

## Lecture 8

### Bias - Variance Decomp

Consider a model  $Y = f(\underline{X}) + \varepsilon$

where  $\varepsilon \perp Y, \underline{X}$

and  $\mathbb{E}\varepsilon = 0$

$\text{Var}\varepsilon = \sigma^2$

Fix  $\underline{X}_0$  at  $\underline{X}_0$  and look at  $\text{Err}(\underline{X}_0) = \mathbb{E}[(Y_0 - \hat{f}(\underline{X}_0))^2]$

*new pt not training* (pointing to  $\underline{X}_0$ )  
*sq. err.* (pointing to the squared term)

we can decompose this

$$= \mathbb{E}[(Y_0 - \underbrace{\mathbb{E}[\hat{f}(\underline{X}_0)]}_a + \underbrace{\hat{f}(\underline{X}_0) - \mathbb{E}[\hat{f}(\underline{X}_0)]}_b)^2]$$

$$\downarrow (a+b)^2 = a^2 + b^2 + 2ab$$

$$= \mathbb{E}[(Y_0 - \mathbb{E}[\hat{f}(\underline{X}_0)])^2] + \mathbb{E}[(\hat{f}(\underline{X}_0) - \mathbb{E}[\hat{f}(\underline{X}_0)])^2]$$

$$+ 2\mathbb{E}[(Y_0 - \mathbb{E}[\hat{f}(\underline{X}_0))](\hat{f}(\underline{X}_0) - \mathbb{E}[\hat{f}(\underline{X}_0)])]$$

①

$$Y_0 = f(\underline{X}_0) + \varepsilon_0$$

$$\textcircled{1} = \mathbb{E}[(f(\underline{X}_0) + \varepsilon_0 - \mathbb{E}[\hat{f}(\underline{X}_0)])^2]$$

$$= \cancel{\mathbb{E}[(f(\underline{X}_0) - \mathbb{E}[\hat{f}(\underline{X}_0)])^2]} + \mathbb{E}\varepsilon_0^2 + 2\mathbb{E}[\varepsilon_0(f(\underline{X}_0) - \mathbb{E}[\hat{f}(\underline{X}_0)])]$$

*constant* (pointing to  $f(\underline{X}_0)$ )

$$\underbrace{(f(\underline{x}_0) - E\hat{f}(\underline{x}_0))^2}_{\text{Bias}(\hat{f})} \quad \sigma^2 \quad + 2(f(\underline{x}_0) - E\hat{f}(\underline{x}_0)) \underbrace{E[\varepsilon]}_0$$

diff betwn truth ( $f$ ) and avg. est ( $E\hat{f}$ )

$$(1) = \text{Bias}(\hat{f})^2 + \sigma^2$$

$$(2) \quad E[(\hat{f}(\underline{x}_0) - E[\hat{f}(\underline{x}_0)])^2] = \text{Var}(\hat{f})$$

$$\text{Var}(z) = E[(z - E z)^2]$$

$$(3) = 0$$

$$2 E[(Y_0 - E\hat{f}(\underline{x}_0))(\hat{f}(\underline{x}_0) - E\hat{f}(\underline{x}_0))]$$

↑ indep      ← can distribute b/c indep

$$= E[Y_0 - E\hat{f}(\underline{x}_0)] \underbrace{E[\hat{f}(\underline{x}_0) - E\hat{f}(\underline{x}_0)]}_0$$

$$E[z - E z] = E z - E z = 0$$

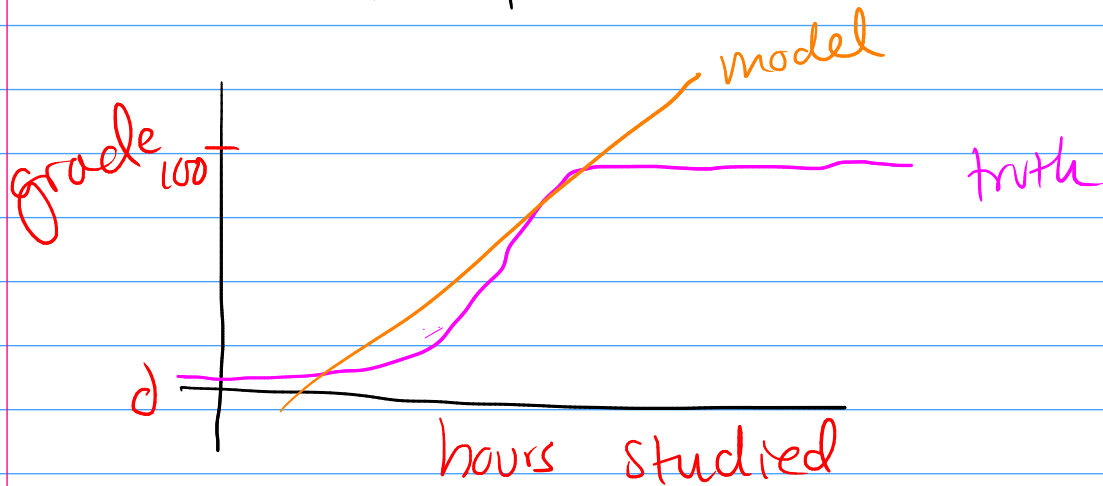
All together:

$$\text{Err}(\underline{x}_0) = \text{Bias}(\hat{f}(\underline{x}_0))^2 + \text{Var}(\hat{f}(\underline{x}_0)) + \sigma^2$$

↑  
irreducible  
var. term

# Laymen's terms

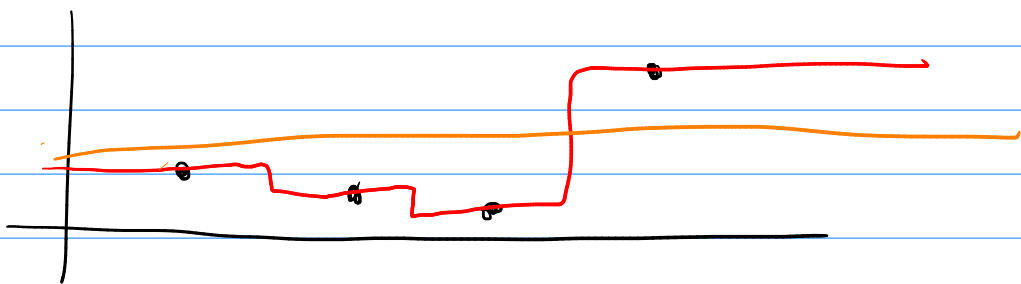
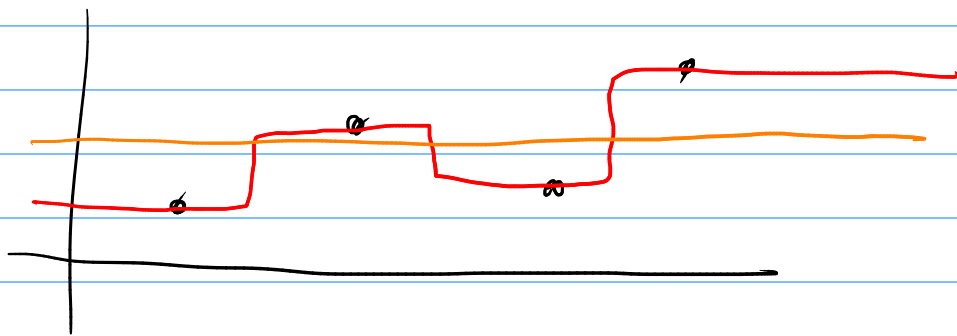
Bias: err. b/c we approx. a complicated real life fn w/ a simpler model



Variance of  $\hat{f}$ : how much my model changes when the training data changes

e.g. 1-NN

$N$ -NN = kNN w/  $k = \# \text{training}$

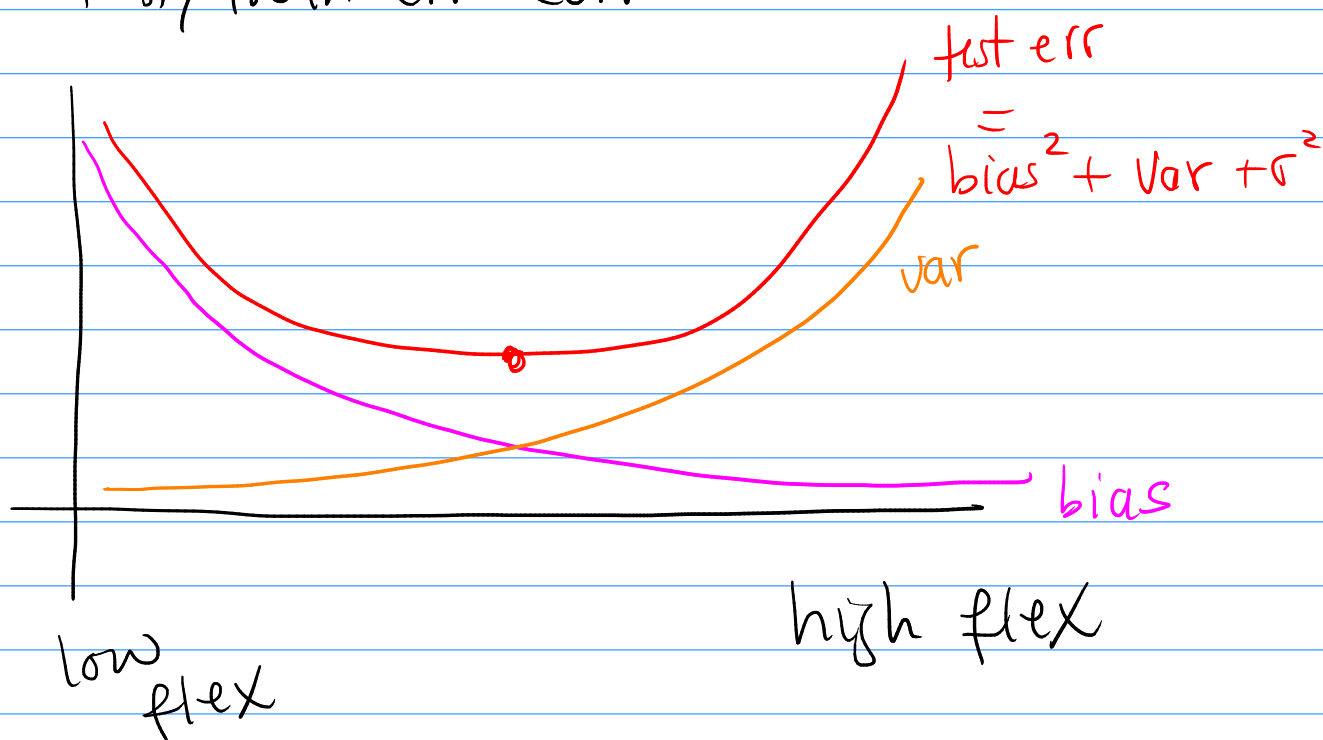


Typically:

low flex  $\Leftrightarrow$  low var, high bias

high flex  $\Leftrightarrow$  high var, low bias

Recall test/train err curves



Risk Minimization

Theoretically want

$$f^* = \underset{f}{\operatorname{argmin}} \mathbb{E}[L(Y, f(\underline{X}))]$$

Theoretical minimum?

Let  $(\underline{X}, Y) \sim p$   $\leftarrow$  joint density

$$\mathbb{E}[L(Y, f(X))] = \iint L(y, f(x)) p(x, y) dx dy$$

and  $p(x, y) = p(y|x) p(x)$

$$\downarrow = \iint L(y, f(x)) p(y|x) dy p(x) dx$$

iterated  
expectation  
→

$$\uparrow \underbrace{\mathbb{E}[L(Y, f(x)) | X=x]}_{\downarrow}$$

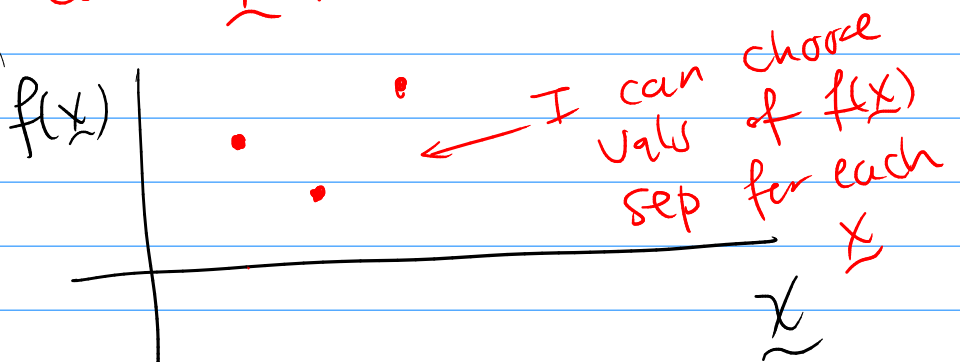
$$= \mathbb{E}_x \left[ \underbrace{\mathbb{E}[L(Y, f(x)) | X=x]}_{A(x)} \right]$$

Want to choose  $f$  to minimize.

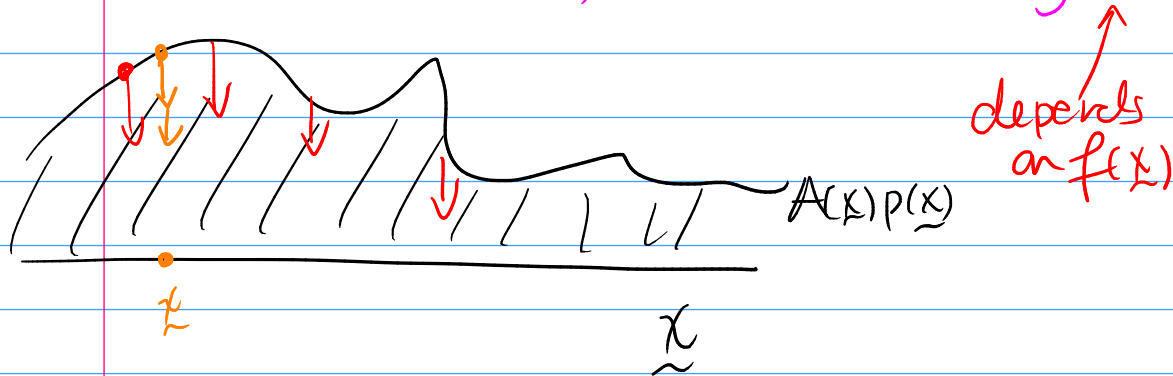
$$= \int A(x) p(x) dx$$

Consider no restrictions on  $f$ .

To choose  $f$  I need to tell you value of  $f(x)$  for each  $x$ .



Want to choose  $f(\underline{x})$  to min.  $\int A(\underline{x}) p(\underline{x}) d\underline{x}$



Game to do this is to sep. for each  $\underline{x}$  choose val. of  $f(\underline{x})$  to make  $A(\underline{x}) p(\underline{x})$  as small as possible.

Since  $f$  doesn't affect  $p$  we can just choose  $f$  to make  $A(\underline{x})$  as small as possible.

Punchline: to min risk I can choose  $f(\underline{x})$  sep. for each  $\underline{x}$  to min

$$A(\underline{x}) = \int L(y, f(\underline{x})) p(y|\underline{x}) dy$$

consider:  $L(y, f(\underline{x})) = (y - f(\underline{x}))^2$

$$f^*(\underline{x}) = \operatorname{argmin}_{f(\underline{x})} E[(Y - f(\underline{x}))^2 | \underline{X} = \underline{x}]$$

$$= \operatorname{argmin}_c E[(Y - c)^2 | \underline{X} = \underline{x}]$$

$$\arg \min_c \mathbb{E}[(z-c)^2] = \mathbb{E}z$$

why?  $\mathbb{E}[z^2 + c^2 - 2zc] = \mathbb{E}[z^2] + c^2 - 2c\mathbb{E}z$

$$\frac{\partial}{\partial c} [\dots] = 2c - 2\mathbb{E}z = 0$$

$$\Rightarrow c = \mathbb{E}z$$

So ...

$$f^*(\underline{x}) = \mathbb{E}[Y | \underline{X} = \underline{x}]$$



What about other losses?

$$L(y, f(x)) = |y - f(x)|$$

$$f^*(\underline{x}) = \text{median}(Y | \underline{X} = \underline{x})$$

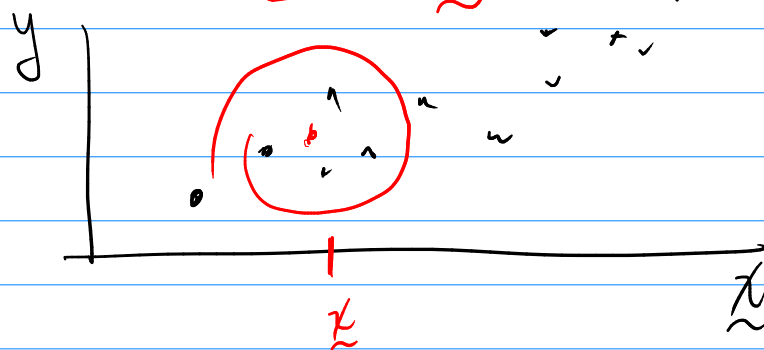
Problem: I don't know  $p(\underline{x}, y)$ .

All I have is training data.

Let's approximate  $E[Y | \underline{X} = \underline{x}]$ .

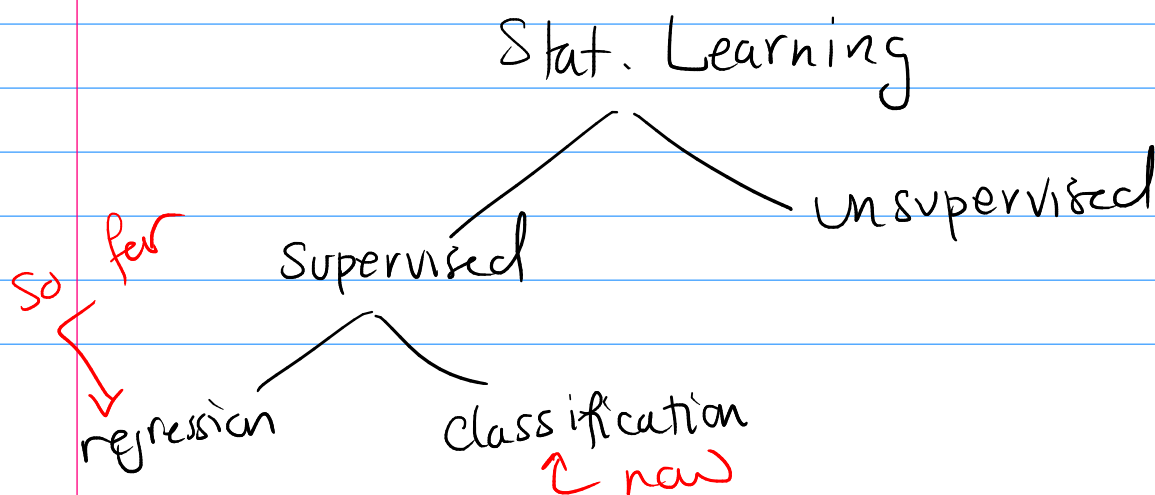
Ex. (1)  $\hat{f}(\underline{x}) = \text{avg. of } y_n \text{ s for } \underline{x}_n$   
near  $\underline{x}$

KNN  
regression



(2) Maybe assume  $E[Y | \underline{X} = \underline{x}] = \underline{x}^T \beta$   
and then build  $\hat{f}$  by est.  $\beta$ .

OLS  
regression





Setup: Classification

$$\underline{x} \in \mathbb{R}^P \quad \text{and} \quad y \in \mathcal{C}$$

$$\mathcal{C} = \{c_1, c_2, \dots, c_k\}$$

↑ set of possible classes.

Goal: find some  $\hat{f}$  so that  $\hat{f}(\underline{x}) \approx y$

---

Need a loss fn 0-1 loss

$$\begin{aligned} L(y, \hat{f}(\underline{x})) &= \mathbb{1}(y \neq \hat{f}(\underline{x})) \\ &= \begin{cases} 0 & \hat{f}(\underline{x}) = y \\ 1 & \hat{f}(\underline{x}) \neq y \end{cases} \end{aligned}$$

What is  $f^*$ ?

$$f^*(\underline{x}) = \arg \min_c \mathbb{E}[L(y, c) | \underline{X} = \underline{x}]$$

$$\mathbb{E}[\mathbb{1}(z \in A)] = P(z \in A)$$

$$= \arg \min_c P(Y \neq c | \underline{X} = \underline{x})$$

$$= \arg \min_c ( -P(Y=c | \underline{X} = \underline{x}) )$$

$$\hat{f}^*(\underline{x}) = \arg \max_c P(Y=c | \underline{X} = \underline{x})$$

← Bayes' classifier

↑ predict class  $c$  w/ highest cond. prob.

e.g.  $K=3$

$$P(Y=1 | \underline{X} = \underline{x}), P(Y=2 | \underline{X} = \underline{x}), P(Y=3 | \underline{X} = \underline{x})$$

↑ pick the largest

---