

Лекция 6

Линейная регрессия

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

Преподаватель: Владимир Омелюсик

25 апреля 2022 г.

$X_1 \dots X_n \sim F_x(x)$ — модель до этого

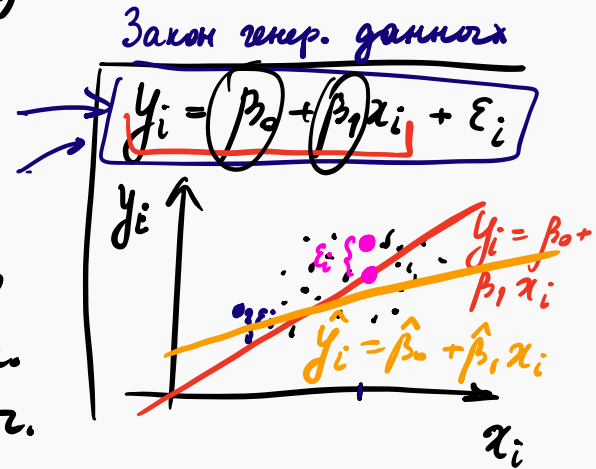
Постановка задачи

$$\rightarrow y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i$$

y_i — зависимая переменная (target, цель.)
 β_0 — константа (коэффициент)
 $x_{1i}, x_{2i}, \dots, x_{ki}$ — независимые переменные (признаки, регрессоры)
 ε_i — случайная ошибка

	x_1	x_2	...	x_k	y
1
2
3
⋮					

Пусть x_{ij} — константы
 β_j — всегда константы
 ε_i — сл. велич.
 y_i — сл. велич.



Хотим: по выборке получить $\hat{\beta}_0$ и $\hat{\beta}_1$ — оценки β_0 и β_1

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

сл. величина

По чему линейная линейная регрессия?

Пример:

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

Доход β_0 Стаж работы (y.e.) ε_i Сл. ошибка

2.9

Δαφν vs Σταμ ραβου
300 2.5

(y.e.) $\hat{\beta}_0 = 2.9$
 $\rightarrow \hat{\beta}_0, \hat{\beta}_1 = 4.3$

$$\hat{y}_i = 2.9 + 4.3 x_i$$

2000 год.

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_K x_{Ki} + \varepsilon_i$$

①
②
③
④

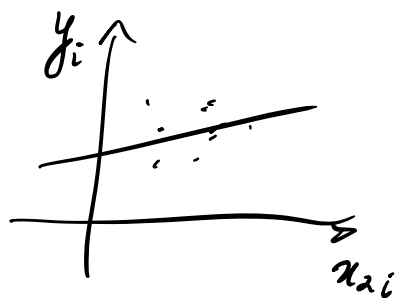
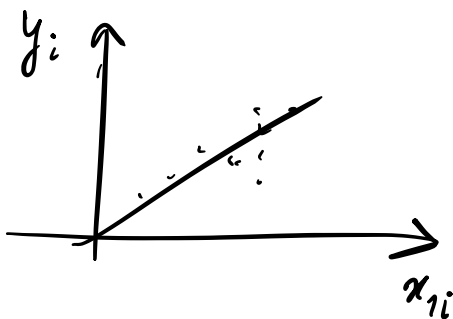
Лич. разр. лич-а по коэф-м

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i - \text{LP}$$

$$\ln y_i = \beta_0 + \beta_1 (x_i^2) + (\varepsilon_i^3) - \text{LP}$$

$$y_i = \beta_0^2 + (\ln \beta_1) x_i + \varepsilon_i - \text{не LP}$$

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \varepsilon_i$$



↙ LP

парная

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

↘
2 пар.

↘ LP

мног

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

> 2 пар.

	1	x_1	x_2	...	y
→	1	⋮			
→	1	⋮			
→	1	⋮			
	⋮				

Метод наименьших квадратов

$$\sum_i (y_i - \underbrace{\hat{y}_i})^2 \rightarrow \min_{\hat{\beta}_0, \dots, \hat{\beta}_k} \quad (\text{МНК})$$