

## Lec 6.1 Least squares Geometry

1. Recall:

- We're try to find  $\underline{w}$  that minimize

$$\min_{\underline{w}} \|\underline{y} - \underbrace{\underline{X}\underline{w}}_{\text{LS}}\|^2$$

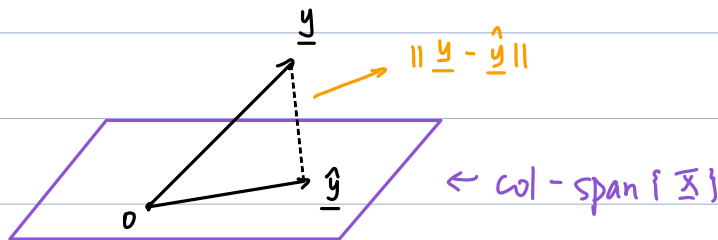
- LS solution:  $\underline{w}_{LS} = \underbrace{(\underline{X}^T \underline{X})^{-1}}_{\underline{X}^+} \underline{X}^T \underline{y}$

2. Moving on:

↳  $\hat{\underline{y}} = \underline{X} \underline{w}$ : linear comb. of col of  $\underline{X}$ .

$\Rightarrow \hat{\underline{y}}$  is a vector in col-span of  $\{\underline{X}\}$ .

$\therefore$  To minimize  $\|\underline{y} - \hat{\underline{y}}\|$  (distance between  $\hat{\underline{y}}$  and  $\underline{y}$ )



- The best  $\hat{\underline{y}}$  (i.e.  $\hat{\underline{y}}_{LS}$ ) is the projection of  $\underline{y}$  onto col-span  $\{\underline{X}\}$ .

$\Rightarrow$  Every col. of  $\underline{X}$  is orthogonal to  $\underline{y} - \hat{\underline{y}}$

$$\hookrightarrow \underline{a} \perp \underline{b} \Leftrightarrow \underline{a}^T \underline{b} = 0$$

$$\Rightarrow \underline{X}^T (\underline{y} - \hat{\underline{y}}_{LS}) = 0$$

$$\Rightarrow \underline{X}^T (\underline{y} - \underline{X} \underline{w}_{LS}) = 0$$

$$\Rightarrow \underline{X}^T \underline{X} \underline{w}_{LS} = \underline{X}^T \underline{y}$$

$$\Rightarrow \underline{w}_{LS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

// Same result, but from geometric view.

## Lec 6.2 Regularized linear regression / Least squares

### 1. Regularized version of Lin. regr:

$$\hookrightarrow \underbrace{\min_{\underline{w}} \|\underline{X} \underline{w} - \underline{y}\|^2}_{\text{simple version}} + \underbrace{\lambda \|\underline{w}\|^2}_{\text{penalty func. (against large } \|\underline{w}\|)} \\ \text{if } \underline{w} \text{ is large, penalty } \uparrow$$

#### Motivation:

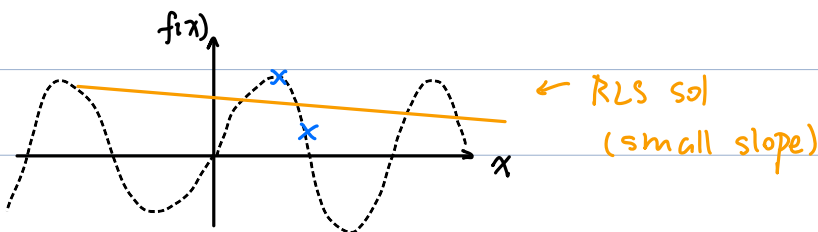
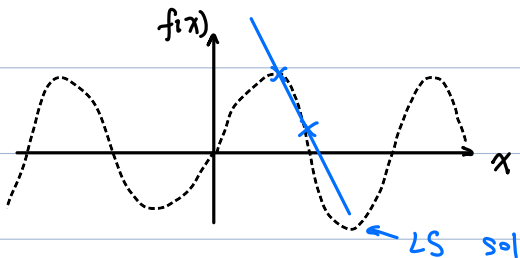
Want  $\|\underline{w}\|$  small to avoid over-fitting

due to ① noisy data ② Not enough data

### 2. Text book example (can skip):

↳ Target func:  $f(x) = \sin(\pi x)$

Training set: 2 points (randomly sampled)



← RLS sol  
(small slope)

$$\hookrightarrow f(\underline{w}) = \|\underline{X} \underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|^2$$

$$\frac{\partial f}{\partial w_k} = \frac{\partial}{\partial w_k} \left( \|\underline{X} \underline{w} - \underline{y}\|^2 + \lambda \sum_{j=0}^d w_j^2 \right) \quad // \text{break into one}$$

$$= 2 [\underline{X}^T (\underline{X} \underline{w} - \underline{y})]_k + 2\lambda w_k \quad \text{↗ k-th item}$$

$$= 2 [\underline{X}^T (\underline{X} \underline{w} - \underline{y}) + \lambda \underline{w}]_k$$

$$\nabla f(\underline{w}) = 2 [\underline{X}^T (\underline{X} \underline{w} - \underline{y}) + \lambda \underline{w}]$$

$$\text{We want } \nabla f(\underline{w}) = 0$$

$$\therefore (\underline{X}^T \underline{X} + \lambda \underline{1}) \underline{w}_{\text{RLS}} = \underline{X}^T \underline{y}$$

$$\underline{w}_{\text{RLS}} = (\underline{X}^T \underline{X} + \lambda \underline{1})^{-1} \underline{X}^T \underline{y}$$

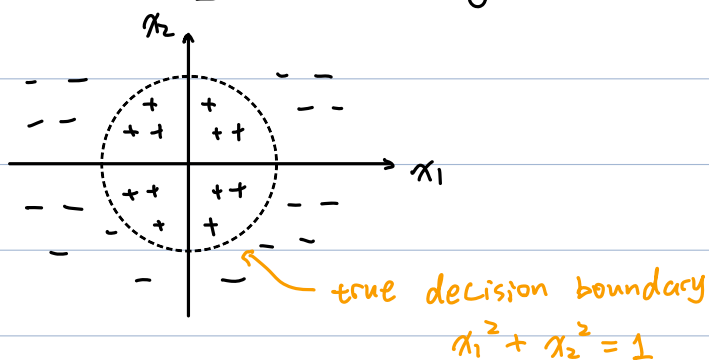
Note: ① If  $\lambda = 0$ , then RLS is same as LS form

② RLS side benefit more numerically stable solution.

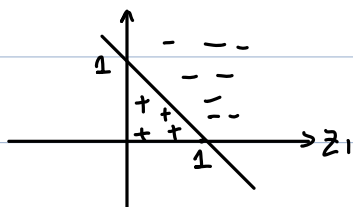
### Lec 6.3 Non-linear transformation

1. linear func may not be enough:

e.g.



↳ We transfer this into  $z_1 = x_1^2$ ,  $z_2 = x_2^2$



$$\Rightarrow z_1 + z_2 = 1$$

↳ Suppose  $h(\underline{z}) = \text{sign}(z_1 + z_2 - 1)$

Then our solution  $g(\underline{x}) = \text{sign}(x_1^2 + x_2^2 - 1)$

∴ We can transfer to another space, and use Lin. Regression

## 2. General non-linear regression:

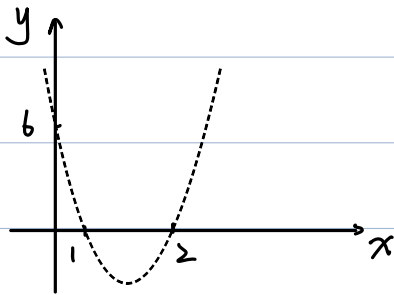
↳ Let  $\underline{z} = \Phi(\underline{x})$  be a non-linear transform

$\underline{h}(\underline{z}) = \underline{w}^T \underline{z}$  be a linear classifier in  $\underline{z}$  space

Then  $\underset{\substack{\uparrow \\ \text{classifier in terms of } x}}{g(\underline{x})} = h(\Phi(\underline{x}))$  is a non-linear classifier in  $\underline{x}$  space

//  $\Phi(\cdot)$  is called a feature transform

## 3. e.g. Quadratic Regression:



$$\underline{z} = (z_0 = 1, z_1 = x, z_2 = x^2) \quad // d=2$$

$$\therefore y = \underline{w}^T \underline{z}$$

$$= w_0 + w_1 z_1 + w_2 z_2$$

$$= w_0 + w_1 x + w_2 x^2$$

$$\therefore \underline{z}_1 = (1, 0, 0), \quad y_1 = 6$$

$$\underline{z}_2 = (1, 1, 1), \quad y_2 = 0$$

$$\underline{z}_3 = (1, 2, 4), \quad y_3 = 0$$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ 1 & x & x^2 \end{matrix}$

$$\therefore \underline{Z} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}, \quad \underline{y} = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$$

$$\therefore \underline{w}_{LS} = \underbrace{(\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T}_{\text{pseudo inverse of } \underline{Z}} \underline{y}$$

$$\therefore \underline{Z} \text{ is in full rank} \quad \therefore \underline{Z}^+ = \underline{Z}^{-1}$$

$$\therefore \underline{w}_{LS} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y} = \underline{Z}^{-1} \underline{y} = \begin{bmatrix} 6 \\ -9 \\ 3 \end{bmatrix}$$

$$\therefore y = 6 - 9z_1 + 3z_2 = 6 - 9x + 3x^2$$