Jan B

Math 650.4                                    Late April

## Theory of Bias-Variance Decomposition

Recall $y = g + e$ ← residual

$$y = g + (f-g) + \delta$$

$\underbrace{\hspace{2cm}}$ error due
to misspecification
& estimation error

$$e = y - g = f - g + \delta$$
$$so \quad e^2 = (f - g + \delta)^2$$

**r.v.**
random
variable

I wan to gauge "mean squared error" (MSE)
for a new observation $\vec{x}^*$ and a model
built from a single data set $D$

────────────────

to do
mean,
as in
expectation,
need
random
variables
from
Math 241

We need to assume a data generating
process (DGP) (random variable model)

### DGP Assumptions
① $D = (X, y)$ is fixed
② $\vec{x}^*$ is fixed
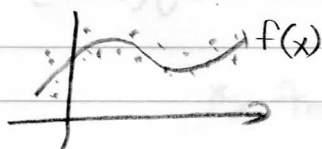③ $\delta$ is drawn from r.v. $\Delta$ with $E[\Delta | \vec{X} = \vec{x}^*] = 0$

$\delta$ is a r.v.

④ $Y = f(\vec{x}^*) + \Delta$
    $Y$ is random through $\Delta$ only
    $E$ is also a r.v. through $\Delta$.

**r.v.  realization**
$E$        $e$
$Y$        $y$

$$\Rightarrow \underset{\uparrow}{E}[Y | \vec{X} = \vec{x}^*] = f(\vec{x}^*) \quad \text{conditional expectation}$$
$$\text{expectation} \qquad\qquad\qquad \text{function}$$

$\left(\begin{array}{c}\text{expectation}\\ \text{of } Y \text{ is } f\end{array}\right)$

Theory of Bias-Variance Decomposition

Only source of error is $\Delta^*$

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

everywhere from now on, but notation suppressed

Imagine $g$ is perfect, i.e. $g = f$

$$= E_{\Delta^*}\left[(Y^* - f(\vec{x}^*))\right] = E_{\Delta^*}[\Delta^2] = \sigma^2$$

irreducible squared error (theoretical best MSE)

Assume $\Delta$ is homoskedastic

$$Var[\Delta \mid X] = Var[\Delta] = E[\Delta^2] - (E[\Delta])^2 = E[\Delta^2]$$

this is 0 $= \sigma^2$

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

$g$ is const so $g$ is indep of $\Delta^*$

$g = \mathcal{A}_{\mathcal{H}}(\mathbb{D})$

$g = \mathcal{A}_{\mathcal{H}}(X, y)$

since $\mathbb{D}$ is fixed, $g$ is fixed

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

$$E\left[f(\vec{x}^*) + 2f(\vec{x}^*) + \Delta^{*2}\right]$$

$$f(\vec{x}^*)^2 + \sigma^2 \qquad \text{by fact}$$

bias $= g - f$ of $g$

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

upward bias if $g > f$

$$= \sigma^2 + Bias[g(\vec{x}^*)]^2 \qquad \leftarrow \quad Bias[g(x)] = g(x) - f(x)$$

downward bias if $g < f$

$$\left(\text{without homoskedicity, would be}\right.$$
$$\left. = \sigma^2(\vec{x}^*) + Bias[g(\vec{x}^*)]^2 \right)$$

different MSE for every $x^*$ $\leftarrow$ MSE is a function of $x^*$

Modify Assumption 1

① X is fixed, y is drawn from $f(\vec{x}) + \Delta$

$$\mathbb{D} = \left( [x], [y] \right) \qquad \vec{y} = f(\vec{x}^*) + \vec{\delta}$$

$\vec{y}$ is drawn from $\vec{Y}_s$ r.v.'s

$\delta_1, \delta_2, ..., \delta_n$ is drawn from $\Delta_1, \Delta_2, ..., \Delta_3$

$g = \mathcal{A}(x, y)$

↳ this is called <u>Data Set — Data Set Variability</u>

Only source of error is $\Delta^*, \Delta_1, \Delta_2, ..., \Delta_n$

$$MSE(\vec{x}^*) = E_{\Delta_1, ... \Delta_n, \Delta^*}\left[ E^2 | \vec{X} = \vec{x}^* \right]$$

$$= E_{\Delta_1, ..., \Delta_n, \Delta^*}\left[ (Y^* - g(\vec{x}^*))^2 | \vec{X} = \vec{x}^* \right]$$

everywhere, but notation dropped

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

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

$$E_{\Delta^*}\left[ f(\vec{x}^*) + 2f(\vec{x}^*)\Delta^* + \Delta^{*2} \right] - 2f(\vec{x}^*) E\left[ g(\vec{x}^*) \right]$$

$$+ Var\left[ g(\vec{x}^*) \right] + E\left[ g(\vec{x}^*) \right]$$

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

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

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

Now, change Assumption 1 again,

made X
random
too
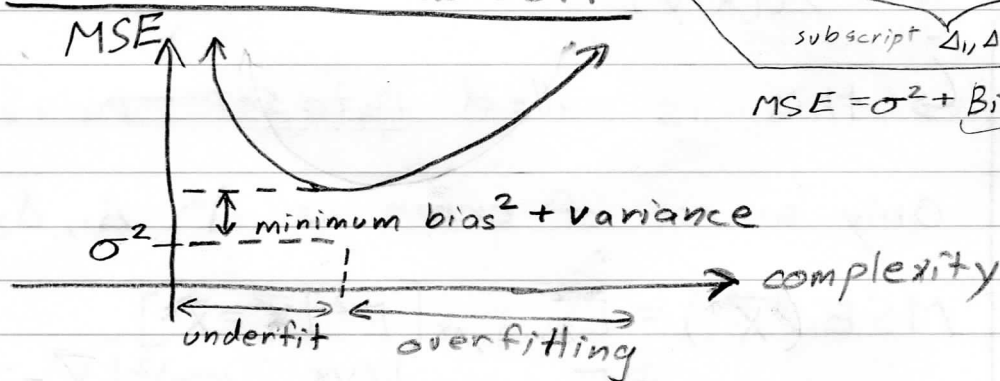
① X is random, drawn from r.v. $\vec{X}$

Y is random, drawn from $f(\vec{x}) + \Delta$

r.v.'s $\Delta_1, ..., \Delta_n$

Now, MSE $= E_x\left[ MSE(\vec{X}) \right]$

↑
generalization
error

mean model

$= \sigma^2 + E_x\left[ Bias\left( E[g(\vec{x})] \right)^2 \right]$

$+ E_x\left[ Var[g(\vec{x})] \right]$

subscript $\Delta_1, \Delta_2, ..., \Delta_n$

MSE $= \sigma^2 + Bias^2 + Variance$

want to
minimize
this part

Bias-Variance Tradeoff



MSE

$\sigma^2$ — minimum bias$^2$ + variance
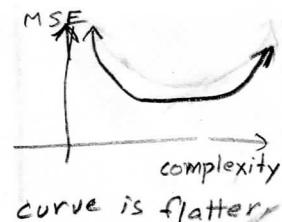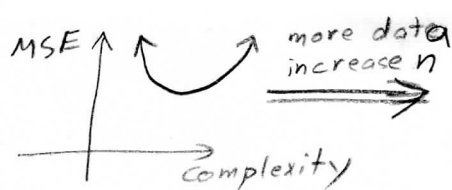
complexity

underfit    overfitting

Underfitting (complexity too low)

E(g) is far from f

⇒ Bias$^2$ term is high
Var term is low

Overfitting (complexity too high)

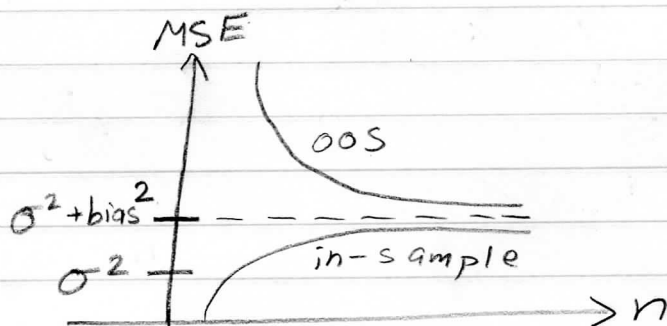E(g) is approximately f

⇒ Bias$^2$ term is low
Var term is high

MSE ↑ ⤸ ⤴ more data increase n → MSE ↑ ↑ ↗

complexity → complexity

curve is flatter

can make bias →0 by overfitting (more complexity)

but as complexity ↑, get variance ↑

Let's say we have a model that's
somewhat overfit now

look at number of data points n

As $n \to \infty$, get $g \to E[g] \Rightarrow Var[g] \to 0$

MSE ↑

OOS

$\sigma^2 + bias^2$ - - - - - - -

$\sigma^2$ in-sample

→ n

if
more complex →

OOS

$\sigma^2 + bias^2$ ↓ - - - - - -

$\sigma^2$ in-sample