

# Linear Regression w/ Polynomial Features

↳ Supervised Learning

TRAINING DATA :  $\{(x_i, t_i)\}_{i=1}^N$  ,  $x_i \in \mathbb{R}$   
 $t_i \in \mathbb{R}$

MODEL OR MAPPER :  $y = f(\phi(x), w)$  ,  $w \equiv$  parameters of the model

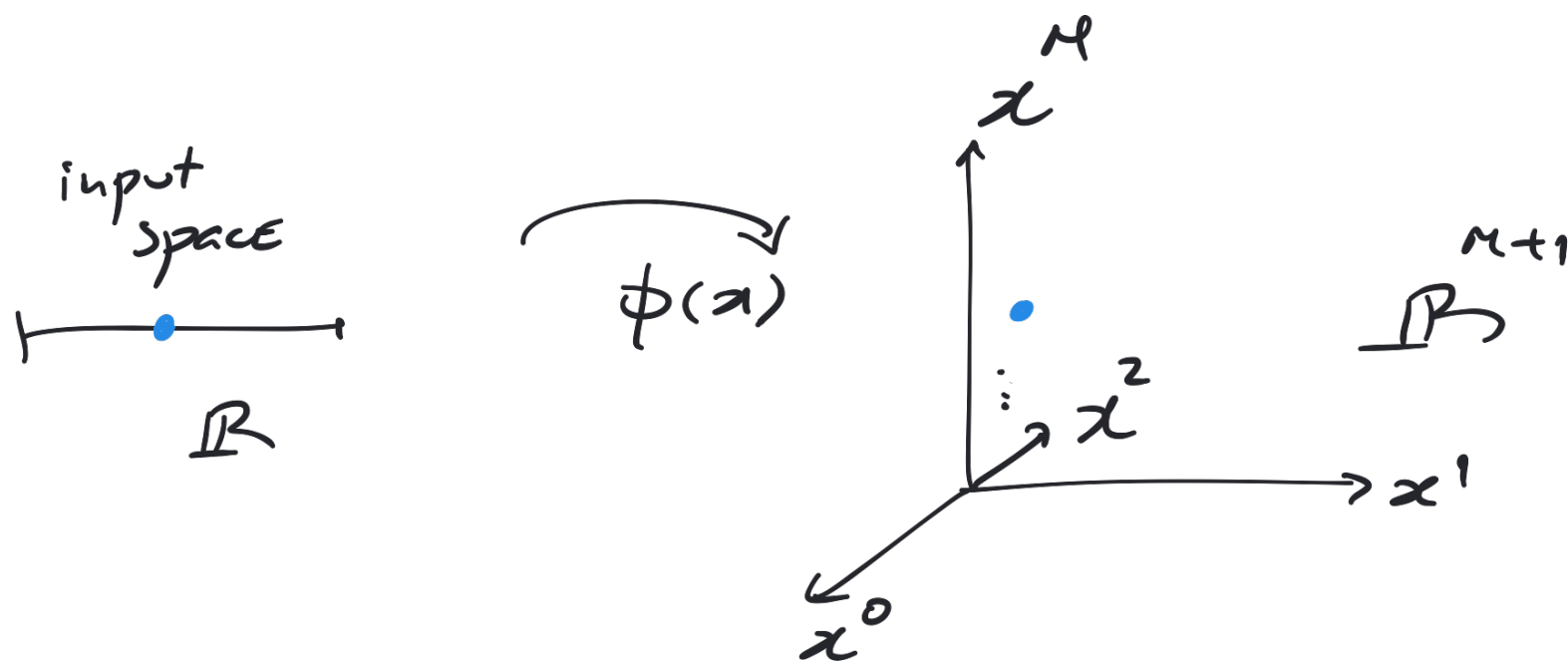
$$y = w_0 + w_1 \cdot x + w_2 \cdot x^2 = w^T \cdot \phi(x) = \phi^T(x) \cdot w$$

Polynomial Basis function :  $\phi(x) = [x^0, x^1, x^2]^T$

PARAMETERS :  $w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$

In general, the dimensionality of feature space is determined by the hyperparameter

$M \equiv$  model order



Hyperparameters  $\leftarrow$  are user-defined  
parameters  $\leftarrow$  are found using the  
hyperparameters for a  
given training set.

For each training sample  $x_i$ :

$$\phi(x_i) = [x_i^0, x_i^1, x_i^2, \dots, x_i^N]^T$$

Mapper:  $y(x_i) = \phi^T(x_i) \cdot w$

$$\begin{cases} y(x_1) = \phi^T(x_1) \cdot w \\ y(x_2) = \phi^T(x_2) \cdot w \\ \vdots \\ y(x_N) = \phi^T(x_N) \cdot w \end{cases}$$

$N \equiv \#$  samples  
in training

$M \equiv$  hyperparameter

FEATURE  
MATRIX :  $X = \begin{bmatrix} \phi^T(x_1) \\ \phi^T(x_2) \\ \vdots \\ \phi^T(x_N) \end{bmatrix}$   $N \times (M+1)$

$$X = \begin{bmatrix} x_1^0 & x_1^1 & x_1^2 & \dots & x_1^M \\ x_2^0 & x_2^1 & x_2^2 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N^0 & x_N^1 & x_N^2 & \dots & x_N^M \end{bmatrix}, \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

$(M+1) \times 1$

$y = X \cdot w$

$y$  is  $N \times 1$

Let  $t$  be the vector with the  
desired / target values :

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \quad N \times 1$$

Error of  
prediction :  $E = t - y$   
 $N \times 1$

OBJECTIVE  
FUNCTION

: ①

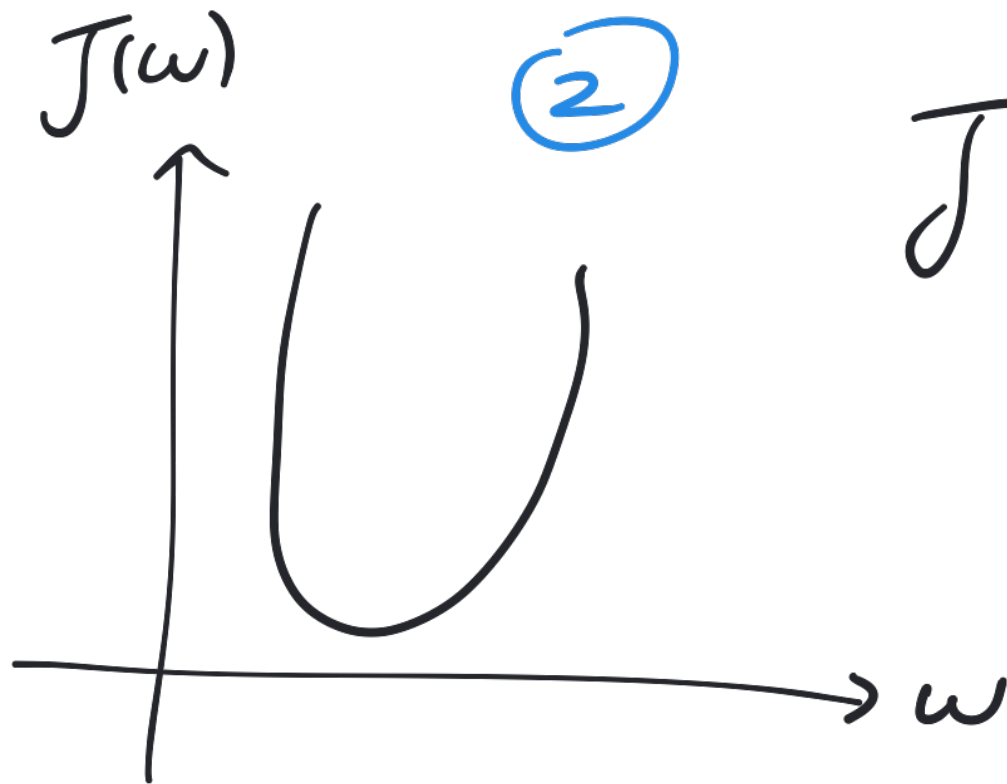
MEAN SQUARED ERROR

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

MEAN ABSOLUTE ERROR (MAE)

②

$$J(w) = \frac{1}{N} \sum_{i=1}^N |\epsilon_i|$$



Consider the objective function:

$$\begin{aligned} J(w) &= \frac{1}{2} \sum_{i=1}^N \epsilon_i^2 \\ &= \frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2 \\ &= \frac{1}{2} \sum_{i=1}^N (t_i - \phi^T(x_i) \cdot w)^2 \end{aligned}$$

$$= \frac{1}{2} (t - Xw)^T (t - Xw)$$

$$= \frac{1}{2} \|t - Xw\|_2^2$$

Reminder:

$$\|x\|_2^2 = x^T \cdot x$$

$$\|x\|_2 = (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2}$$

$$\|x\|_2 \equiv \text{L2-norm}$$



## Learning Algorithm

Pose the question:

$$\arg \min_w J(w) = w^*$$

NECESSARY  
condition:

$$\frac{\partial J(w^*)}{\partial w} = 0$$

$$J(\omega) = \frac{1}{2} \|t - X\omega\|_2^2$$

$$= \frac{1}{2} (t - X\omega)^T (t - X\omega)$$

$$= \frac{1}{2} (t^T - \omega^T X^T) (t - X\omega)$$

$$= \frac{1}{2} (t^T t - t^T X \omega - \omega^T X^T t + \omega^T X^T X \omega)$$

Reminder:

$$(AB)^T = B^T A^T$$

$$(AB^T)^T = B \cdot A^T$$

$$\frac{d(X \cdot \omega)}{d\omega} = X, \quad \frac{d(X \omega^T)}{d\omega} = X^T$$

---


$$\frac{\partial J(\omega)}{\partial \omega} = 0 \Leftrightarrow -t^T X - (X^T t)^T + (X^T X \cdot \omega)^T + \omega^T X^T X = 0$$

$$\Leftrightarrow -t^T X - t^T X + \omega^T X^T X + \omega^T X^T X = 0$$

$$\Leftrightarrow -t^T X + \omega^T X^T X = 0$$

$$\Leftrightarrow -X^T t + X^T X \cdot \omega = 0$$

$\downarrow$   
 applying  
 transposes

$\Rightarrow \textcircled{*}$

$$\textcircled{*} \Leftrightarrow X^T X w = X^T t$$

$\Leftrightarrow$



if  $(X^T X)^{-1}$  exists

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

pseudo-inverse of  $X$   
denoted as  $X^+$  ( $X$  "dagger")

$X$  is (typically) a tall matrix

$\Leftrightarrow$

$$w = X^+ t$$

## Food for thought:

- ① In what scenarios is  $X^T X$  not invertible?
- ② Can you force it to be invertible? How?