

СТАТИСТИЧЕСКИЕ КРИТЕРИИ ДЛЯ ПОРЯДКА РЕГРЕССИОННОЙ МОДЕЛИ ПО ЕДИНСТВЕННОЙ СЕРИИ НАБЛЮДЕНИЙ

We construct an order definition test for a linear regression model on the basis of unique series of observations in the case of normally distributed noise with zero mean and limited dispersion. We propose a new statistics of equilibrium for the model with infinite (but unique) series of observations and, by using this statistics, construct an order definition test for this model.

Побудовано критерій визначення порядку лінійної регресійної моделі за єдиною серією спостережень для випадку нормального розподіленого шуму з пульсовим середнім та обмеженою дисперсією. Запропоновано нову статистику рівноваги для моделі з нескінченною (але єдиною) серією спостережень і на її основі побудовано критерій визначення порядку моделі.

1. Введение. Определение порядка модели регрессии в случае единственной серии наблюдений — довольно сложная задача. При единственной серии остаточные дисперсии необходимо являются зависимыми. Более того, оказывается, при большой выборке коэффициент корреляции стремится к единице. Однако при определенных условиях можно построить критерий для определения порядка такой модели.

В статье рассматривается модель линейной регрессии, ошибки которой имеют нормальное распределение с нулевым математическим ожиданием и положительно определенной матрицей ковариаций. Рассматриваем нулевую гипотезу о порядке регрессионной модели и альтернативную ей простую гипотезу. Для данной модели построен критерий проверки гипотез и предельное распределение этой статистики. Конструктивный метод, примененный в статье, основан на свойствах известных распределений.

2. Постановка задачи и проблематика. Приведем некоторые предположения относительно рассматриваемой модели: имеем единственную (!) серию n наблюдений:

$$(x_i, y_i), \quad i = \overline{1, n}, \quad (1)$$

в модели линейной регрессии

$$y = X\beta + \varepsilon; \quad (2)$$

ошибки ε имеют нормальное распределение с нулевым математическим ожиданием и положительно определенной матрицей ковариаций

$$\begin{aligned} E(\varepsilon) &= 0, \\ \text{Cov}(\varepsilon) &= \sigma^2 \Omega, \quad \sigma^2 < \infty. \end{aligned} \quad (3)$$

Будем рассматривать гипотезу о порядке модели (1)–(3). За основную гипотезу принимаем гипотезу H_p : модель (1)–(3) имеет порядок p , где $p < n$. Предполагаем, что $\text{rank}(X) = p$.

Известно, что для любой положительно определенной матрицы Ω существует такая невырожденная матрица P размера $n \times n$, что

$$PP' = \Omega.$$

Умножая (2) слева на P^{-1} , получаем $P^{-1}y = P^{-1}X\beta + P^{-1}\varepsilon$. Обозначим $u = P^{-1}y$, $Q = P^{-1}X$, $z = P^{-1}\varepsilon$. Тогда (2) можно представить в виде

$$z = Q\beta + u. \quad (4)$$

Ошибки модели (4) имеют нулевое математическое ожидание

$$E(u) = P^{-1}E(\varepsilon) = 0 \quad (5)$$

и ковариационную матрицу

$$\begin{aligned} Cov(u) &= E(uu') = E(P^{-1}\varepsilon\varepsilon'(P^{-1})') = P^{-1}E(\varepsilon\varepsilon')(P')^{-1} = \\ &= P^{-1}\sigma^2 PP'(P')^{-1} = \sigma^2 I. \end{aligned} \quad (6)$$

По методу наименьших квадратов получаем оценку Гаусса – Маркова вектора параметров β :

$$\hat{\beta}_p = (Q'Q)^{-1}Q'z,$$

или в обозначениях исходной модели

$$\hat{\beta}_p = (X'\Omega^{-1}X)^{-1}X'\Omega^{-1}y,$$

где матрица $X'\Omega^{-1}X$ невырождена [1, с. 379].

Оценка $\hat{\beta}_p$ не смешена:

$$E(\hat{\beta}_p) = (Q'Q)^{-1}Q'E(z) = (Q'Q)^{-1}Q'Q\beta + (Q'Q)^{-1}Q'E(z) = \beta,$$

а ее ковариационная матрица имеет вид

$$\begin{aligned} Cov(\hat{\beta}_p) &= E\left(\left(\hat{\beta}_p - E(\hat{\beta}_p)\right)\left(\hat{\beta}_p - E(\hat{\beta}_p)\right)'\right) = \\ &= (Q'Q)^{-1}Q'E(uu')Q(Q'Q)^{-1} = \sigma^2(Q'Q)^{-1} = \sigma^2(X'\Omega^{-1}X)^{-1}. \end{aligned}$$

Учитывая предположение (3), а также используя (5) и (6), получаем $u \sim N(0, \sigma^2 I)$. Обозначим остаточную сумму квадратов S_p :

$$S_p = (z - Q\hat{\beta}_p)'(z - Q\hat{\beta}_p),$$

или в обозначениях исходной модели:

$$S_p = (y - X\hat{\beta}_p)' \Omega^{-1} (y - X\hat{\beta}_p).$$

Тогда несмешенной оценкой для σ^2 будет остаточная дисперсия:

$$\hat{S}_p^2 = \frac{S_p}{v}, \quad \text{где } v = n - p.$$

Рассмотрим альтернативную гипотезу о порядке модели — гипотезу H_{p+k} . Оценку методом наименьших квадратов вектора β для этой гипотезы обозначим $\hat{\beta}_{p+k}$. Соответственно обозначим остаточную сумму квадратов через S_{p+k} и остаточную дисперсию —

$$\hat{S}_{p+k}^2 = \frac{S_{p+k}}{v-k}.$$

Тогда справедлива следующая теорема.

Теорема 1. Остаточные дисперсии альтернативных гипотез \hat{S}_p^2 и \hat{S}_{p+k}^2 модели (1)–(3) коррелированы с коэффициентом корреляции

$$r = \sqrt{1 - k/v}. \quad (7)$$

Доказательство. Легко убедиться, что $\hat{\beta}_p \sim N(\beta, \sigma^2(X'\Omega^{-1}X)^{-1})$. Величина $\hat{S}_p^2 \sim \frac{\sigma^2}{v}\chi_v^2$. $\hat{\beta}_p$ не зависит от \hat{S}_p^2 . Используя свойства χ^2 -распределения, получаем $E(\hat{S}_p^2) = \sigma^2$, $D(\hat{S}_p^2) = 2\sigma^4/v$. Аналогично $E(\hat{S}_{p+k}^2) = \sigma^2$, $D(\hat{S}_{p+k}^2) = 2\sigma^4/(v-k)$. Вычислим величину $E(\hat{S}_p^2 \hat{S}_{p+k}^2)$:

$$\begin{aligned} E(\hat{S}_p^2 \hat{S}_{p+k}^2) &= \frac{2\sigma^4}{v(v-k)} E(\chi_v^2 \chi_{v-k}^2) = \frac{2\sigma^4}{v(v-k)} E\left(\sum_{j=p+1}^n \xi_j^2 \sum_{j=p+k+1}^n \xi_j^2\right) = \\ &= \frac{2\sigma^4}{v(v-k)} E\left(\left(\xi_{p+1}^2 + \xi_{p+2}^2 + \dots + \xi_{p+k}^2 + \sum_{j=p+k+1}^n \xi_j^2\right) \sum_{j=p+k+1}^n \xi_j^2\right) = \\ &= \frac{2\sigma^4}{v(v-k)} (k(v-k) + 3(v-k) + (v-k)^2 - (v-k)) = \sigma^2\left(1 + \frac{2}{v}\right). \end{aligned}$$

Теперь вычислим коэффициент корреляции r :

$$\begin{aligned} r &= \frac{E(\hat{S}_p^2 \hat{S}_{p+k}^2) - E(\hat{S}_p^2)E(\hat{S}_{p+k}^2)}{\sqrt{D(\hat{S}_p^2)D(\hat{S}_{p+k}^2)}} = \frac{\sigma^4\left(1 + \frac{2}{v}\right) - \sigma^4}{\sqrt{\frac{2\sigma^4}{v} \frac{2\sigma^4}{v-k}}} = \\ &= \frac{1}{v} \sqrt{v(v-k)} = \sqrt{1 - k/v}. \end{aligned}$$

Теорема доказана.

Из теоремы 1 следует, что остаточные дисперсии асимптотически линейно зависимы.

3. Распределение t -статистики. Критерий для определения порядка модели.

Определение 1. Статистику отношения остаточных дисперсий будем называть t -статистикой модели (1)–(3):

$$T_{p,k}^{(n)} = \frac{\hat{S}_p^2}{\hat{S}_{p+k}^2}. \quad (8)$$

Теорема 1 явным образом отрицает использование критерия для проверки гипотезы о порядке модели [2]. Попробуем найти распределение статистики (8).

Определение 2. Обратной нормализованной t -статистикой назовем величину

$$\zeta_{p,k}^{(n)} = \frac{r^2}{T_{p,k}^{(n)}}. \quad (9)$$

Лемма 1. Обратная нормализованная t -статистика $\zeta_{p,k}^{(n)}$ модели (1)–(3) имеет β -распределение с параметрами $(v-k)/2$ и $k/2$.

Доказательство. На основании теоремы 1 получаем

$$\begin{aligned} \zeta_{p,k}^{(n)} &= \frac{r^2}{T_{p,k}^{(n)}} = r^2 \frac{v}{v-k} \frac{\hat{S}_p^2}{\hat{S}_{p+k}^2} = \frac{\sum_{j=p+k+1}^n \xi_j^2}{\sum_{j=p+1}^n \xi_j^2} = \end{aligned}$$

$$= \frac{\sum_{j=p+k+1}^n \xi_j^2}{\sum_{j=p+k+1}^n \xi_j^2 + \sum_{j=p+1}^{p+k} \xi_j^2} = \frac{\chi_{v-k}^2}{\chi_{v-k}^2 + \chi_k^2},$$

причем случайные величины в знаменателе независимы, а отсюда из свойств β -распределения имеем в правой части случайную величину с распределением $\beta((v-k)/2, k/2)$.

Теорема 2. Плотность распределения t -статистики имеет вид

$$f_{T_{p,k}^{(n)}}(x) = \frac{\Gamma\left(\frac{v}{2}\right)}{\Gamma\left(\frac{v-k}{2}\right)\Gamma\left(\frac{k}{2}\right)} r^{v-2} \left(\frac{x}{r^2} - 1\right)^{\frac{k}{2}-1} x^{\frac{v}{2}}, \quad \text{где } x \geq r^2.$$

Доказательство. Используем лемму 1. Производя замену переменных $u = r^2/x$, получаем

$$f_{\zeta_{p,k}^{(n)}}(x) = \begin{cases} \frac{\Gamma(v_1 + v_2)}{\Gamma(v_1)\Gamma(v_2)} u^{v_1-1} (1-u)^{v_2-1}, & u \in (0, 1); \\ 0, & u \notin (0, 1). \end{cases} \quad (10)$$

В рассматриваемом случае $v_1 = (v-k)/2$, $v_2 = k/2$. Из (10) получаем утверждение теоремы.

Лемма 2. Критерий проверки гипотезы о порядке регрессионной модели (1)–(3) против простой альтернативы имеет вид

$$T_{p,k}^{(n)} < 1 + \frac{k(F_\alpha(k, v-k) - 1)}{v}, \quad (11)$$

где $F_\alpha(k, v-k)$ — α -квантиль F -распределения с k и $v-k$ степенями свободы.

Доказательство. Из свойств β -распределения известно, что случайная величина

$$\frac{k_1 F(k_1, k_2)}{k_2 + k_1 F(k_1, k_2)}$$

имеет β -распределение $\beta(k_1/2, k_2/2)$. Согласно лемме 1

$$\zeta_{p,k}^{(n)} \sim \beta\left(\frac{v-k}{2}, \frac{k}{2}\right).$$

Из представления (9) получаем

$$T_{p,k}^{(n)} \sim r^2 \left[1 + \frac{k}{(v-k)F(v-k, k)} \right].$$

Учитывая (7) и свойства F -распределения, имеем

$$T_{p,k}^{(n)} \sim \frac{v-k}{v} \left[1 + \frac{kF(k, v-k)}{v-k} \right] = 1 + \frac{k(F(k, v-k) - 1)}{v}. \quad (12)$$

Представление (12) дает нам выражение для распределения t -статистики $T_{p,k}^{(n)}$ через F -распределение. При таких условиях можем построить простой в применении критерий проверки гипотезы о порядке регрессионной модели (1)–(3).

против простой альтернативы. Этот критерий задается неравенством (11). Следствие доказано.

4. Предельное распределение t -статистики. Рассмотрим бесконечную серию наблюдений

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), \dots \quad (13)$$

Построим критерий для проверки гипотезы о порядке модели (1)–(3), (13).

Определение 3. Назовем величину

$$\tau_{p,k}^{(n)} = \nu \left(1 - \zeta_{p,k}^{(n)} \right) \quad (14)$$

статистикой равновесия.

Теорема 3. Статистика равновесия (14) модели (1)–(3), (13) для произвольного p при $n \rightarrow \infty$ сходится по распределению к случайной величине χ_k^2 / σ^2 .

Доказательство. Используя лемму 1, имеем

$$1 - \zeta_{p,k}^{(n)} = 1 - \frac{\chi_{v-k}^2}{\chi_{v-k}^2 + \chi_k^2} = 1 - \frac{\chi_v^2 - \chi_k^2}{\chi_v^2} = \frac{\chi_k^2}{\chi_v^2}.$$

Принимая во внимание, что

$$\frac{\chi_v^2}{\nu} \xrightarrow{n \rightarrow \infty} \sigma^2,$$

получаем

$$\tau_{p,k}^{(n)} \xrightarrow{n \rightarrow \infty} \frac{\chi_k^2}{\sigma^2}.$$

Теорема доказана.

Следствие. Для достаточно больших n критерий проверки гипотезы о порядке регрессионной модели (1)–(3), (13) против простой альтернативы имеет вид

$$T_{p,k}^{(\infty)} < \frac{1}{1 - \chi_k^2(\alpha) / \nu \sigma^2},$$

где $\chi_k^2(\alpha)$ — α -квантиль χ^2 -распределения с k степенями свободы.

1. Себер Дж. Линейный регрессионный анализ. — М.: Мир, 1980. — 466 с.

2. Кендал М., Стьюарт А. Статистические выводы и связи: В 2-х т. — М.: Наука, 1973. — 899 с.

Получено 26.05.99