

Boosting for regression

- Fix: B - number of trees
 λ - shrinking λ - small - requires large B
 d - number of splits in each tree ($d+1$ leaves)

$$\begin{matrix} X_1, X_2, \dots, X_n & \in & \mathbb{R}^p \\ \downarrow & & \downarrow \\ y_1 & y_2 & \dots & y_n & \in & \mathbb{R} \end{matrix}$$

1) Set $\hat{f}(x) = 0$ and $r_i = y_i$, $i \in \{\text{training set}\}$

$$\left\{ \underbrace{\hat{f}(x_i)}_{\text{prediction}} + \underbrace{r_i}_{\text{error}} = \underbrace{y_i}_{\text{actual value}} \right.$$

2) For $b=1, 2, \dots, B$:

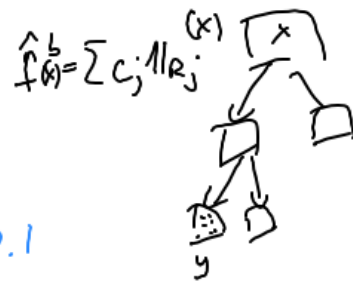
(a) fit a tree \hat{f}^b with d splits to the training data (X, r)

(b) update \hat{f} : $\text{new } \hat{f}(x) := \hat{f}(x) + \lambda \hat{f}^b(x)$

(c) update the residuals:

$$\text{new } r_i := r_i - \lambda \hat{f}^b(x)$$

3) \rightarrow model: $\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x)$



$$\frac{\lambda=1}{\lambda=0.1}$$

x:	135	141	156	132	181
y:	11	12	11	13	14

\hat{f}	\hat{f} :	4	13	11	13	13	= \hat{f}^1
r:	r:	0	-1	0	0	+1	
	r:	9.9	10.7	9.9	11.2	12.2	

Ada Boost for binary classification

Dataset: (X, Y) $X = (x_1, \dots, x_N)$, $Y = (y_1, \dots, y_N)$, $x_i \in \mathbb{R}^P$, $y_i \in \{-1, 1\}$

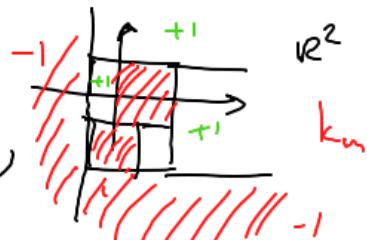
• $C_0(x) = 0$

• $C_m(x) = C_{m-1}(x) + \alpha_m k_m(x)$

$\alpha_m > 0$, k_m - classifier, $k_m: \mathbb{R}^P \rightarrow \{-1, 1\}$

How to choose α_m, k_m ? — To minimize the total error E of C_m ,

$$E = \sum_{i=1}^N e^{-y_i C_m(x_i)} = \sum_{i=1}^N e^{-y_i C_{m-1}(x_i)} e^{-y_i \alpha_m k_m(x_i)} = \sum_{i=1}^N w_i^{(m)} e^{-y_i \alpha_m k_m(x_i)}$$



e.g. $y_i = -1$, $C_m(x_i) = -10 \rightarrow e^{-y_i C_m(x_i)} = e^{-10}$ - very small
 $C_m(x_i) = -0.1 \rightarrow e^{-y_i C_m(x_i)} = e^{-0.1}$ - < 1 , but close to 1
 $C_m(x_i) = 1 \rightarrow e^{-y_i C_m(x_i)} = e^1 = e$
 $C_m(x_i) = 3 \rightarrow e^{-y_i C_m(x_i)} = e^3$

$$E = \sum_{i=1}^N w_i^{(m)} e^{-y_i d_m k_m(x_i)} = \sum_{i: y_i = k_m(x_i)} w_i^{(m)} e^{-d_m} + \sum_{i: y_i \neq k_m(x_i)} w_i^{(m)} e^{d_m} =$$

$$= \underbrace{\left(\sum_i w_i^{(m)} \right)}_{\text{does not depend on } x} e^{-d_m} + \sum_{i: y_i \neq k_m(x_i)} w_i^{(m)} (e^{d_m} - e^{-d_m})$$

$$= \underbrace{(e^{d_m} - e^{-d_m})}_{\text{does not depend on } k \text{ and } \neq 0} \underbrace{\sum_{i: y_i \neq k_m(x_i)} w_i^{(m)}}_{\text{does not depend on } d_m, \text{ only on } k_m}$$



→ k_m should be chosen to minimise $\sum_{i: y_i \neq k_m(x_i)} w_i^{(m)}$, i.e., weighted error.

once k_m is chosen, we find d_m to minimise E

$$\leadsto d_m = \frac{1}{2} \ln \frac{\sum_{i: y_i = k_m(x_i)} w_i^{(m)}}{\sum_{i \neq k_m(x_i)} w_i^{(m)}}$$

Drawback: excessive weight assigned to outliers



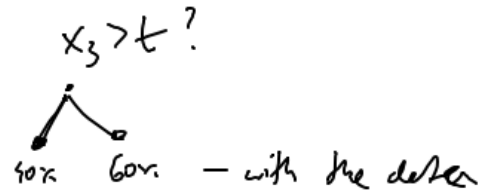
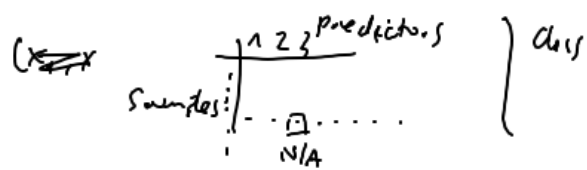
Handling missing data:

- ignore records with missing attribute values
(or ignore an attribute)

drawback:

- could result in bias
- smaller dataset

- 'guess' the missing attribute value
- treat missing values as a special category
- trees: split the sample (used in C4.5 by Ross Quinlan)
C5.0



Trees 'ask' ~~the~~ question about just one attribute at a time, so missing data is less problematic than ~~in~~ in the case of other methods. (but not in scipys)

Trees may be used to find important attributes.

LDA, QDA (Linear / Quadratic Discriminant Analysis) $X = (x_1, \dots, x_N), x_i \in \mathbb{R}^p \quad y_i \in \{1, 2, \dots, K\}$

Model distribution of X separately in each of the response classes (i.e., given Y)
and use Bayes Theorem

$$\Pr(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

↑ predicted
these are probabilities that
given sample is in the class k
($k=1, 2, \dots, K$)

How to estimate f_k ?

we assume that f_k is of some special form
and try to estimate the coefficients →



$f_1 \sim N(\dots)$

$$\pi_k = \Pr(Y=k) = \frac{\#\{i: y_i = k\}}{N}$$

estimate

$f_k(x) = \text{"Pr}(X=x | Y=k)\text{" density function}$

$$\Pr(X \in A | Y=k) = \int_A f_k(x) dx$$