

$$Q(\beta) = \underbrace{(y - X\beta)^T}_{1 \times m} \underbrace{(y - X\beta)}_{m \times 1} + \lambda \beta^T \beta$$

① это скаляр (но вектору)

② взять дифференциал

$$\underline{d(AB) = (dA)B + A d(B)}$$

$$d(Q(\beta)) = - \underbrace{d\beta^T}_{1 \times m} \underbrace{X^T (y - X\beta)}_{m \times 1} -$$

$$- \underbrace{(y - X\beta)^T}_{1 \times m} \underbrace{X d\beta}_{m \times 1} + 2 \beta^T d\beta = 0 \quad |_{\beta = \hat{\beta}}$$

$$+ 2 (y - X\hat{\beta})^T X d\hat{\beta} = 2 \hat{\beta}^T d\hat{\beta}$$

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

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

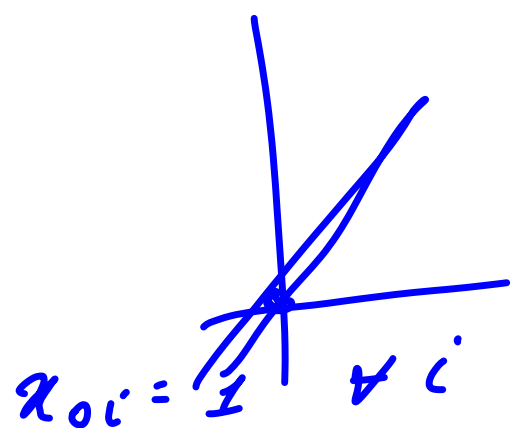
$$\hat{\beta}_{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

$$\underbrace{(x_1 \ x_2 \ x_3)}_{x^T} \underbrace{\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}}_x = \boxed{x^T x}$$

$$\{(x_i, y_i)\}_{i=1}^N$$

$y_i \in \mathbb{R}$ - наблюдения

$x_i \in \mathbb{R}^{d+1}$ - признаки
 ← можно задать y_i



$$a(x_i) = \underline{\langle w, x_i \rangle} = w_0 + w_1 x_{1i} + \dots + w_d x_{di}$$

Обучение ☺

① аналитически

$$(X^T X)^{-1}$$

② градиентный спуск

$$\nabla f(x_1, \dots, x_n) = (f'_{x_1}, \dots, f'_{x_n})$$

• инициал. веса

$$w^{(0)} = (w_0^{(0)}, w_1^{(0)}, \dots, w_d^{(0)})$$

• обновление

$$w^{(k)} = w^{(k-1)} - \eta_k \nabla Q(w^{(k-1)})$$

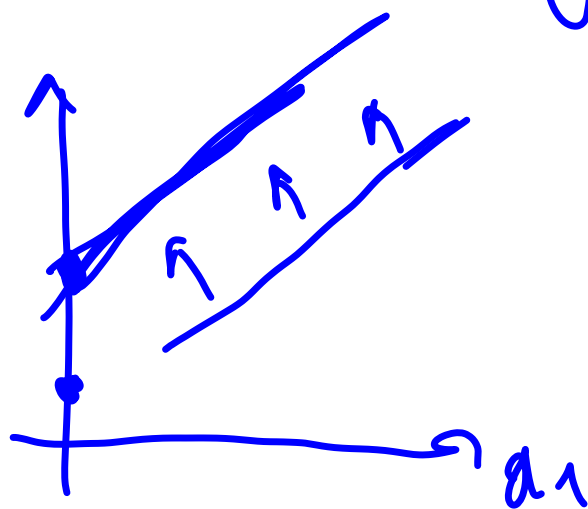
↑ learning rate

$$Q(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \underbrace{\sum_{j=1}^d w_j x_{ij}}_{\text{dot product}})^2 + \lambda \underbrace{\langle w, w \rangle}_{\sum_j x_{ij} w_j}$$

$$\textcircled{1} \nabla_w Q(w) = -\frac{1}{n} \sum 2x_i (y_i - \hat{y}_i) + 2\lambda w$$

$y_1 \dots y_n$

$$\eta > 0$$



$$w_0^{(1)} = w_0^{(0)} - \eta \left[-\frac{2}{n} \sum_{i=1}^n x_{i0} (y_i - \hat{y}_i) \right]$$

⋮

$$w_d^{(1)} = w_d^{(0)} - \eta \left[-\frac{2}{n} \sum_{i=1}^n x_{id} (y_i - \hat{y}_i) + 2\lambda w_d^{(0)} \right]$$

$$w_d^{(k)} = w_d^{(k-1)} - \eta \left[\dots \right]$$

$$w_d^{(1)} = \left(w_d^{(0)} - 2\eta \lambda w_d^{(0)} \right) - \eta \text{ grad } OLS$$

grad MLR (OLS)

$$w_d^{(1)} = w_d^{(0)} - \eta \text{ grad } OLS$$

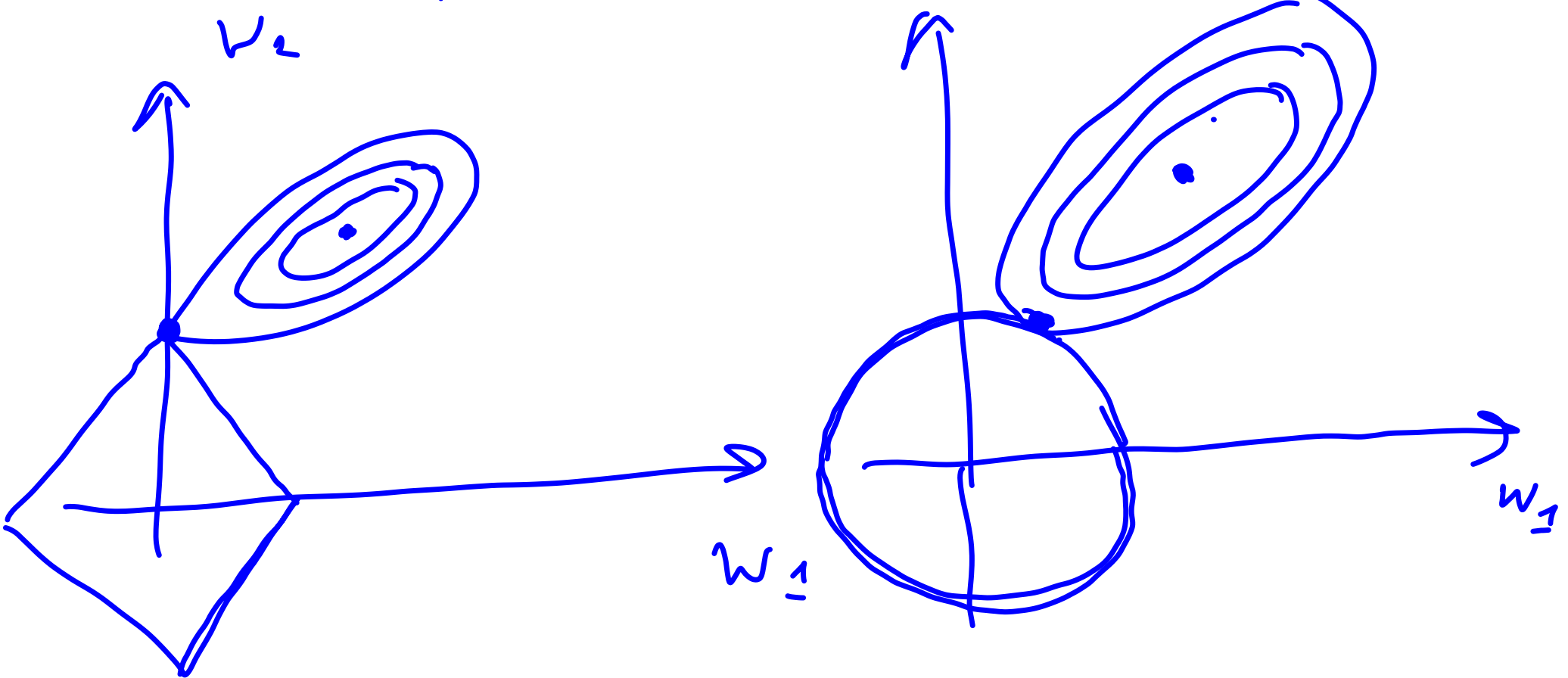
$$w_d^{(1)} = w_d^{(0)} \left(1 - 2\eta \lambda \right) - \eta \text{ grad } OLS$$

$\lambda > 0$
 $\lambda = 0$
 $\lambda < 1$

Ch-b0: shrinking

Lasso

Ridge

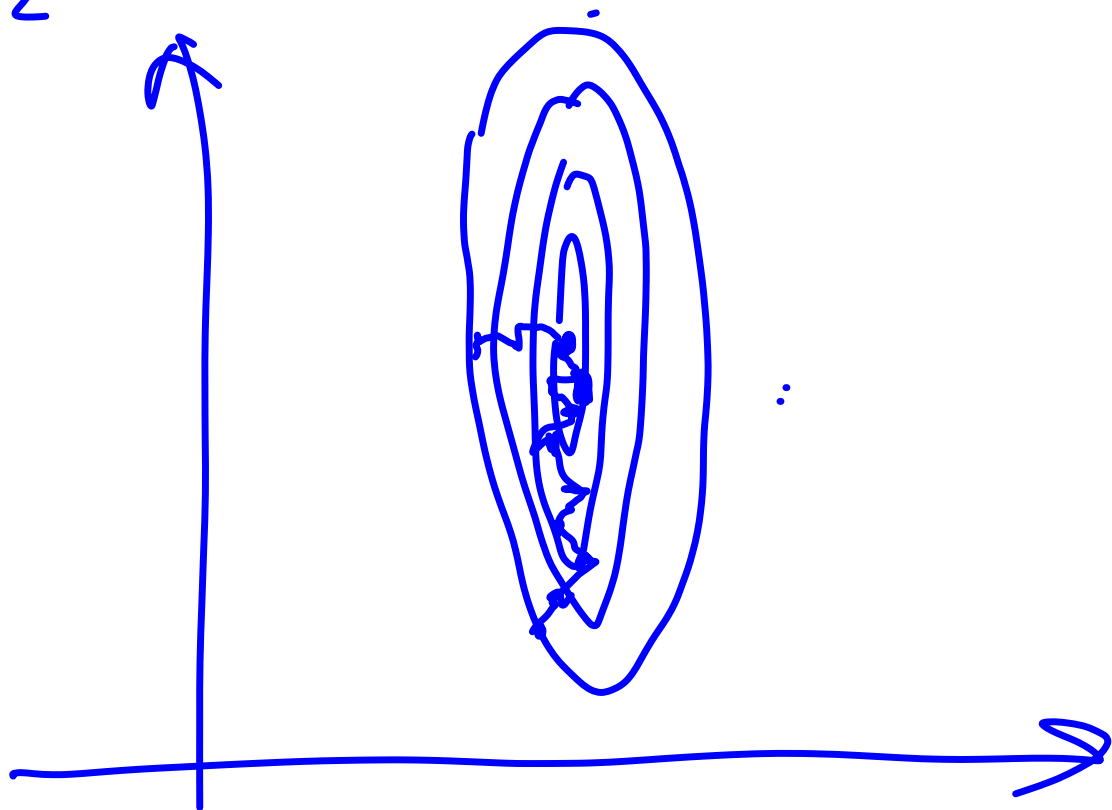


Корреляция

① $x_1 \in [20, 200] \rightarrow W_1$

$x_2 \in [1, 5] \rightarrow W_2$

W_2

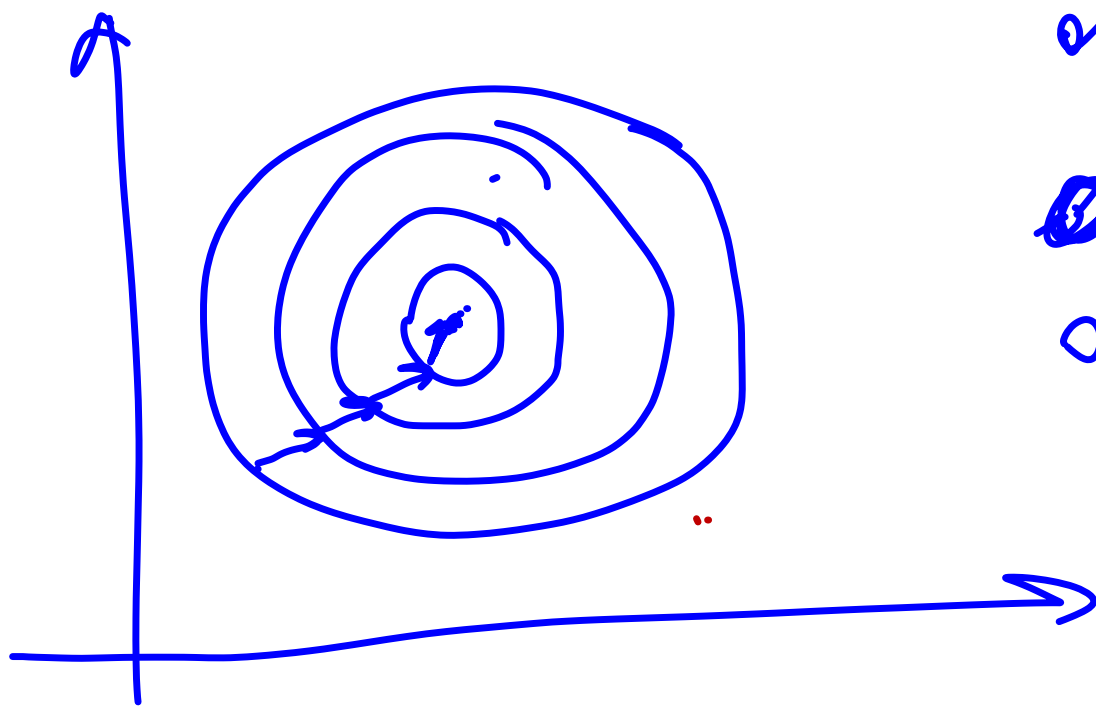


W_1

корреляция x_1, x_2

②

W_2



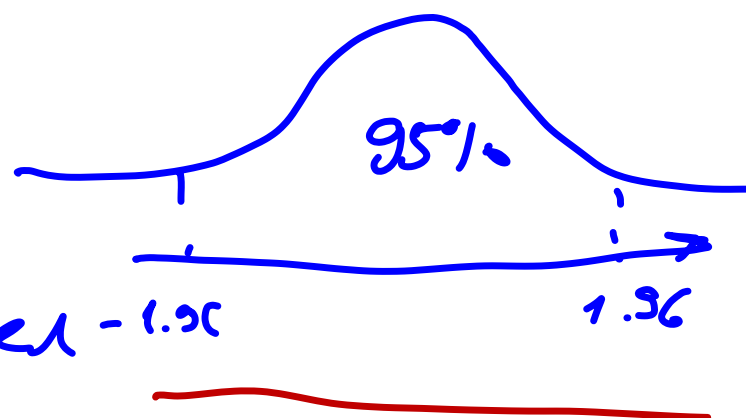
W_1

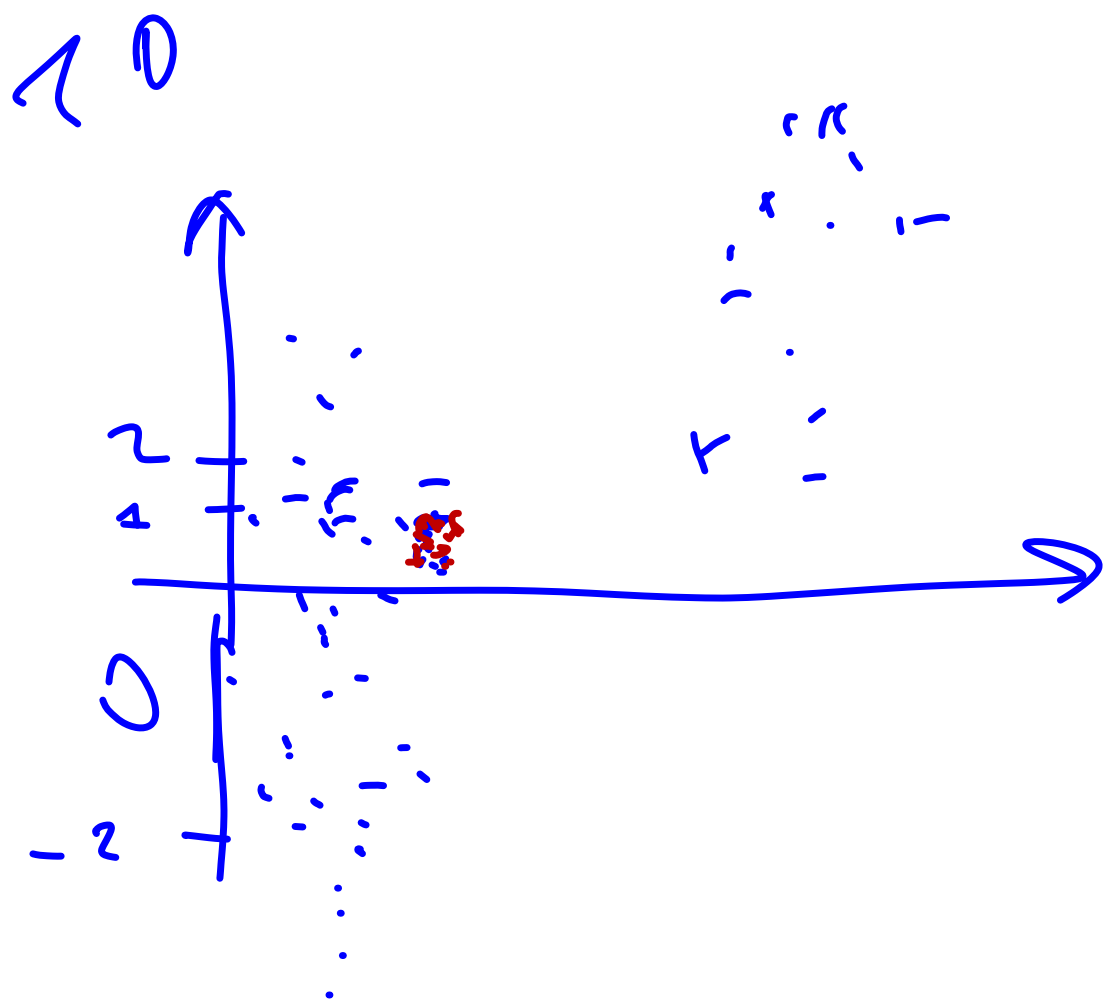
$0 \leq x_1 \leq 1$
 $0 \leq x_2 \leq 1$

③

$$z' = \frac{x - \bar{x}}{\hat{\sigma}}$$

-cf. статистика - 1.96





② min - max

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \rightarrow [0; 1]$$

$$x' (b - a) + a \rightarrow [a, b]$$

Кросс - валидация.

- ① подборовое с перебором
- ② настраивать гиперпараметры

a) k-fold CV

d) LOO CV

Leave-one-out

Набор данных $D = \{D_1, D_2 \dots D_k\}$
k частей равного
размера

$$\bigcup_i^k D_i = D$$

$$D_i \cap D_j = \emptyset \quad \forall i, j = \{1 \dots k\} \quad i \neq j$$

Алгоритм k-fold

① Случайно разбить выборку
на k частей. каждая - "11"

② для k от 1 до k повторить

$$\text{Test set} \leftarrow D_k$$

$$\text{Train} \leftarrow D \setminus D_k$$

Обучить модель на "Train"

Посчитать ошибку для D_k и

назвать Error_k

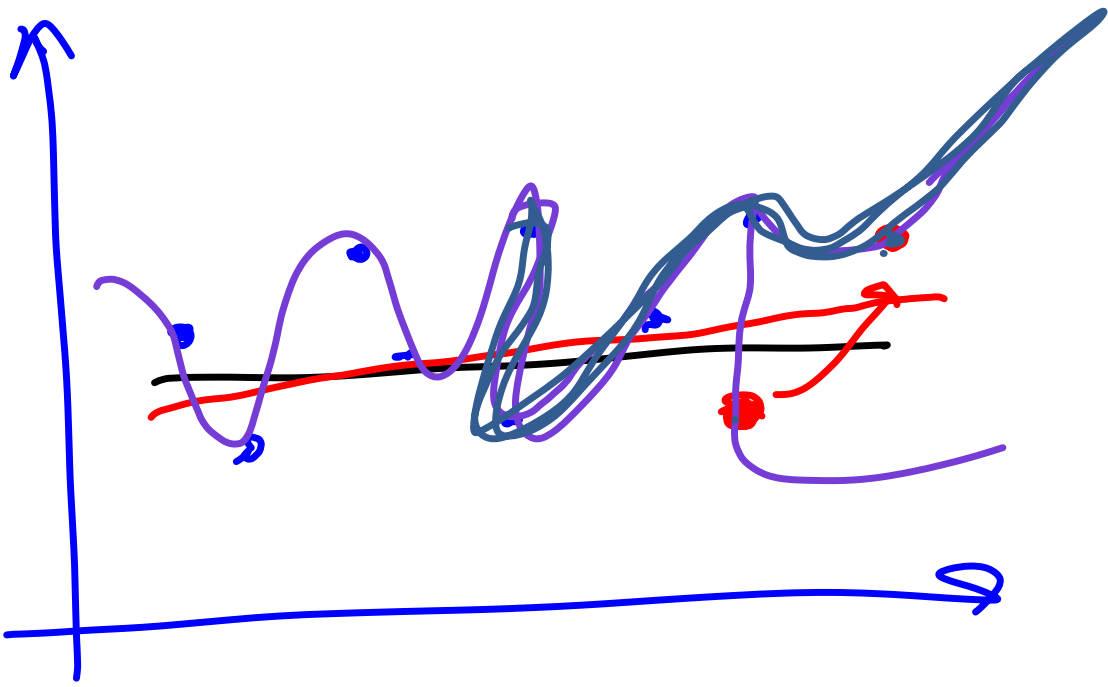
$$\text{Error} = \frac{1}{k} \sum_{i=1}^k \text{Error}_i$$

↑
исходящая ошибка
(cross-val.)
можно

LOOCV (когда надо проверить)
 $|D| = N$
 вынуждены
 проверять
 все

① for k to N
 Test $\leftarrow k$ -й элемент
 Проверка
 Train $\leftarrow D \setminus k$ -й эл.
 Об. обуч. на Train.
 — — — — —

$$\text{Error} = \frac{1}{N} \sum_{k=1}^N \text{Error}_k$$



настройка гиперпар.
 настройка гиперпар.

X_{train} , X_{test} , $X_{\text{validation}}$
 60 20 20