

# Linear Algebra for Computer Science

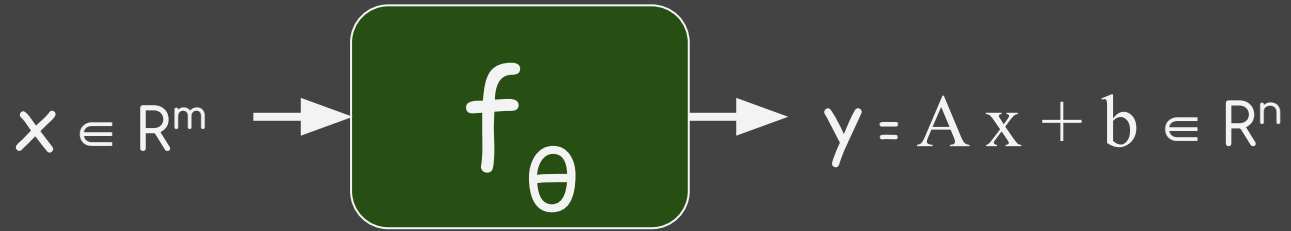
## Lecture 31

Parameter Learning, Cost Function, Overfitting,  
Cross-validation

# Example: Linear Regression



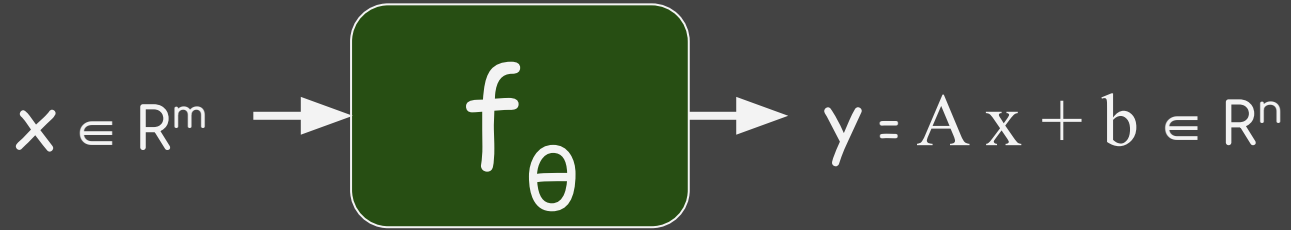
K. N. Toosi  
University of Technology



# Example: Linear Regression



K. N. Toosi  
University of Technology



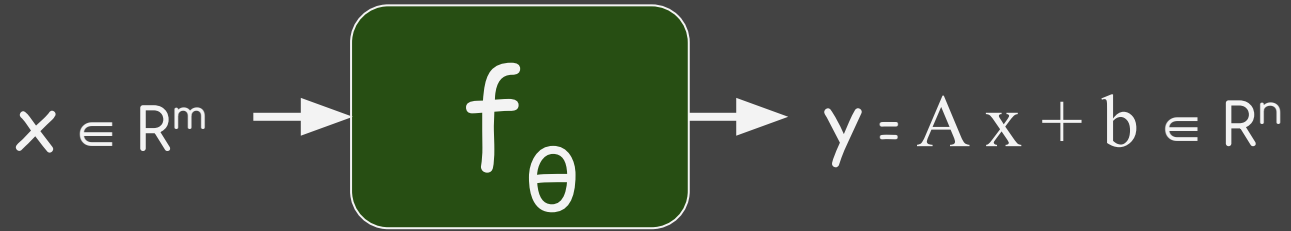
A: ? by ? matrix

b: ?-D vector

# Example: Linear Regression



K. N. Toosi  
University of Technology



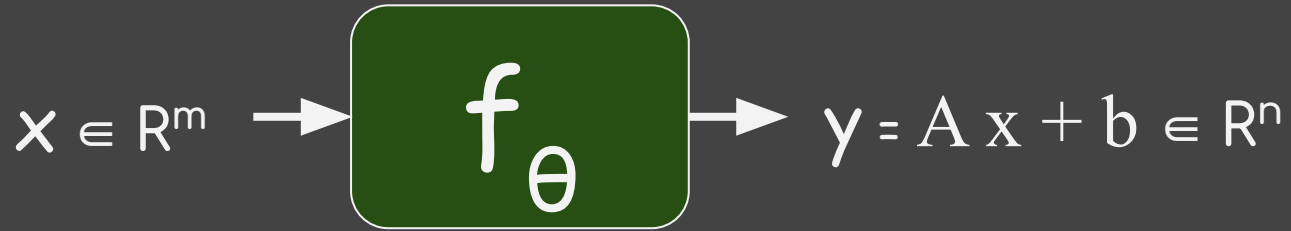
$\mathbf{A}$ :  $n$  by  $m$  matrix

$\mathbf{b}$ :  $n$ -D vector

# Example: Linear Regression



K. N. Toosi  
University of Technology



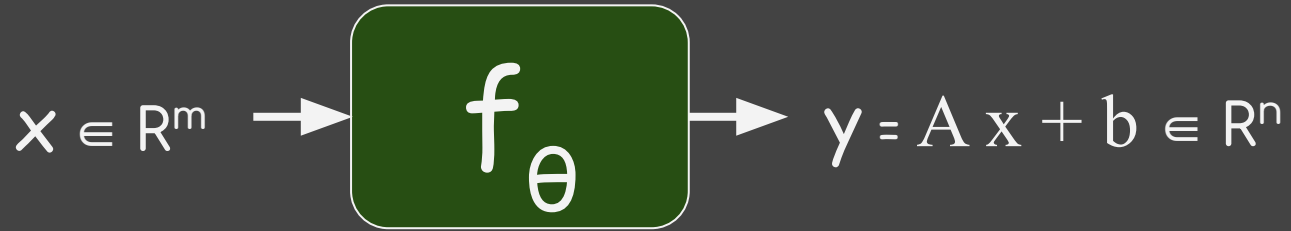
$$\mathbf{y} = \mathbf{f}(\theta, \mathbf{x})$$

$$\theta = ?$$

# Example: Linear Regression



K. N. Toosi  
University of Technology



$$\mathbf{y} = \mathbf{f}(\theta, \mathbf{x})$$

$$\theta = (\mathbf{A}, \mathbf{b})$$

# Example: Linear Regression



$$f(x) = Ax + b$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

is  $f$  linear?

$$f(\alpha x) = A(\alpha x) + b = \alpha Ax + b$$

$$\alpha f(x) = \alpha Ax + \alpha b$$

$\Rightarrow f$  is not linear in general ( $b \neq 0$ )

# Affine maps



K. N. Toosi  
University of Technology

$$\cancel{g(x+x_0)} \quad g(x) = f(x+x_0) - f(x_0) = A(x+x_0) + b - (Ax_0 + b) = Ax$$

$$g(x) = f(x) - f(0) = Ax$$
$$g(x) = f(x) - b = Ax$$

$$f(x) = Ax + b$$

affine function

$$f(x) : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\exists v \in \mathbb{R}^n$$

$$g(x) = f(x) - v$$

linear

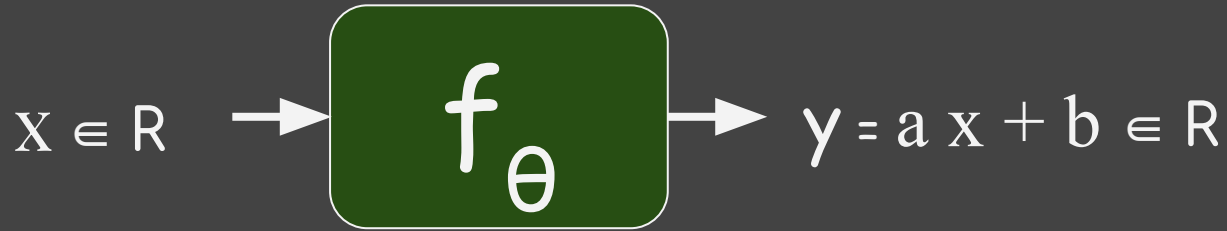
$$f(x-x_0) - f(x_0)$$



# Example: Linear Regression



K. N. Toosi  
University of Technology



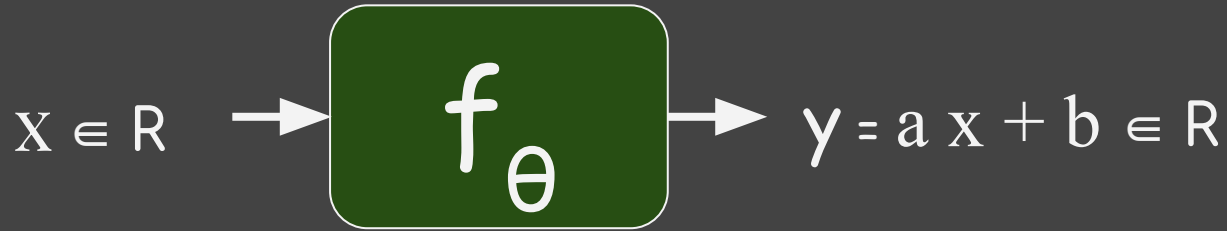
$$y = f(\theta, x)$$

$$\theta = ?$$

# Example: Linear Regression



K. N. Toosi  
University of Technology



$$y = f(\theta, x)$$

$$\theta = (a, b)$$

# Example: Linear Regression



K. N. Toosi  
University of Technology

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(x) = ax + b$$

$f$  is not linear

$$h(x, \theta) = h(\underline{x}, \underline{\theta}) = ax + b \quad \begin{cases} \text{is } h \text{ linear in } x? \text{ NO} \\ \text{is } h \text{ linear in } \underline{\theta}? \text{ YES} \end{cases}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \theta \in \mathbb{R}^2$$

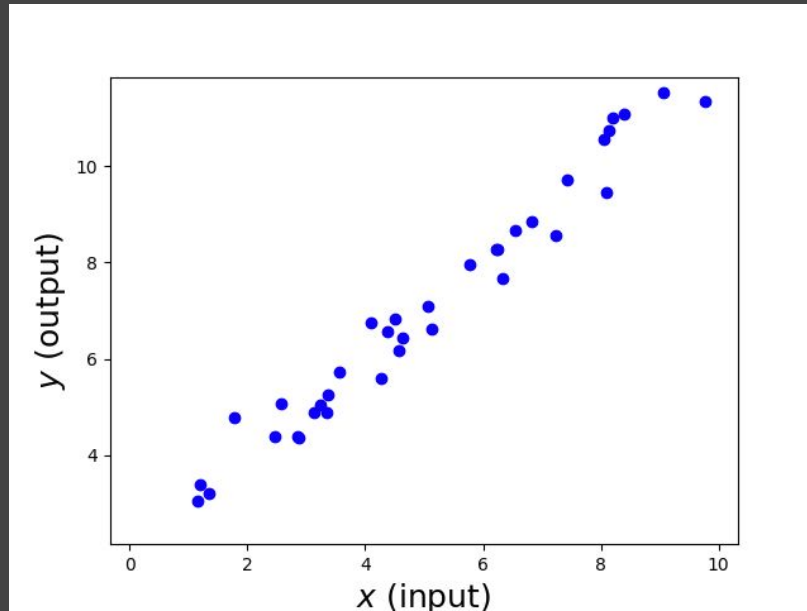
$$h(x, \theta) = ax + b = [x \ 1] \begin{bmatrix} a \\ b \end{bmatrix} = v^T \begin{bmatrix} a \\ b \end{bmatrix} = v^T \theta$$

# Linear Regression: Cost function



K. N. Toosi  
University of Technology

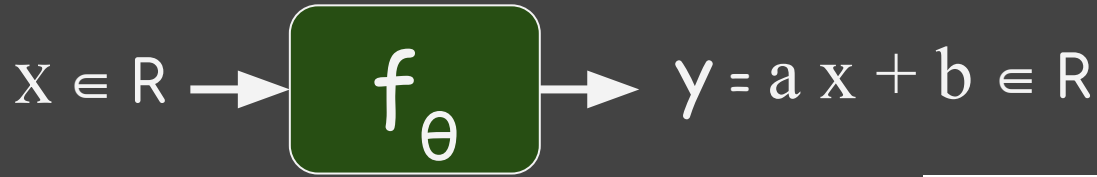
Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



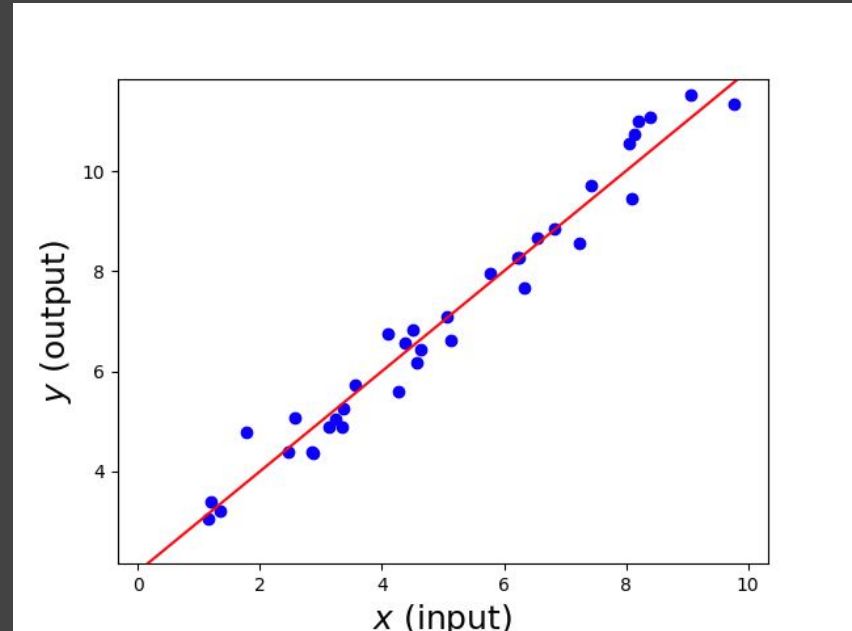
# Linear Regression: Cost function



K. N. Toosi  
University of Technology



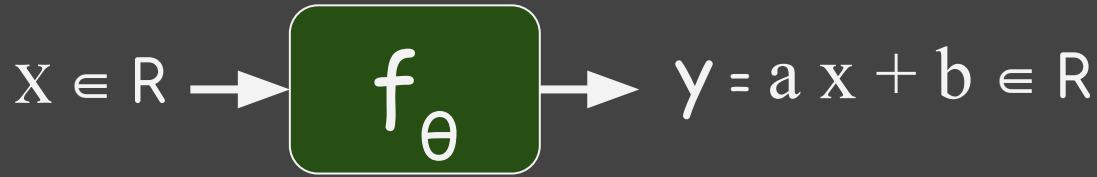
Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



# Linear Regression: Cost function



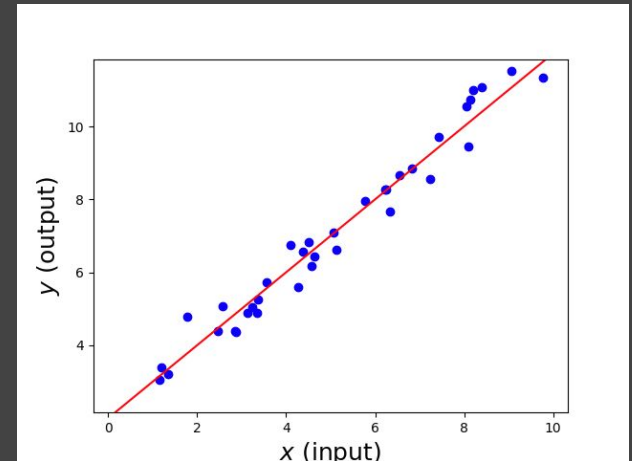
K. N. Toosi  
University of Technology



Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

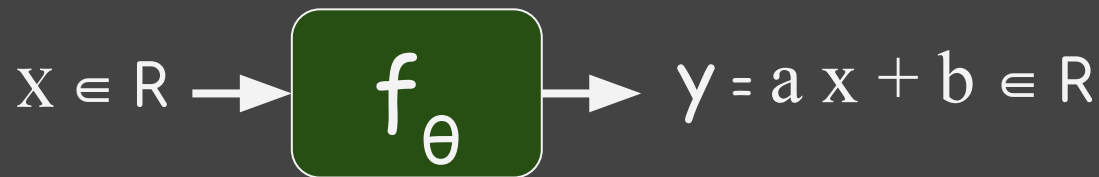
cost function:

$$C(\theta) = \sum_{i=1..N} d(f(\theta, x_i), y_i)$$





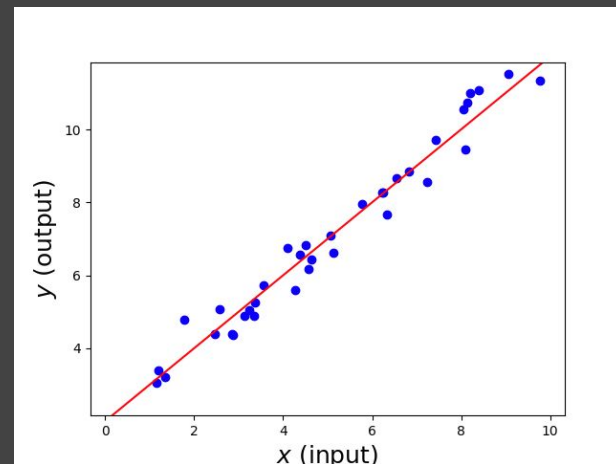
# Linear Regression: Cost function



Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

cost function:

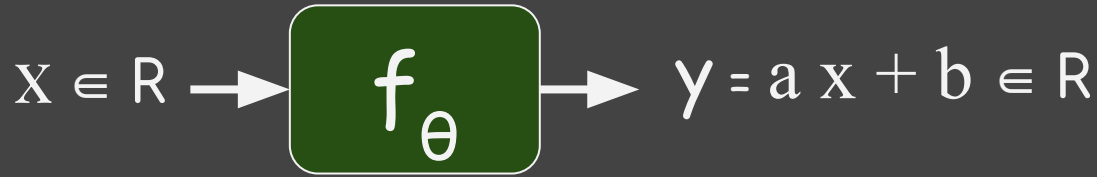
$$C(a,b) = \sum_{i=1..n} d( f(a,b, x_i), y_i )$$



# Linear Regression: Cost function



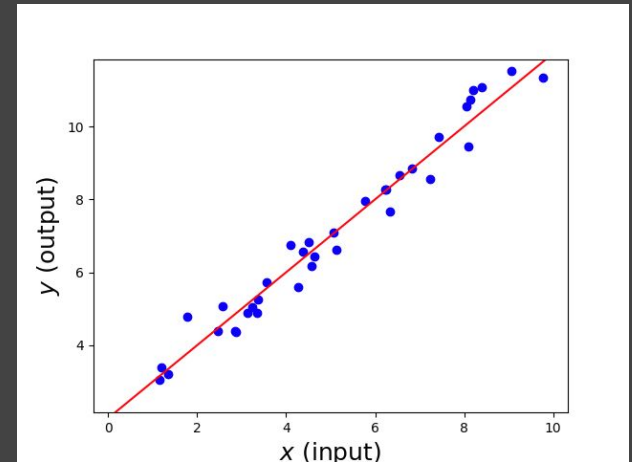
K. N. Toosi  
University of Technology



Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

cost function:

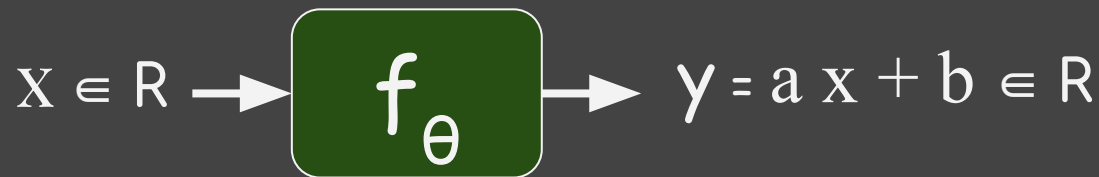
$$\begin{aligned} C(a,b) &= \sum_{i=1..n} d(f(a,b, x_i), y_i) \\ &= \sum_{i=1..n} d(a x_i + b, y_i) \end{aligned}$$







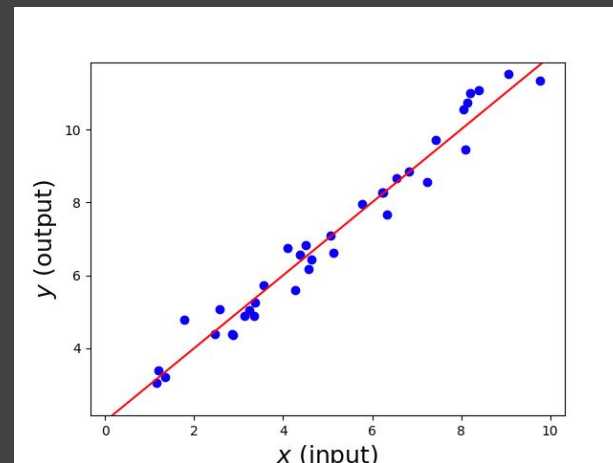
# Linear Regression: Cost function



Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

cost function (sum of squared errors):

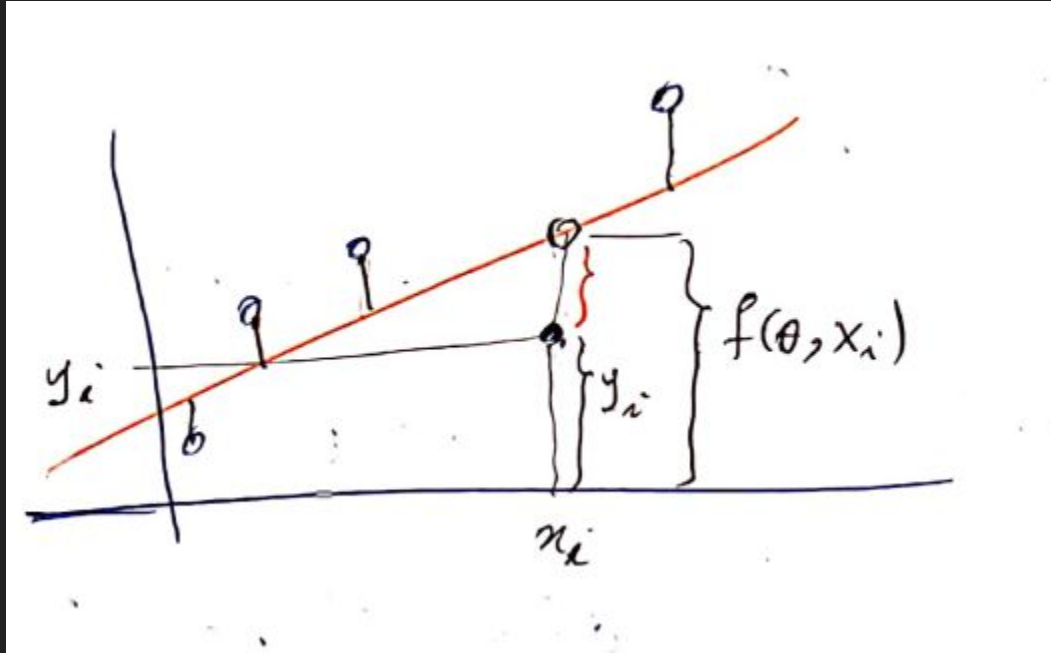
$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$



# The cost function

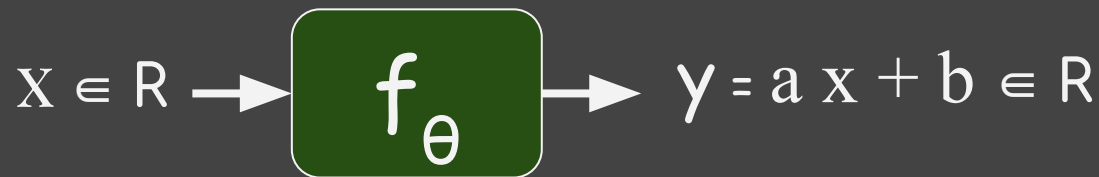


K. N. Toosi  
University of Technology





# Linear Regression: Cost function

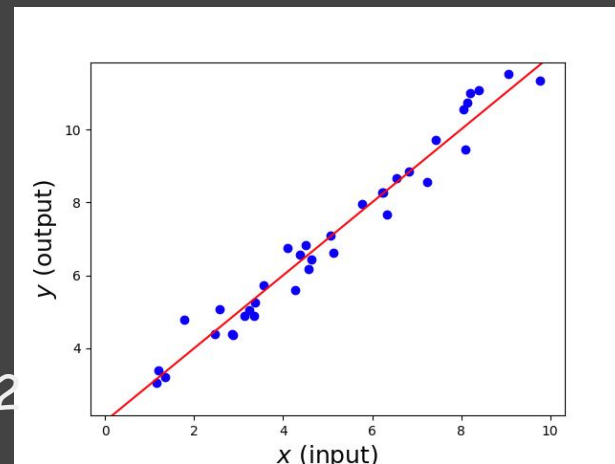


Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

cost function (sum of squared errors):

$$C(a, b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a, b} \sum_{i=1..n} (a x_i + b - y_i)^2$$



# Linear Regression: Cost function



$$f(x, \theta) = f(x, \begin{bmatrix} a \\ b \end{bmatrix}) = ax + b = [x \ 1] \begin{bmatrix} a \\ b \end{bmatrix} = [x \ 1] \theta$$

linear in  $x$ ? Not in general

linear in  $\begin{bmatrix} a \\ b \end{bmatrix}$ ? YES

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$C(\theta) = C(a, b) = \sum_{i=1}^N d(f(x_i, \theta), y_i) = \sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$

$$\begin{aligned} C(a, b) &= \sum_{i=1}^N (ax_i + b - y_i)^2 = \sum_{i=1}^N ([x_i \ 1] \begin{bmatrix} a \\ b \end{bmatrix} - y_i)^2 \\ &= \left\| \begin{bmatrix} ax_1 + b - y_1 \\ ax_2 + b - y_2 \\ \vdots \\ ax_N + b - y_N \end{bmatrix} \right\|^2 = \left\| \underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \vdots & \vdots \\ x_N & 1 \end{bmatrix}}_{N \times 2} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{2 \times 1} - \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{N \times 1} \right\|^2 \end{aligned}$$

# Example: Linear Regression



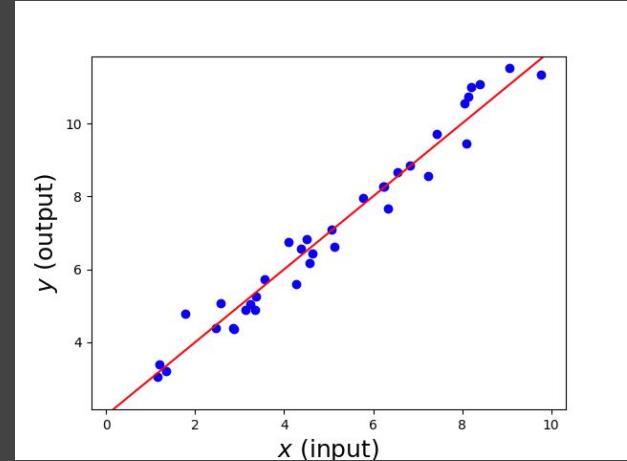
K. N. Toosi  
University of Technology

cost function (sum of squared errors):

$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

How to find  $a^*, b^*$ ?



# Solution 1: Least squares



training data  $(x_1, y_1), (x_2, y_2) \dots, (x_N, y_N)$

$$C(\theta) = \sum_{i=1}^N d(f(x_i, \theta), y_i)$$

$$= \sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$

$$= \sum_{i=1}^N (f(x_i, a, b) - y_i)^2$$

$$C(a, b) = \sum_{i=1}^N (ax_i + b - y_i)^2 \quad \text{cost function}$$

$$= \sum_{i=1}^N \left( \begin{bmatrix} x_i & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} - y_i \right)^2$$

# Solution 1: Least squares



$$C(a, b) = \sum_{i=1}^N (ax_i + b - y_i)^2 \quad \text{cost function}$$

$$= \sum_{i=1}^N \left( [x_i \ 1] \begin{bmatrix} a \\ b \end{bmatrix} - y_i \right)^2$$

$$= \begin{bmatrix} ax_1 + b - y_1 \\ ax_2 + b - y_2 \\ ax_3 + b - y_3 \\ \vdots \\ ax_N + b - y_N \end{bmatrix}^2 = \begin{bmatrix} x_{1,1} \\ x_{2,1} \\ x_{3,1} \\ \vdots \\ x_{N,1} \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix}^2 = \|A\theta - b\|^2$$

$\underbrace{\begin{bmatrix} x_{1,1} \\ x_{2,1} \\ x_{3,1} \\ \vdots \\ x_{N,1} \end{bmatrix}}_{N \times 2} \quad \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{2 \times 1} \quad \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix}}_{N \times 1}$

# Solution 1: Least squares



$$\theta^* = \begin{bmatrix} a^* \\ b^* \end{bmatrix} = \operatorname{argmin} \|A\theta + b\|^2$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \theta = (A^T A)^{-1} A^T b$$

$$= \left( \begin{bmatrix} x_1 & x_2 & \dots & x_N \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_N & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} x_1 & x_2 & \dots & x_N \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & N \end{bmatrix}^{-1} \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

closed-form  
solution



# Example: Linear Regression



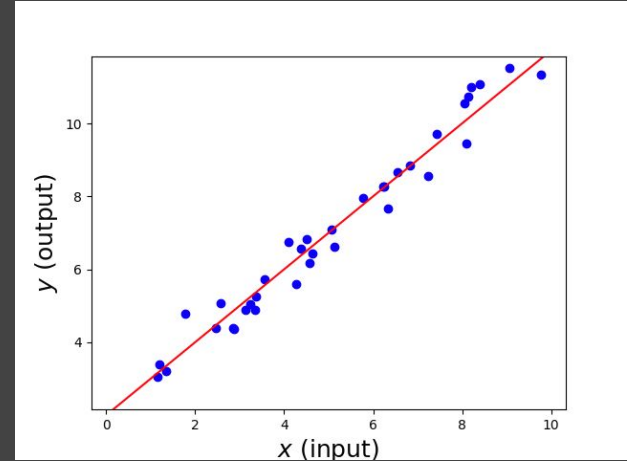
K. N. Toosi  
University of Technology

cost function (sum of squared errors):

$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

How to find  $a^*, b^*$ ?



# Solution 2: partial derivatives



K. N. Toosi  
University of Technology

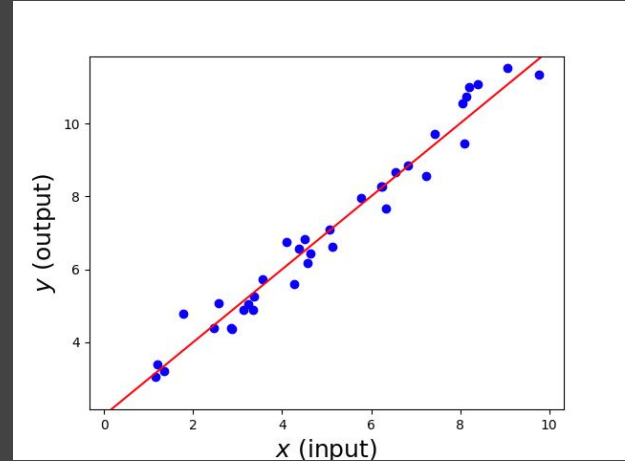
cost function (sum of squared errors):

$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$\partial C(a,b) / \partial a = 0$$

$$\partial C(a,b) / \partial b = 0$$



# Solution 2: partial derivatives

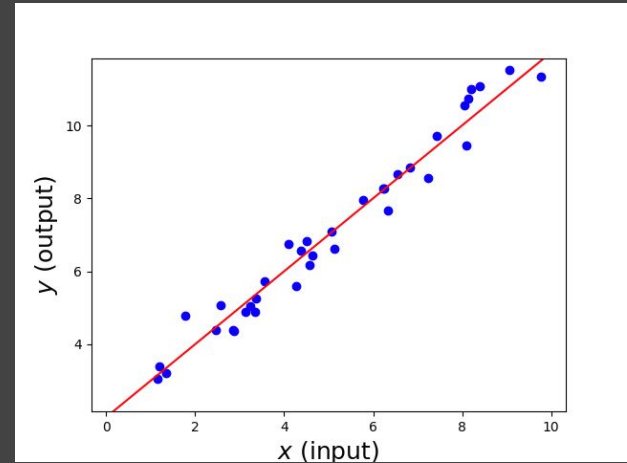


K. N. Toosi  
University of Technology

cost function (sum of squared errors):

$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$



$$\frac{\partial C(a,b)}{\partial a} = 2 \sum_{i=1..n} x_i (a x_i + b - y_i) = 0$$

$$\frac{\partial C(a,b)}{\partial b} = 2 \sum_{i=1..n} (a x_i + b - y_i) = 0$$

# Solution 2: partial derivatives



K. N. Toosi  
University of Technology

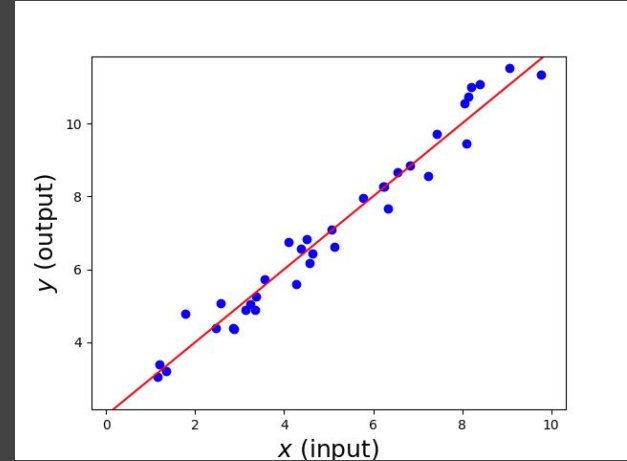
cost function (sum of squared errors):

$$C(a,b) = \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$\sum_{i=1..n} x_i (a x_i + b - y_i) = 0$$

$$\sum_{i=1..n} (a x_i + b - y_i) = 0$$



# Solution 2: partial derivatives



K. N. Toosi  
University of Technology

$$a^*, b^* = \operatorname{argmin}_{a, b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$\sum_{i=1..n} x_i (a x_i + b - y_i) = a \sum_{i=1..n} x_i^2 + b \sum_{i=1..n} x_i - \sum_{i=1..n} x_i y_i = 0$$

$$\sum_{i=1..n} (a x_i + b - y_i) = a \sum_{i=1..n} x_i + b n - \sum_{i=1..n} y_i = 0$$



## Solution 2: partial derivatives

$$a^*, b^* = \operatorname{argmin}_{a,b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$\left( \sum_{i=1..n} x_i^2 \right) a + \left( \sum_{i=1..n} x_i \right) b = \sum_{i=1..n} x_i y_i$$

$$\left( \sum_{i=1..n} x_i \right) a + n b = \sum_{i=1..n} y_i$$



## Solution 2: partial derivatives

$$a^*, b^* = \operatorname{argmin}_{a, b} \sum_{i=1..n} (a x_i + b - y_i)^2$$

$$\left( \sum_{i=1..n} x_i^2 \right) a + \left( \sum_{i=1..n} x_i \right) b = \sum_{i=1..n} x_i y_i$$

$$\left( \sum_{i=1..n} x_i \right) a + n b = \sum_{i=1..n} y_i$$

$a^*, b^* \leftarrow$  solve system of linear equations

# Solution 2: partial derivatives



$$\begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & N \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & N \end{bmatrix}^{-1} \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

closed-form solution

x

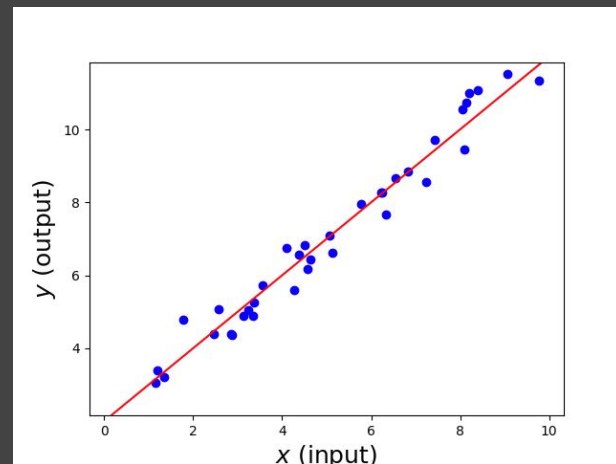




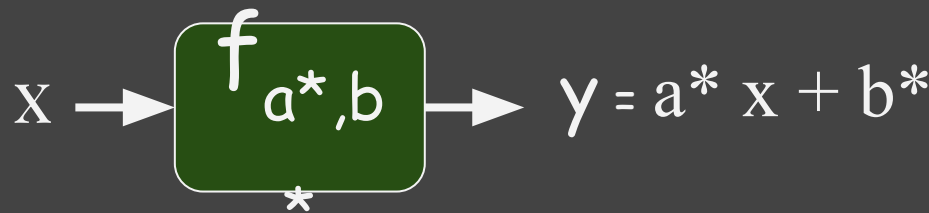
# Example: Linear Regression

$$a^*, b^* = \operatorname{argmin}_{a, b} \sum_{i=1..n} (a x^i + b - y^i)^2$$

$a^*, b^* \Leftarrow$  solve system of linear equations



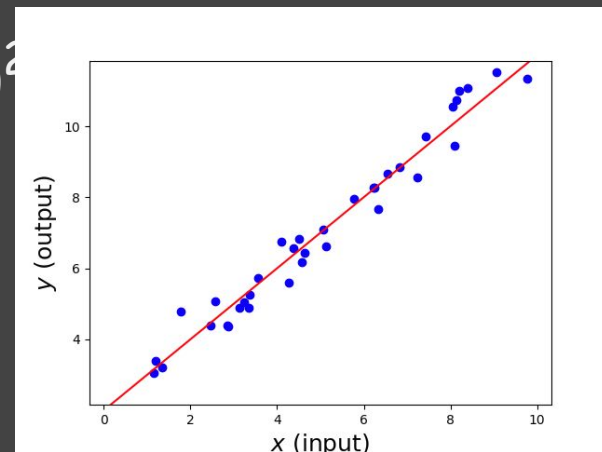
$$y = a^* x + b^*$$





# Evaluation

- Find good parameters  $\theta$ 
  - $\theta^* = \operatorname{argmin}_{\theta} \sum_{i=1..n} (f(\theta, x_i) - y_i)^2$
  - another method
- How good  $\theta^*$  is?
- How well the regressor works?



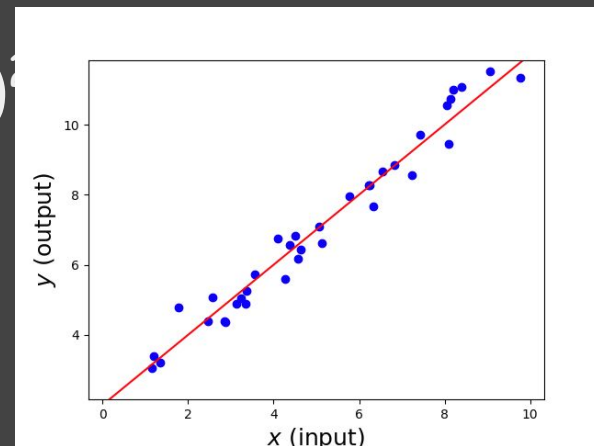
$$y = a^* x + b^*$$





# Evaluation

- Find good parameters  $\theta$ 
  - $\theta^* = \operatorname{argmin}_{\theta} \sum_{i=1..N} (f(\theta, x_i) - y_i)$
  - another method
- How good  $\theta^*$  is?
- How well the regressor works?
- Given Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



$$y = a^* x + b^*$$

$$\text{Error} = C(\theta^*) = \sum_{i=1..N} (f(\theta^*, x_i) - y_i)^2$$

# Learning from data



K. N. Toosi  
University of Technology



- Parameter Learning:
  - A collection of input-output pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ,
  - choose  $\theta$  such that  $y = f(\theta, x)$  is a reasonable output
    - for training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
    - for unseen data

# Learning from data



K. N. Toosi  
University of Technology



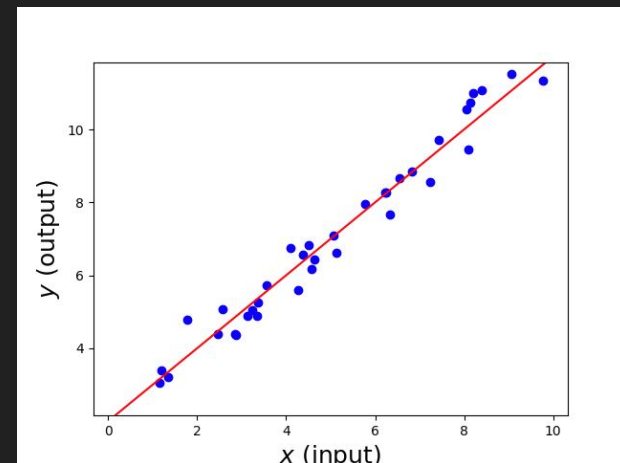
- Parameter Learning:
  - A collection of input-output pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ,
  - choose  $\theta$  such that  $y = f(\theta, x)$  is a reasonable output
    - for training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
    - for unseen data
  - **Generalization:** How well the model works on unseen data



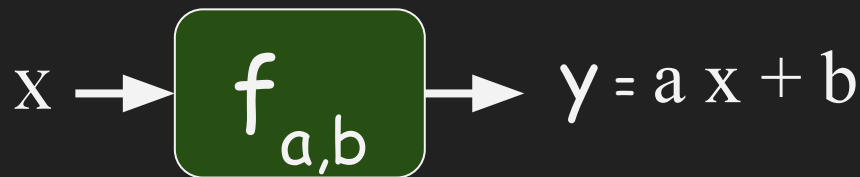
# Example: Linear Regression

$$a^*, b^* = \operatorname{argmin}_{a, b} \sum_{i=1..n} (a x^i + b - y^i)^2$$

$a^*, b^* \Leftarrow$  solve system of linear equations



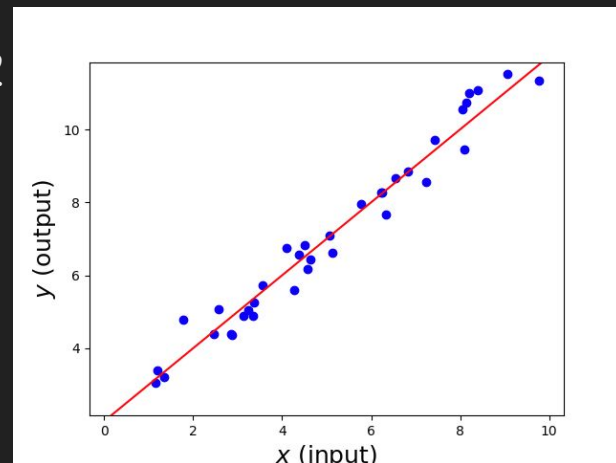
$$y = a^* x + b^*$$





# Evaluation

- Find good parameters  $\theta$ 
  - $\theta^* = \operatorname{argmin}_{\theta} \sum_{i=1..n} (f(\theta, x_i) - y_i)^2$
  - another method
- How good  $\theta^*$  is?
- How well the regressor works?



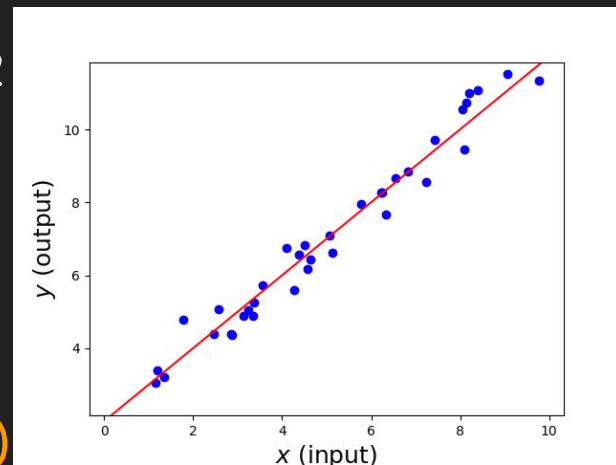
$$y = a^* x + b^*$$





# Evaluation

- Find good parameters  $\theta$ 
  - $\theta^* = \operatorname{argmin}_{\theta} \sum_{i=1..n} (f(\theta, x_i) - y_i)^2$
  - another method
- How good  $\theta^*$  is?
- How well the regressor works?
- Given Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



$$y = a^* x + b^*$$

$$\text{Error} = C(\theta^*) = \sum_{i=1..N} (f(\theta^*, x_i) - y_i)^2$$



# Learning from data



K. N. Toosi  
University of Technology

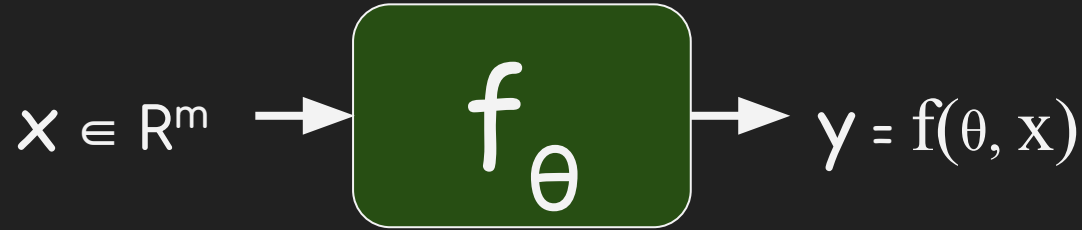


- Parameter Learning:
  - A collection of input-output pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ,
  - choose  $\theta$  such that  $y = f(\theta, x)$  is a reasonable output
    - for training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
    - for unseen data

# Learning from data



K. N. Toosi  
University of Technology



- Parameter Learning:
  - A collection of input-output pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ,
  - choose  $\theta$  such that  $y = f(\theta, x)$  is a reasonable output
    - for training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
    - for unseen data
  - **Generalization:** How well the model works on unseen data

# Example: Polynomial Regression

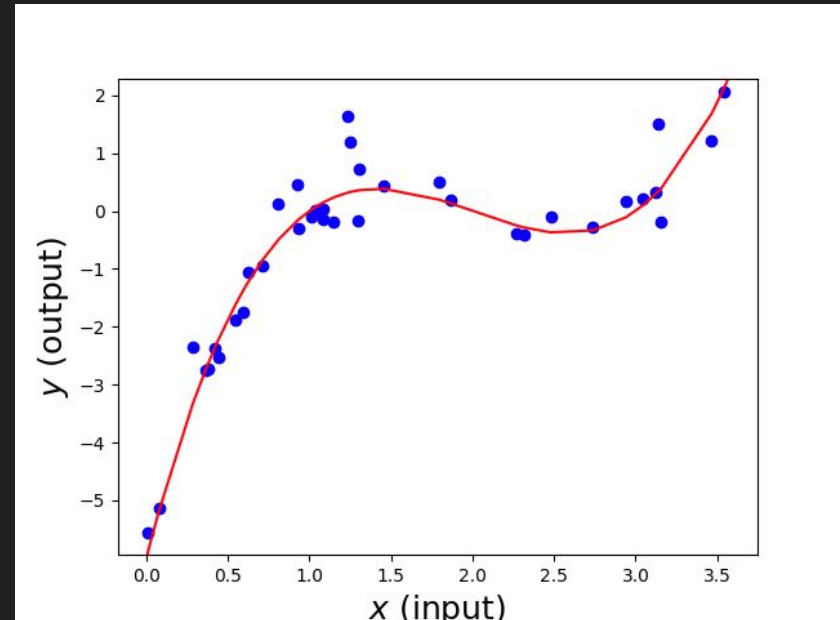


K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

$$\theta = (a_p, a_{p-1}, \dots, a_1, a_0)$$



# Example: Polynomial Regression



Polynomial Regression

$$f(x, \theta) = a_p x^p + a_{p-1} x^{p-1} + \dots + a_2 x^2 + a_1 x + a_0$$

$$\theta = [a_p, a_{p-1}, \dots, a_2, a_1, a_0]^T \in \mathbb{R}^{p+1}$$

$$f: \mathbb{R} \times \mathbb{R}^{p+1} \rightarrow \mathbb{R}$$

$$f(x, \theta) = [x^p \quad x^{p-1} \quad \dots \quad x^2 \quad x \quad 1]$$

$\underbrace{\hspace{10em}}_{V^T}$

$$\begin{bmatrix} a_p \\ a_{p-1} \\ \vdots \\ a_2 \\ a_1 \\ a_0 \\ \theta \end{bmatrix} = V^T \theta$$

$$C(\theta) = \sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$

# Example: Polynomial Regression



$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(\underline{x}) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n \in \mathbb{R}$$

$$\theta = [a_0, a_1, \dots, a_n]^T \in \mathbb{R}^{n+1}$$

$f(x, \theta)$ : linear in  $x$ ? **No**  
linear in  $\theta$ ? **YES**  $\Rightarrow f(x) = M\theta$   
( $x^{(n+1)}$ )

$$f(x, \theta) = [1, x, x^2, \dots, x^n] \theta = [1, x, x^2, \dots, x^n] \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix}$$

# Example: Polynomial Regression



$$f(x, \theta) = [1, x, x^2, \dots, x^n] \theta = [1, x, x^2, \dots, x^n] \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix}$$

→ training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

$$C(\theta) = \sum_{i=1}^N (f(\theta, x_i) - y_i)^2 = \left\| \begin{bmatrix} 1, x_1, x_1^2, \dots, x_1^n \\ 1, x_2, x_2^2, \dots, x_2^n \\ \vdots \\ 1, x_N, x_N^2, \dots, x_N^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \right\|^2$$

$\|A \cdot \theta - Y\|^2$

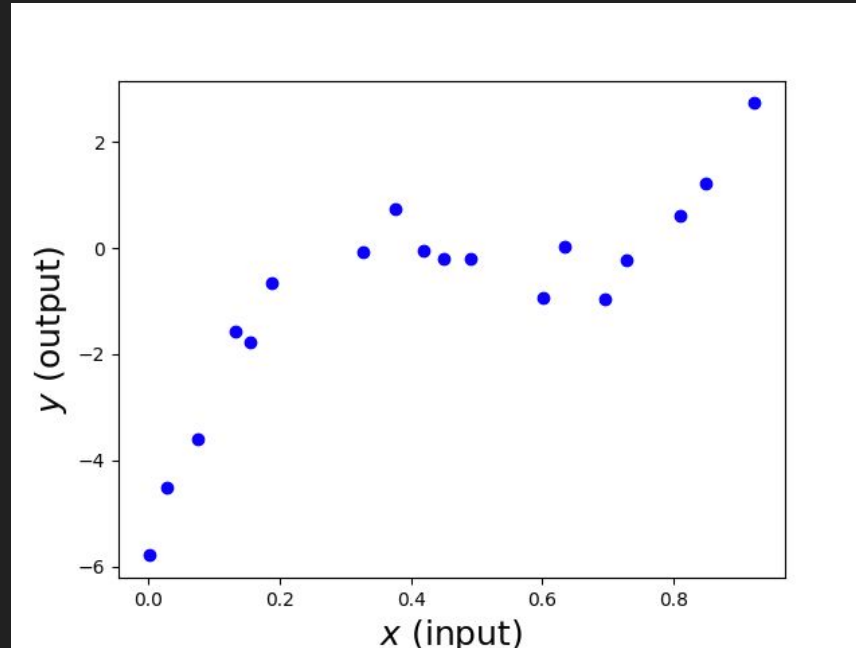
$$\theta^* = (A^T A)^{-1} A^T Y$$

# Vary polynomial degree $p$



K. N. Toosi  
University of Technology

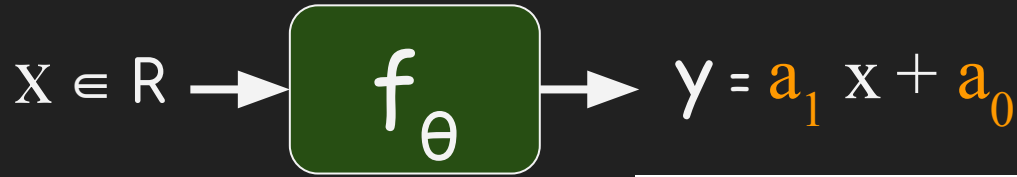
$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$



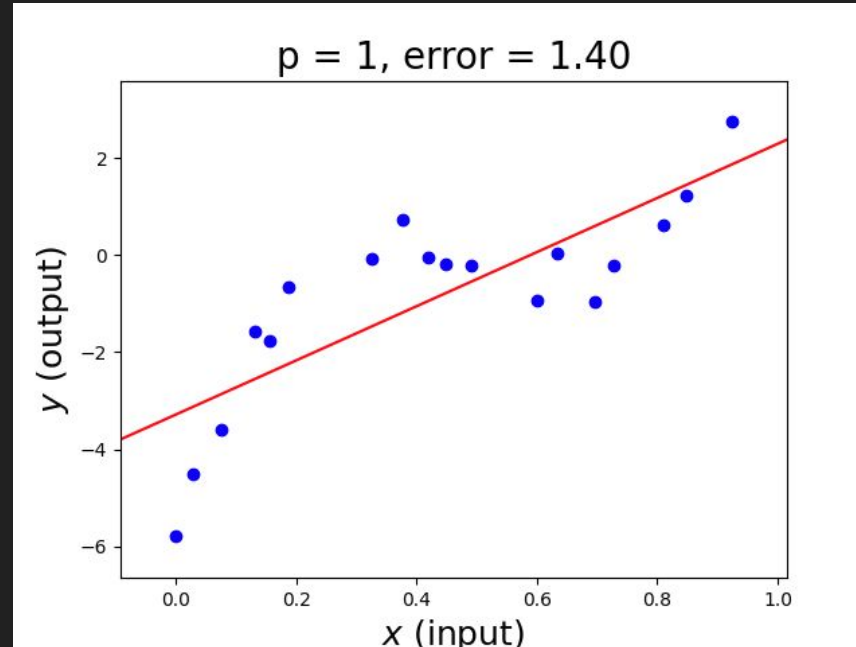
# polynomial degree $p = 1$



K. N. Toosi  
University of Technology

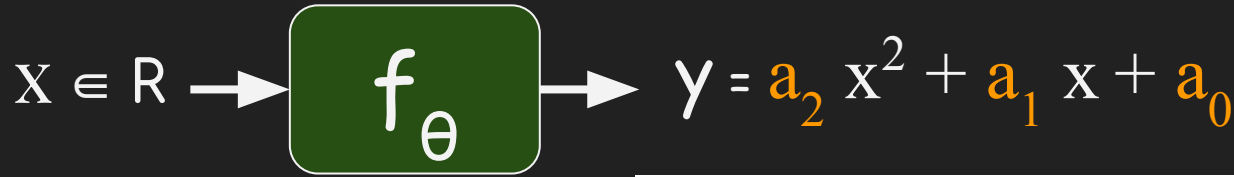


$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$

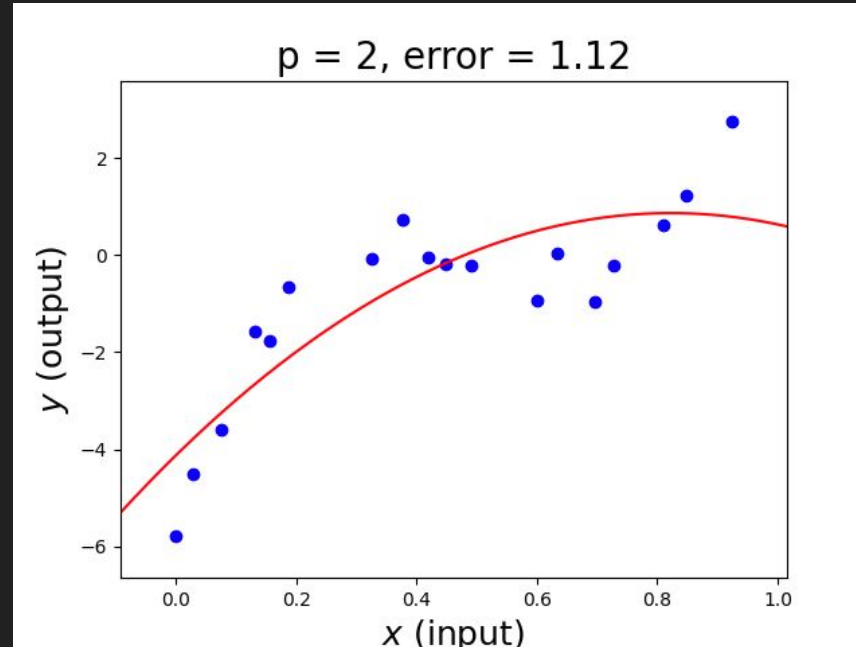




polynomial degree  $p = 2$



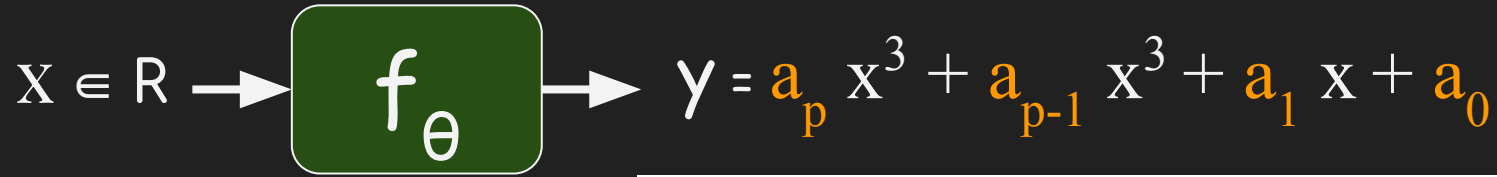
$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$



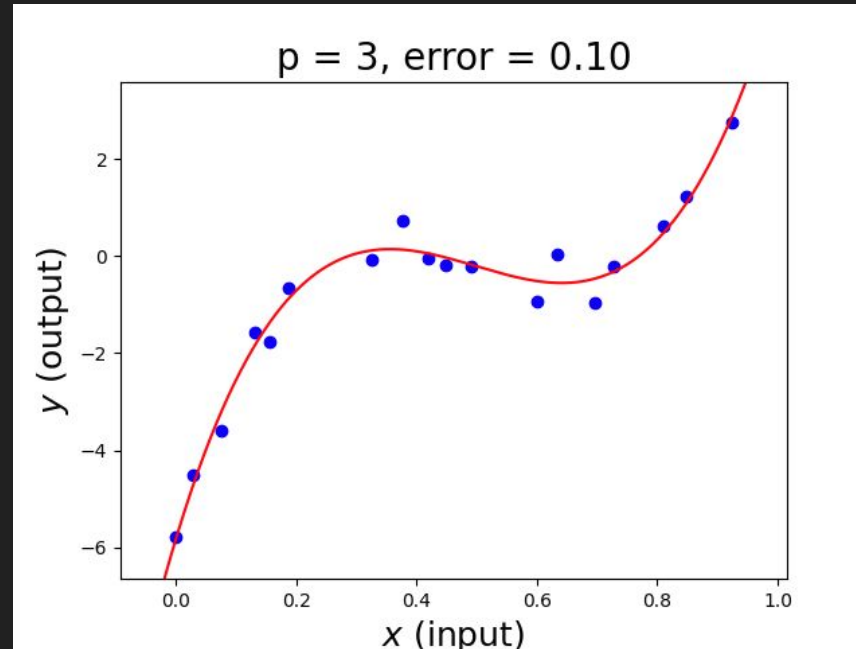
# polynomial degree $p = 3$



K. N. Toosi  
University of Technology



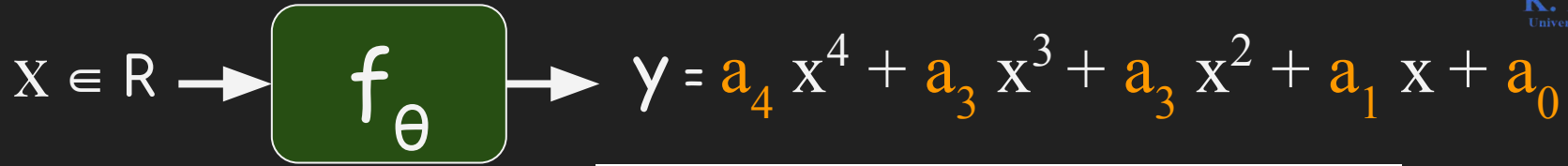
$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$



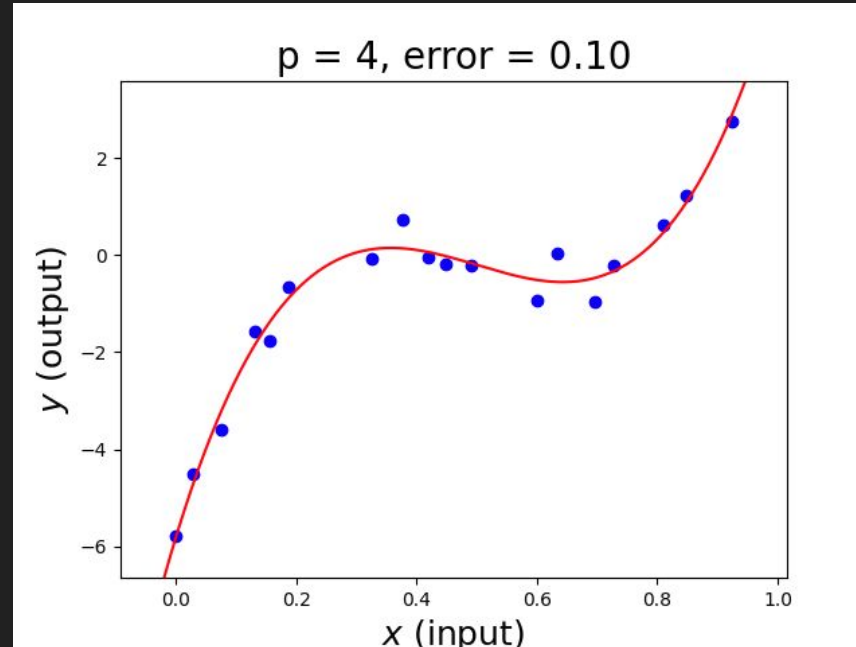
polynomial degree  $p = 4$



K. N. Toosi  
University of Technology



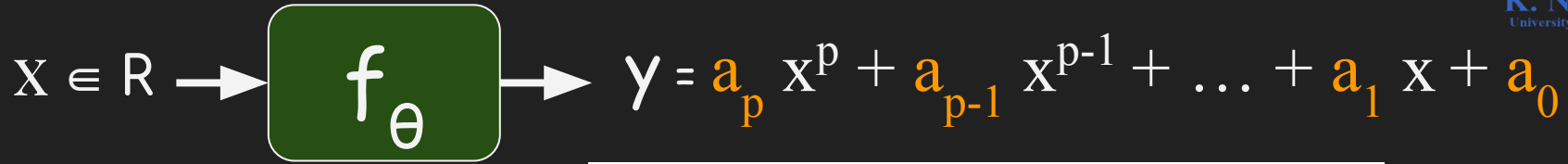
$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$



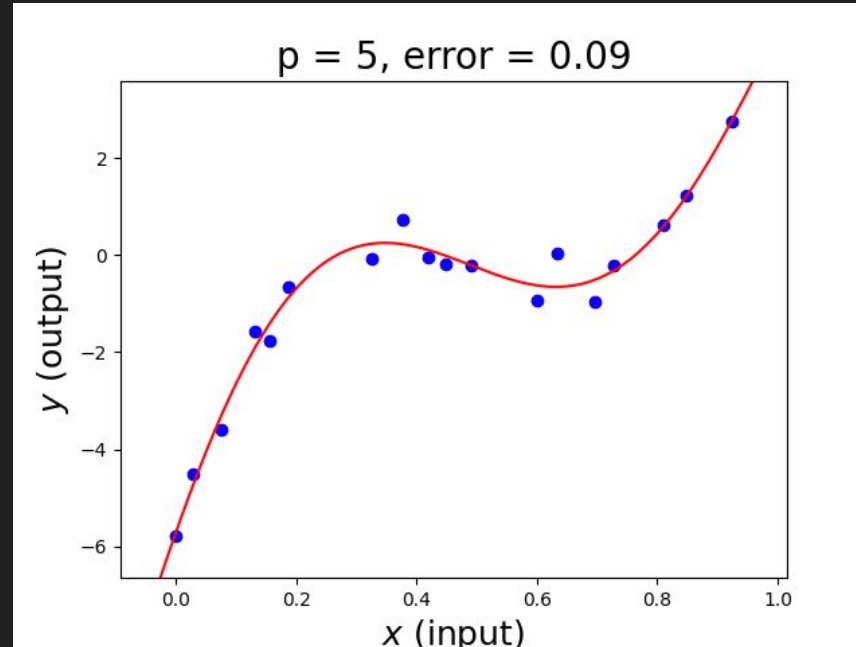
# polynomial degree $p = 5$



K. N. Toosi  
University of Technology



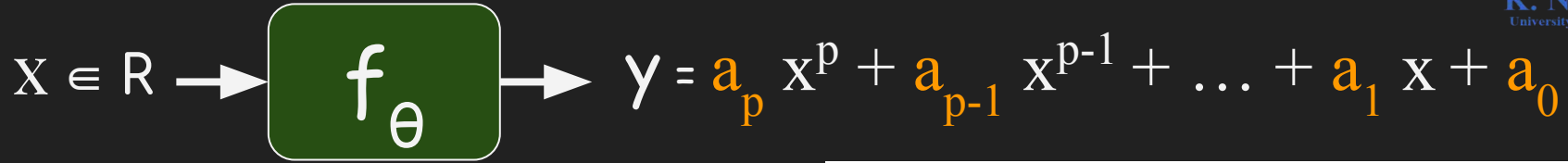
$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$



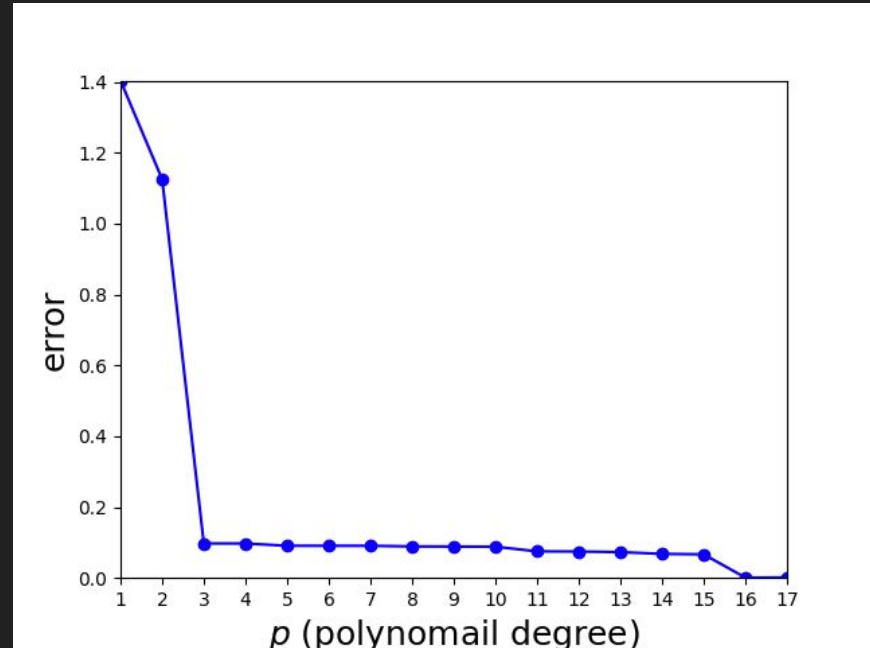
polynomial degree  $p = 5$



K. N. Toosi  
University of Technology



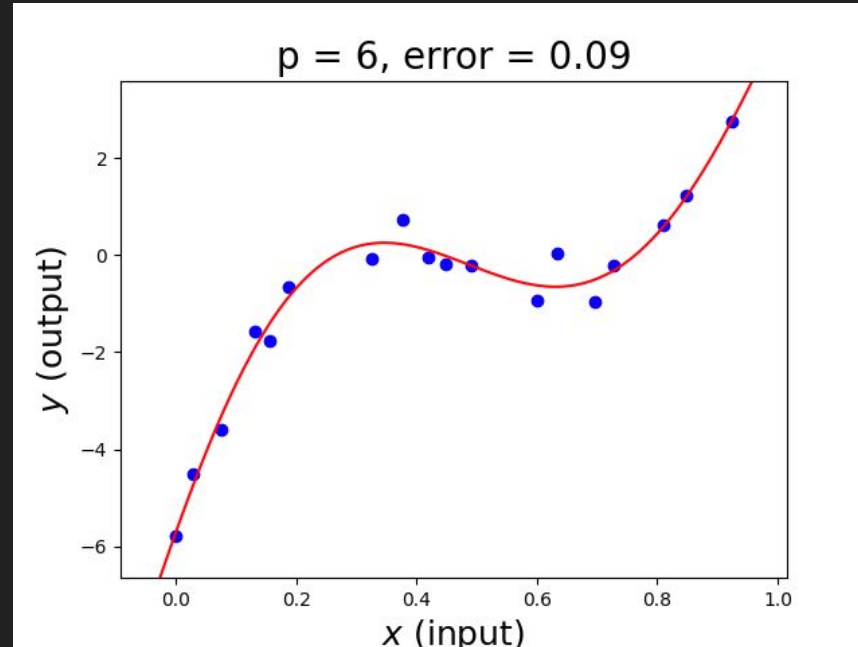
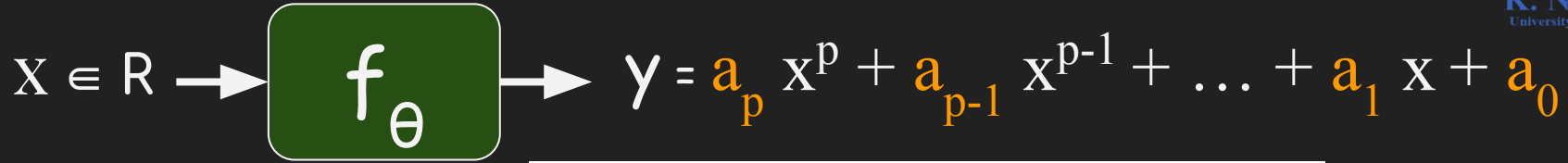
$$\text{Error} = \sum_{i=1..n} (f(\theta^*, x_i) - y_i)^2$$



# polynomial degree $p = 6$



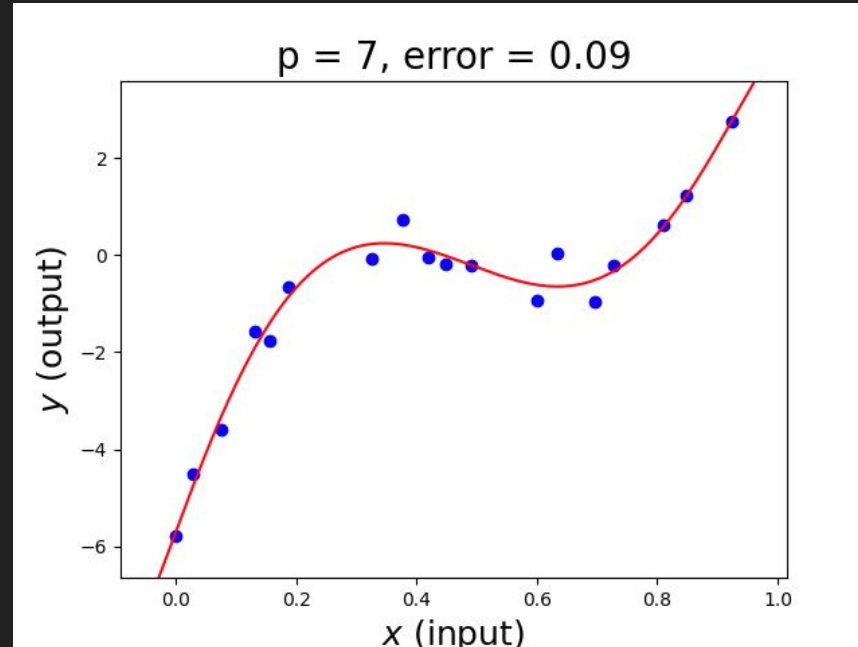
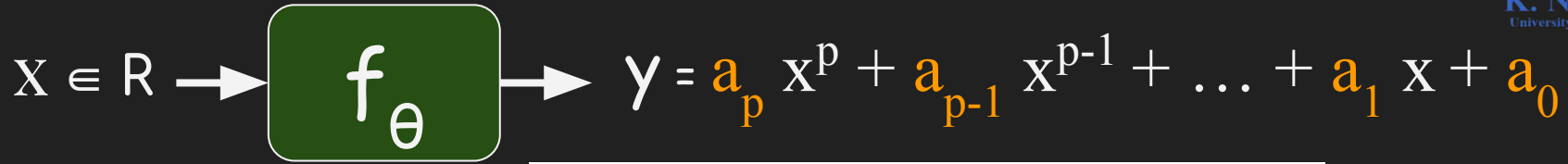
K. N. Toosi  
University of Technology



polynomial degree  $p = 7$




K. N. Toosi  
University of Technology

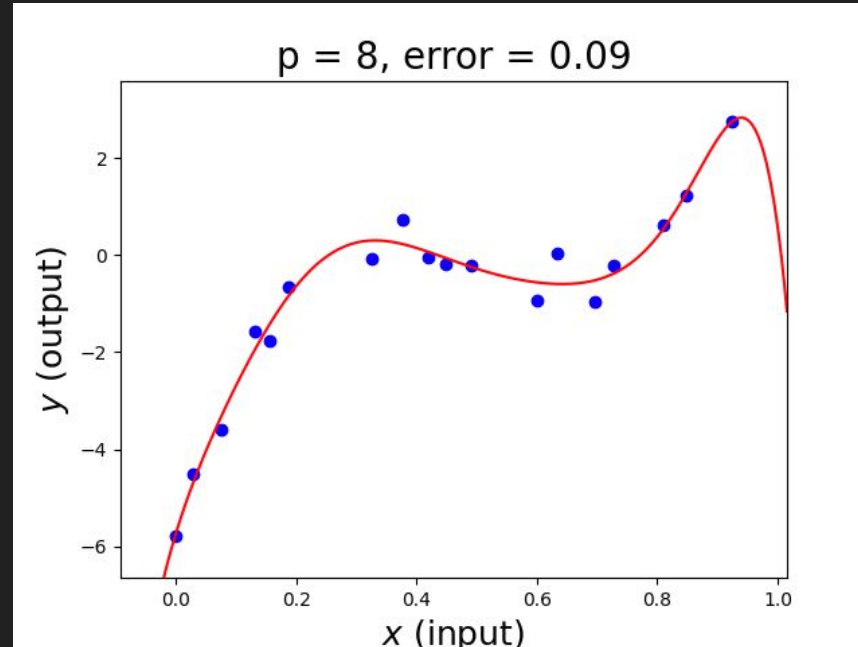


polynomial degree  $p = 8$



K. N. Toosi  
University of Technology

$X \in \mathbb{R} \rightarrow$    $\rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$



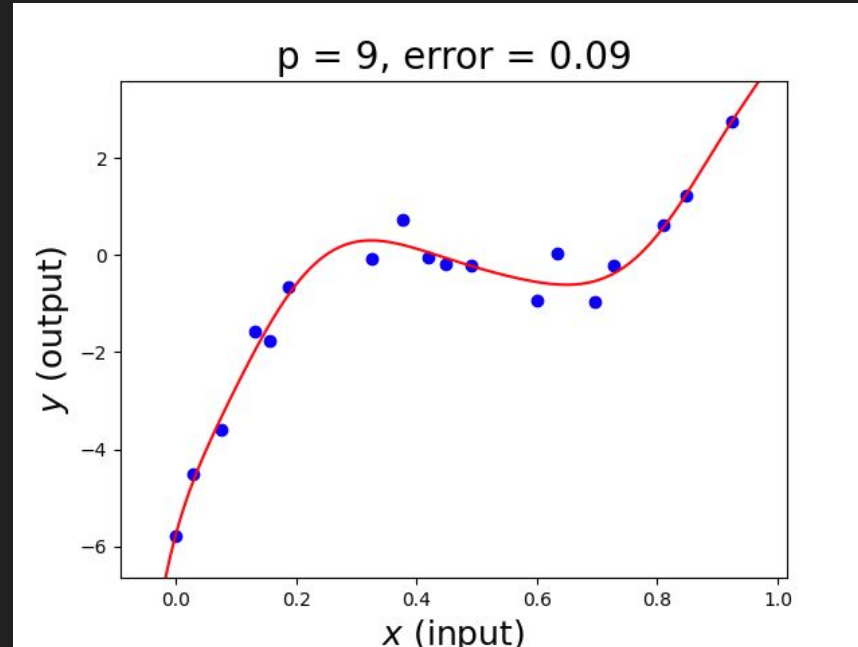


# polynomial degree $p = 9$



K. N. Toosi  
University of Technology


$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

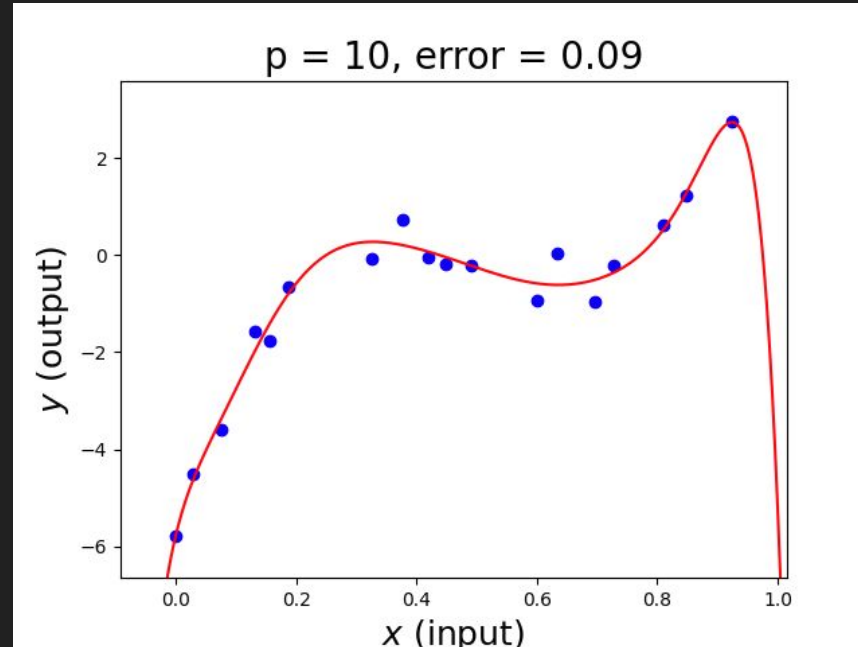


polynomial degree  $p = 10$



K. N. Toosi  
University of Technology

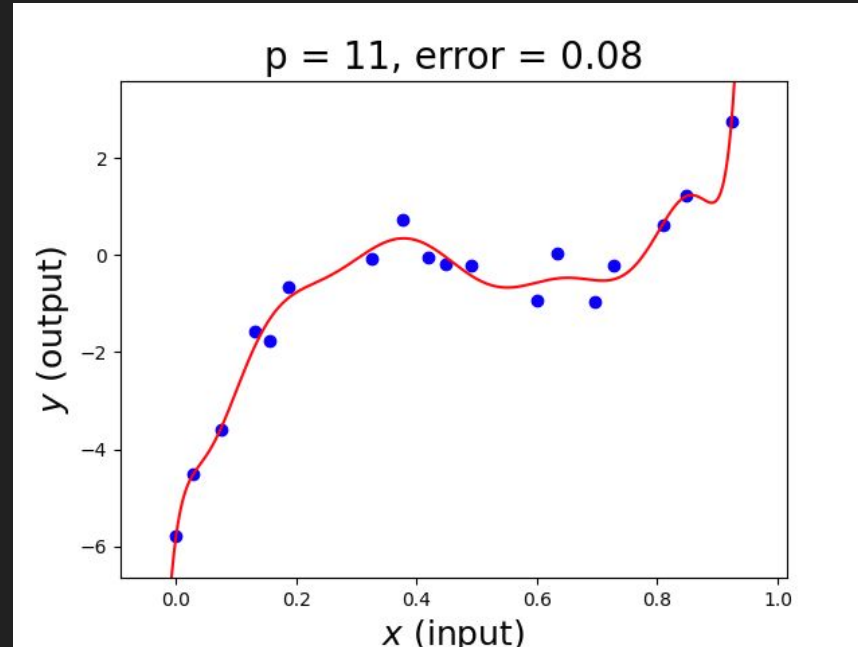
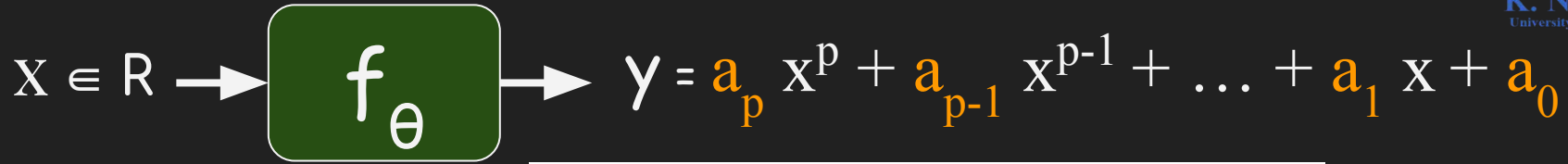
$X \in \mathbb{R} \rightarrow$    $\rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$



polynomial degree  $p = 11$



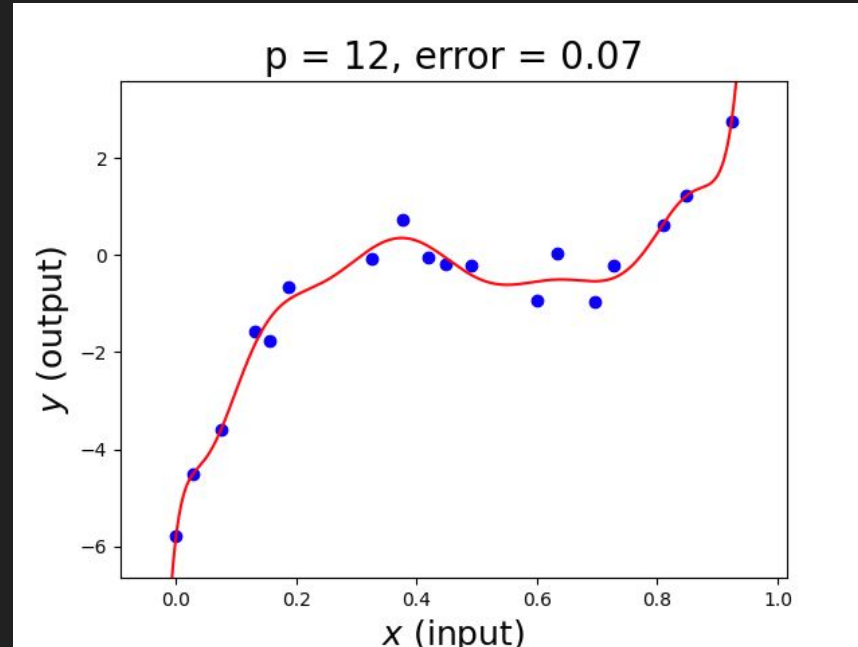
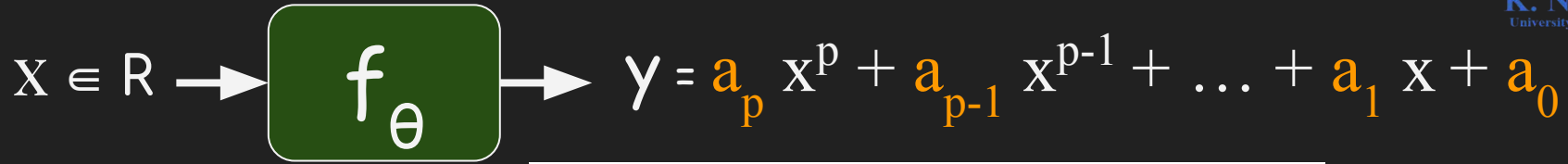
K. N. Toosi  
University of Technology



polynomial degree  $p = 12$




K. N. Toosi  
University of Technology

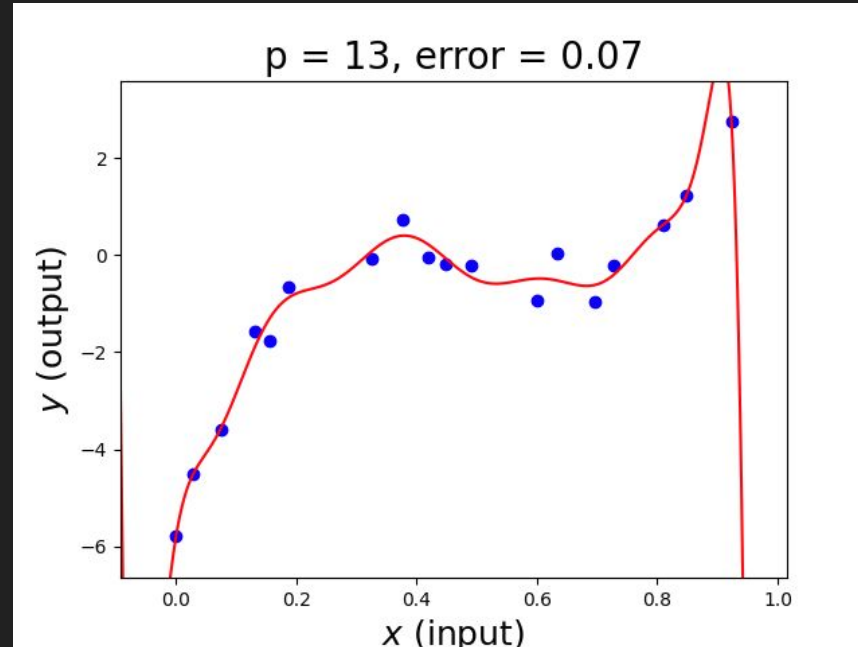


polynomial degree  $p = 13$



K. N. Toosi  
University of Technology

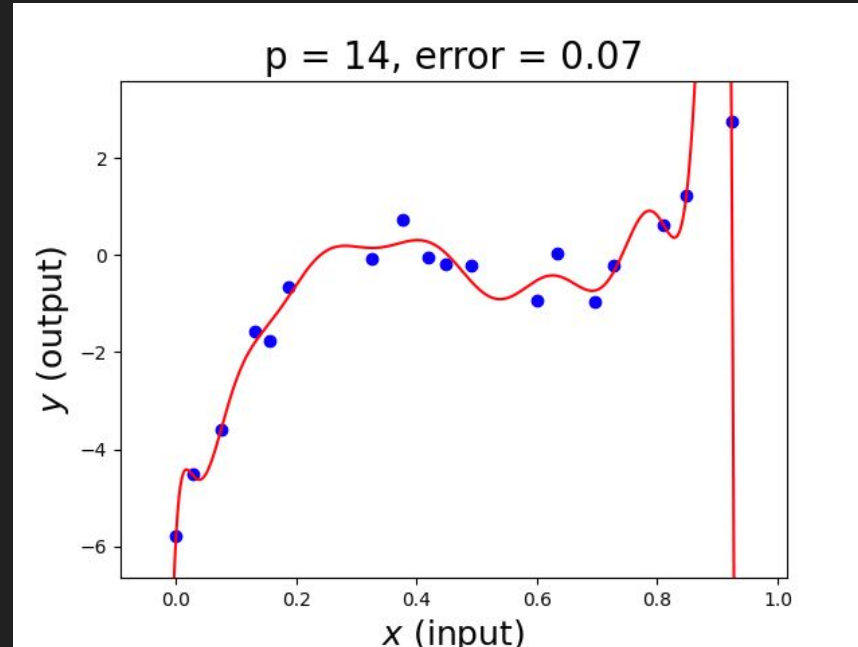
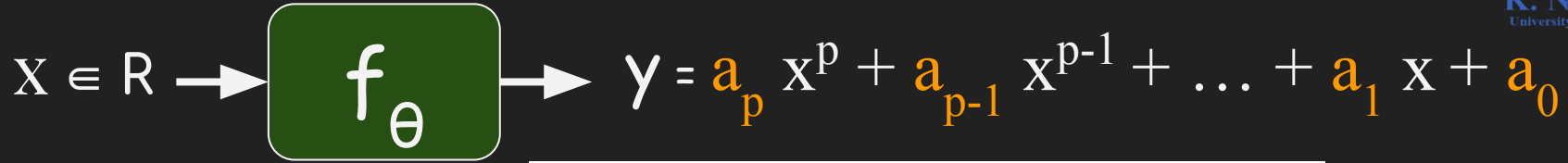
$X \in \mathbb{R} \rightarrow$    $\rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$



polynomial degree  $p = 14$




K. N. Toosi  
University of Technology

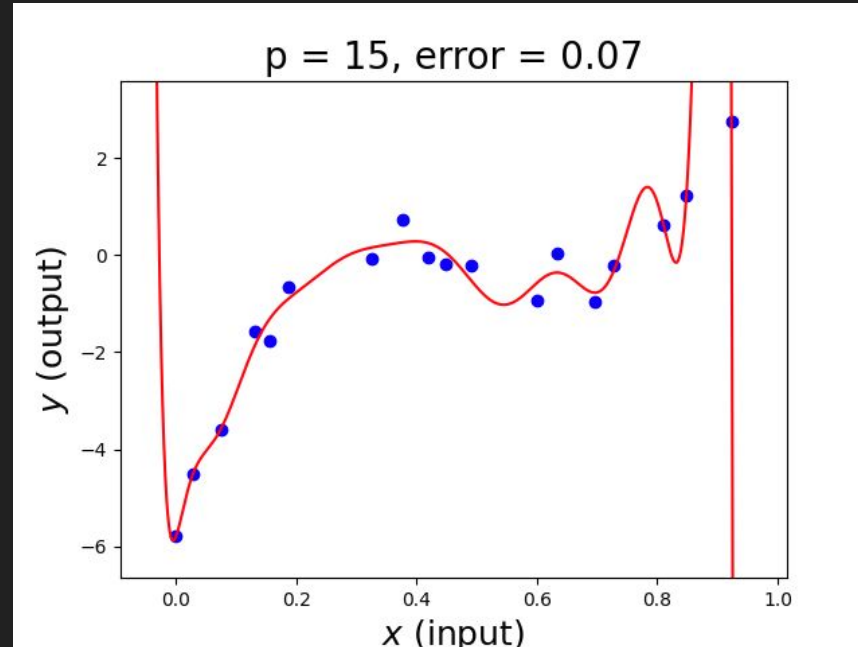


polynomial degree  $p = 15$



K. N. Toosi  
University of Technology


$X \in \mathbb{R} \rightarrow$    $\rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$

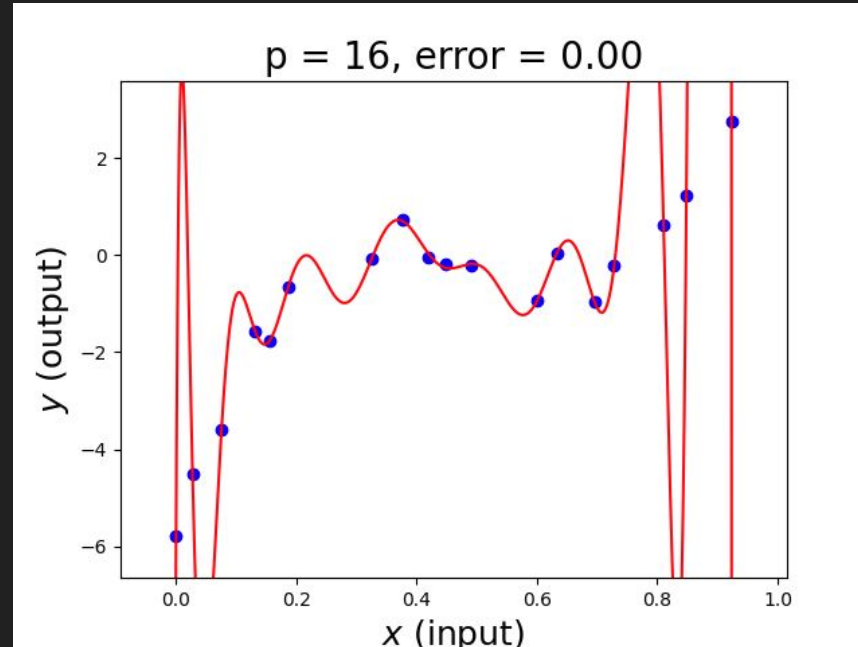


polynomial degree  $p = 16$



K. N. Toosi  
University of Technology

$X \in \mathbb{R} \rightarrow$    $\rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$

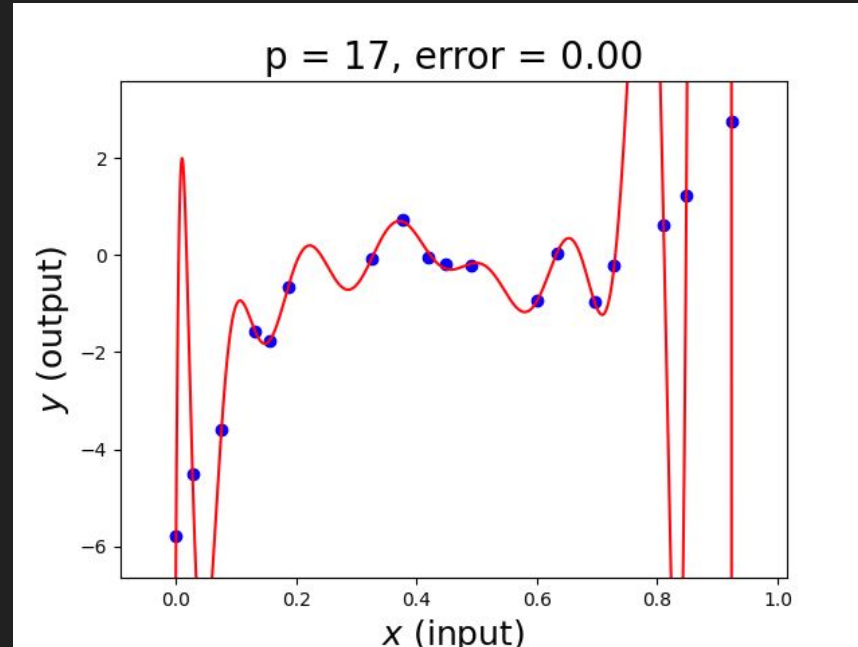
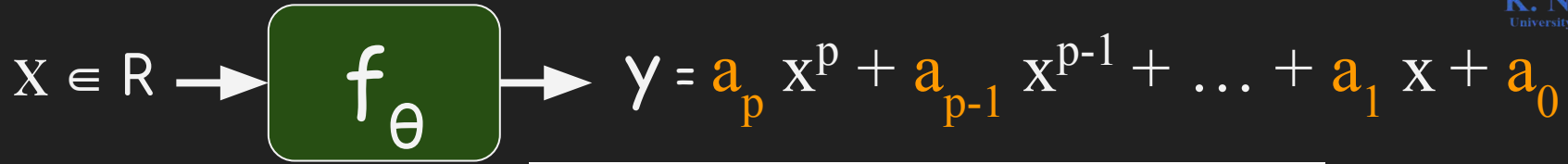




polynomial degree  $p = 17$



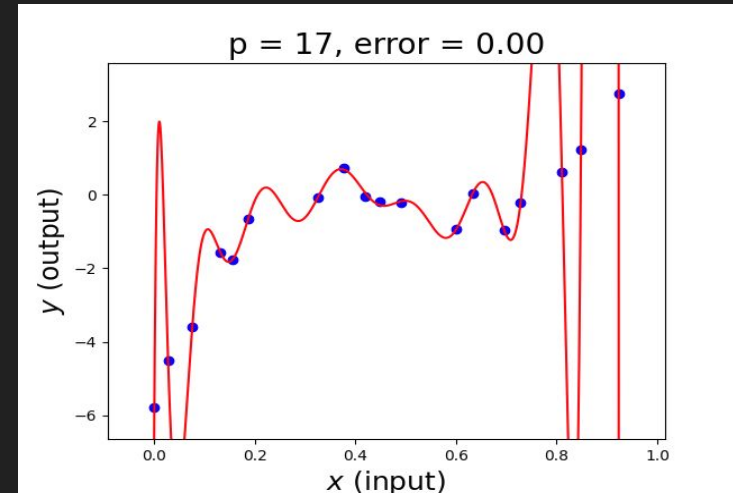
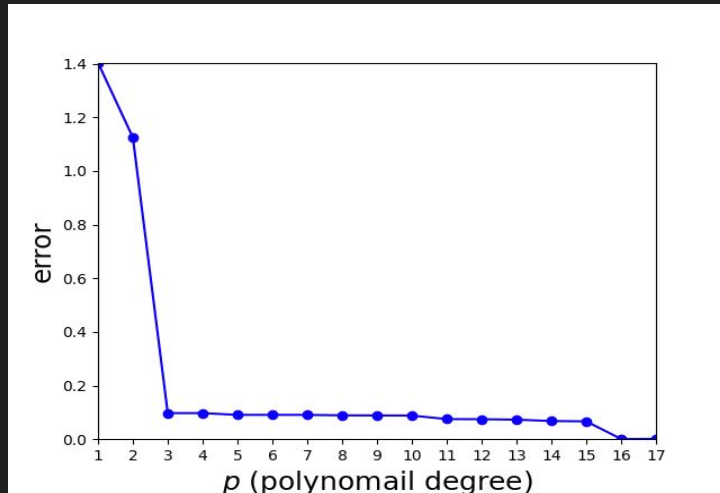
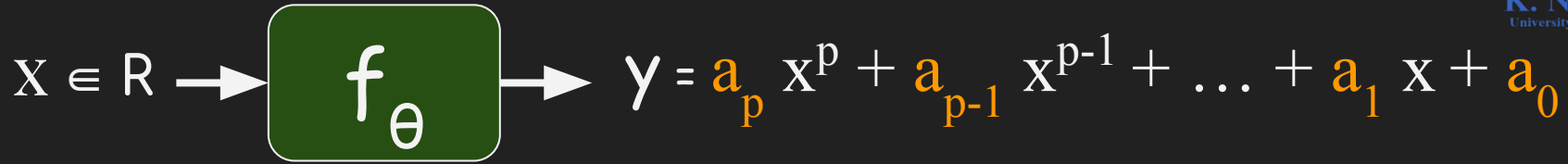
K. N. Toosi  
University of Technology



# polynomial degree $p = 17$



K. N. Toosi  
University of Technology



# Learning from data



K. N. Toosi  
University of Technology

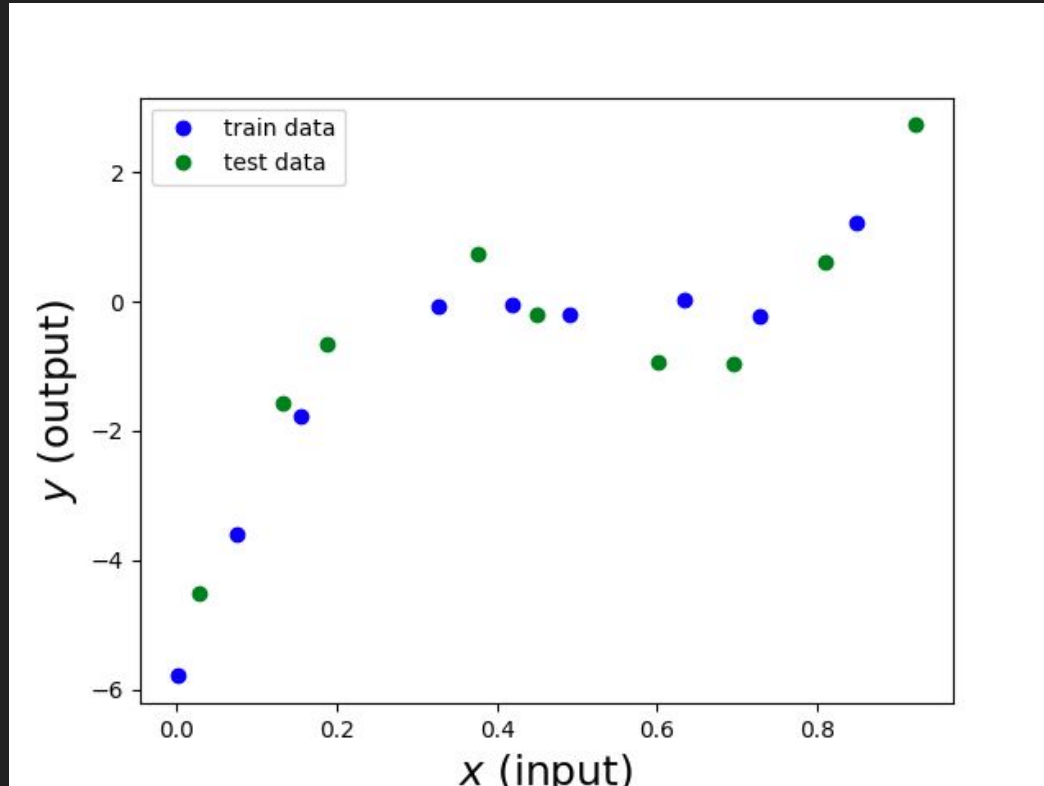


- Parameter Learning:
  - A collection of input-output pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ,
  - choose  $\theta$  such that  $y = f(\theta, x)$  is a reasonable output
    - for training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
    - for unseen data
  - **Generalization:** How well the model work on unseen data
  - How to evaluate error on **unseen data?**

# Train-Test split



K. N. Toosi  
University of Technology



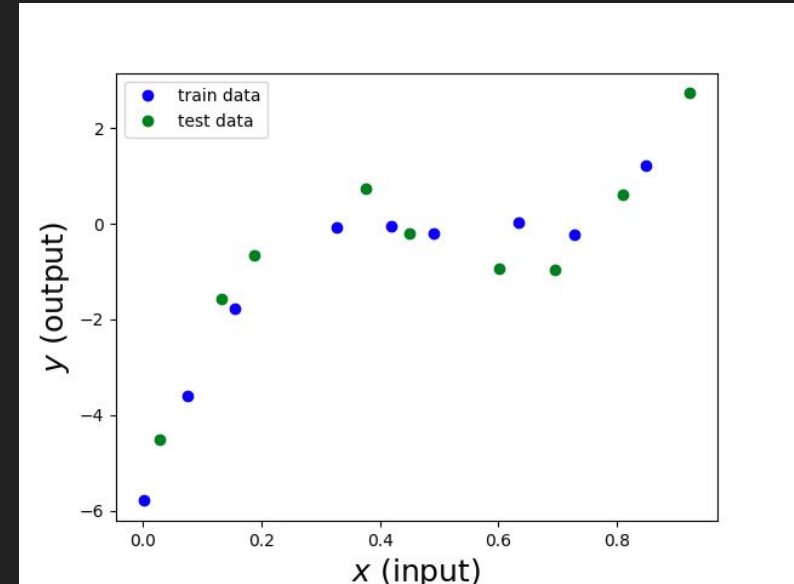
# Train-Test split



- Divide data into train and test sets
- Train the model using train data

Evaluate on test data

$$D = D_{\text{test}} \cup D_{\text{train}}$$
$$\text{Train Error} = \sum_{(x_i, y_i) \in D_{\text{train}}} (f(x_i; \theta) - y_i)^2$$
$$\text{Test Error} = \sum_{(x_i, y_i) \in D_{\text{test}}} (f(x_i; \theta) - y_i)^2$$



# Train-Test split



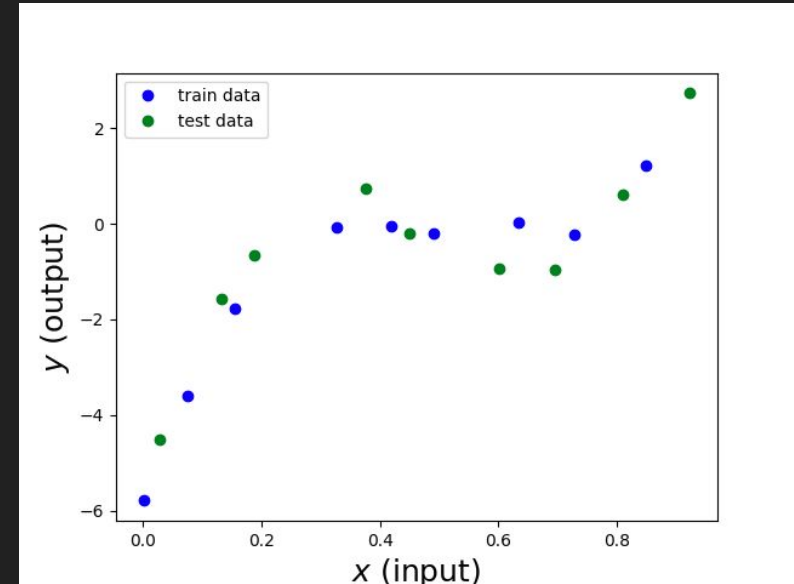
K. N. Toosi  
University of Technology

- Divide data into train and test sets
- Train the model using train data
- Evaluate on test data

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{i \in \text{train}} (f(\theta, x_i) - y_i)^2$$

$$\text{Err}_{\text{train}} = \sum_{i \in \text{train}} (f(\theta^*, x_i) - y_i)^2$$

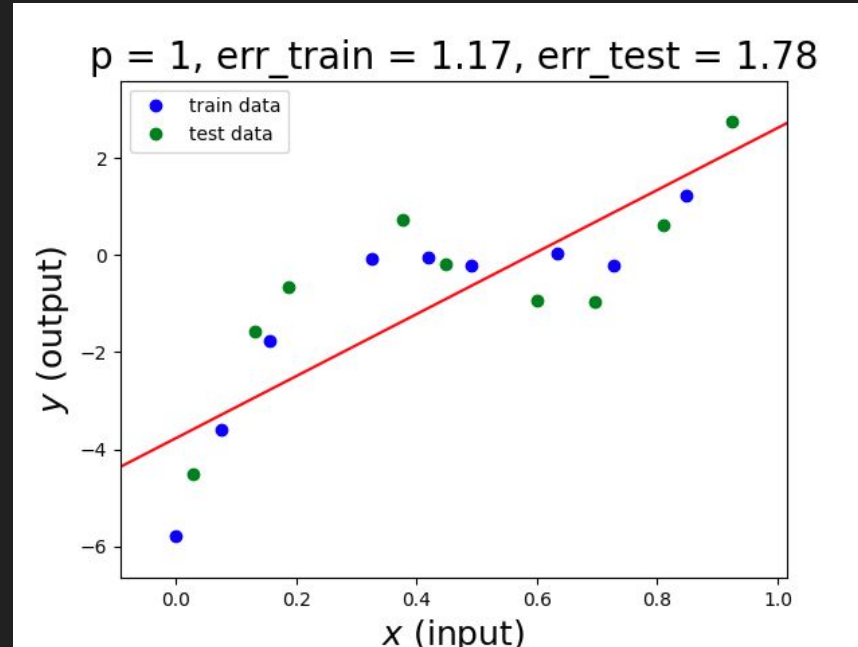
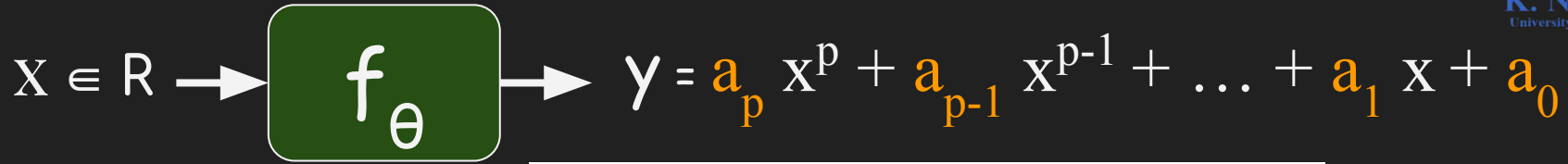
$$\text{Err}_{\text{test}} = \sum_{i \in \text{test}} (f(\theta^*, x_i) - y_i)^2$$



# polynomial degree $p = 1$



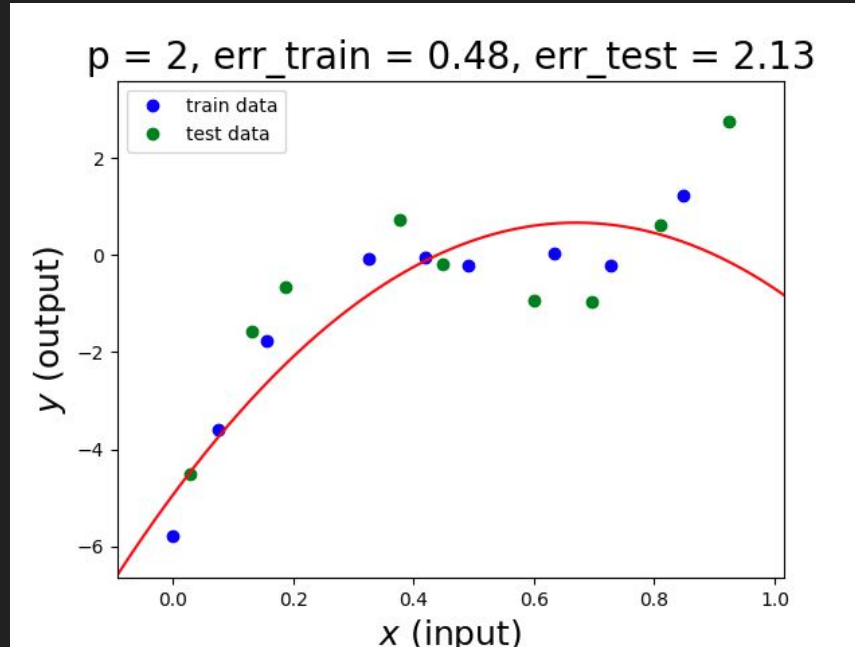
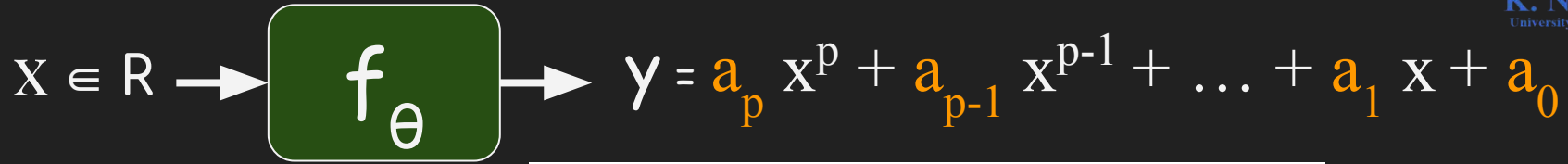
K. N. Toosi  
University of Technology



# polynomial degree $p = 2$



K. N. Toosi  
University of Technology

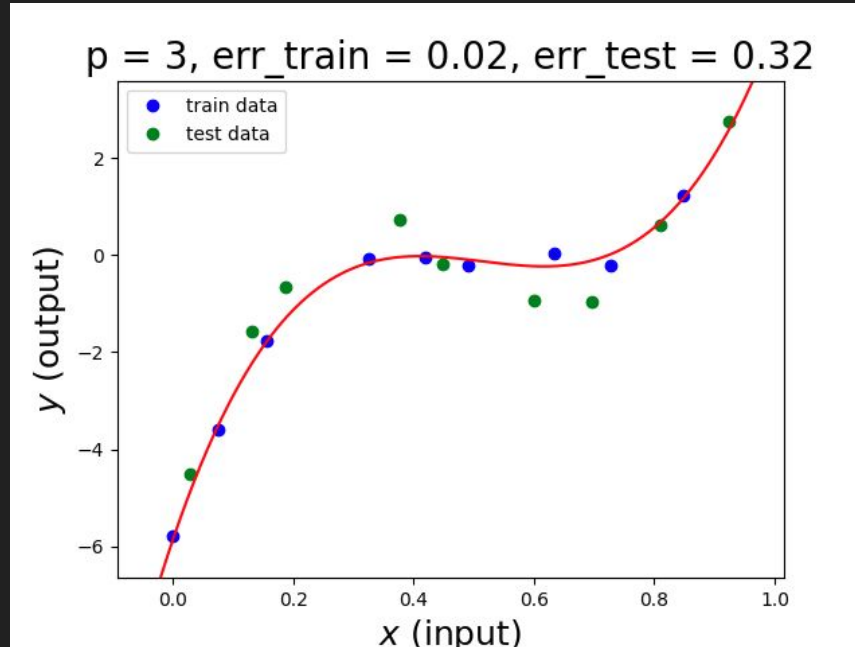
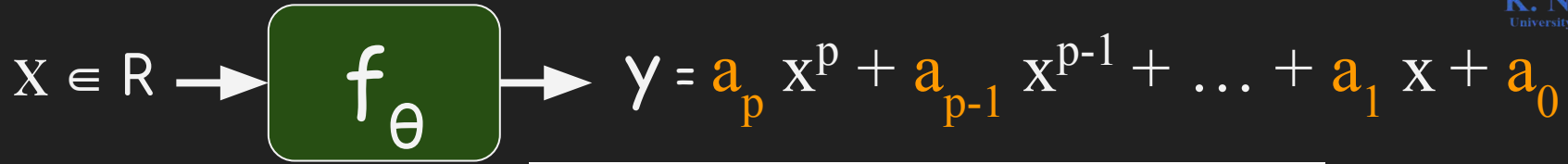




# polynomial degree $p = 3$



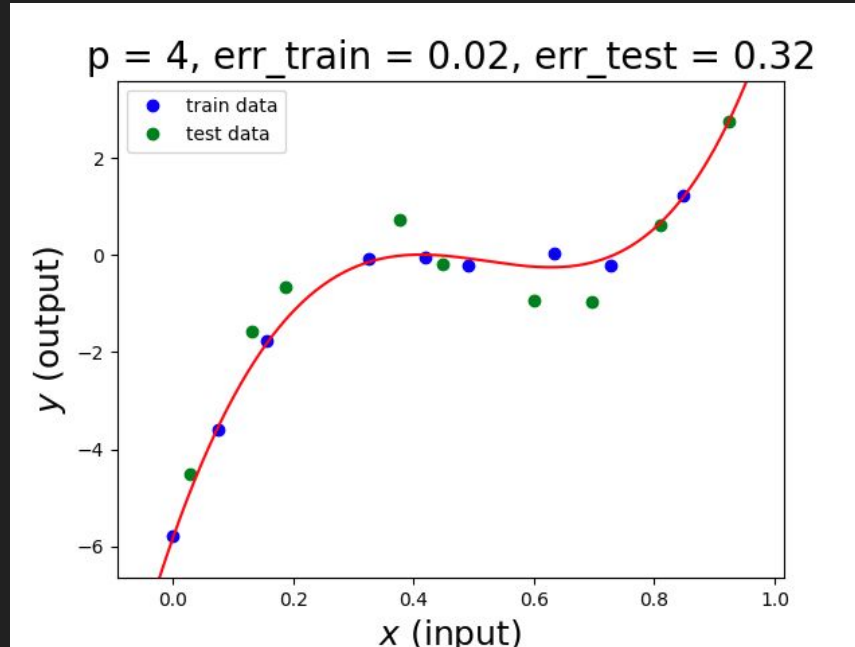
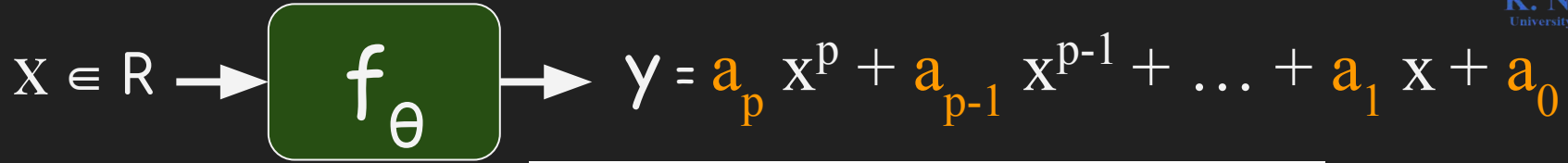
K. N. Toosi  
University of Technology



# polynomial degree $p = 4$



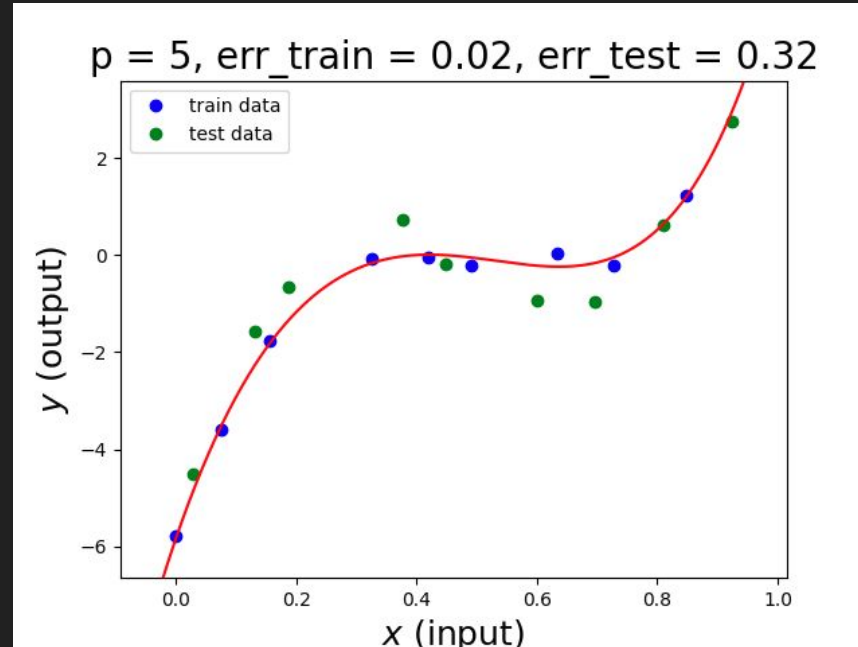
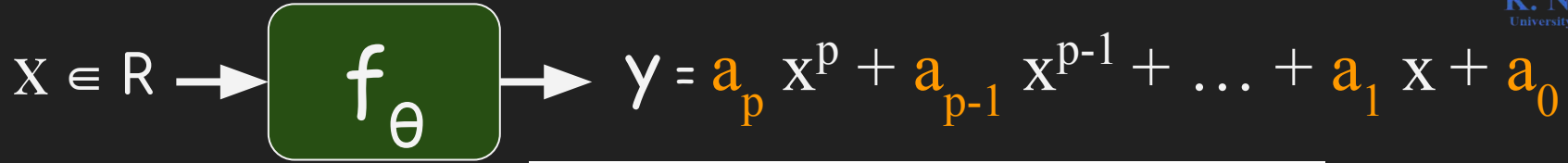
K. N. Toosi  
University of Technology



# polynomial degree $p = 5$



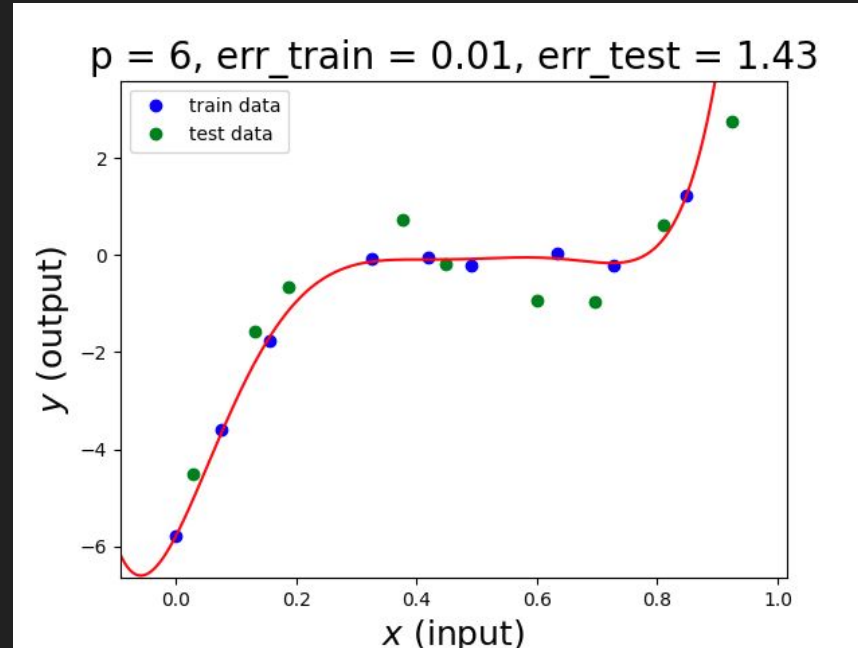
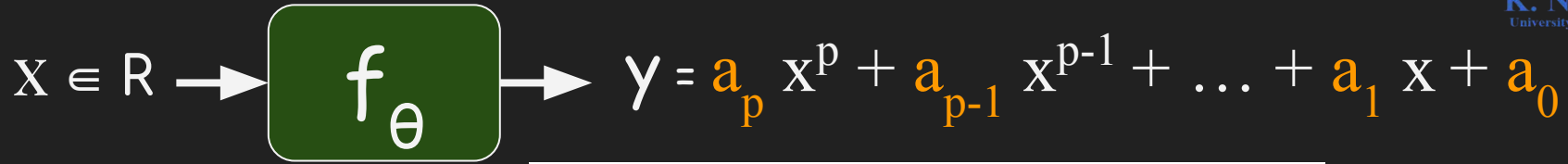
K. N. Toosi  
University of Technology



# polynomial degree $p = 6$



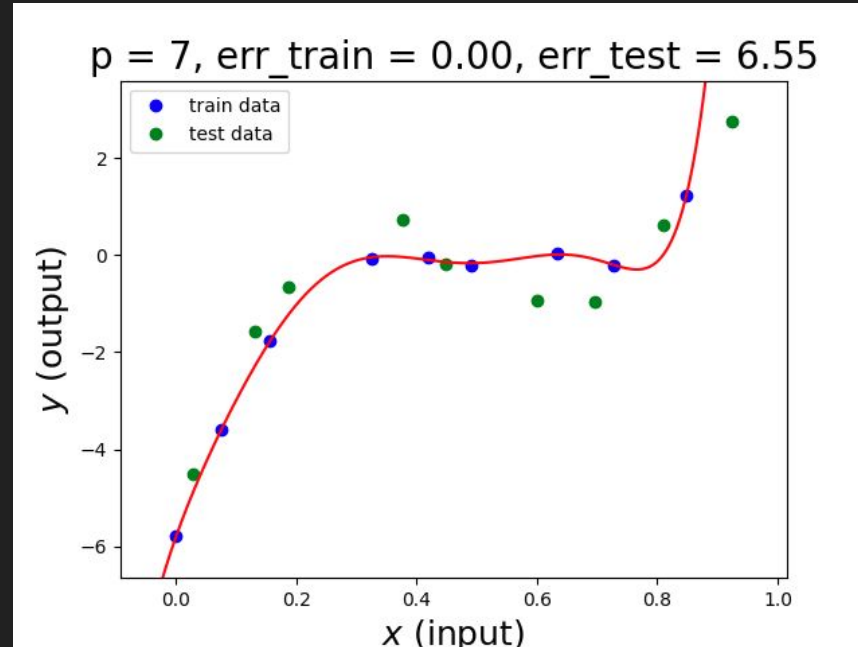
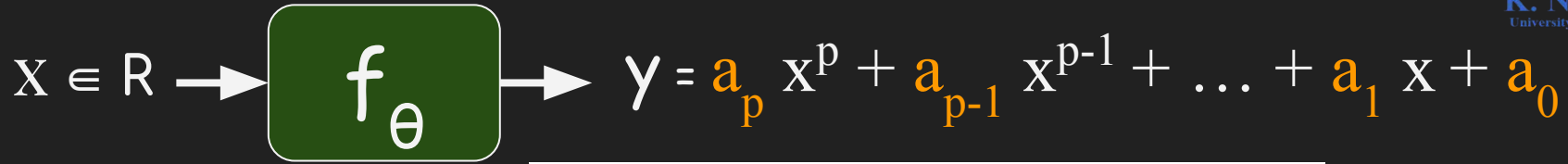
K. N. Toosi  
University of Technology



# polynomial degree $p = 7$



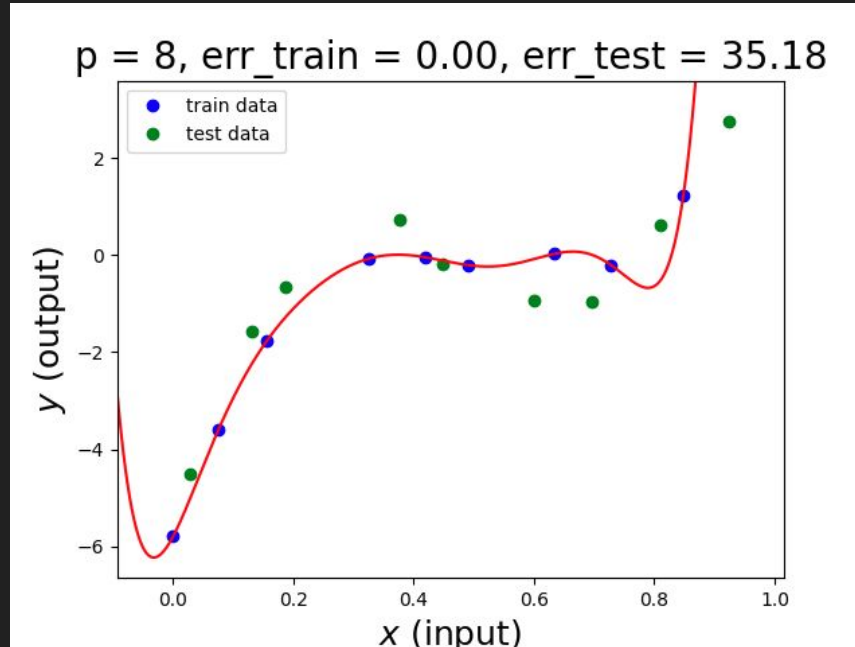
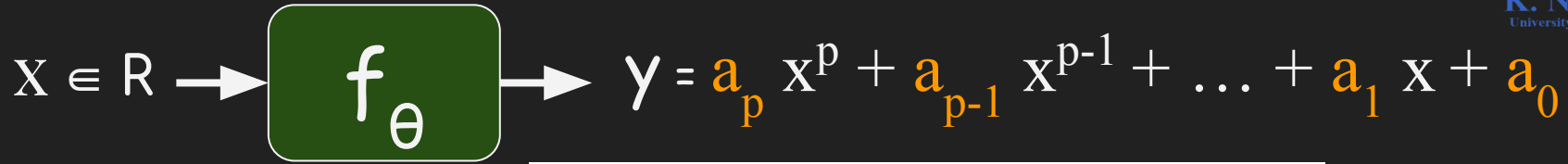
K. N. Toosi  
University of Technology



# polynomial degree $p = 8$



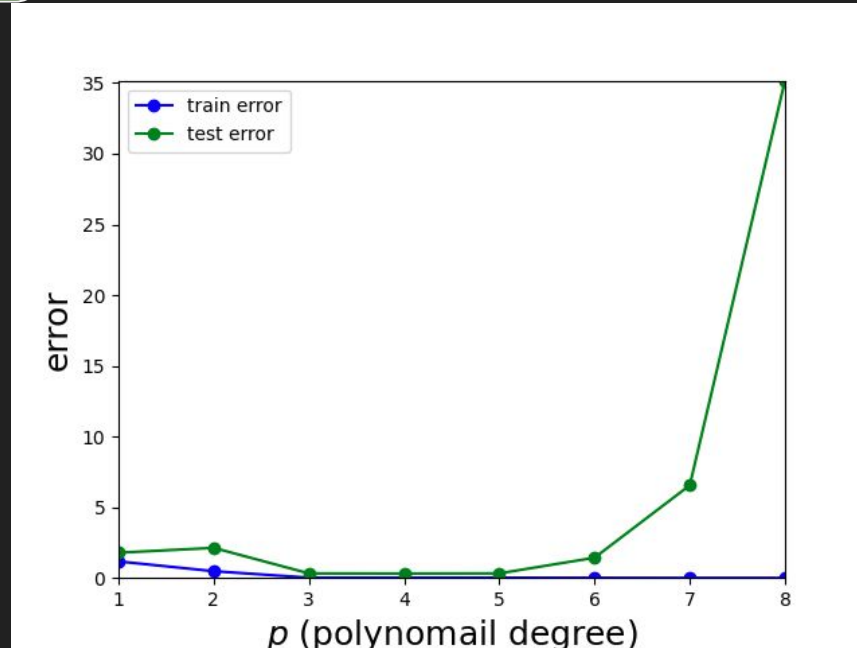
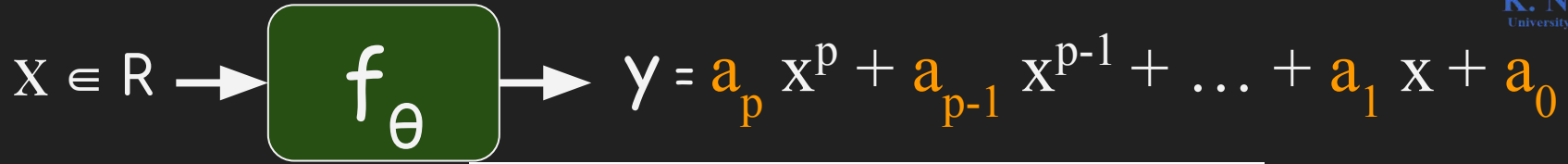
K. N. Toosi  
University of Technology



# Train and test errors



K. N. Toosi  
University of Technology

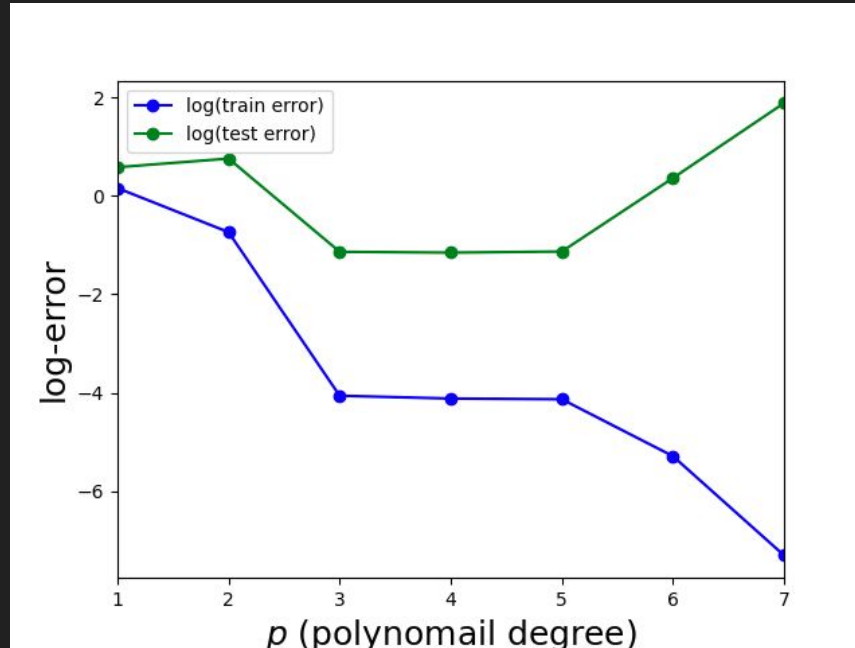


# Train and test errors



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$







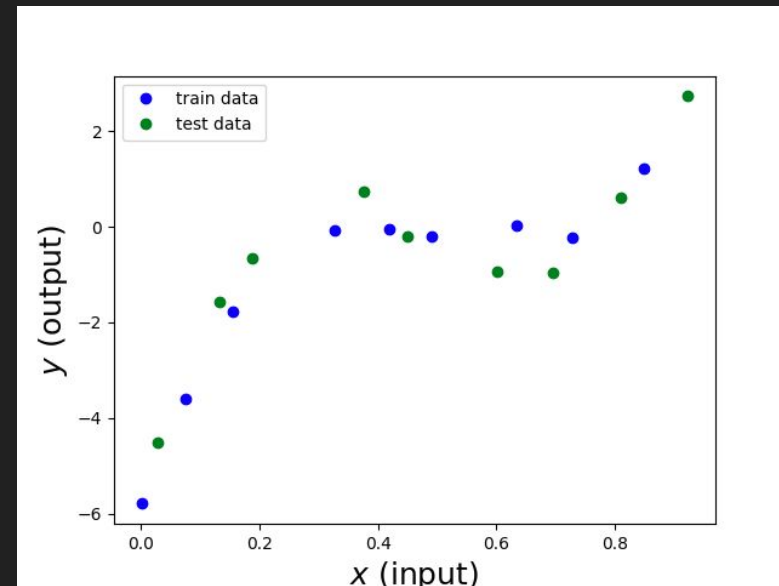
# Cross-validation

How good model work on new data?

- Divide data into train and test sets
- Train the model using train data
- Evaluate on test data

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{i \in \text{train}} (f(\theta, x_i) - y_i)^2$$

$$\text{Err}_{\text{test}} = \sum_{i \in \text{test}} (f(\theta^*, x_i) - y_i)^2$$



- divide the data in different ways and report the average error

# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- Hyperparameters
  - Are not learned during model fitting
  - Used for
    - Learning the structure
    - Choosing the right model
    - ...
- In polynomial regression the degree parameter ( $p$ ) is a hyperparameter.

# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

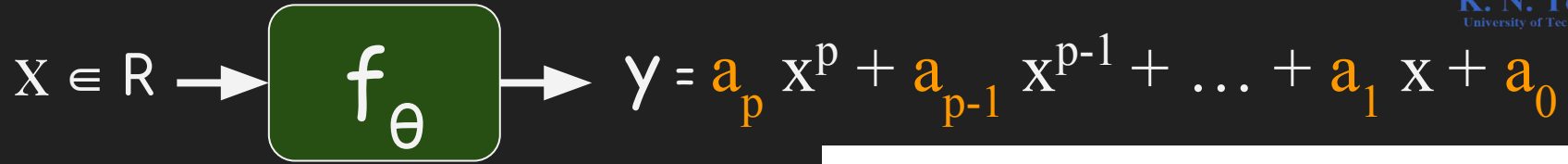
$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- Hyperparameters
  - Are not learned during model fitting
  - Used for
    - Learning the structure
    - Choosing the right model
- In polynomial regression the degree parameter ( $p$ ) is a hyperparameter.
- How to choose hyper parameters?

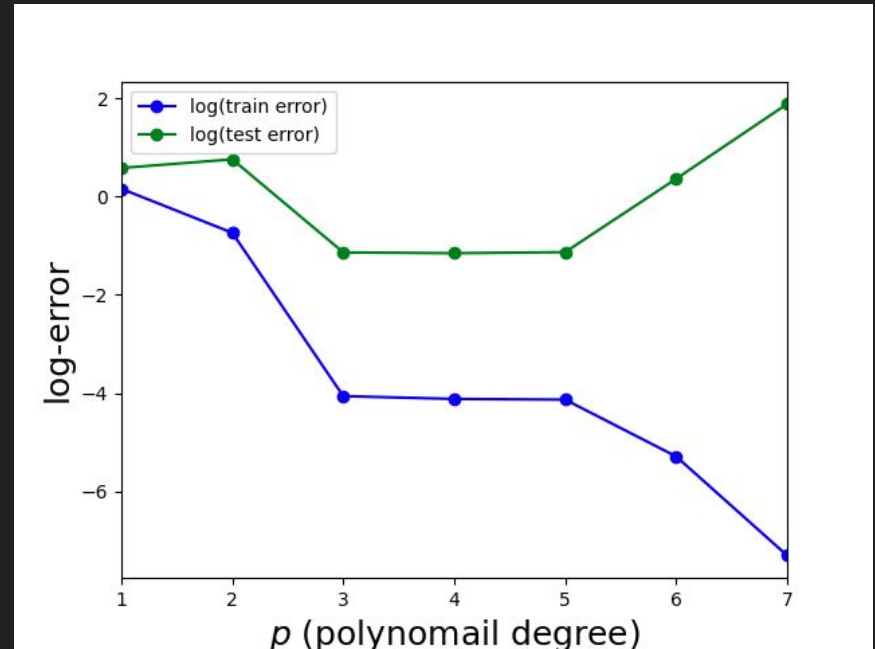
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology



- How to choose  $p$ ?



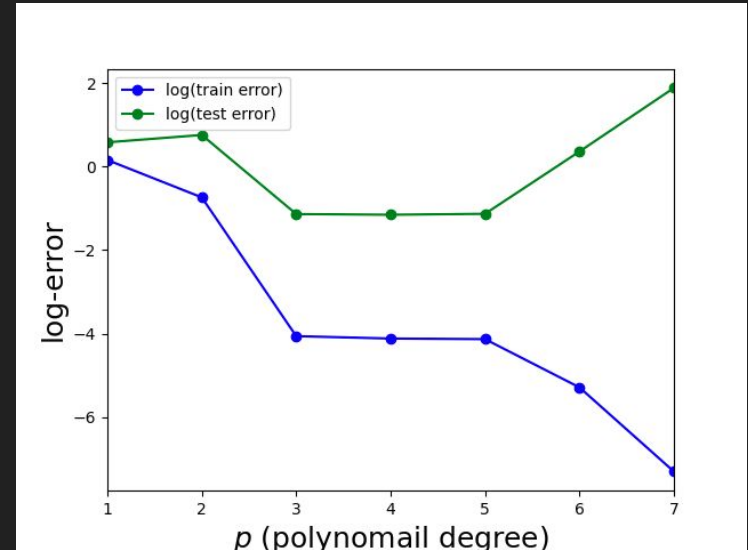
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- How to choose  $p$ ?
  - Split data to train and test sets
  - Train regular parameters using training data
  - Select the  $p$  that minimizes test error



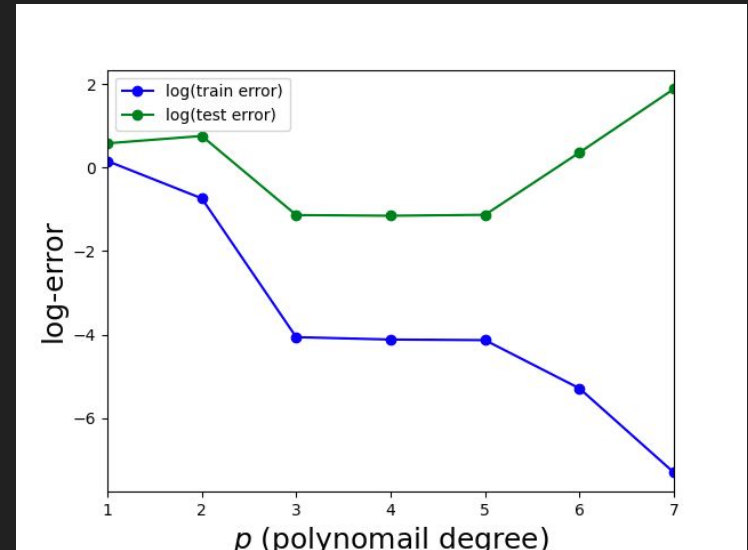
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- How to choose  $p$ ? (cross-validation)
  - Split data to train and test sets
  - Train regular parameters using training data
  - Select the  $p$  that minimizes test error
- [Divide in different ways and use the average error. ]



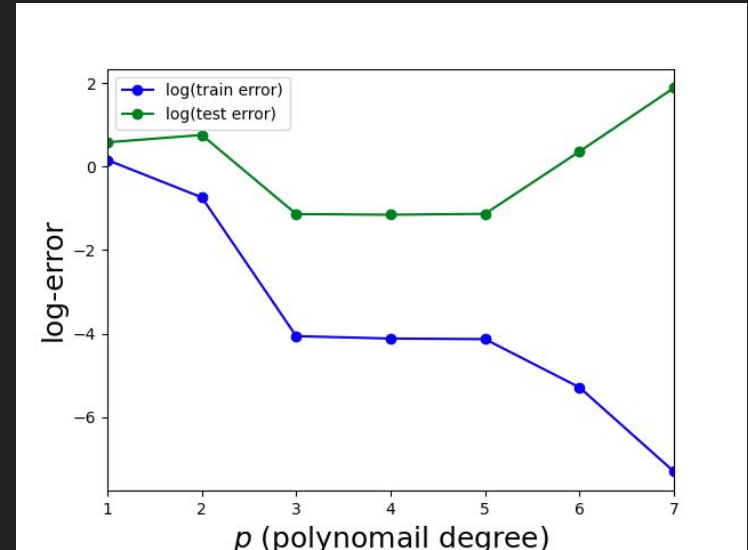
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- How to choose  $p$ ?
  - Split data to train and test sets
  - Train regular parameters using training data
  - Select the  $p$  that minimizes test error
- How to evaluate model on unseen data?



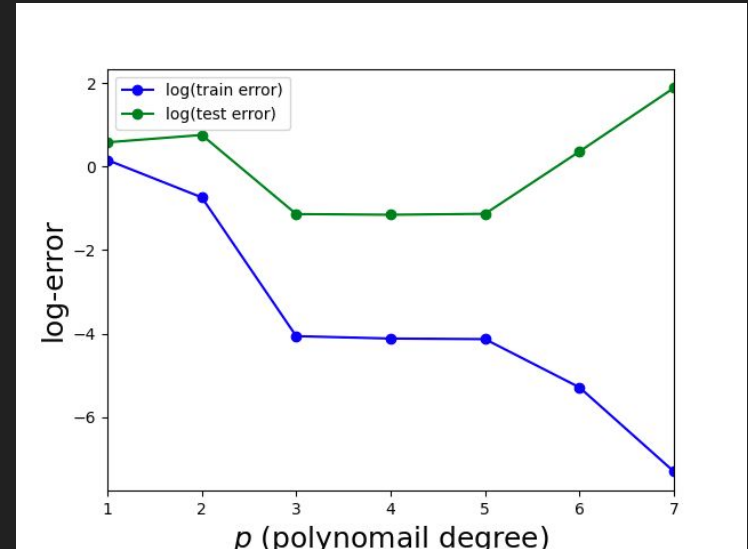
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- How to choose  $p$ ?
  - Split data to train and test sets
  - Train regular parameters using training data
  - Select the  $p$  that minimizes test error
- How to evaluate model on unseen data?
  - Error on train set?
    - train data is used to learn parameters
  - Error on test set?
    - test data is used to choose hyperparameters





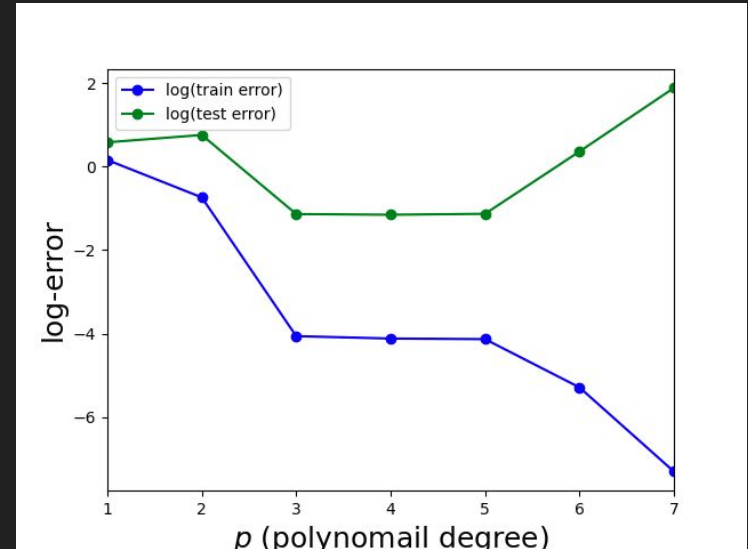
# How to choose hyper-parameters?



K. N. Toosi  
University of Technology

$$X \in \mathbb{R} \rightarrow \boxed{f_{\theta}} \rightarrow y = a_p x^p + a_{p-1} x^{p-1} + \dots + a_1 x + a_0$$

- **How to evaluate model on unseen data?**
  - Split data into **3** subsets: **train**, **validation**, and **test**
  - Train regular parameters using **train set**
  - Choose hyperparameters using **validation set**
  - Report system error on **test set**
- [Divide in different ways and report the average error. ]



# How to avoid overfitting

- Limit model complexity (e.g. degree of polynomial)
- **Regularization**

