

You've seen this before!

Constant in Δ 's

Constant in Δ 's

$$MSE = \sigma^2 + E_X [\text{Var}[f(\vec{x})]] + E_X [\text{Bias}[g(\vec{x})]^2]$$

done

$$Y|\vec{X}=\vec{x} = g(\vec{x}) + (h^*(\vec{x}) - g(\vec{x})) + (f(\vec{x}) - h^*(\vec{x})) + \Delta$$

$$E[(Y - g(\vec{x}) | \vec{X}=\vec{x})^2] = E[(h^*(\vec{x}) - g(\vec{x}) + (f(\vec{x}) - h^*(\vec{x})) + \Delta)^2]$$

$$MSE = E[\underbrace{(h^* - g)^2}_V + \underbrace{(f - h^*)^2}_B + \Delta^2 + 2\underbrace{(h^* - g)\Delta}_V + 2\underbrace{(f - h^*)\Delta}_B + 2\underbrace{(h^* - g)(f - h^*)}_V]$$

$$= E[V^2] + B^2 + \sigma^2 + 2B E[V] = \text{Var}[V] + B^2 + \sigma^2 + 2B E[V] + E[V]^2$$

If algorithm is "unbiased" in its estimate of g then $E[V] = 0$. E.g. OLS is unbiased. (Prove in ECON 387.)

If unbiased...

this should read $\text{Bias}[g]^2$

$$MSE = E_X [\text{Var}[g - h^*]] + E_X [\text{Bias}[h^*]^2] + \sigma^2$$

if not, there's a

fourth component

this should read
 $\text{Var}[g] = E[(g - h^*)^2]$

estimation error component

misspecification error component

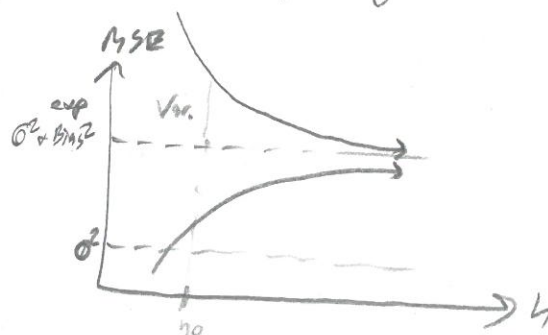
irreducible error component

\Rightarrow the bias-var decomp. you know about the whole semester!

if $h \rightarrow \infty$ $g \rightarrow h^*$ You can prove this for many alg's

\Rightarrow est. error $\Rightarrow 0$

but bias is always the same!

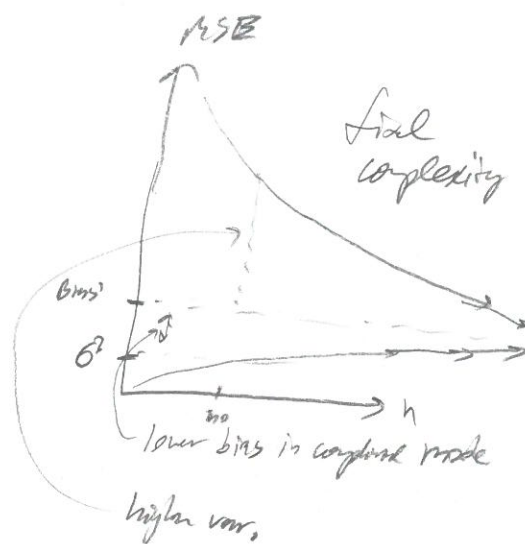
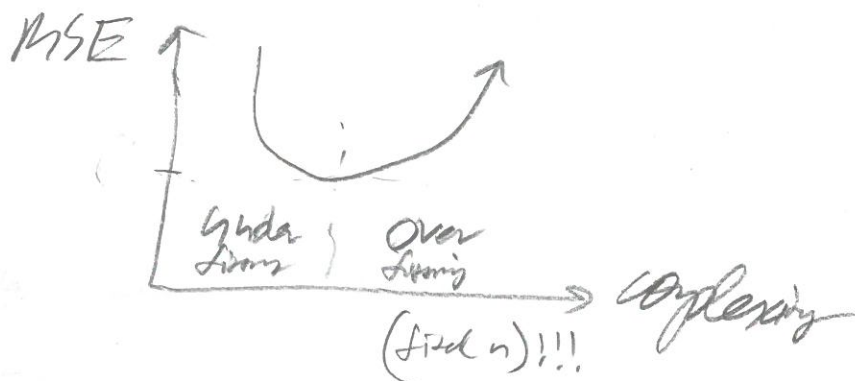


How to minimize bias? Make $h^* \rightarrow f$ by increasing model complexity!

Is there a tradeoff? Yes & No.

the tradeoff is between underfitting (high bias / low var) and overfitting (low bias / high var)

but perfect fitting minimize both



Back to trees...

regression trees have... low bias & high var
 but the ability to get close to capturing f they overfit!

Lemma

A couple of observations

① Many models averaged together is very stable!

We can consider the average model as a model itself!

$$g_{avg} = \frac{g_1 + g_2 + \dots + g_m}{m}$$

3

② The avg of many trees is very low in bias

Can we average many trees together?

But with one $D \Rightarrow$ one tree! It's deterministic! So no...

Unless... we change D or change A to be random.

What if we change D ? Let's sample n rows from D with replacement. This is called a "bootstrap sample". This $D_{(1)}$ is $\approx D$ but a little bit different. It has about $\frac{2}{3}$ the original rows and duplicates of those... regression tree

So let's fit $g_1 = A(H, D_{(1)})$

Let's then draw another bootstrap sample $D_{(2)}$ and build another tree:

$$g_2 = A(H, D_{(2)}), \text{ etc...}$$

g_1, \dots, g_B are all regression trees which are similar. Then...

$g_{\text{avg}} = \frac{g_1 + \dots + g_B}{B}$ is the "average" or "aggregate" tree

Bootstrap + Aggregation = "Bagging" (Breiman, 1994)

"Mean-Algorithm" which leads to "improvements for unstable procedures"

HUGE!

high variance algorithms

demo

(9)

What's going on. Imagine a tree g_t . It has low bias & high variance.

T Many trees averaged $g_{\text{bagged}} = \frac{g_1 + \dots + g_T}{T}$

also has low bias since

$$\begin{aligned} \text{Bias} &:= E[g_{\text{bagged}}] - f = E\left[\frac{g_1 + \dots + g_T}{T}\right] - f \\ &= \frac{1}{T} E[g_1 + \dots + g_T] - \frac{Tf}{T} = \frac{1}{T} (E[g_1 - f + g_2 - f + \dots + g_T - f]) \\ &= \frac{1}{T} \left(\underbrace{(E[g_1] - f)}_{\substack{\text{Bias of} \\ 1^{\text{st}} \text{ tree} \\ (\text{small})}} + \underbrace{(E[g_2] - f)}_{\substack{\text{Bias of} \\ 2^{\text{nd}} \text{ tree} \\ (\text{small})}} + \dots + \underbrace{(E[g_T] - f)}_{\substack{\text{Bias of} \\ \text{last tree} \\ (\text{small})}} \right) = \frac{Tf}{T} (\text{small}) \end{aligned}$$

But the variance term looks like:

$$\text{Var}(g_{\text{bagged}}) = \text{Var}\left[\frac{g_1 + \dots + g_T}{T}\right] = \frac{1}{T^2} \text{Var}\left[\sum_{t=1}^T g_t\right]$$

if g_1, \dots, g_T independent

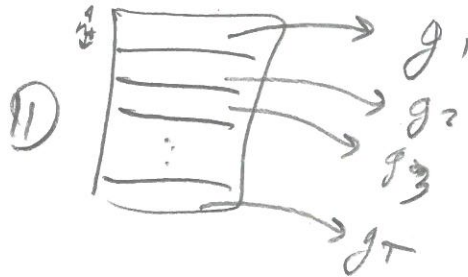
$$= \frac{1}{T^2} \sum \text{Var}(g_i) = \frac{\text{Var}(g_i)}{T}$$

which means as $T \rightarrow \infty$, the variance term vanishes!

$$\Rightarrow \text{MSE} = \sigma^2 + \underbrace{\frac{\text{Var}(g_i)}{T}}_{\text{goes to zero}} + \underbrace{\text{Bias}}_{\text{small}} \approx \sigma^2 \text{ PERFECT!!!}$$

What's the problem?

g_1, \dots, g_T not independent... Why? Usually from same data (bootstrap samples). But if they are from different data, this works!



As long as n is large enough to create an unbiased model then your MSE will be very small.

Usually, you use all D to build g bagged... in which case the ~~models~~ are dependent because bootstrap samples contain a lot of the same data!

Back to Math 241.

If X_1, \dots, X_n identically distr but not independent, what is $\text{Var}[\bar{X}_n]$?

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

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

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

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

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

this is the answer but let's get it into canonical form

Note $\rho_{ij} := \frac{\sigma_{ij}}{\sigma_i \sigma_j} = \frac{\sigma_{ij}}{\sigma^2} \Rightarrow \sigma_{ij} = \rho \sigma^2$

↑↑ same

Note $\rho \geq 0$ $\text{Var}[\bar{X}] \rightarrow \frac{\sigma^2}{n}$ as expected

$$MSE = \sigma^2 + \underbrace{\left(\rho \text{Var}[g_t] + \frac{1-\rho}{T} \text{Var}[g_t] \right)}_{\text{variance term}} + \underbrace{\text{Bias}(g_t)^2}_{\text{low}}$$

$$< \text{Var}[g_t]$$

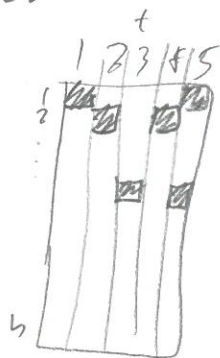
itself since $\rho \neq 1$

this is why bagging works so well

Validation for bagged models (not only trees)

Usually $D = D_{\text{train}} \cup D_{\text{test}}$, build on train, validate on test.
Here, each tree has its own $D_{\text{train},t}$, $D_{\text{test},t}$. why?
 "out of bag (oob)"

Bootstrap sample. This means that each tree can validate itself by predicting on $D_{\text{test},t}$. Now let each tree do so and average by observation.



Validating the first obs is done on trees 1 & 5 and we average the \hat{y} 's. Validation for the second obs. is done on trees 2 & 4 and we avg. the \hat{y} 's.

If T is sufficiently large, all n obs. will be validated. We then compute E_i 's and then an error metric e.g. R^2_{oob} , SE_{oob} , etc.

⇒ Out of bag error est (oob error). It works! Empiric does it. Theoretically it's \approx K-fold CV with $K=2$, but I don't know much about it.