

## Week 02 - Part 3

Review:

- Deriving  $\underline{w}_{ls}$  by finding the gradient and setting it to zero
- Deriving  $\underline{w}_{ls}$  by (pseudo)-solving the system of linear equations  $\underline{y} = X \underline{w}$ .

Today:

- Deriving  $\underline{w}_{ls}$  with Geometric interpretations
- Regularized Least squares
- Non-linear transformation

Recall:

■ Least square solution:  $\underline{w}_{ls} = X^T \underline{y} = (X^T X)^{-1} X^T \underline{y}$

■ Prediction by  $\underline{w}_{ls}$ :  $\hat{\underline{y}}_{ls} = X \underline{w}_{ls} = X X^T \underline{y}$

■ It's like we take  $\underline{y}$  and with a projection matrix transforming it into  $\hat{\underline{y}}_{ls}$ .

■  $XX^T$  is a projection matrix.

■ This observation leads us into geometric interpretation of Least squares.

## Geometric Interpretation of Least Squares

■ Observe that  $\hat{\underline{y}} = X \underline{w} =$

$$\begin{bmatrix} x_{10} & x_{11} & \dots & x_{1d} \\ x_{20} & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \dots & \vdots \\ x_{N0} & x_{N1} & \dots & x_{Nd} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$= ?$

■ So,  $\hat{\underline{y}}$  is linear combination of ...

■ Thus,  $\hat{y}$  is in the space of all possible linear combinations of columns of  $X$

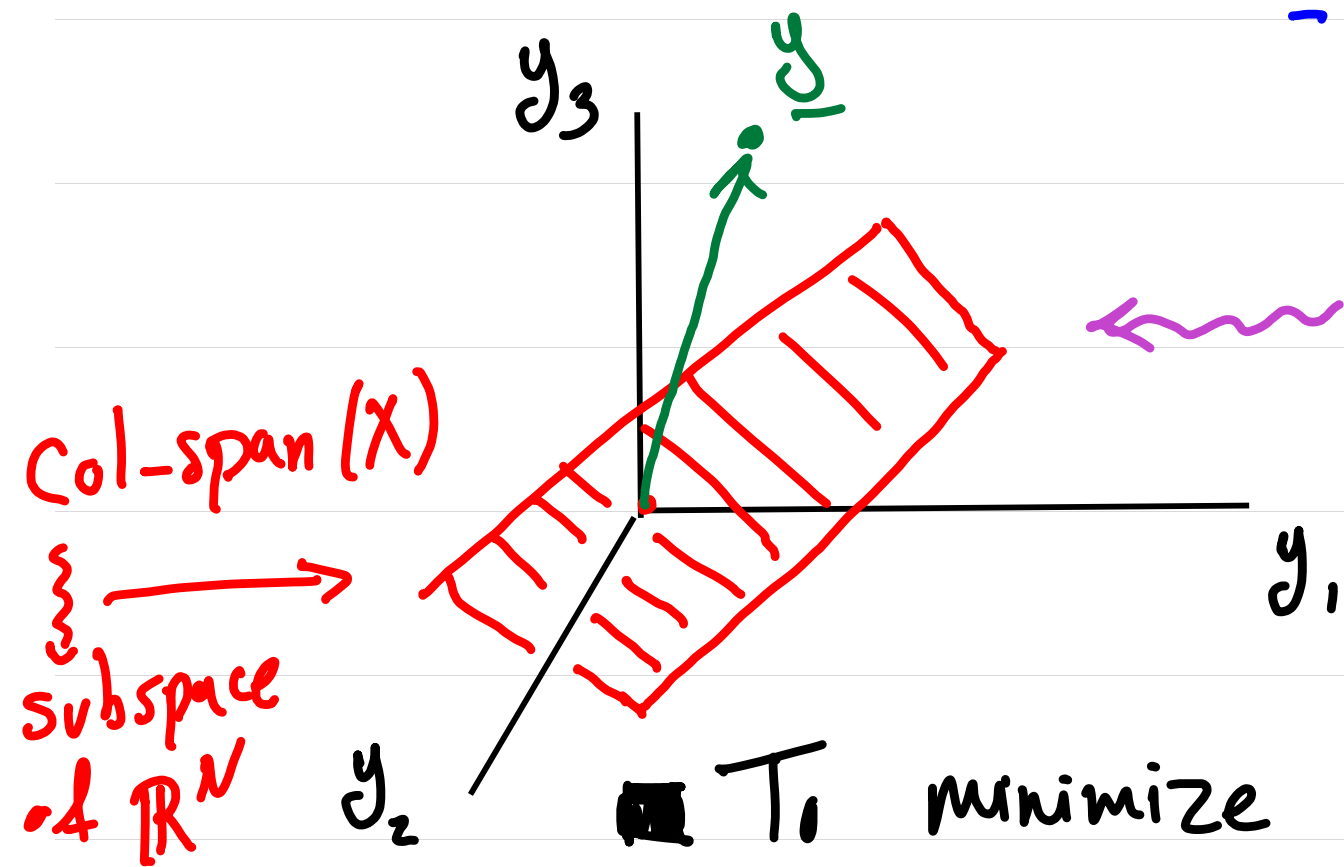
■ The space of all possible linear combination of columns of  $X$  is called  $\text{Col-span}\{X\}$

■ Let's illustrate  $\text{Col-span}(X)$  for  $N=3$ ,  $d=1$ , and  $X = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix}$ .

$$\text{Col-span}(X) =$$

$$\blacksquare \text{Col-span}(X) = \left\{ \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} : -y_1 + 2y_2 - y_3 = 0 \right\}$$

- A plane that goes through origin.
- It is a subspace of  $\mathbb{R}^N$ .



← the space of possible  $\hat{y}$ .

$\blacksquare$  To minimize  $\|\underline{y} - \hat{\underline{y}}\|$  : Euclidean distance b/w  $\underline{y}$  &  $\hat{\underline{y}}$

$\blacksquare$  must find  $\hat{\underline{y}}$  on  $\text{Col-span}(X)$  that is closest to  $\underline{y}$ .

■ The best  $\hat{\underline{y}}$  (i.e.  $\hat{\underline{y}}_{ls}$ ) is the projection of  $\underline{y}$  onto  $\text{Col-span}\{X\}$ .

■ That means  $(\underline{y} - \hat{\underline{y}}_{ls})$  must be orthogonal to any vector in  $\text{Col-span}\{X\}$ .

■ Thus,  $(\underline{y} - \hat{\underline{y}}_{ls})$  is orthogonal to every column of  $X$ .

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

■ Thus,  $X^T (\underline{y} - \hat{\underline{y}}_{ls}) = \underline{0} \Rightarrow$

# Regularized Linear Regression / Least squares

■ Previously, we tried to minimize  $\|X\underline{w} - \underline{y}\|^2$

■ In regularized version, we minimize  $\|X\underline{w} - \underline{y}\|^2 + \underbrace{\lambda \|\underline{w}\|^2}_{\text{penalty function (against large weight)}}$

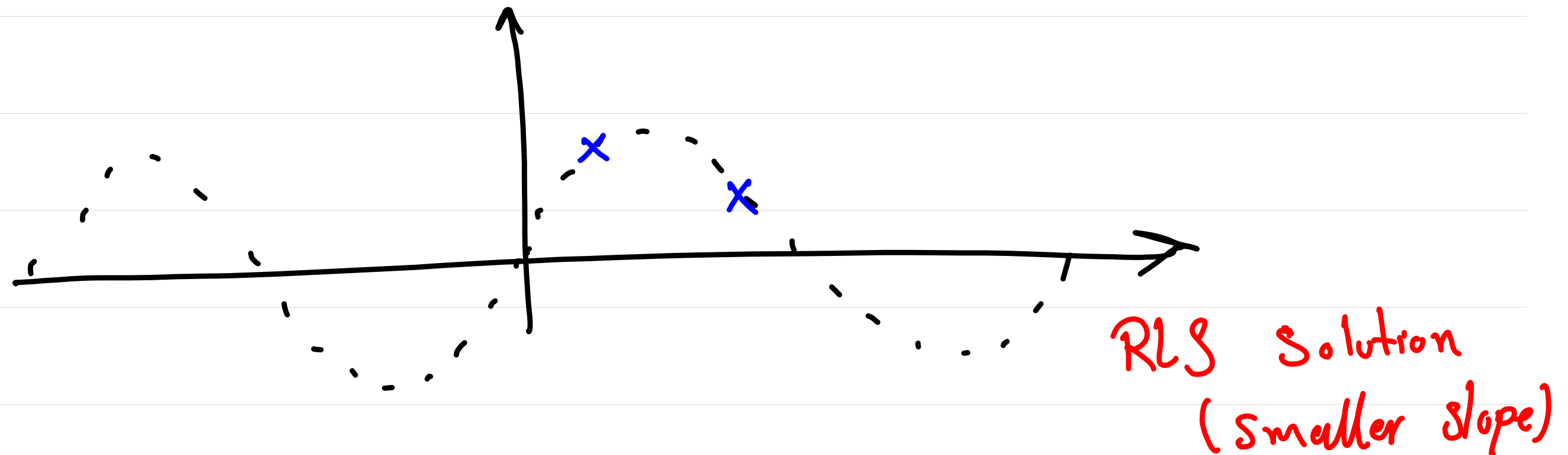
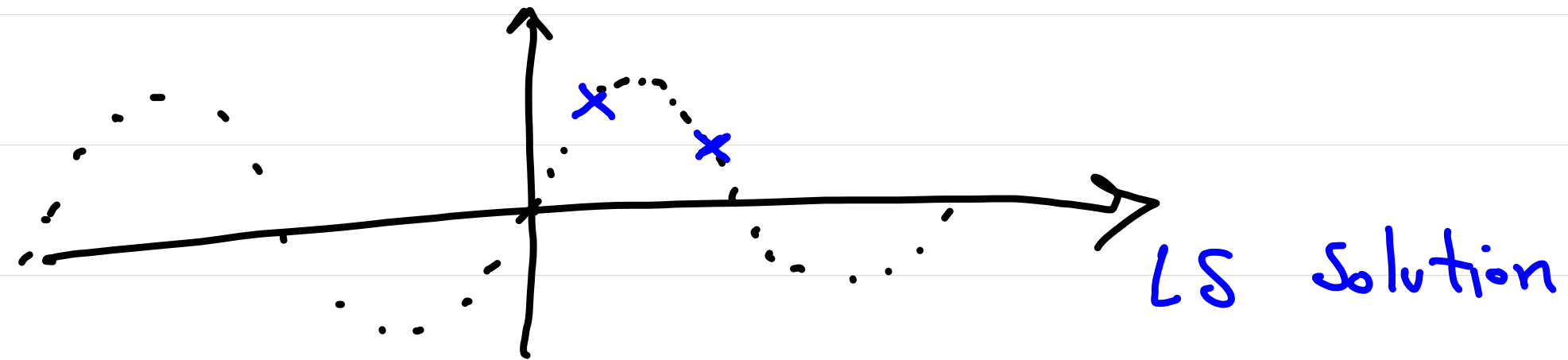
■ The motivation is to avoid overfitting

☞ your data is noisy

☞ you do not have enough data (compared to the complexity of the target function)



■ E.g. target:  $f(x) = \sin(\pi x)$



Note: ①  $\lambda = 0 \Rightarrow LS$

② How to choose  $\lambda$ ? validation

How do we solve this Problem?

■ We want to  $\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|^2$

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

■ Observe that  $\nabla_{\underline{w}} f(\underline{w}) = 2\underline{X}^T(\underline{X}\underline{w} - \underline{y}) + 2\lambda \underline{w}$

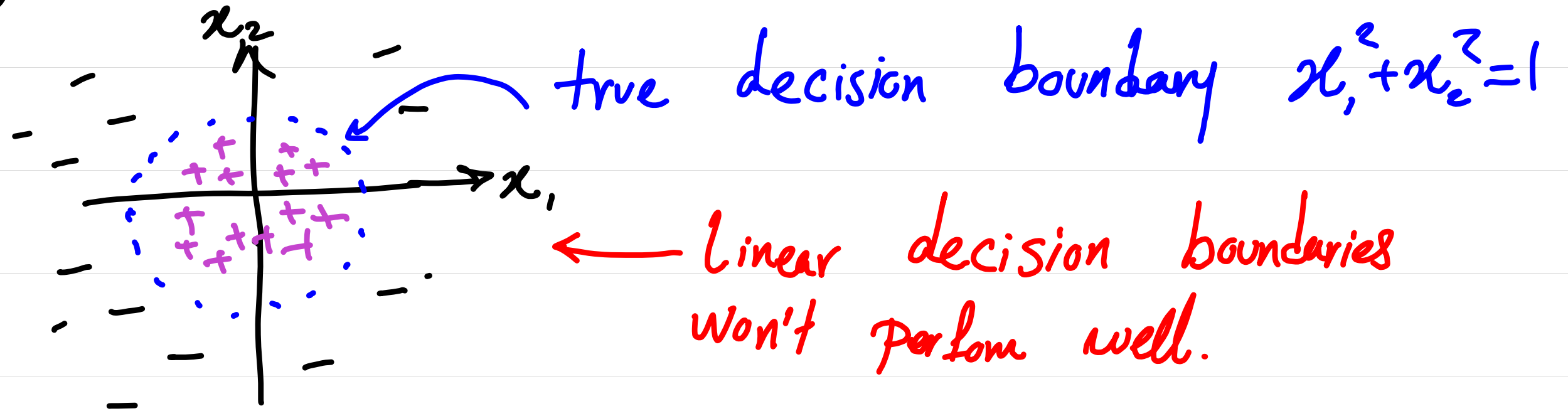
■ We want  $\nabla_{\underline{w}} f(\underline{w}) = 0 \Rightarrow (\underline{X}^T \underline{X} + \lambda \underline{I})\underline{w} = \underline{X}^T \underline{y}$

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

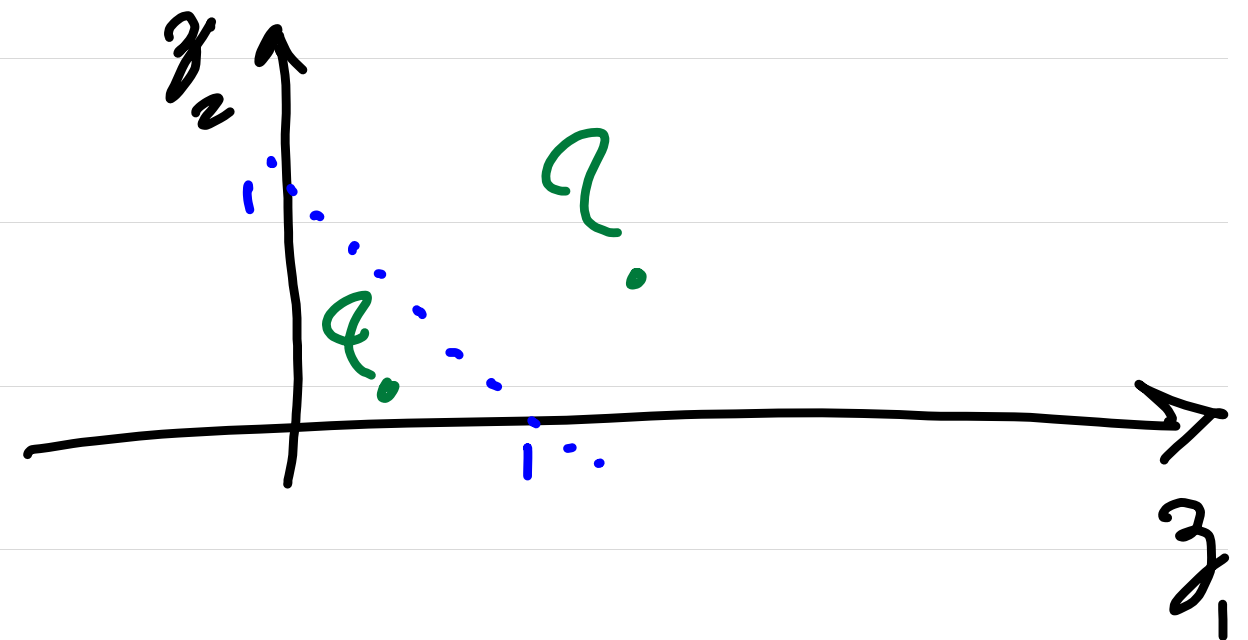
So far: We studied Linear Models

But in many cases linear Models are not good enough.

E.g.



Then define  $z_1 = x_1^2$  and  $z_2 = x_2^2$



The points are

linearly separable in  $Z$ -space.

Suppose PLA gives you  $h(\underline{z}) = \text{Sign}(z_1 + z_2 - 1)$ . Then, we know  $g(\underline{x}) = \text{Sign}(x_1^2 + x_2^2 - 1)$

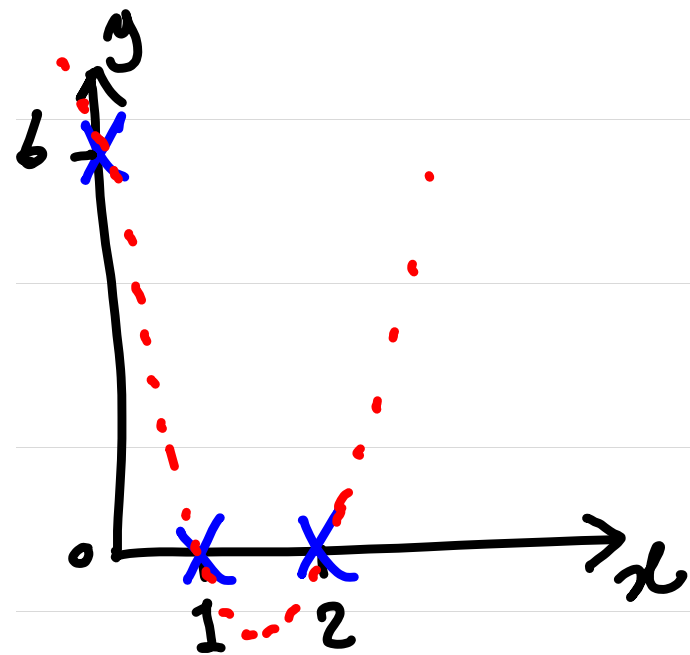
In general:

■ Let  $\underline{z} = \Phi(\underline{x})$  be non-linear transformation  
(~feature transformation~)

■ Let  $h(\underline{z})$  be a linear classifier/regression function  
in  $\underline{z}$  space ( $h(\underline{z}) = \text{Sign}(\underline{w}^T \underline{z})$  or  $h(\underline{z}) = \underline{w}^T \underline{z}$ )

■ Then  $g(\underline{x}) = h(\Phi(\underline{x}))$  is non-linear classifier in  
 $\underline{x}$  space

## E.g. Quadratic Regression



■ Define  $\underline{z} = (z_0=1, z_1=x, z_2=x^2)$

■  $\underline{y} = \underline{w}^T \underline{z} = w_0 + w_1 z_1 + w_2 z_2$  ← Linear in  $\underline{z}$

$= w_0 + w_1 x + w_2 x^2$  ← Quadratic in  $\underline{x}$

■ Let's find the  $\underline{w}_{LS}$ :  $\underline{x}_1 =$  ,  $\underline{x}_2 =$  ,  $\underline{x}_3 =$   
 $\underline{z}_1 =$  ,  $\underline{z}_2 =$  ,  $\underline{z}_3 =$  →  $\underline{Z} =$

$\underline{y} =$

$\hat{\underline{y}} = \underline{w}_{LS}^T \underline{z} =$

$\underline{w}_{LS} =$

$= \dots = \Rightarrow$