

Gradient Boosting

Suppose we have a family of weak classifiers \mathcal{F} .

$$\text{Set } f_0 = \arg \inf_{f \in \mathcal{F}} \sum_{i=1}^n L(f(x_i), y_i)$$

Update the estimate for $m=1, 2, \dots, M$

$$f_m = \sum_{i=0}^m h_i \quad h_i \in \mathbb{R}^N; \quad h_0 = f_0(x).$$

Steepest descent: $h_m = \overset{\text{length}}{-g_m} \overset{\text{direction}}{g_m} \quad g_m \in \mathbb{R}^N$

$$g_m = \nabla_{f_{m-1}} L(f) = \nabla_{f_{m-1}} \sum_{i=1}^n L(y_i, f(x_i)) \quad (\text{steepest direction})$$

$$f_m = \arg \min_f L(f_{m-1} - f g_m) \quad (\text{best length})$$

$$f_m = f_{m-1} - f_m g_m \quad (\text{step downwards})$$

Random Forests

Suppose we have a tree method $h(\Theta, x)$ where Θ is the parameter set.

Algorithm:

for $b=1, \dots, B$

a. Draw a bootstrap sample of data Z^b of size N .

b. Grow a tree T_b on Z^* by the following procedure.

- (i) Select $m \ll p$ variables
- (ii) Pick the best variable/split point
- (iii) Split node

2. Output $\{T_b\}_{b=1}^B$.

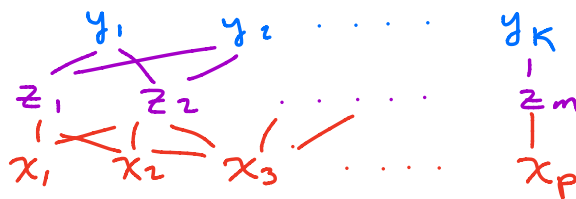
Prediction Rule: $\hat{f}_{RF}(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$ (regression)

$\hat{f}_{RF}(x) = \text{majority vote } \{T_b\}_{b=1}^B$ (classification).

Suppose $(X, Y) \sim P_{X,Y}$ and define the margin

$$\begin{aligned} \text{mg}(X, Y) = \text{margin} &= \text{Avg}_m [I(h_m(x) = y)] \\ &= \max_{k \neq y} \text{Avg} (I(h_m(x) = k)) \end{aligned}$$

Neural Networks



Three layer neural network.

Assumption information is passed $(x) \rightarrow (z) \rightarrow (y)$

In the second layer

$$z_m = \underbrace{\sigma}_{\text{activation}} \left(\alpha_{0m} + \underbrace{\alpha_m^T x}_{\text{first layer information}} \right) \quad m=1, \dots, M$$

$$T_k = \beta_{0k} + \beta_k^T z, \quad k=1, \dots, K$$

$$y_k = \underbrace{g_k}_{\text{Chosen to match response } y_k.}(T_k)$$

e.g. $g(T) = T$ (regression)

$$g(T) = \frac{e^{T_k}}{\sum_{k=1}^K e^{T_k}} \quad (\text{classification}).$$

For a single y

$$y = g \left(\beta_0 + \sum_{i=1}^M \beta_i \sigma(\alpha_{0i} + \alpha_i^T x) \right)$$

Training: Uses Gradient descent to find local optimum

$$(\hat{\alpha}_{0m}, \hat{\alpha}_m, \hat{\beta}_{0k}, \hat{\beta}_k)_{m=1, n=1}^{M, K}$$

Called "Back Propagation"