



We fit M models using same algorithm:

$$g_{(1)} = A(D_{(1)}, \mathcal{H}), g_{(2)} = A(D_{(2)}, \mathcal{H}), \dots, g_{(M)} = A(D_{(M)}, \mathcal{H})$$

we average all of these models
$$g_{BAG} := \frac{g_{(1)} + g_{(2)} + \dots + g_{(M)}}{M}$$

$$MSE = \sigma^2 + E_x[Bias[g_{BAG}]^2] + E_x[Var[g_{BAG}]] \approx \sigma^2 + E_x[Var[g_{BAG}]] \approx \sigma^2 + \frac{1}{M} E_x[Var[g_{(1)}]]$$

if \mathcal{H} sufficiently complex relative to f
e.g. trees M large $\epsilon(0,1)$

Another benefit of bagging: validation for free using "out-of-bag (oob) validation".

$$D_{oob(1)} := D \setminus D_{(1)} \text{ is a set of about } \frac{1}{3} n. \text{ The } g_{(1)}(D_{oob(1)}) \text{ will give honest predictions.}$$

\uparrow \uparrow
 $D_{test(1)}$ $D_{train(1)}$

$$D_{oob(2)} := D \setminus D_{(2)} \text{ is a different set of } \approx \frac{1}{3} n. \text{ The } g_{(2)}(D_{oob(2)}) \text{ is honest.}$$

$$\vdots$$

$$D_{oob(M)} := D \setminus D_{(M)} \dots \dots \dots$$

How do we get validation for g_{BAG} ?

$$\hat{y}_{i,oob} := \text{Avg (only } g_{(m)} \text{ predictions where } i \text{ is oob).}$$

Each obs is oob $\approx \frac{1}{3} M$. Since M large, each $\hat{y}_{i,oob}$ will be accurate.

OOB validation approximately $\approx K=2$ -fold CV.

Advantages to bagging

- ① Obliterates Bias if an A with complex \mathcal{H} is employed (e.g. trees).
- ② Reduces Variance Substantially.
- ③ Free Validation during the fitting step.

$F \propto$

Assume we are using trees

$$MSE = \sigma^2 + \underbrace{\rho}_{\text{Correlation}} E_x[Var[g_{(1)}]]. \text{ How can we make MSE smaller?}$$

ρ = Avg Correlation between the trees each built with a different bootstrap sample.
 How can we further de-correlate the trees during tree construction?
 What if during each node's construction you only split on a subset of features
 of size $p_{try} < p$ i.e. $\{j_1, j_2, \dots, j_{p_{try}}\} \subset \{1, 2, \dots, p\}$?

This would make the trees more different, hence ρ would decrease.
 Amazingly this doesn't increase bias too much. Random Forests, RF (Breiman, 2001).