## Bias-Var Decomp

One $(X, \vec{y})$, one $\vec{x}^\alpha$:    $MSE(\vec{x}^\alpha) = \sigma^2 + Bias[g(\vec{x}^\alpha)]^2$

$$\overbrace{(g(\vec{x}^\alpha) - f(\vec{x}^\alpha))^2}$$

One $X$, many $\vec{y}$, one $\vec{x}^\alpha$:   $MSE(\vec{x}^\alpha) = \sigma^2 + Bias[g(\vec{x}^\alpha)]^2 + \underbrace{Var(g(\vec{x}^\alpha))}_{\text{bias-var var}}$

Most general situation:    $\underbrace{(E[g(\vec{x}^\alpha)] - f(\vec{x}^\alpha))^2}$

Many $X$, many $\vec{y}$, many $\vec{x}^\alpha$    $MSE = \sigma^2 + E_X\left[Bias[g(\vec{x})]^2\right] + E_X\left[Var(g(\vec{x}))\right]$

---

How to make $g$ better.  ("Bagging" 1998)

Imagine many models $\overset{g_1, \dots, g_m}{\text{averaged together:}}$

$$g_{avg} := \frac{g_1 + \dots + g_m}{M} \qquad \text{What is } MSE?$$

$$MSE = \sigma^2 + E_X\left[Bias[g_{avg}]^2\right] + E_X\left[Var[g_{avg}]\right]$$

$$= \sigma^2 + E_X\, E\left[\frac{g_1 + \dots + g_m}{M} - f\right]^2 + E_X\, Var\left[\frac{g_1 + \dots + g_m}{M}\right]$$

$$= \sigma^2 + E_X\, E\left[\frac{1}{M}\underbrace{\left((g_1 - f) + (g_2 - f) + \dots + (g_m - f)\right)}_{Bias(g_1) + Bias(g_2) + \dots + Bias(g_m)}\right]^2 + E_X\, \frac{1}{M^2} Var\left[g_1 + \dots + g_m\right]$$

Assume ① Bias of all $g_1, \dots, g_m$ same and ② $g_1, \dots, g_m$ independent

$$= \sigma^2 + E_X\, Bias[g_1]^2 + E_X\left[\frac{Var[g_1]}{M}\right] \longrightarrow = \sigma^2 + E_X\, Bias[g_1]^2 = \sigma^2$$

If $M \to \infty$

If Bias is low i.e. overfit model

Can we do this?   A little...

(I) It is easy to drive bias → 0. You just overfit. It will fit $f$ perfectly on average. But you introduce large variance.

~~The variance~~ can be driven to ~~zero~~ if ... $g_1$ ... are indep.

(II) How to get $g_1, ..., g_m$ independent given one dataset $D$?

You can't!

Ever Beginning in 1998... [Leo]

Imagine sampling $D$ with replacement of size $n$.

$D_{(1)} = $ sample $(D)$ (non-parametric bootstrap sample) → we will explain the meaning of this term shortly

$D_{(1)}$ has about $\frac{2}{3}$ of the rows of $D$ and $\frac{1}{3}$ duplicates

then do it again ... and again

$D_{(1)}, D_{(2)}, \ldots, D_{(M)}$

Each $D_{(b)}$ is a little bit different from the others since it has slightly different data.

Now build a model on each bootstrap sample

$$g_b = A_{zt}(D_{(b)}) \qquad b = 1, \ldots, M$$

and average...

$$g_{avg} = \frac{1}{M} \sum_{b=1}^{M} g_b \qquad \text{does the regression via averaging}$$

this is called bootstrap + aggregation = "bagging"

It is a "meta-algorithm" as in its something done on top of an $A$ and can yield results for any $A$.

Is $\text{Var}[g_{avg}] = \dfrac{\text{Var}(g_1)}{M}$ s.t. $M \to \infty$ $\text{Var}[g_{avg}] \to 0$?

No! $g_1, \ldots, g_m$ are not independent

↙ since a lot of the same data, therefore $g(1), g(2)$ will be similar

How similar?



$D_{(1)}$ , $D_{(2)}$

March 28? exercise...

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n}\sum x_i\right) = \frac{1}{n^2}\left(\text{Var}(X_1) + \text{Var}(X_2) + \ldots + \text{Var}(X_n) + \sum_{i \neq j} \text{Cov}(X_i, X_j)\right)$$

$\sigma^2 = \text{Var}(X_i)$ Same for all $i$

Assume $\sigma_{ij} = \text{Cov}(X_i, X_j)$ is the same for all $i, j$

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

let $\rho = \dfrac{\sigma_{ij}}{\sigma_i \sigma_j} = \dfrac{\sigma_{ij}}{\sigma^2}$ (correlation) ↗

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

$$= \rho\sigma^2 + \frac{1-\rho}{n}\sigma^2$$

If $\rho \to 0 \Rightarrow \sigma_{ij} \to 0$ and all $X_i$'s indep $\Rightarrow \text{Var}(\bar{X}) \to \dfrac{\sigma^2}{n}$

otherwise $\text{Var}(\bar{X}) > \dfrac{\sigma^2}{n}$

In expansion, each $g(b)$ has some corr. with another $g(b)$.

Variance term

$$\Rightarrow MSE = \sigma^2 + E_x\left[bias(\hat{g})^2\right] + E_x\left[\rho \, Var(\hat{g}) + \frac{1-\rho}{m} Var(\hat{g})\right]$$

Now with $m \to \infty$ bootstrap samples

$$MSE = \sigma^2 + E_x\left[bias(\hat{g})^2\right] + E_x\left[\rho \, Var(\hat{g})\right]$$

If $\rho < 1$, the gang does better than a single $\hat{g}$

And $\rho < 1$ for bootstrapping $\Rightarrow MAGIC!$

Bootstrap: raise yourself up by pulling on your bootstrap.

Get someone for another!!

## Validation for models that implement bagging

Usually $D = D_{train} \cup D_{test}$

Here, you can imagine

$\frac{2}{3}$ of runs in $D_1 \longrightarrow D_{test,(1)}$

$D = D_{(1)} \cup \left(D \setminus D_{(1)}\right) \longleftarrow$ the $\frac{1}{3}$ left over! "Out of bag" (OOB)

$D = D_{(2)} \cup D_{(2)}^{test}$

$\vdots$

etc.

these are $n$ rows and $M \gg 4$ models trained on bootstrap samples

where each row has a $\frac{1}{3}$ chance of being left out of each model

⟹ each row has $\approx \frac{M}{3}$ models that are not built with that row.

Thus, we can predict on these $\frac{M}{3}$ models to get an "OOB" estimate. ∧(with $M$ large, no problem to get est's.) ∧It's similar to leaving out the dataset.

Surprisingly, theoretically oob is $\approx$ ∧(cv with) $K=2$.

---

Let's say we bag trees s.t. nodesize $=$ small. and $M=$ large

⟹ $\text{Bias}[\hat{g}] \approx 0$

⟹ $\text{MSE} \approx \sigma^2 + E_x\left[\rho\, \text{Var}(\hat{g})\right]$

### Advantages of Bagged trees

① Low MSE Due to low bias since trees are complex enough to fit $f$.

② No need to specify a model since trees are non-parametric

③ Validation for free!

### DEMO

How can we do better? Make $\rho$ as small as possible!

How can we "decorrelate" the trees?

Modify the regression tree algorithm. Instead of spliting by trying all $p$ features use $\{i_1, i_2, \ldots, i_{p_{or}}\} \subset \{1, 2, \ldots, p\}$

7

same as any split

This will bring $\rho \downarrow$ but increase $bias(g)$ but not too much!

Each tree is much more "random" and there are lots of trees

$\Rightarrow$ a "random forest"! (Brennan, 2001)

<u>Concludes the material needed for your project</u>