

# LINEAR REGRESSION & REGULARIZATION

Narges Norouzi

# EXTENDING LINEAR REGRESSION TO MORE COMPLEX MODELS

- The inputs  $X$  for linear regression can be:
  - Original quantitative inputs
  - Transformation of quantitative inputs (log, exp, square root, square, etc.)
  - Polynomial transformation (example:  $y = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$ )
  - Dummy encoding of categorical inputs
  - Interactions between variables (example:  $x_3 = x_1 \times x_2$ )
- This allows use of linear regression techniques to fit non-linear datasets.

# LINEAR BASIS FUNCTION MODEL

- Generally,

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j \phi_j(x)$$

Basis Function

- Typically,  $\phi_0(x) = 1$  so that  $\theta_0$  acts as a bias.
- In the simplest case, we can use linear basis function:

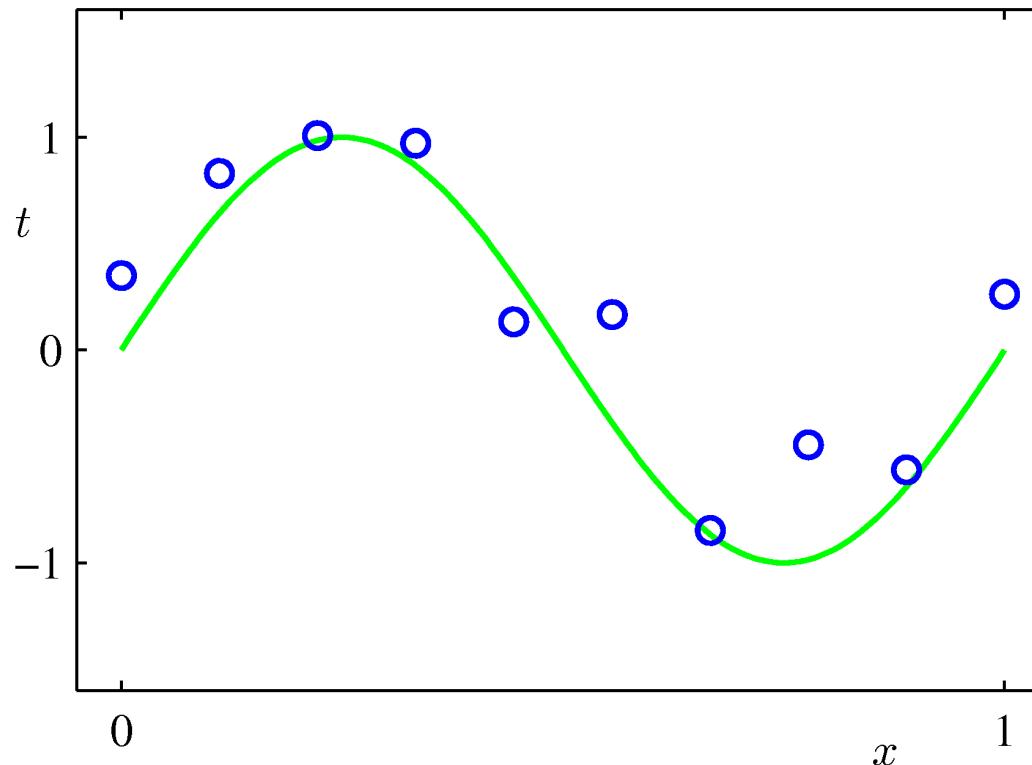
$$\phi_j(x) = x_j$$

- Polynomial basis function:  $\phi_j(x) = x^j$

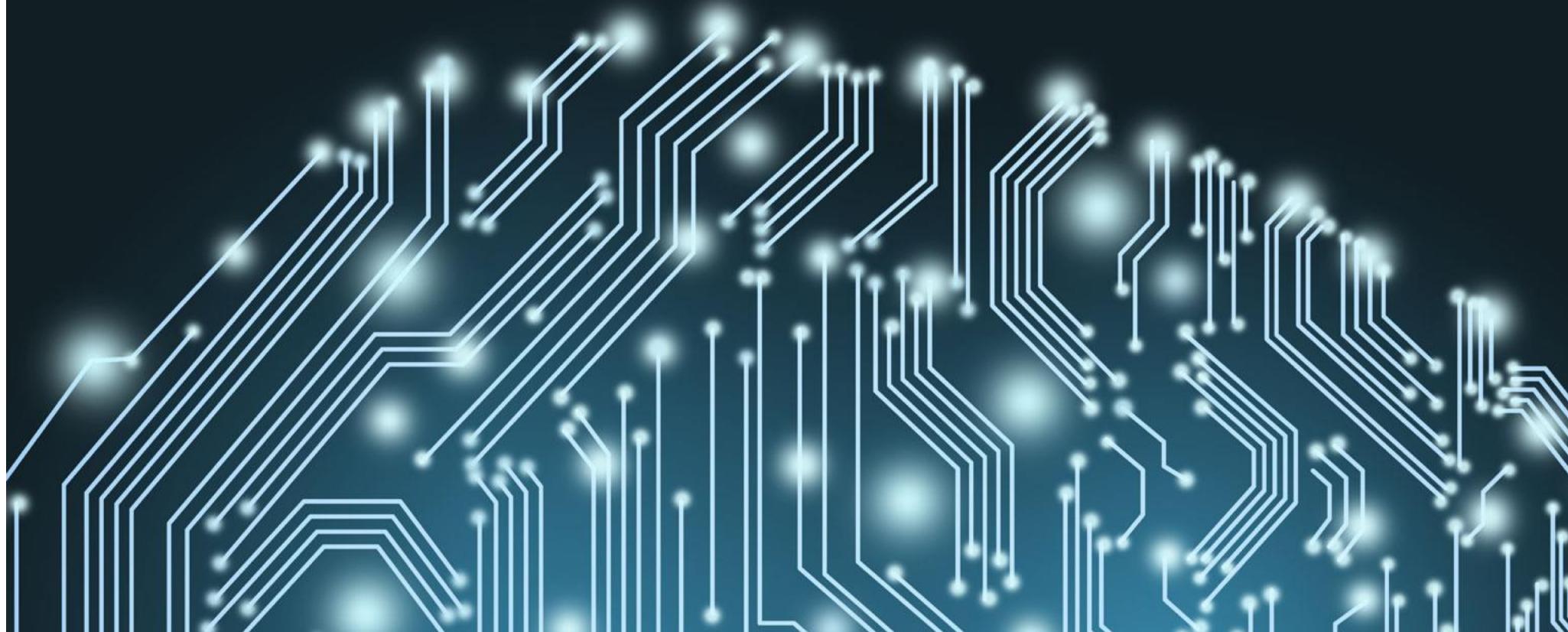
$$\frac{(x - \mu_j)^2}{2s^2}$$

- Gaussian basis function:  $\phi_j(x) = e^{-\frac{(x - \mu_j)^2}{2s^2}}$

# EXAMPLE OF FITTING A POLYNOMIAL CURVE WITH A LINEAR MODEL



$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d = \sum_{j=0}^d \theta_j x^j$$



# VECTOR CALCULATION

# LINEAR ALGEBRA CONCEPTS

- Vector in  $\mathbb{R}^d$  is an ordered set of  $d$  real numbers.
- $v = [1, 2, 3, 4]$  is in  $\mathbb{R}^4$  and is a column vector.
- An m-by-n matrix is an object with  $m$  rows and  $n$  columns, where each entry is a real number

$$\begin{bmatrix} 1 & 0 & 2 \\ 4 & 0.5 & 7.8 \end{bmatrix}$$

- Transposing the matrix

$$\begin{bmatrix} a \\ b \end{bmatrix}^T = [a \ b], \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

Note:  $(Ax)^T = x^T A^T$

# LINEAR ALGEBRA CONCEPTS

- Vector norms:

- $L_p$  norm of  $v = (v_1, v_2, \dots, v_k) = (\sum_i |v_i|^p)^{\frac{1}{p}}$
- Common norms  $L_1$  and  $L_2$
- Length of the vector  $v$  is  $L_2(v)$

# LINEAR ALGEBRA CONCEPTS

- Vector dot product:  $u \cdot v = (u_1 \ u_2) \cdot (v_1 \ v_2) = u_1 \times v_1 + u_2 \times v_2$ 
  - The dot product of  $u$  with itself =  $\text{length}(u)^2 = \|u\|_2^2$
- Matrix product:

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}$$

$$A \times B = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

# LINEAR ALGEBRA CONCEPTS

- Vector products:

- Dot product:

$$u \cdot v = u^T v = [u_1 \ u_2] \cdot \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = u_1 \times v_1 + u_2 \times v_2$$

- Outer product:

$$uv^T = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} [v_1 \ v_2] = \begin{bmatrix} u_1 v_1 & u_1 v_2 \\ u_2 v_1 & u_2 v_2 \end{bmatrix}$$

# VECTORIZATION

- Benefits of vectorization
  - More compact equations
  - Faster code (using optimized matrix libraries)
- Consider our model:  $h(x) = \sum_{j=0}^d \theta_j x_j$

- Let  $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$

We can write our model in vectorized form as  $h(x) = x^T \theta$

# VECTORIZATION

- Consider our model for  $n$  instances:  $h(x^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$

- Let

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \in \mathbb{R}^{d+1 \times 1}, \quad X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \ddots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d+1}$$

- We can write the vectorized form as  $h_\theta(x) = X\theta$

# VECTORIZATION

- For the linear regression cost function:

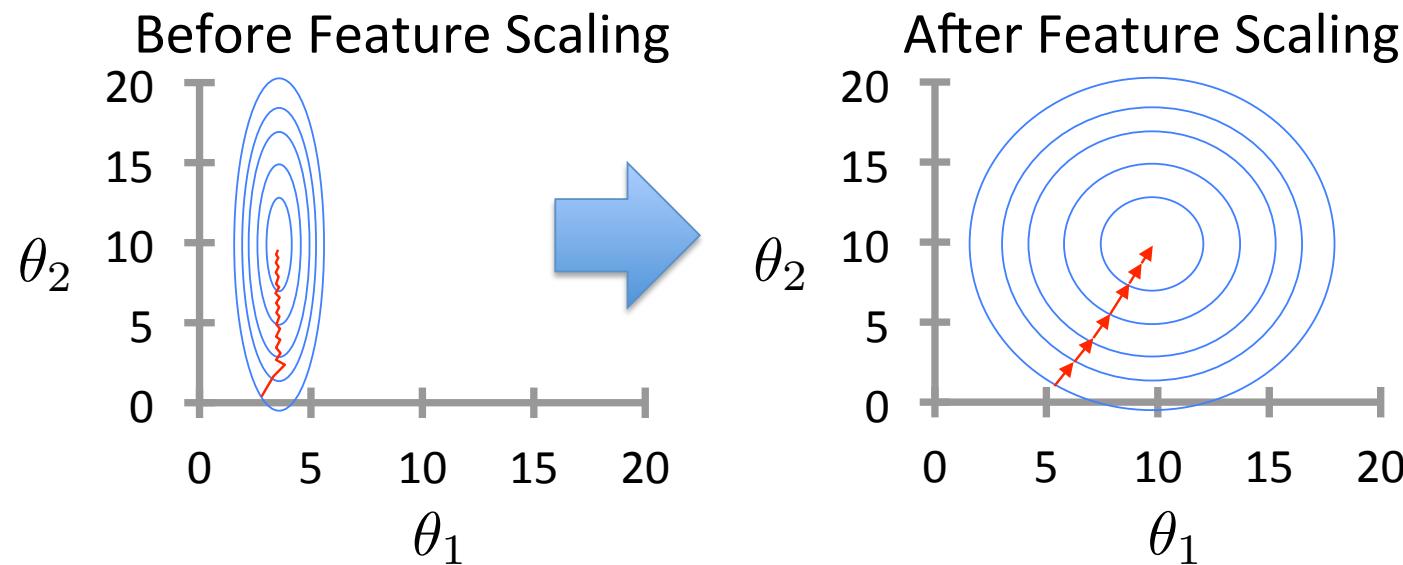
$$\begin{aligned}Cost(\theta) &= \frac{1}{2 \times n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \\&= \frac{1}{2 \times n} \sum_{i=1}^n (x^{(i)T} \theta - y^{(i)})^2 \\&= \frac{1}{2 \times n} (X\theta - Y)^T (X\theta - Y)\end{aligned}$$

- Note that:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^{n \times 1}, \quad X \in \mathbb{R}^{n \times d+1}, \quad \theta \in \mathbb{R}^{d+1 \times 1} \rightarrow X\theta \in \mathbb{R}^{n \times 1}$$
$$X\theta - Y \in \mathbb{R}^{n \times 1}, \quad (X\theta - Y)^T \in \mathbb{R}^{1 \times n} \rightarrow Cost(\theta) \in \mathbb{R}$$

# IMPROVING LEARNING: FEATURE SCALING

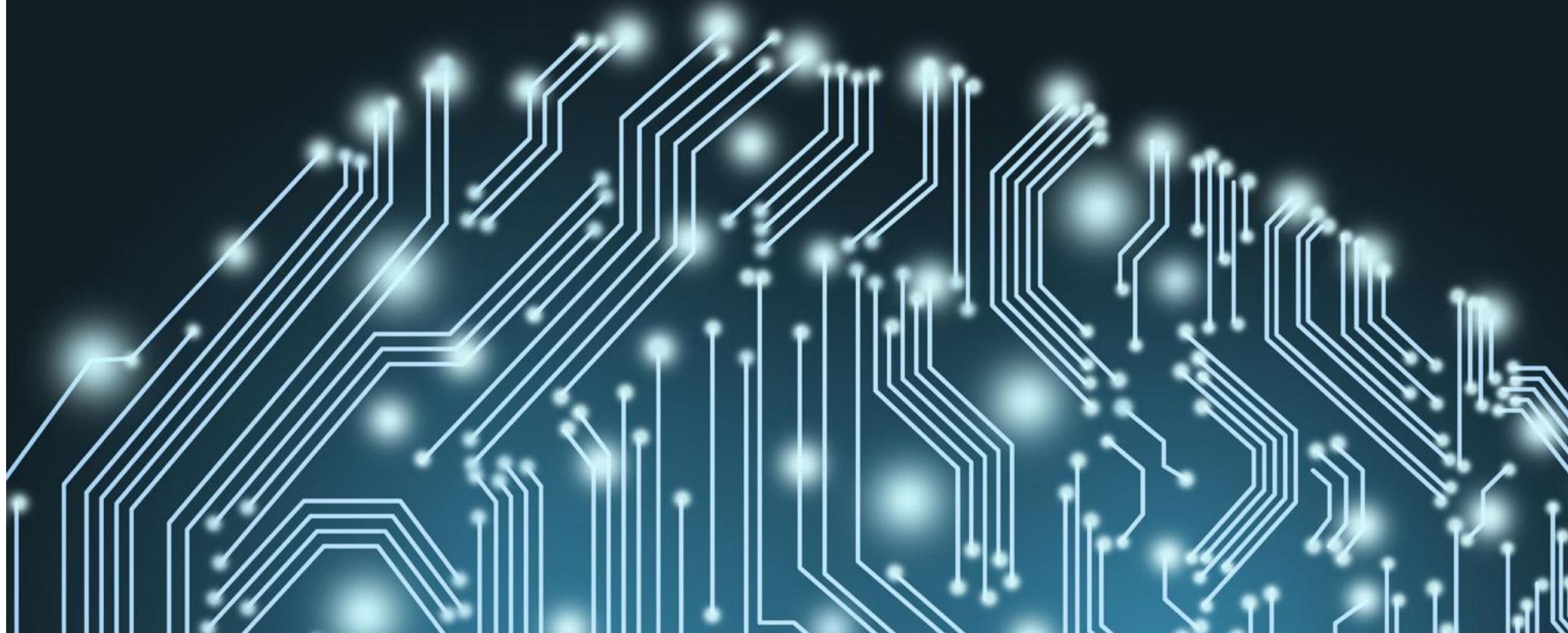
- Idea: ensure that features have similar scales



- Makes gradient descent converge much faster

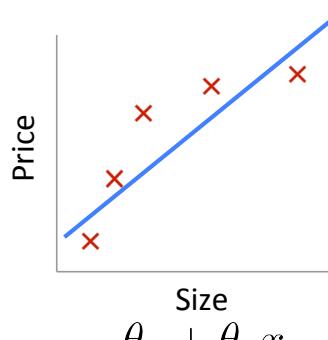
# CLASS EXERCISE

[bit.ly/ce-6](https://bit.ly/ce-6)

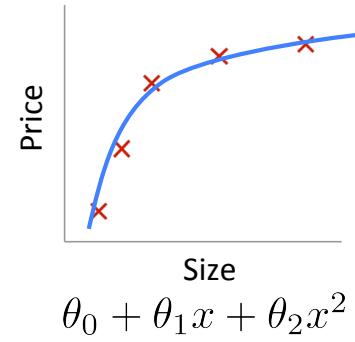


# REGULARIZATION AND OVERFITTING

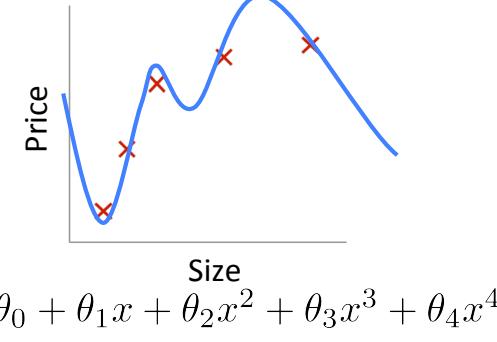
# QUALITY OF FIT



$\theta_0 + \theta_1 x$   
Underfitting  
(high bias)



Correct fit



Overfitting  
(high variance)

- Overfitting: The learned hypothesis may fit the training set very well
  - but fails to generalize to new examples

# REGULARIZATION

- A method for automatically controlling the complexity of the learned hypothesis
- Idea: penalize large values of  $\theta_j$
- There are other techniques we will learn later such as dropout technique in neural networks or dimensionality reduction.

# REGULARIZATION

- Linear regression objective function

$$Cost(\theta) = \left\{ \frac{1}{2 \times n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \right\} + \left\{ \lambda \sum_{j=1}^d \theta_j^2 \right\}$$

Model fit to data

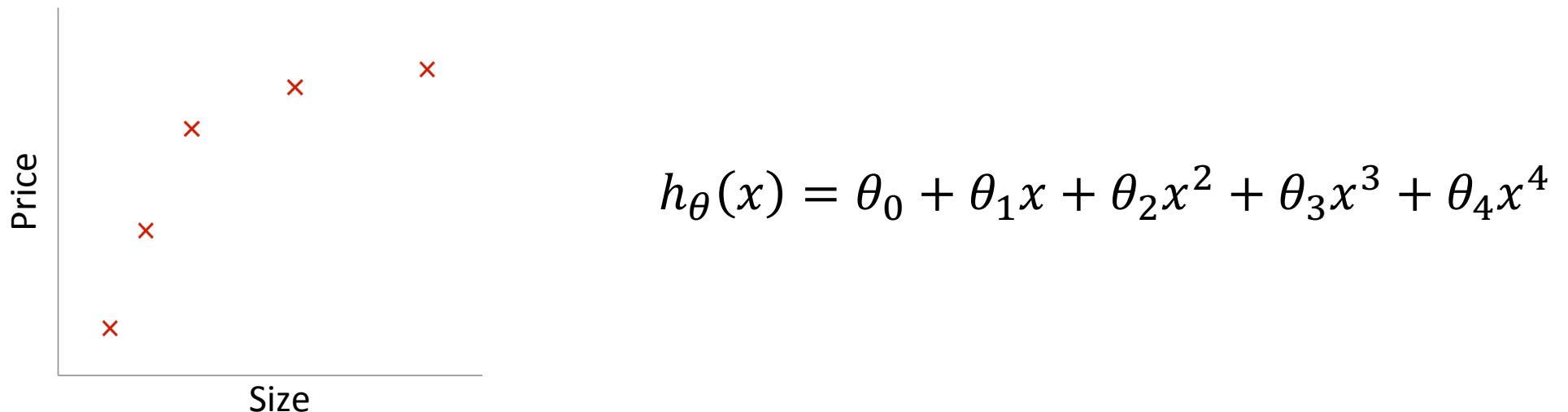
Regularization

- $\lambda$  is the regularization parameter ( $\lambda \geq 0$ )
- No regularization on  $\theta_0$
- This exact regularizer pulls coefficients/parameters to 0

# UNDERSTANDING REGULARIZATION

- What happens if we set  $\lambda$  too huge?

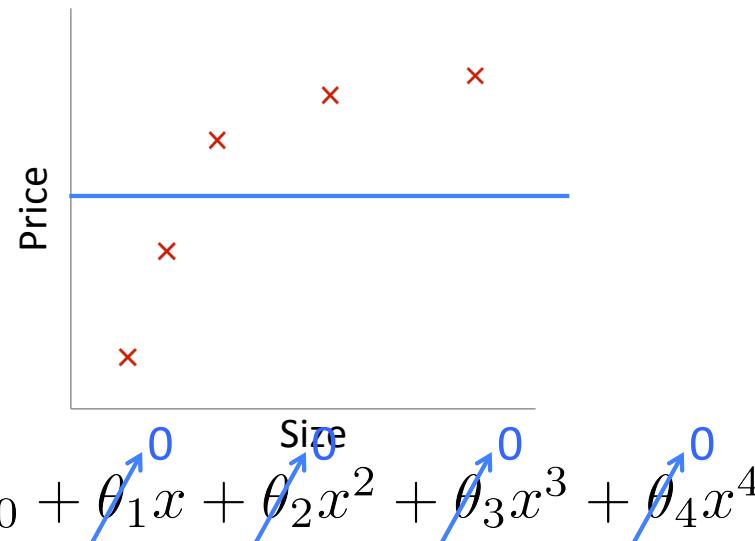
$$Cost(\theta) = \frac{1}{2 \times n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^d \theta_j^2$$



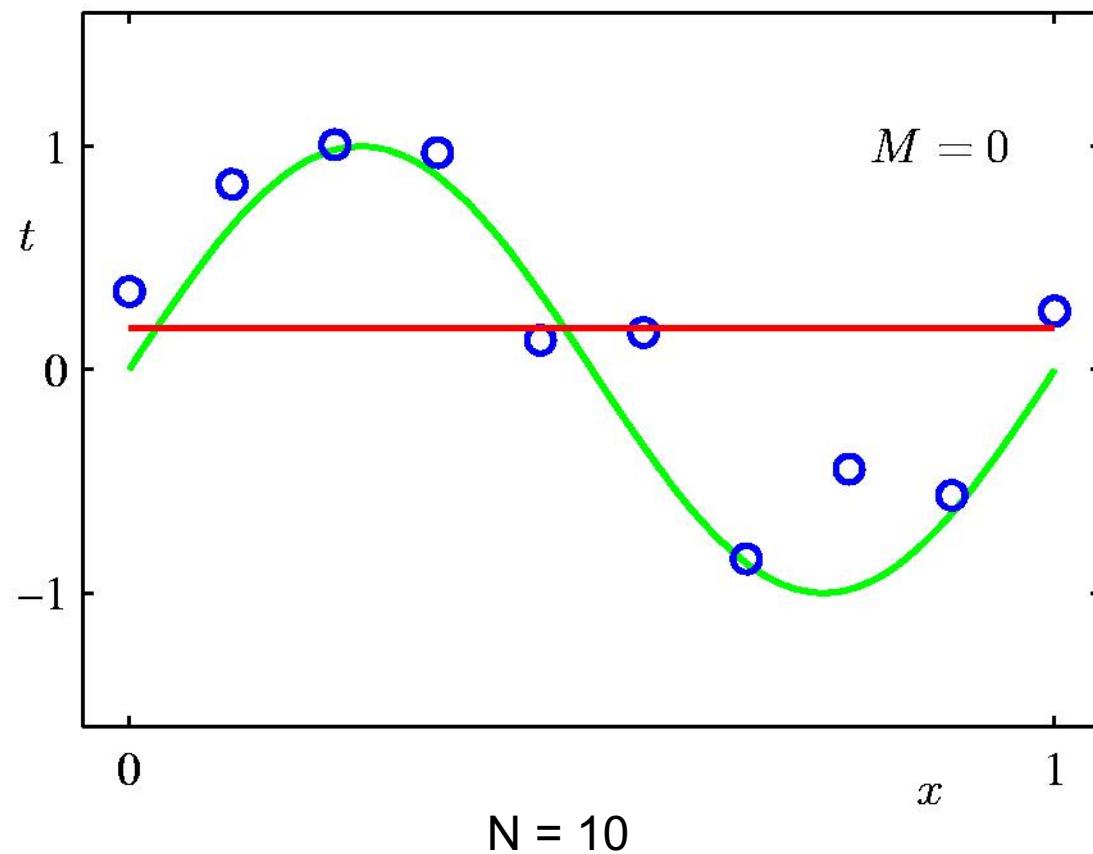
# UNDERSTANDING REGULARIZATION

- What happens if we set  $\lambda$  too huge?

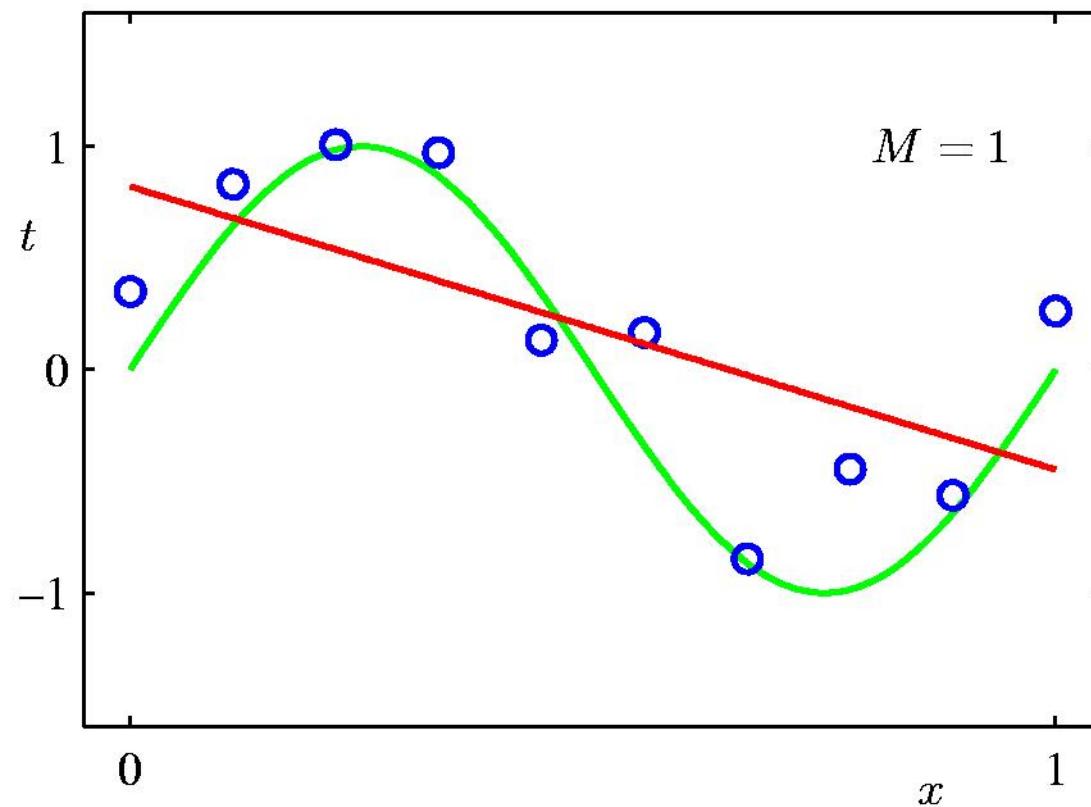
$$Cost(\theta) = \frac{1}{2 \times n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^d \theta_j^2$$



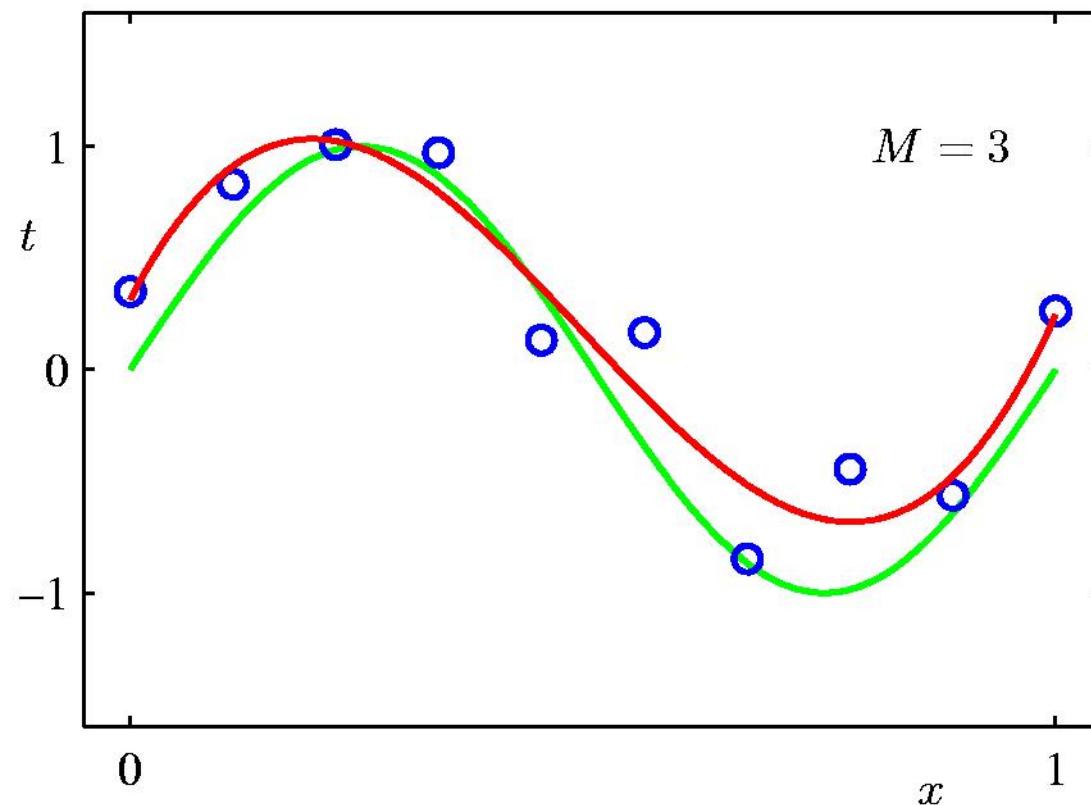
# 0<sup>TH</sup> ORDER POLYNOMIAL



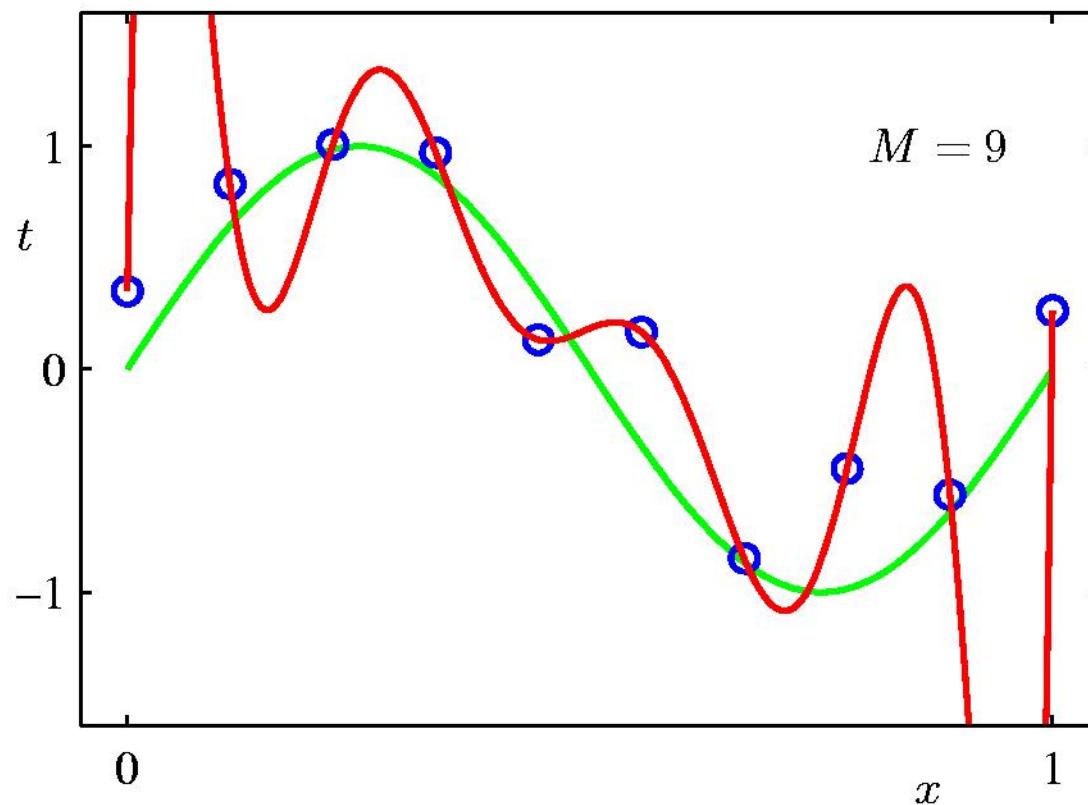
# 1<sup>ST</sup> ORDER POLYNOMIAL



# 3<sup>RD</sup> ORDER POLYNOMIAL



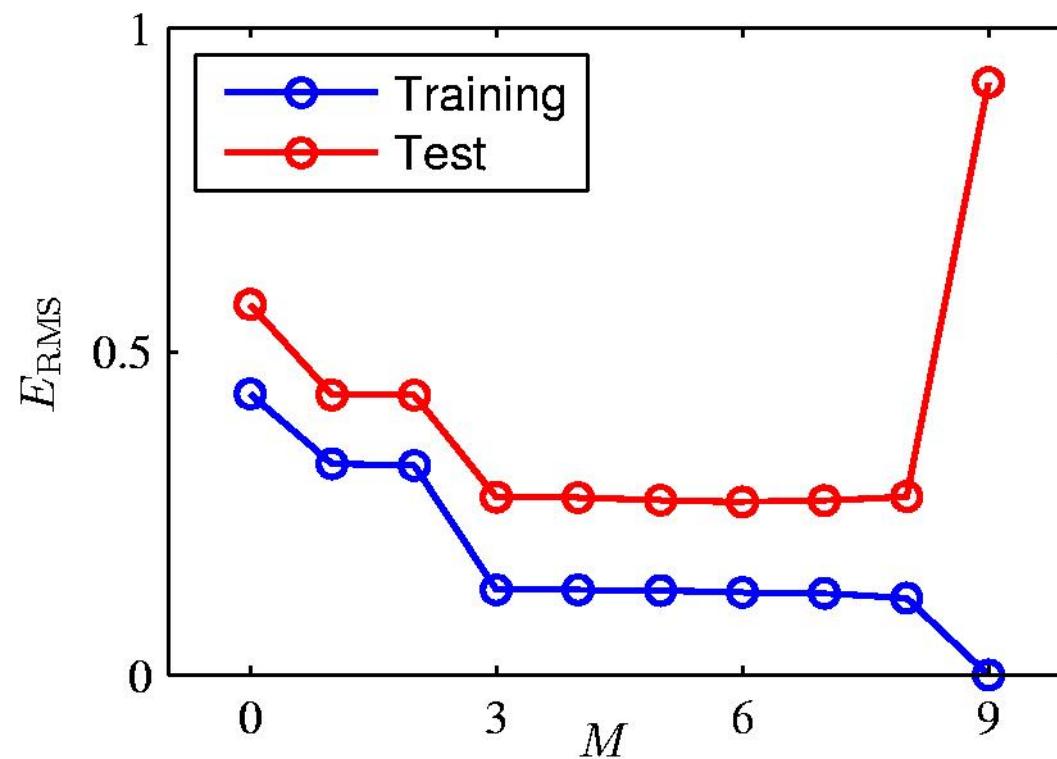
# 9<sup>TH</sup> ORDER POLYNOMIAL



# POLYNOMIAL COEFFICIENTS

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$\theta_0$	0.19	0.82	0.31	0.35
$\theta_1$		-1.27	7.99	232.37
$\theta_2$			-25.43	-5321.83
$\theta_3$			17.37	48568.31
$\theta_4$				-231639.30
$\theta_5$				640042.26
$\theta_6$				-1061800.52
$\theta_7$				1042400.18
$\theta_8$				-557682.99
$\theta_9$				125201.43

# OVER-FITTING





QUESTIONS?