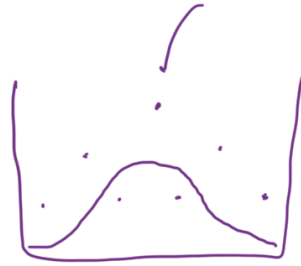
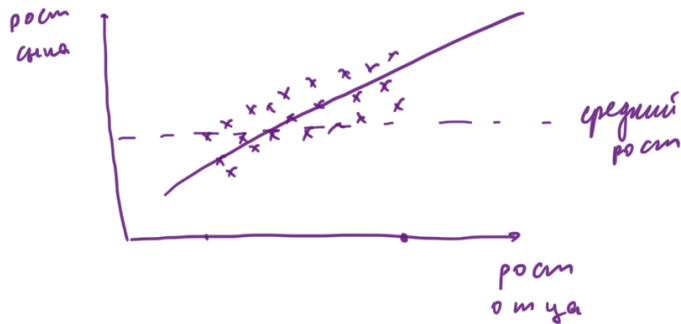


# линейная регрессия

$$Y = R$$



## ① Модель линейной регрессии

$$a(x) = w_0 + w_1 x_1 + \dots + w_d x_d = w_0 + \sum_{j=1}^d w_j x_j \quad \text{①}$$

свободный коэффициент (bias)
веса коэффициентов параметров

$(d+1)$  параметров

$$\textcircled{2} \quad \underline{w_0} + \langle w, x \rangle$$

предположение: 1-й параметр всегда = 1

$w_1 \cdot x_1$   
 $\quad \quad \quad \parallel$   
 $w_1 \cdot 1$

$$\textcircled{a(x) = \langle w, x \rangle}$$

## ② Область применения

$x$  - квартира в Москве

$y$  - рыночная стоимость

$$a(x) = w_0 + w_1 \cdot (\text{площадь}) + w_2 \cdot (\text{кол-во комнат}) +$$

$$+ w_3 \cdot (\text{рост до метро}) + \cancel{w_4 (\text{район})}$$

1) категориальные признаки

$x_1$  - район

$C = \{c_1, \dots, c_n\}$  - мн-во значений признака

$b_1(x), \dots, b_n(x)$  - кодовые признаки

$$b_j(x) = [x_1 = c_j]$$

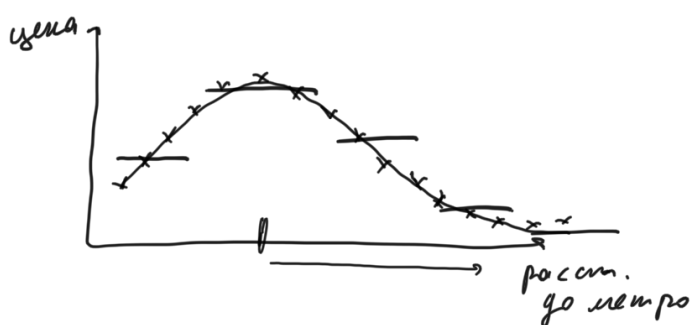
$$b_1(x) + \dots + b_n(x) = 1$$

one-hot  
кодирование

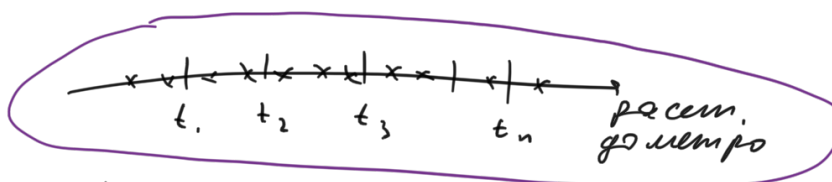
$$a(x) = w_0 + \underbrace{w_1 \cdot [x_1 = c_1]} + \underbrace{w_2 \cdot [x_1 = c_2]} + \dots +$$

$$+ w_n \cdot [x_1 = c_n] + \dots$$

2) дискретизация числовых признаков  
дискретизация



$$w_3 \cdot (\text{рост до метро})$$



$$t_0 = -\infty$$

$$t_{n+1} = +\infty$$

$$b_i(x) = [t_{i-1} \leq x_1 < t_i]$$

$$a(x) = \underbrace{w_1 [t_0 \leq x_1 < t_1]} + \dots + \underbrace{w_{n+1} [t_n \leq x_1 < t_{n+1}]}$$

Для линейных моделей нужно готовить такие данные, чтобы модель была осмысленной

③ Измерение ошибки в заданной регрессии

1)  $L(y, z) = (y - z)^2$  - квадратич. ф-ция потерь

$$MSE(a, X) = \frac{1}{l} \sum_{i=1}^l \underbrace{(\underbrace{a(x_i)}_{\text{предм}} - \underbrace{y_i}_{\text{предм}})^2}_{\text{квадратич. предм}}$$

$$RMSE(a, X) = \sqrt{MSE(a, X)}_{\text{предм}}$$

коэф. детерминации:

$$R^2(a, X) = 1 - \frac{\sum_{i=1}^l (a(x_i) - y_i)^2}{\sum_{i=1}^l (y_i - \bar{y})^2}$$

$a(x)$  идеальная  $\Rightarrow R^2 = 1$

$a(x)$  константа  $\Rightarrow R^2 = 0$   
( $a(x) = \bar{y}$ )

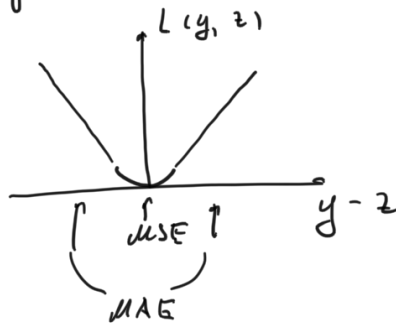
2)  $L(y, z) = |y - z|$

$$MAE(a, X) = \frac{1}{l} \sum_{i=1}^l |a(x_i) - y_i|$$

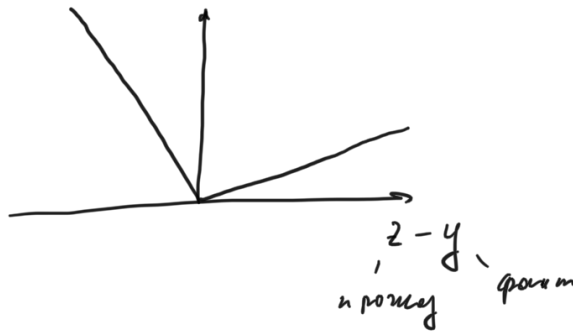
$y$	$a(x)$	$(y - a(x))^2$	$ y - a(x) $
1	2	1	1
1000	2	996004	998
1	1	0	0
1000	3	994009	997

$\Rightarrow$  MAE устойчиво к выбросам

3) Функция Хубера (Huber loss)



4) не самые практичные ф-ции потерь



$$L(y, z) = \begin{cases} \dots \end{cases}$$

5)  $L(y, z) = (\log(z+1) - \log(y+1))^2$  - MSLE

6)  $L(y, z) = \left| \frac{y-z}{y} \right|$

2      1  
1000    999

$$MAPE(a, X) = \frac{1}{l} \sum_{i=1}^l \frac{|y_i - a(x_i)|}{|y_i|}$$

④ Обучение

$$MSE: \frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 \rightarrow \min_w$$

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \vdots & & \vdots \\ x_{l1} & x_{l2} & \dots & x_{ld} \end{pmatrix} \quad \text{- матрица объектов - признаков}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_l \end{pmatrix} \quad w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} \quad \text{вектор весов}$$

$$D \quad | \quad y \in \mathbb{R}^n \quad | \quad w_d$$

$$MSE = \frac{1}{2} \| Xw - y \|_2^2 \rightarrow \min_w$$

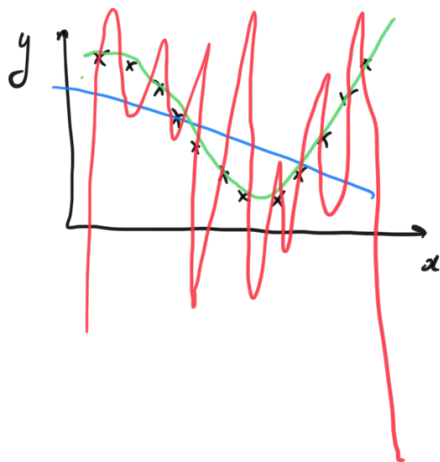
$$X = \begin{pmatrix} \langle w, x_1 \rangle \\ \langle w, x_2 \rangle \\ \vdots \\ \langle w, x_n \rangle \end{pmatrix}$$

$$\nabla_w MSE = 0$$

$$w = \underbrace{(X^T X)^{-1}}_{\substack{d \times d \quad l \times d \\ d \times d}} X^T y \quad (\text{только если } X \text{ не вырождена})$$

$O(d^3)$  - сложность обращения

⑤ Обобщающая способность модели



$$a(x) = w_0 + w_1 x$$

$$a(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$a(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_{n-1} x^{n-1}$$

Бли переобучение (overfitting) - ошибка на обучающих данных близка к нулю, тем на тест. выборке

Обозначение: 1) Обуч. выборка

$X$  - обучающая выборка

$X_t$  - тестовая выборка

... ..

$$U(\theta, X) \geq U(\theta, X_t)$$

2) иross-важация