



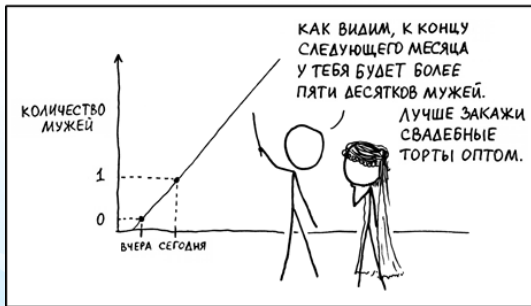
# Статистика ФИВТ ПМИ

## Прикладной поток

Лекция 12

# 7. Линейная регрессия

МОЁ ХОББИ: ЭКСТРАПОЛИРОВАТЬ





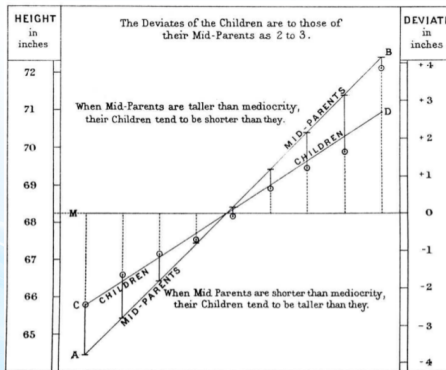
# Первое упоминание регрессии

Впервые регрессия упоминается в работе Гальтона

"Регрессия к середине в наследственности роста", 1885 г.

$x$  — рост родителей,  $y$  — рост детей

Установлена зависимость  $y - \bar{y} \approx \frac{2}{3}(x - \bar{x})$ , т.е. регрессия к середине.



# 7. Линейная регрессия

## 7.1. Постановка задачи линейной регрессии



# Пример

Пусть  $x$  — рост песика, а  $y$  — его вес.

Что мы знаем?

- ▶ чем крупнее песик, тем больший вес он имеет;
- ▶ песики одинакового роста могут иметь разный вес.

Выводы:

- ▶ для фиксированного роста песика  $x$  его вес  $y = f(x)$  является случайной величиной;
- ▶ в среднем вес  $f(x)$  возрастает при увеличении роста песика  $x$ .



# Пример

Простая зависимость:

$$y = \theta_0 + \theta_1 x + \varepsilon,$$

$x$  — рост песика,

$y$  — вес песика,

$\theta_0, \theta_1$  — неизвестные параметры,

$\varepsilon$  — случайная составляющая с нулевым средним.

Зависимость **линейна по параметрам**, линейна по аргументу.



# Пример

Более сложная зависимость:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2^2 + \varepsilon,$$

$x_1$  — рост песика,

$x_2$  — обхват туловища песика,

$y$  — вес песика,

$\theta_0, \theta_1, \theta_2, \theta_3$  — неизвестные параметры,

$\varepsilon$  — случайная составляющая с нулевым средним.

Зависимость **линейна по параметрам**, квадратична по аргументам.



# Модель линейной регрессии

Рассматриваем функциональную зависимость вида

$$y = y(x) = \theta_1 x_1 + \dots + \theta_d x_d$$

$x_1, \dots, x_d$  — признаки ,

$\theta = (\theta_1, \dots, \theta_d)^T$  — вектор параметров.

Для оценки  $\theta$  производится  $n$  испытаний вида

$$Y_i = \theta_1 x_{i1} + \dots + \theta_d x_{id} + \varepsilon_i, \quad i = 1, \dots, n,$$

$x_i = (x_{i1}, \dots, x_{id})$  — признаковые описания объекта  $i$   
(обычно неслучайные),

$\varepsilon_i$  — случайная ошибка измерений.



# Модель линейной регрессии

Введем обозначения

$$Y = \begin{pmatrix} Y_1 \\ \dots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ \dots & & \\ x_{n1} & \dots & x_{nd} \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \dots \\ \varepsilon_n \end{pmatrix}.$$

Матричная форма записи проведенных испытаний

$$Y = X\theta + \varepsilon.$$

$X \in \mathbb{R}^{n \times d}$  — регрессоры (или матрица плана эксперимента),

$Y \in \mathbb{R}^n$  — отклик.

Матричный вид зависимости:  $y(x) = x^T \theta$ .



## Замечание

Зависимость  $y = y(x)$  должна быть **линейна по параметрам**, но не обязана быть линейной по признакам.

Пусть  $z_1, \dots, z_k$  — набор "независимых" переменных.

Можно рассматривать модель

$$y(x) = \theta_1 x_1(z_1, \dots, z_k) + \dots + \theta_d x_d(z_1, \dots, z_k),$$

где  $x_j(z_1, \dots, z_k)$  — некоторые функции (м.б. нелинейные).

Примеры:

▶  $x(z_1, \dots, z_k) = 1;$

▶  $x(z_1, \dots, z_k) = z_1;$

▶  $x(z_1, \dots, z_k) = \ln z_1;$

▶  $x(z_1, \dots, z_k) = z_1^2 z_2.$

## Пример: Потребление мороженого

Предполагается линейная зависимость потребления мороженого в литрах на человека от среднесуточной температуры воздуха:  $ic = \theta_0 + \theta_1 t$ .

Проведена серия наблюдений

$$IC_i = \theta_0 + \theta_1 t_i + \varepsilon_i,$$

$t_i$  — среднесуточная температура воздуха,

$IC_i$  — потребление мороженого в литрах на чел.,

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



## Пример: Потребление мороженого

Наблюдения:  $IC_i = \theta_0 + \theta_1 t_i + \varepsilon_i$ .

В данном примере  $x_1(t) = 1, x_2(t) = t$ ,

$$X = \begin{pmatrix} 1 & t_1 \\ \dots & \\ 1 & t_n \end{pmatrix}, Y = \begin{pmatrix} IC_1 \\ \dots \\ IC_n \end{pmatrix}, \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}.$$

Пусть  $w = I\{\text{выходной день}\}$ , зависимость  $ic = \theta_0 + \theta_1 t + \theta_2 t^2 w$ .

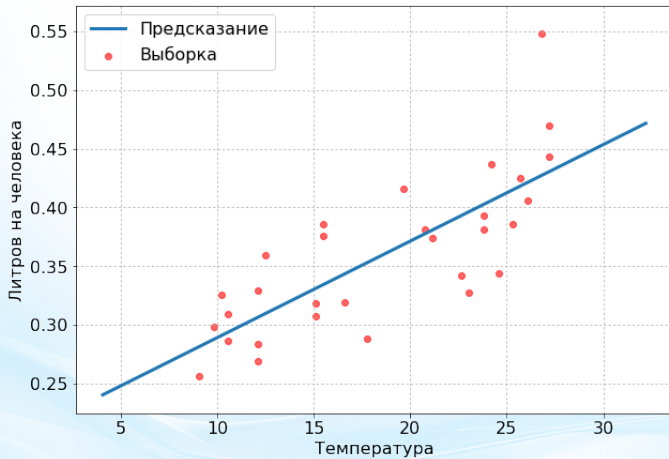
Наблюдения:  $IC_i = \theta_0 + \theta_1 t_i + \theta_2 t_i^2 w_i + \varepsilon_i$ .

В данном примере  $x_1(t, w) = 1, x_2(t, w) = t, x_3(t, w) = t^2 w$ ,

$$X = \begin{pmatrix} 1 & t_1 & t_1^2 w_1 \\ \dots & & \\ 1 & t_n & t_n^2 w_n \end{pmatrix}, Y = \begin{pmatrix} IC_1 \\ \dots \\ IC_n \end{pmatrix}, \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix}.$$



## Пример 3: Потребление мороженого





# Категориальные переменные

$x$  — id должности сотрудника (натуральное число),

$y$  — его зарплата.

Предположим, что должности занумерованы следующим образом:

- ▶  $x = 1$  — простой рабочий;
- ▶  $x = 2$  — сисадмин, присваивающий id;
- ▶  $x = 3$  — директор.

Сисадмин предложит рассмотреть модель  $y = \theta_0 + \theta_1 x$  :))

Если  $x \in \{1, \dots, k\}$ , то рассматриваются **dummy-переменные**:

$$x_j = I\{x = j\}, \quad j = 1, \dots, k - 1,$$

$$\text{модель } y = \theta_0 + \theta_1 x_1 + \dots + \theta_{k-1} x_{k-1}.$$

# 7. Линейная регрессия

## 7.2. Метод наименьших квадратов

# Материал на доске





# Реализация в sklearn

```
m = sklearn.linear_model.LinearRegression(fit_intercept=True)
```

Обучение модели:

```
m.fit(X, Y)
```

Вектор коэффициентов:

```
m.coef_
```

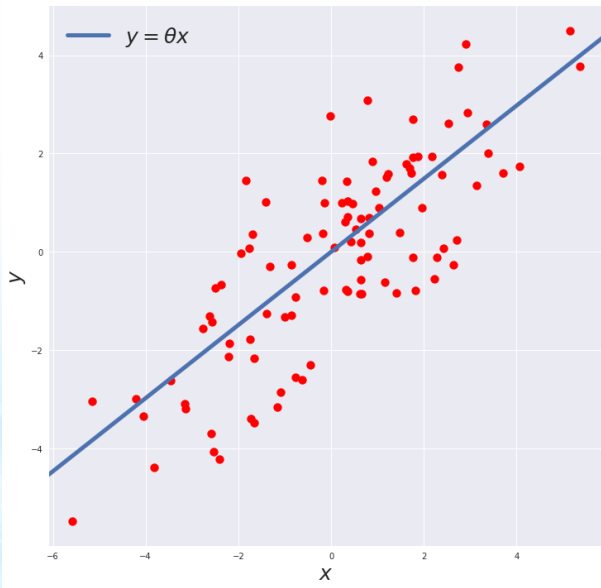
Свободный коэффициент:

```
m.intercept_
```

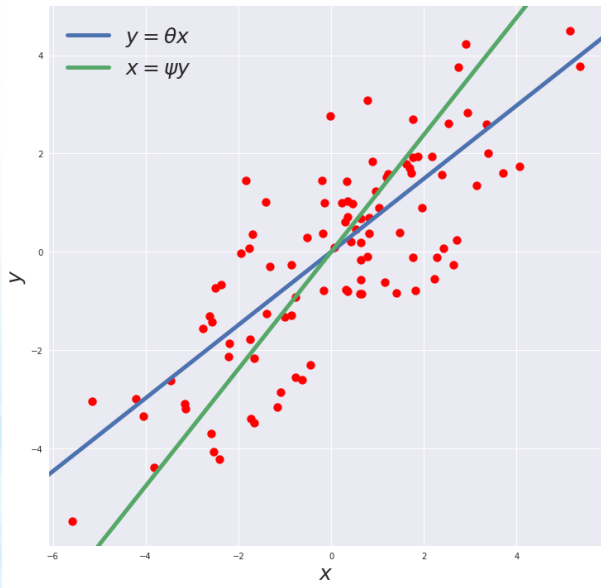
Предсказания:

```
m.predict(X)
```

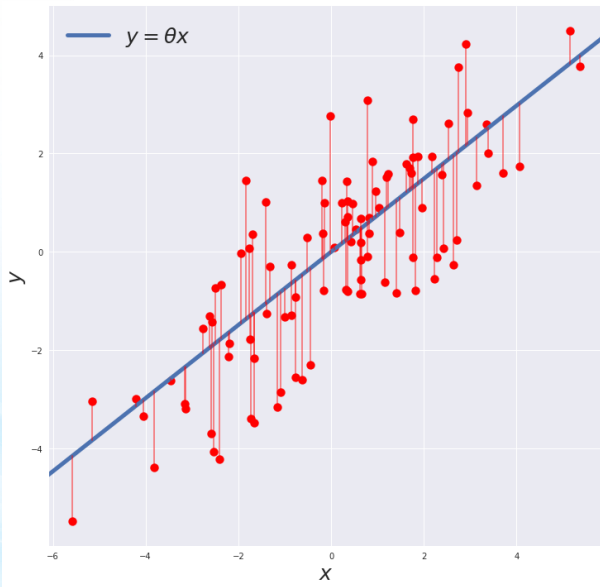
# Инверсия



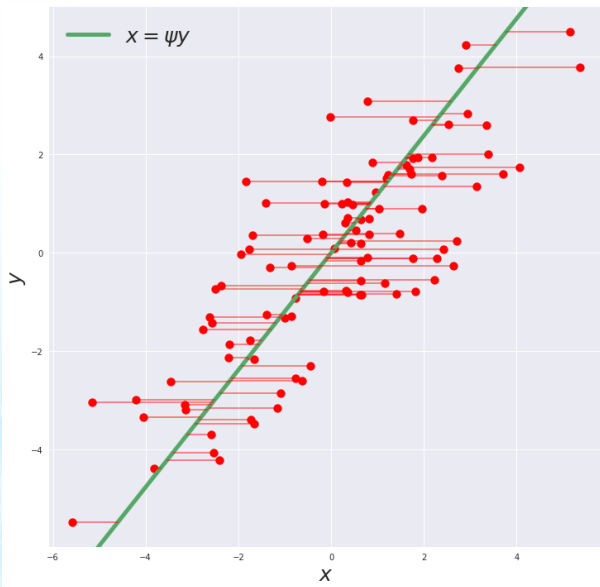
# Инверсия



# Инверсия



# Инверсия





**BCÈ!**