

## Lecture 7: Model Evaluation

Need to be precise about what is random.

Consider training data  $(X_n, Y_n) \stackrel{iid}{\sim} p$

So  $T = \{(X_n, Y_n)\}_{n=1}^N$

$\nwarrow$  is random

Use this to fit

$$\hat{f} = \hat{f}_T \leftarrow \text{random}$$

Let  $X_0, Y_0$  be an independent (from training data) sample from  $p$ .

and let  $M$  be my perf. metric

$$\text{e.g. } M(A, B) = (A - B)^2$$

then let

$$\text{Err}_T = \mathbb{E}[M(Y_0, \hat{f}_T(X_0)) | T]$$

$\nearrow$   
probably what  
we want to  
estimate

= given some  $T$  expected  
err. on new data

However, in practice often difficult to estimate.

Instead, easier to estimate

$$\begin{aligned}\text{Err} &= \mathbb{E}_T[\text{Err}_T] \\ &= \mathbb{E}[M(Y_0, \hat{f}(x_0))]\end{aligned}$$

= expected gen. perf. across  
all possible  $T$

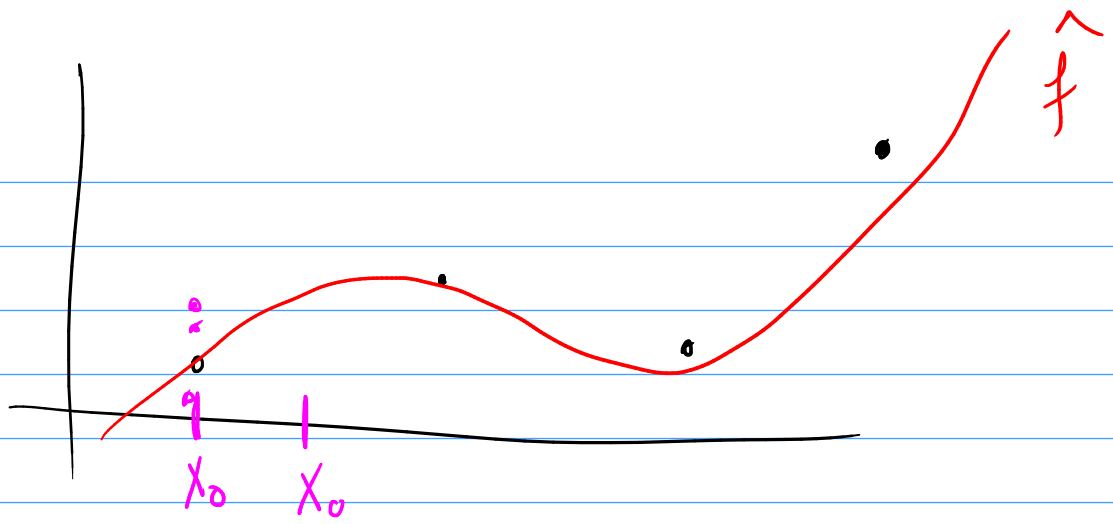
Also define

$$\begin{aligned}\overline{\text{err}} &= \text{training error} \\ &= \frac{1}{N} \sum_n M(y_n, \hat{f}(x_n))\end{aligned}$$

Typically  $\overline{\text{err}} < \text{Err}_T$  b/c we over-fit  
both eval. for fixed  $T$

Part of the problem:

- (1)  $x_0$  is different than training  $x_n$
- (2)  $y_0$  is different than training  $y_n$   
(even if  $x_0 = x_n$ )



Partially simplify by considering only  $Y_0$  to be random.

In-sample error: 
$$\text{Err}_{\text{in}} = \frac{1}{N} \sum_n \mathbb{E}_{Y_0} [M(Y_0, \hat{f}(X_n))]$$

fix  $X_0$ 's at training  $\underline{X}_n$ 's  
and calc. gen. err.

- like  $\text{Err}_T$  but fixing  $X_0$ 's at training.

Define optimism

$$op = \text{Err}_{\text{in}} - \overline{\text{err}}$$

in-sample - training.

typically  $op > 0$  as  $\overline{\text{err}}$  underestimates  $\text{Err}_{\text{in}}$

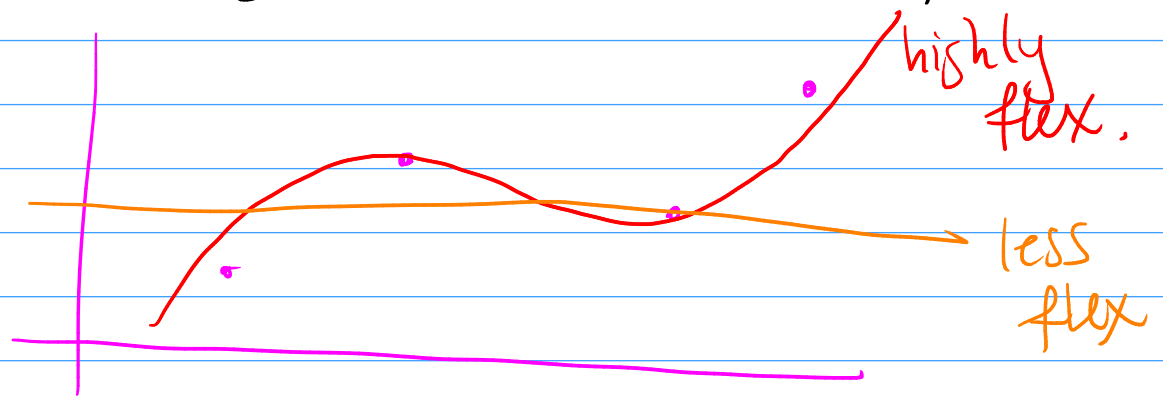
can calculate expected  $op$  over potential training data ( $\underline{X}_n$ 's)

$$\omega = \mathbb{E}[op]$$

In general

$$\omega = \frac{2}{N} \sum_n \text{cov}(\hat{Y}_n, Y_n)$$

So  $\omega$  the avg. amount that  $Y_n$  affects  $\hat{Y}_n$



Notice

$$E[\text{Err}_n] = E[\overline{\text{err}}] + \omega$$

in some cases I can estimate  $\omega$  as  $\hat{\omega}$  and then est. my  $\text{Err}_n$

$$\text{Err}_n \approx \overline{\text{err}} + \hat{\omega}$$

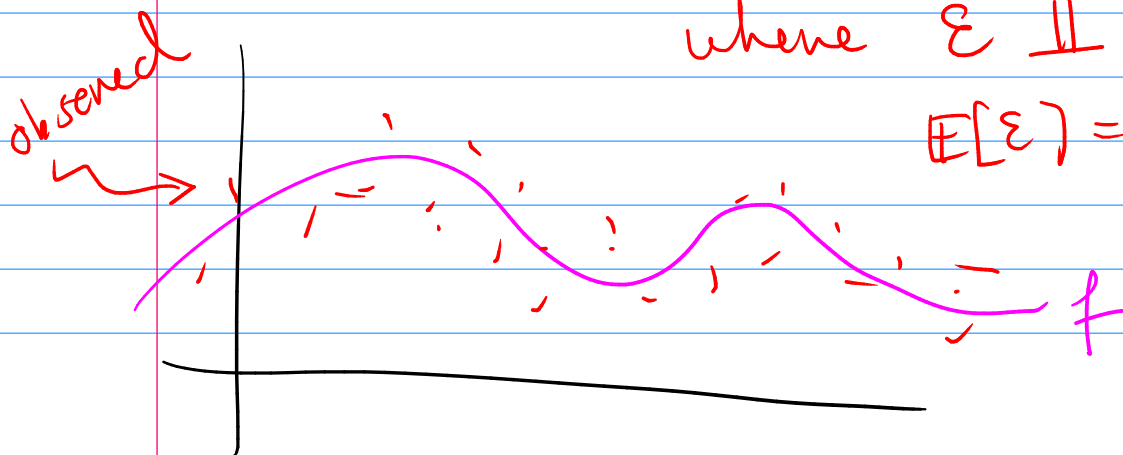
Ex.

Assume

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

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

$$E[\varepsilon] = 0, \text{Var}(\varepsilon) = \sigma^2$$



Then if I fit a lin. regression w/  
P variables, then

$$\omega = \frac{2P}{N} \sigma^2.$$

In practice what we can do is estimate  $\sigma^2$  as  $\hat{\sigma}^2$  and then calc.

$$\hat{\text{Err}}_{\text{in}} = \overline{\text{err}} + \frac{2P}{N} \hat{\sigma}^2.$$

↑ mallows's  $C_p$   
generally AIC.



---

### Tuning model complexity

① CV estimates Err

② estimate  $\text{Err}_{\text{in}}$

) neither estimating  
 $\text{Err}_T$

Downside of CV is that it trains on a slightly smaller data set.

Do a  $K$ -fold  $x$ -val.

Can choose  $K$ . Can set  $K=N$ .

Called Leave-one-out  $x$ -val (LOO).

---

## Bias-Variance Tradeoff.

Consider a model

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

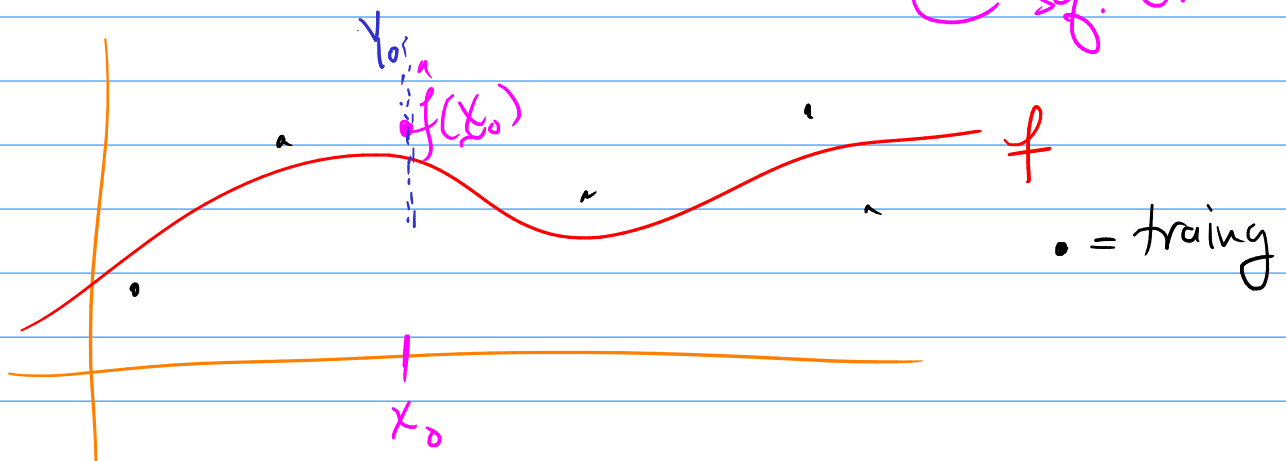
$$\varepsilon \perp \underline{X}, Y$$

$$\mathbb{E}\varepsilon = 0, \text{Var } \varepsilon = \sigma^2 < \infty.$$

Fix  $\underline{X}_0$  at some fixed  $\underline{x}_0 \in \mathbb{R}$  and calc

$$\text{Err}(\underline{x}_0) = \mathbb{E}[(Y_0 - \hat{f}(\underline{x}_0))^2] = (*)$$

↑ sq. err.



Can decompose:

expected over  $T$

$$(*) = \mathbb{E} \left[ \underbrace{(Y_0 - \mathbb{E}[\hat{f}(x_0)])}_a + \underbrace{(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))}_b \right]^2$$

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

$$= \underbrace{\mathbb{E}[(Y_0 - \mathbb{E}[\hat{f}(x_0)])^2]}_{a^2} + \underbrace{\mathbb{E}[(\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)])^2]}_{b^2}$$

$$+ 2 \mathbb{E}[(Y_0 - \mathbb{E}[\hat{f}(x_0)])(\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)])]$$

+ 2ab

Claim:  $\mathbb{E}[a^2] = \text{Bias}(\hat{f})^2 + \sigma^2$

$$\mathbb{E}[b^2] = \text{Var}(\hat{f})$$

$$\mathbb{E}[2ab] = 0$$

i.e.  $\text{Err}(x_0) = \text{Bias}^2 + \text{Var} + \sigma^2$

irreducible noise

Bias =  $\mathbb{E}[\hat{f}(x_0) - f(x_0)]$  = expected diff  
btwn est. and truth

Var =  $\text{Var}(\hat{f}(x_0))$  = amt. my  $\hat{f}$  changes due to  
changing  $T$

Typically:

low flex  $\Leftrightarrow$  low var, high bias

high flex  $\Leftrightarrow$  high var, low bias

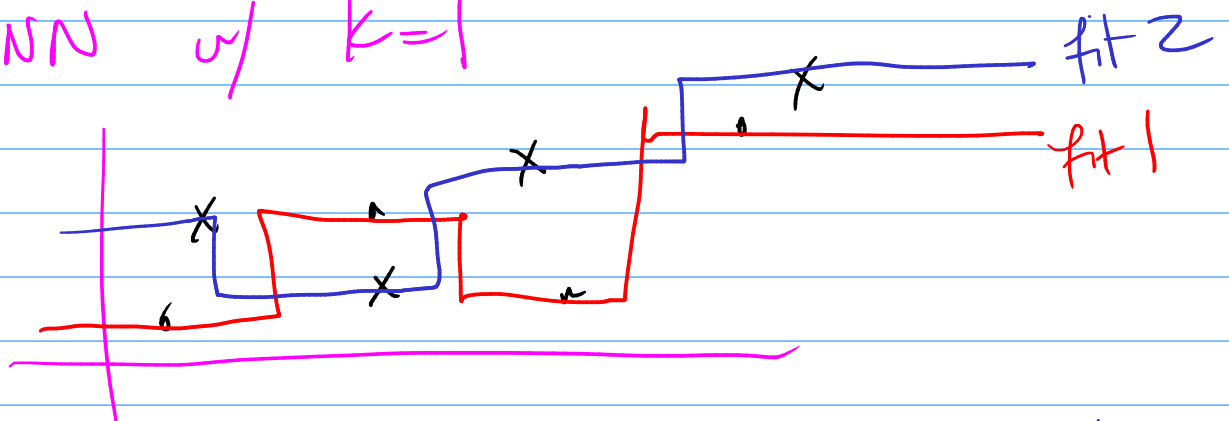
In laymens terms

bias = err incurred b/c approx. complicated  
real life w/ simple model



variance = sensitivity of fit

e.g. KNN w/  $k=1$



KNN w/  $k=N$

