

# Введение в АД

## Лекция 6

### Кросс-валидация, градиентный бустинг



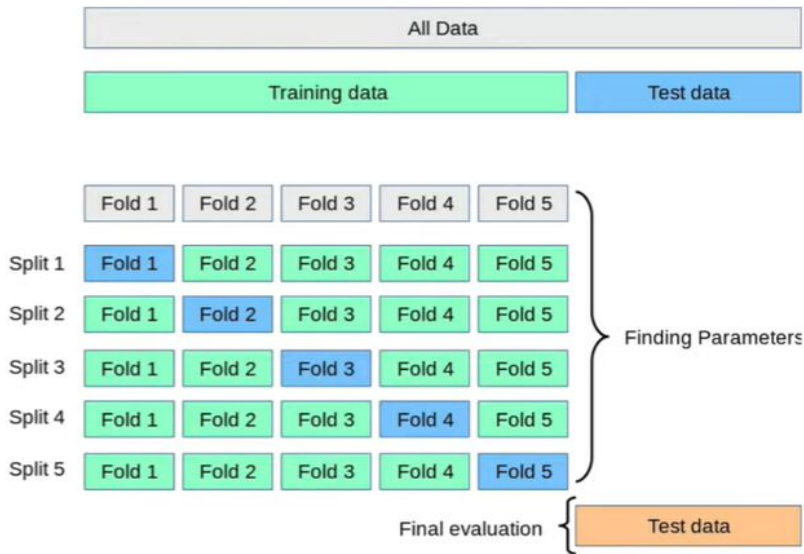
# Cross-validation

По какому принципу нужно делить выборку на train/test?

Что делать, если данных мало, а хочется обучиться хорошо?

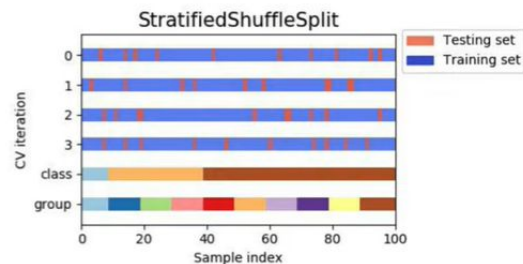
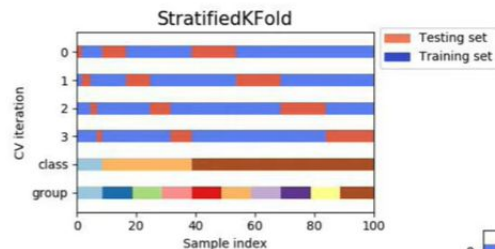
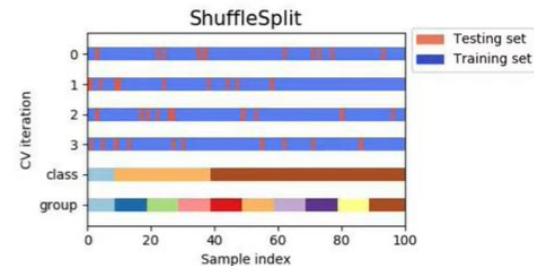
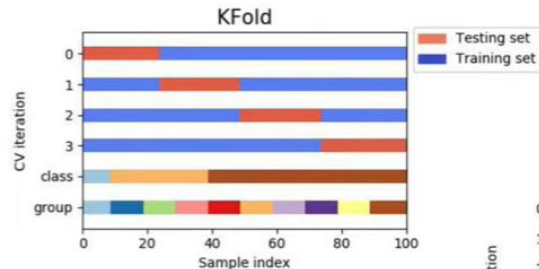
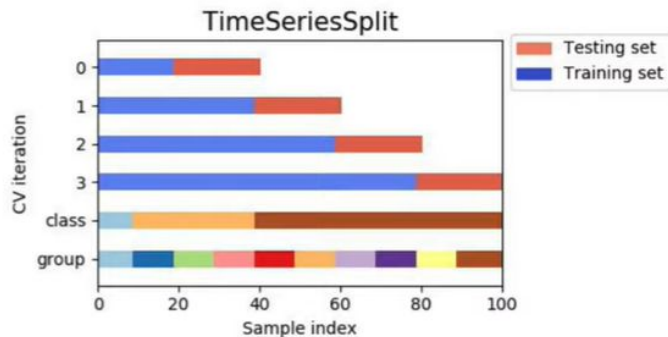
Стратегии валидации:

- hold-out
- k-fold
- leave-one-out (LOO)
- etc.



# Cross-validation

**Важно:** кросс-валидация не дает оценку качества обучения, она дает оценку качества подобранной **модели**

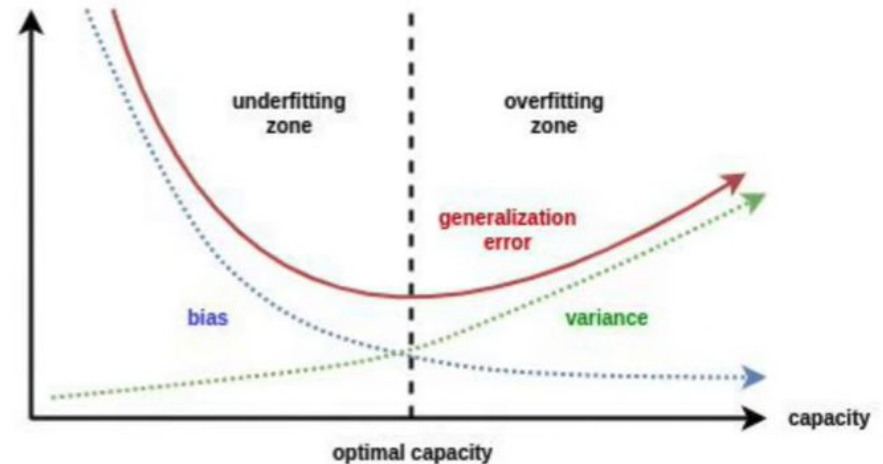




# Bias-variance tradeoff

При обобщении результатов работы модели на генеральную совокупность, ошибку модели можно разложить на три части:

- **bias**, неправильные или недостаточные предположения модели о данных
- **variance**, излишняя чувствительность к флуктуациям в данных, “обучение на шум”
- **шум** в исходных данных, также называемый **irreducible error**





# Bias-variance decomposition

The dataset  $X = (x_i, y_i)_{i=1}^{\ell}$  with  $y_i \in \mathbb{R}$  for regression problem.

Denote loss function  $L(y, a) = (y - a(x))^2$ .

The corresponding risk estimation is

$$R(a) = \mathbb{E}_{x,y} \left[ (y - a(x))^2 \right] = \int_{\mathbb{X}} \int_{\mathbb{Y}} p(x, y) (y - a(x))^2 dx dy.$$

Let's show that  $a_*(x) = \mathbb{E}[y | x] = \int_{\mathbb{Y}} yp(y | x)dy = \arg \min_a R(a)$ .

$$\begin{aligned} L(y, a(x)) &= (y - a(x))^2 = (y - \mathbb{E}(y | x) + \mathbb{E}(y | x) - a(x))^2 = \\ &= (y - \mathbb{E}(y | x))^2 + 2(y - \mathbb{E}(y | x))(\mathbb{E}(y | x) - a(x)) + (\mathbb{E}(y | x) - a(x))^2. \end{aligned}$$

Let's return to the risk estimation:

$$\begin{aligned} R(a) &= \mathbb{E}_{x,y} L(y, a(x)) = \\ &= \mathbb{E}_{x,y} (y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y | x) - a(x))^2 + \\ &\quad + 2\mathbb{E}_{x,y} (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)). \end{aligned}$$

$$\begin{aligned} \mathbb{E}_x \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)) \mid x \right] &= \\ = \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) \mid x \right] \right) &= \\ = \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) (\mathbb{E}(y | x) - \mathbb{E}(y | x)) \right) &= \\ = 0 \end{aligned}$$

So the risk takes form:

$$R(a) = \mathbb{E}_{x,y} (y - \mathbb{E}(y|x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y|x) - a(x))^2.$$

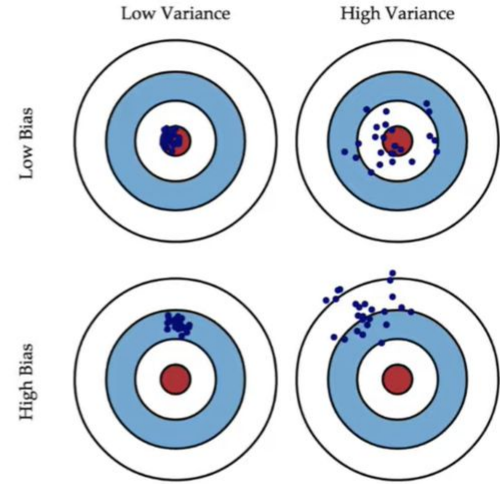
Does not depend on  $a(x)$

The minimum is reached when  $a(x) = \mathbb{E}(y|x)$ .

So the optimal regression model with square loss is

$$a_*(x) = \mathbb{E}(y|x) = \int_{\mathbb{Y}} yp(y|x)dy.$$

$$L(\mu) = \underbrace{\mathbb{E}_{x,y} [(y - \mathbb{E}[y|x])^2]}_{\text{noise}} + \underbrace{\mathbb{E}_x [(\mathbb{E}_X [\mu(X)] - \mathbb{E}[y|x])^2]}_{\text{bias}} + \underbrace{\mathbb{E}_x [\mathbb{E}_X [(\mu(X) - \mathbb{E}_X [\mu(X)])^2]]}_{\text{variance}}.$$



# Gradient boosting

Бустинг – подход, при которой для обучения следующих моделей мы используем данные о предыдущих. Реализован в CatBoost, XGBoost

Плюсы:

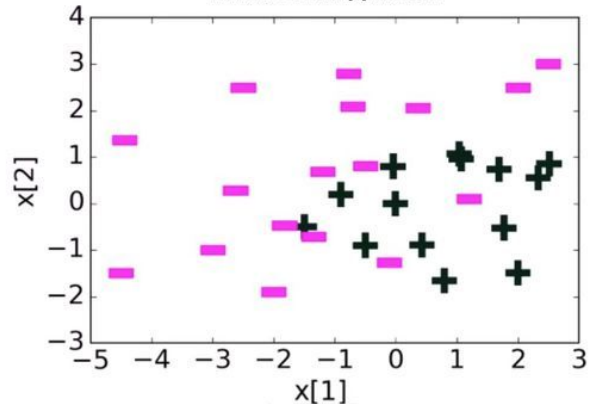
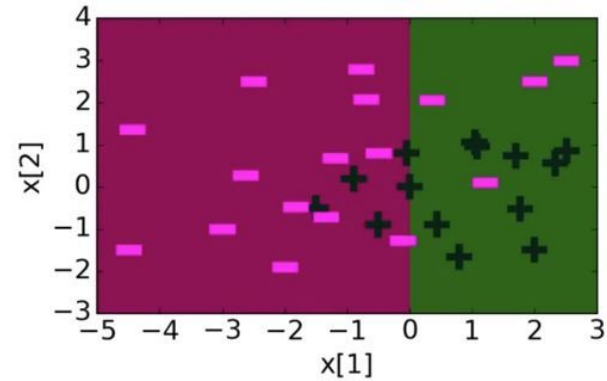
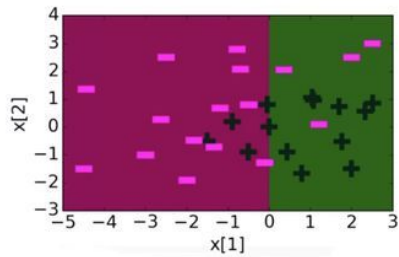
- легкое построение модели
- быстрое обучение, если за базовую модель берем быстрые алгоритмы
- можно хорошо идентифицировать выбросы

Минусы:

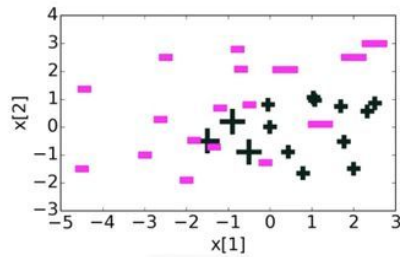
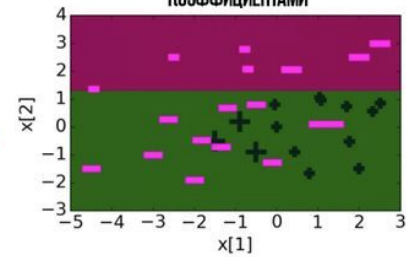
- теряется интерпретируемость
- слабая параллелизация



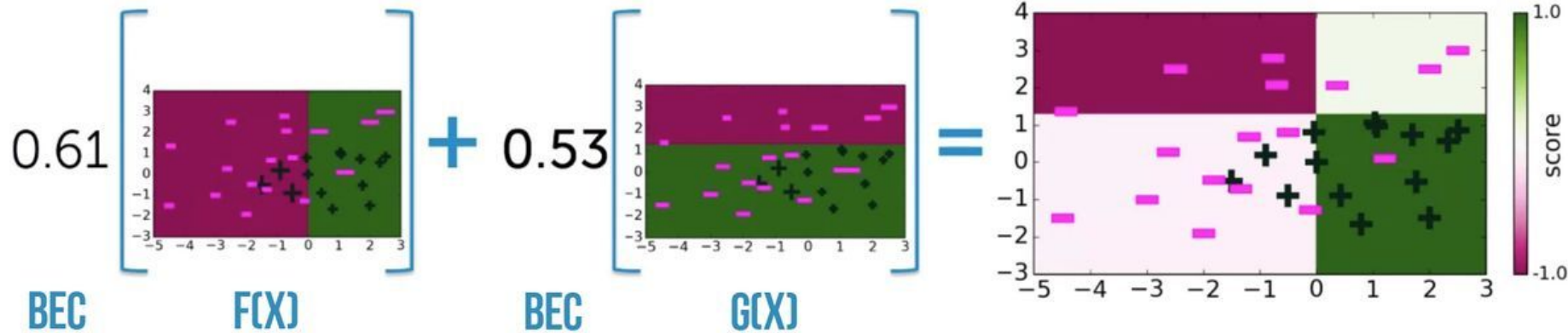
НАЧАЛЬНЫЕ ДАННЫЕ

РЕШЕНИЕ  $f(x)$ РЕШЕНИЕ  $f(x)$ 

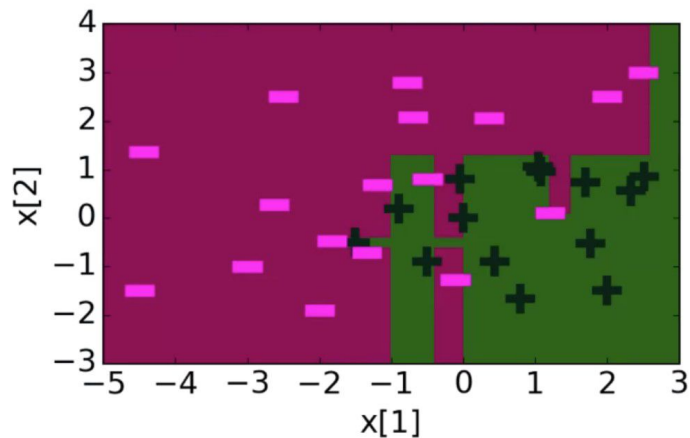
ПЕРЕСЧЕТ ВЕСОВЫХ КОЭФФИЦИЕНТОВ

РЕШЕНИЕ  $g(x)$  С НОВЫМИ ВЕСОВЫМИ КОЭФФИЦИЕНТАМИ





после 30 итераций:



# Откуда берутся коэффициенты?

Optimal model:

$$\hat{f}(x) = \arg \min_{f(x)} L(y, f(x)) = \arg \min_{f(x)} \mathbb{E}_{x,y}[L(y, f(x))]$$

Let it be from parametric family:

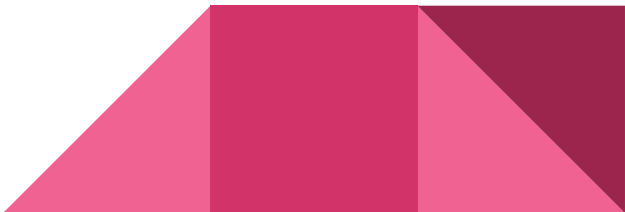
$$\hat{f}(x) = f(x, \hat{\theta}),$$

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E}_{x,y}[L(y, f(x, \theta))]$$

$$\hat{f}(x) = \sum_{i=0}^{t-1} \hat{f}_i(x),$$

$$(\rho_t, \theta_t) = \arg \min_{\rho, \theta} \mathbb{E}_{x,y}[L(y, \hat{f}(x) + \rho \cdot h(x, \theta))],$$

$$\hat{f}_t(x) = \rho_t \cdot h(x, \theta_t)$$



Попробуем использовать градиентный спуск *в пространстве функций*

$$r_{it} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x)=\hat{f}(x)}, \quad \text{for } i = 1, \dots, n,$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n (r_{it} - h(x_i, \theta))^2, \quad \rho_t = \arg \min_{\rho} \sum_{i=1}^n L(y_i, \hat{f}(x_i) + \rho \cdot h(x_i, \theta_t))$$

In linear regression case with MSE loss:

$$r_{it} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x)=\hat{f}(x)} = -2(\hat{y}_i - y_i) \propto \hat{y}_i - y_i$$