

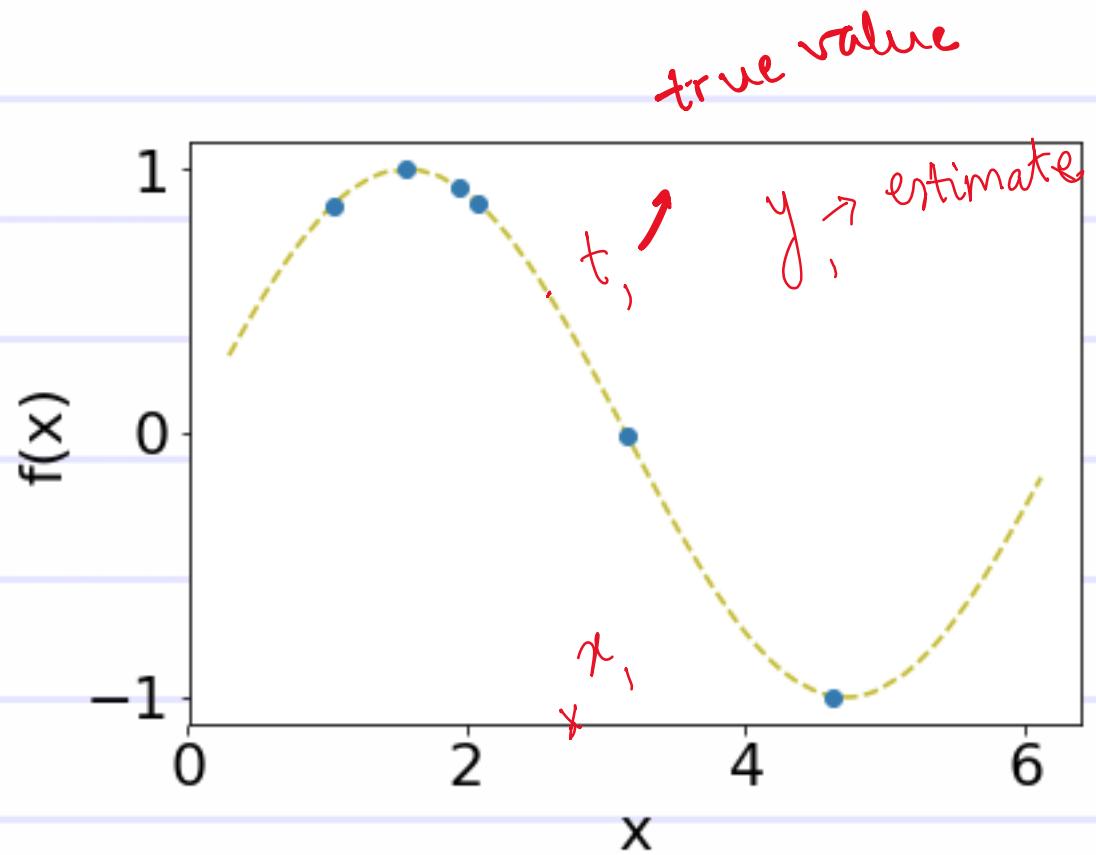
# Linear Regression

---

*EE698V - Machine Learning for Signal Processing*

Vipul Arora





$x \rightarrow F \rightarrow y$

training data

$x_1, y_1, t_1$

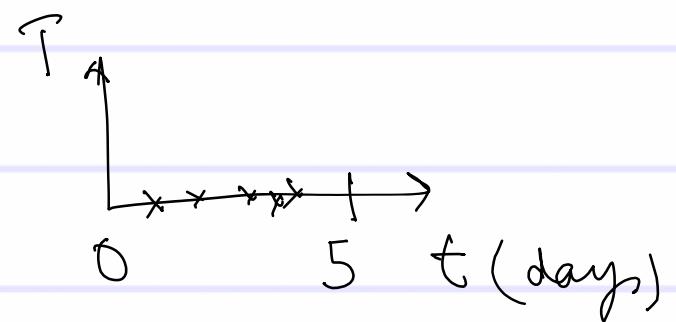
$x_2, y_2, t_2$

⋮

$x_5, y_5, t_5$

Given the blue points, how to estimate  $f(x)$  for any  $x$

?



Actual value = t

Let us take a linear model

*estimate*  $y = w_0 + w_1x + w_2x^2 + \dots + w_Dx^D$

$$x \rightarrow [f] \rightarrow y = \begin{bmatrix} 1 & x & x^2 & \dots & x^D \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

$D+1$  unknowns  $\rightarrow w_0, \dots, w_D$

5 eq's

$$\begin{array}{ll} x_1 & y_1 \\ \vdots & \vdots \\ x_5 & y_5 \\ \hline y_t = & \sum w_i x_i \end{array} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^D \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_5 & x_5^2 & \dots & x_5^D \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

$$y = \Phi w$$

if  $D+1 = S$ , then

$$w = \Phi^{-1} y \quad \text{unique.}$$

if  $D+1 < S$ , then

$$\text{minimize } \|t - y\|$$

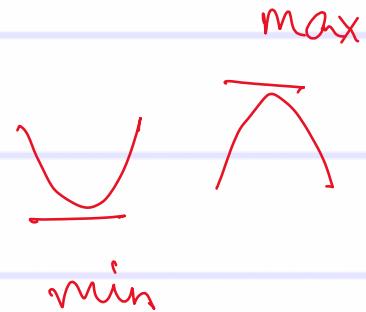
$$y_1 = \Phi w$$

$$y_1 = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

least squares solution

$$\min E = (\mathbf{t} - \Phi w)^T (\mathbf{t} - \Phi w) = \|\mathbf{t} - y\|^2$$

$$\frac{\partial E}{\partial w} = \left( \frac{\partial}{\partial w} (\mathbf{t} - \Phi w) \right) (\mathbf{t} - \Phi w) \times 2$$



$$0 = -\Phi^T (\mathbf{t} - \Phi w) \times 2$$

$$0 = -\Phi^T \mathbf{t} + \Phi^T \Phi w$$

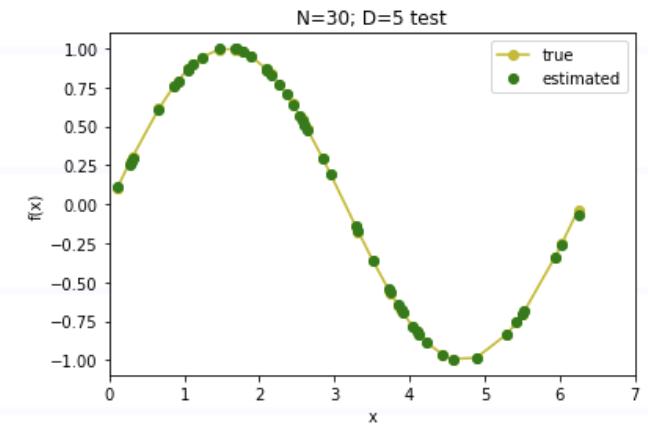
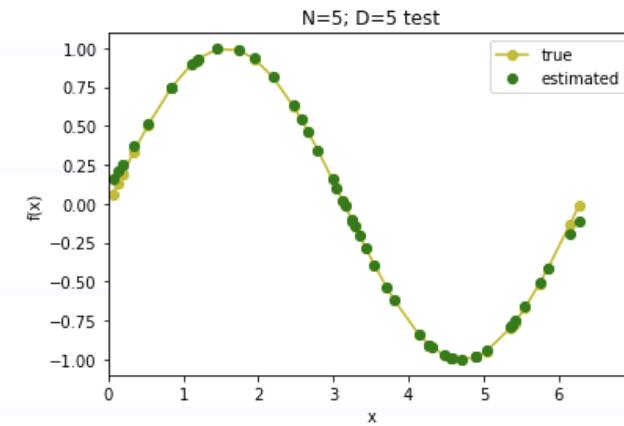
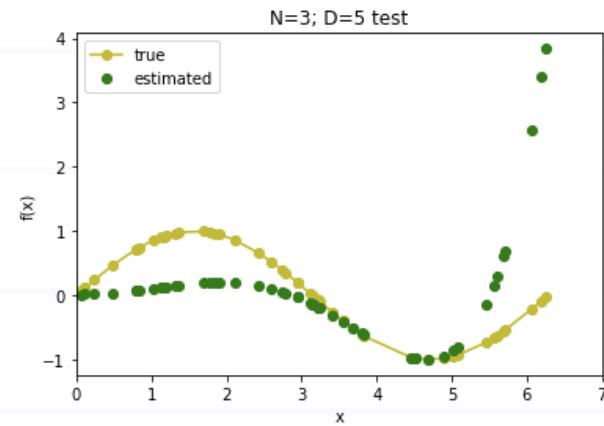
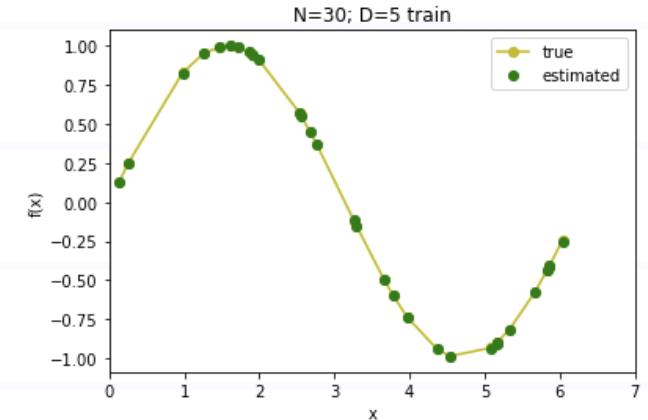
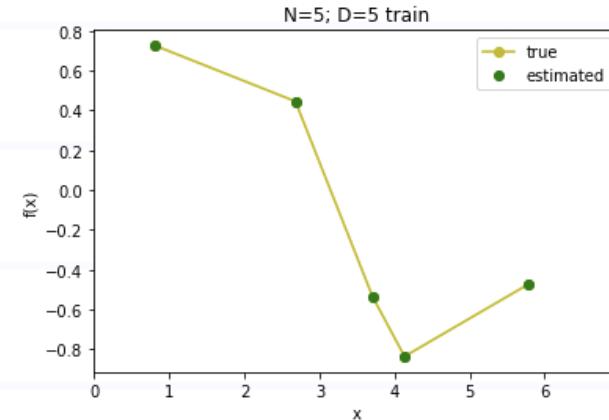
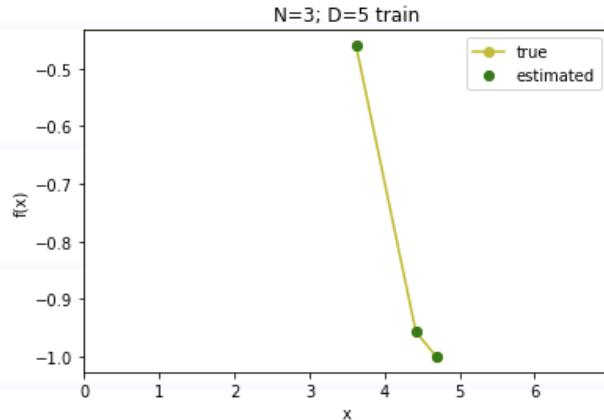
$$w = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\frac{\partial (u^T v)}{\partial w} = \frac{\partial u}{\partial w} v$$

$$+ \frac{\partial v}{\partial w} u$$

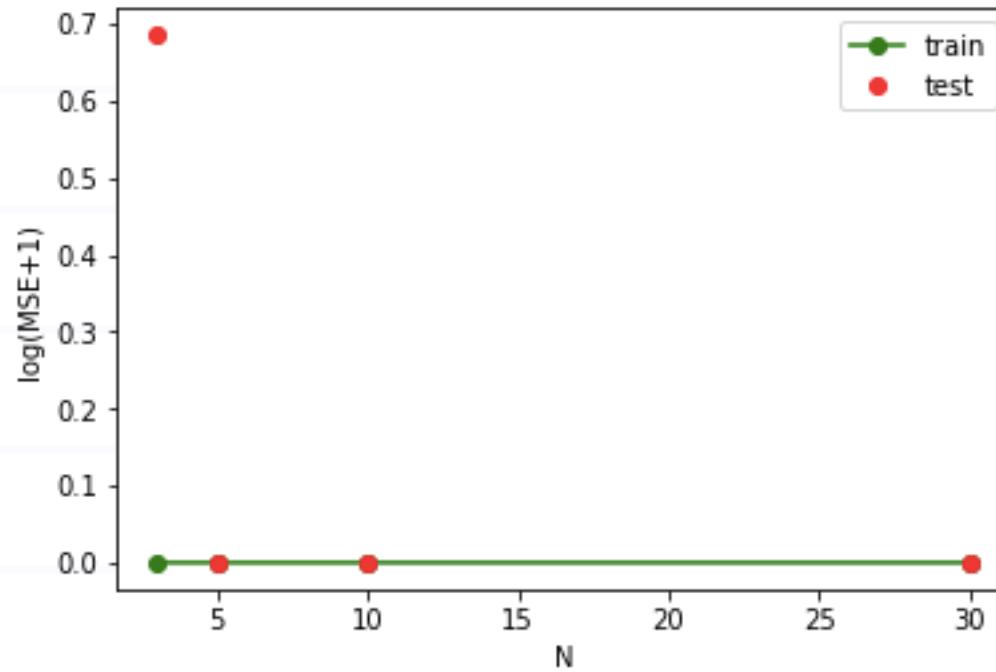
$$\frac{\partial \Phi u}{\partial w} = \frac{\partial u}{\partial w} \Phi^T$$

# Effect of increasing data points (fixed D)

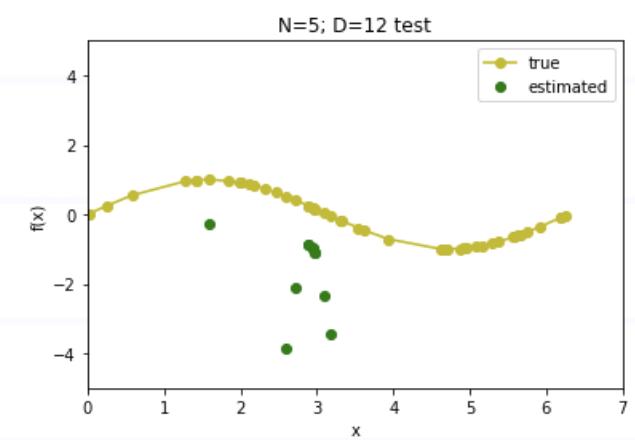
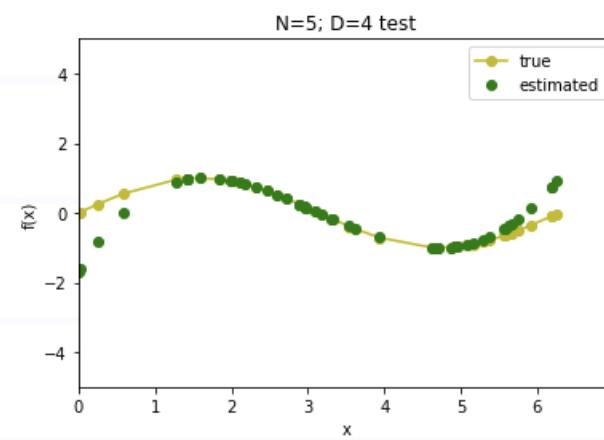
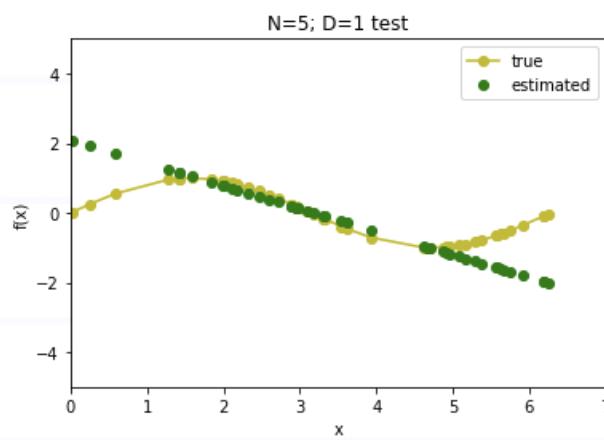
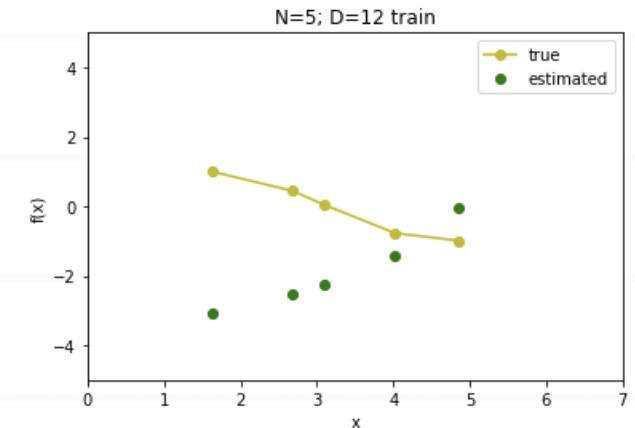
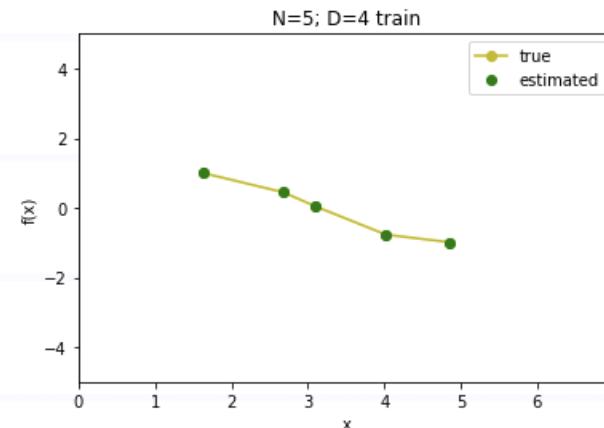
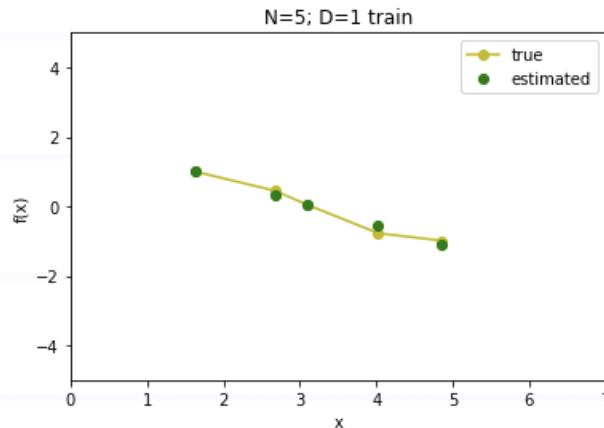


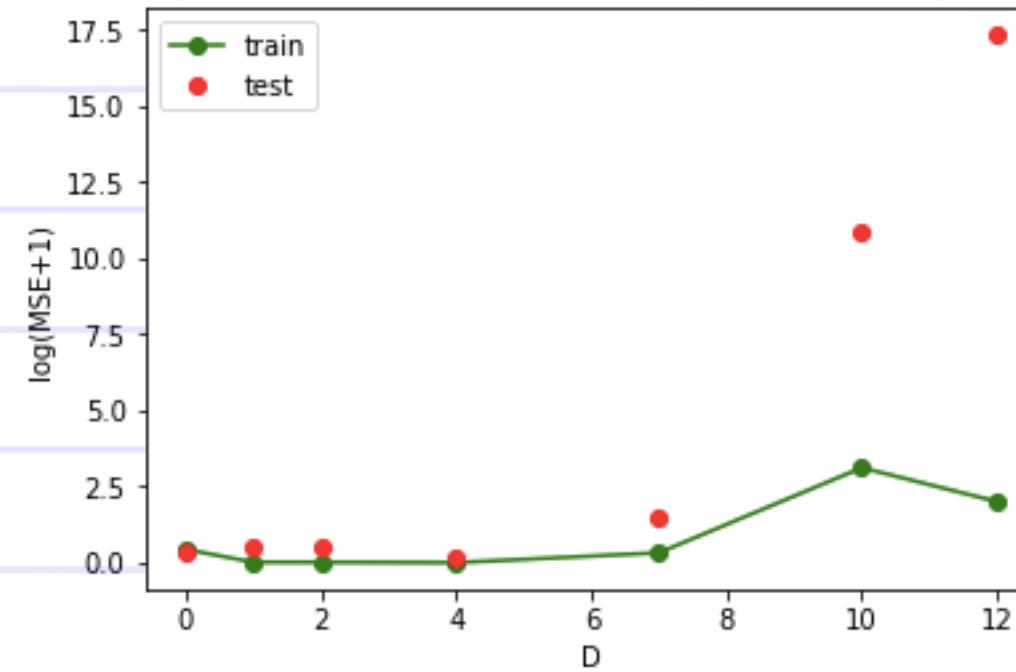
# Effect of increasing data (fixed model size)

---



# Effect of changing model size (same data)





# L2 Regularization or Ridge Regression

$$y = \Phi \omega$$

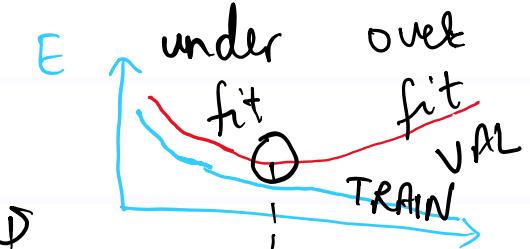
$$E = (\mathbf{t} - \Phi \omega)^T (\mathbf{t} - \Phi \omega) + \lambda \omega^T \omega$$

$$\frac{\partial E}{\partial \omega} = (-\Phi^T \mathbf{t} + \Phi^T \Phi \omega + \lambda \omega) \times 2$$

$$\Rightarrow \omega = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

RIDGE REGRESSION

$$\lambda \omega^T \omega$$



model  
capacity

Validation

## Another L2 normalization: Min-norm

$$D+1 > N$$

$$E = \|t - \Phi w\| + \lambda w^T w$$

$$\frac{\partial E}{\partial w} = -\Phi^T \operatorname{sgn}(t - \Phi w) + 2\lambda w$$

$\parallel_0$

$$w = \frac{1}{2\lambda} \Phi^T \operatorname{sgn}(t - \Phi w) \quad \dots \quad (i)$$

Consider,  $\underline{t = \Phi w} = \frac{\Phi \Phi^T}{2\lambda} \operatorname{sgn}(t - \Phi w)$

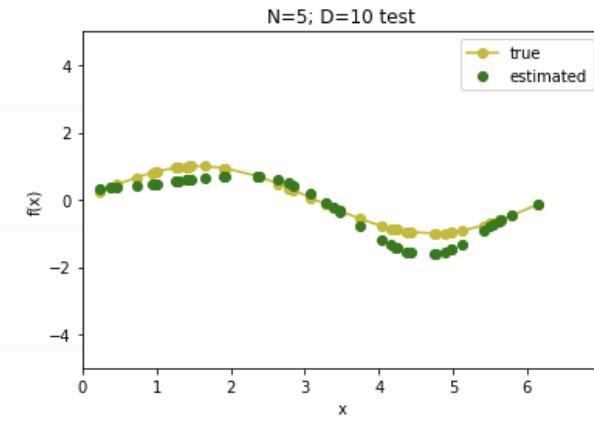
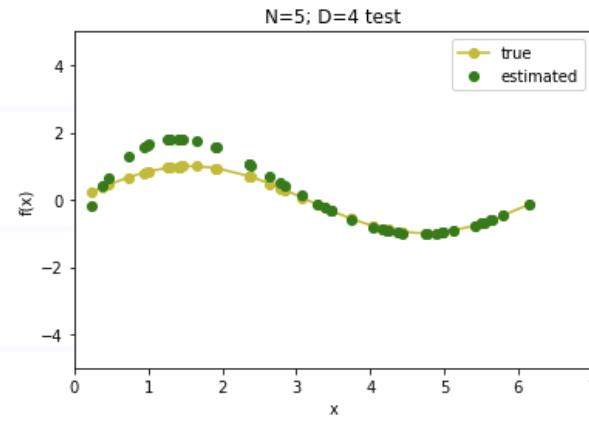
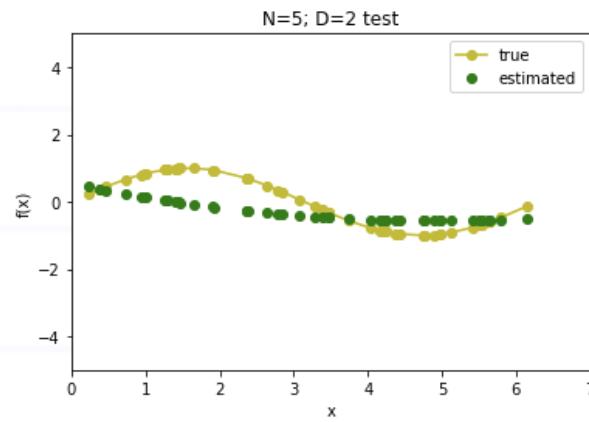
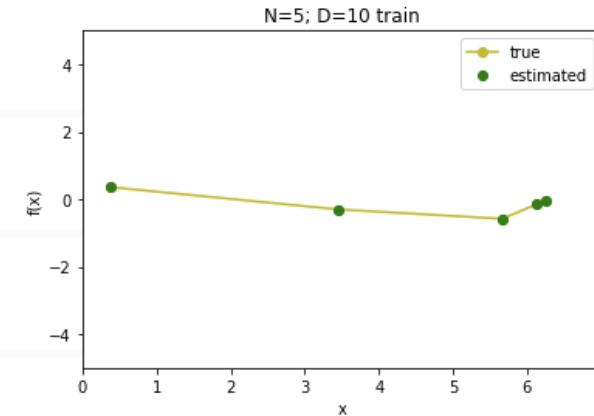
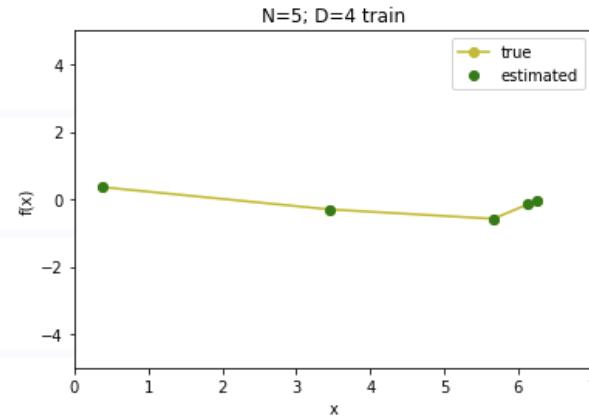
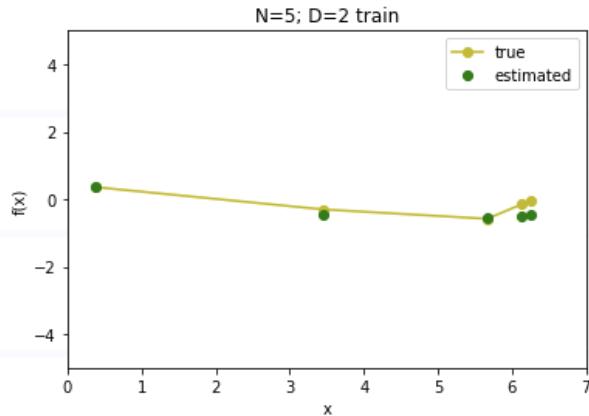
$$\Rightarrow \frac{1}{2\lambda} \operatorname{sgn}(t - \Phi w) = (\Phi \Phi^T)^{-1} t$$

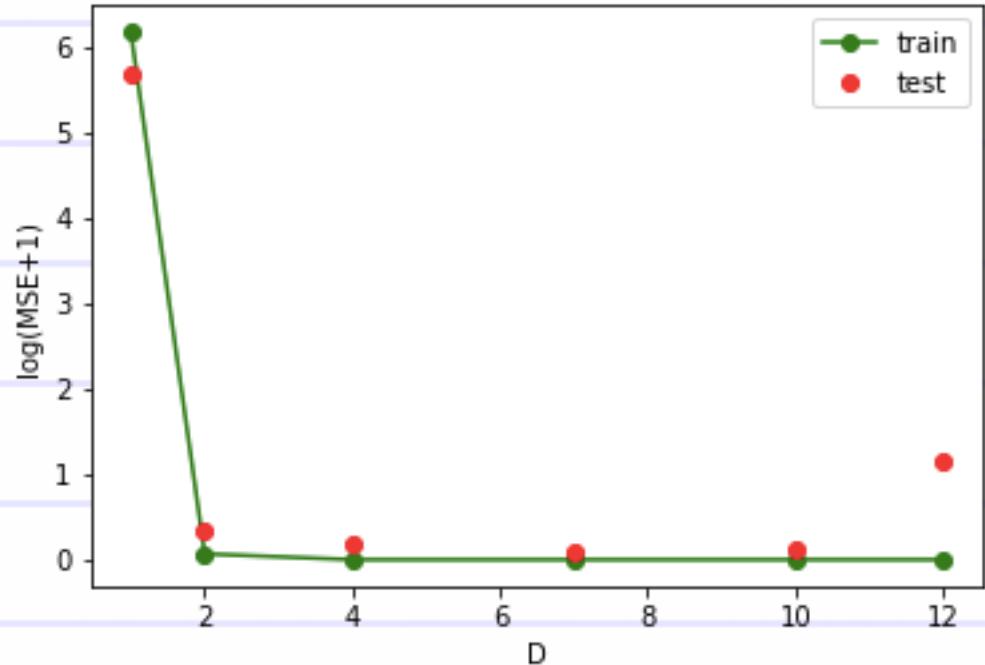
$\therefore$  from (1),

$$\hat{w} = \bar{\Phi}^T (\bar{\Phi} \bar{\Phi}^T)^{-1} t$$

Minimum norm  
solution

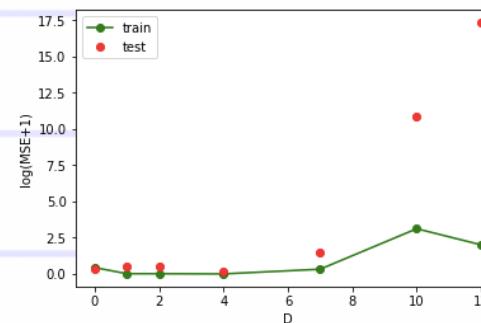
# Effect of changing model size (fixed data)





- Does better for large  $D$ , but worse for small  $D$

Recall:



- Caution: I have used single run, not averaged over many runs
- A hybrid approach may help

# Regularization

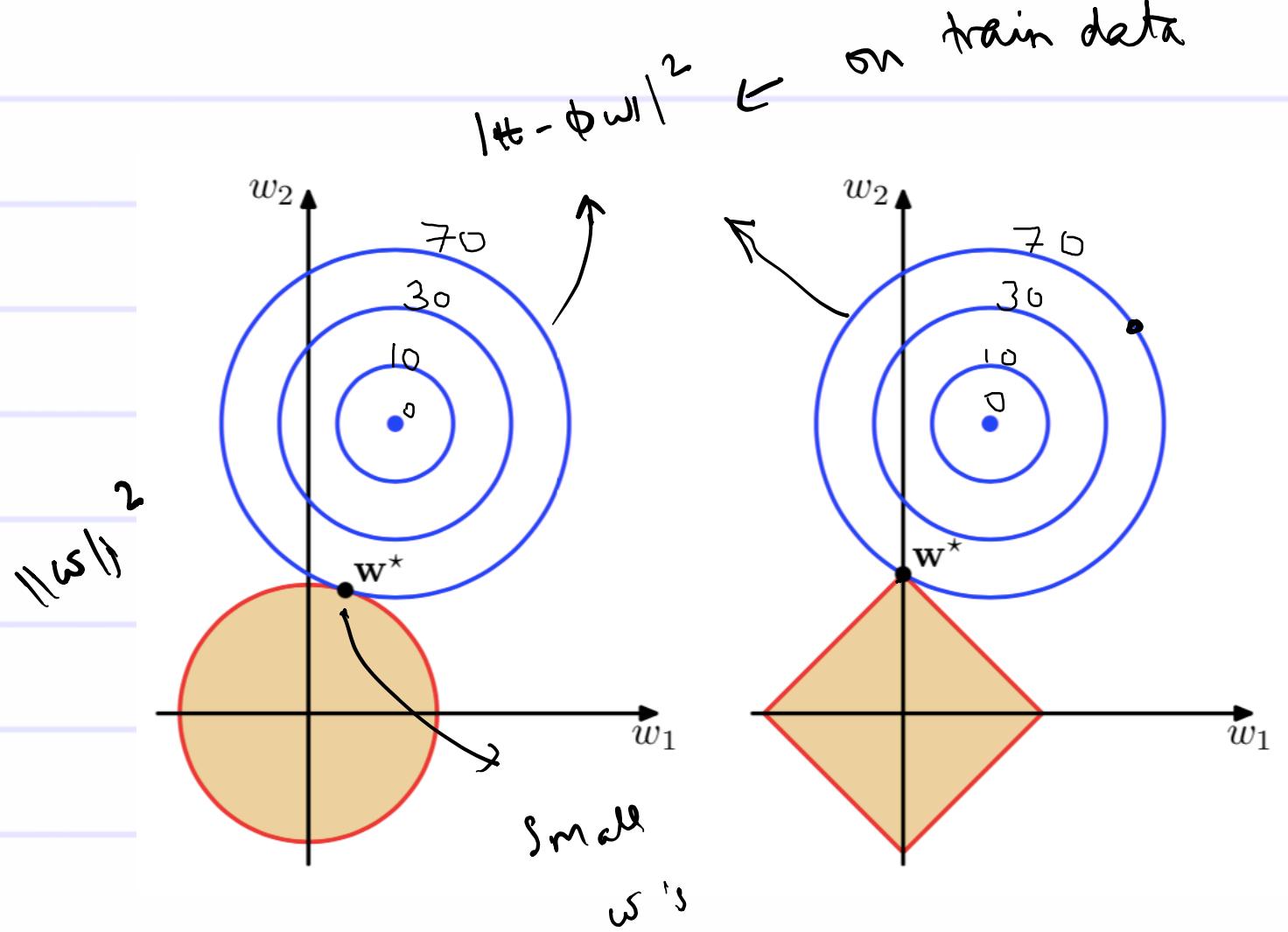
---

In general, with least squares loss,

$$E = \underbrace{(t - \Phi w)^T(t - \Phi w)}_{\text{least sq. loss}} + \lambda \sum_j |w_j|^q$$

$q=1$ , LASSO

$q=2$ , Ridge Regression



PRML: Fig 3.4

feature selection

$$\begin{aligned}
 y &= \Phi w \\
 &= \begin{bmatrix} 1 & x_1 & x_2 & \dots \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}
 \end{aligned}$$

# Multiple Outputs

Read PRML Sec. 3.1.5

$$y_i \rightarrow y_i \quad \forall i = 1, 2, \dots, 5$$

$y_i$  is  $K \times 1$  vector

$$\begin{bmatrix} y_1^\top \\ y_2^\top \\ \vdots \\ y_5^\top \end{bmatrix} = \Phi \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^D \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_5 & x_5^2 & \dots & x_5^D \end{bmatrix} \begin{bmatrix} w_{11} & \dots & w_{1K} \\ w_{21} & & w_{2K} \\ \vdots & & \vdots \\ w_{D1} & & w_{DK} \end{bmatrix}$$

$\textcolor{red}{Y} \quad 5 \times K$        $\textcolor{red}{\Phi} \quad 5 \times D$        $\textcolor{red}{W} \quad D \times K$

$$W = (\bar{\Phi}^\top \bar{\Phi})^{-1} \bar{\Phi}^\top T$$

same derivation  
as seen before

## Multiple Inputs

$$y, \Phi = [1 \ e^{(x-1)^2} \ e^{(x-2)^2} \dots]$$

- $\Phi$  can be a function of  $x$ , instead of just  $x^d$   
*polyn.*

- In general,  
for single sample  $y^\top = \Phi(x) w$

i.e.  $y_j = \sum_i \phi_i(x) w_{ij}$   $j=0, 1, 2, \dots, N_y - 1$   
 $i=0, 1, 2, \dots, D$

for multiple  
samples

$$y = \Phi(x) w$$

$$y_{j,t} = \sum_i \phi_i(x_t) w_{ij}$$

$j = 0, 1, \dots, N_y - 1$   
 $i = 0, 1, \dots, D$   
 $t = 0, 1, \dots, T - 1$   
samples

# References

---

- "Pattern Recognition and Machine Learning", C.M. Bishop, 2nd Edition, Springer, 2011.
  - **Section 1.1 (highly recommended)**
  - **Section 3.1 (highly recommended)**
- Matrix Calculus – wikipedia

[https://en.wikipedia.org/wiki/Matrix\\_calculus](https://en.wikipedia.org/wiki/Matrix_calculus)