Lec 21  4/30/19  Math 390.4

~~Regression~~
Trees $(y \subseteq \mathbb{R})$  review with $\mathbb{R}$ dense

## Theory of Bias-Variance Decomposition

Recall... $\quad y = g + e = g + (f - g) + \delta$

error due to miss + estimation

error due to ignorance

$$e = y - g = f - g + \delta$$

$$\Rightarrow e^2 = (f - g + \delta)^2$$

What if I want to measure the "mean squared error" (MSE) for a new observation $\vec{x}_*$?

I need to assume a r.v. model somewhere!

Let's assume $\Delta$ is the r.v. where $\delta$ is realized from.
$\vec{x}$, and $f(\vec{x})$ be constant.

$$Y = f(\vec{x}) + \Delta \quad \leftarrow \text{r.v.}$$

this is equivalent to assuming

$$\boxed{I} \quad E[Y | \vec{X} = \vec{x}] = f(\vec{x}) \quad \leftarrow \text{conditional expectation function}$$

$$\Rightarrow E[Y | \vec{X} = \vec{x}] = E[f(\vec{x}) + \Delta \,|\, \vec{X} = \vec{x}] = E[f(\vec{x}) \,|\, ] + E[\Delta \,|\, ] = f(\vec{x})$$

$$\Rightarrow E[\Delta \,|\, \vec{X} = \vec{x}] = 0.$$

the expected $Y$ at any $x$

is $f(\vec{x}) \Rightarrow$ the error due to ignorance must have mean center of $0$.

(II) Also assume variance does not depend on $\vec{x}$. $Var(\Delta | \vec{X} = \vec{x})$

$$= Var(\Delta) = \sigma^2$$
$$\Rightarrow E[\Delta^2] = \sigma^2$$

Not strictly necessary! But makes life easy here.

---

Back to MSE for a new obs. $x^*$. If we knew $f$...

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

$$= E(\Delta^{*2}) = \sigma^2 \quad \text{expectation taken over } \Delta^*$$

If we don't know $f$... then $MSE(\vec{x}^*) \geq \sigma^2$. Proof:

drop this notation...

$$E[(Y - g(\vec{x}))^2 | \downarrow] = E[Y^2 - 2Y g(\vec{x}) + g(\vec{x}^*)^2] \quad \text{the only randomness is in } \Delta^*!$$

$$= E[Y^2] - 2E[Y g(\vec{x})] + E[g(\vec{x})^2]$$

$$= E[(f + \Delta)^2] - 2E[Y]E[g] + E[g^2]$$

$$= E[f^2 + 2f\Delta + \Delta^2] - 2fg + g^2$$

$$= f^2 + \sigma^2 - 2fg + g^2 \quad \longrightarrow \quad \text{expected sqd. errors additive}$$

$$= (f(\vec{x}) - g(\vec{x}^*))^2 + \sigma^2 \geq \sigma^2$$

Now... instead of taking Expectation over just $\Delta^*$, we

take expectation over $\Delta_1, \Delta_2, \ldots, \Delta_u, \Delta^*$ which means

the randomness in $D$ itself.

"Dataset-dataset variability"

Assume!

$\vec{X}, Y$ are r.v.'s

and $(\vec{X}_1, Y_1), (\vec{X}_2, Y_2),$ are all iid relations.

$$MSE(\vec{x}^\alpha) = E_{\Delta_1, \ldots, \Delta_u, \Delta^*}\left[\left(Y^* - g(\vec{x}^\alpha)\right)^2 \mid \vec{X} = \vec{x}^\alpha\right]$$

just a function of $\Delta^*$

is a r.v. now based on $\Delta_1, \ldots, \Delta_u$

independent! why $D \perp (\vec{X}^*, Y^*)$!

$$= E_{\Delta_1, \ldots, \Delta_u, \Delta^*}\left[Y^{*2}\right] - 2E_{\Delta_1, \ldots, \Delta_u, \Delta^*}\left[Y^* g(\vec{x}^\alpha)\right] + E_{\Delta_1, \ldots, \Delta_u, \Delta^*}\left[g(\vec{x}^\alpha)^2\right]$$

$$= E_{\Delta^*}\left[Y^{*2}\right] - 2 E_{\Delta^*}\left[Y^*\right] E_{\Delta_1, \ldots, \Delta_u}\left[g(\vec{x}^\alpha)\right] + E_{\Delta_1, \ldots, \Delta_u}\left[g(\vec{x}^\alpha)^2\right]$$

subscripts dropped everywhere

$$= \left(f(\vec{x}^\alpha)^2 + \sigma^2\right) - 2 f(\vec{x}^\alpha) E(g(\vec{x}^\alpha)) + Var[g(\vec{x}^\alpha)] + E[g(\vec{x}^\alpha)]^2$$

$$= \sigma^2 + \left(E[g(\vec{x}^\alpha)] - f(\vec{x}^\alpha)\right)^2 + Var[g(\vec{x}^\alpha)]$$

$$= \sigma^2 + Bias[g(\vec{x}^\alpha)]^2 + Var[g(\vec{x}^\alpha)]$$

↑
irreducible gen. err.
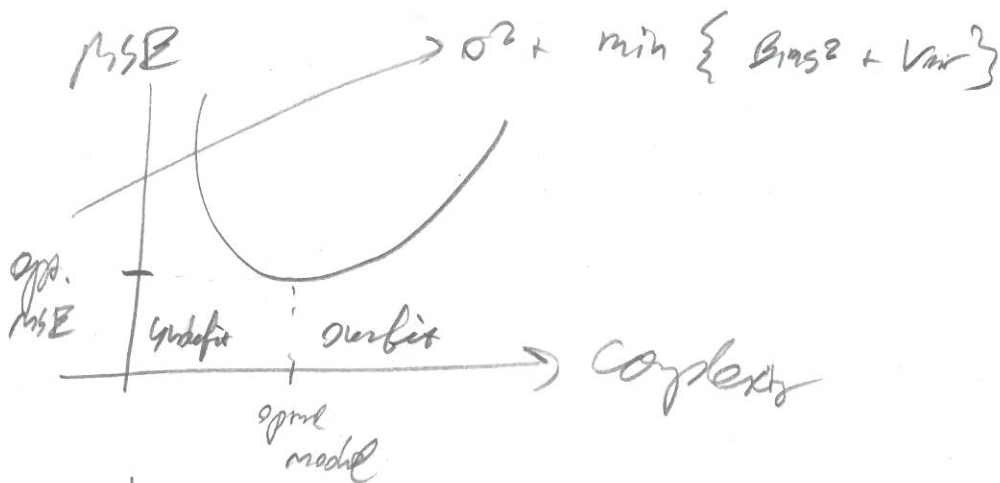
↑
How far is $g$ from $f$ on average?

↑
How variable is $g$ from its center measure, is mean sqd. distance.

One more small point... this is all for one new obs. $\vec{X}^*$.
Let's assume $\vec{X}^*$ is realized from $P(X)$. Then...

$$MSE := E_X\left[MSE(\vec{X}^*)\right] = E_X\left[\sigma^2 + Bias(g(\vec{X}^*))^2 + Var(g(\vec{X}^*))\right]$$

$$= \underbrace{\sigma^2}_{\substack{\text{irred.}\\\text{err.}}} + \underbrace{E_X\left[Bias(g(\vec{X}^*))\right]^2}_{\substack{\text{expected sqd. model}\\\text{bias}}} + \underbrace{E_X\left[Var(g(\vec{X}^*))\right]}_{\substack{\text{expected model}\\\text{variance}}}$$

$$DEMO$$

---

Is there a "bias-variance tradeoff"? Yes and no.
Yes in the extremes. No in not the extremes

MSE $\longrightarrow \sigma^2 + \min\{Bias^2 + Var\}$



opt.
MSE |  underfit | overfit $\longrightarrow$ Complexity
           optimal
           model

What $\wedge$ underfitting? $\underbrace{E_X\left[Bias(g(\vec{X}))\right]^2}_{}\left[E\left[g(\vec{X}) - f(\vec{X})\right]\right]^2$ is high since

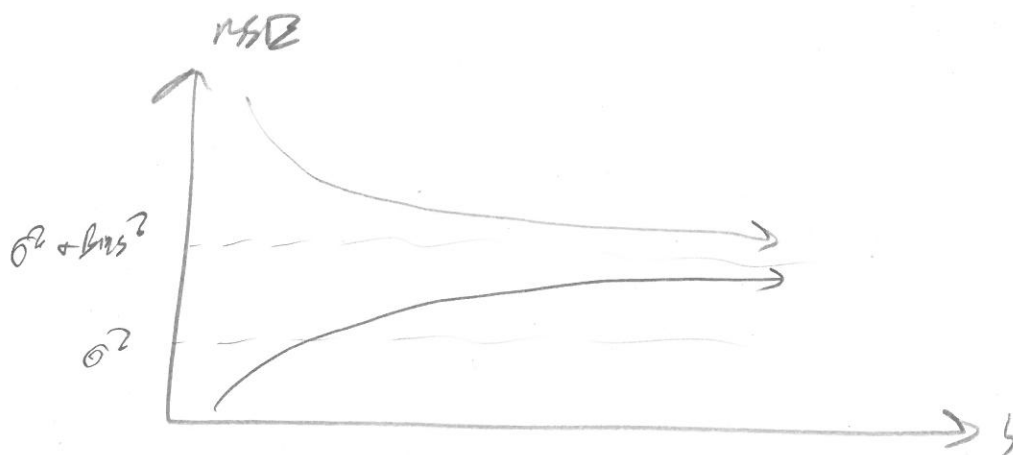$f(x)$ likely way more complex than $g(x)$

$E_X\left[Var(g(\vec{X}))\right]$ is low since if $D$ change, fit doesn't change
too much

What happens during overfitting?

$E_x\left[bias(g(x))^2\right]$ is low since $g(x)$ will be complicated enough to locate $f(x)$ low

$E_x\left[Var(g(x))\right]$ is high since it is tightly fitting $\int$ which change dataset — dataset.

For any fixed complexity, (fixed algorithm)



as $n \to \infty \Rightarrow Var(g(x)) \to 0$ why? But bias does not budge. Why?

So there is a bias variance tradeoff but there is a spot where both together are minimized. Each algorithm has a different curve.