

Ada-Boost

· Error rate: $\overline{err} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq G(x_i))$

AdaBoost M_1

1. Initialize weights $w_i = 1/N, i=1, 2, \dots, N$
2. for $m = 1:M$
 - (a) Fit a classifier $G_m(x)$ to the training data
 - (b) $err_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i}$
 - (c) $\alpha_m = \frac{1}{2} \log \frac{1 - err_m}{err_m}$
 - (d) $w_i \leftarrow w_i \cdot \exp(\alpha_m \cdot I(y_i \neq G_m(x_i)))$
3. $G(x) = \text{sign}[\sum_{m=1}^M \alpha_m G_m(x)]$

· Additive Model: $f(x) = \sum_{m=1}^M \beta_m \cdot b(x; \gamma_m)$

Forward Stagewise

1. Initialize $f_0(x) = 0$
2. For $m = 1:M$
 - (a) $(\beta_m, \gamma_m) = \argmin_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma_m))$
 - (b) $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$

Exponential Loss: $(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N \exp[-y_i(f_{m-1}(x_i) + \beta G(x_i))]$

$$= \arg \min_{\beta, G} \sum_{i=1}^N w_i^{(m)} \exp(-\beta y_i G(x_i))$$

$$G_m = \arg \min_G \sum_{i=1}^N w_i^{(m)} \exp(-\beta) \cdot I(y_i = G(x_i))$$

$$+ \sum_{i=1}^N w_i^{(m)} \exp(\beta) \cdot I(y_i \neq G(x_i))$$

$$= \arg \min_G \sum_{i=1}^N w_i^{(m)} (\exp(\beta) - \exp(-\beta)) I(y_i \neq G(x_i))$$

$$+ \sum_{i=1}^N w_i^{(m)} \exp(-\beta)$$

$$= \arg \min_G \sum_{i=1}^N w_i^{(m)} I(y_i \neq G(x_i))$$

$$\beta_m = \arg \min_{\beta} \sum_{i=1}^N w_i^{(m)} (\exp(\beta) - \exp(-\beta)) I(y_i \neq G(x_i))$$

$$+ \sum_{i=1}^N w_i^{(m)} \exp(-\beta)$$

$$= \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}$$

$$w_i^{(m+1)} = w_i^{(m)} \cdot e^{-\beta_m G_m(x_i) \cdot y_i}$$

· Error Bound

$$\text{err} = \frac{1}{N} \sum_{i=1}^N I(G(x_i) \neq y_i) \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i \sum_{m=1}^M \beta_m G_m(x_i))$$

$$= \sum_{i=1}^N w_i^{(1)} \prod_{m=1}^M \exp(-\beta_m y_m G_m(x_i))$$

$$= \sum_{i=1}^N w_i^{(2)} \prod_{m=2}^M \exp(-\beta_m y_m G_m(x_i))$$

$$= \dots = \prod_{m=1}^M Z_m$$

$$\text{since } Z_m = \sum_{i=1}^N w_i^{(m)} \exp(-\beta_m y_i G_m(x_i))$$

$$= \sum_{i=1}^N w_i^{(m)} \exp(-\beta_m) I(y = G_m(x_i)) + \sum_{i=1}^N w_i^{(m)} e^{\beta_m} I(y_i \neq G_m(x_i))$$

$$= (1 - e_m) e^{-\beta_m} + e_m e^{\beta_m}$$

$$\alpha_m = \frac{1}{2} \log \frac{1 - e_m}{e_m} = 2 \sqrt{e_m(1 - e_m)} = \sqrt{1 - 4r_m^2} \quad \text{where } r_m = \frac{1}{2} - e_m$$

$$\text{then } \prod_{m=1}^M Z_m = \prod_{m=1}^M \sqrt{1 - 4r_m^2} \leq \exp\left(-2 \sum_{m=1}^M r_m^2\right)$$

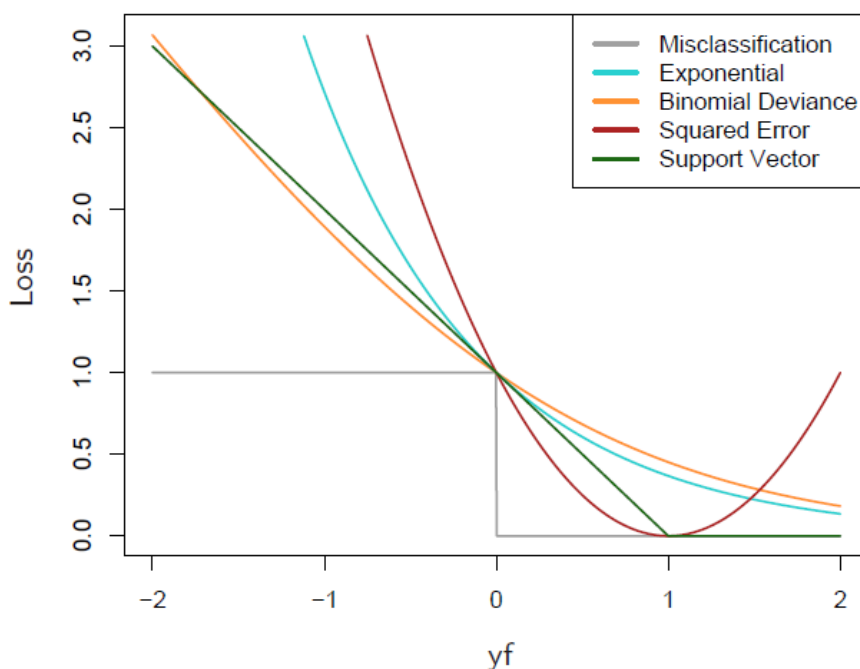


TABLE 10.1. *Some characteristics of different learning methods. Key: ▲ = good, ◆ = fair, and ▼ = poor.*

Characteristic	Neural Nets	SVM	Trees	MARS	k-NN, Kernels
Natural handling of data of “mixed” type	▼	▼	▲	▲	▼
Handling of missing values	▼	▼	▲	▲	▲
Robustness to outliers in input space	▼	▼	▲	▼	▲
Insensitive to monotone transformations of inputs	▼	▼	▲	▼	▼
Computational scalability (large N)	▼	▼	▲	▲	▼
Ability to deal with irrelevant inputs	▼	▼	▲	▲	▼
Ability to extract linear combinations of features	▲	▲	▼	▼	◆
Interpretability	▼	▼	◆	▲	▼
Predictive power	▲	▲	▼	◆	▲

Boosting Tree

- Tree rule: Partition the space into disjoint regions $R_j, j=1, 2, \dots, J$

$$x \in R_j \Rightarrow f(x) = \gamma_j, \quad T(x; \theta) = \sum_{j=1}^J \gamma_j I(x \in R_j)$$

$$\text{where } \hat{\theta} = \arg \min_{\theta} \sum_{j=1}^J \sum_{x_i \in R_j} L(y_i, \gamma_j)$$

- Tree Optimization:

- given R_j : $\hat{\gamma}_j = \overline{y_j}$

- find R_j : $\tilde{\theta} = \arg \min_{\theta} \sum_{i=1}^N \hat{L}(y_i, T(x_i; \theta))$

- Boosting Tree: $f_M(x) = \sum_{m=1}^M T(x; \theta_m)$

- stagewise $\tilde{\theta}_m$: $\hat{\theta}_m = \arg \min_{\theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \theta_m))$

Numerical Boosting

- $\hat{f} = \arg\min_f L(f)$, where $f = \{f(x_1), f(x_2), \dots, f(x_N)\}$
- $f_m = \sum_{n=0}^m h_n$, $h_n \in \mathcal{R}^N$ is a descent step.

Steepest Descent:

$$g_m = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i) = f_{m-1}(x_i)}$$

$$f_m = \arg\min_p L(f_{m-1} - p g_m)$$

then $h_m = -p_m g_m$, $f_m = f_{m-1} - p_m g_m$

Gradient Boosting

- $\hat{\Theta}_m = \arg\min_{\Theta} \sum_{i=1}^N (-g_{im} - T(x_i, \Theta))^2$, a tree closest to gradient

Gradient Tree Boosting

1. Initialized $f_0(x) = \arg\min_r \sum_{i=1}^N L(y_i, r)$

2. For $m=1:M$

(a) For $i=1, 2, \dots, N$, compute

$$r_{im} = - \left[\frac{\partial L(y_i, f_{m-1}(x_i))}{\partial f_{m-1}(x_i)} \right]$$

(b) Fit a regression tree to fit r_{im} with regions

$$R_{jm}, j=1, 2, \dots, J_m$$

(c) For $j=1, 2, \dots, J_m$ compute

$$r_{jm} = \arg\min_r \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + r)$$

(d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} r_{jm} I(x \in R_{jm})$

3. Output $\hat{f}(x) = f_M(x)$