

# Supervised Learning for Regression

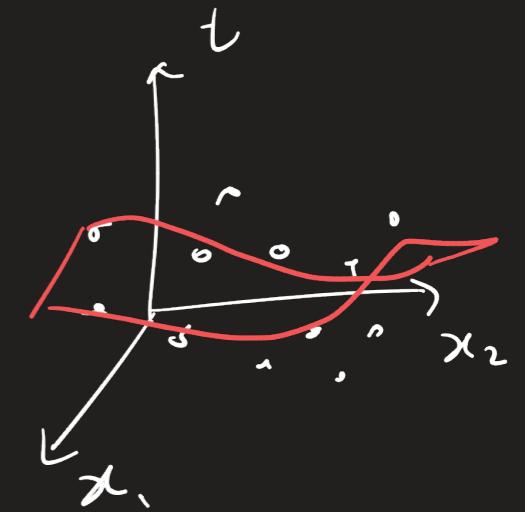
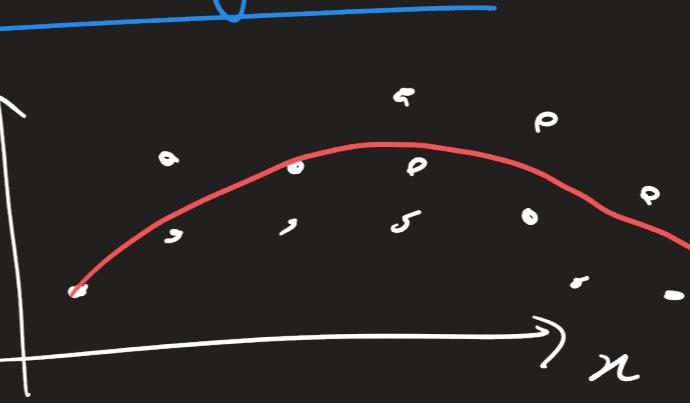
## Tasks

### Linear Regression

① INPUT SPACE  
Input data:  $\{x_i\}_{i=1}^N$ ,  $N = \# \text{ samples}$   
 $x_i \in \mathbb{R}$

Target values:  $\{t_i\}_{i=1}^N$

DATASET:  $\{(x_i, t_i)\}_{i=1}^N$



② FEATURE SPACE

Polynomial basis function:

$$\phi(x_i) = [x_i^0, x_i^1, x_i^2, \dots, x_i^M]^T \in \mathbb{R}^{(M+1)}$$

### ③ MODEL / MAPPER :

$$y_i = f(\phi(x_i), w)$$

$$\boxed{y_i = w_0 \cdot x_i^0 + w_1 \cdot x_i^1 + \dots + w_M \cdot x_i^M}$$

$f = \text{linear regression model}$

$$= w^T \phi(x_i)$$

$$= \phi(x_i)^T w$$

$w$  = generalized notation for model parameters

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

$$\phi(x_i) = \begin{bmatrix} x_i^0 \\ x_i^1 \\ \vdots \\ x_i^M \end{bmatrix}$$

$y_i$  = model prediction for  $i^{th}$  sample

$$\left\{ \begin{array}{l} y_1 = w_0 \cdot x_1^0 + w_1 \cdot x_1^1 + \dots + w_M \cdot x_1^M \\ y_2 = w_0 \cdot x_2^0 + w_1 \cdot x_2^1 + \dots + w_M \cdot x_2^M \\ \vdots \\ y_N = w_0 \cdot x_N^0 + w_1 \cdot x_N^1 + \dots + w_M \cdot x_N^M \end{array} \right.$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_1^0 & x_1^1 & \dots & x_1^M \\ x_2^0 & x_2^1 & \dots & x_2^M \\ \vdots & \vdots & \ddots & \vdots \\ x_N^0 & x_N^1 & \dots & x_N^M \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

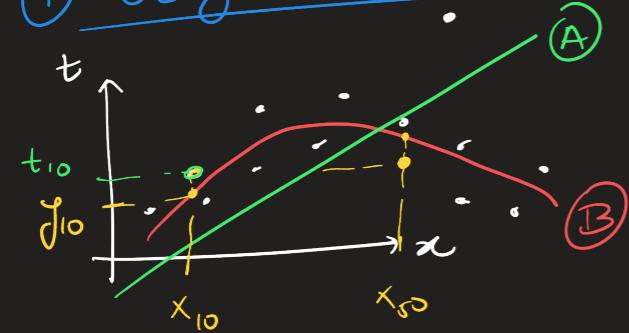
FEATURE MATRIX  $\equiv X$   
 $N \times (M+1)$

output vector  $\equiv y$   
 Model parameter  $\equiv w$

MODEL /  
MAPPER:

$$y = X \cdot w$$

#### ④ OBJECTIVE FUNCTION



$$\epsilon_i = t_i - y_i$$

$$J(\omega) = \frac{1}{N} \sum_{i=1}^N \epsilon_i^2$$

$$= \frac{1}{N} \sum_{i=1}^N (t_i - y_i)^2$$

$$= \frac{1}{N} \sum_{i=1}^N (t_i - \omega^T \phi(x_i))^2$$

$$= \frac{1}{N} (t - \mathbb{X} \omega)^T (t - \mathbb{X} \omega)$$

$$= \frac{1}{N} \|t - \mathbb{X} \omega\|_2^2$$

MEAN Least  
squares  
(MSE)

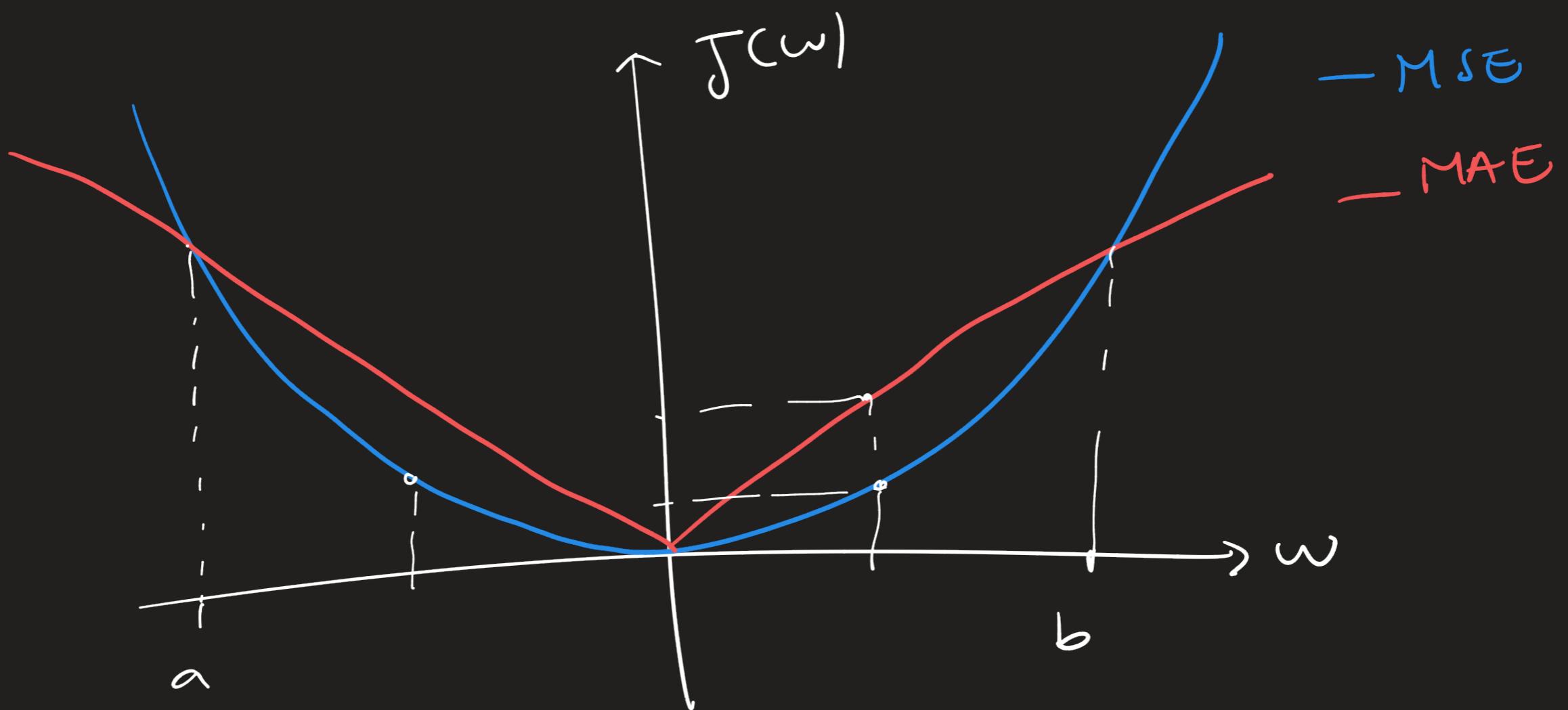
Reminder:  
 $\|\mathbf{x}\|_2 = \|\mathbf{x}\|$   
 $= (\mathbf{x}_1^2 + \mathbf{x}_2^2 + \dots + \mathbf{x}_n^2)^{1/2}$

$$\|\mathbf{x}\|_1 = |x_1| + |x_2| + \dots + |x_n|$$

An alternative, consider sum of absolute error for  $f_i$

$$J(\omega) = \frac{1}{N} \|t - \mathbb{X} \omega\|_1$$

MEAN ABSOLUTE ERROR  
(MAE)



## ⑤ Learning Algorithm

Finds the best solution for model parameters,  $w$ .

$$\frac{\partial J}{\partial w} = 0$$

$$J(w) = \frac{1}{N} \cdot (t - Xw)^T(t - Xw)$$

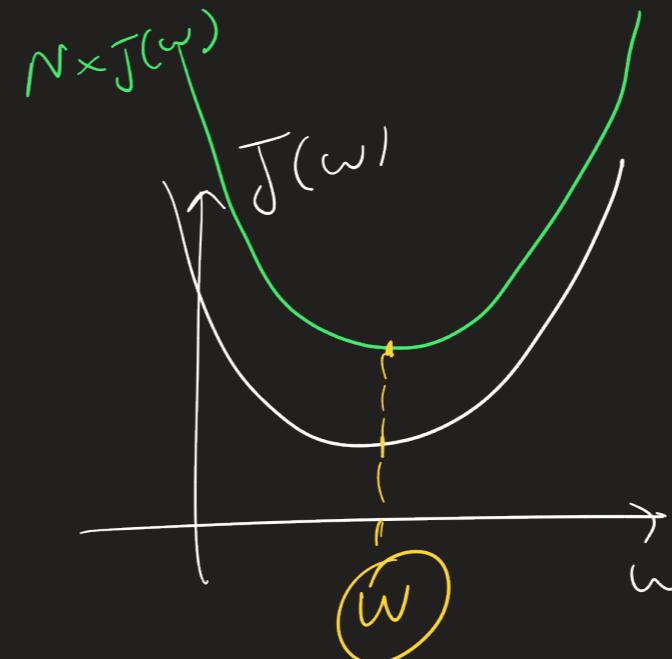
$$= \frac{1}{N} (t^T - w^T X^T)(t - Xw)$$

$$= \frac{1}{N} (t^T t - t^T X w - w^T X^T t + \underline{w^T X^T X w})$$

$$\frac{\partial J}{\partial w} = -t^T X - \cancel{(X^T t)} + \cancel{(X^T X w)} + \underline{w^T X^T X}$$

$$= -t^T X - t^T X + w^T X^T X + w^T X^T X$$

$$= -2t^T X + 2 \cdot w^T X^T X = 0$$



$$-\cancel{t^T \cancel{X}} + \cancel{w^T \cancel{X}^T \cancel{X}} = 0$$

$$w^T \cancel{X}^T \cancel{X} = t^T \cancel{X}$$

$$\cancel{X}^T \cancel{X} w = \cancel{X}^T \cdot t$$

Reminder  
cyclic  
operation  
 $(X^T a)^T$   
 $= a^T \cdot (X^T)^T$

$$w = (\cancel{X}^T \cancel{X})^{-1} \cdot \cancel{X}^T \cdot t$$

if  $(\cancel{X}^T \cancel{X})^{-1}$  (matrix inverse)  
Exists.

If  $X^T X$  is invertible:

1)  $\det(X^T X) \neq 0$

2) Full rank matrix

3) Eigenvalues of  $X^T X$  they

are all  $\neq 0$ .

4) Features are all linearly independent.

$$X = \left[ \begin{array}{c|c|c} f_1 & f_2 & f_3 \end{array} \right]$$

$f_1 \equiv \text{male} \in \{0, 1\}$

$f_2 \equiv \text{female} \in \{0, 1\}$

$f_3 \equiv \text{length} \in \mathbb{R}$

$$= \left[ \begin{array}{ccc} 1 & 0 & 3.5 \\ 0 & 1 & 4.2 \\ \vdots & \vdots & \vdots \end{array} \right]$$

If we have 3 features

then we want them to  
exist in a 3-D basis.

$$B = \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$$

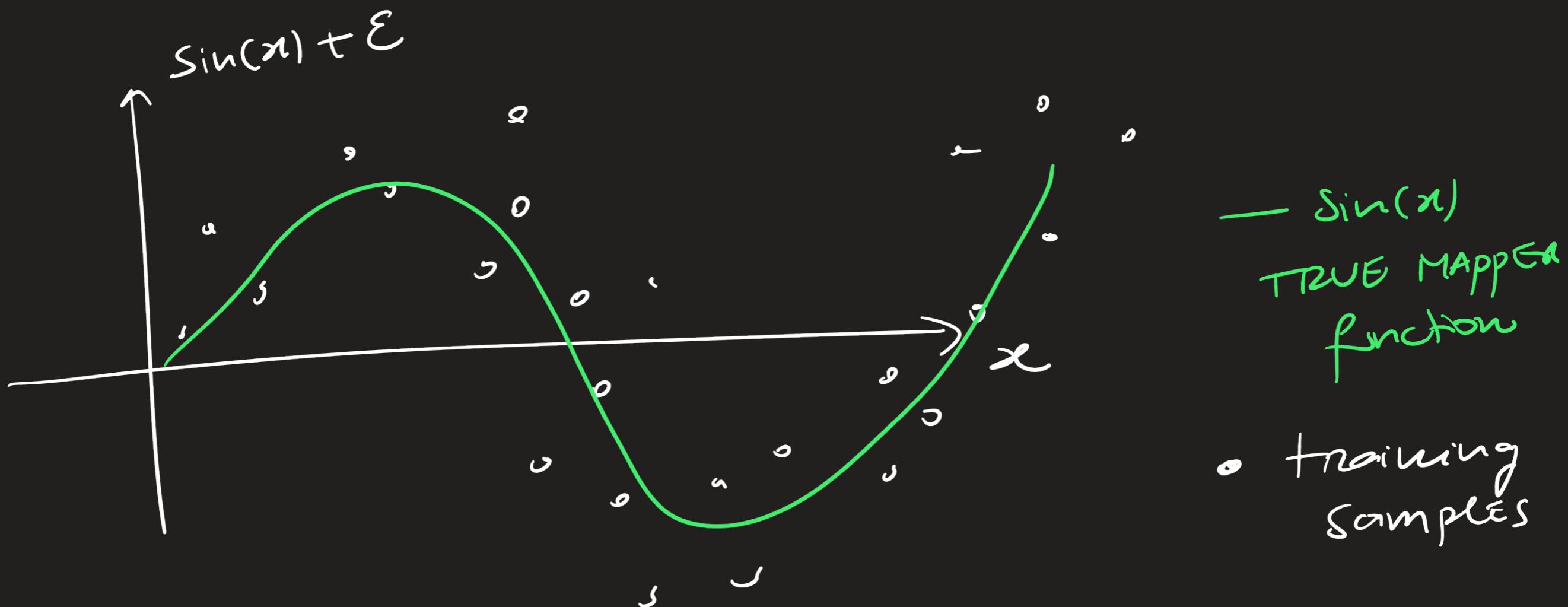
$$\text{NEW } X^T X + \lambda \cdot I$$

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Diagonally-loading

Hyperparameter ← provided by the user.

$M = \text{degree of polynomial}$



## OVERfitting

① Model is NOT able to generalize to new and unseen data.

② The parameters ↑ in magnitude

## Regularization

$$R_{L_2}(w) = \sum_{j=0}^M w_j^2$$

$$= w^T w$$

$$= \|w\|_2^2$$

L<sub>2</sub>-regularizer

or

RIDGE regularizer

$$R_{L_1}(w) = \sum_{j=0}^M |w_j|$$

$$= \|w\|_1$$

L<sub>1</sub>-regularizer

or

Lasso regularizer

→ "prefers" to have small values for all w's

→ promotes sparsity of parameters vector

Include the regularizer in  
the objective function.

$$J(w) = \|t - Xw\|_2^2 + \lambda \cdot \|w\|_2^2$$

squared error  
regularizer  
 $\lambda$  hyperparameter  
ridge regularizer

$\lambda \rightarrow 0$ : does NOT regularize model

$\lambda \rightarrow \infty$ : it will, in the limit, only consider range of values for  $w$ .

$$J(\omega) = (t - X\omega)^T(t - X\omega) + \lambda \cdot \omega^T \omega$$

$$= \underbrace{t^T t - t^T X \omega - \omega^T X^T t + \omega^T X^T X \omega}_{+ \lambda \omega^T \omega}$$

$$\frac{\partial J}{\partial \omega} = -2t^T X + 2\omega^T X^T X + 2\lambda \cdot \omega^T = 0$$

$$\lambda \omega^T + \omega^T X^T X = t^T X \quad | \omega \cdot I = \omega$$

$$\lambda \omega^T + X^T X \omega = X^T t$$

$$\underbrace{(\lambda I + X^T X)}_{(I + X^T X)}, \omega = X^T t$$

$$\boxed{\omega = (\lambda I + X^T X)^{-1} X^T t}$$

$$\begin{bmatrix} \lambda & X^T X \\ 1 \times 1 & (M+1) \times (N+1) \end{bmatrix}$$

diagonally loading  
 $X^T X$