

Ensembling methods

Stacking, bagging, boosting

Data Analytics and Mining, 2024

Majid Sohrabi

National Research University Higher School of Economics



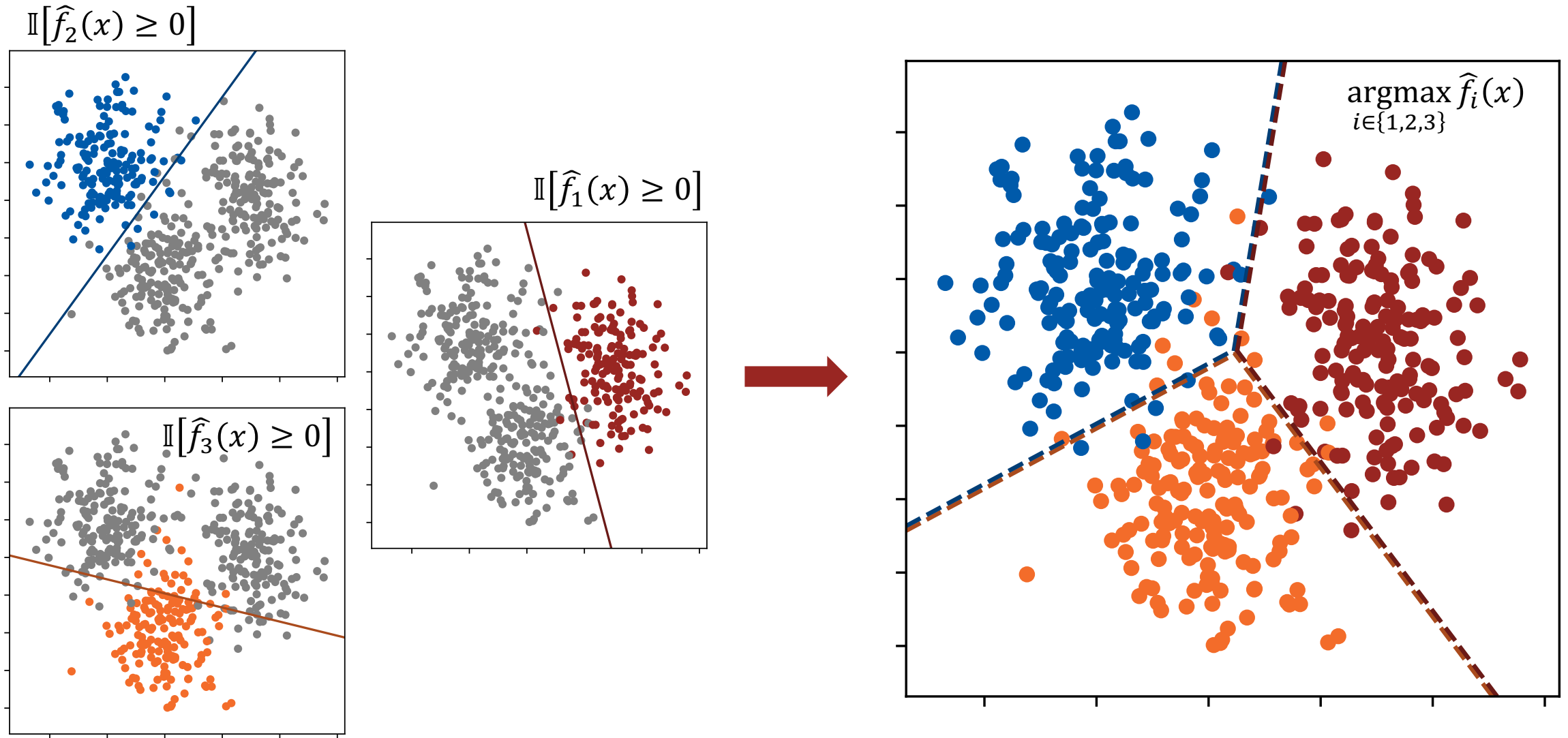
MMCP

November 29, 2024

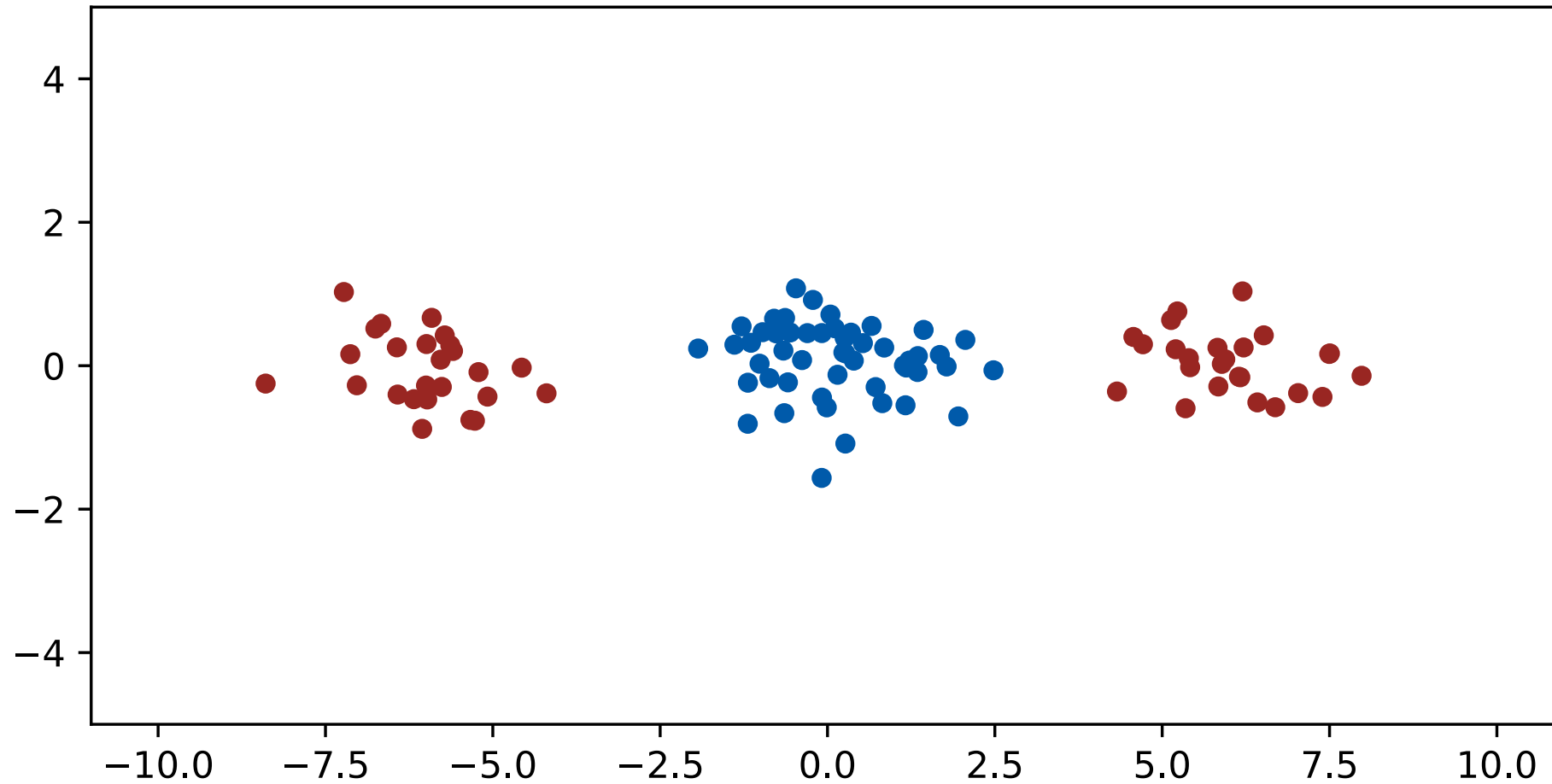
Why ensembles?



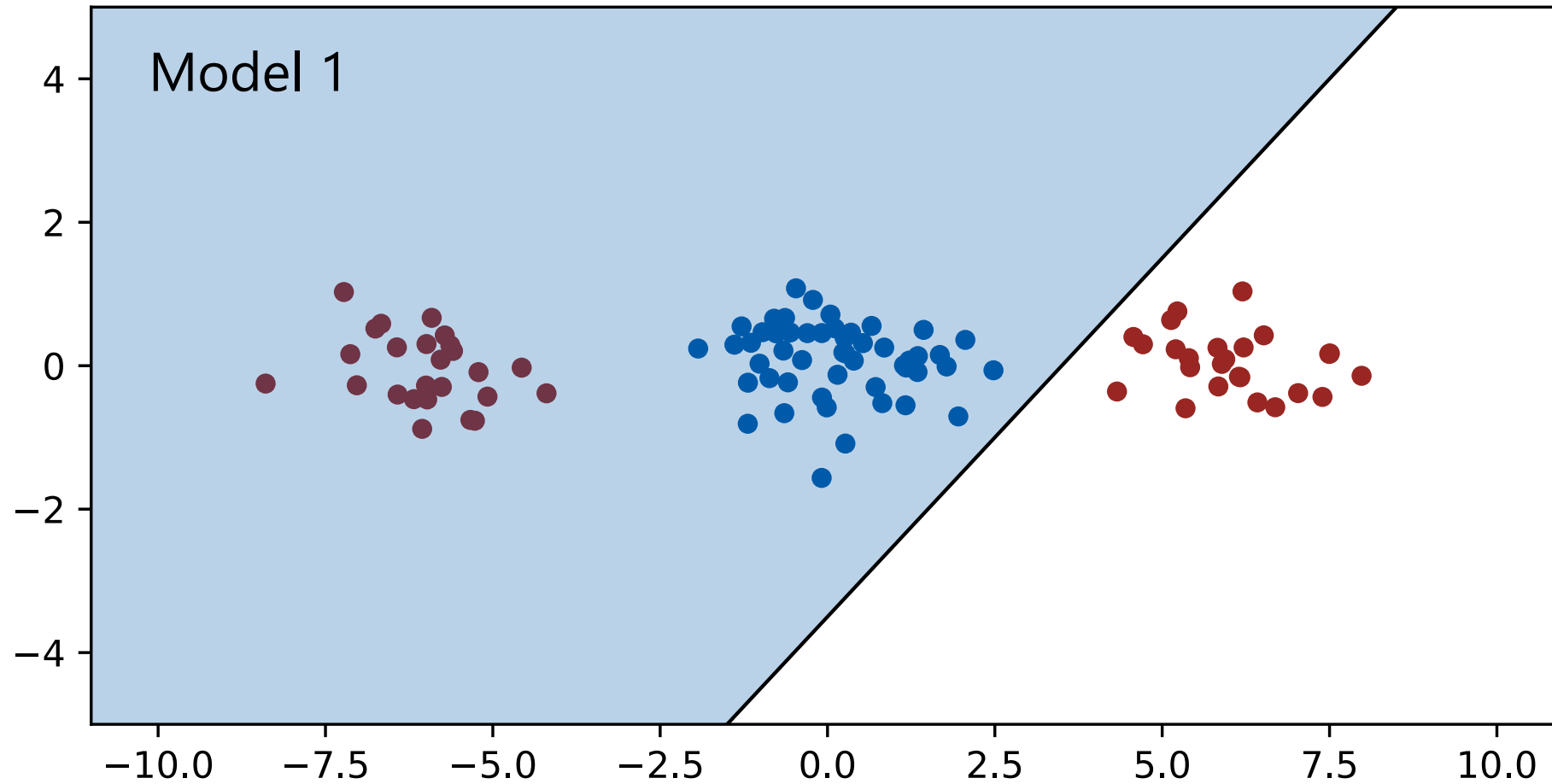
One-vs-rest classification scheme



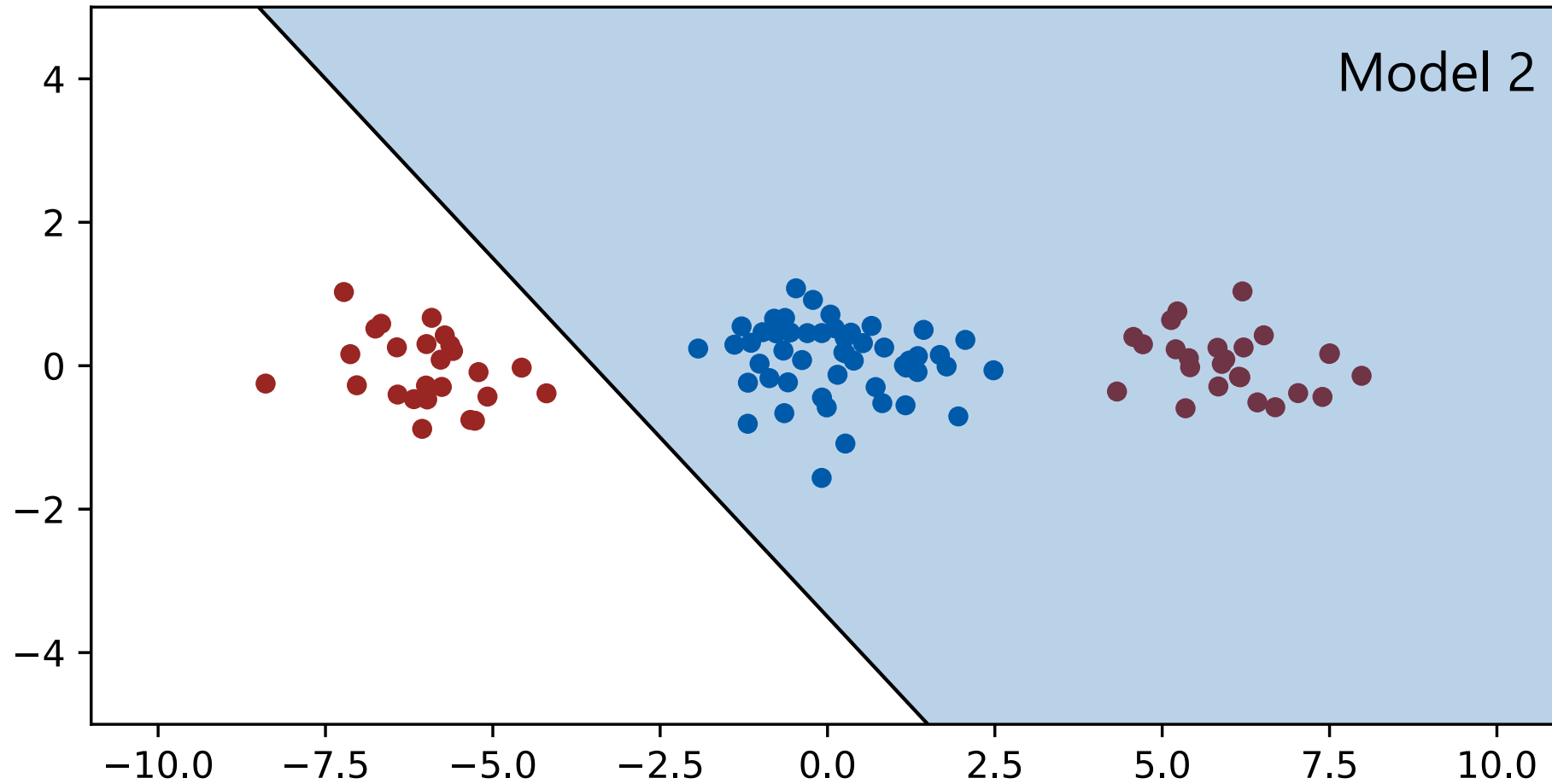
Increasing complexity of simple underfitting models



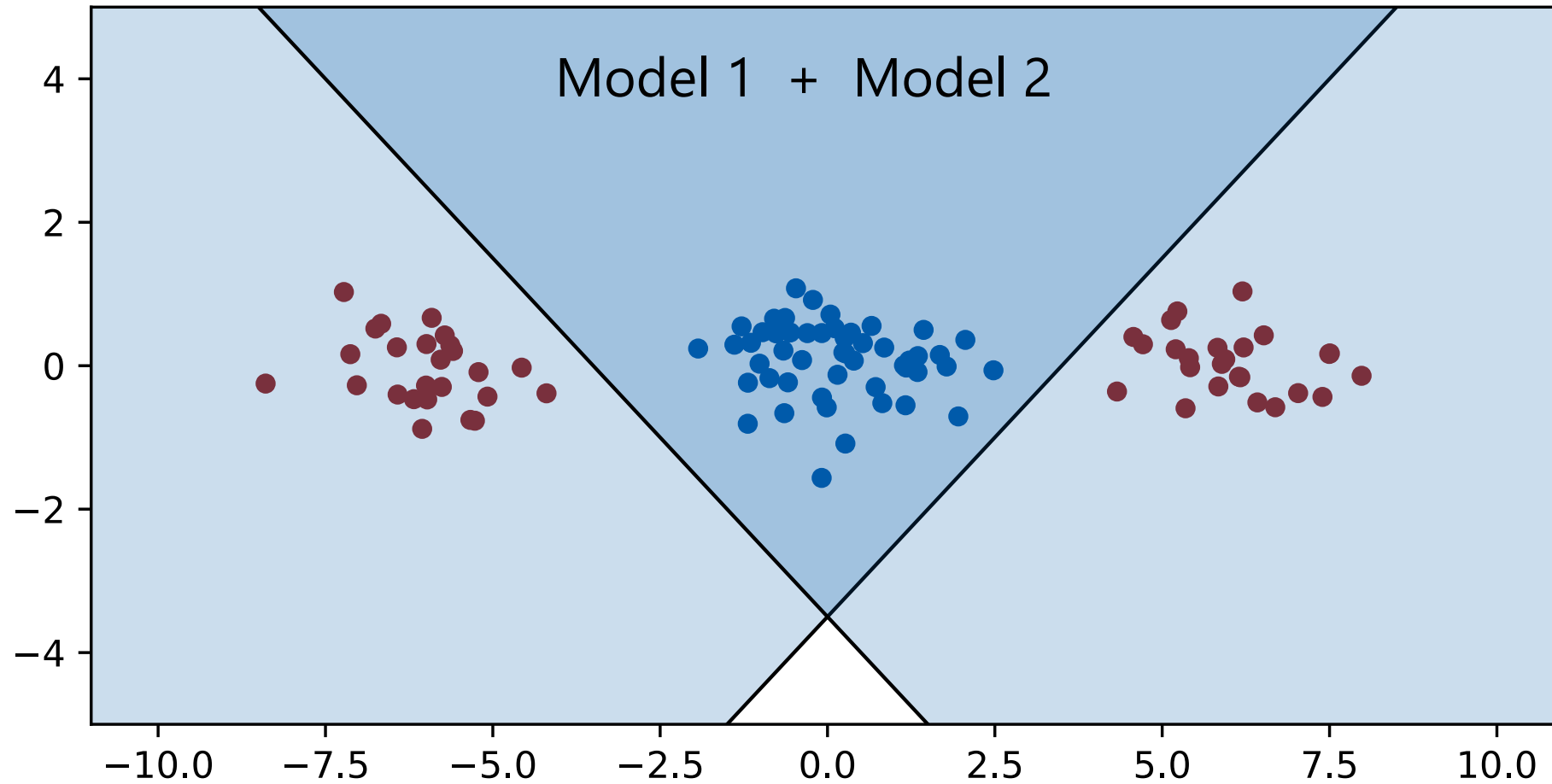
Increasing complexity of simple underfitting models



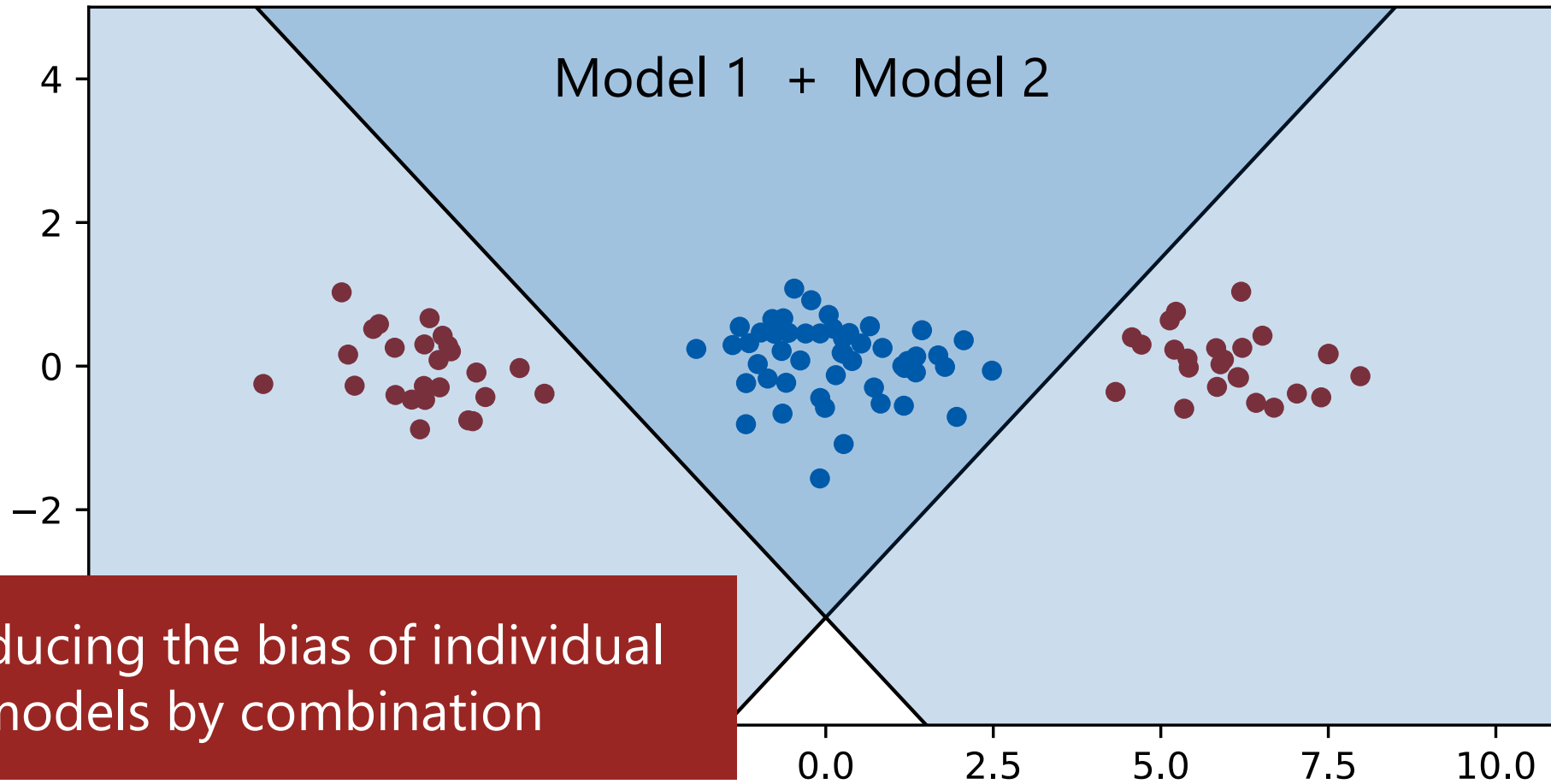
Increasing complexity of simple underfitting models



Increasing complexity of simple underfitting models



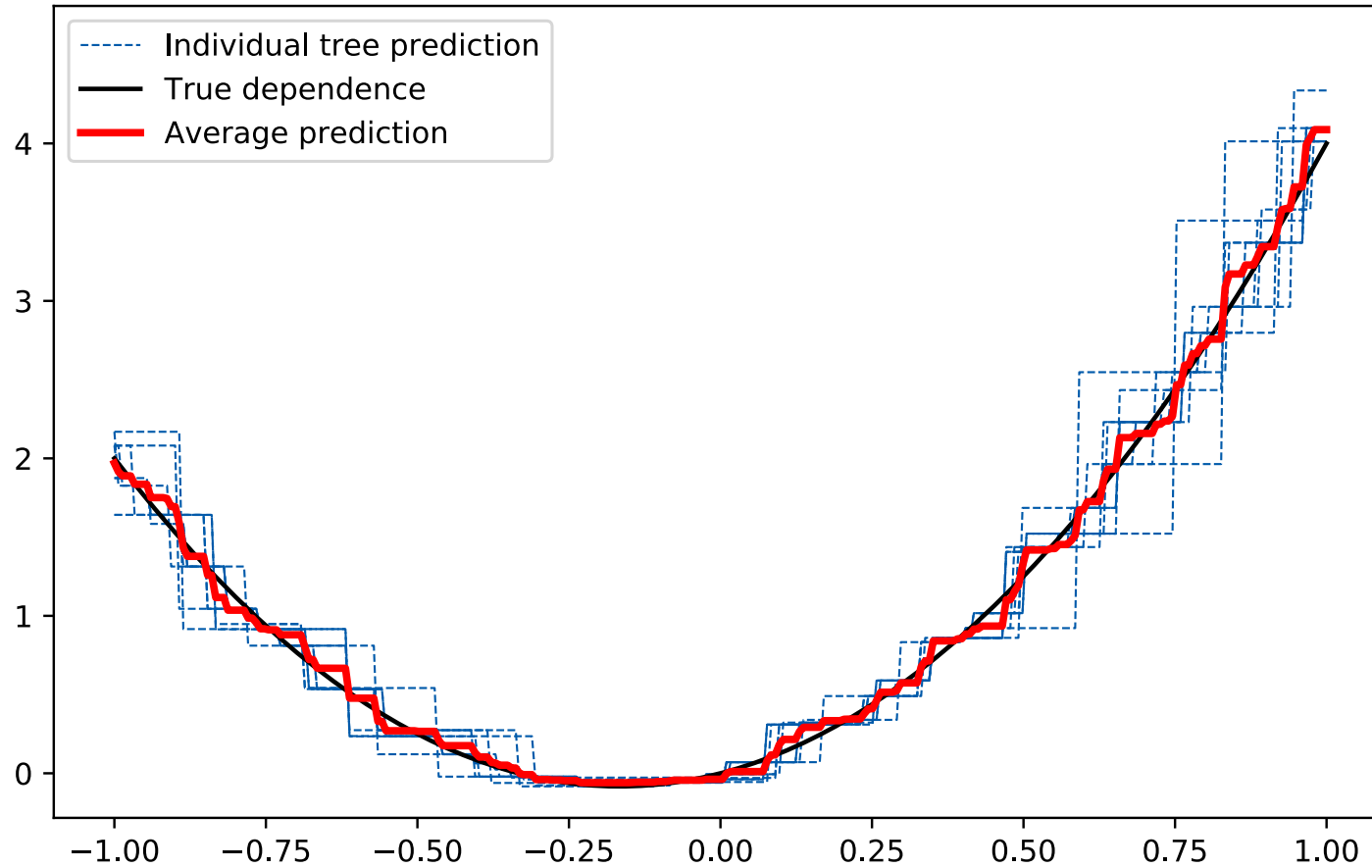
Increasing complexity of simple underfitting models



I.e. reducing the bias of individual models by combination

May increase variance

Reducing variance by averaging



Take a set of overfitting (high-variance) models $\hat{f}_1(x), \dots, \hat{f}_M(x)$

- E.g. decision trees trained on different training sets

For regression – average their predictions:

$$F(x) = \frac{1}{M} \sum_{i=1}^M \hat{f}_i(x)$$

For classification – majority voting

Ambiguity decomposition

Consider the following ensemble:

$$F(x) = \sum_i w_i f_i(x), \quad w_i \geq 0, \quad \sum_i w_i = 1$$

Then it's easy to show that:

$$\underbrace{(F(x) - y)^2}_{\text{Ensemble error}} = \underbrace{\sum_i w_i (f_i(x) - y)^2}_{\text{Base learner error}} - \underbrace{\sum_i w_i (f_i(x) - F(x))^2}_{\text{Ambiguity}}$$

Disagreement reduces the error!

[Ambiguity decomposition proof]

$$F(x) = \sum_i w_i f_i(x), \quad w_i \geq 0, \quad \sum_i w_i = 1$$

$$(F(x) - y)^2 = \left(\sum_i w_i f_i(x) - y \right)^2 = \sum_i w_i f_i(x) \cdot F(x) - 2 \sum_i w_i f_i(x) \cdot y + \sum_i w_i y^2$$

$$= \sum_i w_i [f_i(x) \cdot F(x) - 2f_i(x) \cdot y + y^2]$$

$$= \sum_i w_i [f_i^2(x) - 2f_i(x) \cdot y + y^2 - f_i^2(x) + 2f_i(x) \cdot F(x) - F^2(x) + F^2(x) - f_i(x) \cdot F(x)]$$

$$= \sum_i w_i \left[(f_i(x) - y)^2 - (f_i(x) - F(x))^2 + F^2(x) - f_i(x) \cdot F(x) \right]$$

$$= \sum_i w_i (f_i(x) - y)^2 - \sum_i w_i (f_i(x) - F(x))^2$$

Majority voting classifier error

Assume we have M classifiers with prediction error probability $p < \frac{1}{2}$ each

Suppose the guesses are wrong or correct independently of each other

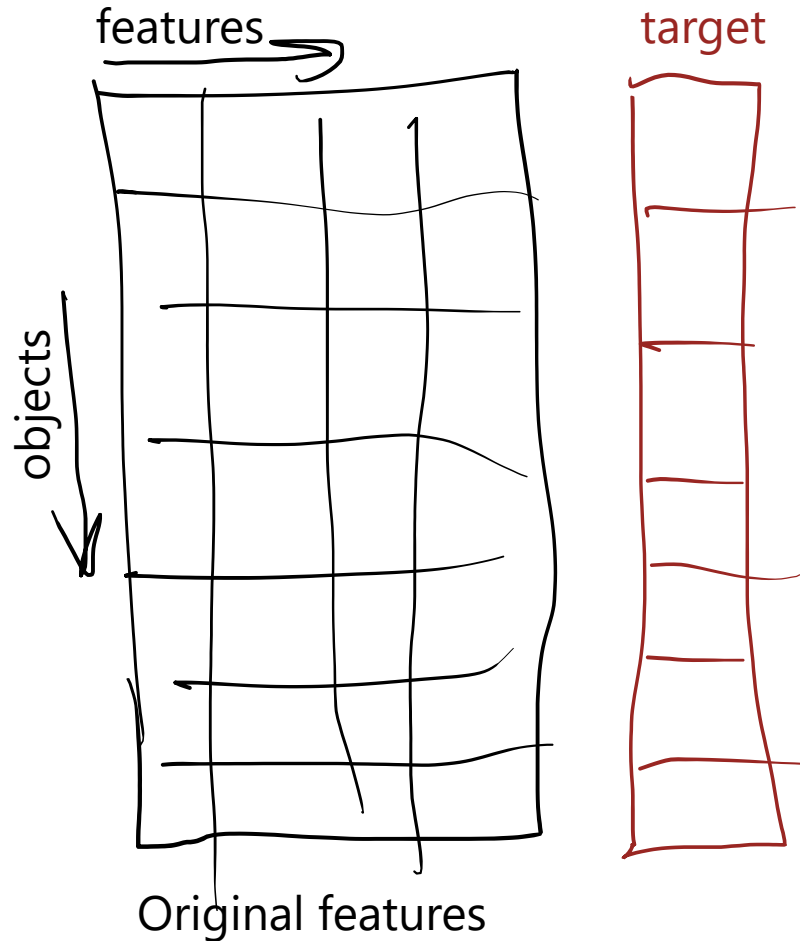
Then major vote error probability $\rightarrow 0$ as $M \rightarrow +\infty$.

Stacking



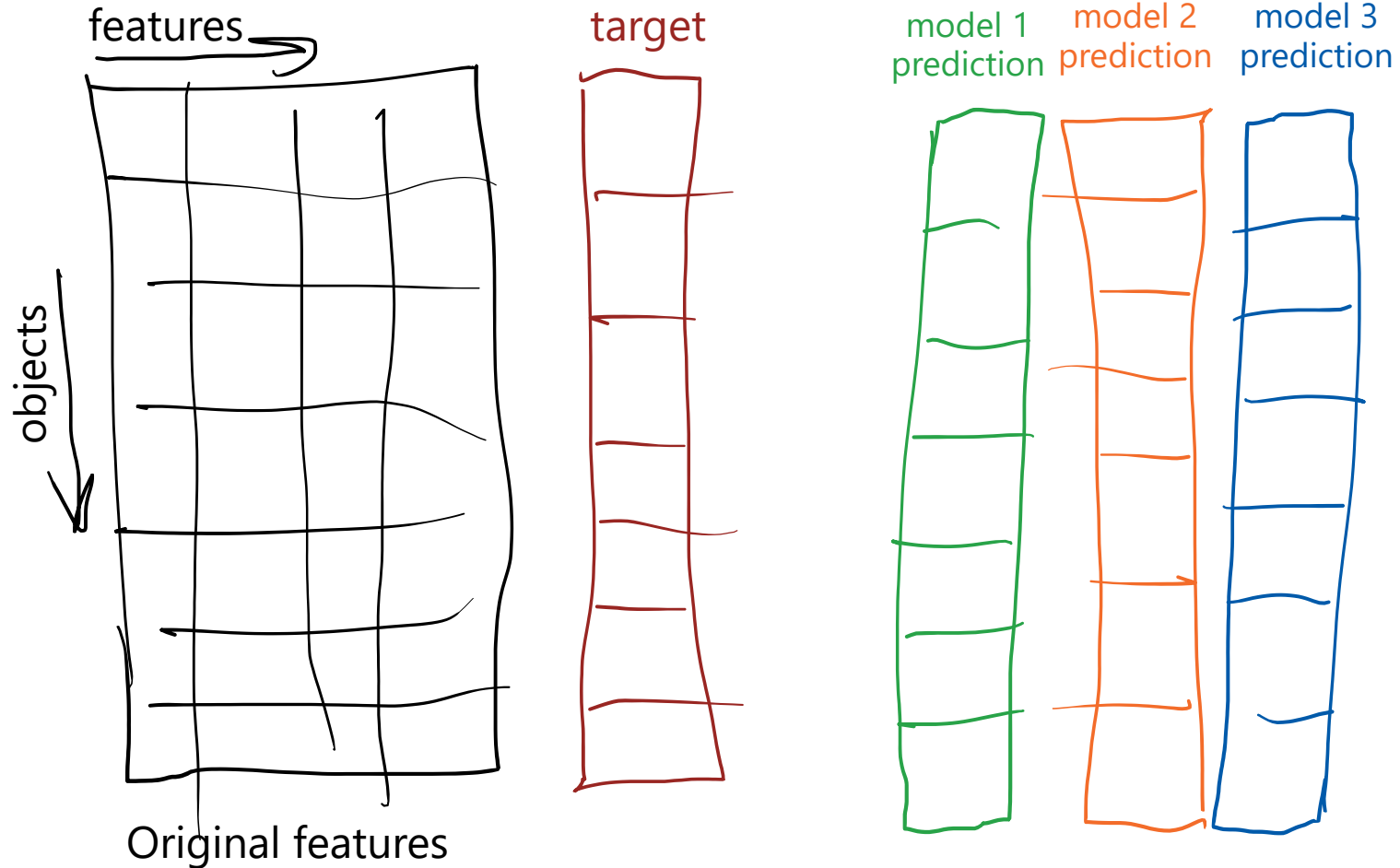
Stacking

Using output of initial models as features for the aggregation model



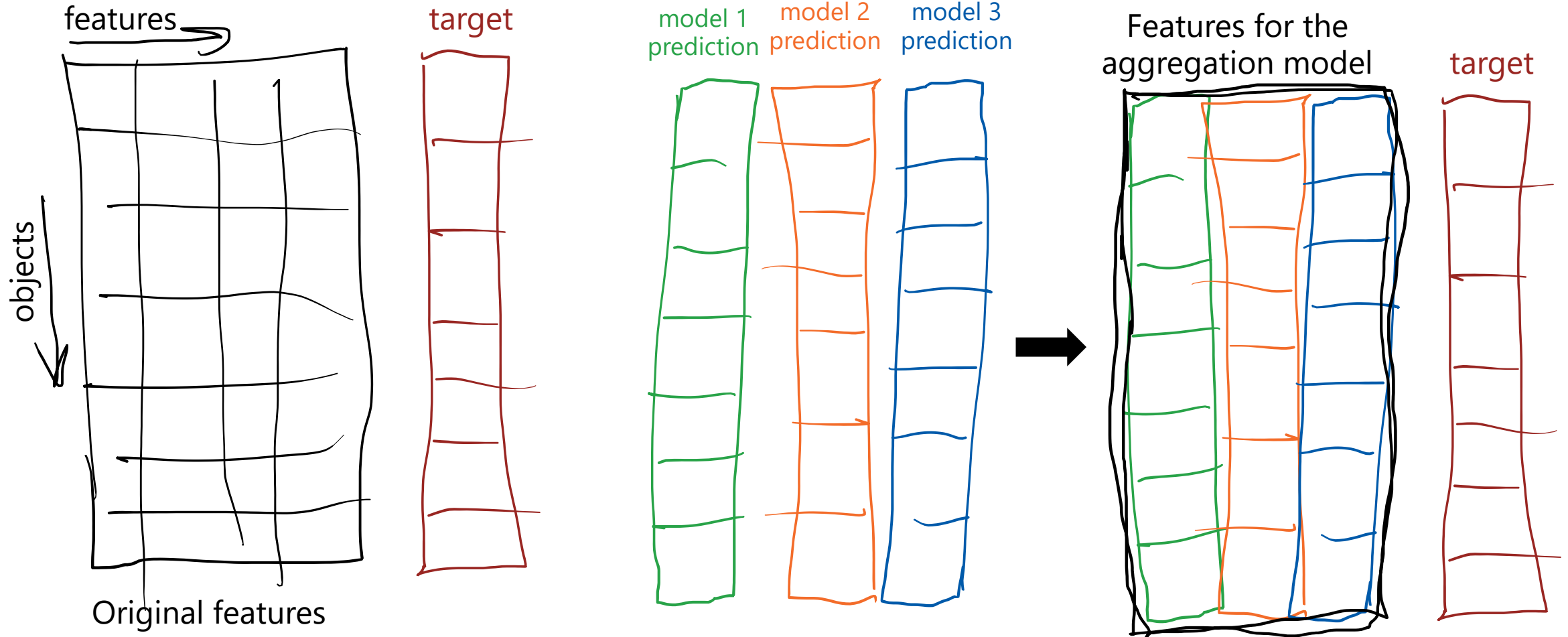
Stacking

Using output of initial models as features for the aggregation model



Stacking

Using output of initial models as features for the aggregation model



Stacking

Using output of initial models as features for the aggregation model

Easy to overfit if using the base models' predictions from the train data

Solution: use cross-validated estimates

Bagging



Bagging

Bagging = bootstrap aggregating

Generate bootstrap samples (re-sample your data with replacement)

Train independent base models on those samples

Combine their predictions

Example

- Bagged decision trees
- Random Forest (= bagged trees + feature subsampling)
- Extra random trees (= Random Forest + randomized splits)

Boosting



Forward stagewise additive modeling (FSAM)

Loss function: $L(f(x), y)$

Base learners \widehat{f}_m

Approximate the output as:

$$F_M(x) = \widehat{f}_0(x) + \sum_{m=1}^M \widehat{c}_m \cdot \widehat{f}_m(x)$$

Do so in steps:

- Start from 0, constant or just fit \widehat{f}_0 to data
- At each step solve:

$$(\widehat{c}_m, \widehat{f}_m) = \operatorname{argmin}_{c, f} \left[\sum_{n=1}^N L(F_{m-1}(x_n) + c \cdot f(x_n), y_n) \right]$$

AdaBoost

AdaBoost = FSAM with exponential loss

$$L(f(x), y) = \exp[-y \cdot f(x)]$$

$$y \in \{-1, +1\},$$

- and base learners being binary classifiers: $\widehat{f}_m(x) \in \{-1, +1\}$

Minimization can be done analytically

- if individual learners allow for weighted samples

AdaBoost

Loss at step m :

$$\sum_{n=1}^N \exp[-y_n(F_{m-1}(x_n) + c \cdot f(x_n))]$$

AdaBoost

Loss at step m :

$$\sum_{n=1}^N \exp[-y_n(F_{m-1}(x_n) + c \cdot f(x_n))]$$
$$= \sum_{n=1}^N \exp[-y_n \cdot F_{m-1}(x_n)] \cdot \exp[-y_n \cdot c \cdot f(x_n)] = \sum_{n=1}^N w_n \cdot \exp[-y_n \cdot c \cdot f(x_n)]$$

AdaBoost

Loss at step m :

$$\begin{aligned} & \sum_{n=1}^N \exp[-y_n(F_{m-1}(x_n) + c \cdot f(x_n))] \\ &= \sum_{n=1}^N \exp[-y_n \cdot F_{m-1}(x_n)] \cdot \exp[-y_n \cdot c \cdot f(x_n)] = \sum_{n=1}^N w_n \cdot \exp[-y_n \cdot c \cdot f(x_n)] \\ &= \sum_{y_n=f(x_n)} w_n \cdot e^{-c} + \sum_{y_n \neq f(x_n)} w_n \cdot e^c \end{aligned}$$

AdaBoost

Loss at step m :

$$\begin{aligned} & \sum_{n=1}^N \exp[-y_n(F_{m-1}(x_n) + c \cdot f(x_n))] \\ &= \sum_{n=1}^N \exp[-y_n \cdot F_{m-1}(x_n)] \cdot \exp[-y_n \cdot c \cdot f(x_n)] = \sum_{n=1}^N w_n \cdot \exp[-y_n \cdot c \cdot f(x_n)] \\ &= \sum_{y_n=f(x_n)} w_n \cdot e^{-c} + \sum_{y_n \neq f(x_n)} w_n \cdot e^c \\ &= \sum_{n=1}^N w_n \cdot e^{-c} + \sum_{y_n \neq f(x_n)} w_n \cdot (e^c - e^{-c}) \end{aligned}$$

AdaBoost

So at step m we optimize:

$$(\widehat{c}_m, \widehat{f}_m) = \operatorname{argmin}_{c, f} \left[\sum_{n=1}^N w_n \cdot e^{-c} + \sum_{y_n \neq f(x_n)} w_n \cdot (e^c - e^{-c}) \right]$$

$$w_n = \exp[-y_n \cdot F_{m-1}(x_n)]$$

In other words, find \widehat{f}_m , s.t. $\sum_{\widehat{f}_m(x_n) \neq y_n} w_n$ is minimized

Easy to show that for c :

$$c_m = \frac{1}{2} \ln \frac{\sum_{\widehat{f}_m(x_n)=y_i} w_n}{\sum_{\widehat{f}_m(x_n) \neq y_i} w_n}$$

AdaBoost

So at step m we optimize:

$$(\widehat{c}_m, \widehat{f}_m) = \operatorname{argmin}_{c, f} \left[\sum_{n=1}^N w_n \cdot e^{-c} + \sum_{y_n \neq f(x_n)} w_n \cdot (e^c - e^{-c}) \right]$$

$$w_n = \exp[-y_n \cdot F_{m-1}(x_n)]$$

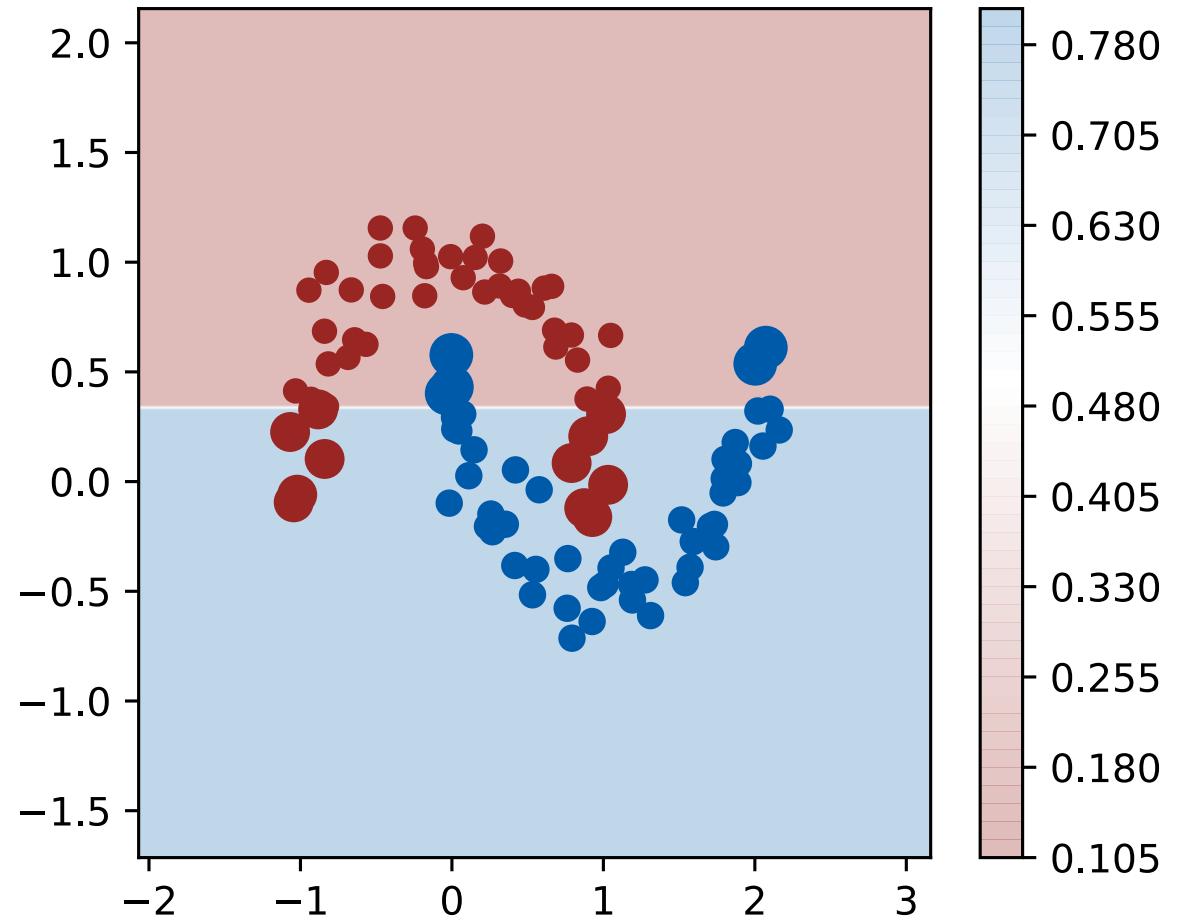
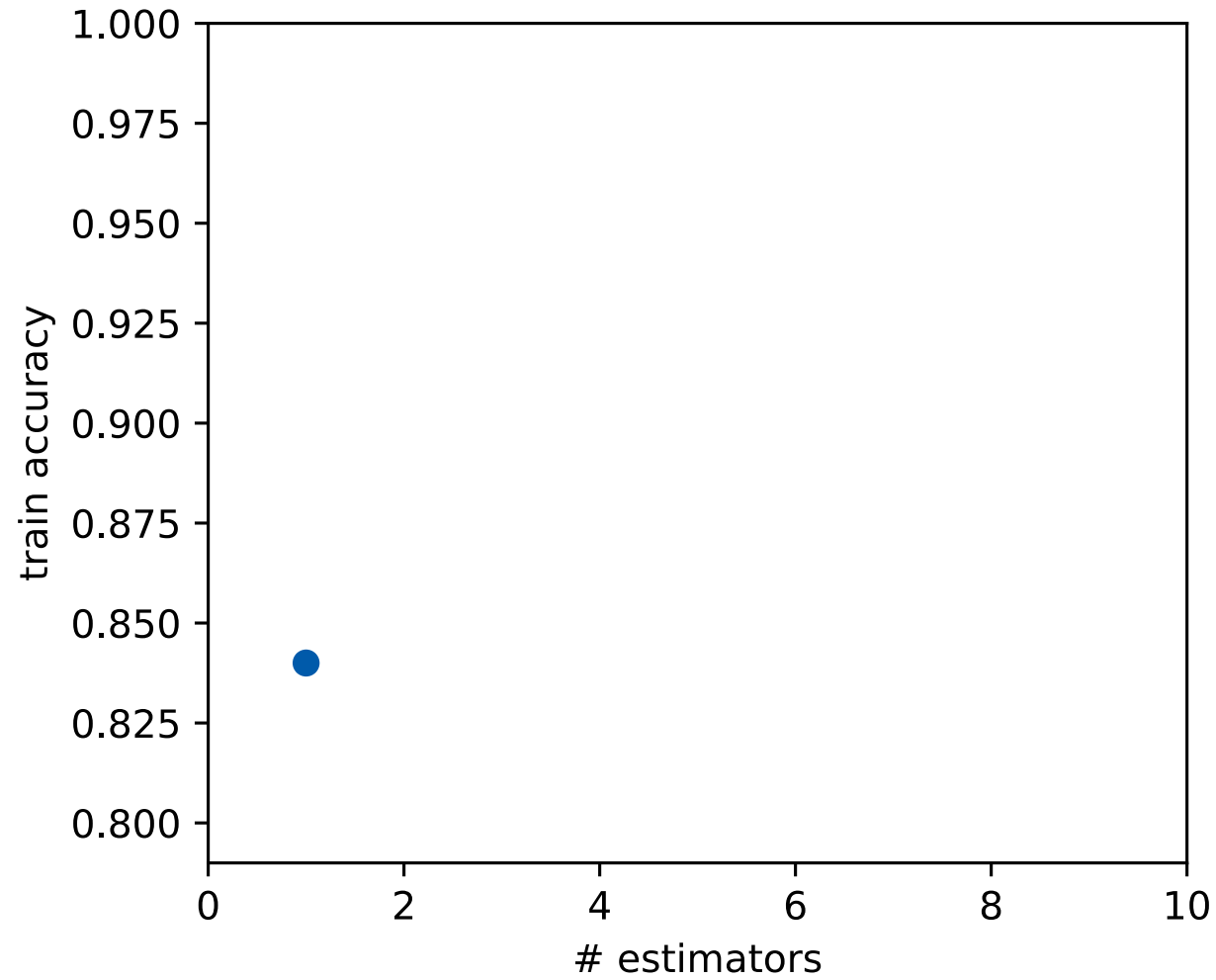
In other words, find \widehat{f}_m , s.t. $\sum_{\widehat{f}_m(x_n) \neq y_n} w_n$ is minimized

Easy to show that for c :

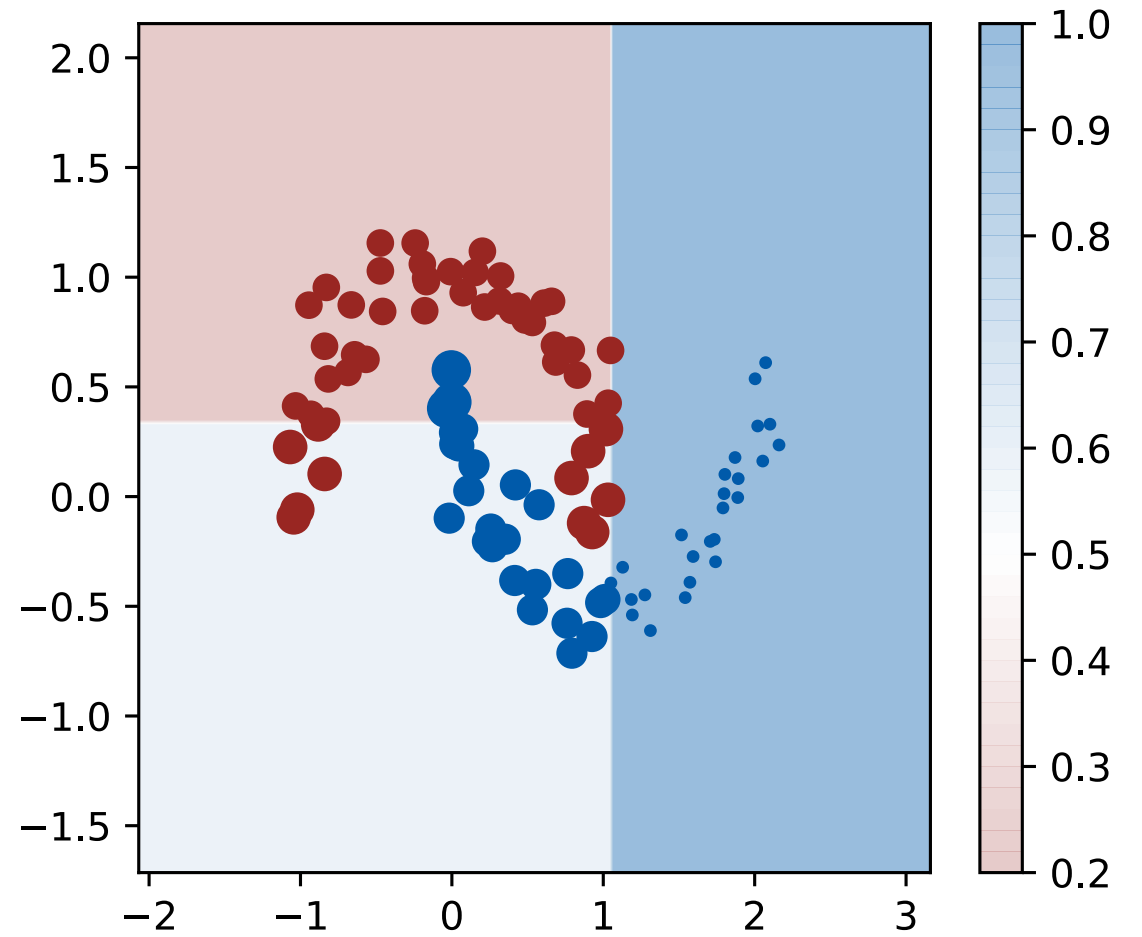
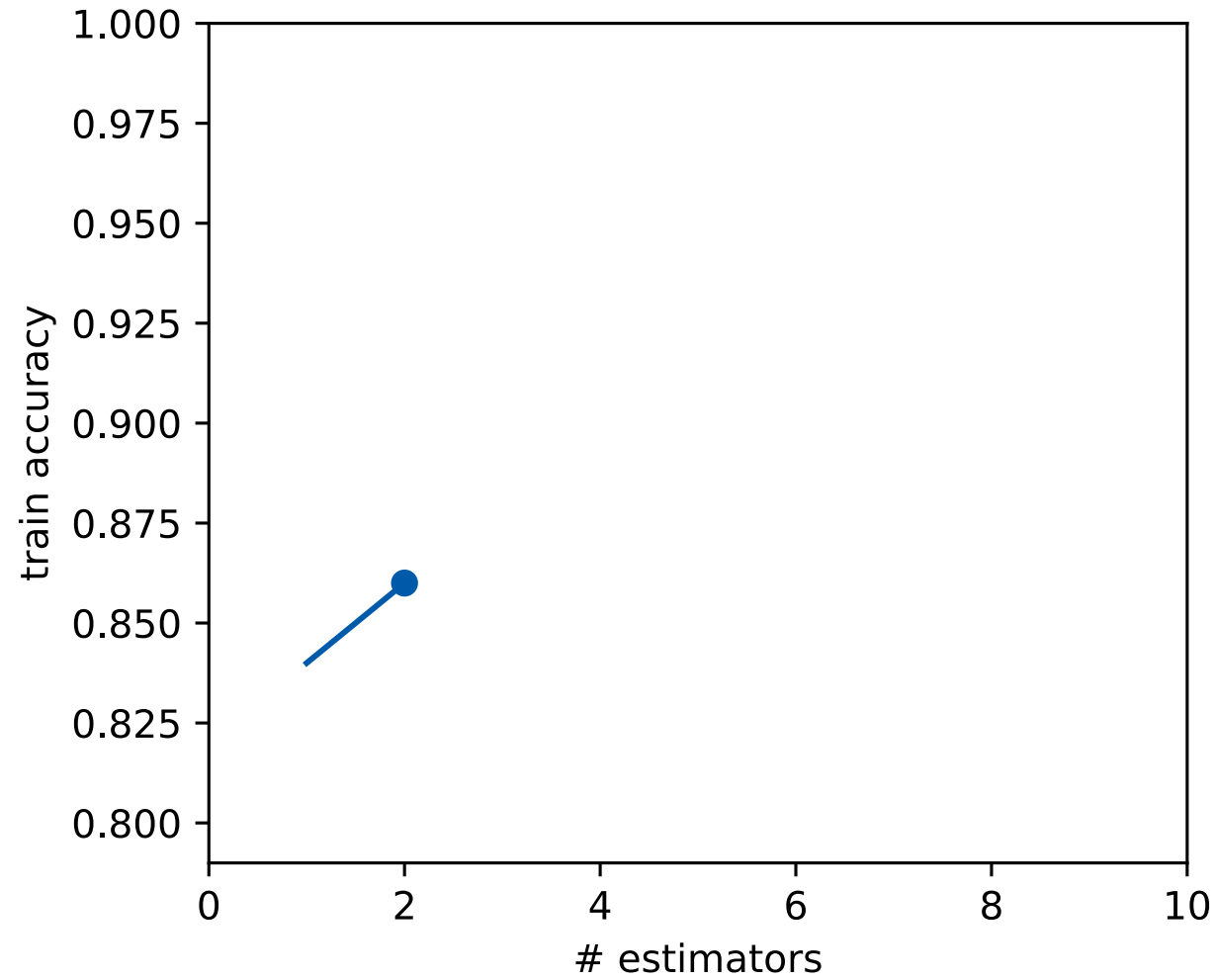
$$c_m = \frac{1}{2} \ln \frac{\sum_{\widehat{f}_m(x_n)=y_i} w_n}{\sum_{\widehat{f}_m(x_n) \neq y_i} w_n}$$

**pay more attention to
objects predicted
wrongly**

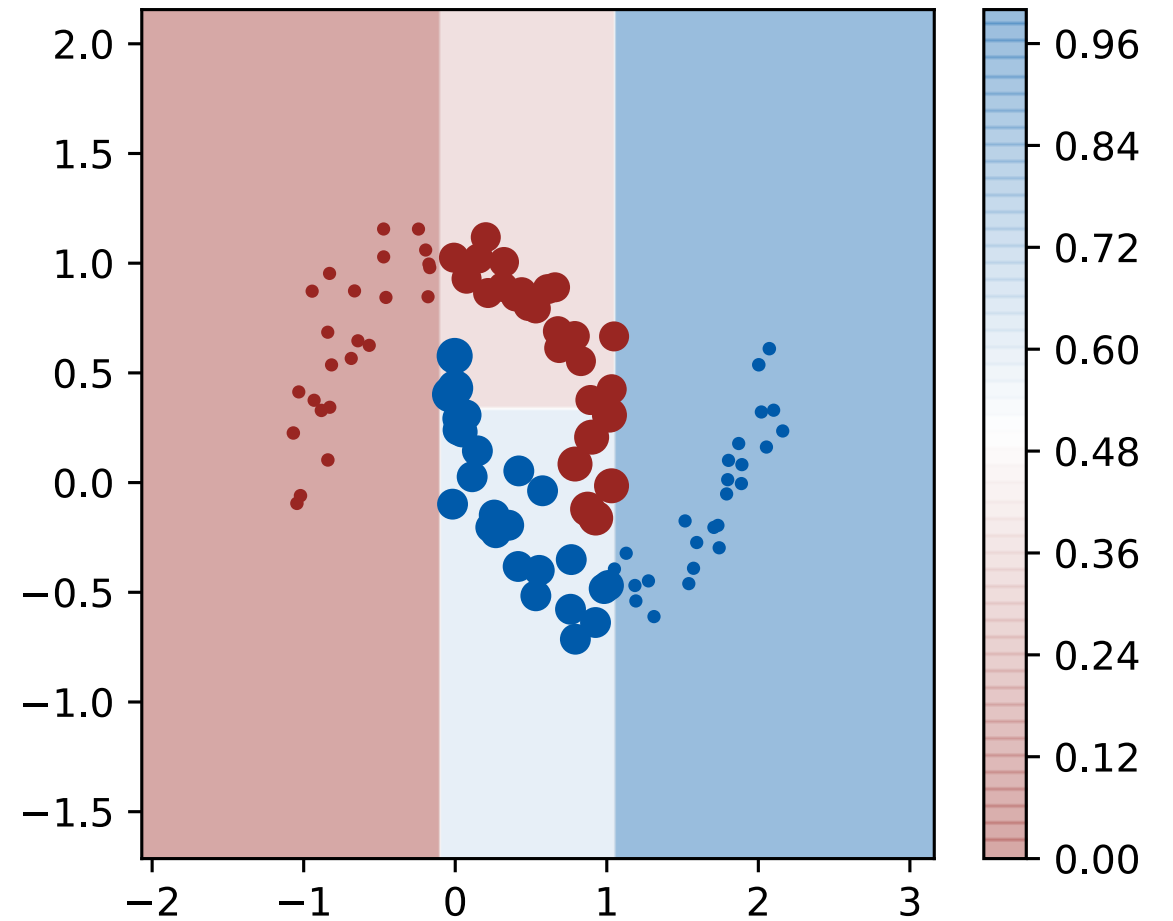
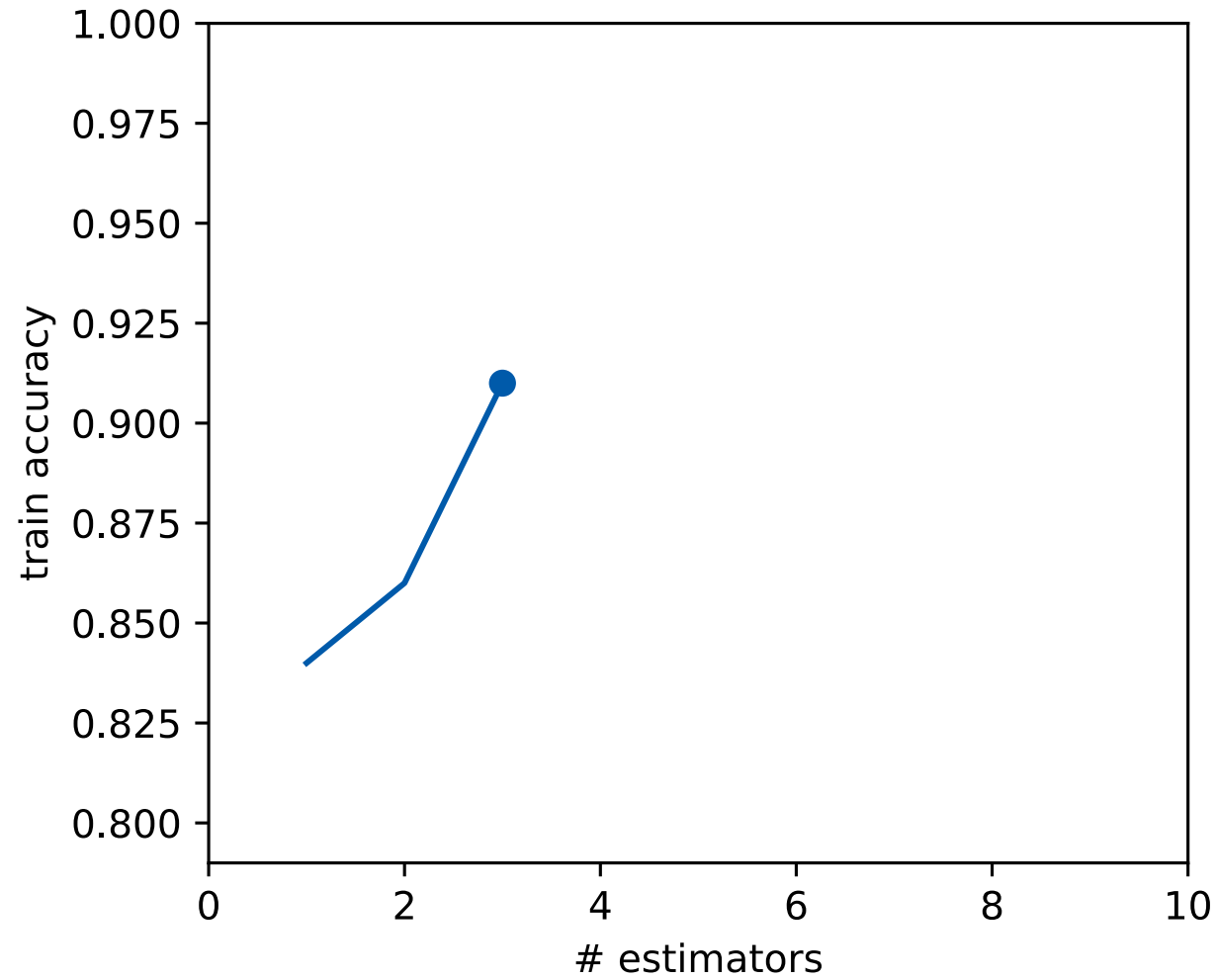
Example



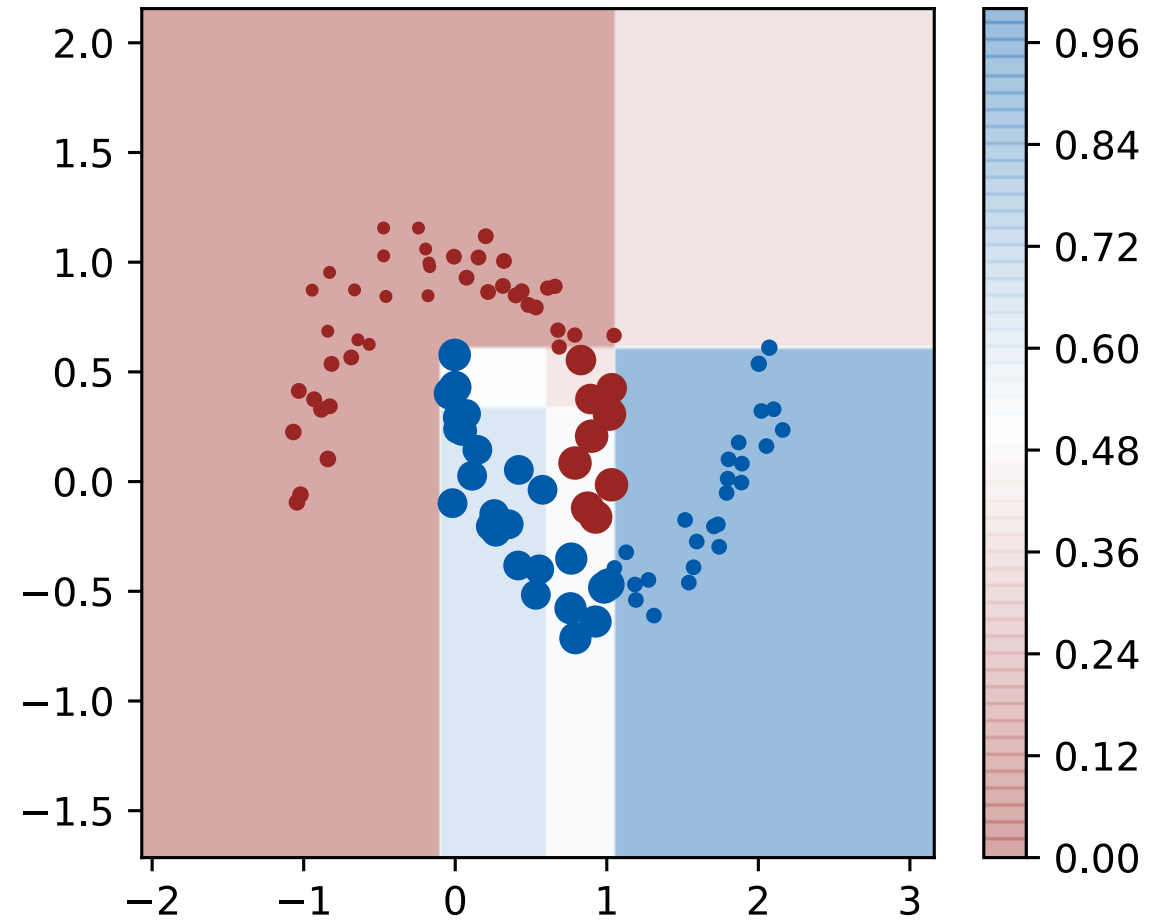
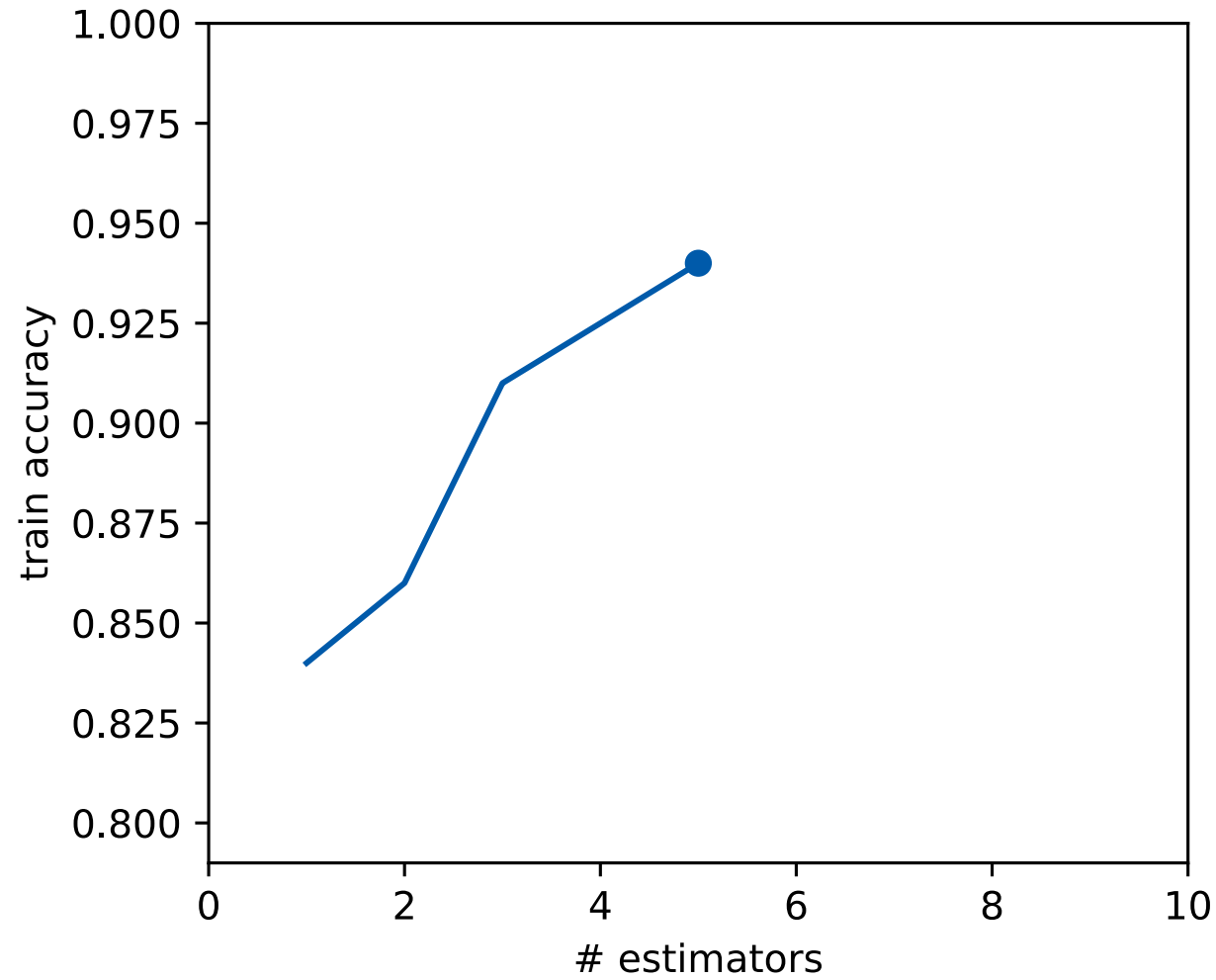
Example



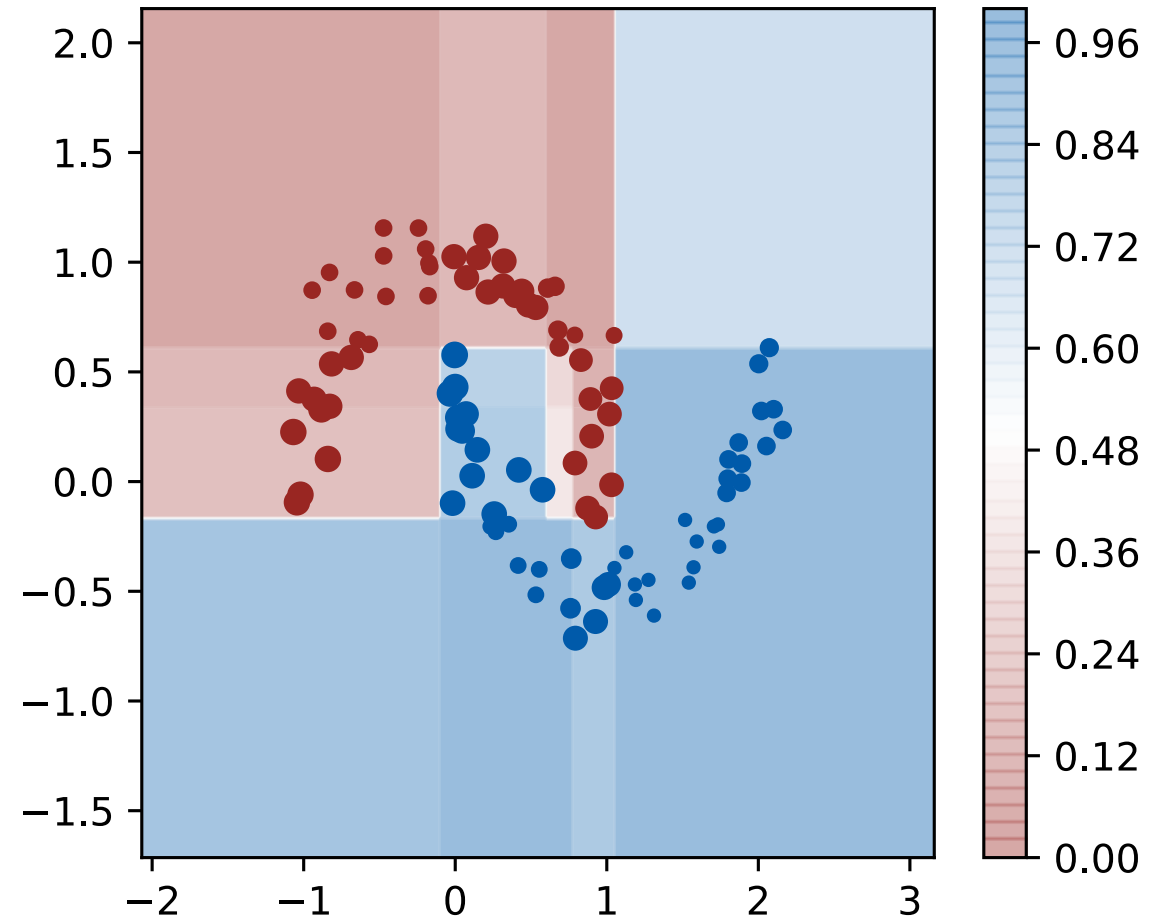
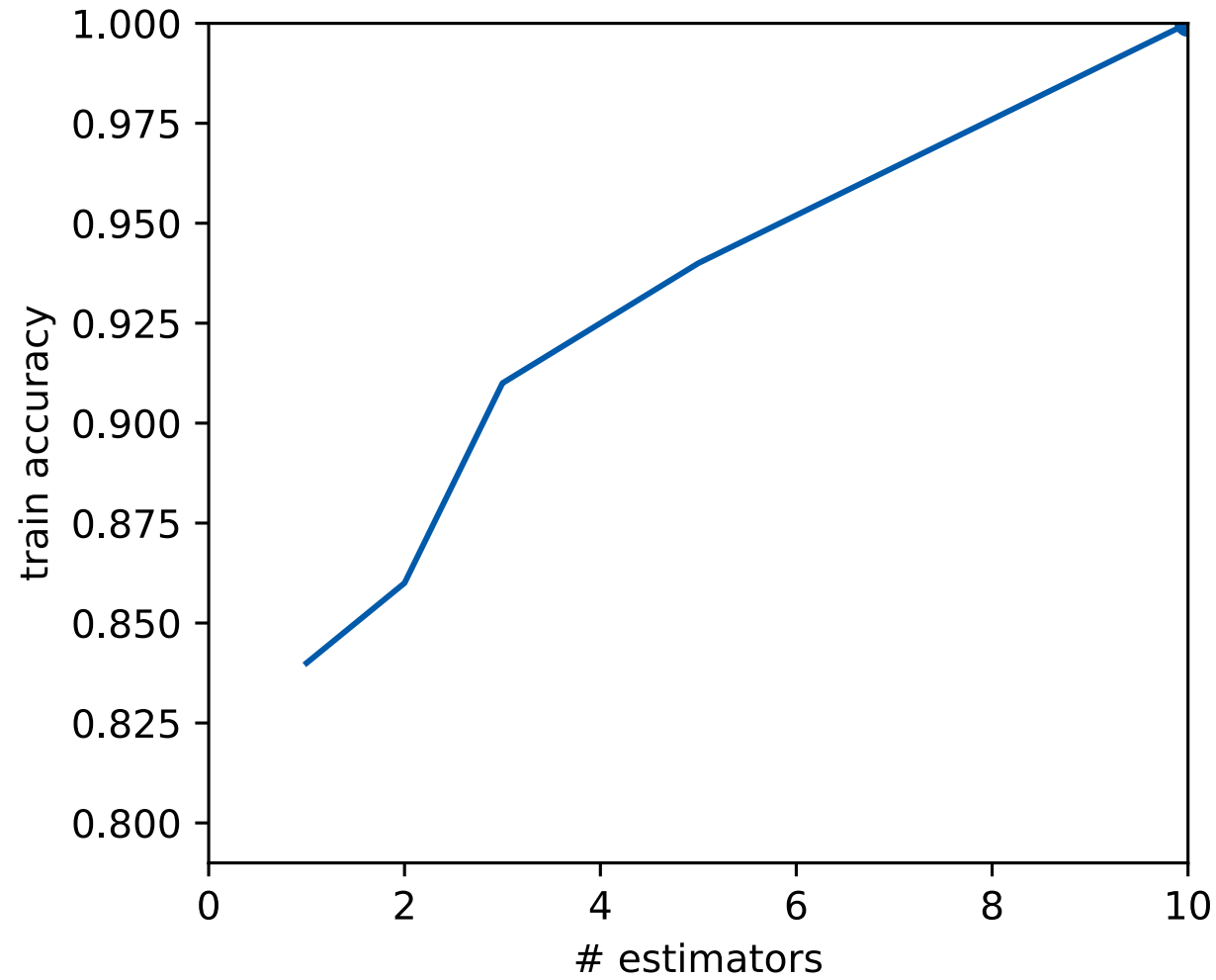
Example



Example



Example



Summary

Ensembling may allow to reduce variance and/or bias of the base learners

Averaging the prediction of **independent** high-variance models reduces the variance

With bagging, the base learners are made (quasi-) independent using bootstrapping

- Can be done in parallel

With boosting, the base learners are built in sequence, each next one trying to improve upon the mistakes of the previous steps

- One of the most powerful classes of models

Question to you: does it make sense to boost linear regression?

Thank you!

Majid Sohrabi



msohrabi@hse.ru



@MSOHRABI_CS