

Gradient Boosted Decision Trees Classification

Classification

Given a training set: $Z=\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
 \mathbf{x}_i - features, y_i - class labels (0, 1)

Goal is to find $f(\mathbf{x})$ using training set, such as

$$\min \sum_{(\mathbf{x}, y) \in T} [f(\mathbf{x}) \neq y]$$

at test set $T=\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

How to build $f(\mathbf{x})$?

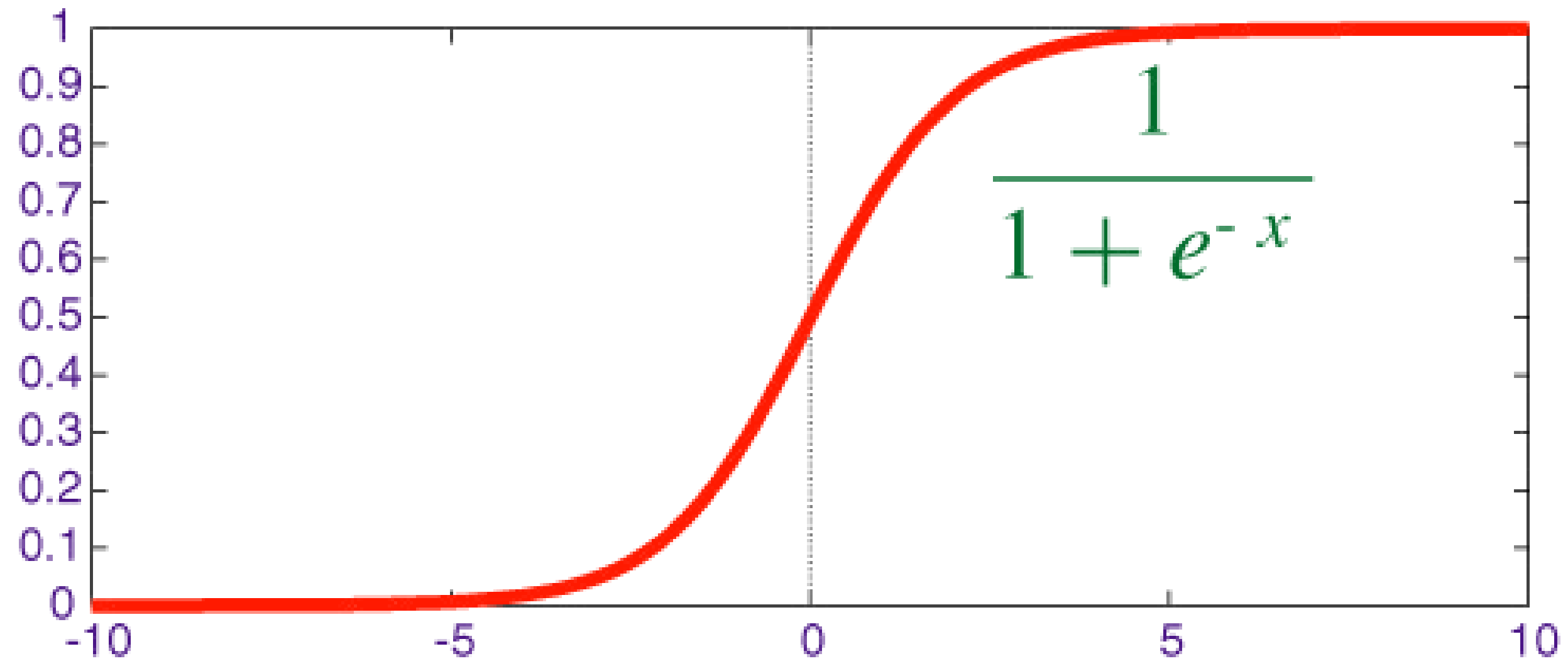
Gradient Boosted Trees for Classification

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\sum_{m=1}^M h_m(\mathbf{x}))}$$

$h_m(\mathbf{x})$ - a decision tree

$$0 < P(y = 1|\mathbf{x}) < 1$$

Sigmoid function



$$f(\mathbf{x}) = \sum_{m=1}^M h_m(\mathbf{x})$$

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

Likelihood:

$$\prod_{i=1}^n P(y_i|\mathbf{x}_i) = P(y_1|\mathbf{x}_1) \cdot \dots \cdot P(y_n|\mathbf{x}_n)$$

“The principle of maximum likelihood”

Algorithm: find a function $f(x)$ maximizing the **likelihood**

Equivalent: find a function $f(x)$ maximizing the logarithm of the **likelihood**
(since logarithm is a monotone function)

$$Q[f] = \sum_{i=1}^n \log(P(y_i | \mathbf{x}_i))$$

$$\max Q[f]$$

$$L(y_i, f(\mathbf{x}_i)) = \log(P(y_i|\mathbf{x}_i))$$

$$Q[f] = \sum_{i=1}^n L(y_i, f(\mathbf{x}_i))$$

Algorithm: Gradient Boosted Trees for Classification

Input: training set $Z = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$,
M – number of iterations

1. $f_0(x) = \log \frac{p_1}{1-p_1}$ p_1 - part of objects of first class
2. For $m=1 \dots M$:
3. $g_i = \frac{dL(y_i, f_m(\mathbf{x}_i))}{df_m(\mathbf{x}_i)}$
4. Fit a decision tree $h_m(\mathbf{x}_i)$ to the target g_i
(auxiliary training set $\{(\mathbf{x}_1, g_1), \dots, (\mathbf{x}_n, g_n)\}$)
5. $\rho_m = \underset{\rho}{\operatorname{argmax}} Q[f_{m-1}(\mathbf{x}) + \rho h_m(\mathbf{x})]$
6. $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \mathbf{v} \rho_m h_m(\mathbf{x}_i)$
7. Return: $f_M(\mathbf{x})$

\mathbf{v} - regularization (learningRate), recommended ≤ 0.1