

## Линейный классификатор

$X \subset \mathbb{R}^d$  - пр-во объектов

$Y = \{-1, +1\}$

$$X^l = (x_i, y_i)_{i=1}^l$$

$$a(\vec{x}, \vec{w}) = \text{sign}(\langle \vec{w}, \vec{x} \rangle + b)$$

$w \in \mathbb{R}^d$ ;  $b \in \mathbb{R}$  - сдвиг

Обучение линейного классификатора заключается в поиске вектора весов, на котором достигается минимум

$$w = \arg \min_{w \in \mathbb{R}^d} Q(w, X^l)$$

Наиболее удобными являются числовые функции потерь классиф.

$$Q(w, X^l) = \sum_{i=1}^l [y_i (\langle \vec{w}, \vec{x}_i \rangle + b) < 0] \rightarrow \min_w$$

$$\sum_{i=1}^l [y_i (\langle \vec{w}, \vec{x}_i \rangle + b) < 0] \leq \sum_{i=1}^l L(y_i (\langle \vec{w}, \vec{x}_i \rangle + b)) \rightarrow \min_w$$

$$L(M) = \log(1 + e^{-M}) - \text{логистическая функция потерь}$$

Замечание:

$$\frac{\partial f}{\partial \vec{v}} = \frac{d}{dt} f(x_{0,1} + t v_1, \dots, x_{0,d} + t v_d) \Big|_{t=0}$$

$$\frac{\partial f}{\partial \vec{v}} = \sum_{j=1}^d \frac{\partial f}{\partial x_j} \frac{d}{dt} (x_{0,j} + t v_j) = \sum_{j=1}^d \frac{\partial f}{\partial x_j} v_j = \langle \nabla f, \vec{v} \rangle$$

$$\langle \nabla f, \vec{v} \rangle = \|\nabla f\| \|\vec{v}\| \cos \varphi = \|\nabla f\| \cos \varphi$$



$$S(x_0) = \{x \in \mathbb{R}^d \mid f(x) = f(x_0)\}$$

Разложим  $f$  в ряд Тейлора на этой линии уровня

$$f(x_0 + \varepsilon) = f(x_0) + \langle \nabla f, \varepsilon \rangle + o(\|\varepsilon\|)$$

, где  $x_0 + \varepsilon \in S(x_0)$ . Поскольку  $f(x_0 + \varepsilon) = f(x_0)$

$$\langle \nabla f, \varepsilon \rangle = o(\|\varepsilon\|)$$

$$\left\langle \nabla f, \frac{\varepsilon}{\|\varepsilon\|} \right\rangle = o(1)$$

Устремим  $\varepsilon \rightarrow 0 \Rightarrow \frac{\varepsilon}{\|\varepsilon\|} \rightarrow$  касательная в точке  $x_0$ . В пределе получим, что ортогонален этой касательной.

Лемма 1

$$\nabla_x \langle a, x \rangle = a$$

$$\frac{\partial}{\partial x_j} \langle a, x \rangle = \frac{\partial}{\partial x_j} \sum_{k=1}^d a_k x_k = a_j$$

Лемма 2

$$\nabla_x \|x\|_2^2 = 2x$$

$$\frac{\partial}{\partial x_j} \sum_{k=1}^d x_k^2 = 2x_j$$

Лемма 3

$$\nabla_x \langle Ax, x \rangle = (A + A^T)x$$

$$= \nabla_x \langle A^T x, x \rangle = ?$$

$$\langle Ax, x \rangle = \sum_{j=1}^n (Ax)_j x_j = \sum_{j=1}^n \left( \sum_{k=1}^n a_{jk} x_k \right) x_j =$$

$$\sum_{j=1}^n \sum_{k=1}^n a_{jk} x_k x_j = \sum_{j=1}^n a_{jj} x_j^2 + \sum_{j \neq k} a_{jk} x_j x_k$$

$$\frac{\partial}{\partial x_i} \sum_{j=1}^n a_{jj} x_j^2 + \frac{\partial}{\partial x_i} \left( \sum_{j \neq i} a_{ij} x_i x_j + \sum_{j \neq i} a_{ji} x_i x_j \right) =$$

$$= 2a_{ii} x_i + \sum_{j \neq i} a_{ij} x_j + \sum_{j \neq i} a_{ji} x_j = \sum_{j \neq i} a_{ij} x_j + \sum_{j \neq i} a_{ji} x_j =$$

$$= (Ax)_i + (A^T x)_i = (A + A^T) x$$



$$\|Ax + b\|^2 = \langle Ax + b, Ax + b \rangle =$$

$$\langle Ax, Ax \rangle + 2\langle Ax, b \rangle + \langle b, b \rangle =$$

$$= \langle A^T A x, x \rangle + 2\langle x, A^T b \rangle + \langle b, b \rangle.$$

$$\nabla_x \|Ax + b\|_2^2 = \nabla_x \langle A^T A x, x \rangle + \nabla_x 2\langle x, A^T b \rangle +$$

$$+ \nabla_x \langle b, b \rangle = (A^T A + x^T A) x + 2A^T b =$$

$$= 2A^T A x + 2A^T b = 2A^T (Ax + b)$$

## Errors in Regression Tasks

### 1 Mean Squared Error

$$\frac{1}{n} \sum_{i=1}^n (a - y_i)^2$$

### 2 Root Mean Squared Error due to interpretation

$$\sqrt{\text{MSE}}$$

### 3 R MSE = 80

if  $a \in [0, 1]$ , RMSE - bad

elif  $a \in [10^4, 10^6]$ , RMSE - good

### 3. Determination coefficient + $R^2$

$$R^2 = \frac{\sum_{i=1}^n (a - y(x_i))^2}{\sum_{i=1}^n (y(x_i) - \bar{y})^2} \quad ?$$

Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |a_i - y(x_i)|$$

Let's consider the following task

$$\frac{1}{n} \sum_{i=1}^n (a - y_i)^2 \rightarrow \min$$

$$a_{mse}^* = \frac{1}{n} \sum_{i=1}^n y_i - \text{it is optimal}$$

$$\frac{1}{n} \sum_{i=1}^n |a - y_i| \rightarrow \min$$

$$a_{mas}^* = \text{median } \{y_i\}_{i=1}^n \rightarrow \text{or mode?}$$



## MEAN SQUARED LOGARITHMIC ERROR

$$L(y, a) = (\log(a+1) - \log(y+1))^2$$

1)  $a > 0; y > 0$

2) штраф за нарушение порядка, чем за значение

3) больше штраф за заниженные значения чем за завышенные.

## Mean Absolute Percentage Error

$$L(y, a) = \frac{|y - a|}{|y|}$$

Symmetric

$$L(y, a) = \frac{|y - a|}{\frac{(y + |a|)}{2}}$$

## Learning of linear regression

$$\text{Task: } \frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2 \rightarrow \min$$

$$\frac{1}{n} \|Xw - y\|^2 \rightarrow \min_w$$

$$w = (X^T X)^{-1} X^T y$$

Problems:  $A^{-1} = O(n^3)$

$X^T X$  — м.с. вычисления

$X$  - objects,  $X \in \mathbb{R}^n$

$Y$  - answers,  $Y \in \mathbb{R}$  or  $Y \in \mathbb{R}^n$

$y_i = y(x_i)$   $y: X \rightarrow Y$  unknown function;

$a(x) = f(x, \alpha)$  - model of dependency

$\alpha \in \mathbb{R}^p$  - vector of model params

Least Squares method

$$Q(\alpha, X^l) = \sum_i w_i (f(x_i, \alpha) - y_i)^2 \rightarrow \min$$

$w_i$  - weight of importance of  $i$  object

$Q(\alpha^*, X^l)$  - residual sum of squares



Task of Least Square method is connected to Principle of maximum likelihood

The model of data with non-correlated gaussian noise

$$y(X_i) = f(X_i, \alpha) + \varepsilon_i \quad (1)$$

, where  $\varepsilon_i \sim N(0, \sigma_i^2)$   $i = 1, \dots, l$

Method of Max. Likelihood

$$(2) L(\varepsilon_1, \dots, \varepsilon_l | \alpha) = \frac{1}{\prod_{i=1}^l \sqrt{2\pi} \sigma_i} \exp\left(-\frac{1}{2\sigma_i^2} \varepsilon_i^2\right) \Rightarrow \max_{\alpha}$$

$$-\ln L(\varepsilon_1, \dots, \varepsilon_l | \alpha) = \text{const}(\alpha) + \frac{1}{2} \sum_{i=1}^l \frac{1}{\sigma_i^2} (f(X_i, \alpha) - y_i)^2 \rightarrow \min_{\alpha}; \quad (3)$$

The application of least square method suppose that data has ① normal distribution and ② independent ③ dispersions are the same

Weights of objects  $w_i = \sigma_i^{-2}$



## Nadaraya - Watson

The approx. with the help of const  
 $a(x) = a$  at the vicinity  $x \in X$

$$Q(a; X^l) = \sum_{i=1}^l w_i(x) (a - y_i)^2 \quad (4)$$

$w_i(x) = K\left(\frac{p(x, x_0)}{h}\right)$  - beca observed  $x_0$   
relatively to  $x_0$

$K$  - kernel, non-increasing, limited, smooth

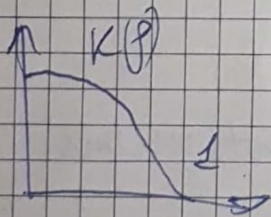
$h$  - width of window

$$a_h(x, X^l) = \frac{\sum_{i=1}^l y_i \widehat{w_i(x)} = 1}{\sum_{i=1}^l \underbrace{w_i(x)}_{=1}} = \frac{\sum_{i=1}^l y_i K\left(\frac{p(x_i, x)}{h}\right)}{\sum_{i=1}^l K\left(\frac{p(x, x_i)}{h}\right)}$$

Question 1: What's the solution of  
LSM if the solution is constant.

Average.

Example of kernel





Let the following conditions be satisfied:

1) Samples  $X^l = (x_i, y_i)_{i=1}^l$  is simple & belongs to distribution  $p(x, y)$ ;

2) Kernel  $K(r)$  is limited;

$$\int_0^\infty K(r) dr < \infty, \quad \lim_{r \rightarrow \infty} r K(r) = 0;$$

3) Dependence  $E(y|x)$  has no vertical asymptotes;

$$E(y^2|x) = \int y^2 p(y|x) dy < \infty \text{ for any } x \in X$$

4) The sequence  $h_l$  decreases, but not too fast.

$$\lim_{l \rightarrow \infty} h_l = 0, \quad \lim_{l \rightarrow \infty} l h_l = \infty$$

Then there is convergence in probability

$$a_{h_l}(x, X^l) \xrightarrow{P} E(y|x) \text{ for any}$$

point  $x \in X$ , where

$$E(y|x), p(x), D(x) \text{ is continuous}$$

$$D(x) > 0$$



Kernel  $K(r)$

- influence on the smoothness of  $a_h(x)$  significantly affects smoothness
- slightly affects the quality of approx

Width of  $h$

- significantly affects the quality of approx

- with an uneven grid  $\{x_i\}$  there is variable width of window  
 $h(x) = \rho(x, x^{(k+2)}), x^{(k+1)}$  -  $k$ -neighbour of  $x_i$

- Optimization of window width by sliding control

$$LOO(h, X^L) = \sum_{i=1}^L (a_h(x_i; X^L \setminus \{x_i\}) - y_i)^2 \rightarrow \min_h$$

How to handle outliers

Main idea:

Bigger the error  $\varepsilon_i = |a_h(x_i; X^L \setminus \{x_i\}) - y_i|$

the more a precedent  $(x_i, y_i)$  is an outlier and less its weight should be.

Input:  $X^L$  - learning set      Algo

Output: coefficients  $\gamma_i$   $i=1, \dots, l$

1 Init  $\gamma_i = 1$   $i=1 \dots l$

2 repeat

3    for all  $i=1 \dots l$

4    calc estimates of sliding control  
$$a_i = a_h(x_i; X^L \setminus \{x_i\}) = \frac{\sum_{j=1, j \neq i}^l y_j \gamma_j K(\frac{f}{h})}{\sum_{j=1, j \neq i}^l \gamma_j K(\frac{f}{h})}$$

5    for all objects  $i=1, \dots, l$

$$\gamma_i = \bar{\kappa}(|a_i - y_i|);$$

6     $\|\gamma_i\| > \tilde{\gamma}$