ind
bootstrap
sample $\to$ $\boxed{D_{(1)}}$          this has $\approx 2/3$ obs. of $D$, $\approx$ left

$D$ $\xrightarrow{\text{ind bootstrap}}$ $\boxed{D_{(2)}}$

$\xrightarrow{\text{ind bootstrap}}$ $\boxed{D_{(m)}}$

We fit M models using same algorithm;
$$g_{(1)} = A(D_{(1)}, H), \; g_{(2)} = A(D_{(2)}, H), \cdots, g_{(m)} = A(D_{(m)}, H)$$

We all of these models $\quad g_{BAG} = \dfrac{g_{(1)} + g_{(2)} + \cdots + g_{(m)}}{M}$

$$MSE = \sigma^2 + E_x[Bias[g_{BAG}]^2] + E_x[Var[g_{BAG}]]$$

$\quad$ if H sufficiently complex relative to f e.g.
$$\approx \sigma^2 + E_x[Var[g_{BAG}]] \qquad \text{trees.}$$
$\quad$ M large $\quad \to G[0,1]$
$$\approx \sigma^2 + \rho \, E_x[Var[g_{BAG}]]$$

---

Another benefit of bagging : Validation for free using "Out-of-bag (oob) Validation."

$\quad$ ① test(1) $\qquad\qquad$ ① train(1)

$D_{oob(1)} := D / D_{(1)}$ is a set of above $\frac{1}{3}n$. Thus
$\quad g_{(1)}(D_{oob(1)})$ will give honest prediction.

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

$\vdots$

$D_{oob(m)} := D / D_{(m)}$ "   "   "

How do we get validation for $g_{BAG}$ ?

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

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

OOB validation $\qquad \approx k = 2 - fold$ CV

Advantage to bagging

① Obliterates Bias if an A wish complex H is employed (e.g trees)

② Reduce variance substantially.

③ free validation during the fitting step.

Assume we are using tree.

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

$\rho$ = Avg correlation between two trees each built with a different bootstrap sample.

How can we further de-correlate the tree during tree construction ?

What if during each nodes construction, you only split on a subset of features of size $P_{try} < p$. i.e. $[j_1, j_2, ..., j_{A_{try}}] \subset [1, 2, ..., p]$ ?

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