

Прикладная статистика. Регрессия II

Леонид Иосипой

Программа «Математика для анализа данных»
Центр непрерывного образования, ВШЭ

11 марта 2021

Повторение

Регрессионный анализ решает задачу выявления искаженной случайным «шумом» зависимости некоторого показателя Y от измеряемых переменных X_1, \dots, X_k .

Повторение

Мы будем изучать **линейную регрессию**.

В линейной регрессии мы делаем предположение, что

$$y_i \approx \beta_0 + \beta_1 x_{i1} + \dots \beta_k x_{ik}, \quad i = 1, \dots, n,$$

где

- ▶ $y_i, x_{i1}, \dots, x_{ik}$ — отклик и значения k признаков для этого отклика (нам известные);
- ▶ $\beta_0, \beta_1, \dots, \beta_k$ — константы, которые не зависят от номера отклика (нам неизвестные).

Задача состоит в том, чтобы оценить $\beta_0, \beta_1, \dots, \beta_k$.

Повторение

Регрессионное равенство можно переписать в матричном виде как

$$y \approx X\beta,$$

где

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} x_{10} = 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} = 1 & x_{n1} & \dots & x_{nk} \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}.$$

Здесь мы добавили в матрицу X единичный столбец, чтобы больше не думать про коэффициент β_0 .

Повторение

Мы будем изучать свойства **метода наименьших квадратов** без использования каких-либо регуляризаторов:

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 = \|y - X\beta\|^2 \rightarrow \min_{\beta}$$

Точное решение $\hat{\beta}$ этой задачи известно и равно

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

Повторение

Чтобы исследовать качество решения метода наименьших квадратов, определим величину **TSS (Total Sum of Squares)** — разброс y относительно своего среднего:

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2.$$

Повторение

Оказывается, что (если в модель включен коэффициент β_0) TSS можно представить в виде суммы:

$$\text{TSS} = \text{RSS} + \text{ESS},$$

- ▶ **RSS (Residual Sum of Squares)** — это сумма квадратов отклонений предсказанных y от их истинных значений:

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

- ▶ **ESS (Explained Sum of Squares)** — это сумма квадратов отклонений среднего y от предсказанных y :

$$\text{ESS} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

Повторение

По величинам RSS и ESS можно составить меру R^2 , которая называется коэффициентом детерминации:

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

По сути, это доля объясненной дисперсии отклика во всей дисперсии отклика.

Повторение

Сделаем следующие предположения:

(П1) Истинная модель действительно является «зашумленной» линейной:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, \dots, n,$$

для некоторых (неизвестных) коэффициентов $\beta_0, \dots, \beta_k \in \mathbb{R}$ и некоторой случайной ошибки ε_i с $\mathbb{E}[\varepsilon_i] = 0$.

(П2) Наблюдения действительно случайны, то есть $(y_i, x_{i1}, \dots, x_{ik})$ для $i = 1, \dots, n$ образуют независимую выборку.

Повторение

(ПЗ) Матрица X является матрицей полного (столбцового) ранга:

$$\text{rank } X = k + 1.$$

То есть ни один из признаков не должен являться линейной комбинацией других. Поскольку среди столбцов есть константа, никакой из признаков в выборке не должен быть константой.

Повторение

Уже из этих трех предположений можно вывести, что оценки, получаемые методом наименьших квадратов, являются **несмещенными и состоятельными**:

$$\mathbb{E}[\hat{\beta}_j] = \beta_j \quad \text{и} \quad \hat{\beta}_j \xrightarrow{\mathbb{P}} \beta_j, \quad j = 0, \dots, k.$$

Повторение

Более того, предположим еще что:

(П4) Ошибки $\varepsilon_1, \dots, \varepsilon_n$ имеют одинаковую дисперсию, которая не зависит от значений признаков (гомоскедастичность ошибок):

$$\text{Var}(\varepsilon_i) = \sigma^2, \quad i = 1, \dots, n,$$

где $\sigma^2 > 0$ — неизвестный параметр.

Тогда можно показать, что дисперсия оценок, получаемых методом наименьших квадратов, является наименьшей в классе всех оценок, линейных по y (теорема Гаусса-Маркова).

То есть оценки метода наименьших квадратов являются в некотором смысле оптимальными.

Повторение

Рассмотрим еще одно предположение:

(П5) Ошибки $\varepsilon_1, \dots, \varepsilon_n$ имеют нормальное распределение

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2).$$

Если выполняются (П1)-(П5), то оценки метода наименьших квадратов совпадают с оценками максимального правдоподобия.

Это означает, что оценки метода наименьших квадратов обладают всеми свойствами, которыми обладают оценки максимального правдоподобия.

Повторение

Более того, при выполнении (П1)-(П5) мы можем посчитать распределения всех случайных объектов в модели:

$$y \sim \mathcal{N}(X\beta, \sigma^2 I_n),$$

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1}),$$

$$\hat{y} \sim \mathcal{N}(X\beta, \sigma^2 X(X^\top X)^{-1}X^\top).$$

Повторение

Эти факты позволяют нам построить следующие доверительные интервалы уровня доверия $1 - \alpha$, $\alpha \in (0, 1)$:

- ▶ для неизвестной дисперсии шума σ^2 :

$$\mathbb{P} \left(\frac{\text{RSS}}{c_{1-\alpha/2}} \leq \sigma^2 \leq \frac{\text{RSS}}{c_{\alpha/2}} \right) = 1 - \alpha,$$

где c_α — квантиль уровня α распределения χ_{n-k-1}^2 .

- ▶ для регрессионных коэффициентов β_0, \dots, β_k :

$$\mathbb{P} \left(\hat{\beta}_j - c_{1-\alpha/2} \hat{\sigma} \sqrt{(X^\top X)_{jj}^{-1}} \leq \beta_j \leq \hat{\beta}_j + c_{1-\alpha/2} \hat{\sigma} \sqrt{(X^\top X)_{jj}^{-1}} \right) = 1 - \alpha,$$

где $(X^\top X)_{jj}^{-1}$ — j, j элемент матрицы $(X^\top X)^{-1}$ и c_α — квантиль уровня α распределения T_{n-k-1} .

Повторение

Аналогично построению доверительных интервалов, можно проверить гипотезу о том, что признак j незначим, то есть что $\beta_j = 0$, $j = 0, \dots, k$.

Критерий Стьюдента

нулевая гипотеза: $H_0 : \beta_j = 0$

альтернатива: $H_1 : \beta_j \neq 0$ или $\beta_j > 0$ или $\beta_j < 0$

статистика:
$$T = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{(X^\top X)^{-1}_{jj}}}$$

нулевое распределение: $T \sim T_{n-k-1}$

Повторение

Можно также проверить гипотезу о том, что сразу несколько коэффициентов β_j равны 0.

Критерий Фишера

нулевая гипотеза: $H_0 : \beta_{j_1} = \dots = \beta_{j_m} = 0$
для некоторых $0 \leq j_1 < \dots < j_m \leq k$

альтернатива: $H_1 : \beta_{j_1}, \dots, \beta_{j_m} \neq 0$ одновременно

статистика: $F = \dots$

нулевое распределение: $F \sim F_{m, n-k-1}$ — распределение Фишера

Повторение

Обратите внимание, что доверительные интервалы и критерии строятся в предположениях (П1)-(П5).

Если ошибки имеют разную дисперсию и/или распределены не нормально, то доверительные интервалы будут неверными!

Повторение

Есть несколько типичных ошибок, которые следует иметь в виду, применяя регрессионный анализ. Сами по себе они достаточно очевидны. Тем не менее, о них часто забывают при работе с реальными данными и в результате приходят к неверным выводам.

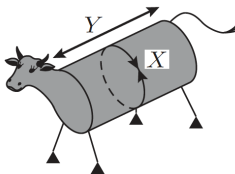
*Существуют три вида лжи: ложь, наглая ложь и статистика.
(Марк Твен)*

Повторение

Пример

Рассмотрим в качестве отклика Z вес коровы, а в качестве предикторов — окружность ее туловища X и расстояние от хвоста до холки Y . Сравнительному анализу были подвергнуты три регрессионные модели:

- (1) линейная: $Z = \theta_1 + \theta_2 X + \theta_3 Y$;
- (2) степенная: $Z = \theta_1' X^{\theta_2'} Y^{\theta_3'}$;
- (3) учитывающая содержательный смысл задачи $Z = \theta_0 X^2 Y$.



Повторение

Модель	По всем наблюдениям		По части наблюдений	
	$\hat{\theta}$	$\hat{\sigma}$	$\hat{\theta}_{\text{тяж}}$	$\hat{\sigma}_{\text{лег}}$
1	$\hat{\theta}_1 = -984,7$ $\hat{\theta}_2 = 4,73$ $\hat{\theta}_3 = 4,70$	25,9	$\hat{\theta}_1 = 453,2$ $\hat{\theta}_2 = 0,62$ $\hat{\theta}_3 = -0,22$	81
2	$\hat{\theta}'_1 = 0,0011$ $\hat{\theta}'_2 = 1,556$ $\hat{\theta}'_3 = 1,018$	24,5	$\hat{\theta}'_1 = 266,4$ $\hat{\theta}'_2 = 0,203$ $\hat{\theta}'_3 = -0,072$	79
3	$\hat{\theta}_0 = 1,13 \cdot 10^{-4}$	26,6	$\hat{\theta}_0 = 1,11 \cdot 10^{-4}$	28

Подробнее про пример можно посмотреть в книге Лагутина.

Спасибо за внимание!