

Лекция 7

Линейная регрессия: продолжение

Курс: Введение в DS на УБ и МиРА (весна, 2022)

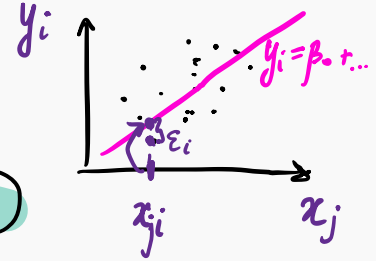
Преподаватель: Владимир Омелюсик

16 мая 2022 г.

- Параметрическое тестирование гипотез (статистики, p-value, доверительные интервалы).
- Непараметрическое тестирование гипотез (χ^2 -критерий согласия Пирсона).
- Линейная регрессия: начало.

Линейная регрессия: напоминание

- Верим, что данные пришли из модели



предп. \rightarrow $y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i} + (\epsilon_i)$

- y_i – целевая (зависимая переменная), $x_{j,i}$ – регрессоры (независимые переменные), β_j – коэффициенты, ϵ_i – случайная ошибка.
(параметры)

- Предполагаем, что $x_{j,i}$ и β_j – константы. – обяз-о

- Линейная регрессия линейна по β_j .

вн $y_i = \beta_0 + \beta_1 x_i^2 + \epsilon_i$ / ЛР

- Хотим по данным оценить $\hat{\beta}_0, \dots, \hat{\beta}_k$, чтобы делать предсказания как

$y_i = \beta_0 + \beta_1^2 x_i^2 + \epsilon_i$
– не ЛР

оценив. \rightarrow $\hat{y}_i = \hat{\beta}_0 + \dots + \hat{\beta}_k x_k$

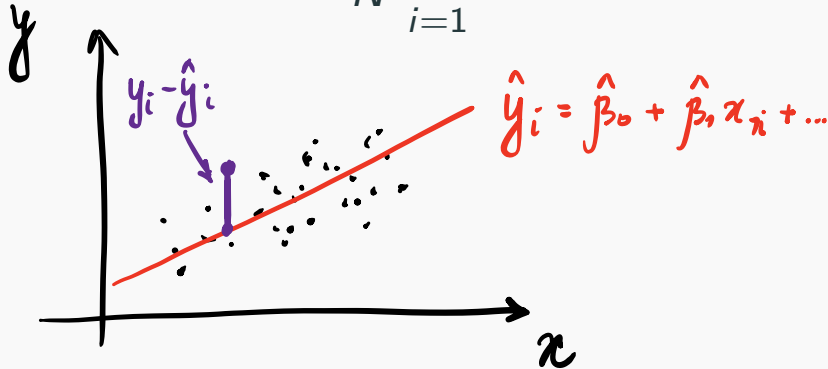
Метод наименьших квадратов

- Будем минимизировать усреднённую сумму квадратов отклонений истинного y_i от предсказанного:

$$MSE = \left[\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right] =$$

— среднекв. ошибки (mse, mean squared error)

$$= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \dots - \hat{\beta}_k x_{k,i})^2 \rightarrow \min_{\hat{\beta}_0, \dots, \hat{\beta}_k}$$



Пример: регрессия на константу

- Модель:

$$y_i = \beta_0 + \varepsilon_i$$

$$\hat{y} = \hat{\beta}_0$$

- Ищем оценку:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}_0)^2 \rightarrow \min_{\hat{\beta}_0}$$

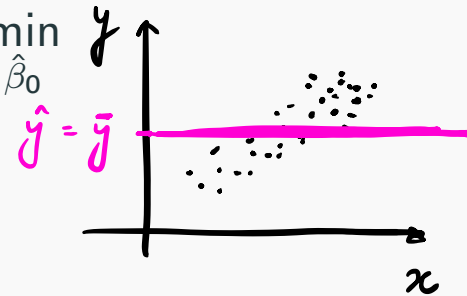
$$MSE'_{\hat{\beta}_0} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}_0) = 0$$

Другой

$$y_i = \beta_1 x_i + \varepsilon_i$$
$$\hat{y} = \hat{\beta}_1 x_i$$

$$\hat{\beta}_1 = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

$$\sum_{i=1}^N y_i = \sum_{i=1}^N \hat{\beta}_0 = N \cdot \hat{\beta}_0$$
$$\Rightarrow \hat{\beta}_0 = \frac{\sum_{i=1}^N y_i}{N} = \bar{y}$$



$$\hat{y} = \hat{\beta}_0 = \bar{y}$$

Пример: парная регрессия

- Модель:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \varepsilon_i$$

- Ищем оценку:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1,i})^2 \rightarrow \min_{\hat{\beta}_0, \hat{\beta}_1}$$

- Есть готовые формулы:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2},$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

$$\bar{y} = \frac{\sum_{i=1}^N y_i}{N}$$
$$\bar{x} = \frac{\sum_{i=1}^N x_i}{N}$$

Формула оценок для множественной регрессии

$$X = \begin{bmatrix} 1 & | & | & | & | \\ 1 & x_1 & x_2 & \dots & x_k \\ \vdots & | & | & & | \\ 1 & | & | & & | \end{bmatrix} \downarrow$$

$$N \times (k+1)$$

$$y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i} + \varepsilon_i$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{bmatrix}$$

- Запишем модель в матричном виде

$$\beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix} \quad y_i - \hat{y}_i \approx \hat{\varepsilon}_i$$

$$\begin{cases} y = X\beta + \varepsilon, \\ \hat{y} = X\hat{\beta} \end{cases}$$

$$N \times 1 \quad (k+1) \times 1$$

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

$$N \times 1 \quad N \times (k+1) \quad (k+1) \times 1 \quad N \times 1$$

$$\|y - \hat{y}\|_2^2 = \sum_i (y_i - \hat{y}_i)^2$$

- Запишем задачу в матричном виде

$$MSE = \frac{1}{N} \|y - X\hat{\beta}\|_2^2 = (y - X\hat{\beta})^T (y - X\hat{\beta}) \rightarrow \min_{\hat{\beta}}$$

- Тогда

v - вектор
 $p \times 1$

$$\|v\|_2^2 = \sqrt{\sum_{i=1}^p v_i^2}$$

- евклид. норма v

$$\hat{\beta}_1 = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

$$\hat{\beta} = \underbrace{(X^T X)^{-1} X^T y}_{(k+1) \times 1}$$

Интерпретация оценок коэффициентов

! При прочих равных

1. Lin-lin модель:

$$y_i = \beta_0 + \dots + \beta_k x_{k,i} + \dots + \varepsilon_i$$

$$\hat{y}_i = \hat{\beta}_0 + \dots + \hat{\beta}_k x_{ki} + \dots$$

- При увеличении $x_{k,i}$ на единицу, y_i увеличится на β_k

2. Log-lin модель:

$$\log y_i = \beta_0 + \dots + \beta_k x_{k,i} + \dots + \varepsilon_i$$

$$y_i = e^{\beta_0 + \dots + \beta_k x_{k,i} + \dots + \varepsilon_i}$$

- При увеличении $x_{k,i}$ на единицу, y_i увеличится в e^{β_k} раз.
- $e^{\beta_k} \approx (1 + \beta_k)$ для маленьких $\beta_k \Rightarrow$ при увеличении $x_{k,i}$ на единицу, y_i увеличится на β_k процентов.

Интерпретация оценок коэффициентов

3. Lin-log модель:

$$y_i = \beta_0 + \beta_1 x_i^2 + \varepsilon_i \longrightarrow y_i = \beta_0 + 2\beta_1 \log x_i + \varepsilon_i$$

// \tilde{x}_{ki}

$$y_i = \beta_0 + \dots + \beta_k (\log x_{k,i}) + \dots + \varepsilon_i$$

$y_i = \beta_0 + \tilde{\beta}_1 \log x_i + \varepsilon_i$
 $\log y_i = \log \beta_0 + \tilde{\beta}_1 \log x_i + \varepsilon_i$

- При увеличении $\log x_{k,i}$ на единицу, y_i увеличится на β_k .
- При увеличении $x_{k,i}$ на единицу, y_i увеличится на

$$\left\{ \beta_k [\log(x_{k,i} + 1) - \log(x_{k,i})] = \beta_k \log \frac{x_{k,i} + 1}{x_{k,i}} \right\} \tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 \log x_i + \varepsilon_i$$

- При увеличении $x_{k,i}$ на один процент, y_i увеличится на

$$\beta_k \log \frac{1.01 x_{k,i}}{x_{k,i}} = \beta_k \log 1.01$$

- При увеличении $x_{k,i}$ на $p\%$, y увеличится на

$$\beta_k \log([100 + p]/100).$$

4. Log-log модель:

$$\log y_i = \beta_0 + \dots + \beta_k \log x_{k,i} + \dots + \varepsilon_i$$

- При увеличении $x_{k,i}$ на $p\%$, y увеличится в $e^{\beta_k \log([100+p]/100)}$ раз \Rightarrow на $\beta_k \log([100+p]/100)$ процентов.



Качество подгонки

- Среднеквадратичная ошибка:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

	M1 (4 ряд)	M2 (8 ряд)
800		300

- Средняя абсолютная ошибка:

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

	7.5	5.4
M2 лучше		

- Коэффициент детерминации:

$R^2 \in [0, 1]$
! только если есть β_0 !

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

среднее по \hat{y}_i

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{\hat{y}})^2 / N-1}{\sum_{i=1}^N (y_i - \bar{y})^2 / N-1} = \frac{sVar(\hat{y})}{sVar(y)}$$

среднее по y_i

(доля вар. дан. целевой перемен, объясн. нашей моделью)

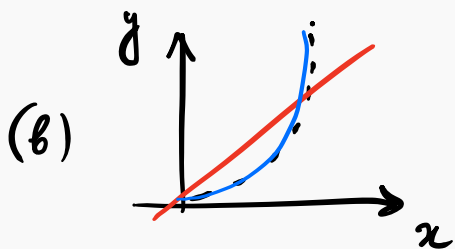
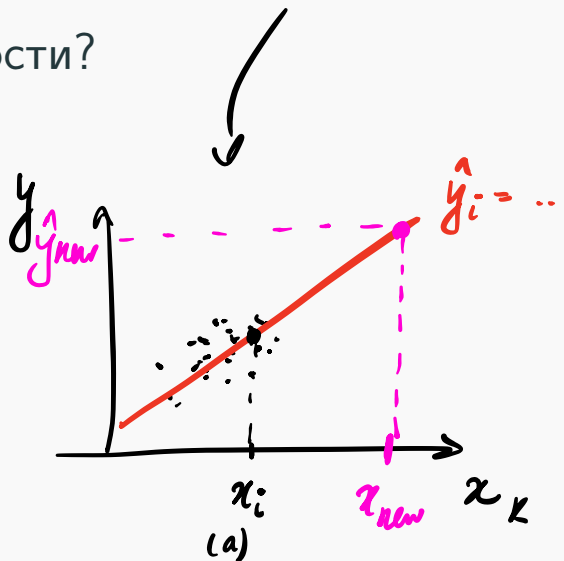
Miscellaneous

а) • Как делать предсказания на имеющейся выборке? $\hat{y}_i = \hat{\beta}_0 + \dots + \hat{\beta}_k x_{ik}$

(б) • Как делать предсказания на новой выборке?

(в) • Как добавлять нелинейности?

$$\hat{y}_{new} = \hat{\beta}_0 + \hat{\beta}_1 x_{new,1} + \dots + \hat{\beta}_k x_{new,k}$$



$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$y_i = \beta_0 + (\beta_1 x_i) + (\beta_2 x_i^2) + \varepsilon_i$$

Проверка гипотезы о значимости отдельного коэффициента

$$y_i = \beta_0 + \dots + \beta_k x_{ki} + \dots$$

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

\downarrow
 $H_0: \beta_k = 0$ не
 отверг.

Гипотеза о значим.
 β_j :

$$\begin{cases} H_0: \beta_j = 0, \\ H_1: \beta_j \neq 0 \end{cases}$$

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

- Предполагаем, что $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ и все независимы.
- Можно показать, что $\begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}$

$$\hat{\beta} \sim \mathcal{MN}(\beta, \text{Var}(\hat{\beta}))$$

- Также можно показать, что $\text{Var}(\hat{\beta}) = (X^T X)^{-1} \hat{\sigma}^2$ — точно

$$\hat{\sigma}^2 = \frac{N \times \text{MSE}}{N - (k + 1)}$$

$$\hat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} \hat{\text{Var}}(\hat{\beta}_0) & \overset{k \times k}{\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \dots} & \text{Cov}(\hat{\beta}_0, \hat{\beta}_k) \\ \vdots & \hat{\text{Var}}(\hat{\beta}_1) & \vdots \\ & & \hat{\text{Var}}(\hat{\beta}_k) \end{pmatrix}$$

Z-тест, t-тест и пример

- Используем Z-тест или t-тест:

$$\frac{\hat{\beta}_j - 0}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}} \sim t_{N-(k+1)} \rightarrow \mathcal{N}(0, 1) \quad t = \frac{4 - 0}{0.8} = \frac{4}{0.8} = 5$$

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$$

	const	x	z
coef ($\hat{\beta}_j$)	3	4	8
se	0.5	0.8	1

standard error

$$se = \sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}$$

std

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases}$$



Но отверг. \Rightarrow коэф. значим

$$\hat{y}_i = 3 + 4x_i + 8z_i$$

Можем пров.

$$1) \begin{cases} H_0: \beta_j = \beta_0 \\ H_1: \beta_j \neq \beta_0 \end{cases}$$

$$2) \begin{cases} H_0: \beta_j = 0 \\ H_1: \beta_j \neq 0 \end{cases}$$

Но не отверг. $\Rightarrow \beta_j(x_j)$ не значим

Проверка гипотезы о значимости регрессии в целом

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + \xi_i$$

- Формулировка гипотезы о значимости в целом:

*Гипотеза о
знач. в
целом*

$$\left\{ \begin{array}{l} H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0, \\ H_1 : \beta_1^2 + \dots + \beta_k^2 > 0 \end{array} \right\}$$

- Тестовая статистика сложная, имеет F -распределение.
- Легко проверять в софте.