

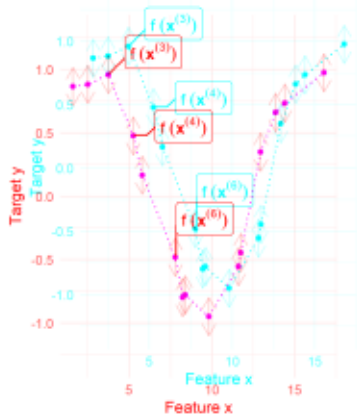
# Introduction to Machine Learning



## Boosting

### Boosting: Concept

### Gradient Boosting: Concept



### Learning goals

#### Learning goals

- Understand idea of forward stagewise modelling
- Understand idea of forward stagewise modelling
- Understand fitting process of gradient boosting for regression problems
- Understand fitting process of gradient boosting for regression problems



## FORWARD STAGewise ADDITIVE MODELING / 3

Hence, we add additive components in a greedy fashion by sequentially minimizing the risk only w.r.t. the next additive component:

$$\min_{\alpha, \theta} \sum_{i=1}^n L \left( y^{(i)}, \hat{f}^{[m-1]}(\mathbf{x}^{(i)}) + \alpha b(\mathbf{x}^{(i)}, \theta) \right)$$



Doing this iteratively is called **forward stagewise additive modeling**.

---

### Algorithm Forward Stagewise Additive Modeling.

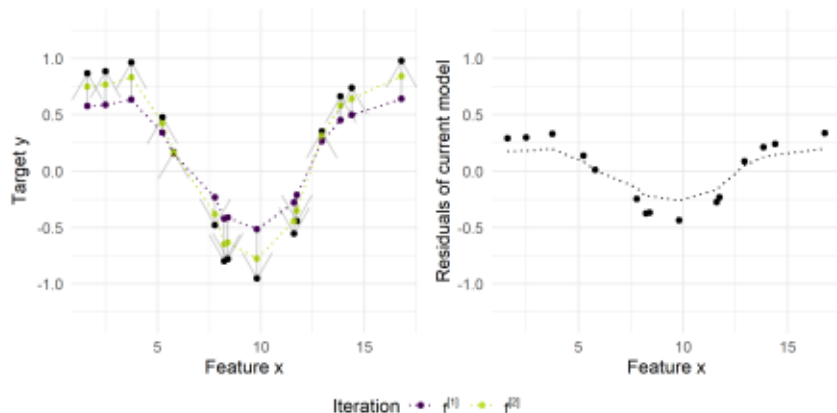
---

- 1: Initialize  $\hat{f}^{[0]}(\mathbf{x})$  with loss optimal constant model
  - 2: **for**  $m = 1 \rightarrow M$  **do**
  - 3:  $(\hat{\alpha}^{[m]}, \hat{\theta}^{[m]}) = \arg \min_{\alpha, \theta} \sum_{i=1}^n L \left( y^{(i)}, \hat{f}^{[m-1]}(\mathbf{x}^{(i)}) + \alpha b(\mathbf{x}^{(i)}, \theta) \right)$
  - 4: Update  $\hat{f}^{[m]}(\mathbf{x}) \leftarrow \hat{f}^{[m-1]}(\mathbf{x}) + \hat{\alpha}^{[m]} b(\mathbf{x}, \hat{\theta}^{[m]})$
  - 5: **end for**
-

# GRADIENT BOOSTING / 2

## Iteration 2:

Let's move our function  $f(\mathbf{x}^{(i)})$  a fraction towards the pseudo-residuals with a learning rate of  $\alpha = 0.6$ .



# GRADIENT BOOSTING / 3

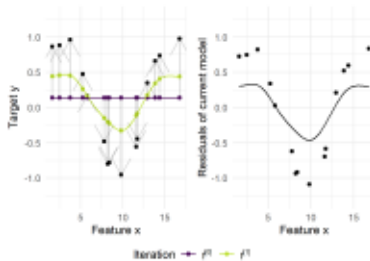
To parameterize a model in this way is pointless, as it just memorizes the instances of the training data.

So, we restrict our additive components to  $b(\mathbf{x}, \theta^{[m]}) \in \mathcal{B}$ .

The pseudo-residuals are calculated exactly as stated above, then we fit a simple model  $b(\mathbf{x}, \theta^{[m]})$  to them:

$$\hat{\theta}^{[m]} = \underset{\theta}{\operatorname{arg\,min}} \sum_{j=1}^m \left( \hat{r}^{[m](j)} - b(\mathbf{x}^{(j)}, \theta) \right)^2$$

So, evaluated on the training data, our  $b(\mathbf{x}, \theta^{[m]})$  corresponds as closely as possible to the negative loss function gradient and generalizes over the whole space.

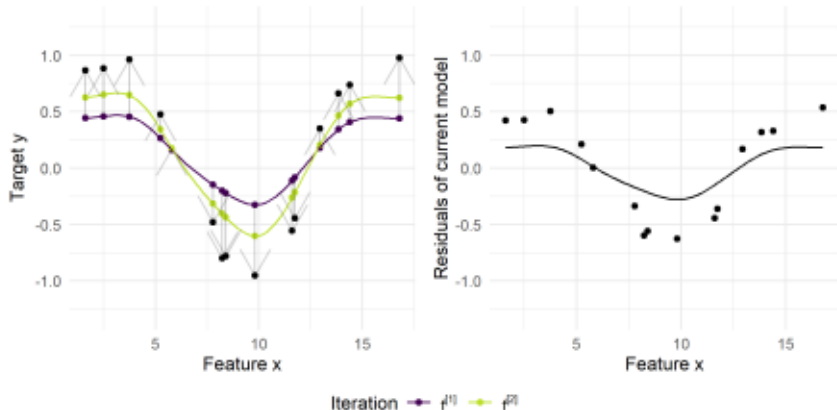




# GRADIENT BOOSTING / 5

Instead of moving the function values for each observation by a fraction closer to the observed data, we fit a regression base learner to the pseudo-residuals (right plot).

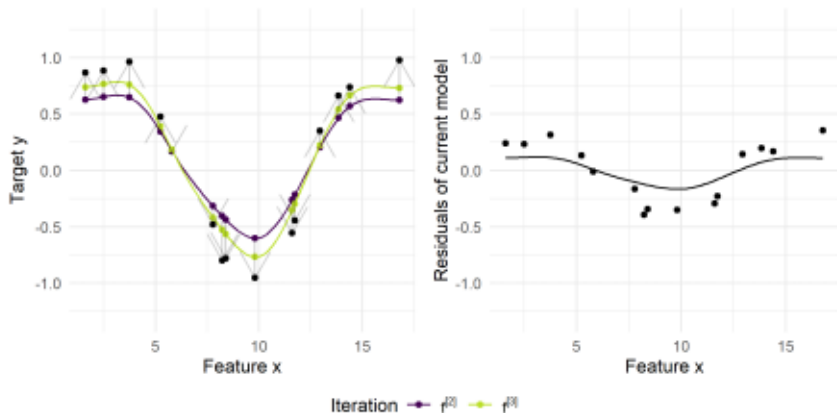
Iteration 2:



# GRADIENT BOOSTING / 6

This base learner is then added to the current state of the ensemble weighted by the learning rate (here:  $\alpha = 0.4$ ) and for the next iteration again the pseudo-residuals of the adapted ensemble are calculated and a base learner is fitted to them.

**Iteration 3:**





# GRADIENT BOOSTING ALGORITHM

---

## Algorithm Gradient Boosting Algorithm.

---

- 1: Initialize  $\hat{f}^{[0]}(\mathbf{x}) = \arg \min_{\theta_0 \in \mathbb{R}} \sum_{i=1}^n L(y^{(i)}, \theta_0)$
  - 2: **for**  $m = 1 \rightarrow M$  **do**
  - 3: For all  $i$ :  $\tilde{r}^{[m](i)} = - \left[ \frac{\partial L(y, f)}{\partial f} \right]_{f=\hat{f}^{[m-1]}(\mathbf{x}^{(i)}, y=y^{(i)}}$
  - 4: Fit a regression base learner to the vector of pseudo-residuals  $\tilde{r}^{[m]}$ :
  - 5:  $\hat{\theta}^{[m]} = \arg \min_{\theta} \sum_{i=1}^n (\tilde{r}^{[m](i)} - b(\mathbf{x}^{(i)}, \theta))^2$
  - 6: Set  $\alpha^{[m]}$  to  $\alpha$  being a small constant value or via line search
  - 7: Update  $\hat{f}^{[m]}(\mathbf{x}) = \hat{f}^{[m-1]}(\mathbf{x}) + \alpha^{[m]} b(\mathbf{x}, \hat{\theta}^{[m]})$
  - 8: **end for**
  - 9: Output  $\hat{f}(\mathbf{x}) = \hat{f}^{[M]}(\mathbf{x})$
- 



Note that we also initialize the model in a loss-optimal manner.

## LINE SEARCH

The learning rate in gradient boosting influences how fast the algorithm converges. Although a small constant learning rate is commonly used in practice, it can also be replaced by a line search.



Line search is an iterative approach to find a local minimum. In the case of setting the learning rate, the following one-dimensional optimization problem has to be solved:

$$\hat{\alpha}^{[m]} = \underset{\alpha}{\operatorname{arg\,min}} \sum_{i=1}^m \mathcal{L}(y^{(i)}, f^{[m-1]}(\mathbf{x}_i) + \alpha b(\mathbf{x}_i, \hat{\theta}^{[m]}))$$

Optionally, an (inexact) backtracking line search can be used to find the  $\alpha^{[m]}$  that minimizes the above equation.