

Ensemble methods

Herman Kamper

2023-04

Bagging

We could overcome problems because of using a single dataset (e.g. overfitting) by collecting more datasets:

$$\underline{x} = \begin{bmatrix} -(\underline{x}^{(1)})^T- \\ -(\underline{x}^{(2)})^T- \\ \vdots \\ -(\underline{x}^{(n)})^T- \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \begin{matrix} \uparrow \\ \downarrow \\ \uparrow \\ \downarrow \\ \vdots \\ \uparrow \\ \downarrow \end{matrix} \begin{matrix} \left[\begin{matrix} \underline{x}_1 \end{matrix} \right]; \left[\begin{matrix} y_1 \end{matrix} \right] \\ \left[\begin{matrix} \underline{x}_2 \end{matrix} \right]; \left[\begin{matrix} y_2 \end{matrix} \right] \\ \vdots \\ \left[\begin{matrix} \underline{x}_B \end{matrix} \right]; \left[\begin{matrix} y_B \end{matrix} \right] \end{matrix} \quad \begin{matrix} \Rightarrow f_1(\underline{x}; \underline{\theta}_1) \\ \Rightarrow f_2(\underline{x}; \underline{\theta}_2) \\ \vdots \\ \Rightarrow f_B(\underline{x}; \underline{\theta}_B) \end{matrix}$$

Aggregate model predictions:

- Regression:

$$f(\underline{x}; \underline{\theta}) = \frac{1}{B} \sum_{b=1}^B f_b(\underline{x}; \underline{\theta}_b)$$

- Classification: Majority or weighted voting

But in most cases we cannot just go out and collect more data. So we simulate this process using *bootstrap* samples:

$$\underline{X} = \begin{bmatrix} \text{blue bar} \\ \text{orange bar} \\ \text{green bar} \\ \text{red bar} \\ \vdots \end{bmatrix} \quad \underline{y} = \begin{bmatrix} \text{blue square} \\ \text{orange square} \\ \text{green square} \\ \text{red square} \\ \vdots \end{bmatrix}$$

Random forests

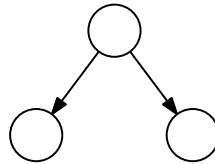
Specifically used with decision trees (regression and classification).

Algorithm:

- Bagging: Train each tree on a different bootstrap sample. But then also ...
- Every time you split, only consider $M < D$ random features.
- Common choice: $M = \sqrt{D}$

Input: $\mathbf{x} \in \mathbb{R}^D$

$x_1, x_2, x_3, \dots, x_D$



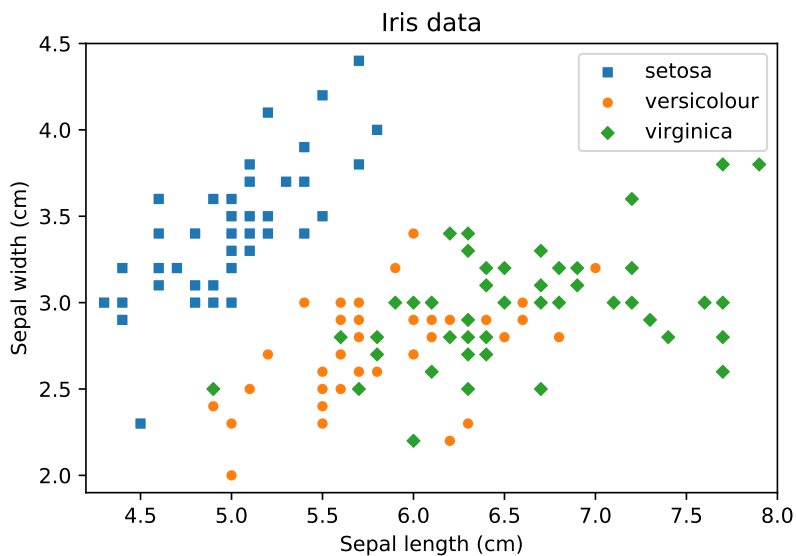
Aggregate model predictions:

- Regression:

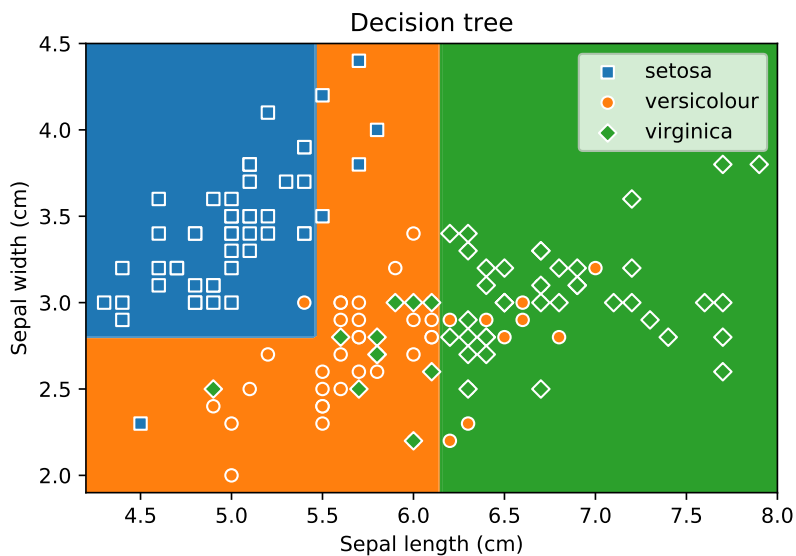
$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{B} \sum_{b=1}^B f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

- Classification: Majority or weighted voting

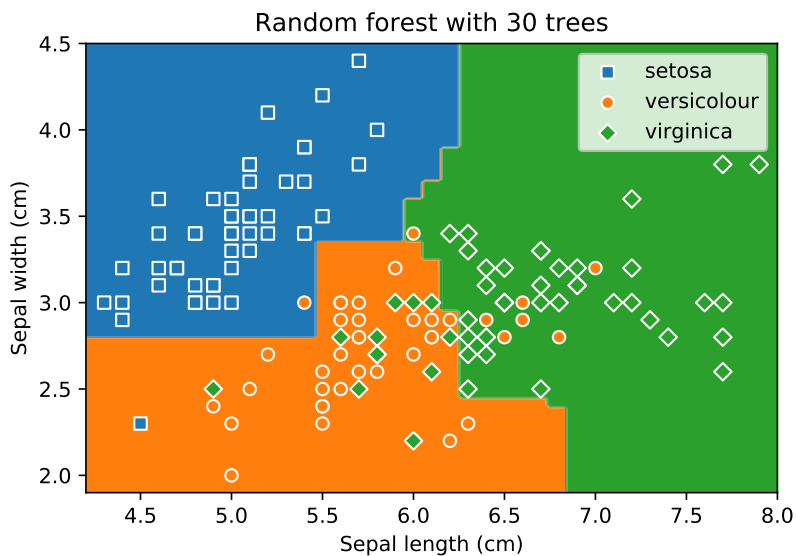
Random forest on iris dataset



Random forest on iris dataset



Random forest on iris dataset



Boosting for regression

Can we combine multiple weak models (just a little bit better than random) to get a “strong” resulting model?

Boosting is different from bagging in that we train a model, look at its mistakes, and then train the next model (instead of just training all the models in one go).

It is different from random forests since it can be used with any weak model (instead of working with decision trees specifically).

Boosting for regression algorithm

Algorithm

- Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow 0$.
- for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunk version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b)$$

- Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

What are we doing?

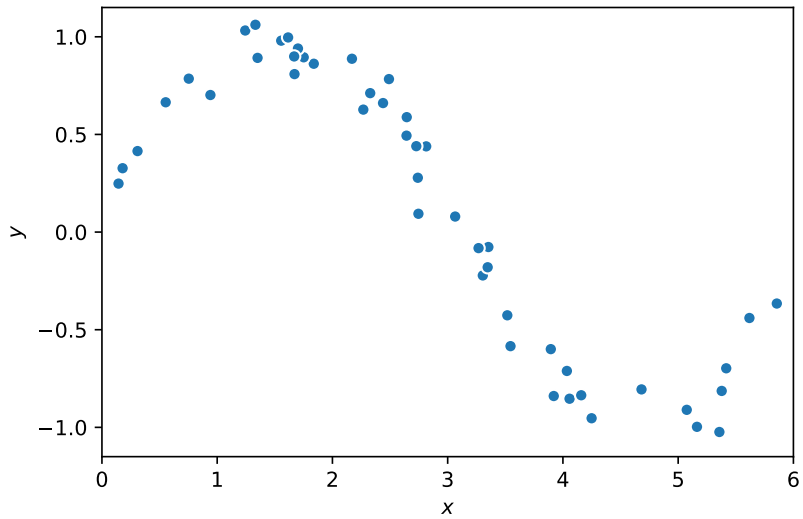
- At $b = 1$:

$\mathbf{r} = \mathbf{y}$, so we are just fitting a model to inputs \mathbf{X} , outputs \mathbf{y} .

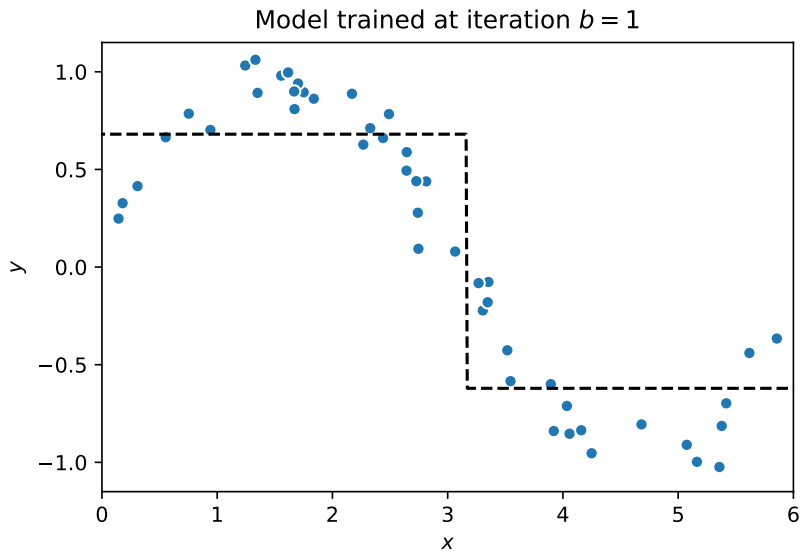
- At $b > 1$:

Fitting a model to the residuals.

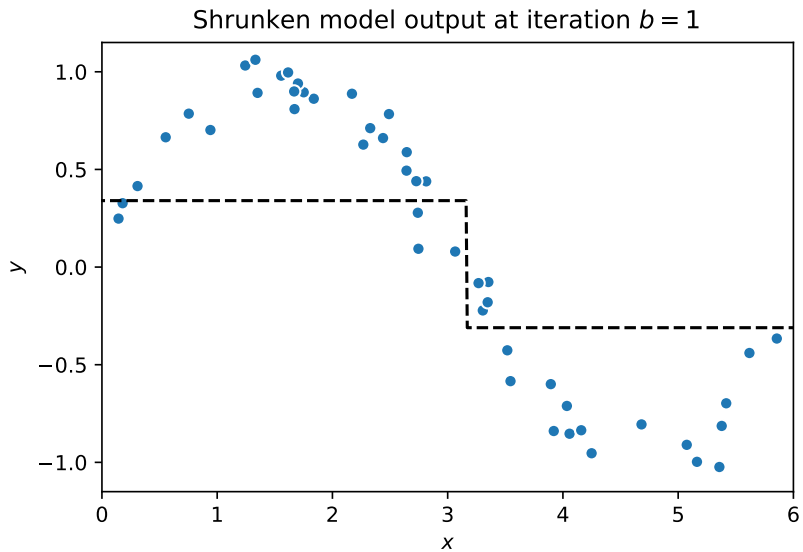
Boosting for regression example



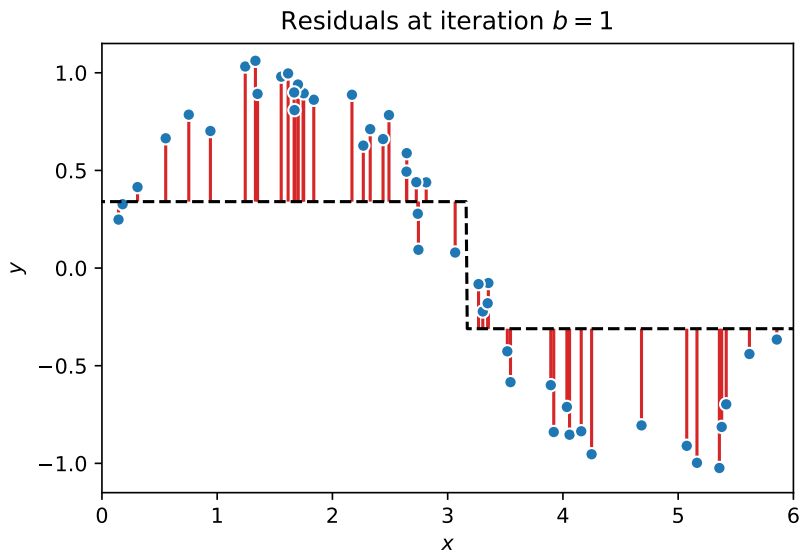
Boosting for regression example



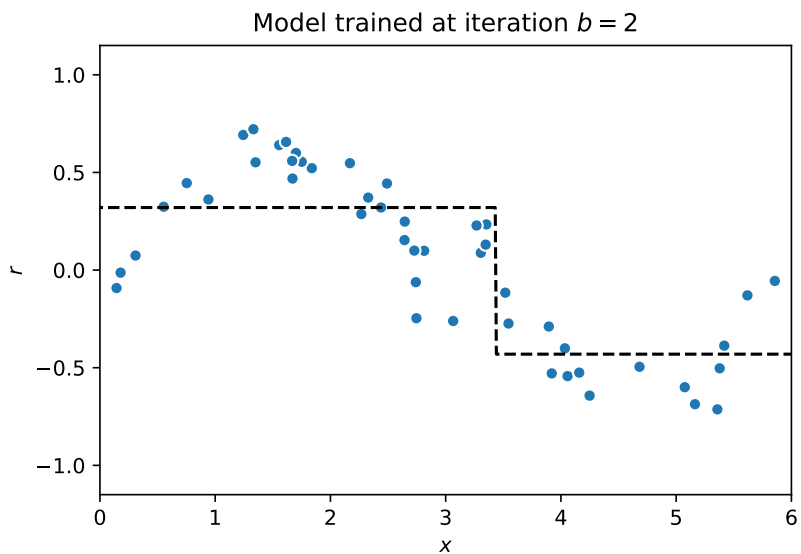
Boosting for regression example



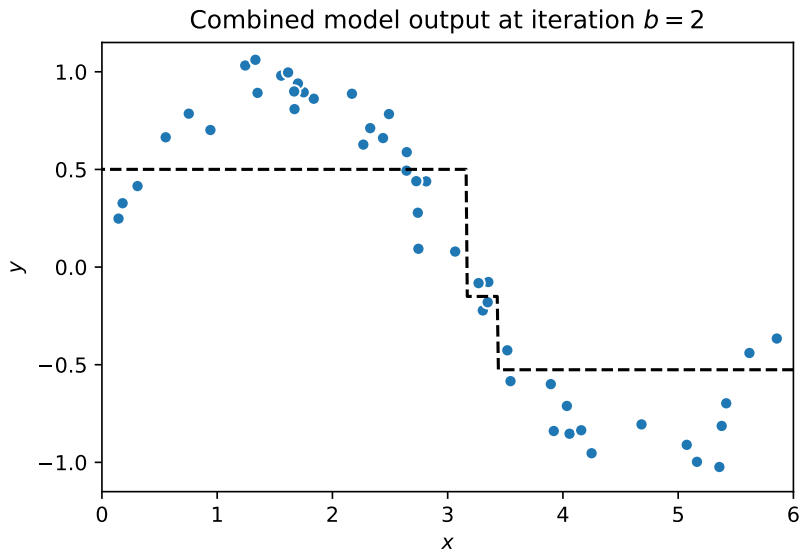
Boosting for regression example



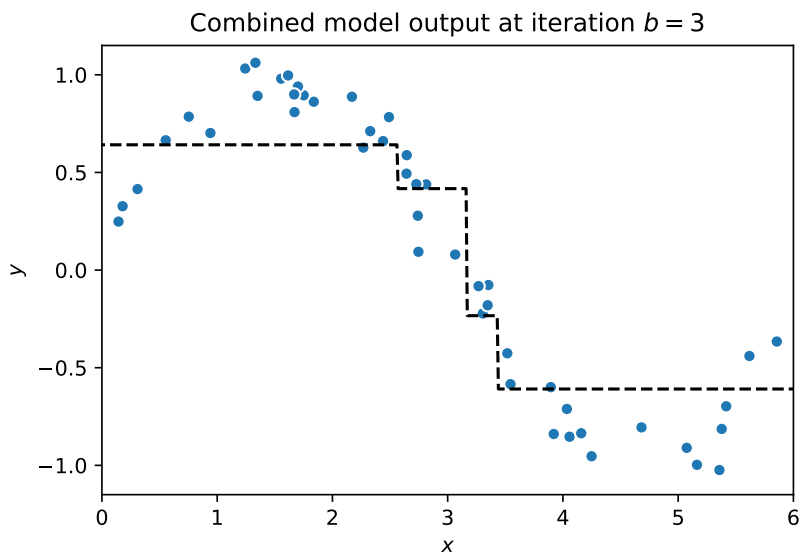
Boosting for regression example



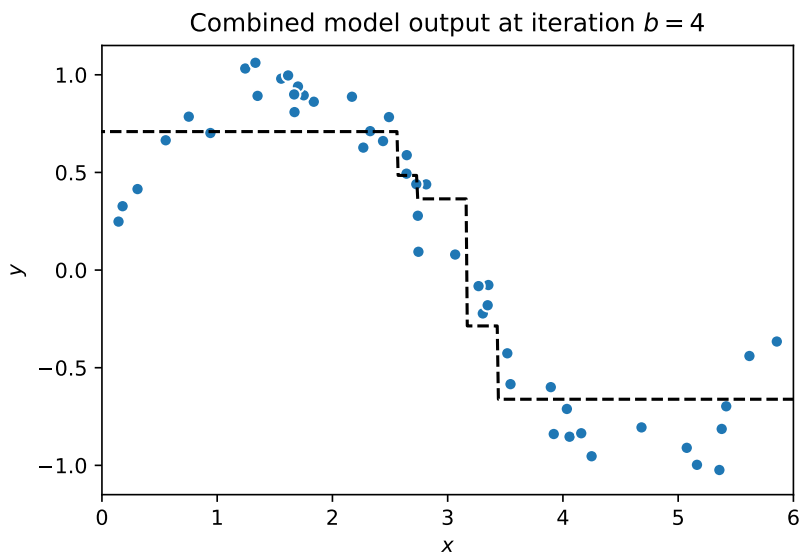
Boosting for regression example



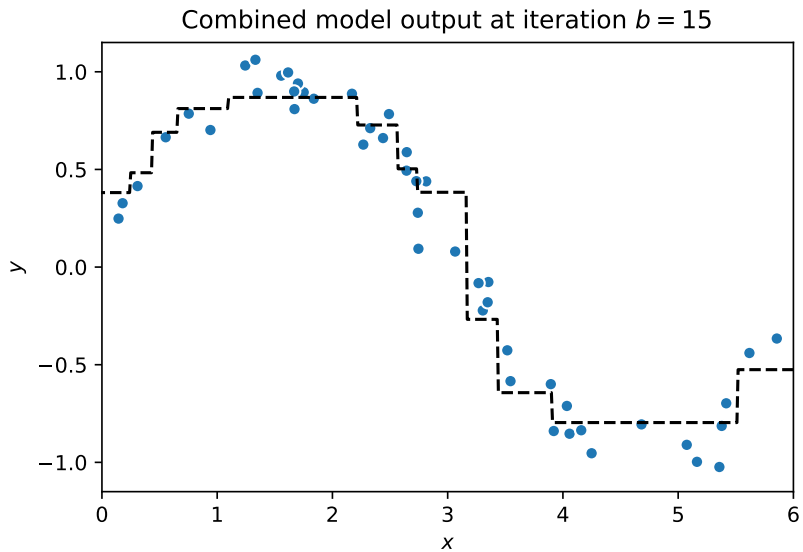
Boosting for regression example



Boosting for regression example



Boosting for regression example



AdaBoost: Boosting for classification

We look at binary classification: $y \in \{-1, 1\}$

Train B models, each $f_b(\mathbf{x}; \boldsymbol{\theta}_b) \in \{-1, 1\}$.

Combine weighted votes:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$$

AdaBoost: Boosting for classification

Algorithm

- Initialise training item weights $w^{(n)} = 1$ for all N training items.
- for iteration $b = 1$ to B :
 - (a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ so that it minimises classification error weighted by $w^{(n)}$.
 - (b) Set model weight using error ϵ :

$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

where ϵ is the normalised error term.

- (c) Update training item weights:

$$\begin{aligned} w^{(n)} &\leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ correct} \\ w^{(n)} &\leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ incorrect} \end{aligned}$$

- Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$

Details

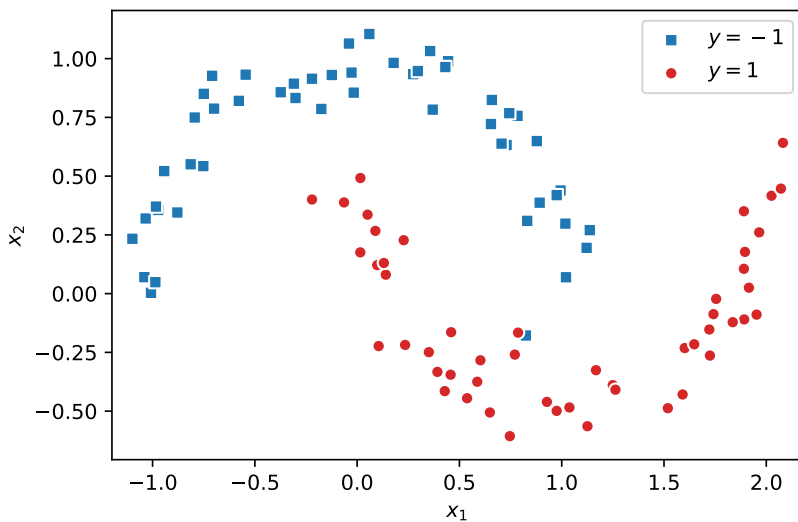
- (a) Goal: Minimise $\sum_{n=1}^N w^{(n)} \mathbb{I} \{y^{(n)} \neq f_b(\mathbf{x}; \boldsymbol{\theta}_b)\}$

- (b) Error term: $\epsilon = \frac{\sum_{n=1}^N w^{(n)} \mathbb{I} \{y^{(n)} \neq f_b(\mathbf{x}; \boldsymbol{\theta}_b)\}}{\sum_{n=1}^N w^{(n)}}$

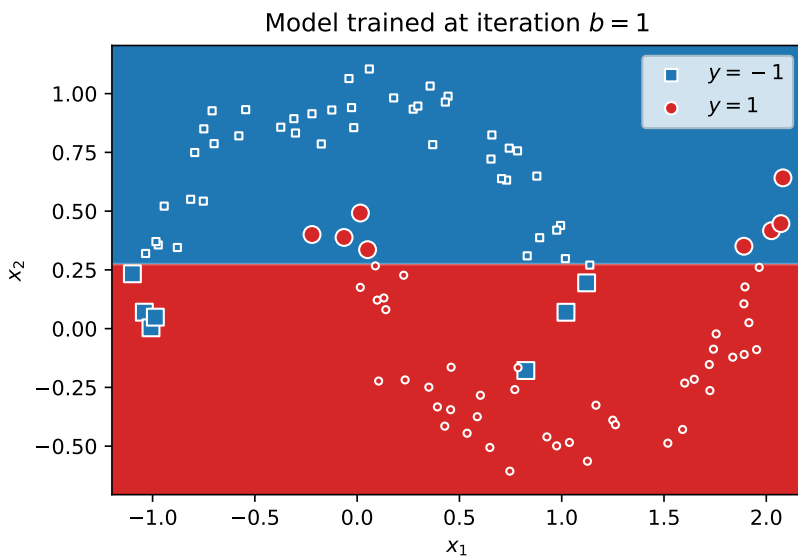
- (c) Correct: $w^{(n)} \leftarrow w^{(n)} \sqrt{\frac{\epsilon}{1 - \epsilon}}$

Incorrect: $w^{(n)} \leftarrow w^{(n)} \sqrt{\frac{1 - \epsilon}{\epsilon}}$

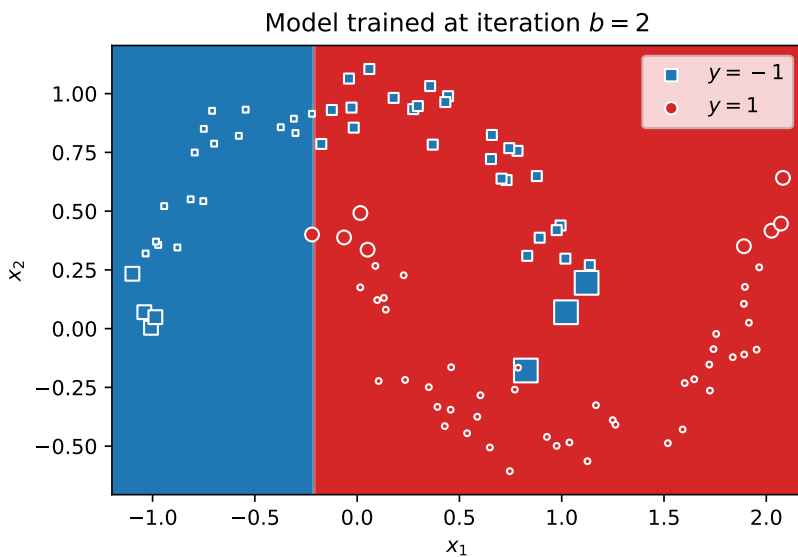
AdaBoost example



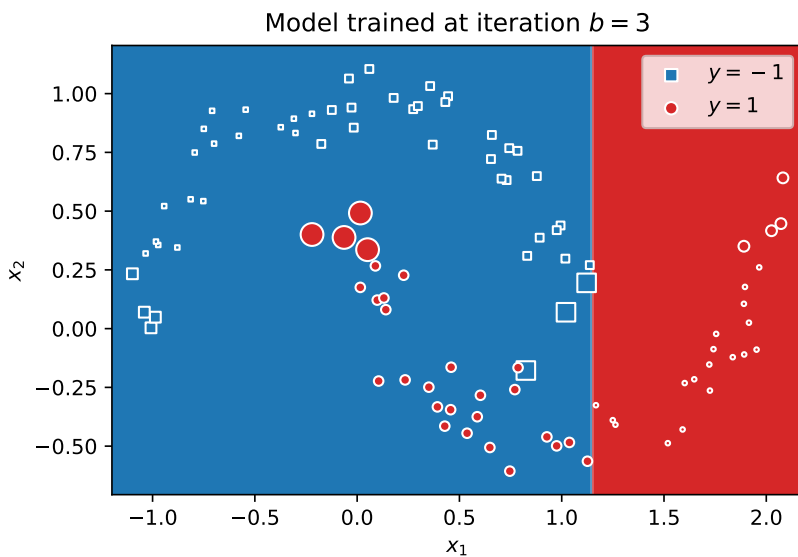
AdaBoost example



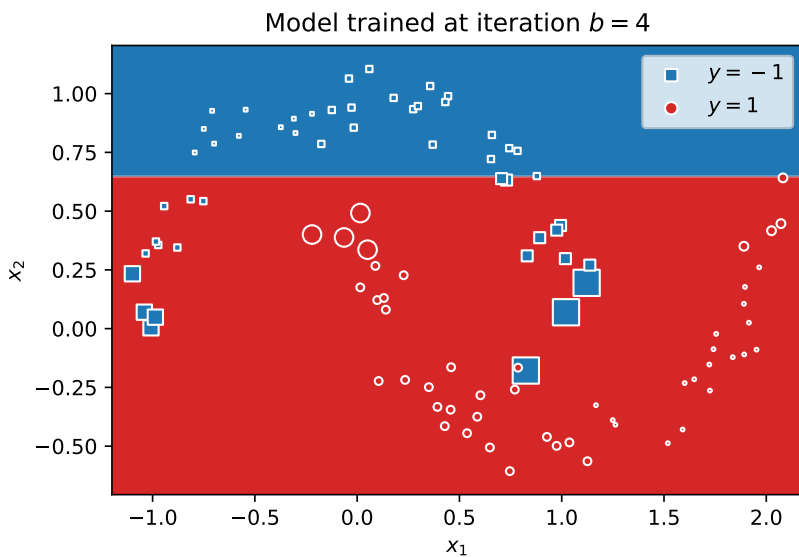
AdaBoost example



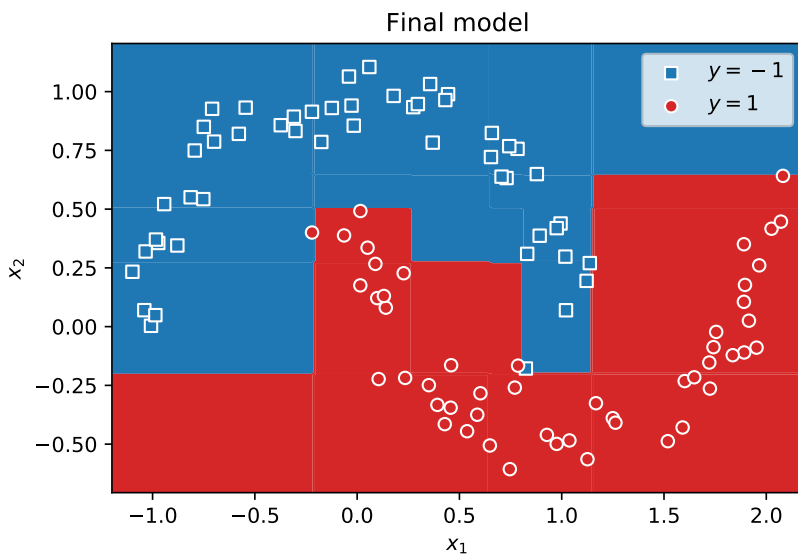
AdaBoost example



AdaBoost example



AdaBoost example



In practice

Many of the machine learning competitions on Kaggle are won with boosting-based methods, especially on problems with tabular data.

Videos covered in this note

- [Ensemble methods 1: Bagging](#) (13 min)
- [Ensemble methods 2: Random forests](#) (7 min)
- [Ensemble methods 3: Boosting for regression](#) (21 min)
- [Ensemble methods 4.1: AdaBoost for classification - Setup](#) (10 min)
- [Ensemble methods 4.2: AdaBoost for classification - Step-by-step](#) (15 min)
- [Ensemble methods 4.3: AdaBoost for classification - Details](#) (11 min)

Reading

- ISLR 8.2.1
- ISLR 8.2.2
- ISLR 8.2.3

Further reading

R. Rojas, "[AdaBoost and the super bowl of classifiers: A tutorial introduction to adaptive boosting](#)," *Freie Universität Berlin*, 2009.

T. Hastie, R. Tibshirani, and J. H. Friedman, *[The Elements of Statistical Learning: Data Mining, Inference, and Prediction](#)*, 2nd ed., 2009, Chapter 10.