

Back to Math 241/368... X_1, \dots, X_n dependent rv's

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n)$$

Assume $\sigma^2 = \text{Var}[X_i]$ same

$$\text{Var}[\bar{X}] = \frac{1}{n^2} \left(\text{Var}\left[\sum X_i\right] \right) = \frac{1}{n^2} \left(\text{Var}[X_1] + \dots + \text{Var}[X_n] + \sum_{i \neq j} \text{Cov}[X_i, X_j] \right)$$

$$\text{Cov}[X_i, X_j] = E[X_i X_j] - E[X_i]E[X_j] \stackrel{\text{assume}}{=} \sigma_{ij} \text{ (same)}$$

$$= \frac{1}{n^2} \left(n \sigma^2 + (n^2 - n) \sigma_{ij} \right) = \frac{1}{n} (\sigma^2 + (n-1) \sigma_{ij})$$

$$\rho := \text{Cov}[X_i, X_j] = \frac{\text{Cov}[X_i, X_j]}{\text{sd}[X_i] \text{sd}[X_j]} = \frac{\sigma_{ij}}{\sigma \sigma} \Rightarrow \sigma_{ij} = \sigma^2 \rho \in [-1, 1]$$

but in our case, $\in (0, 1)$

$$= \frac{1}{n} (\sigma^2 + (n-1) \sigma^2 \rho) = \frac{1}{n} (\sigma^2 + n \sigma^2 \rho - \sigma^2 \rho) \quad \text{Check: } \rho = 0$$

$$= \frac{1}{n} (\sigma^2 (1 - \rho) + n \sigma^2 \rho) = \underbrace{\sigma^2 \rho + \frac{1 - \rho}{n} \sigma^2}_{\rho > 0} \stackrel{\checkmark}{=} \frac{\sigma^2}{n}$$

$$\text{Also } > \frac{\sigma^2}{n}$$

Now let's apply this to the MSE decomposition formula for a model average g_{avg} where the constituent models are zero bias i.e. overfit:

$$\text{MSE} = \sigma^2 + E_x [\text{Bias}[g_1]^2] + E_x [\text{Var}[g_{\text{avg}}]]$$

$$= \sigma^2 + E_x [\text{Var}[g_{\text{avg}}]]$$

$$\stackrel{n \rightarrow \infty}{=} \sigma^2 + E_x \left[\rho \text{Var}[g_1] + \frac{1 - \rho}{n} \text{Var}[g_1] \right]$$

$$\stackrel{\downarrow}{=} \sigma^2 + \rho E_x [\text{Var}[g_1]] \quad \rho \in (0, 1)$$

$$< \sigma^2 + E_x [\text{Var}[g_1]] \quad \text{which means we get a discount almost for free}$$

"Raise oneself by their Bootstrap" means getting something for nothing. The slight cost is that you're model-averaging meaning you lose visibility into how your model is coming up with its predictions.

There is another bonus to this bagging algorithm: free validation.

$$\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{test}}$$

In a bootstrap sample:

$$\mathcal{D} = \mathcal{D}_{(1)} \cup (\mathcal{D} \setminus \mathcal{D}_{(1)})$$

rows that are missing in the bootstrap sample provides a natural Dtest for g_1

$$\mathcal{D} = \mathcal{D}_{(2)} \cup (\mathcal{D} \setminus \mathcal{D}_{(2)})$$

...

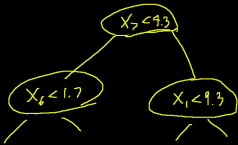
$$\mathcal{D} = \mathcal{D}_{(m)} \cup (\mathcal{D} \setminus \mathcal{D}_{(m)})$$

Since M is large, $M \gg n$. Since each bootstrap sample has $\sim 1/3$ left out, each observation in \mathcal{D} has $\sim M/3$ models that are built without seeing it. We can predict that observation on these $M/3$ models. Then we do this for all $1 \dots n$ observations. These "out of sample" predictions \hat{y} -hats are called "out of bag" (OOB). This procedure is called "bootstrap validation". Theoretically, they say bootstrap validation is approximately equivalent to $K=2$ cross fold validation.

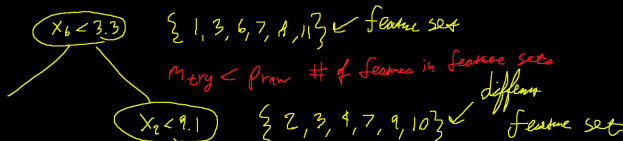
Why use trees in the bagging procedure? They have no bias. There is no need to specify transformations of the p_{raw} features. Then the bagging gives you a discount on the variance. And the bagging gives you validation without further work.

$$\text{MSE} = \sigma^2 + \underbrace{\rho}_{< 1} E_x [\text{Var}[g_1]]$$

Leo Breiman hits the scene again 7 yr later. In 2001, he said "we can do better" i.e. we can make ρ even smaller. In the context of bagged trees, we can do the following during the individual tree construction:



In the original 1984 CART algorithm, at every split point, every feature $1 \dots p$ is searched to find the best split point. What if at every split point, we take a random *subset* of the p features?



If you build trees like this on the bootstrapped samples, you decorrelate the tree models even more! You get a further discount on the MSE via the variance term. And... bias doesn't suffer that much!

Default $mtry$ for regression is $\text{floor}(\text{praw} / 3)$. Default $mtry$ for classification is $\text{floor}(\text{sqrt}(\text{praw}))$. The hyperparameter $mtry$ is definitely cross-validated over because it really matters.

How to name this procedure? There's lots of trees. They're all random... "Random Forests" (RF).

