# Lecture 7: Bias, Var, Model Complexity

To precisely discuss evaluation need to keep track of what is random and what is fixed.

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

Let $T = \{(X_n, Y_n)\}_{n=1}^{N}$ be training data

↳ this is random

We use this training data to build a model

$$\hat{f} = \hat{f}_T$$

← also random

Let $X_0, Y_0$ be indep (from training) sample from $p$

If $M$ is some perf. metric (like $L$)

then let

$T$ is fixed

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

what we'd really like to estimate →

how close $\hat{f}(X_0)$ is to $Y_0$

= given specific $T$ what is expected err. if we predict using $\hat{f}$ on independent sample $X_0, Y_0$

= gen. perf. using this training data $T$

This is quite difficult in practice to estimate.

Instead, typically easier to estimate

$$Err = E_T[Err_T]$$

$$= E[M(Y_0, \hat{f}(X_0))]$$

= exp. gen. perf. over new samples

$\underline{AND}$
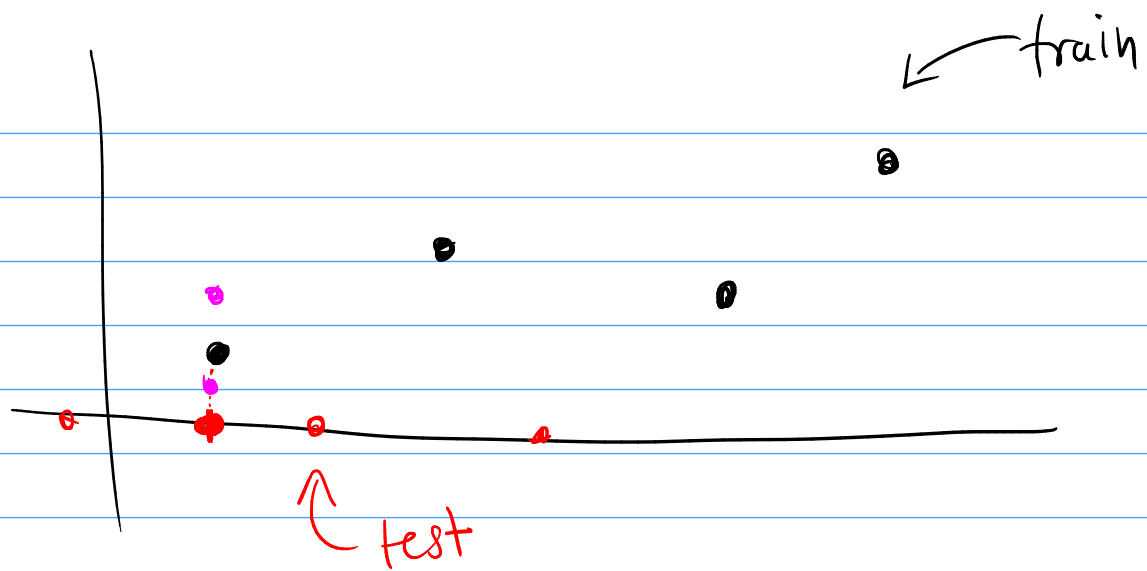
over all possible training data

= exp. perf. of my model building process

Let $\overline{err}$ = training error = $\frac{1}{N} \sum_n M(y_n, \hat{f}(X_n))$

Typically $\overline{err} < Err_T$ b/c we overfit

$\uparrow$ both calc. w/ fixed training

Parts of the problem

extra

truly

① $\underset{\sim}{X_0}$ is random but $\underset{\sim}{X_n}$s are fixed

So new $\underset{\sim}{X_0}$ may not be exactly the same as $\underset{\sim}{X_n}$s.

② $Y^0$ is random so might not match $Y_n$s.

$\leftarrow$ train

$\uparrow$ test

To simplify analysis, consider only $Y^o$ being random and define the in-sample error as

$$Err_{in} = \frac{1}{N} \sum_n \mathbb{E}[M(Y_o, \hat{f}(\underset{\sim}{x_n})) | T]$$

$Y^o$ — $\uparrow$ fixed at train pts

$\uparrow$ like $Err_T$ but fixing $\underset{\sim}{x_n}$ at train pts

let the optimism be

$\leftarrow$ random (deps on $Y^o$)

$$Op = Err_{in} - \overline{err}$$
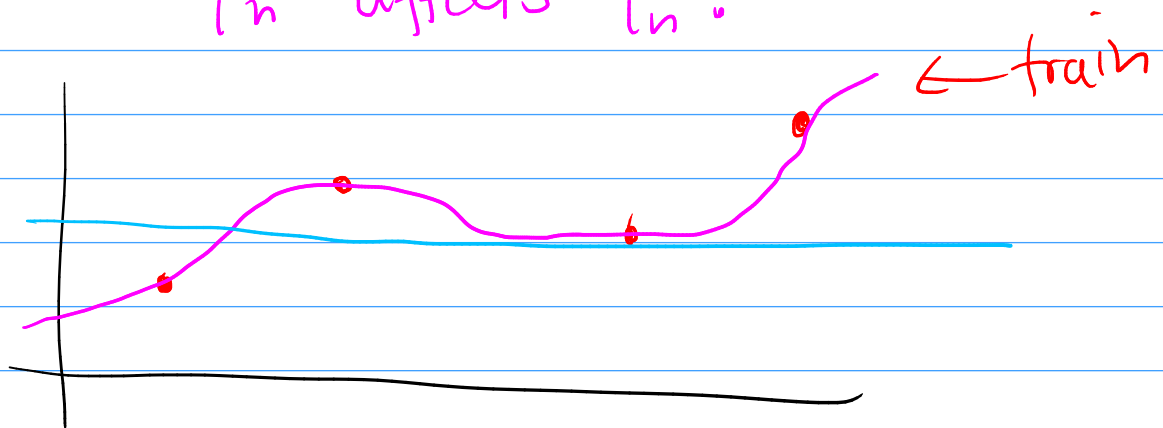
in sample — training

Typically $op > 0$ as $\overline{err}$ underestimates $Err_{in}$ and if

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

= avg optimism

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

So $\omega$ is rel. to the avg. amt that $Y_n$ affects $\hat{Y}_n$.


← train

So

$$\mathbb{E}[\text{Err}_{in}] = \mathbb{E}[\overline{\text{err}}] + \omega$$

In some cases one can directly estimate $\omega$ as $\hat{\omega}$ in which care

$$\text{est. of Err}_{in} = \overline{\text{err}} + \hat{\omega}$$

↑ know     ↑ estimate

e.g. If we assume $Y = f(\underset{\sim}{X}) + \varepsilon$

where $\varepsilon \perp\!\!\!\perp \underset{\sim}{X}, Y$ and

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

# covs.

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

and so $\hat{\omega} = \frac{2P}{N} \hat{\sigma}^2$

est. of in-sample err $= \overline{err} + \frac{2P}{N}\hat{\sigma}^2$

in-sample est. penalty



err

$\overline{err}$

P

Many ways to do this dep. on assumptions,
$C_p$, AIC, BIC, adj. $R^2$

Hold on, this only estimates in-sample err —
doesn't take into account randomness in $\underset{\sim}{X}_0$.

Yes, but often its good enough for model selection,
b/c it still allows me to eval. rel. perf.
among models.

If want to est. gen. perf. still need to
do some kind of x-validation.
$\hookrightarrow$ estimating Err
not $Err_T$