

LINEAR REGRESSION REMINDERS

MODEL  $f(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x})$

Have code to minimize  $\sum_n (y^{(n)} - \underline{w}^T \underline{\phi}(x^{(n)}))^2$  wrt  $\underline{w}$   
 $= (\underline{y} - \Phi \underline{w})^T (\underline{y} - \Phi \underline{w})$

$$\underline{\phi}(\underline{x}) = [\phi_1(\underline{x}) \ \phi_2(\underline{x}) \ \dots \ \phi_k(\underline{x})]^T$$

$\phi_k(\underline{x})$  any scalar function:

- Monomial, eg  $x_2, x_3 x_4^3, \dots, 1, \dots$

- Radial Basis Function

"  $\underline{w}^T$   " = 

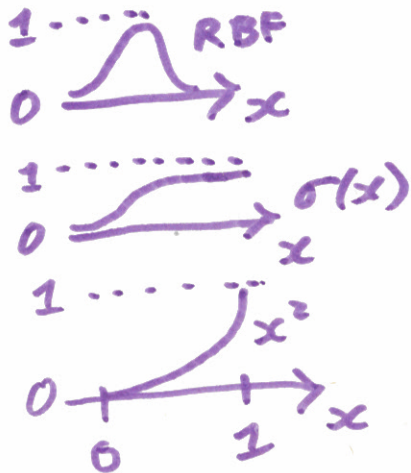
- Logistic sigmoid



We noticed:

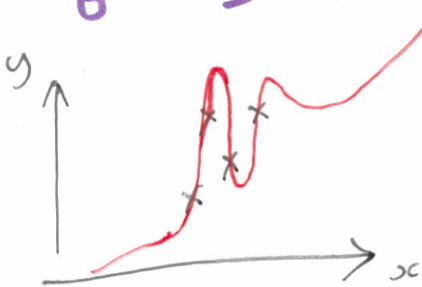
- a) Extrapolate differently; b) Placement matters.  
 Might need large  $K$ .

# Why are large weights bad? MLPR 2016 L4 ©



If basis  $f_i$ 's  
are bounded

Large  $f \Rightarrow$  large  $w$

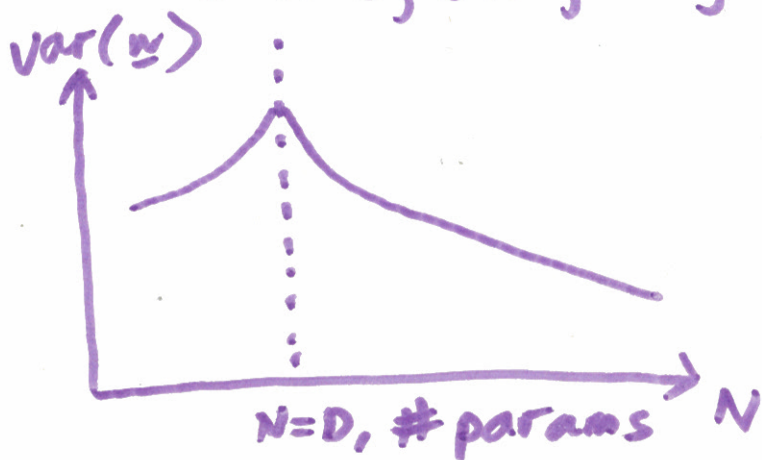


Large derivatives

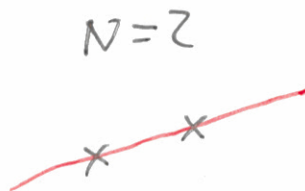
$\Rightarrow$  we've taken  
extreme differences  
between  $\phi$ 's

$\Rightarrow$  large  $w^T w$

Is amount of overfitting =  $\underline{w}^T \underline{w}$  ?



$$\text{var}: \frac{1}{K} \sum w_k^2 - \left[ \frac{1}{K} \sum w_k \right]^2$$



Pseudo-Data Trick

$$y' = \begin{bmatrix} y \\ \underline{0}_k \end{bmatrix} \left. \vphantom{\begin{bmatrix} y \\ \underline{0}_k \end{bmatrix}} \right\} \begin{array}{l} k \text{ fake} \\ \text{observations} \end{array}$$

$$y_{N+1} = 0$$

$$\vdots$$

$$y_{N+k} = 0$$

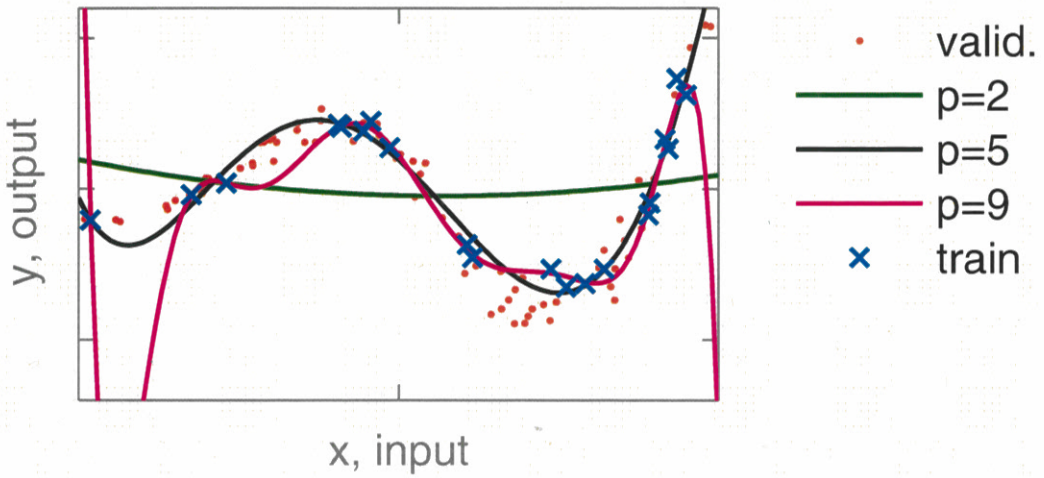
$$\Phi' = \begin{bmatrix} \Phi \\ \sqrt{\lambda} \begin{matrix} 0 & 0 & 0 & 0 \\ 0 & \sqrt{\lambda} & 0 & 0 \\ 0 & 0 & \sqrt{\lambda} & 0 \\ 0 & 0 & 0 & \sqrt{\lambda} \end{matrix} \end{bmatrix}$$

$$\sqrt{\lambda} I$$

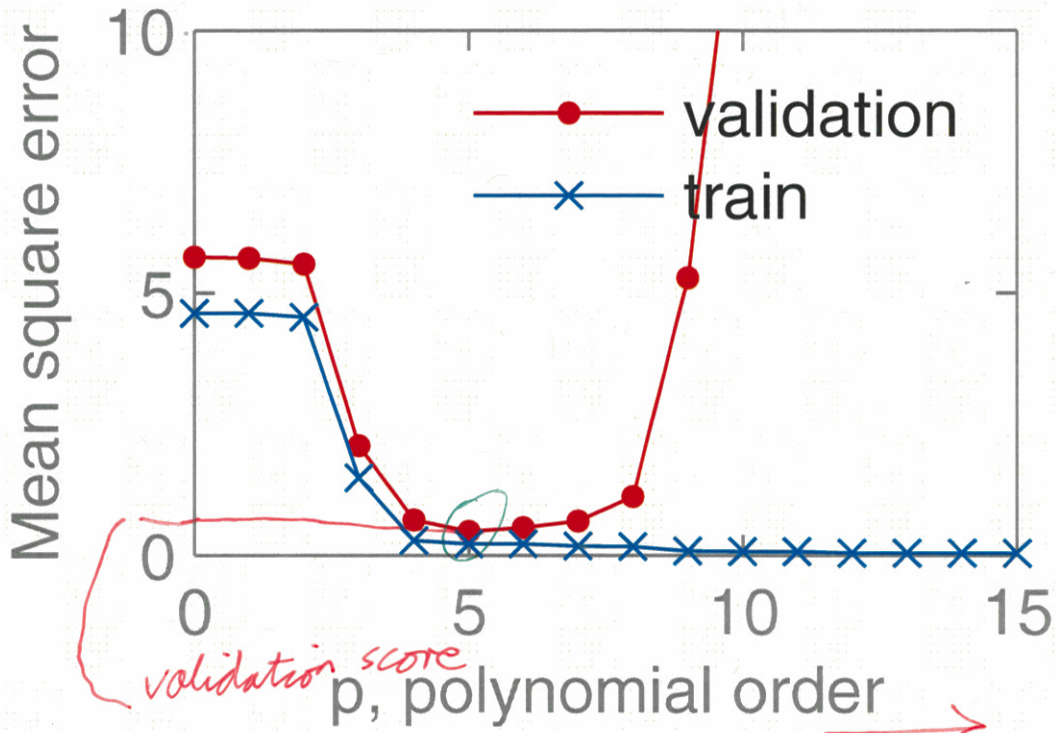
$$\begin{aligned} E_{\lambda}(w) &= (y' - \Phi' w)^T (y' - \Phi' w) \\ &= (y - \Phi w)^T (y - \Phi w) + \lambda w^T w \end{aligned}$$

$(\lambda > 0)$

Can't use  $E_{\lambda}$  to set  $\lambda$



should be a log scale



Polynomials are "nested"

or  $-\log \lambda$

TRAIN ERROR MONOTONIC BECAUSE MODELS ARE NESTED

# Generalisation

$$\mathbb{E}_{p(x,y)} [L(y, f(x))] = E_{\text{gen}}$$

↑  
Loss, e.g.  $(y - f(x))^2$

We assume there is some fixed distribution on future inputs + outputs

$$E_{\text{gen}} = \int L(y, f(x)) p(x, y) dx dy$$

$$\approx \frac{1}{M} \sum_{m=1}^M L(y^{(m)}, f(x^{(m)})) = E_{\text{test}}$$

$$x^{(m)}, y^{(m)} \sim p(x, y)$$

Drawn Test set  
from a

Unbiased:

$$\begin{aligned} \mathbb{E}[E_{\text{test}}] &= \frac{1}{M} \sum_{m=1}^M \underbrace{\mathbb{E}[L_m]}_{E_{\text{gen}}} \\ &= \frac{1}{M} M E_{\text{gen}} = E_{\text{gen}} \end{aligned}$$

But not if model was selected  $E_{\text{test}}$  small!

## Data splits

Training set: fit  $w$ , parameters like those

Validation set: used to pick models  
eg. polynomial order  
or regularizers  $\lambda$

Test set: to report how well we  
do in future.