

Дельта-метод и его некоторые применения в статистике

Вирко Елизавета Петровна, гр. 422

Санкт-Петербургский государственный университет
Математико-механический факультет
Кафедра статистического моделирования

Научный руководитель: к.ф.-м.н., доцент В. В. Некруткин
Рецензент: исследователь Е. А. Советкин



Санкт-Петербург
2019г.

Одномерная модель линейной регрессии с фиксированными регрессорами

Модель:

при $i = 1, \dots, n$ $y_i = ax_i + b + \varepsilon_i$,

- x_i — известные постоянные числа (“регрессоры”),
- ε_i — независимые одинаково распределенные случайные величины, $E \varepsilon_i = 0, D \varepsilon_i = \sigma^2$,
- y_i — результаты наблюдений.

Применение МНК приводит к оценкам

$$\hat{a}_n = \frac{\text{COV}_n}{\bar{s}_{xn}^2}, \quad \hat{b}_n = \bar{y}_n - \hat{a}_n \bar{x}_n, \quad \hat{\sigma}_n^2 = \bar{s}_{yn}^2 - \hat{a}_n^2 (\bar{s}_{xn}^2)^2.$$

Некоторые свойства:

- $\hat{a}_n, \hat{b}_n, (n-2)\hat{\sigma}_n^2/n$ — несмещенные ([Ивченко, Медведев, 2010]).
- Если $\varepsilon_i \sim N(0, \sigma^2)$, то $\hat{\sigma}_n^2$ не зависит от $(\hat{a}_n, \hat{b}_n)^T$ и известны распределения $(\hat{a}_n, \hat{b}_n)^T, \hat{\sigma}_n^2$.

Одномерная модель линейной регрессии с условными математическими ожиданиями

Обычно имеем дело с повторной независимой выборкой.

Модель: пусть x, y — случайные величины с конечными вторыми моментами, $Dx = \sigma_x^2 > 0$, такие, что $E(y | x) = ax + b$. Тогда

$$a = \frac{\text{cov}(x, y)}{\sigma_x^2}, \quad b = E y - a E x, \quad \sigma^2 = \sigma_y^2 - \text{cov}^2(x, y) / \sigma_x^2.$$

$(x_1, y_1)^T, \dots, (x_n, y_n)^T$ — повторная независимая выборка из $(x, y)^T$.
Оценки параметров формально совпадают с полученными для модели с фиксированными регрессорами.

Некоторые свойства:

- \hat{a}_n, \hat{b}_n — несмещенные ([Демиденко, 1981]).
- Если $E(\varepsilon^2 | x) = \text{const}$, то $E \hat{\sigma}_n^2 = n\sigma^2 / (n - 2)$, известно асимптотическое распределение $(\hat{a}_n, \hat{b}_n)^T$.

На предыдущих двух моделях регрессии основан **стандартный критерий**.

Нулевая гипотеза: $H_0: b = 0$. Статистика:

$$\tau_{\text{std}} = \sqrt{n} \hat{b}_n / \hat{\sigma}_{\text{std}}, \quad \text{где } \hat{\sigma}_{\text{std}} = \frac{n}{n-2} \frac{\hat{\sigma}_n^2}{\overline{(x^2)}_n}.$$

Распределение статистики при $b = 0$:

- x_i фиксированы, $\varepsilon_i \sim N(0, \sigma^2)$. Тогда $\mathcal{L}(\tau_{\text{std}}) = t(n-2)$ ([Rao et al., 2007], [Ивченко, Медведев, 2010]).
- $E(y|x) = ax + b$ и $E((y - ax - b)^2 | x) = \text{const}$. Тогда $\mathcal{L}(\tau_{\text{std}}) \Rightarrow N(0, 1)$ [Демиденко, 1981].
- $(x, y)^T$ имеет невырожденное гауссовское распределение. Тогда $\mathcal{L}(\tau_{\text{std}}) = t(n-2)$.

Согласно этим распределениям строится критерий, называемый **стандартным**. Именно он реализован в пакетах STATISTICA, SPSS, STATGRAPHICS и stats для языка R.

Одномерная модель линейной регрессии с линейной аппроксимацией в $\mathbb{L}^2(dP)$. Постановка задачи

Проблемы:

- При каких еще условиях стандартный критерий применим (асимптотически точен)?
- Что делать, если стандартный критерий не применим?

Чтобы ответить на эти вопросы — более общая модель регрессии — линейная аппроксимация в $\mathbb{L}^2(dP)$. Снова повторная независимая выборка. Параметры и их МНК оценки те же, что в модели с УМО. Но **без условия** $E(y | x) = ax + b$.

Задача:

- Найти условия применимости стандартного критерия в случае линейной аппроксимации в $\mathbb{L}^2(dP)$.
- Построить критерии, применимые в более общей ситуации.

Обозначим $x^* = (x - Ex)/\sigma_x$, $\varepsilon = y - ax - b$.

Предложение

Пусть $(x, y)^T$ обладает конечными четвертыми моментами, распределение x непрерывно. Тогда $\mathcal{L}\left(\sqrt{n}((\hat{a}_n, \hat{b}_n, \hat{\sigma}_n^2)^T - (a, b, \sigma^2)^T)\right) \Rightarrow N(\mathbf{0}, \Sigma)$, где Σ имеет вид

$$\frac{1}{\sigma_x^2} \begin{pmatrix} D(\varepsilon x^*) & E(\varepsilon^2 x^*)\sigma_x - ExD(\varepsilon x^*) & E(\varepsilon^3 x^*)\sigma_x \\ E(\varepsilon^2 x^*)\sigma_x - ExD(\varepsilon x^*) & D(\varepsilon(\sigma_x - x^*Ex)) & \sigma_x^2 E\varepsilon^3 - E(\varepsilon^3 x^*)\sigma_x \\ E(\varepsilon^3 x^*)\sigma_x & \sigma_x^2 E\varepsilon^3 - E(\varepsilon^3 x^*)\sigma_x & \sigma_x^2 D(\varepsilon^2) \end{pmatrix}.$$

Следствие

Пусть $E(\varepsilon^2 | x) = \text{const}$. Тогда распределение $\sqrt{n}(\hat{a}_n, \hat{b}_n)^T$ совпадает с известным для модели с УМО: $\mathcal{L}\left(\sqrt{n}((\hat{a}_n, \hat{b}_n)^T - (a, b)^T)\right) \Rightarrow N(\mathbf{0}, \Sigma)$, где

$$\Sigma = \frac{\sigma^2}{\sigma_x^2} \begin{pmatrix} 1 & -Ex \\ -Ex & Ex^2 \end{pmatrix}.$$

Проверка гипотезы $b = 0$. Два критерия

Гипотеза $H_0: b = 0$. Обозначим $\sigma_b^2 = D((y - ax - b)(Ex^2 - xEx))/\sigma_x^4$.

Статистика критерия: $\tau_n = \sqrt{n}\hat{b}_n/\hat{\sigma}_{0n}$, где $\hat{\sigma}_{0n}^2 \xrightarrow{P_{H_0}} \sigma_0^2$ и $P(\hat{\sigma}_{0n} > 0) = 1$.

Критерий

Пусть $\alpha \in (0, 1)$. Критерий отвергает H_0 , если $|\tau_n| \geq C_{1-\alpha/2}$, где $C_{1-\alpha/2}$ — квантиль $N(0, 1)$ уровня $1 - \alpha/2$.

Два критерия: если (общий критерий)

$$\hat{\sigma}_{0n} = \hat{\sigma}_{\text{gen}}^2 = \frac{1}{n\bar{s}_{xn}^4} \sum_{i=1}^n \left((y_i - \hat{a}_n x_i - \hat{b}_n)(\bar{x}_n^2 - x_i \bar{x}_n) \right)^2,$$

или (модифицированный критерий)

$$\hat{\sigma}_{0n} = \hat{\sigma}_{\text{mod}}^2 = \frac{1}{n\bar{s}_{xn}^4} \sum_{i=1}^n \left((y_i - \hat{a}_n x_i)(\bar{x}_n^2 - x_i \bar{x}_n) \right)^2,$$

тогда $\mathcal{L}(\tau_n) \xrightarrow{H_0} N(0, 1)$.

Оба критерия асимптотически точны и состоятельны против альтернативы $H_{b^*}: b = b^* \neq 0$.

Положим

$$\sigma_{0,\text{std}}^2 = \frac{1}{\sigma^2 \sigma_x^2 \mathbb{E} x^2} D((y - ax - b)(\mathbb{E} x^2 - x \mathbb{E} x)).$$

Утверждение

Пусть x, y имеют конечные четвертые моменты, распределение x непрерывно. Стандартный критерий асимптотически точен тогда и только тогда, когда $\sigma_{0,\text{std}}^2 = 1$.

Следствие

- Стандартный критерий асимптотически точен тогда и только тогда, когда $(y - ax - b)^2$ и $(x - \mathbb{E} x)^2$ некоррелированы.*
- В частности, это выполняется, если $\mathbb{E}((y - ax - b)^2 | x) = \text{const}$ или если $\mathbb{E} x = 0$.*

Ошибки первого рода. Равномерное распределение в эллипсе

$$E(y|x) = ax + b, E((y - ax - b)^2|x) \neq \text{const}, b = 0.$$

а) $Ex = 0, \sigma_{0,\text{std}}^2 = 1$: стандартный критерий асимптотически точен,

б) $Ex = 2, \sigma_{0,\text{std}}^2 \approx 0.85^2$: стандартный критерий асимптотически консервативен.

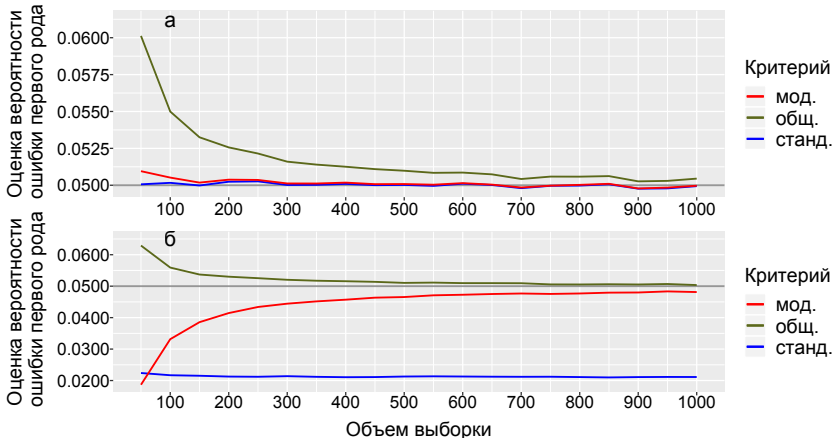


Рис.: Оценки вероятностей ошибок первого рода для равномерного распределения в эллипсе, $n = 50$ (50) 1000, $\alpha = 0.05$.

Поведение общего и модифицированного критериев при сдвигах

Пусть $E x = E y = 0$, тогда $b = 0$.

Для $c \neq 0$ положим $x' = x + c$, $y' = y + ac$. Тогда $b' = 0$.

Зафиксируем n и $(x_1, y_1)^T, \dots, (x_n, y_n)^T$, и рассмотрим статистики критериев, построенные по $(x'_1, y'_1)^T, \dots, (x'_n, y'_n)^T$ как функции от $c = E x'$.

Лемма

1. Статистика общего критерия при $|c| \rightarrow \infty$ имеет вид

$$-\sqrt{n} \operatorname{sgn}(c)(\hat{a}_n - a)/\hat{\theta} + O(1/|c|),$$

где

$$\hat{\theta} = \frac{1}{\sqrt{n\bar{s}_{xn}^2}} \left(\sum_{i=1}^n (y_i - \hat{a}_n x_i - \hat{b}_n)^2 (\bar{x}_n - x_i)^2 \right)^{1/2}.$$

2. Статистика модифицированного критерия при $|c| \rightarrow \infty$ имеет вид

$$-\sqrt{n\bar{s}_{xn}}/c + O(1/c^2).$$

Ошибки первого рода. Влияние сдвига

$$E(y|x) = ax + b, \quad E((y - ax - b)^2|x) \neq \text{const}, \quad b = 0.$$

$$\mathcal{L}(\tau_{\text{std}}) \Rightarrow N(0, \Delta^2), \quad \Delta^2 = 2/3 + O(1/c^2), \quad |c| \rightarrow \infty.$$



Рис.: Оценки вероятностей ошибок первого рода для равномерного распределения в эллипсе с различными сдвигами $c = E x = -7 (0.5) 7$ при фиксированном объеме выборки $n = 300$, $\alpha = 0.05$.

- $y_i = a_1x_{i1} + a_2x_{i2} + b + \varepsilon_i$,
 x_{i1}, x_{i2} фиксированы, ε_i — i.i.d., $E\varepsilon_i = 0$, $D\varepsilon_i = \sigma^2$.
 $(\hat{a}_{1n}, \hat{a}_{2n}, \hat{b}_n, \hat{\sigma}_n^2)^T$ — оценки МНК.
- x_1, x_2, y — случайные величины, $E(y | x_1, x_2) = a_1x_1 + a_2x_2 + b$,
 $(x_{11}, x_{12}, y_1)^T, \dots, (x_{n1}, x_{n2}, y_n)^T$ — повторная независимая
выборка, оценки МНК формально те же.

Свойства первых двух моделей переносятся с одномерного случая.

- Снова повторная независимая выборка, отказ от
 $E(y | x_1, x_2) = a_1x_1 + a_2x_2 + b$, параметры и оценки те же.

Проверка гипотезы $c_1 a_1 + c_2 a_2 = 0$. Стандартный критерий

Гипотеза $H_0: c_1 a_1 + c_2 a_2 = 0$, где c_1, c_2 — фиксированные числа, не равные 0 одновременно.

Для обычно используемых моделей строится **стандартный критерий**.

Он точен, если

- x_{i1}, x_{i2} фиксированы, $\varepsilon_i \sim N(0, \sigma^2)$ ([Ивченко, Медведев, 2010]),
- $(x_1, x_2, y)^T$ имеет невырожденное гауссовское распределение.

Асимптотически точен, если

- $E(y | x_1, x_2) = a_1 x_1 + a_2 x_2 + b$,
 $E((y - a_1 x_1 - a_2 x_2 - b)^2 | x_1, x_2) = \text{const}$ ([Демиденко, 1981]).

Аналогично одномерному случаю интересно провести анализ применимости стандартного критерия в случае линейной аппроксимации в $\mathbb{L}^2(dP)$, а также построить другие критерии.

Гипотеза $H_0: c_1 a_1 + c_2 a_2 = 0$.

В работе доказывается вариант ЦПТ для $\sqrt{n}((\hat{a}_{1n}, \hat{a}_{2n})^T - (a_1, a_2)^T)$. Он применяется к построению критериев для проверки этой гипотезы.

Статистика критерия: $\tau_n = \sqrt{n}(c_1 \hat{a}_{1n} + c_2 \hat{a}_{2n}) / \hat{\sigma}_{0n}$, где $\hat{\sigma}_{0n}^2$ — состоятельная оценка соответствующей асимптотической дисперсии (снова ЦПТ, но при выполнении H_0).

Как и ранее, получаем две оценки — *общую* (выборочную оценку) и *модифицированную* (состоятельную, только если верна H_0), а также два соответствующих критерия.

В обоих случаях $\mathcal{L}(\tau_n) \xrightarrow{H_0} N(0, 1)$, и **оба критерия асимптотически точны и состоятельны против альтернативы $H_c: c_1 a_1 + c_2 a_2 = c \neq 0$.**

Пусть x_1, x_2, y имеют конечные четвертые моменты, распределение $(x_1, x_2)^T$ непрерывно.

Явно выписывается необходимое и достаточное условие применимости стандартного критерия, оно состоит в равенстве единице определенной величины $\sigma_{0,\text{std}}^2$.

Оно выполняется, если $E((y - a_1x_1 - a_2x_2 - b)^2 | x_1, x_2) = \text{const}$.

Гипотеза $H_0: a_1 = a_2$. Ошибки первого рода

Равномерное распределение в шаре.

$$E(y | x_1, x_2) = a_1 x_1 + a_2 x_2 + b,$$

$$E((y - a_1 x_1 - a_2 x_2 - b)^2 | x_1, x_2) \neq \text{const}, a_1 = a_2,$$

$\sigma_{0,\text{std}}^2 = 5/7$: стандартный критерий не применим.

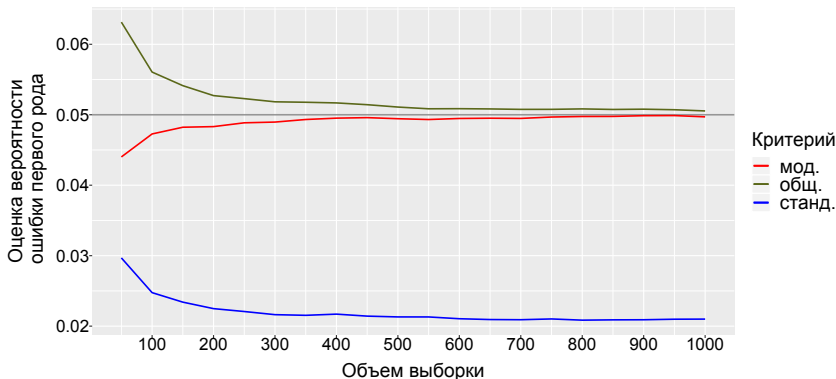


Рис.: Оценки вероятностей ошибок первого рода, $n = 50$ (50) 1000, $\alpha = 0.05$ для равномерного распределения в шаре.

- Найдены асимптотические распределения оценок метода наименьших квадратов для одномерной и двумерной регрессии.
- На основе этих распределений построено по два асимптотически точных и состоятельных критерия для проверки гипотез $H_0: b = 0$ (в одномерном случае), $H_0: c_1 a_1 + c_2 a_2 = 0$ (в двумерном случае), изучены некоторые свойства этих критериев.
- Полученные критерии были сравнены при помощи вычислительных экспериментов между собой и с критериями, обычно применяемыми при проверке таких гипотез.
- Найдены условия применимости стандартных критериев при отказе от предположений модели с УМО.

При нахождении асимптотического распределения в общем случае использовался “дельта-метод” [Shao Jun, 1999].

Теорема

Пусть $\boldsymbol{\eta}_n$ — последовательность d -мерных векторов, $U \subset \mathbb{R}^d$ — открытое, и $\mathbf{a} \in U$. Предположим, что отображение $f = (f_1, \dots, f_k)^T: U \rightarrow \mathbb{R}^k$ является дифференцируемым в точке \mathbf{a} . Обозначим $\Delta f_i \in \mathbb{R}^d$ — градиент функции f_i и положим $\Delta f = (\Delta f_1 : \dots : \Delta f_k)$.

Пусть для некоторой последовательности $c_n \rightarrow +\infty$, $n \rightarrow \infty$ имеет место сходимость

$$\mathcal{L}(c_n(\boldsymbol{\eta}_n - \mathbf{a})) \Rightarrow N(\mathbf{0}, \Sigma).$$

Пусть, к тому же, $P(\boldsymbol{\eta}_n \in U) = 1$. Тогда при $n \rightarrow \infty$

$$\mathcal{L}(c_n(f(\boldsymbol{\eta}_n) - f(\mathbf{a}))) \Rightarrow N(\mathbf{0}, \Sigma_f), \quad \text{где } \Sigma_f = \Delta f^T(\mathbf{a}) \Sigma \Delta f(\mathbf{a}).$$