Обозначим $C = X'X$, $\Psi = X'Y$, тогда СЛУ запишется в виде

$$(7) \quad C\beta = \Psi$$

Свойства матрицы C :

а) так как регрессоры x_1, \dots, x_q линейно независимы, то матрица C не вырождена;

б) C положительно определена и ранг ее в точности равен q ;

с) C – симметричная матрица, $C = C'$, т. е. является эрмитовой.

Поэтому СЛУ (6) имеет единственное решение

$$(8) \quad \hat{\beta} = C^{-1}\Psi = (X'X)^{-1}X'Y$$

Оценка называется оценкой метода наименьших квадратов (МНК-оценкой). Так как в решении (8) матрица $A = (X'X)^{-1}X'$ неслучайная, то является линейной комбинацией наблюдений Y . В соответствии с теоремой Гаусса-Маркова МНК-оценка имеет наименьшую дисперсию среди всех возможных несмещенных линейных оценок.

Свойства оценок

1. Несмещенность

МНК-оценка $\hat{\beta}$ является случайной величиной

$$\begin{aligned} M[\hat{\beta}] &= M[(X'X)^{-1} X'Y] = M[(X'X)^{-1} X'(X\beta + \varepsilon)] = \\ &= (X'X)^{-1} X'X\beta + (X'X)^{-1} M[\varepsilon] = \beta. \end{aligned}$$

2. Распределение

$M[\varepsilon] = 0$, $D[\varepsilon] = \sigma_\varepsilon^2 I_n$, или $\varepsilon : N_n(0, \sigma_\varepsilon^2 I_n)$, то $Y : N_n(X\beta, \sigma_\varepsilon^2 I_n)$

$$1) \hat{\beta} \sim N_n(\beta, \sigma_\varepsilon^2 (X'X)^{-1}) \quad 2) (\hat{\beta} - \beta)' X'X (\hat{\beta} - \beta) / \sigma_\varepsilon^2 \sim \chi_{(q+1)}^2$$

$$3) \hat{\beta} \text{ не зависит от } S_\varepsilon^2 \quad 4) SSE / \sigma_\varepsilon^2 = (n - q - 1) S_\varepsilon^2 / \sigma_\varepsilon^2 \sim \chi_{(n-q-1)}^2$$

$$S_\varepsilon^2 = \frac{e'e}{n-q-1} = \frac{(Y - X\hat{\beta})'(Y - X\hat{\beta})}{n-q-1} = \frac{\sum_{i=1}^n e_i^2}{n-q-1} = \frac{SSE}{n-q-1}$$

где $e_i = Y_i - \hat{Y}_i$ – оценка остатка ε , а SSE (Sum of Squares of Errors) – сумма квадратов (оцененных) остатков

Списки использованной литературы и источников:

- А.А.Большаков, Р.Н.Каримов «Методы обработки многомерных данных и временных рядов» Москва 2007 г.
- Электронный учебник StatSoft по анализу данных.