# Introduction to Machine Learning
# Lecture 6: Gradient boosting

Harbour.Space University
February 2021
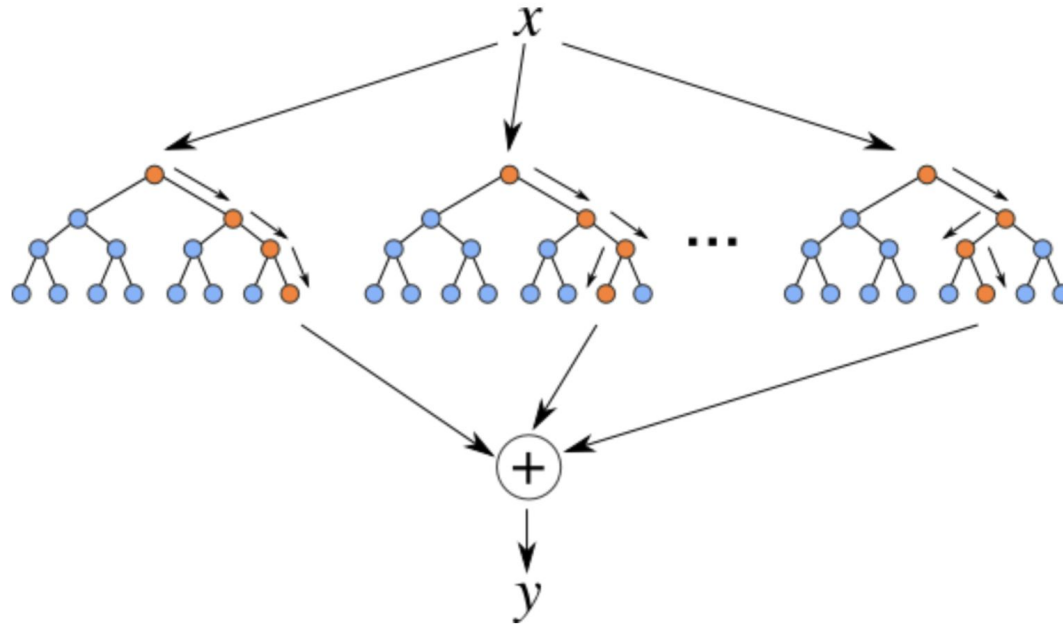
**Nikolay Karpachev**

# Outline

1. Random Forest recap. Out-of-bag error.
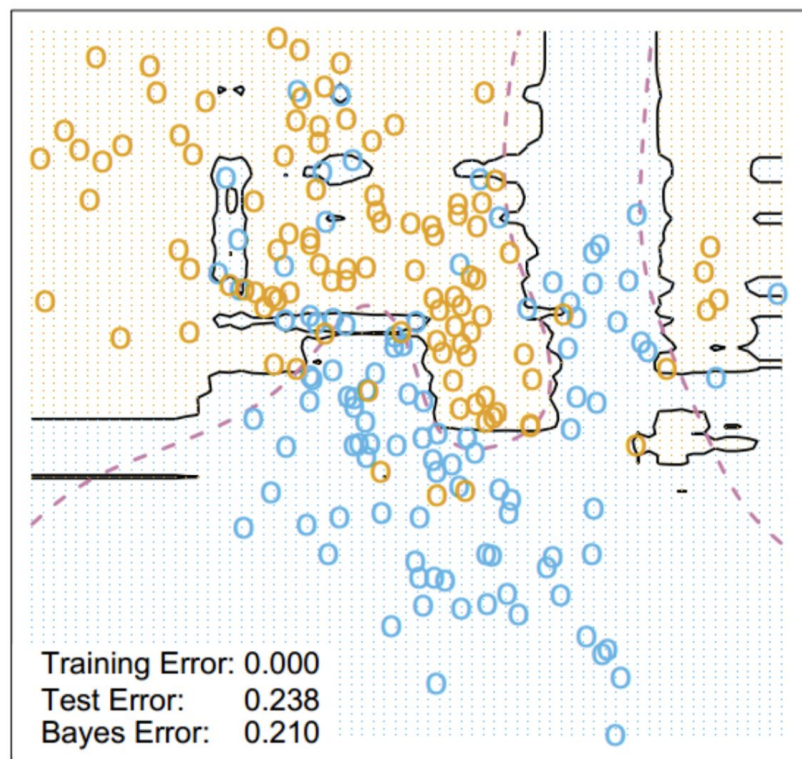2. Boosting technique
3. Gradient Boosting

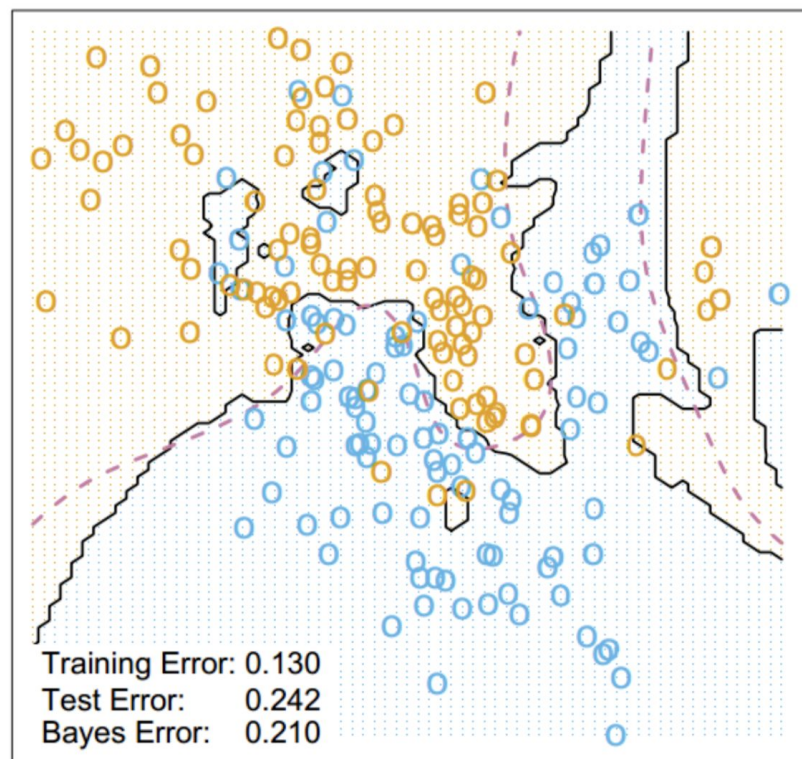## Bagging + RSM = Random Forest

# Random Forest

- One of the greatest "universal" models.
- There are some modifications: Extremely Randomized Trees, Isolation Forest, etc.
- Allows to use train data for validation: OOB

$$\text{OOB} = \sum_{i=1}^{\ell} L\left(y_i, \frac{1}{\sum_{n=1}^{N}[x_i \notin X_n]} \sum_{n=1}^{N}[x_i \notin X_n] b_n(x_i)\right)$$
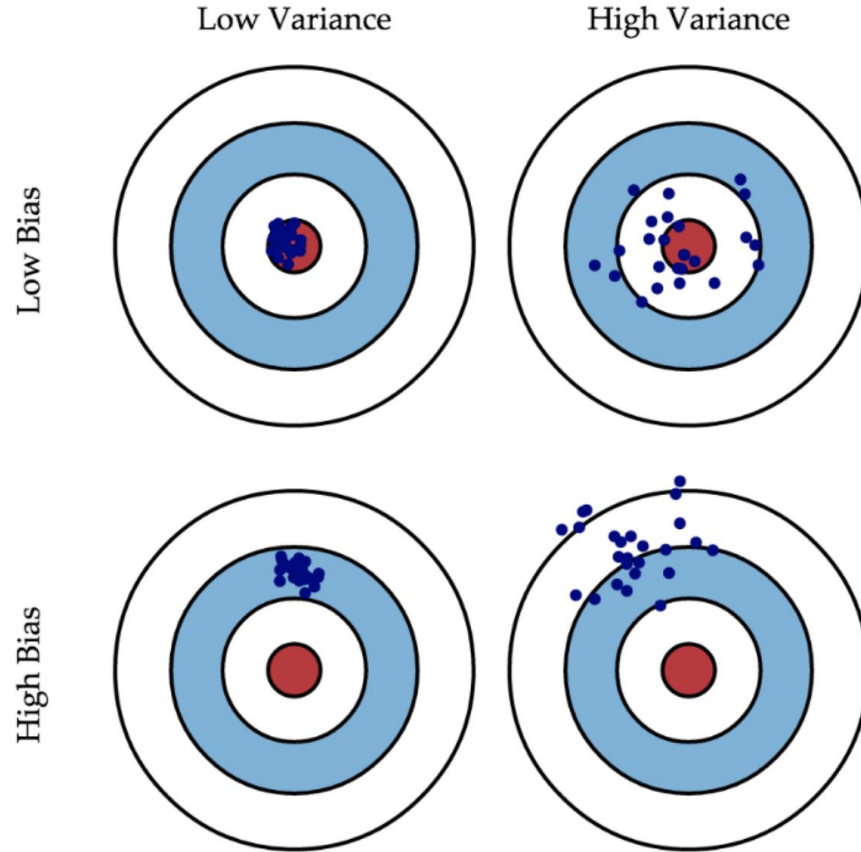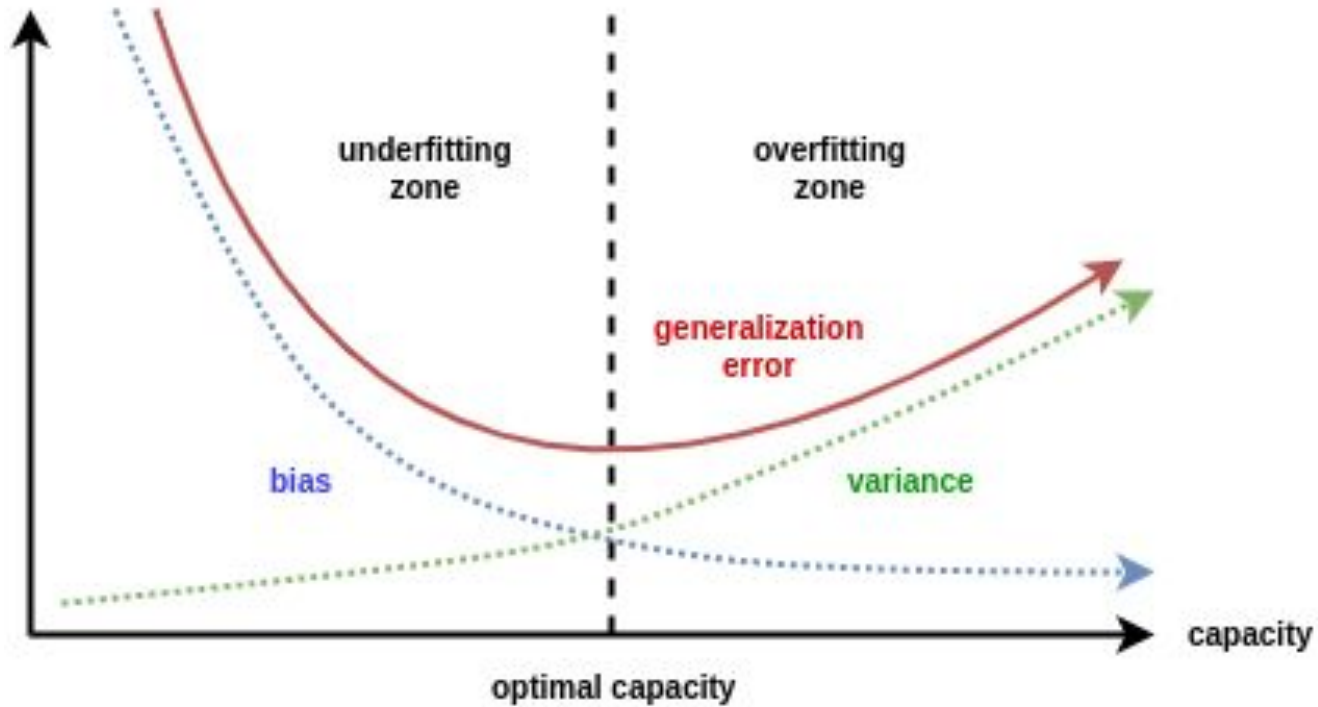
## Random Forest Classifier

Training Error: 0.000
Test Error: 0.238
Bayes Error: 0.210

## 3−Nearest Neighbors

Training Error: 0.130
Test Error: 0.242
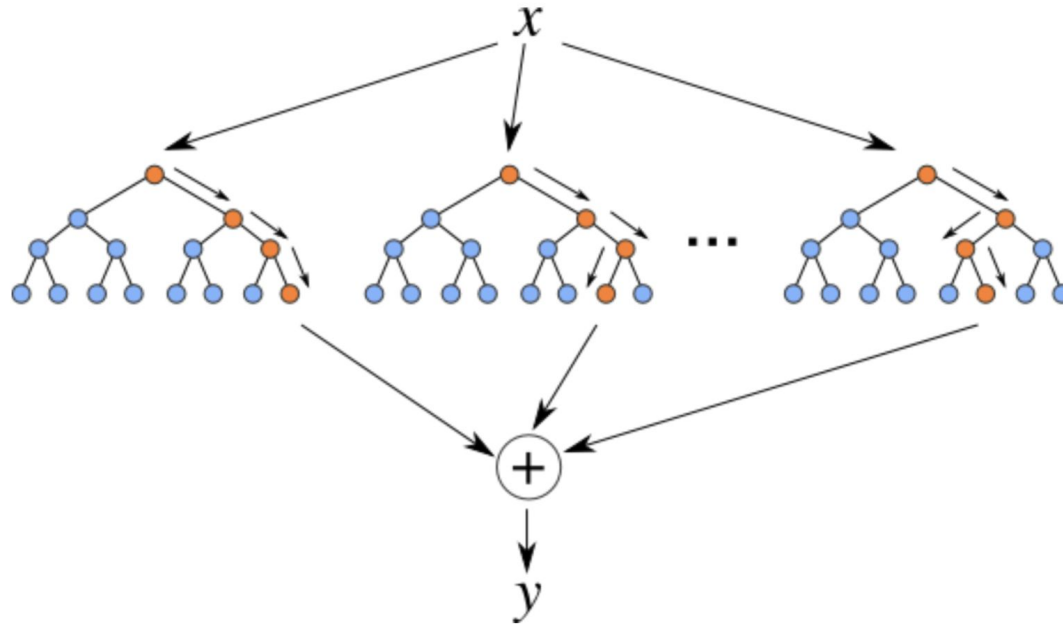Bayes Error: 0.210

# Bias-variance tradeoff

# Bias-variance tradeoff

## Is Random Forest decreasing bias or variance by building the trees ensemble?
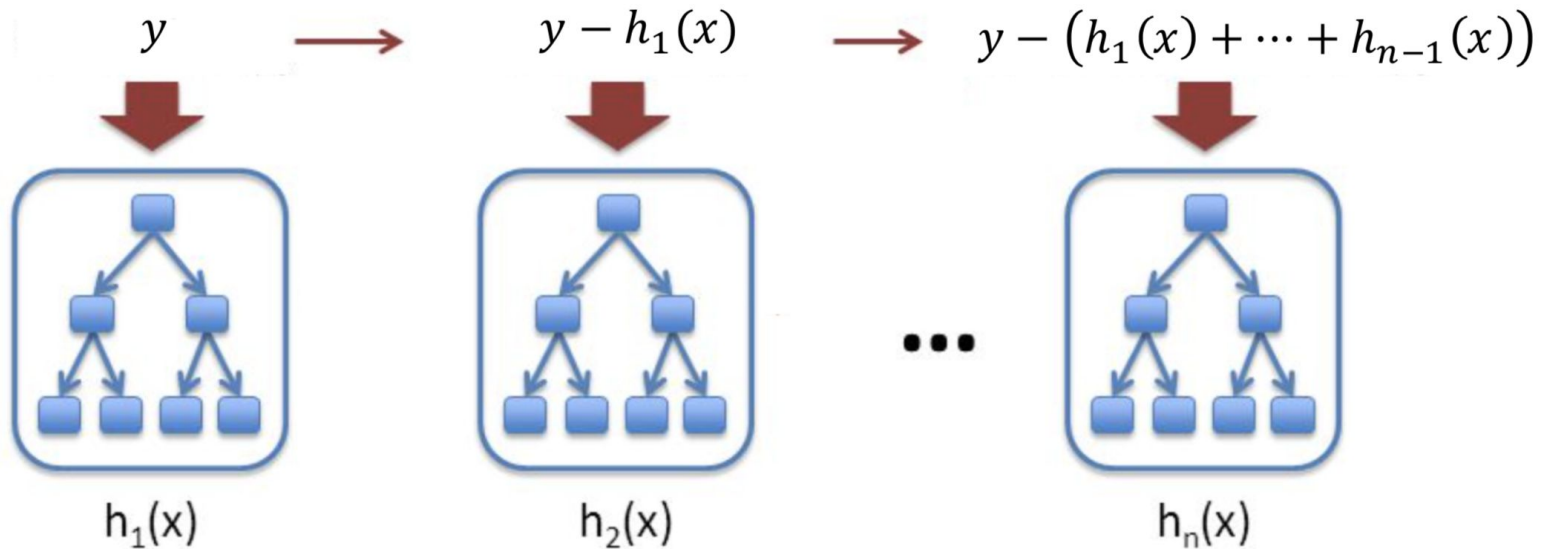
# Boosting

# Boosting

**Boosting** is a technique to combine multiple weak algorithms into one strong ensemble.

- Boosting adds algorithms incrementally
- At the current step, ensemble consists of N algorithms
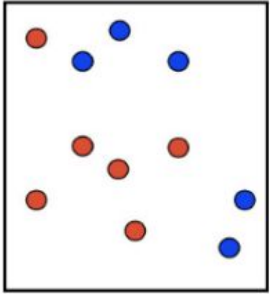- Algorithm N+1 is trained with respect to already built composition to minimize the total error
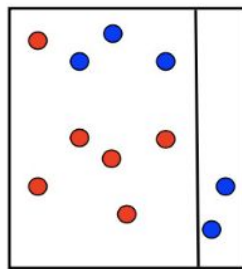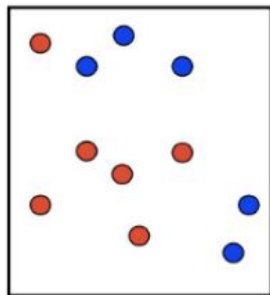
Example: regression

$$a_n(x) = h_1(x) + \cdots + h_n(x)$$



$y \quad \longrightarrow \quad y - h_1(x) \quad \longrightarrow \quad y - \big(h_1(x) + \cdots + h_{n-1}(x)\big)$

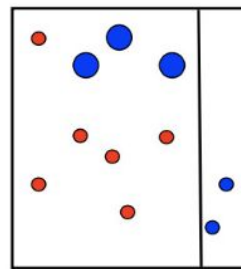$h_1(x)$ $\qquad\qquad$ $h_2(x)$ $\qquad\qquad$ $h_n(x)$
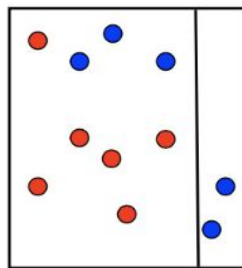
Binary classification problem.
Models - decision stumps.

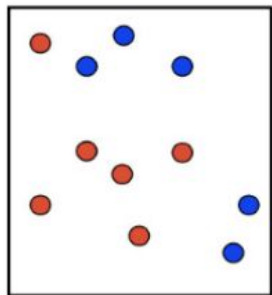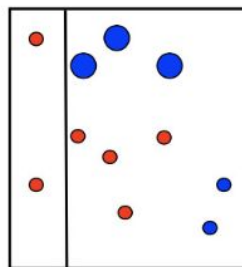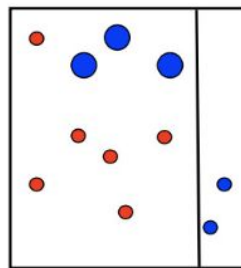# Boosting: intuition

t = 1

Boosting: intuition

t = 1

t = 2

Boosting: intuition

t = 1

t = 2

t = 3

Binary classification problem.
Models - decision stumps.



$$\alpha_1 \quad + \alpha_2 \quad + \alpha_3$$

$$=$$

Denote dataset $\{(x_i, y_i)\}_{i=1,\ldots,n}$, loss function $L(y, f)$.

Denote dataset $\left\{(x_i, y_i)\right\}_{i=1,\ldots,n}$, loss function $L(y, f)$.

Optimal model:

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

Denote dataset $\{(x_i, y_i)\}_{i=1,\ldots,n}$, loss function $L(y, f)$.

Optimal model:

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

Let it be from parametric family: $\hat{f}(x) = f(x, \hat{\theta})$,

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

# Gradient boosting: theory

- Let's assume real-valued y
- Gradient boosting builds a single model as a weighted sum of *"base learners"*

$$\hat{F}(x) = \sum_{i=1}^{M} \gamma_i h_i(x) + \text{const.}$$

# Gradient boosting: theory

- The goal is to minimize the loss function on the training set

**Boosting procedure:**

1. Start with the constant function

$$F_0(x) = \arg\min_{\gamma} \sum_{i=1}^{n} L(y_i, \gamma),$$

2. Incrementally expand the composition by minimizing the **residual error**

$$F_m(x) = F_{m-1}(x) + \arg\min_{h_m \in \mathcal{H}} \left[ \sum_{i=1}^{n} L(y_i, F_{m-1}(x_i) + h_m(x_i)) \right],$$

We select the optimal *h_m*

Loss of the new composition

# Time for your questions and a coffee break

$$F_m(x) = F_{m-1}(x) + \arg\min_{h_m \in \mathcal{H}} \left[ \sum_{i=1}^{n} L(y_i, F_{m-1}(x_i) + h_m(x_i)) \right],$$

We select the
optimal *h_m*

Loss of the new composition

**How to find the optimal *h_m*?**

# Gradient boosting: theory

**Key idea: apply gradient descent in the space of model predictions**

$$F_m(x) = F_{m-1}(x) - \gamma_m \sum_{i=1}^{n} \nabla_{F_{m-1}} L(y_i, F_{m-1}(x_i)),$$

Loss gradient at current prediction point

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^{n} L\left(y_i, F_{m-1}(x_i) - \gamma \nabla_{F_{m-1}} L(y_i, F_{m-1}(x_i))\right),$$

Find the optimal "learning rate"

"Learning rate"

- We fit the next base learner to the anti-gradient of the error
- Thus, the composition resembles gradient descent
- Each base learner addition = gradient step in functional space

## Complete algorithm:

- Initialize the model with a constant value

$$F_0(x) = \arg\min_{\gamma} \sum_{i=1}^{n} L(y_i, \gamma).$$

Incrementally:
- Compute *pseudo-residuals* for each data point

$$r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \ldots, n.$$

- Fit new base learner to pseudo-residuals
- Find the optimal "learning rate" coefficient
- Update the model

$$\gamma_m = \arg\min_{\gamma} \sum_{i=1}^{n} L\left(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)\right).$$

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

# Gradient boosting: beautiful demo

Great demo:

http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html

# Gradient boosting: theory

What we need:

- Data.
- Loss function and its gradient.
- Family of algorithms (with constraints on hyperparameters if necessary).
- Number of iterations M.
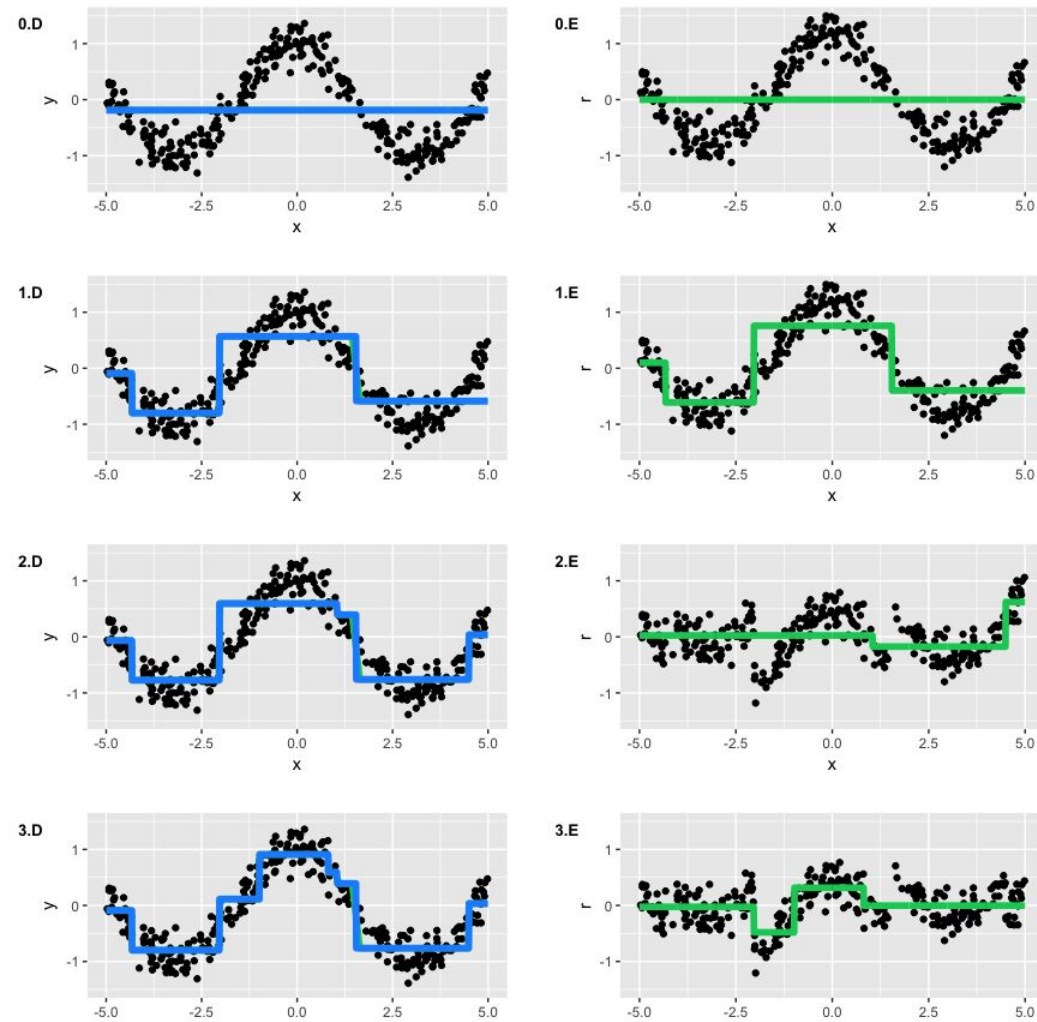- Initial value (GBM by Friedman): constant.

# Gradient boosting: example

What we need:

- Data: toy dataset  $y = cos(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \frac{1}{5}), x \in [-5, 5]$
- Loss function: MSE
- Family of algorithms: decision trees with depth 2
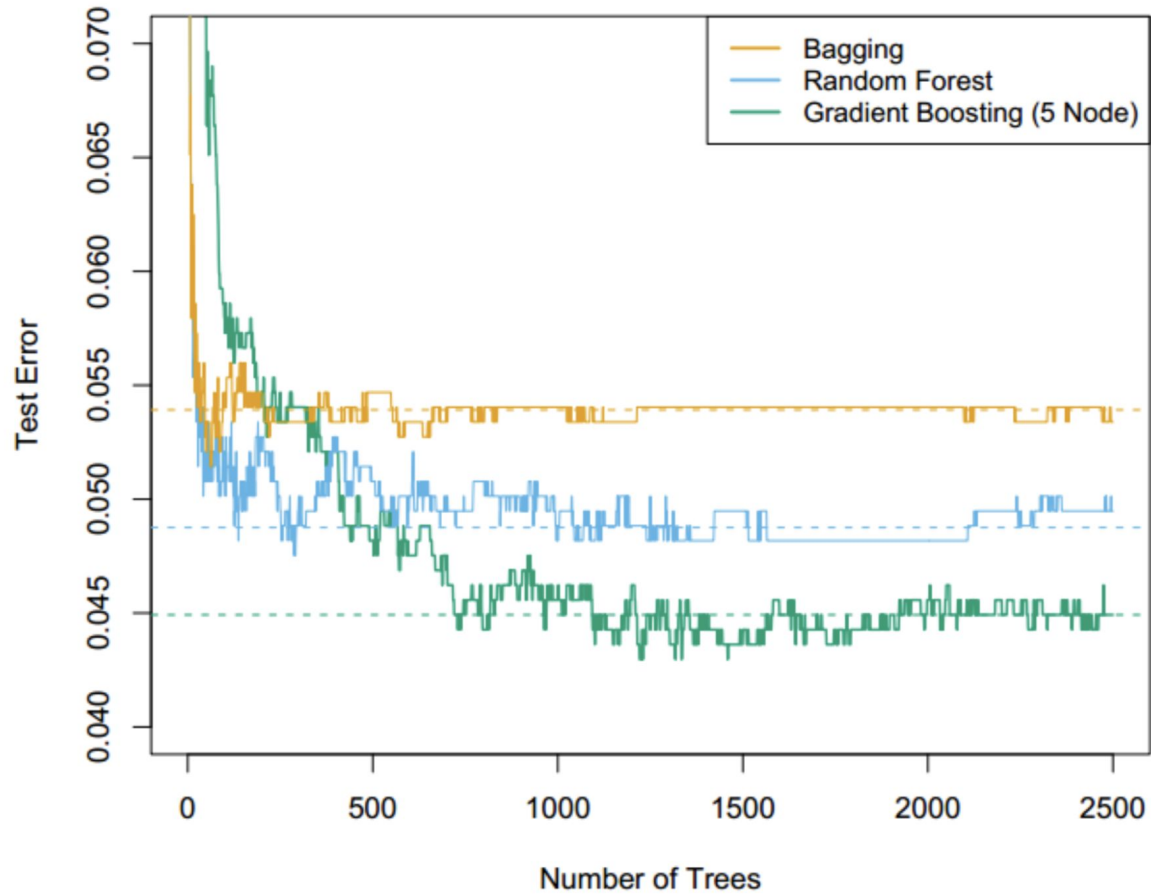- Number of iterations M = 3
- Initial value: just mean value

# Gradient boosting: example

Left: full ensemble on each step.

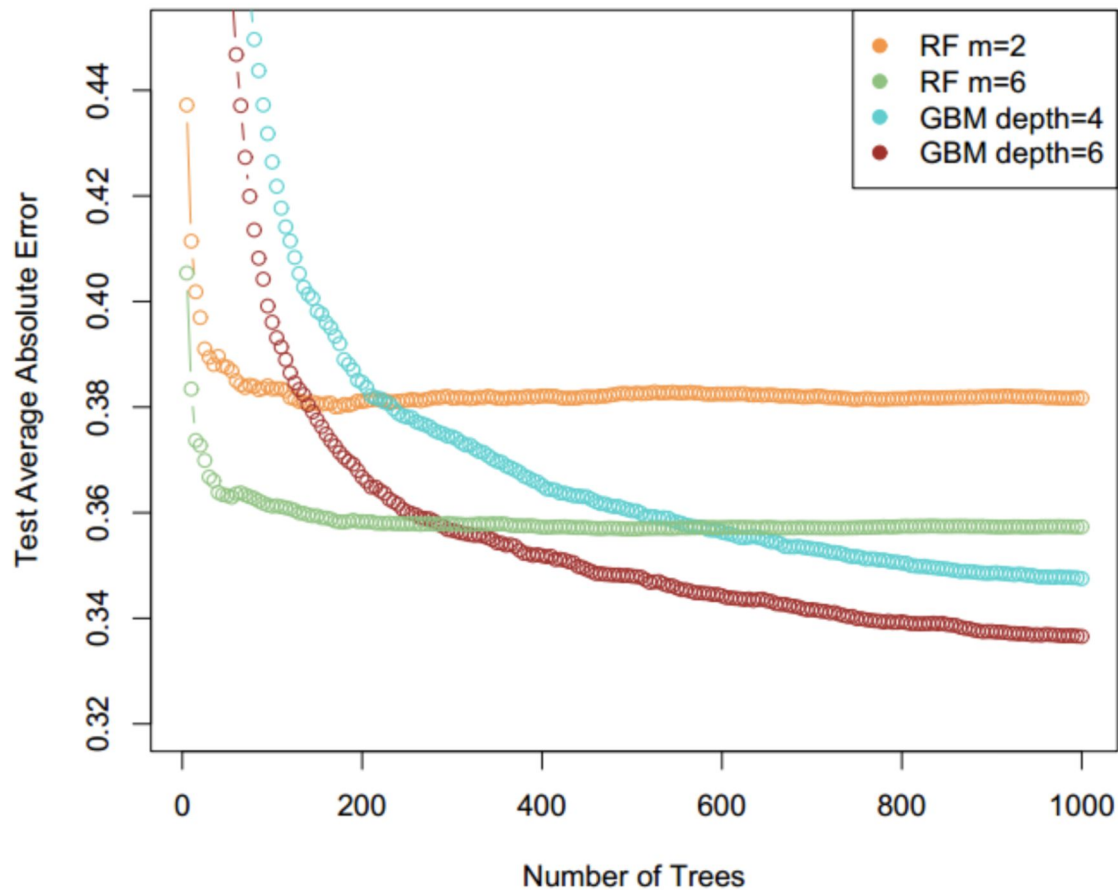Right: additional tree decisions.

Example by ODS; source: https://habr.com/ru/company/ods/blog/327250/

**Spam Data**

California Housing Data

Legend:
- RF m=2
- RF m=6
- GBM depth=4
- GBM depth=6

Y-axis: Test Average Absolute Error
X-axis: Number of Trees

31

# Gradient boosting: regularization

Boosting is a very powerful algorithm

But it can overfit the training data and lose the generalization ability
We need to regularize the model

**Regularization options:**
- Constrain number of trees
- Set max. depth of the trees

# Thanks for the attention!