

# Machine Learning

## Lecture 6: Gradient boosting

Young & Yandex



Radoslav Neychev



# Outline

01 Boosting intuitions

02 Gradient boosting

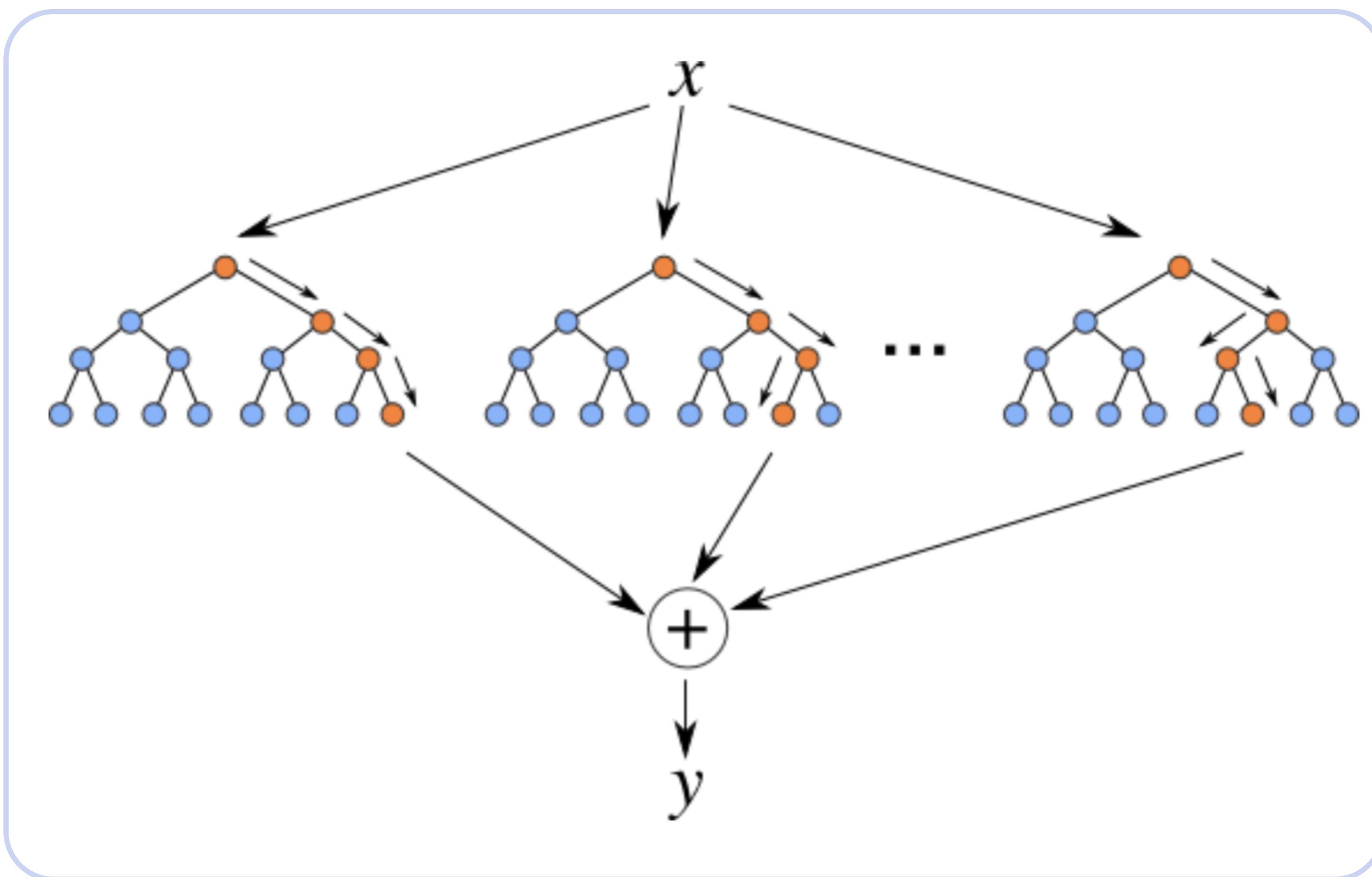
03 Blending

04 Stacking

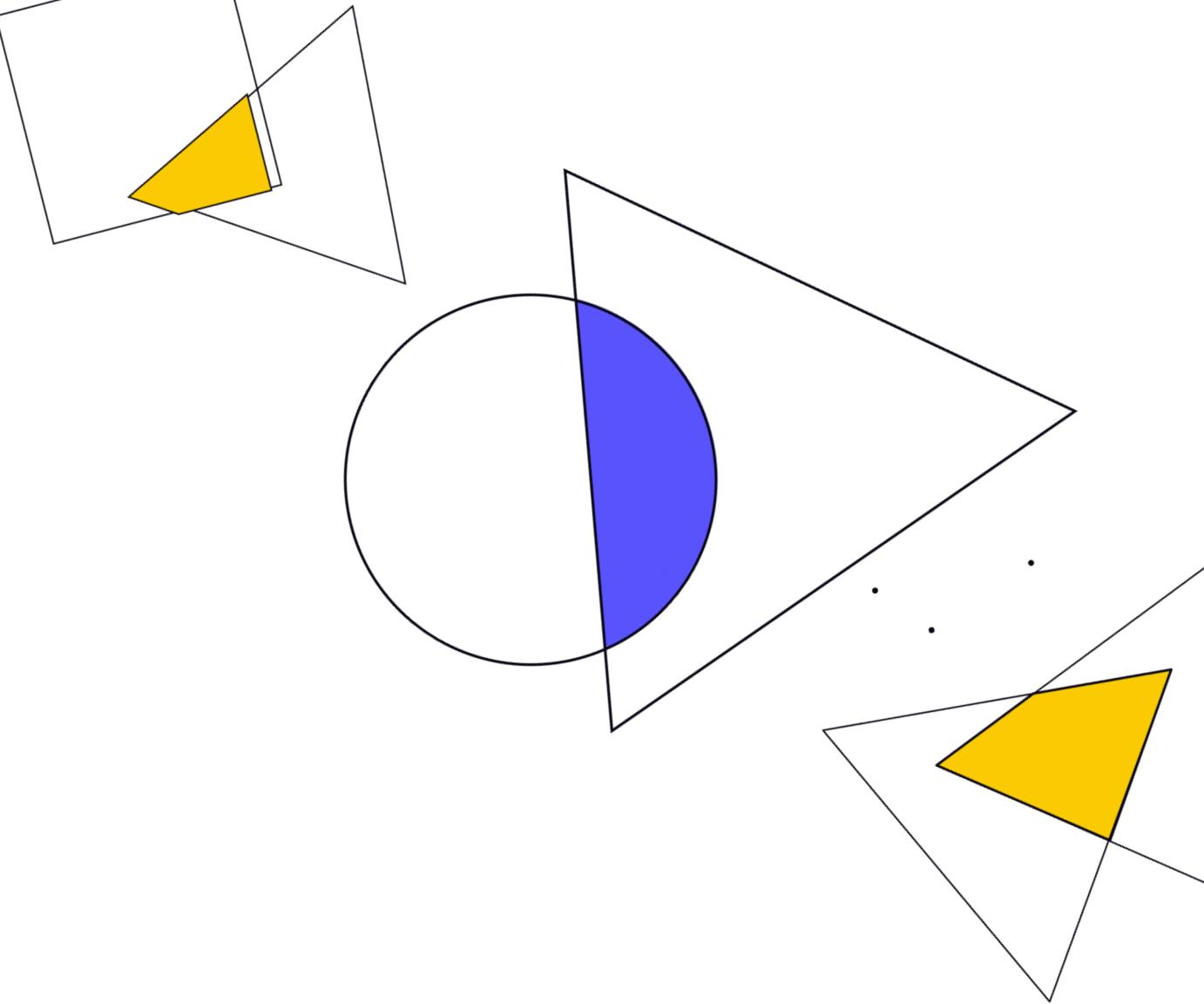


# Random Forest

Bagging + RSM = Random Forest

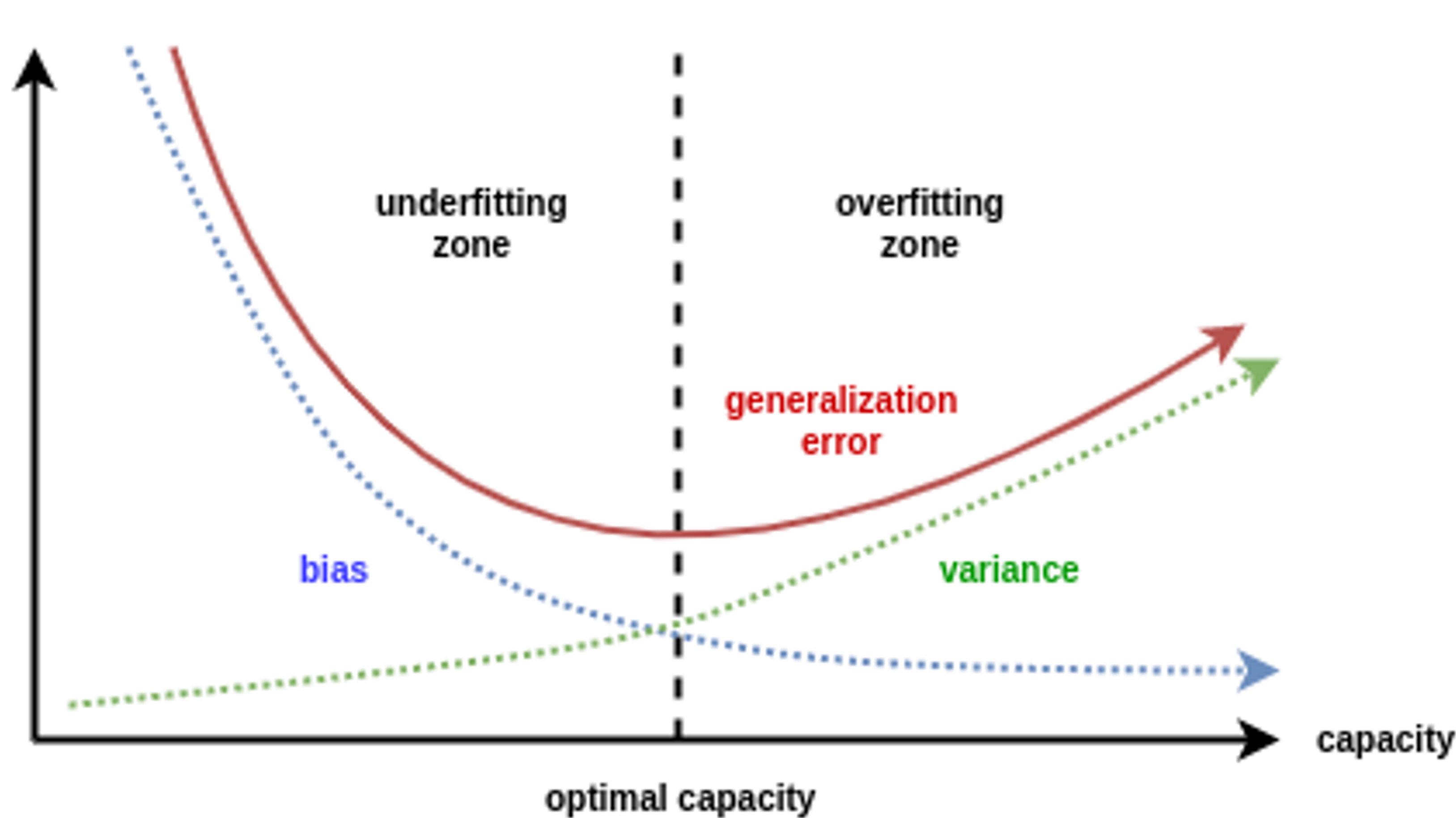


# Random Forest

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

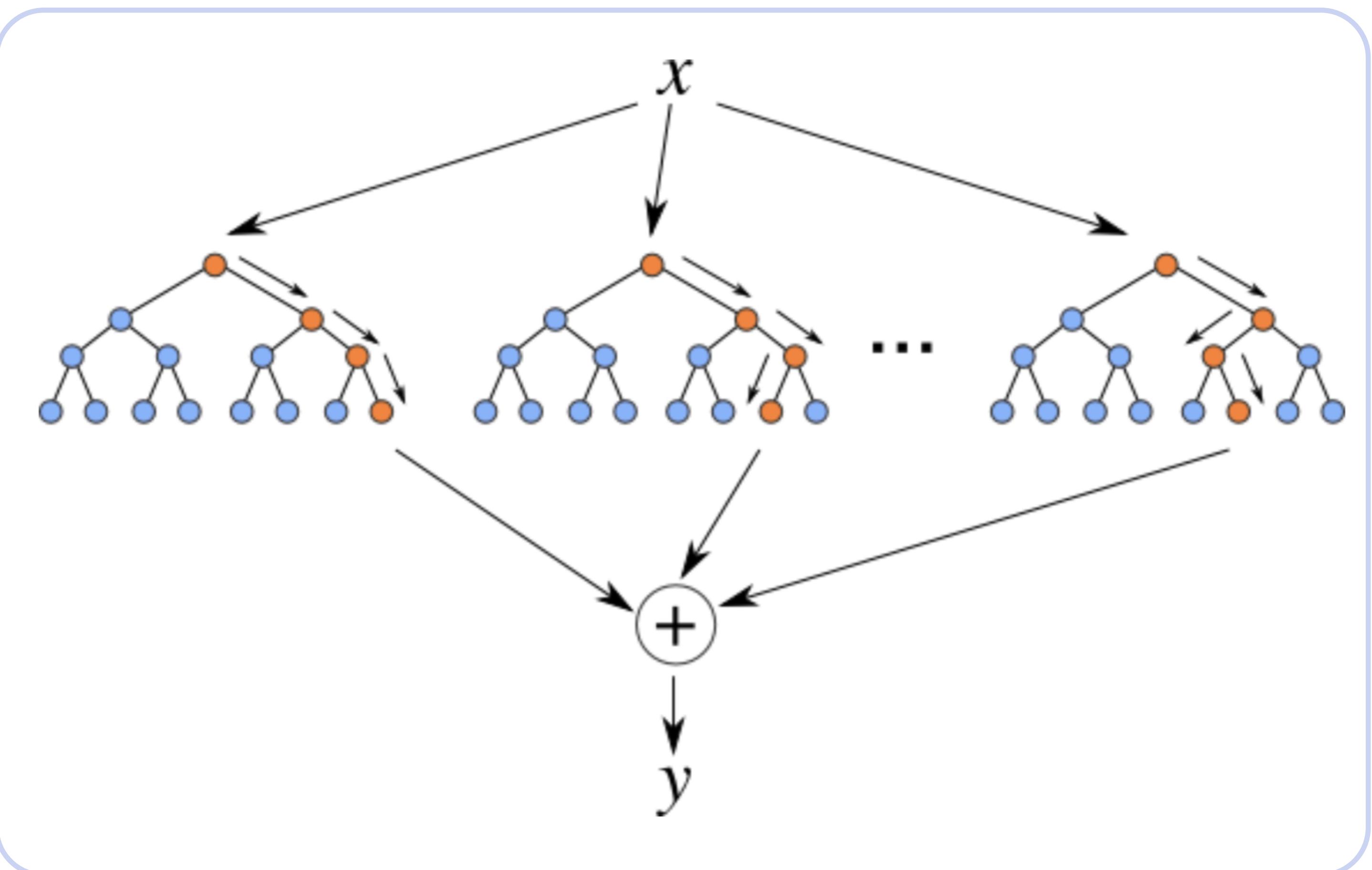
$$\text{OOB} = \sum_{i=1}^{\ell} L \left( y^{(i)}, \frac{1}{\sum_{n=1}^N [\mathbf{x}^{(i)} \notin X_n]} \sum_{n=1}^N [\mathbf{x}^{(i)} \notin X_n] b_n(\mathbf{x}^{(i)}) \right)$$

# Bias-variance tradeoff



# Random Forest

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

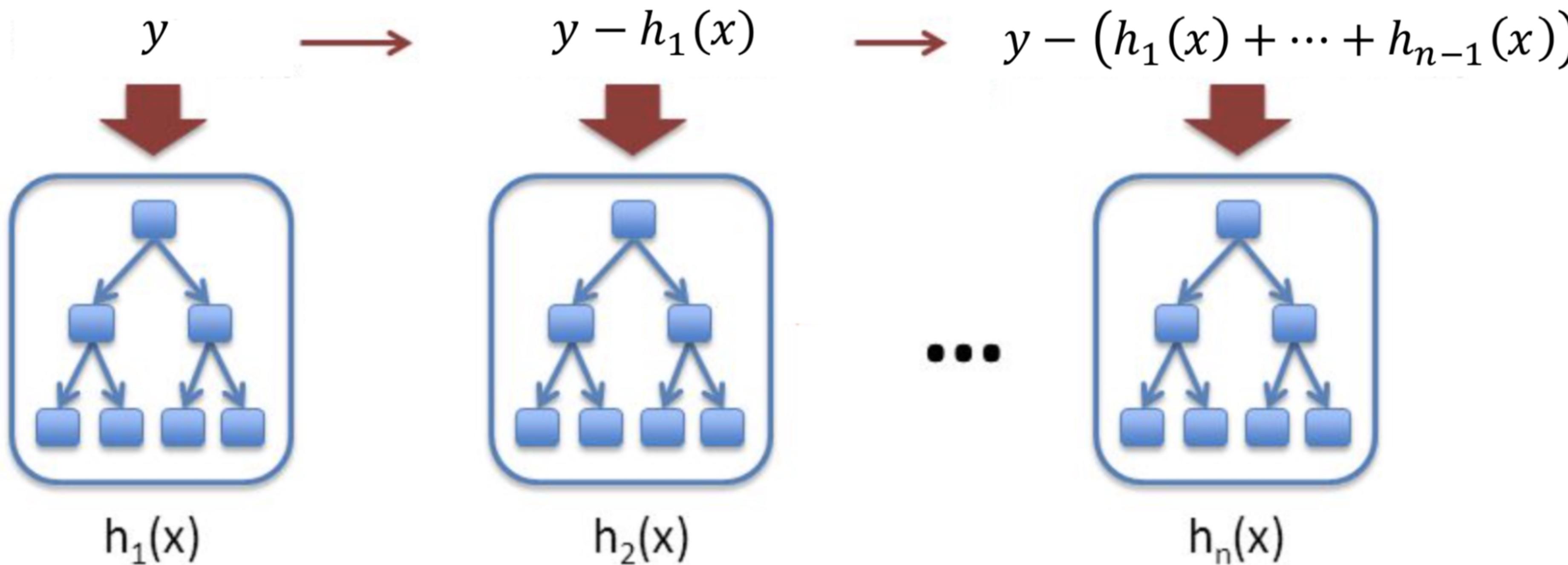


# Boosting



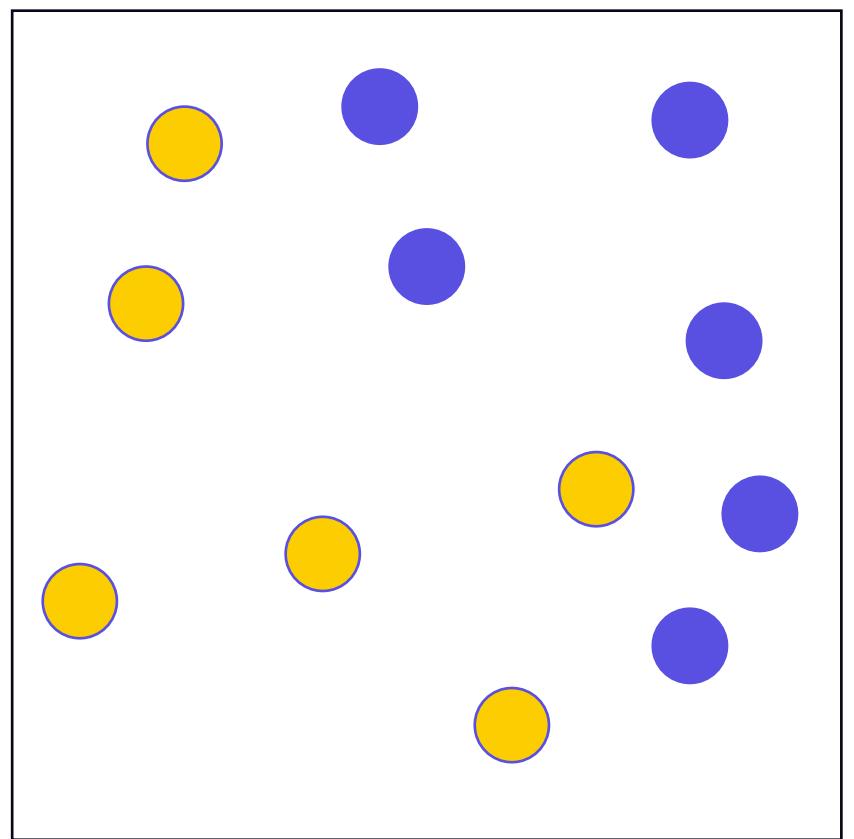
# Boosting

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



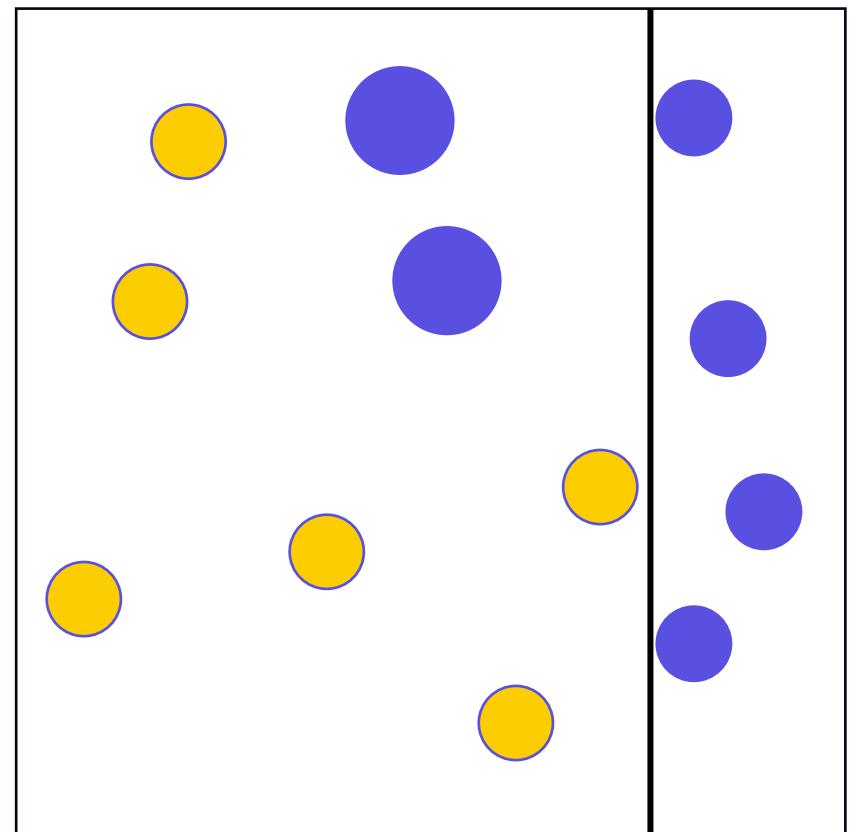
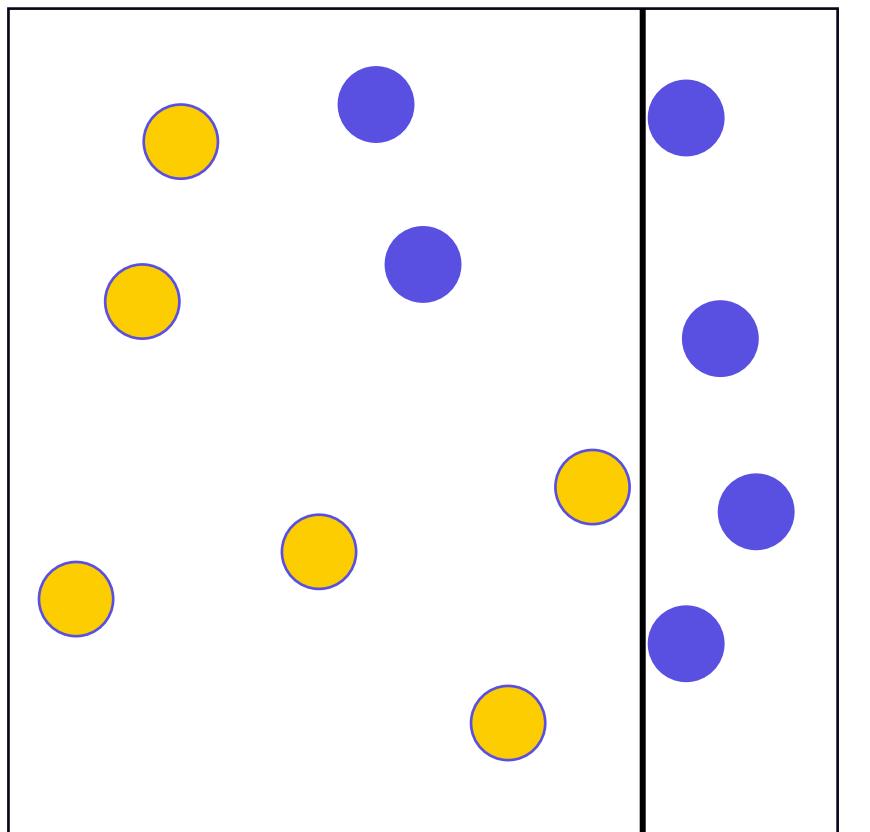
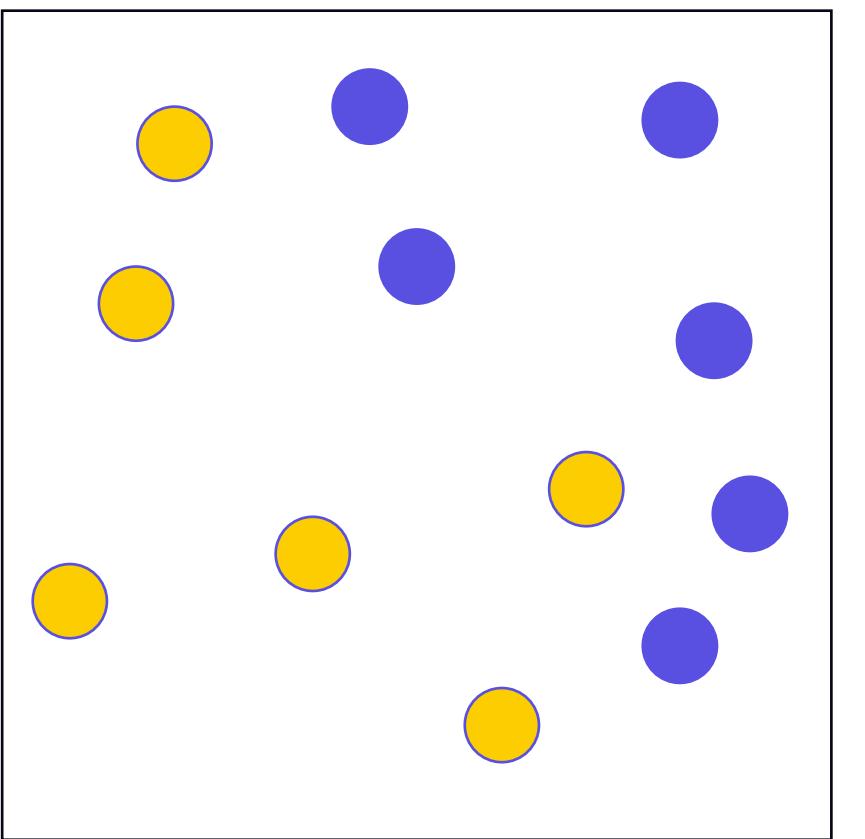
# Boosting: intuition

Binary classification  
Use decision stumps.



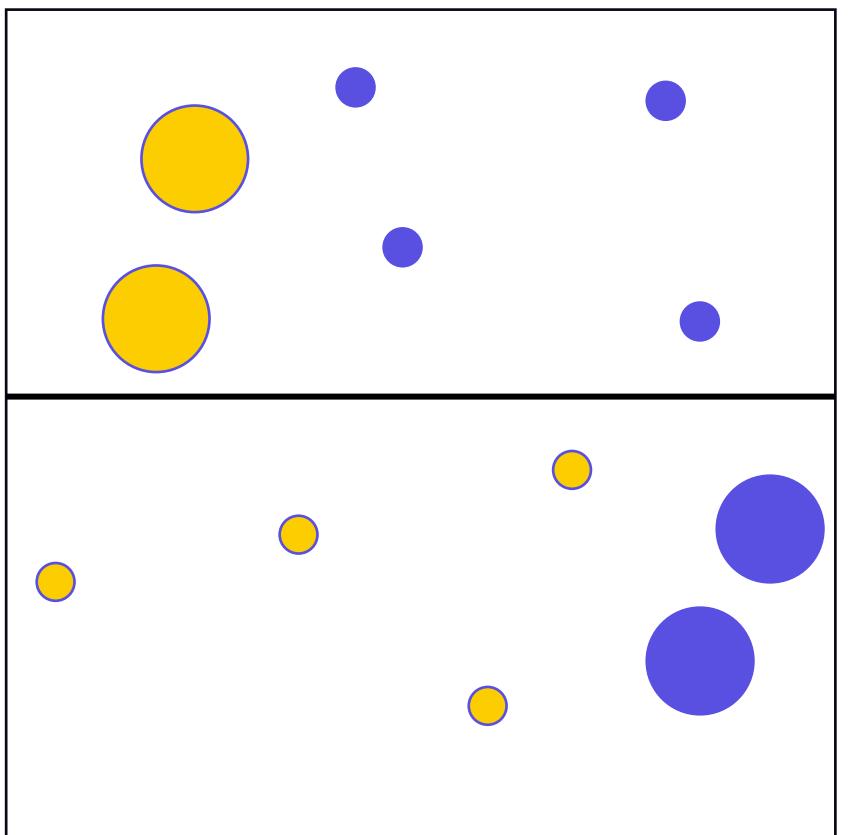
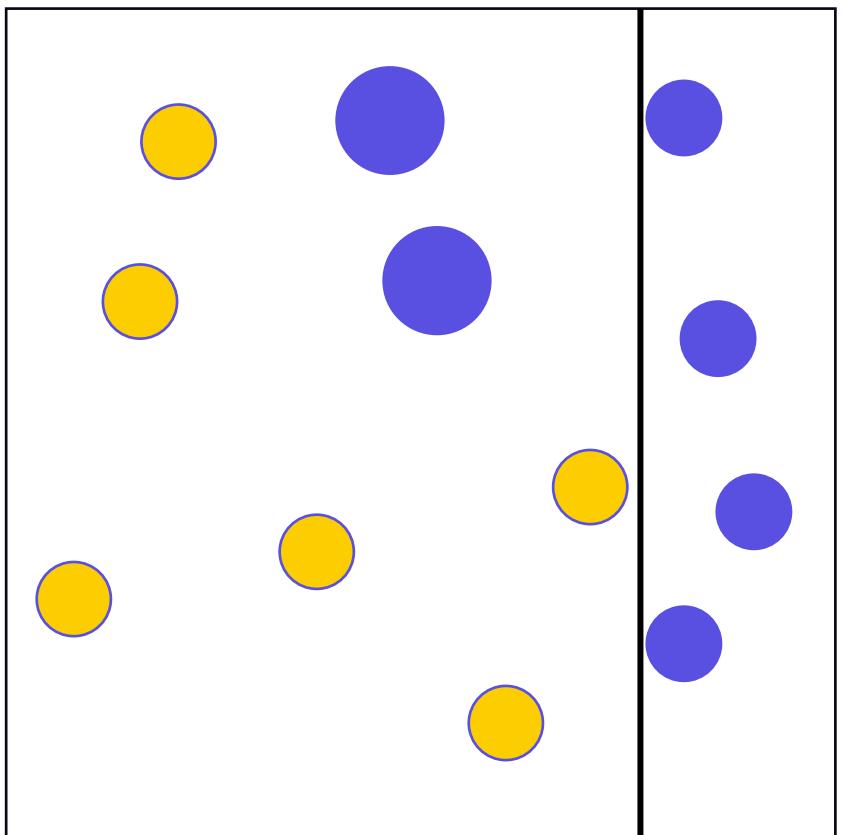
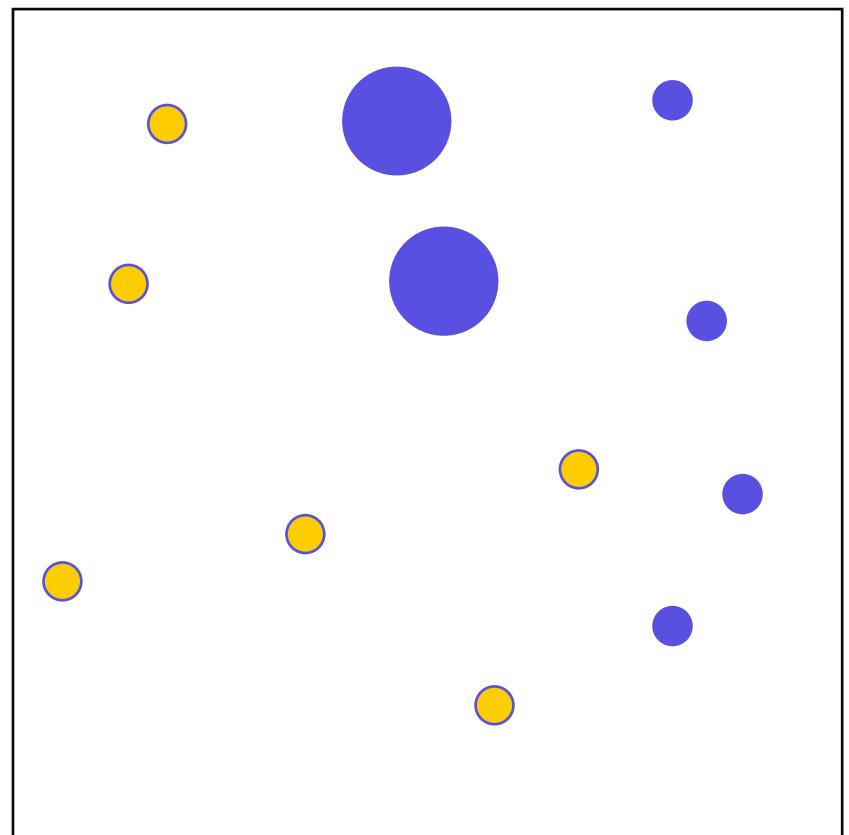
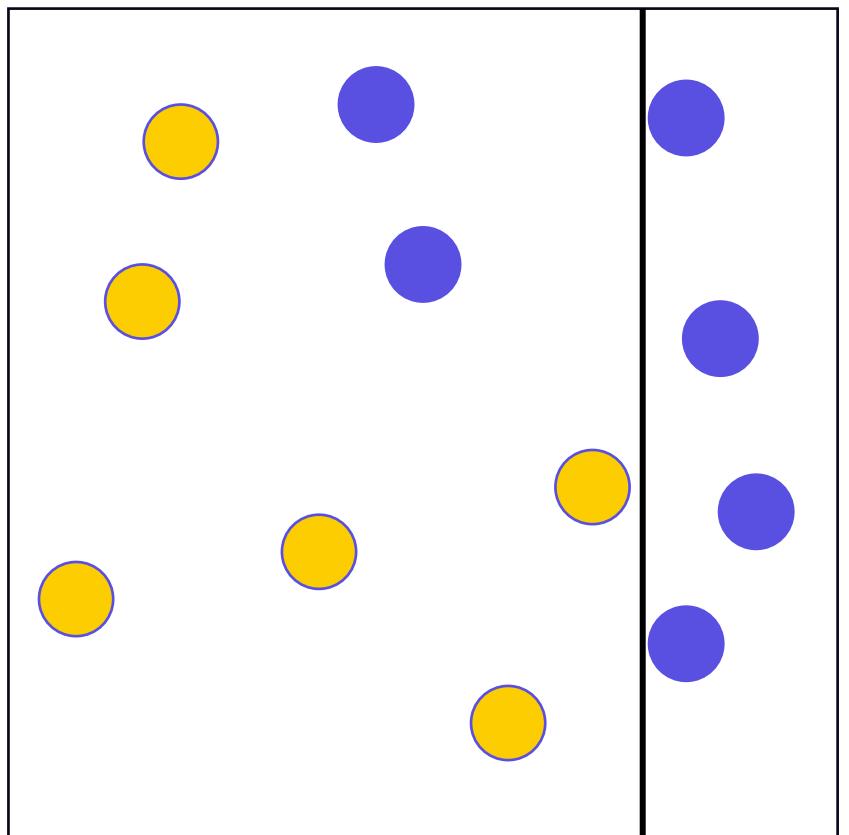
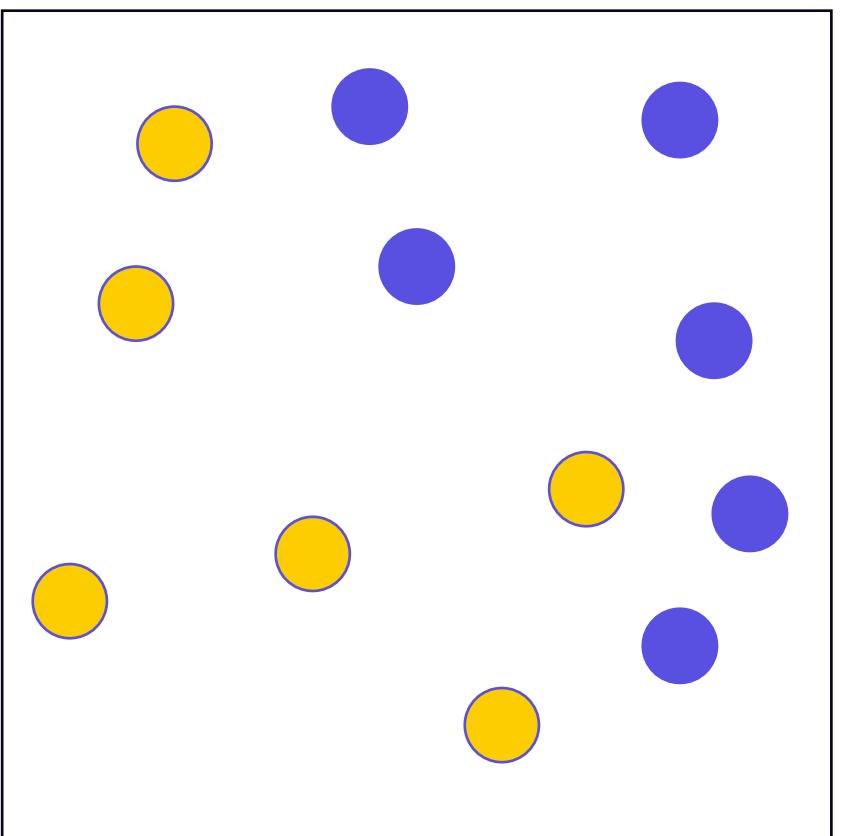
# Boosting: intuition

Binary classification  
Use decision stumps.



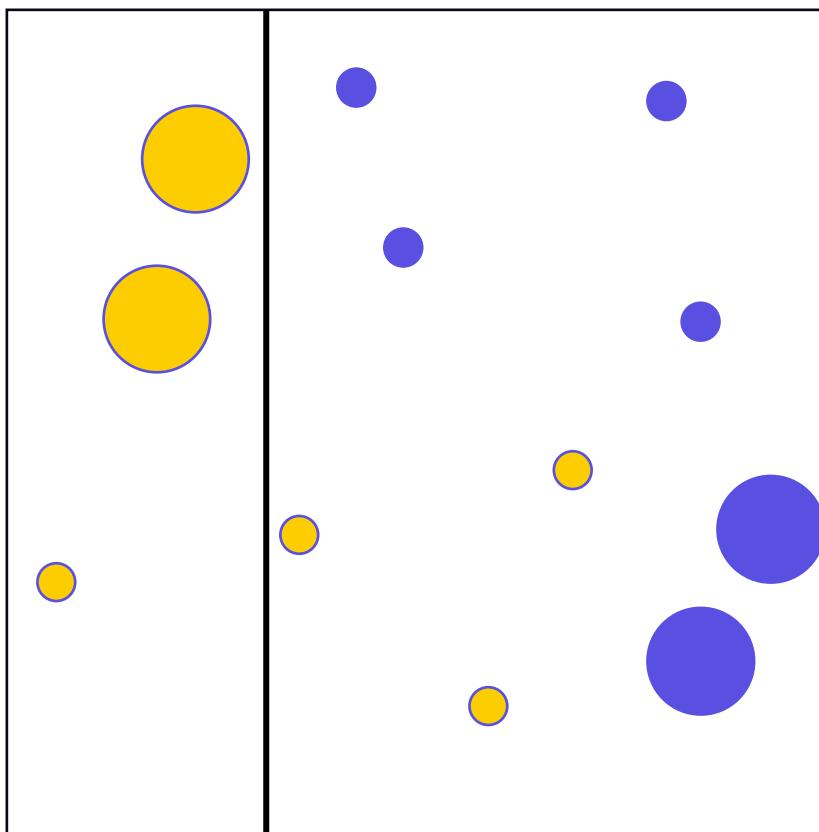
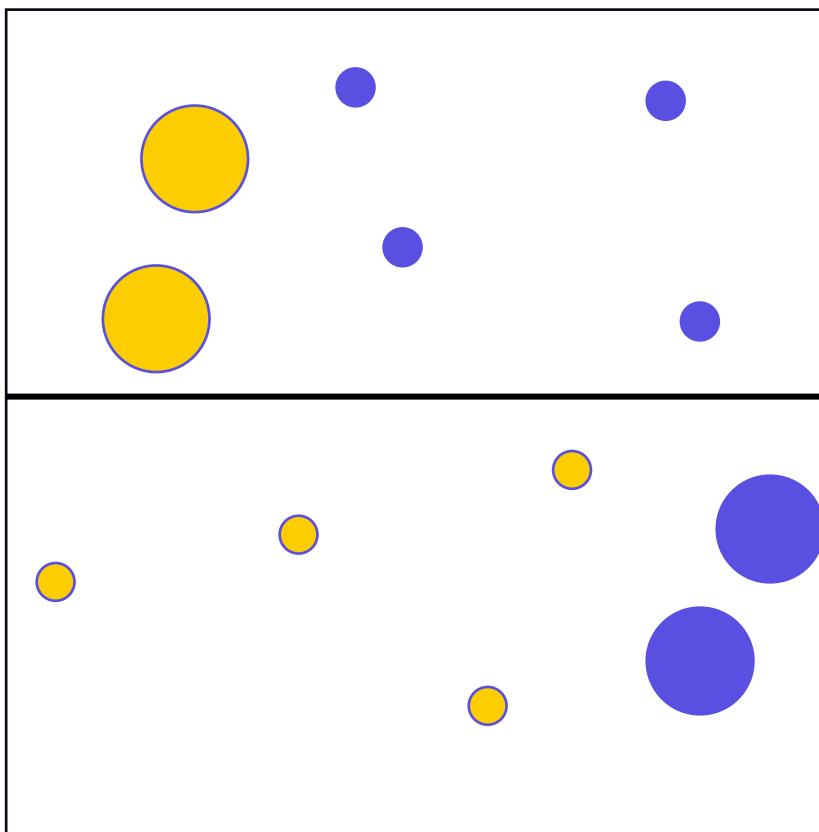
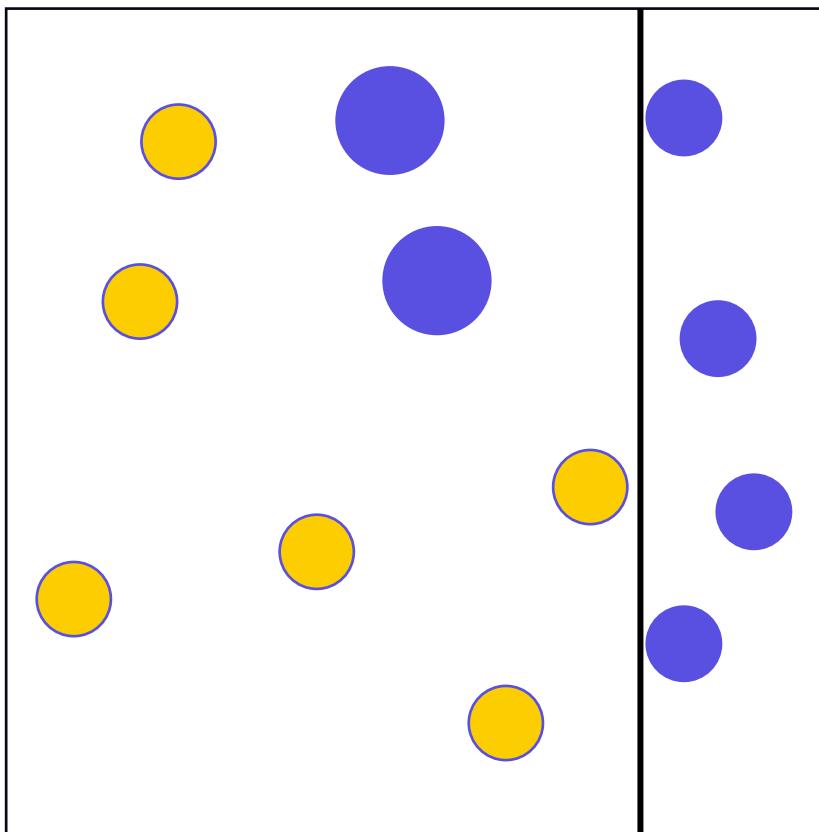
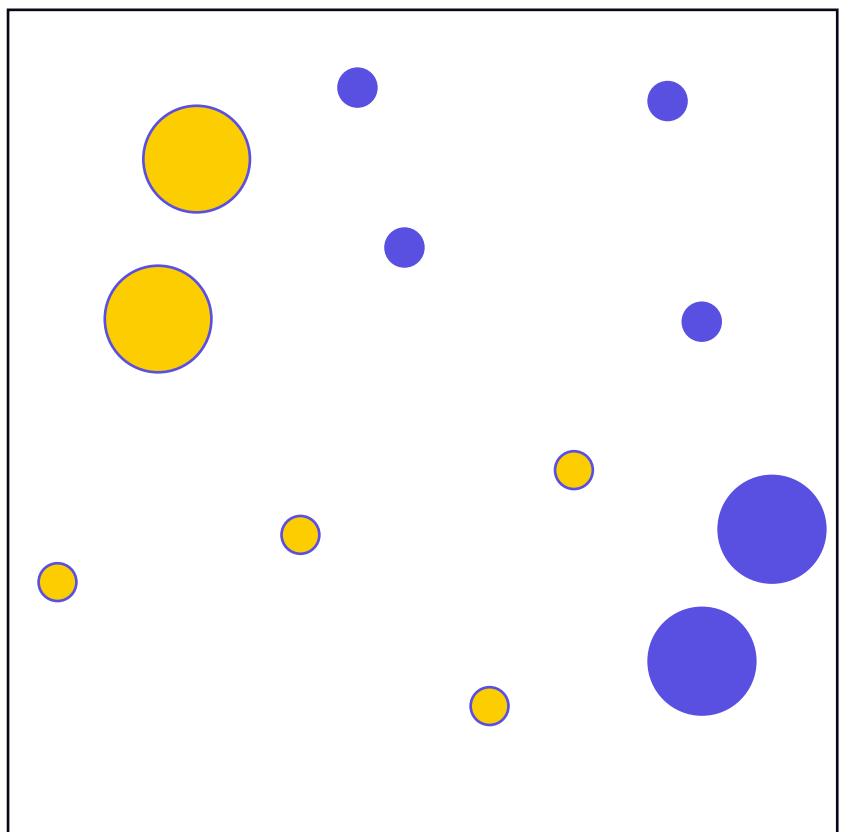
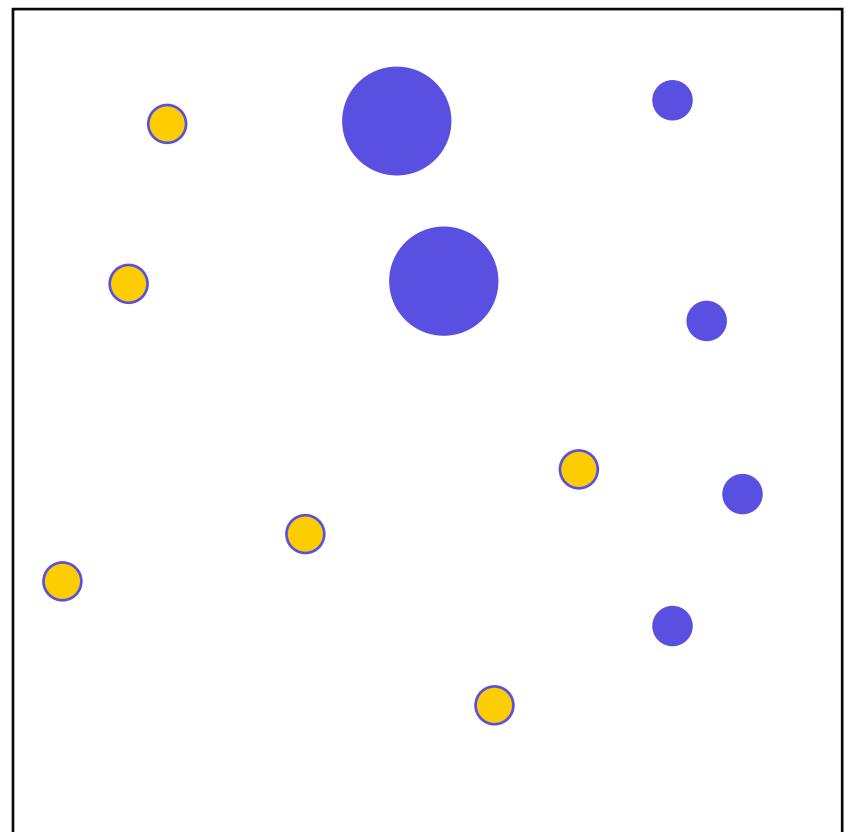
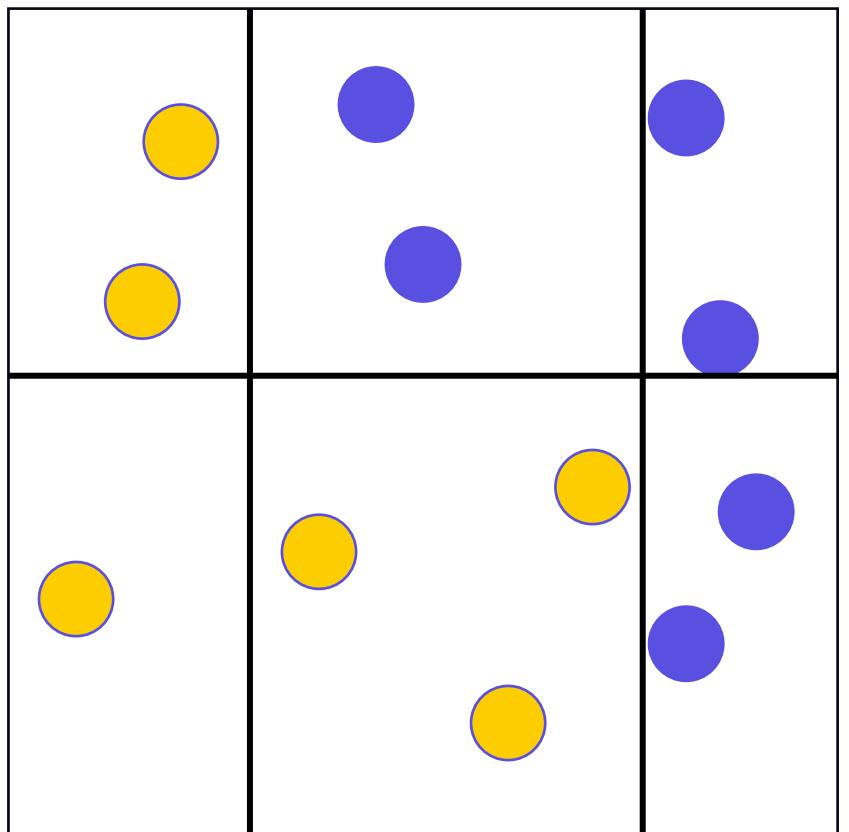
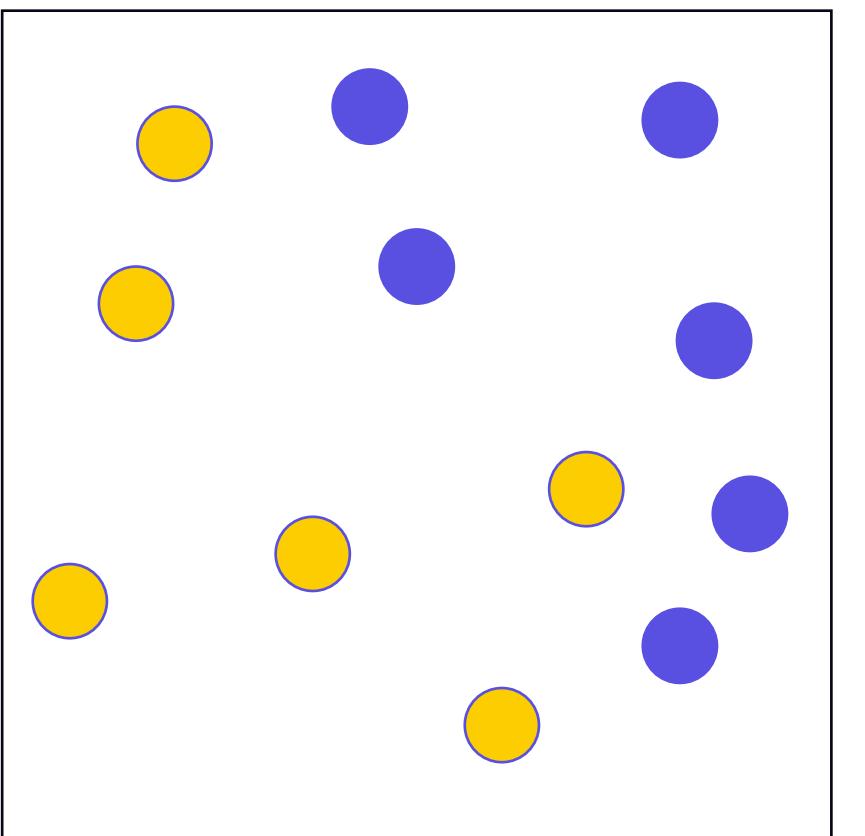
# Boosting: intuition

Binary classification  
Use decision stumps.



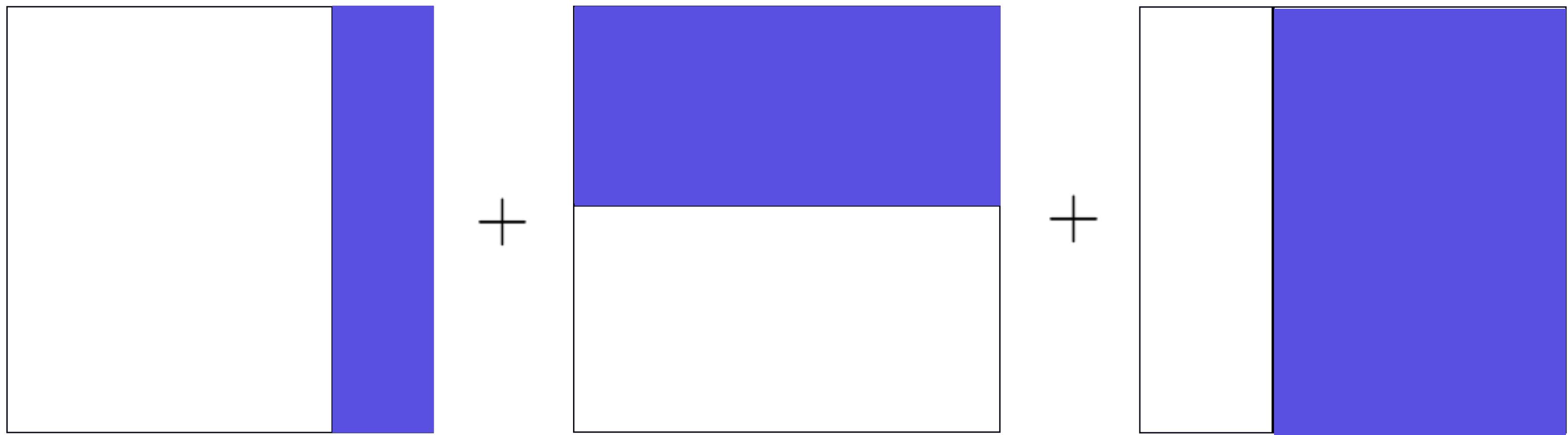
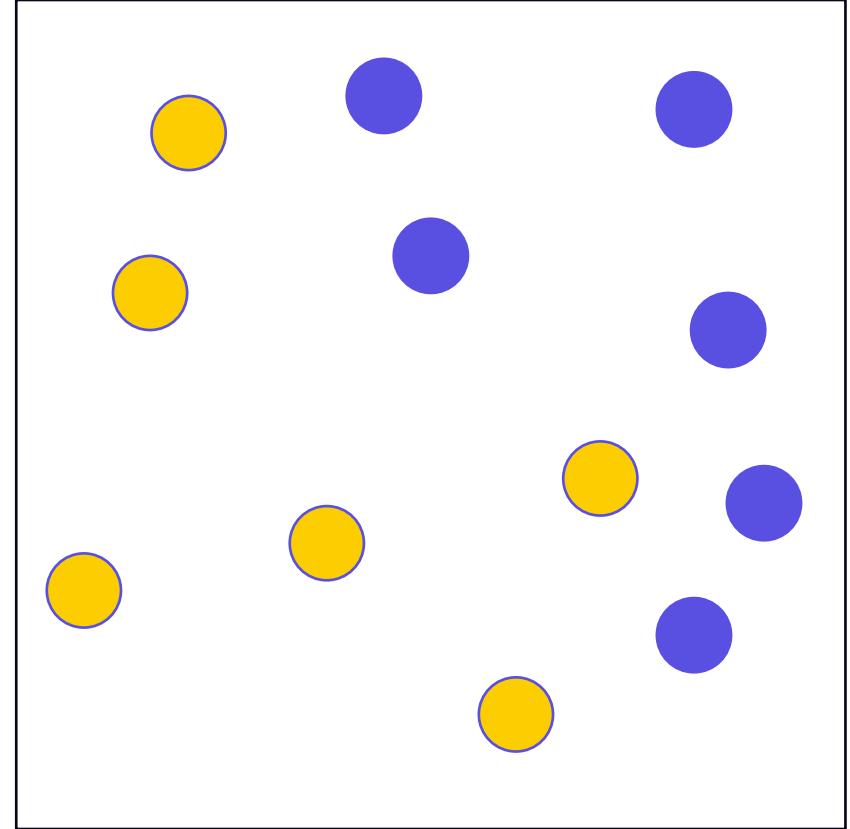
# Boosting: intuition

Binary classification  
Use decision stumps.

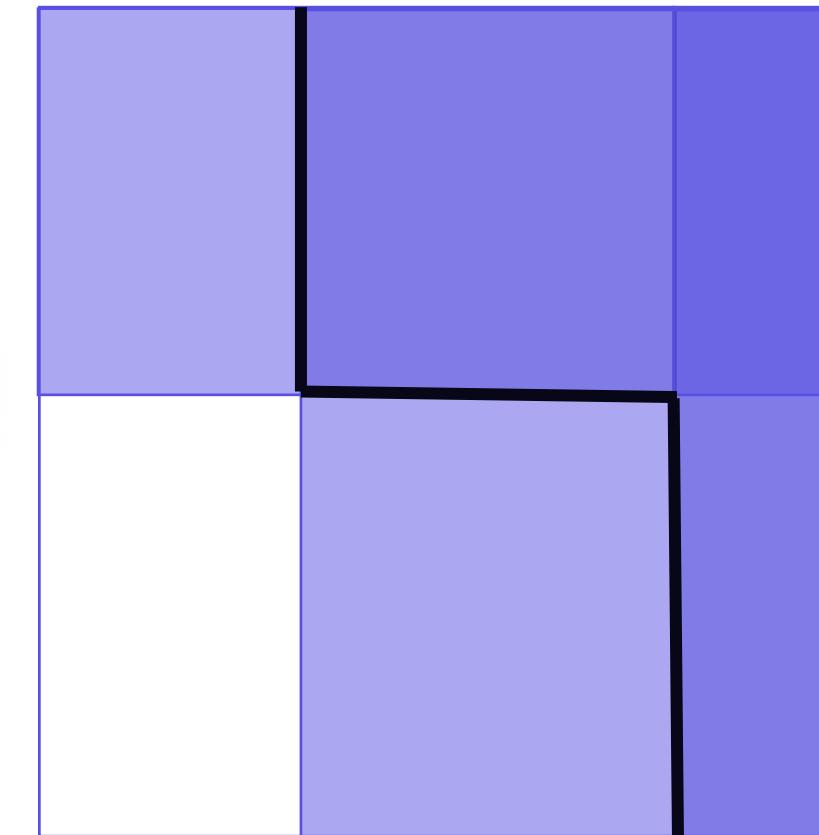


# Boosting: intuition

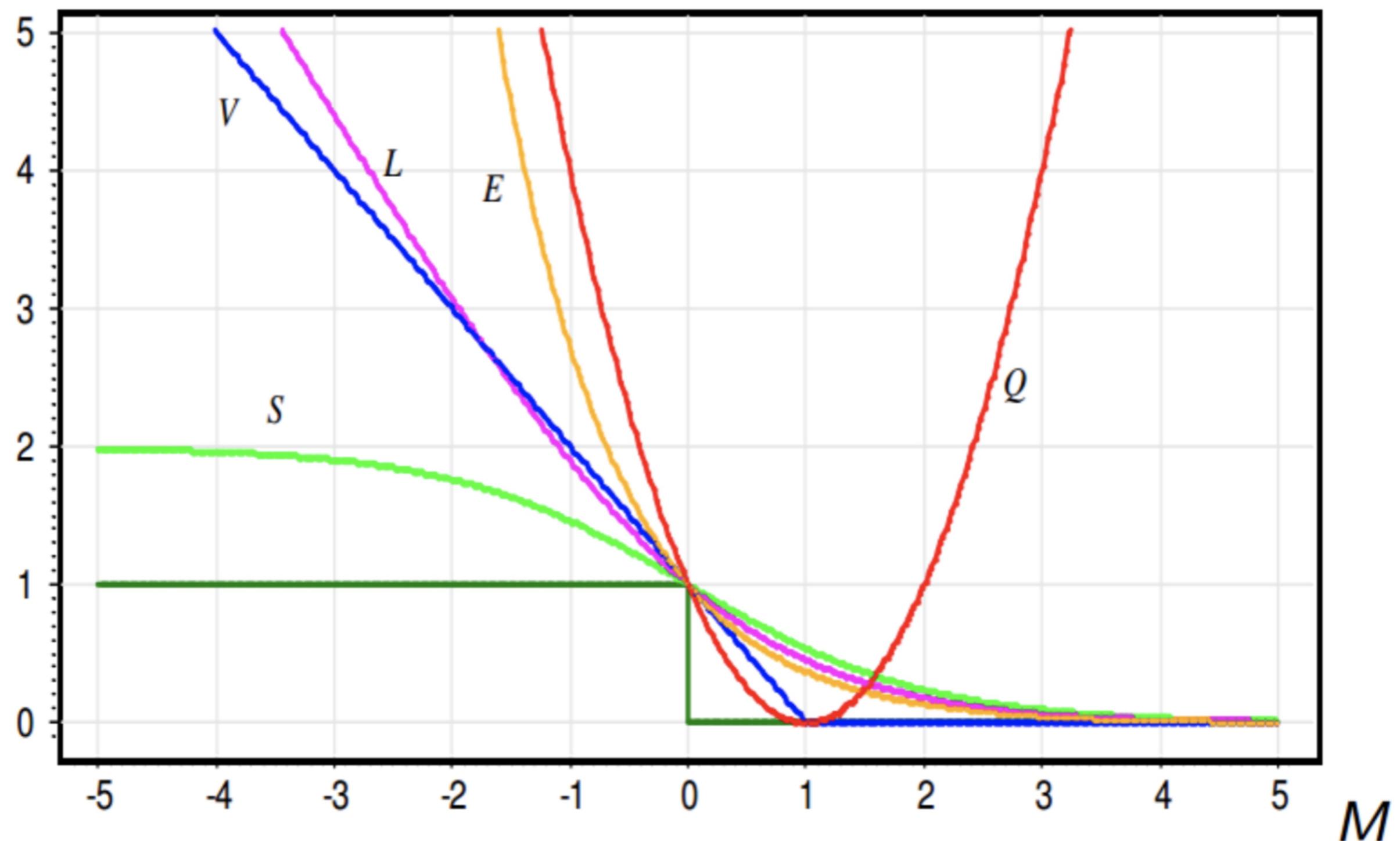
Binary classification  
Use decision stumps.



$$\hat{f}_T(x) = \sum_{t=1}^T \rho_t h_t(x) =$$



# Recap: loss functions for classification



**square loss**

$$Q(M) = (1 - M)^2$$

**hinge loss**

$$V(M) = (1 - M)_+$$

**savage loss**

$$S(M) = 2(1 + e^M)^{-1}$$

**logistic loss**

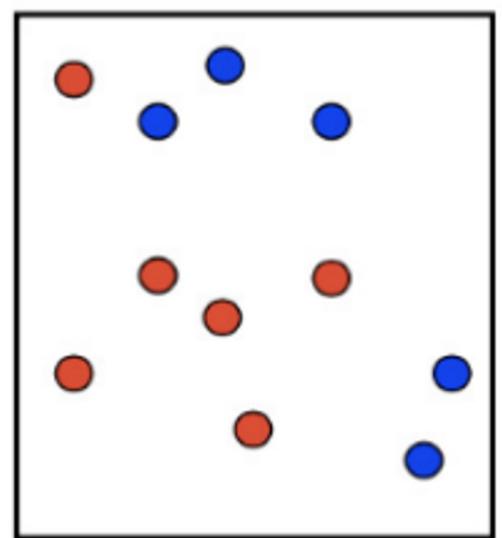
$$L(M) = \log(1 + e^{-M})$$

**exponential loss**

$$E(M) = e^{-M}$$

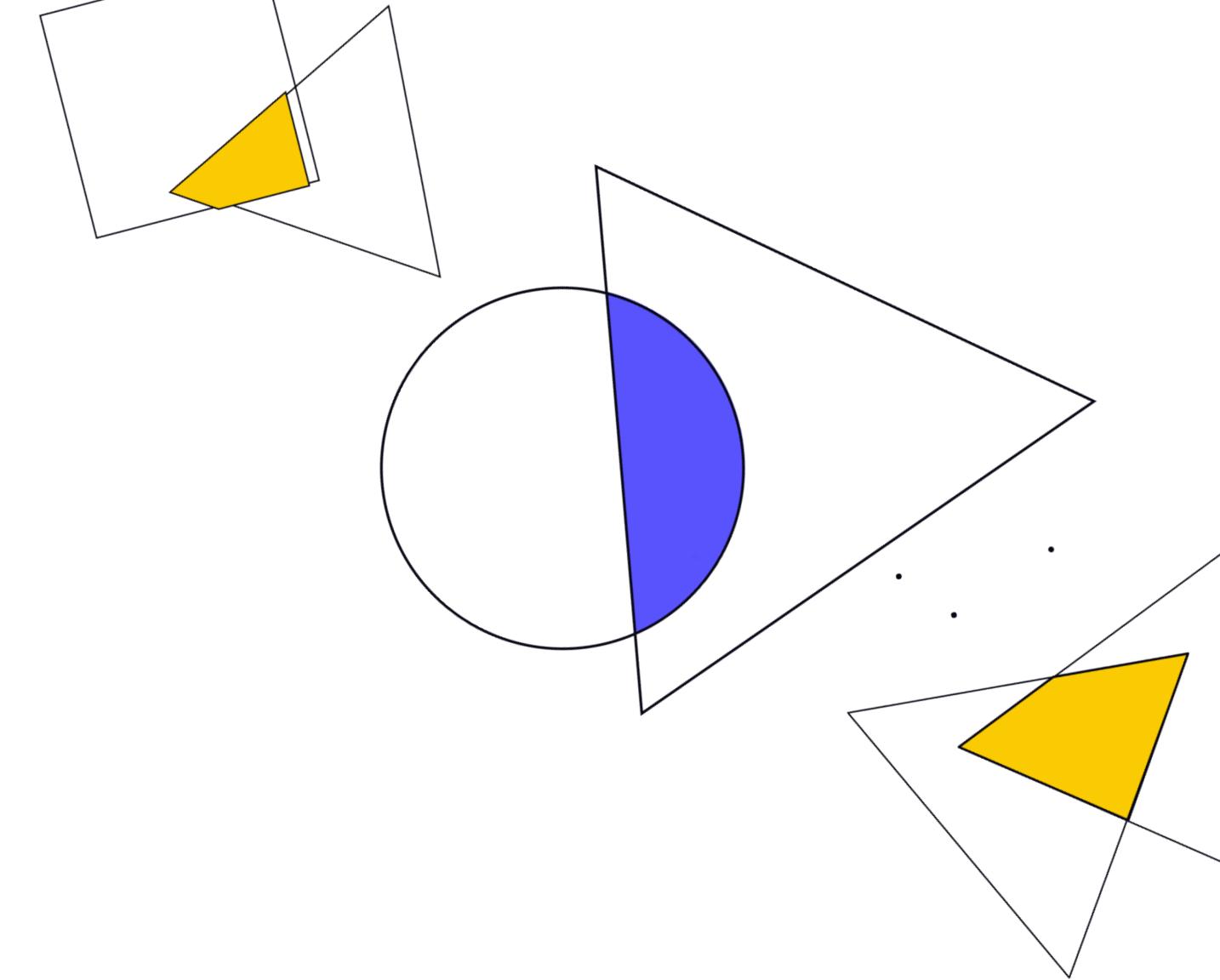
Loss functions for classification

# Boosting: AdaBoost



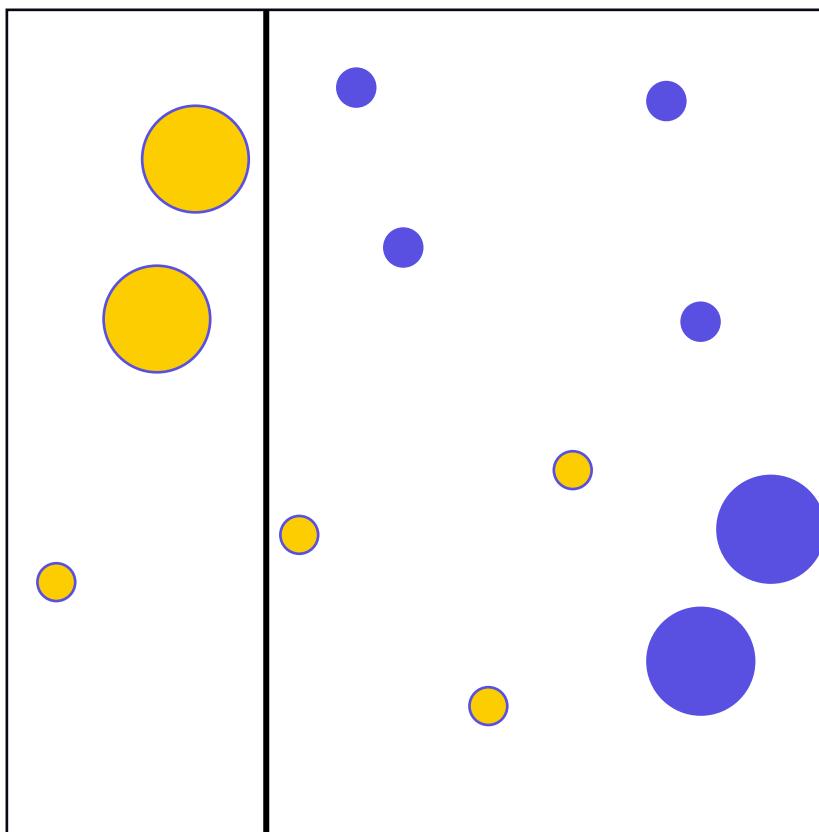
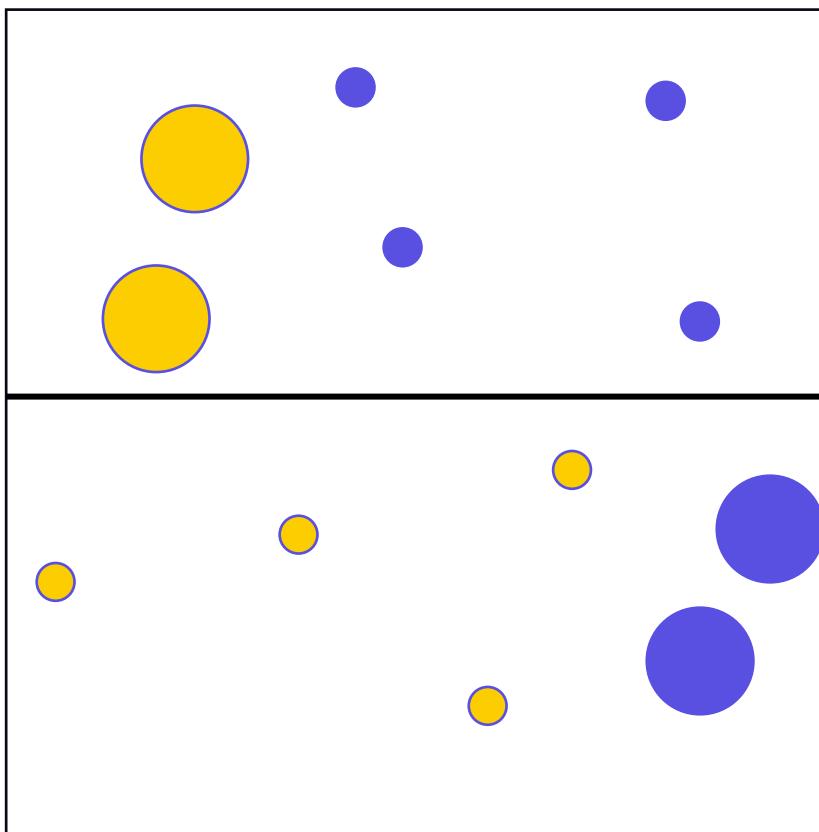
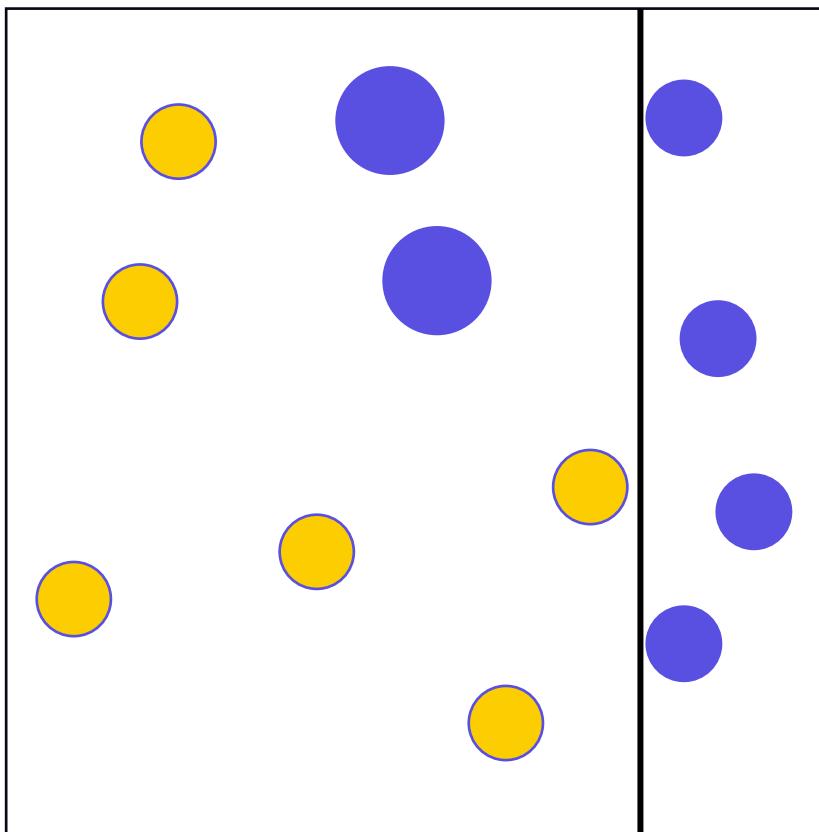
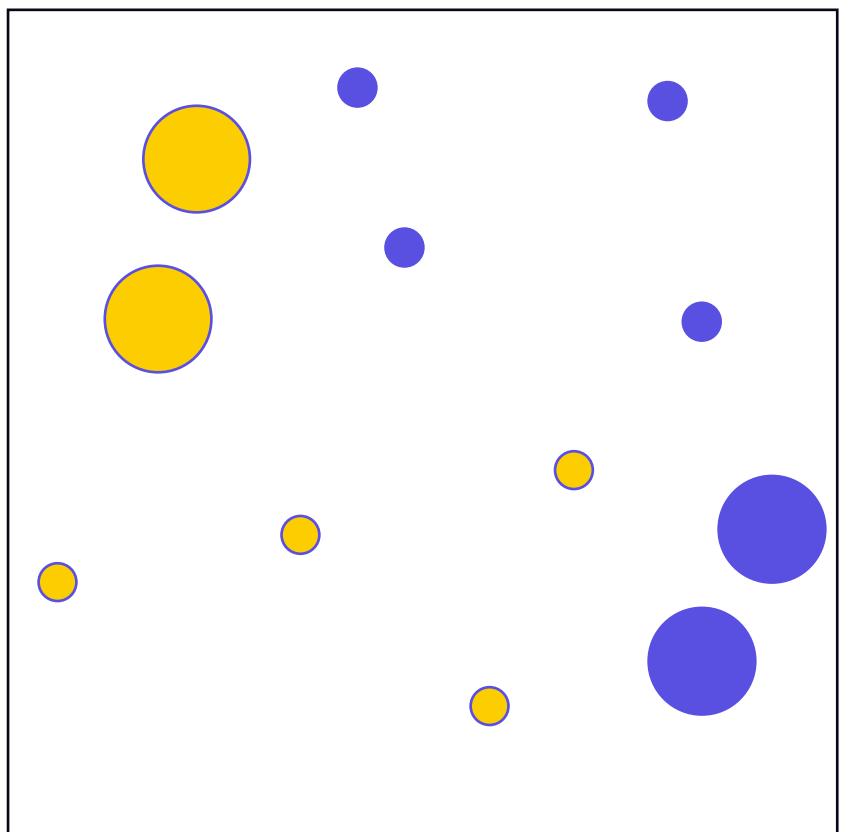
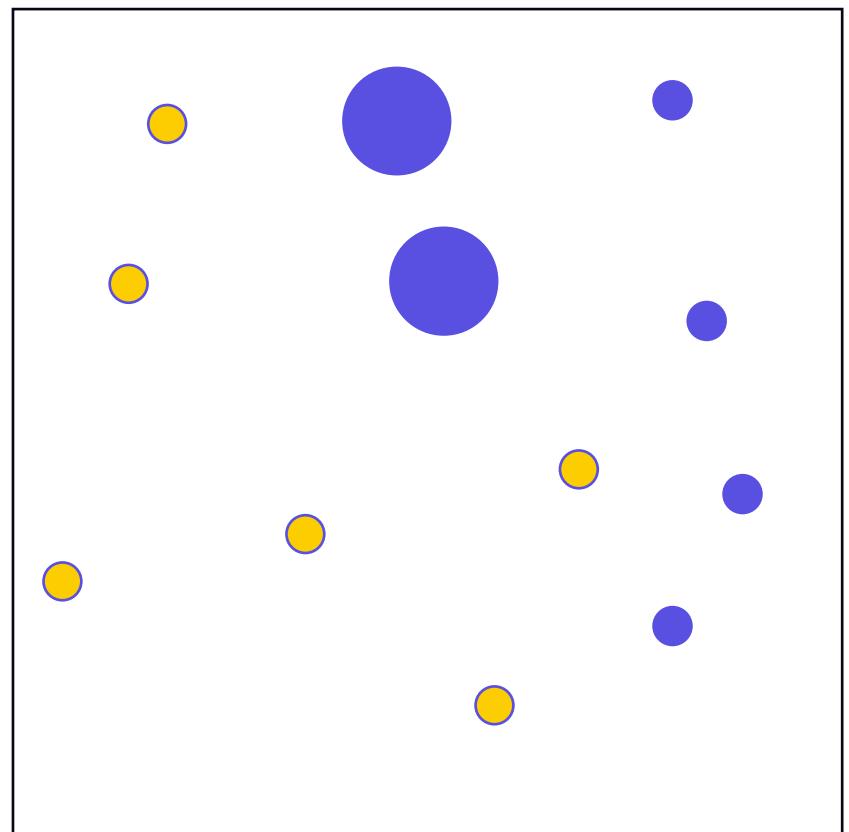
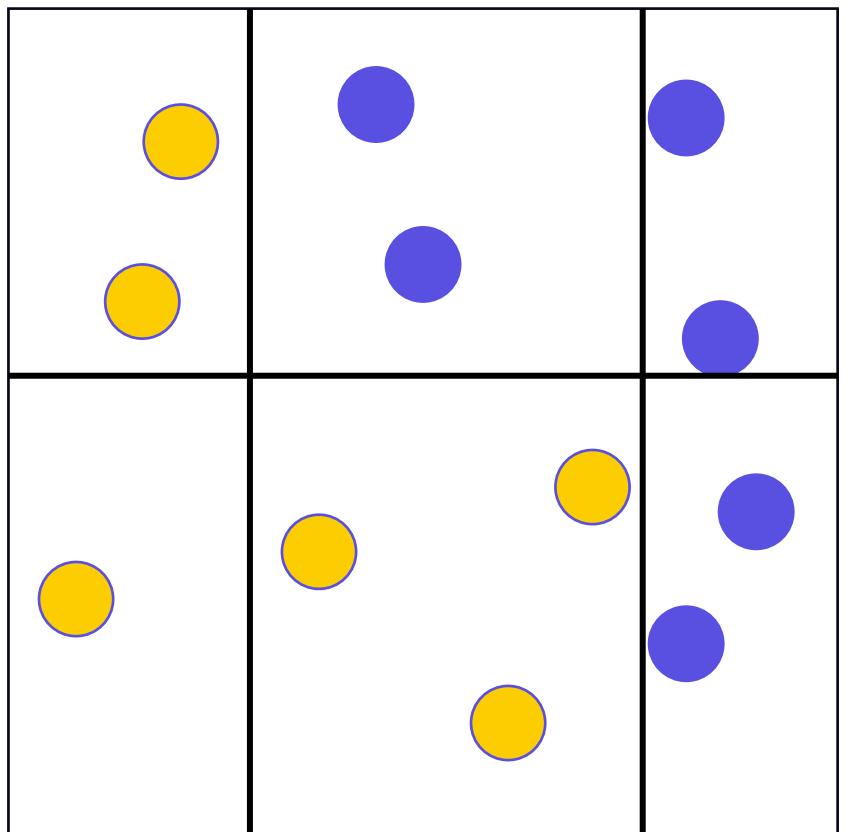
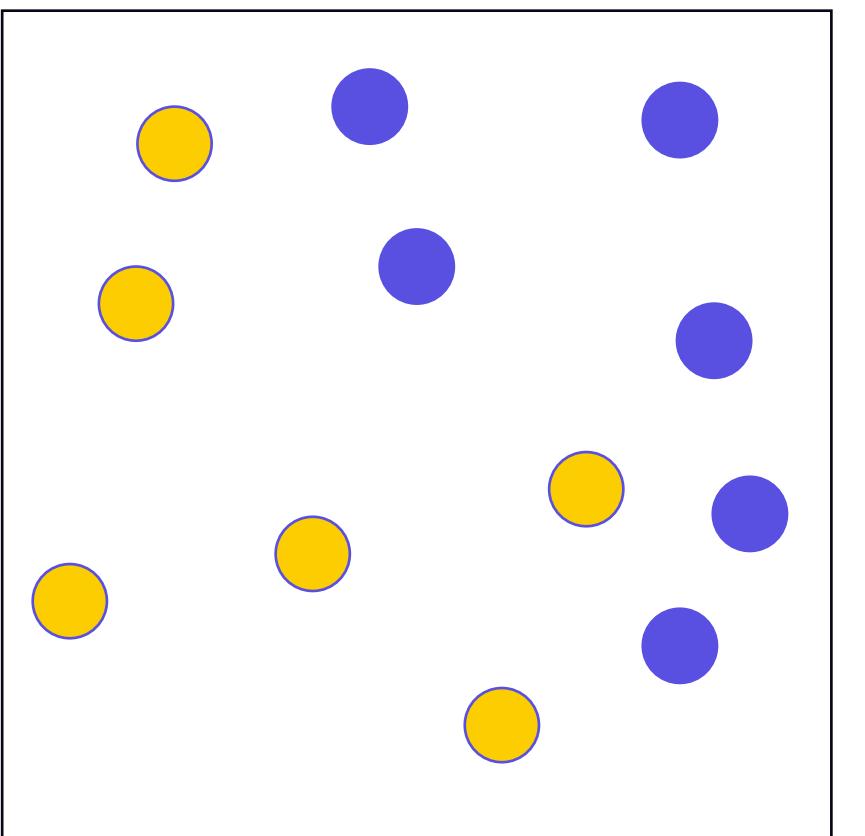
$$\hat{f}_T(\mathbf{x}) = \sum_{t=1}^T \rho_t h_t(\mathbf{x})$$

$$\begin{aligned} L(y^{(i)}, \hat{f}_T(\mathbf{x}^{(i)})) &= \exp\left(-y^{(i)} \hat{f}_T(\mathbf{x}^{(i)})\right) = \exp\left(-y^{(i)} \sum_{t=1}^T \rho_t h_t(\mathbf{x}^{(i)})\right) \\ &\quad \text{const on step T} \\ &= \exp\left(-y^{(i)} \sum_{t=1}^{T-1} \rho_t h_t(\mathbf{x}^{(i)})\right) \cdot \exp\left(-y^{(i)} \rho_T h_T(\mathbf{x}^{(i)})\right) \\ &= w^{(i)} \cdot \exp\left(-y^{(i)} \rho_T h_T(\mathbf{x}^{(i)})\right) \end{aligned}$$

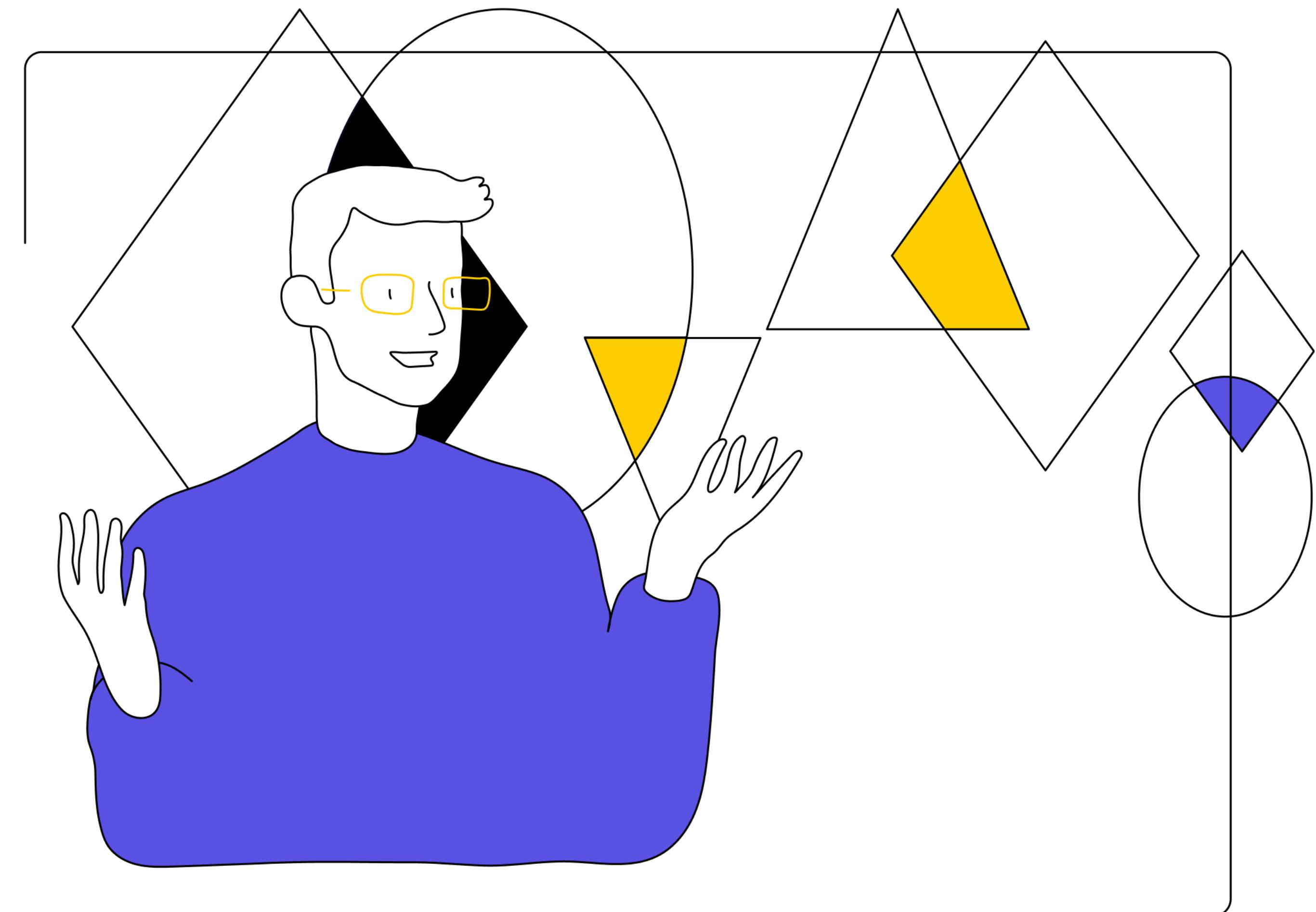


# Boosting: intuition

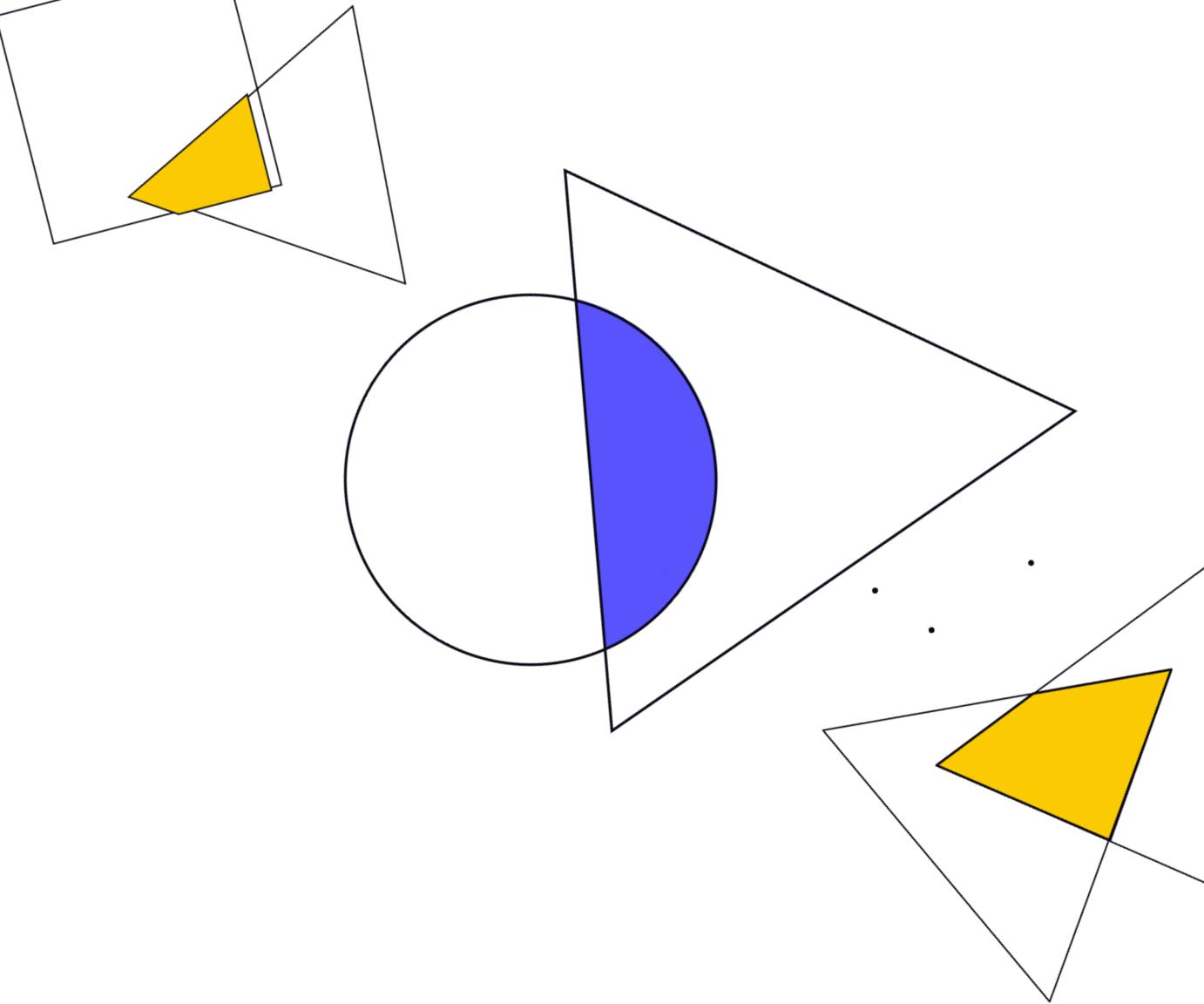
Binary classification  
Use decision stumps.



# Gradient boosting



# Gradient boosting: theory



Denote dataset  
and Loss function

$$\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$$

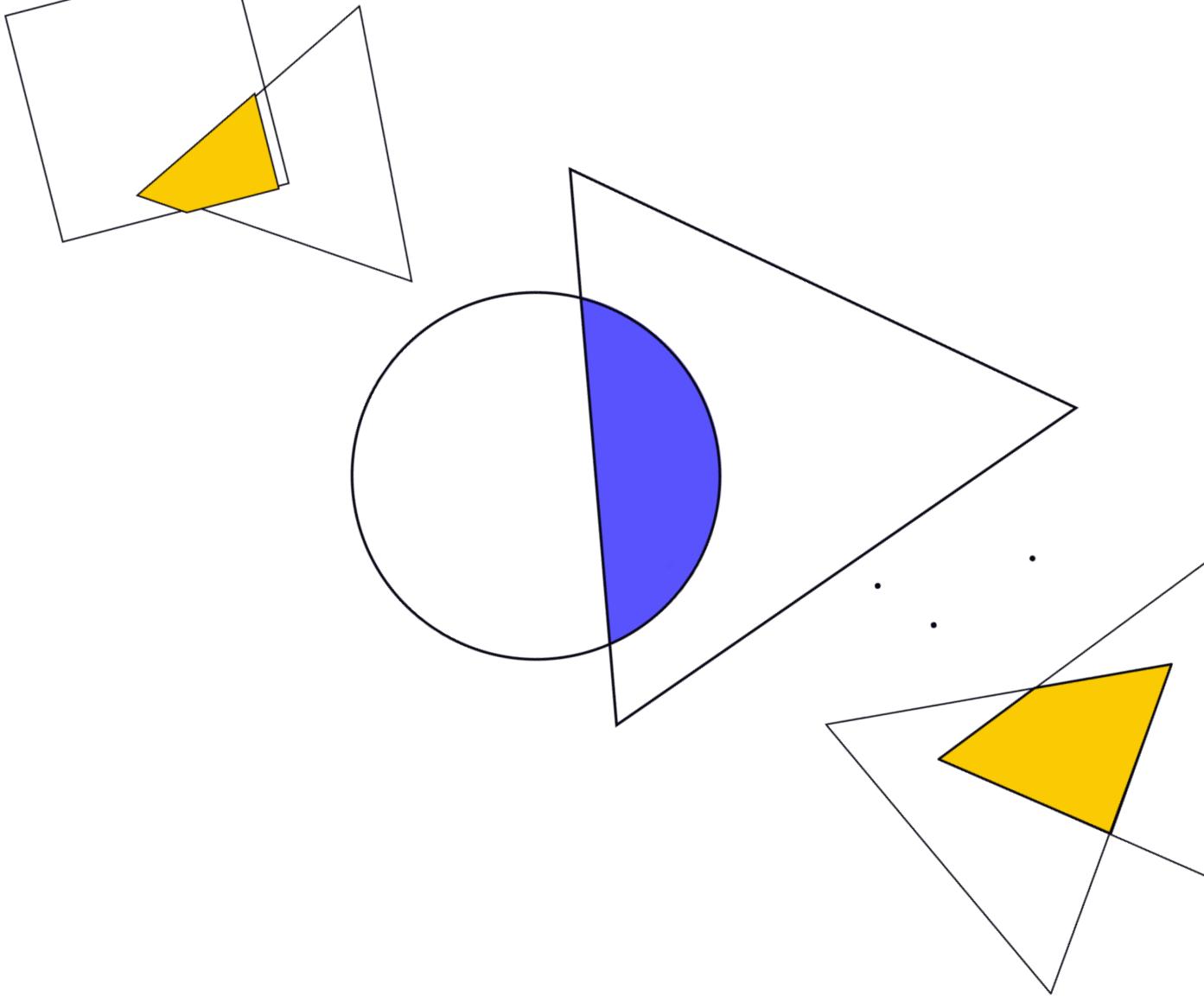
Optimal model:

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

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

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

# Gradient boosting: theory



What if we could use  
gradient descent in space  
of our models?

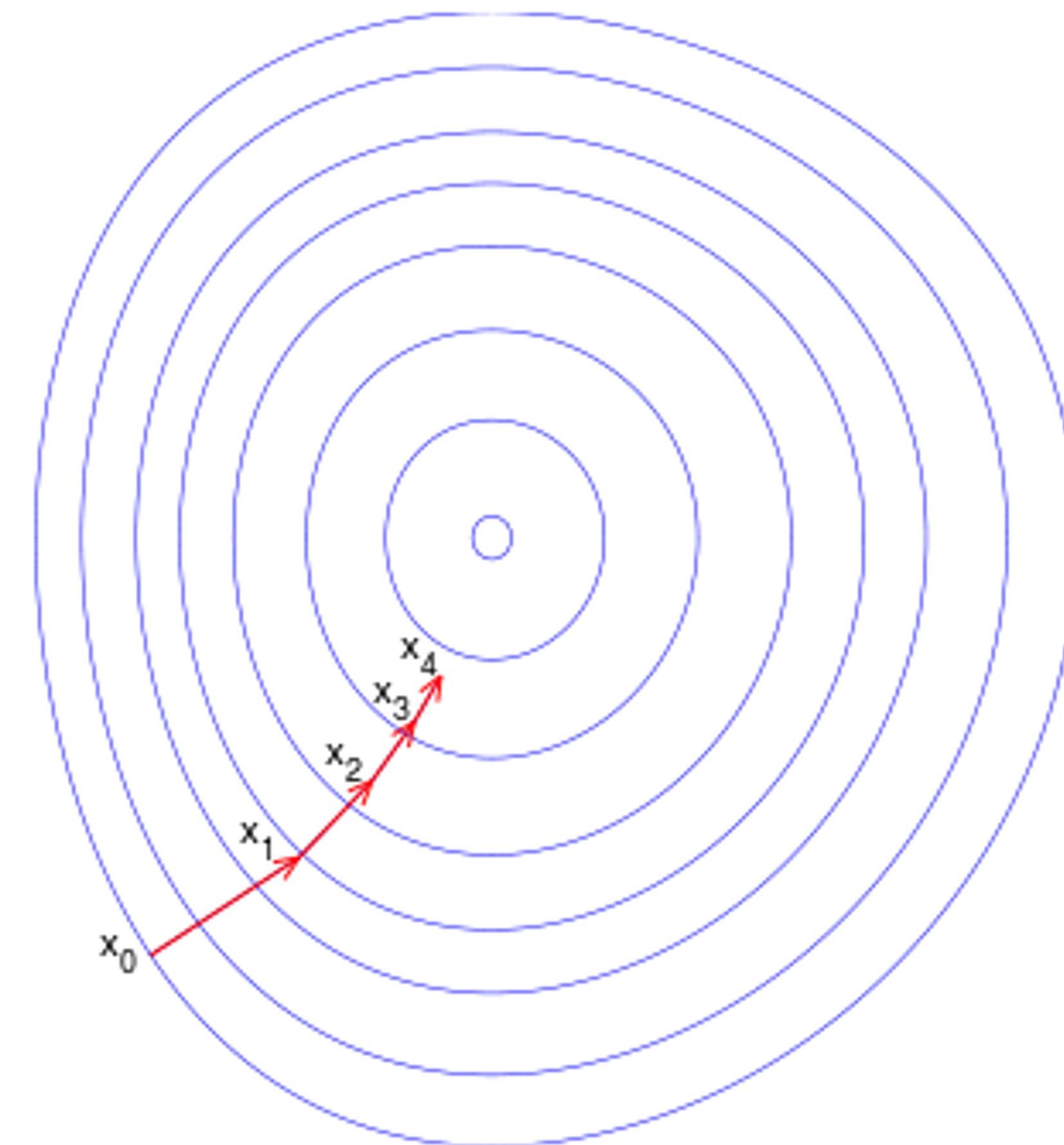
$$\hat{f}_{T-1}(\mathbf{x}) = \sum_{t=0}^{T-1} g_t(\mathbf{x}),$$

$$(\rho_T, \boldsymbol{\theta}_T) = \arg \min_{\rho, \boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}, y} \left[ L \left( y, \hat{f}_{T-1}(\mathbf{x}) + \rho_T \cdot h(\mathbf{x}, \boldsymbol{\theta}_T) \right) \right],$$

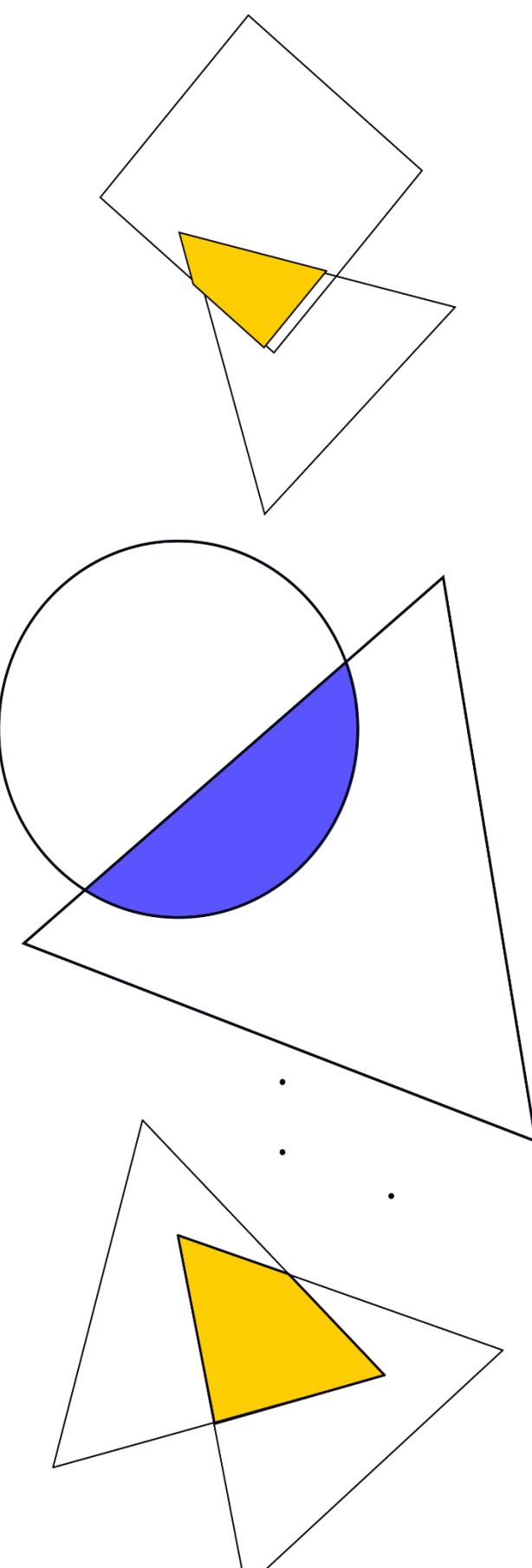
$$g_t(\mathbf{x}) = \rho_t \cdot h(\mathbf{x}, \boldsymbol{\theta}_t)$$

# Gradient boosting: theory

What if we could use  
gradient descent in space  
of our models?



# Gradient boosting: theory



$$\hat{f}_{T-1}(\mathbf{x}) = \sum_{t=0}^{T-1} g_t(\mathbf{x}),$$
$$r_t^{(i)} = - \left[ \frac{\partial L(y^{(i)}, f(\mathbf{x}^{(i)}))}{\partial f(\mathbf{x}^{(i)})} \right]_{f(\mathbf{x})=\hat{f}_T(\mathbf{x})}, \quad \text{for } i = 1, \dots, n,$$

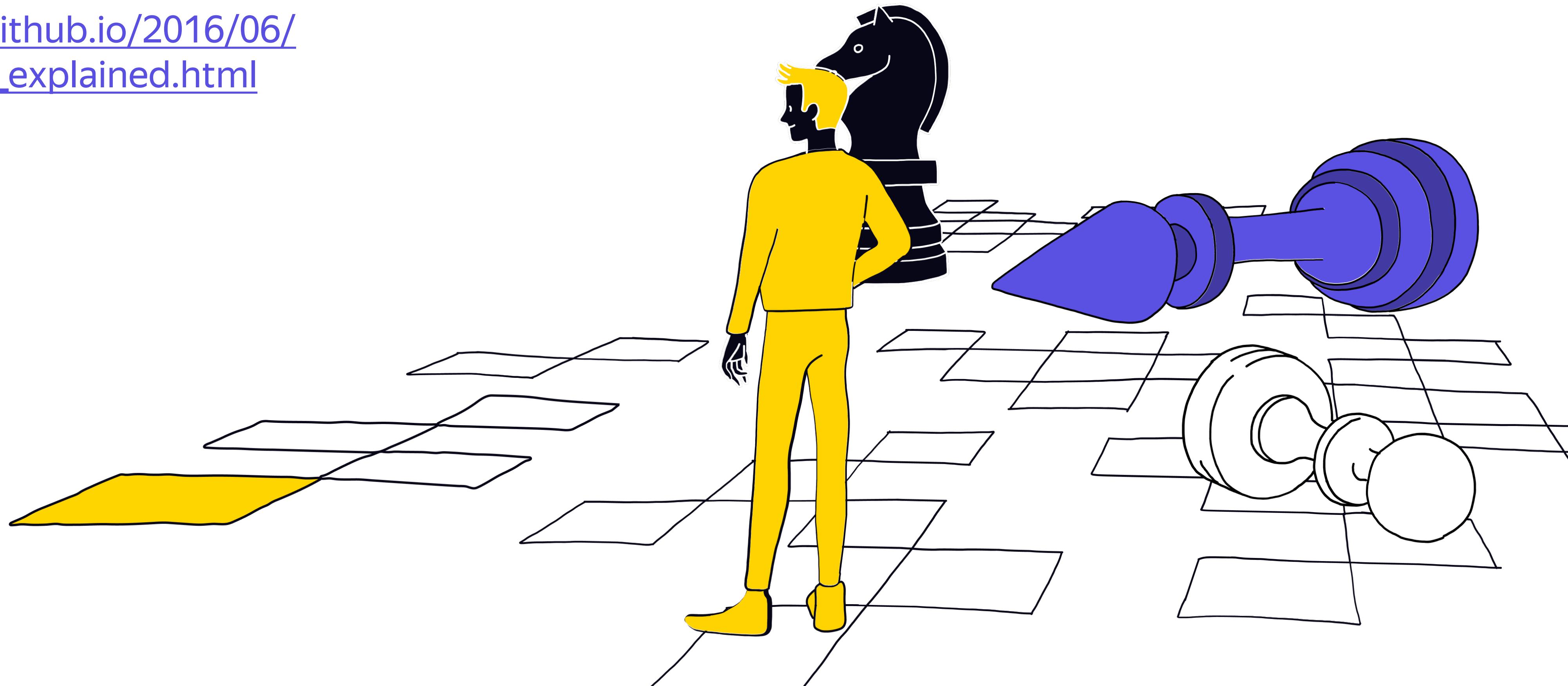
$$\boldsymbol{\theta}_T = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \left( r_t^{(i)} - h(\mathbf{x}^{(i)}, \boldsymbol{\theta}) \right)^2,$$

$$\rho_t = \arg \min_{\rho} \sum_{i=1}^n L \left( y^{(i)}, \hat{f}_{t-1}(\mathbf{x}^{(i)}) + \rho \cdot h(\mathbf{x}^{(i)}, \boldsymbol{\theta}_T) \right)$$

# Gradient boosting: beautiful demo

Great demo:

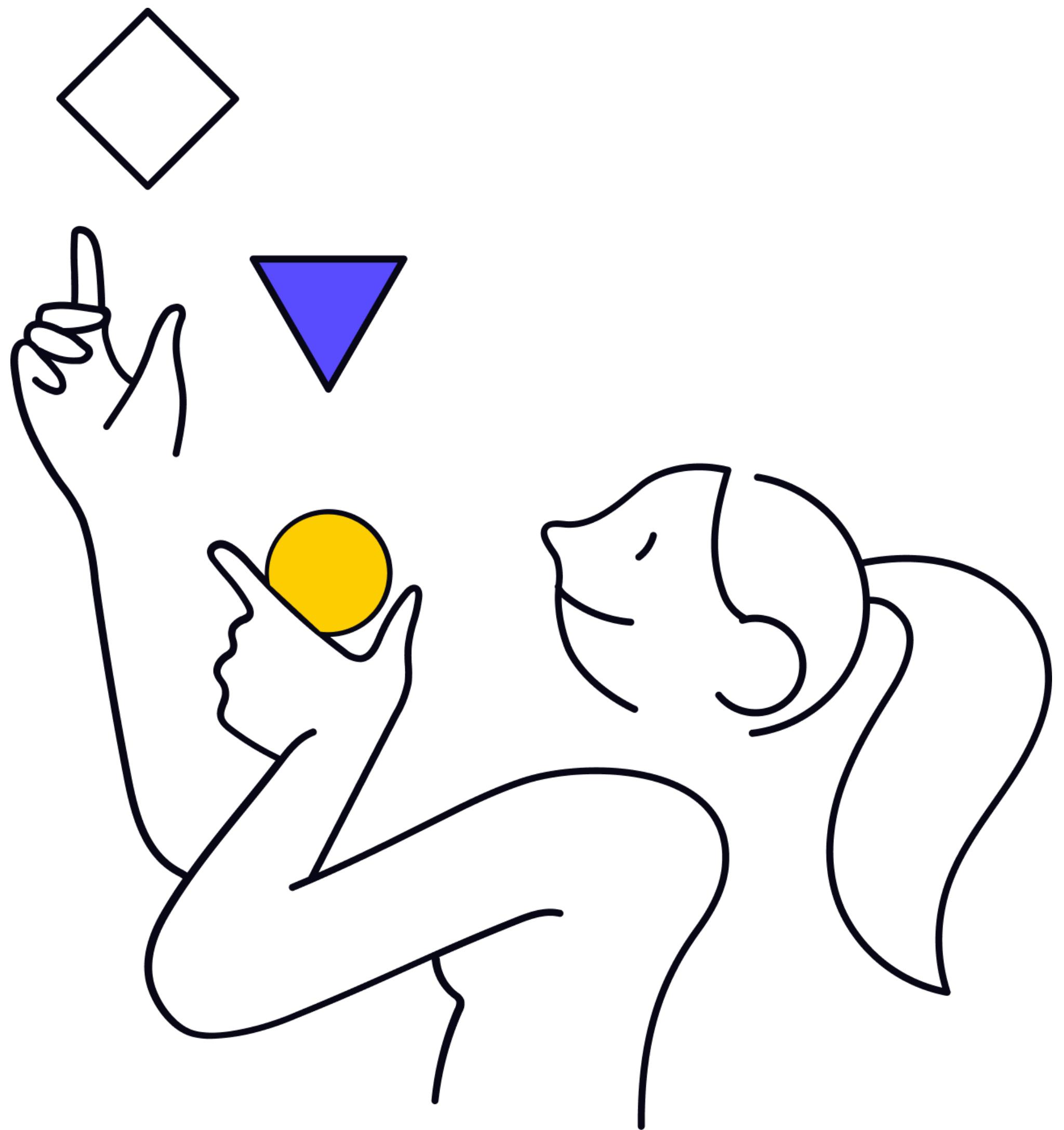
[http://arogozhnikov.github.io/2016/06/24/gradient\\_boosting\\_explained.html](http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html)



# Gradient boosting

What we need:

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

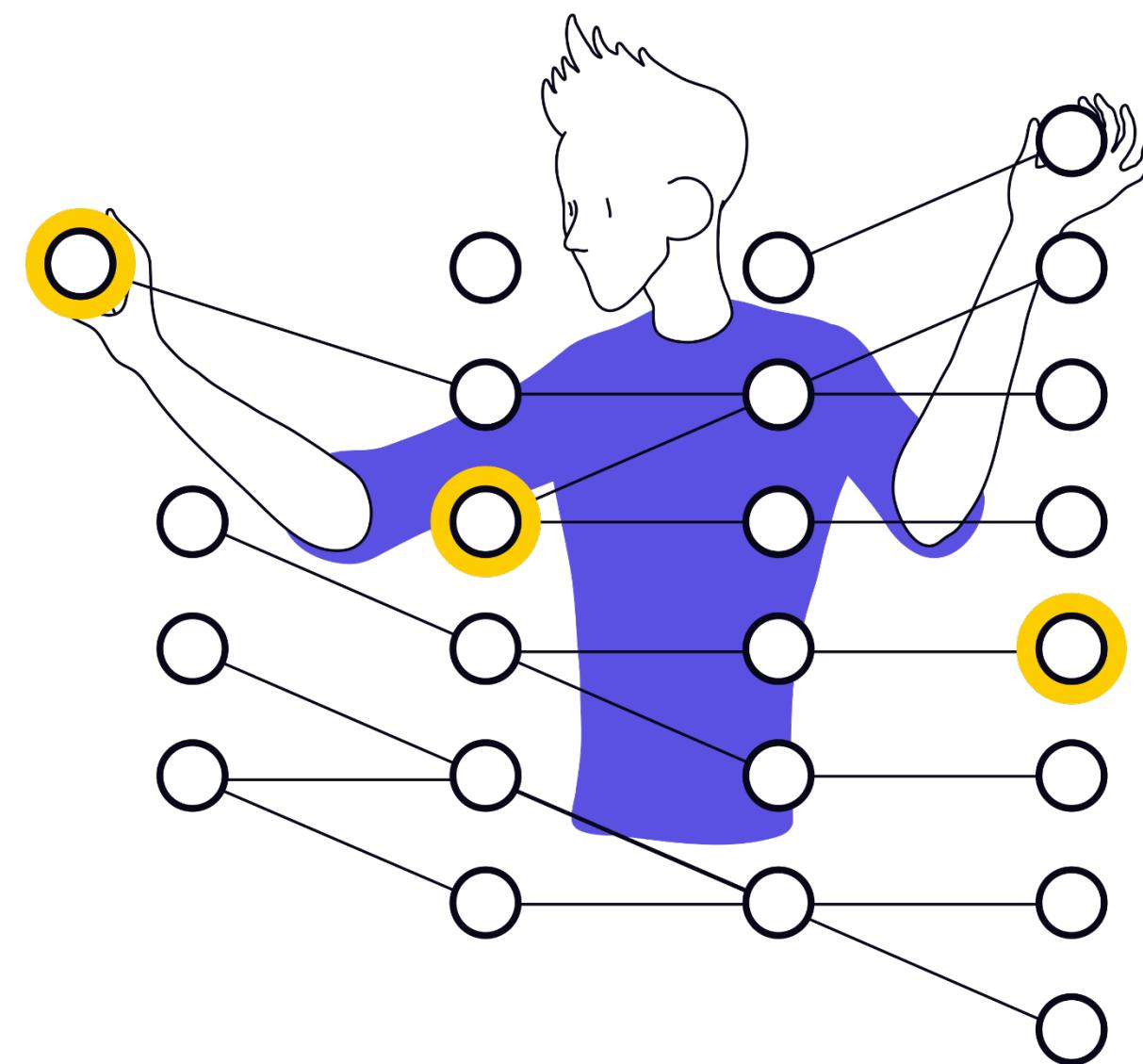


# Gradient boosting: example

What we need:

- Data: toy dataset
- Loss function: MSE
- Family of algorithms: decision trees with depth 2
- Number of iterations M = 3
- Initial value: just mean value

$$y = \cos(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \frac{1}{5}), x \in [-5, 5]$$

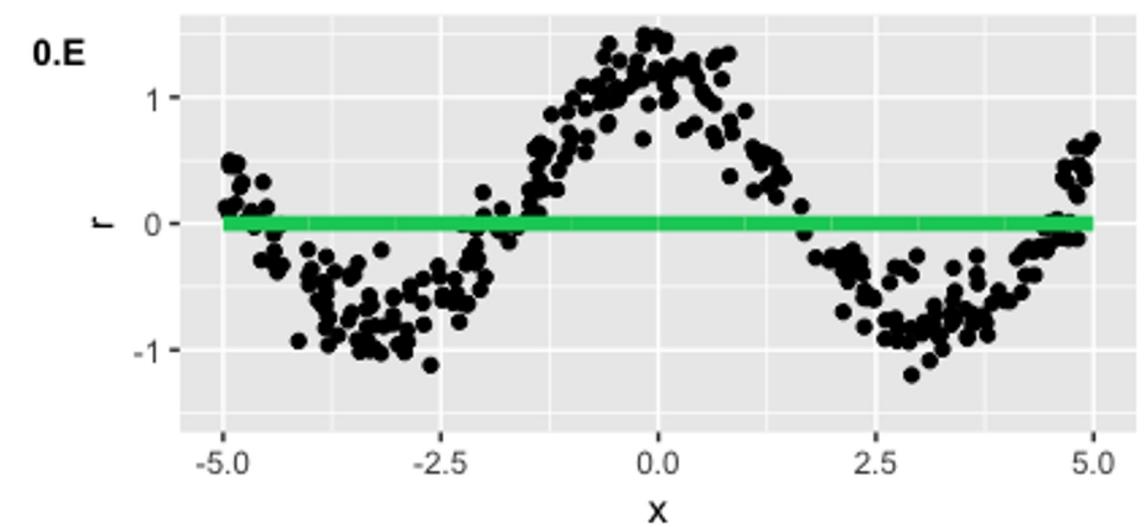
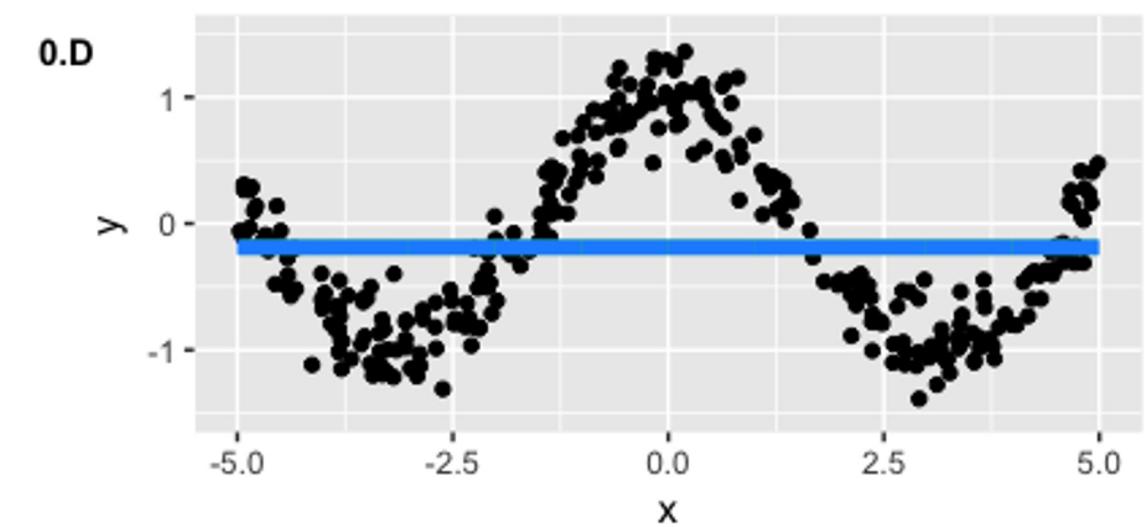


Example by ODS; source:

<https://habr.com/ru/company/ods/blog/327250/>

# Gradient boosting: example

Left: full ensemble on each step.  
Right: additional tree decisions.

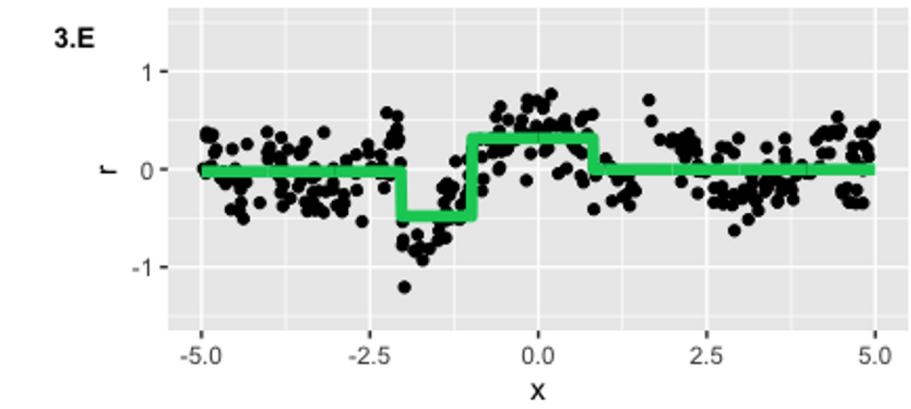
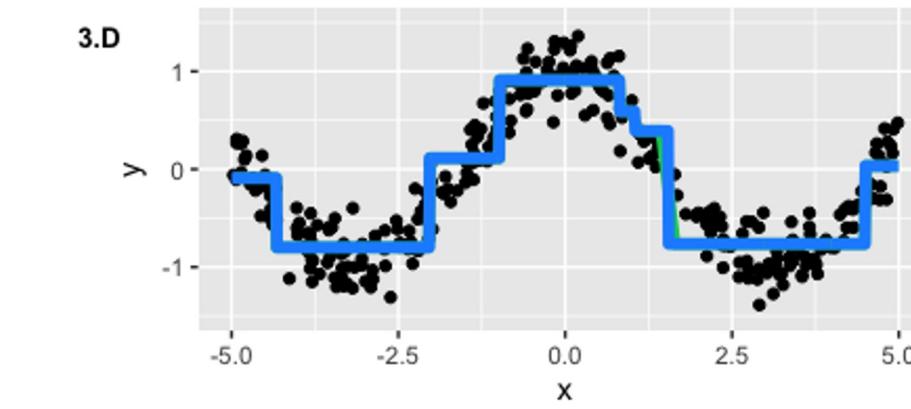
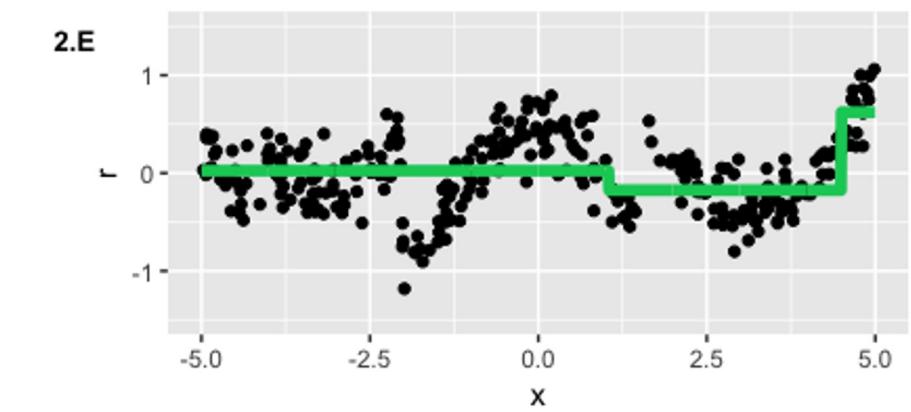
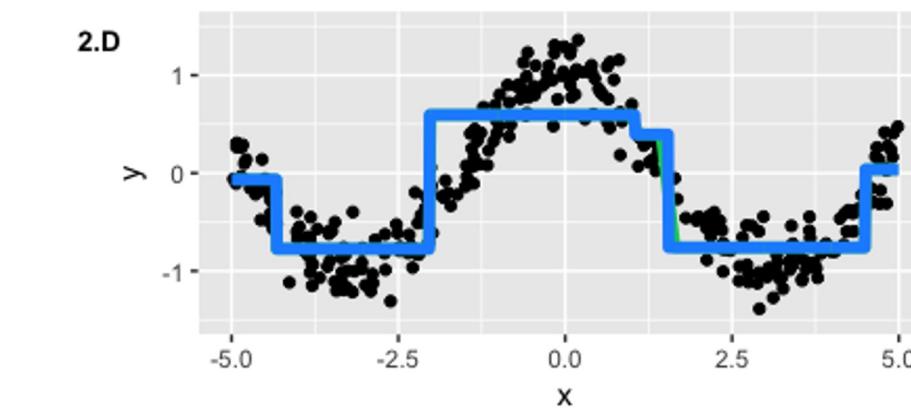
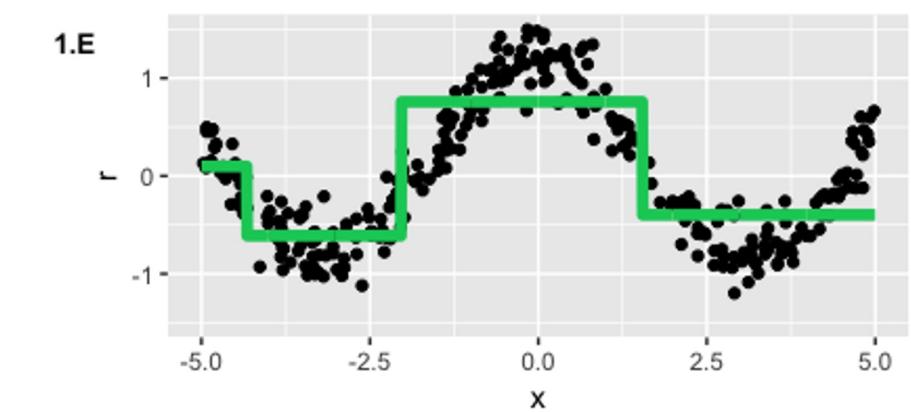
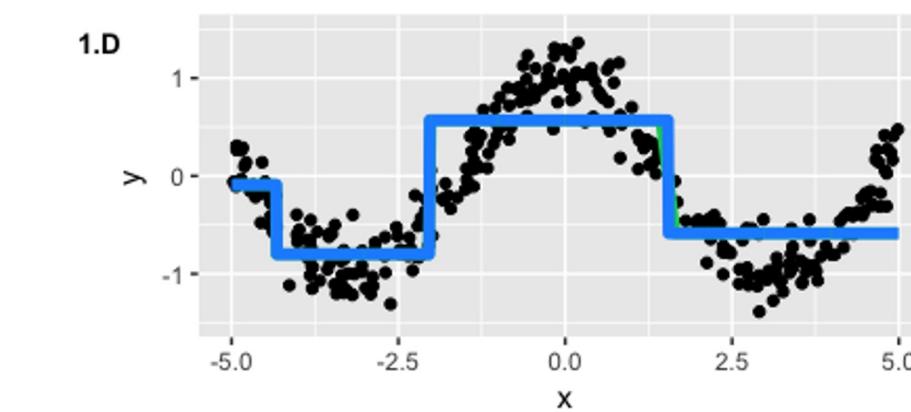
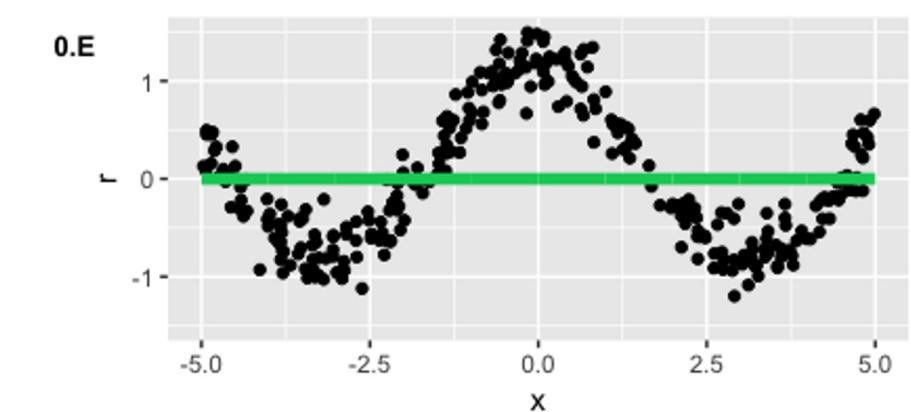
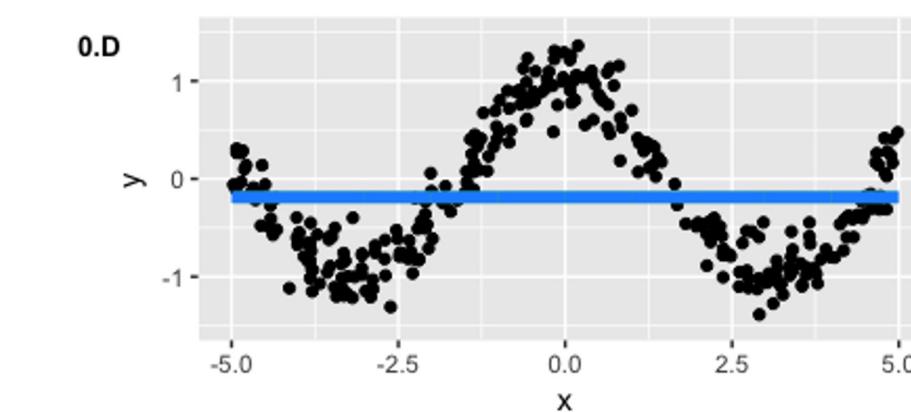


Example by ODS; source:

<https://habr.com/ru/company/ods/blog/327250/>

# Gradient boosting: example

Left: full ensemble on each step.  
Right: additional tree decisions.



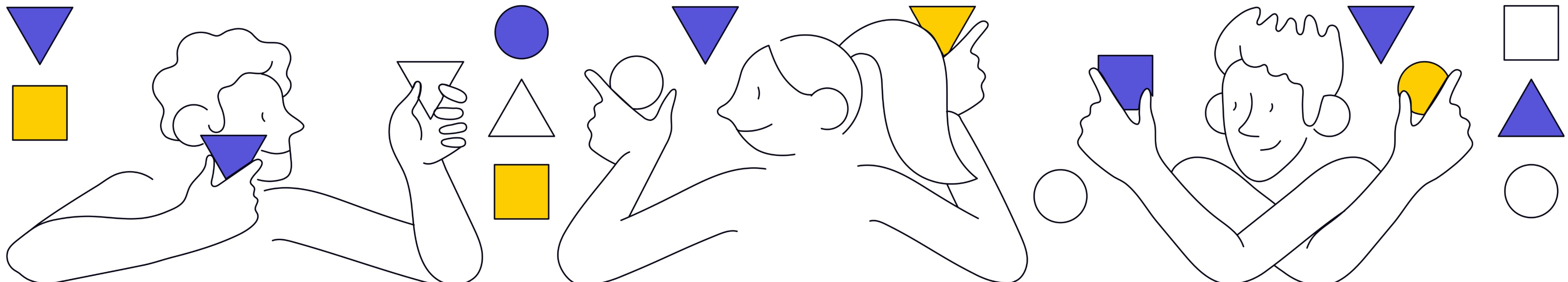
Example by ODS; source:

<https://habr.com/ru/company/ods/blog/327250/>

# Technical side: training in parallel

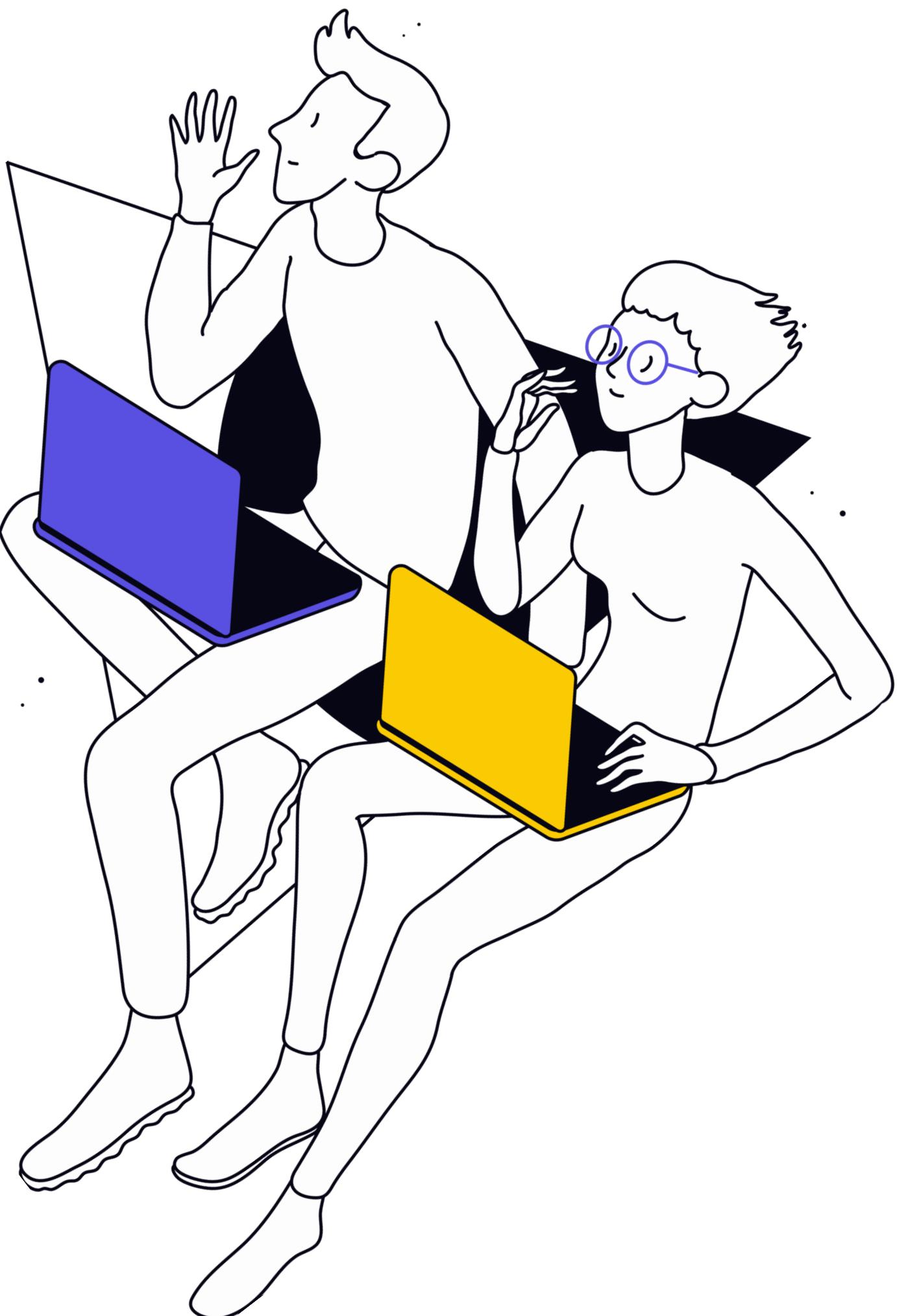
Which of the ensembling methods could be parallelized?

- Random Forest: parallel on the forest level (all trees are independent)
- Gradient boosting: parallel on one tree level



# Revise

- 01** Boosting intuitions
- 02** Gradient boosting
- 03** Blending
- 04** Stacking



# Thanks for attention!

Questions?

