

# Bagging (Breiman, 1996)

## Bias-Variance Tradeoff

$$\mathbb{E}[(h_0(x) - y)^2] = \underbrace{\mathbb{E}[(h_0(x) - \bar{h}(x))^2]}_{\text{variance}} + \underbrace{\mathbb{E}[(\bar{h}(x) - \bar{y}(x))^2]}_{\text{bias}} + \underbrace{\mathbb{E}[(\bar{y}(x) - y)^2]}_{\text{noise}}$$

Bagging Reduces Variance

$$h_0 \rightarrow \bar{h}$$

## Weak Law of Large Numbers

i.i.d. RVs  $x_i$   $\mathbb{E}[x_i] = \bar{x}$

$$\frac{1}{m} \sum_{i=1}^m x_i \rightarrow \bar{x} \text{ as } m \rightarrow \infty$$

i.i.d.  $D_1, \dots, D_m \sim P(X, Y)$

$$\frac{1}{m} \sum_{i=1}^m h_i \rightarrow \bar{h} \text{ as } m \rightarrow \infty$$

## Bagging Algo.

1. Sample  $m$  datasets:  $D_1, \dots, D_m$  from u.a.r. w/ replacement
2. For each  $D_i$  train  $h_i(\cdot)$
3.  $\hat{h}(x) = \frac{1}{m} \sum_{i=1}^m h_i(x)$

Adv.

- Predicted var  $\rightarrow$  decreases uncertainty
- unbiased of test error:

$$(x_i, y_i) \in D \quad S_i = \{D_k \mid (x_i, y_i) \notin D_k\} = \frac{1}{m} P_i \cap$$

$$\hat{h}_i(x) = \frac{1}{|S_i|} \sum_{k \in S_i} h_k(x)$$

$$\text{out-of-bag error} \quad \epsilon_{\text{OOB}} = \frac{1}{m} \sum_{(x_i, y_i) \in D} (\hat{h}_i(x_i) - y_i)^2$$

## Solution: Bagging (Bootstrap Aggregating)

$$Q(X, Y | D) \approx P(X, Y)$$

$$Q((x_i, y_i) | D) = \frac{1}{n} \quad \forall (x_i, y_i) \in D \quad n = |D|$$

$$D_i \sim Q$$

Claim: Samples from  $Q$  aren't i.i.d.

are from  $P$

(Proof)

$$P(X=x_i) = P_i \quad \Omega = x_1, \dots, x_n$$

$$Q(X=x_i | D) =$$

$$\sum_{k=1}^n \underbrace{\binom{n}{k} p_i^k (1-p_i)^{n-k}}_{\substack{\text{prob. of } k \\ \text{copies of} \\ x_i \text{ in } D}} \times \underbrace{\frac{k}{n}}_{\substack{\text{prob. of} \\ \text{picking} \\ x_i}}$$

$$= \frac{1}{n} \sum_{k=1}^n \binom{n}{k} p_i^k (1-p_i)^{n-k} \frac{k}{n} = \frac{1}{n} \sum_{k=1}^n \binom{n}{k} p_i^k (1-p_i)^{n-k} k$$

$$\mathbb{E}[B(p_i, n)] = p_i n$$

## Random Forests

Algo.

1. Sample  $D_1, \dots, D_m$  u.a.r. w/ replacement
2. Train  $h_i$  (decision tree) for  $D_i$ 
  - randomly subsample  $k \in d$  features (w/o replacement)
3.  $\hat{h}(x) = \frac{1}{m} \sum_{j=1}^m h_j(x)$

Adv.

1. 2 hyperparams:  $\textcircled{m}$   $\textcircled{k}$ 
  - $\downarrow$
  - as large as comp. possible
2. Not much preprocessing necessary for DTs

## Variants

$$1. D_j = \{D_j^A, D_j^B\}$$

build the tree  $\downarrow$  make prediction at leaves  
\* statistically consistent

2. Don't grow every tree fully
  - prone based on leave-out