

## Bagging (Bootstrap aggregation)

Given: TS  $S = \{(x_1, y_1) \dots (x_n, y_n)\}$  iid  $\sim (x, y)$

Goal: Estimate  $f(x) = E(y|x)$

$\hat{f}(x, S)$  estimate for  $f(x)$  based on training sample  $S$  (for example CART tree)

$S_1^*, \dots, S_B^*$  Bootstrap samples drawn from  $S$   
(Bootstrap sample = random sample of size  $n$  drawn with replacement)

### Bagged estimate

$$\hat{f}_{\text{Bag}}(x, S) = \frac{1}{B} \sum_i \hat{f}(x, S_i^*)$$

### (Bogus) motivation

If we had  $B$  training samples  $S_1, \dots, S_B$  of size  $n$  all drawn from  $(x, y)$  then we could generate  $B$  prediction rules and average them

$$\bar{f}(x) = \frac{1}{B} \sum_i \hat{f}(x, S_i)$$

- $\bar{f}$  would have lower variance than  $\hat{f}(x, S)$
- Lacking additional training samples we use Bootstrap samples.

## Empirical observation

Bagging can lead to improved performance for prediction rules that depend on the training sample in non-smooth way.

(for example CART)

**Conjecture:** Bagging CART trees helps because bagged trees better approximate smooth responses.

**Note:** Interpretability is lost

**Note:** Bagging is often applied to fully grown trees — no pruning

**Note:** Can use "out-of-bag" error estimation instead of CV

$$\hat{f}_{\text{OOB}}(x_i, S) = \text{ave}_j (\hat{f}(x_i, S_j^*) \mid (x_i, y_i) \notin S_j^*)$$

$$\text{EPE} = \frac{1}{n} \sum (y_i - \hat{f}_{\text{OOB}}(x_i, S))^2$$

## Boosting regression trees

**Given:** TS  $(x_1, y_1) \dots (x_n, y_n) \sim \text{iid} (x, y)$

**Goal:** Estimate  $f(x) = E(y \mid x)$

Initially:  $y_i = y_i \quad i=1 \dots n$

Define  $T(\underline{x}, \underline{y})$  CART tree trained on  $(\underline{x}_1, y_1) \dots (\underline{x}_n, y_n)$  with  $K$  splits ( $K$  usually 1...3)

Boosting algorithm

For  $i = 1 \dots M$  {

$$T_m(\underline{x}) = T(\underline{x}, \underline{y}) \quad (1)$$

$$\underline{y} = \underline{y} - \alpha (T_m(\underline{x}_1) \dots T_m(\underline{x}_n)) \quad (2)$$

}

$$\hat{f}(\underline{x}) = \alpha \sum_{m=1}^M T_m(\underline{x}) \quad \alpha \text{ between } 0.1 \text{ and } 0.01$$

- (1) Apply "base learner" to current residuals  
(2) update residuals

→ Build tree  $T_1$  on  $(\underline{x}_1, y_1) \dots (\underline{x}_n, y_n)$

→ Compute predicted values  $T_1(\underline{x}_1) \dots T_1(\underline{x}_n)$

→ Form new "training data"

$$(\underline{x}_1, y_1 - \alpha T_1(\underline{x}_1)) \dots (\underline{x}_n, y_n - \alpha T_1(\underline{x}_n)).$$

## Notes

$K=1$  Base learner generates a tree with one split and two leaves ("stump")

$\Rightarrow \hat{f}(\underline{x})$  will be additive with piece wise constant coordinate functions

⇒ Alternative way of learning additive model with piece wise constant coordinate functions:

- Define dictionary  $\mathcal{B} = \mathcal{B}_1 \cup \mathcal{B}_2 \cup \dots \cup \mathcal{B}_p$   
 $\mathcal{B}_j$  set of basis functions that depend only on predictor  $j$

$$\mathcal{B}_j = \{ B_{jk}(z) \quad k=1 \dots n \}$$

$$B_{jk}(z) = \mathbb{I}(z_j \geq x_{kj})$$

Note: These functions are not linearly independent but that does not matter

We have such models before (for piece wise linear basis functions (using forward stepwise regression).

What's different

- Boosting uses "stagewise", not "stepwise" regression.

When a new basis function is added to the model, coefficients of previous basis functions are **not** readjusted.

- The least squares coef of the new basis function is multiplied by the learning rate.