

Lec 16 Math 390.4

11

Polynomial regression has a type of "non-linear" linear modeling.

Here is another "non-linear" linear modeling: interactions.

lap?
Approx?
Yes?

Consider the data with $p=2$. $X = \begin{pmatrix} 1 & x_1 & x_2 \\ \vdots & \vdots & \vdots \end{pmatrix}$

hypothesis set

How to fit?

$$\mathcal{H} = \{ w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 : \vec{w} \in \mathbb{R}^4 \}$$

$$\Rightarrow X = \begin{pmatrix} 1 & x_{11} & x_{21} & x_{11}x_{21} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & x_{1n}x_{2n} \end{pmatrix}$$

def. = 4

Why would you do this? Well, remember -

$$\mathcal{H} = \{ w_0 + w_1 x_1 + w_2 x_2 \} \xrightarrow{A} g(x) = b_0 + b_1 x_1 + b_2 x_2$$

Interp: if x_1 increases by one unit, \hat{y} increases by b_1 units on average regardless of the value of x_2

$$\mathcal{H} = \{ w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 \} \xrightarrow{A} g(x) = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_1 x_2 \\ = b_0 + x_1 (b_1 + b_3 x_2) + b_2 x_2$$

Interp: if x_1 increases by one unit, \hat{y} increases by $b_1 + b_3 x_2$ units on average

Now \Rightarrow You're given the model d.o.f. to have different slopes based on the value of the other predictors \Rightarrow much more expressive!

If x_2 e.g. is binary $x_2 = \text{Is male}$
then slope of x_1 is b_1 if female & $b_1 + b_3$ if male.

CMAD

Nar problem...

given the same X , there are many H 's (and A 's) that produce different g 's.

e.g.

$$g_1 = b_0 + b_1 x_1$$

$$g_2 = b_0 + b_1 x_1 + b_2 x_2$$

$$g_3 = b_0 + b_1 \log(x_1) + b_2 x_2$$

$$g_4 = b_0 + b_1 x_1^2 + b_2 x_2$$

$$g_5 = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_1 x_2$$

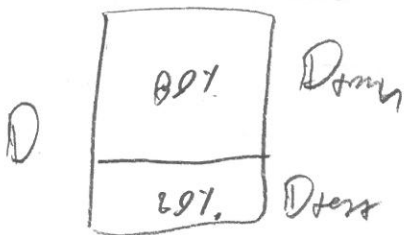
\vdots

infinite models!!

How to "choose" which one to employ?

Model selection - one of the most fundamental questions in all of statistics (maybe even science - is it the "best" model?)

Ideas? Planning



Purpose? Use D_{train} to estimate error (a conservative honest estimate of future performance)

Why? $D_{train} \subset D$
if we only $D \Rightarrow$ ^{very} small

Why not test g_1, \dots, g_m by fitting them on D_{train} , testing on D_{test} and picking the one with the least error?

You can! But there is a problem!

What do we use now to estimate future prediction?

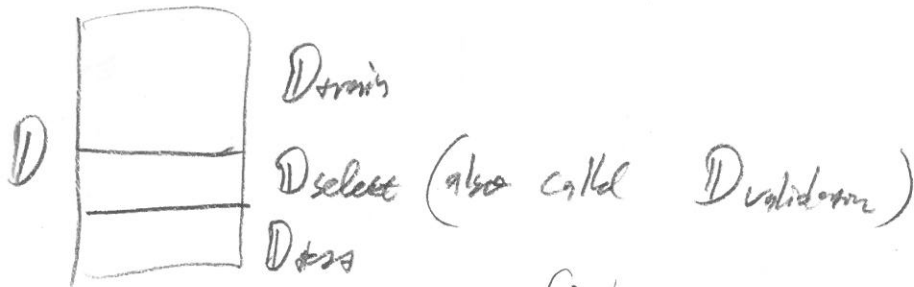
Why can't we use D_{test} again?

Training

g_1, \dots, g_M where M is large.

I can find g_j s.t. error on D_{test} is near 0 by coincidence! This is not a reliable estimate of future performance!

Edison?



Randomly create
3 splits!

$K=5$ typical

60% train
20% select
20% test

Protocol: \rightarrow defining \mathcal{A} and \mathcal{A} .

for each model $j \in \{1, \dots, M\}$

build it $g_j = \mathcal{A}_j(\mathcal{A}_j, D_{train})$

get error $oos e_j = \text{error}(y_{test}, g_j(x_{test}))$ (usually se)

Select $j^* := \argmin \{oos e_1, \dots, oos e_M\}$

Compute $oos e_{j^*} = \text{error}(y_{test}, g_{j^*}(x_{test}))$

Steps 1 & 2
are like a
given $\mathcal{A}_{combined}$

test results
 $\mathcal{A}_{combined}$

Ship model built with steps 1-4 on full dataset!

This is called model selection by oos error.

This is one such method. Other methods are "analytical"
i.e. they rely on statistical/prob. models - we didn't talk about
those! AIC, BIC, Cp, etc....

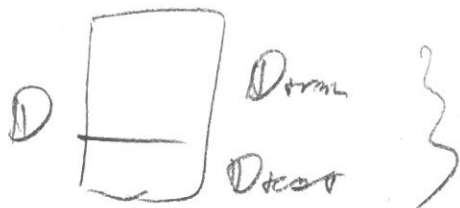
probably won't get to
it this semester
unfortunately.

DEMO

↗ Midterm 2

↘ FINAL

There is an extra step in validation that is usually employed.
Previously before validation we were looking at one model



Random 80% / 20% split of n observations

What if you get an
"outlucky" split?

e.g. all the "noisy" observations in Train
⇒ your ^{error} estimate is much higher
than it should be

all the "noisy" observations in Test
⇒ your error is much lower
" " " " " " " " " " " "

⇒ Variance
in the estimate
of future
performance