



# Machine Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>

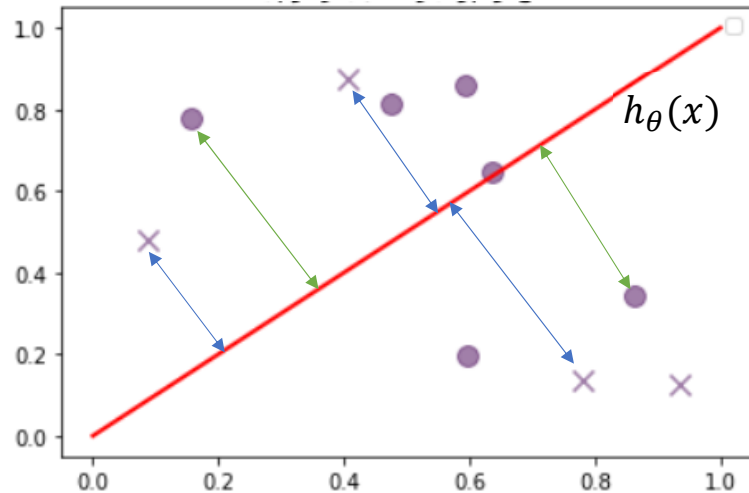


[https://github.com/safayani/machine\\_learning\\_course](https://github.com/safayani/machine_learning_course)



Department of Electrical and computer engineering, Isfahan university of technology, Isfahan, Iran

# Underfitting and Overfitting

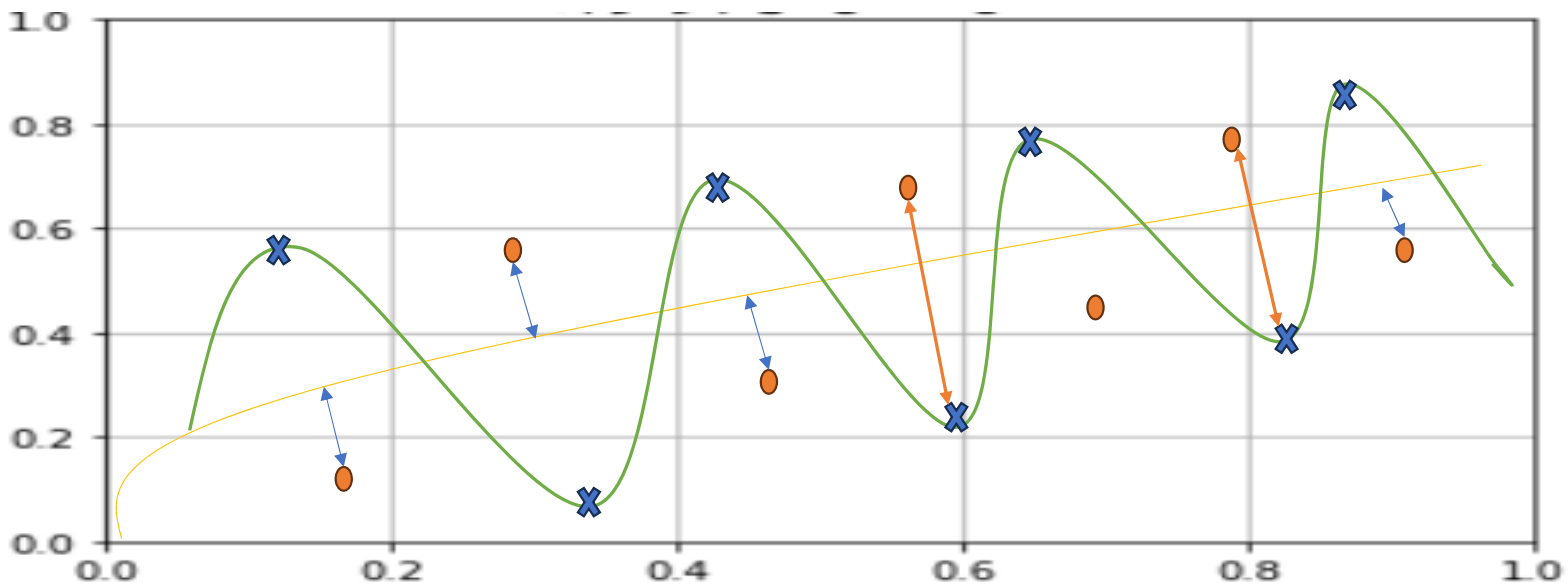


$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

خطای زمان آموزش بسیار زیاد

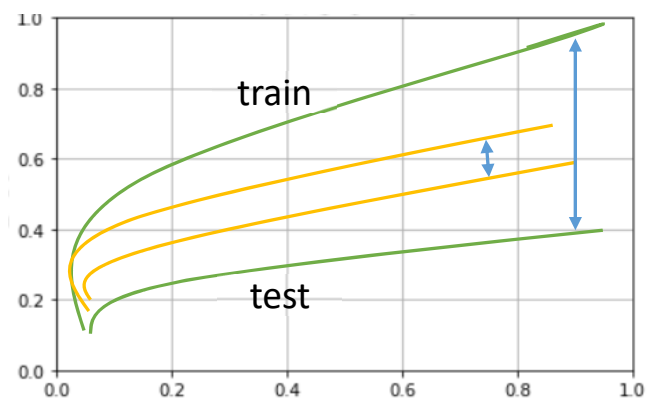
خطای زمان آزمون بسیار زیاد

**underfit**



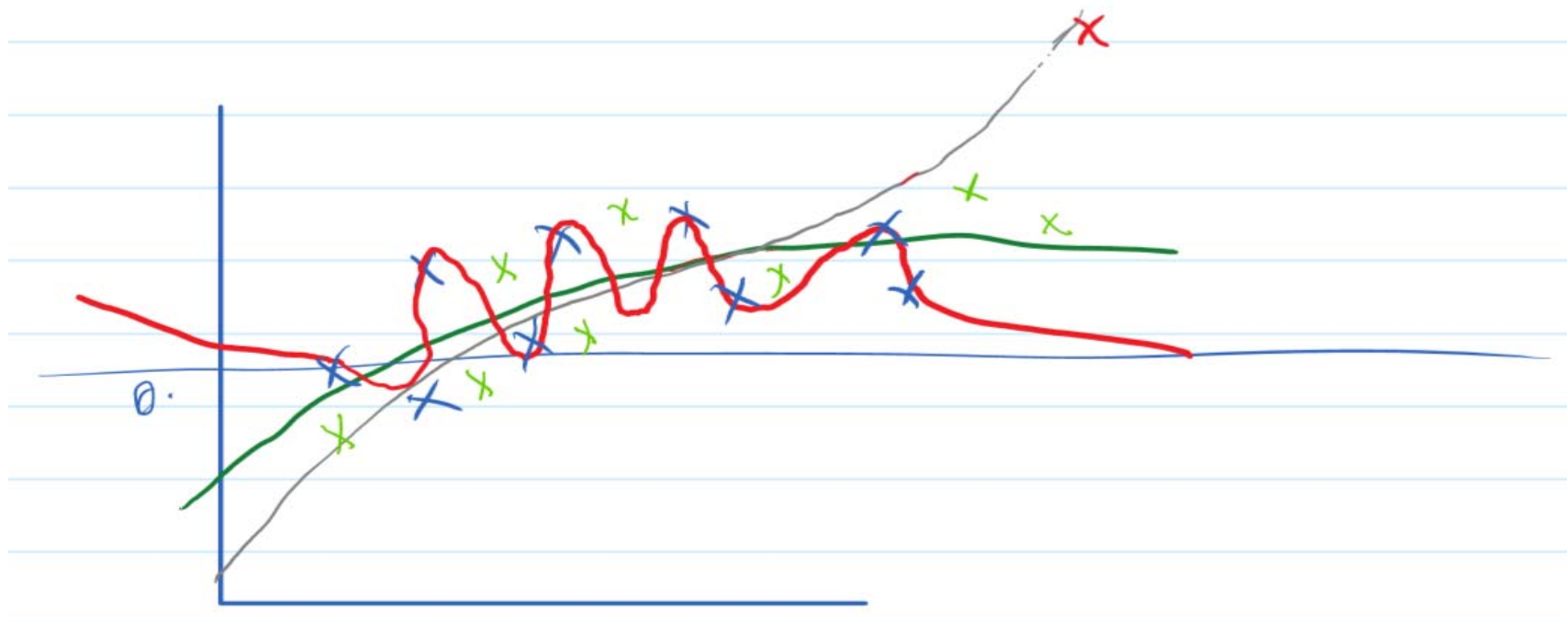
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

خطای آموزش بسیار کم  
خطای آزمون بسیار زیاد  
**overfit**



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

**Just right**



Right :  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

Overfit:  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

# Regularization

$$\text{Minimize}_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

$\theta_3 = 0, \theta_4 = 0$

Features:  $x_1, x_2, \dots, x_{100}$

Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

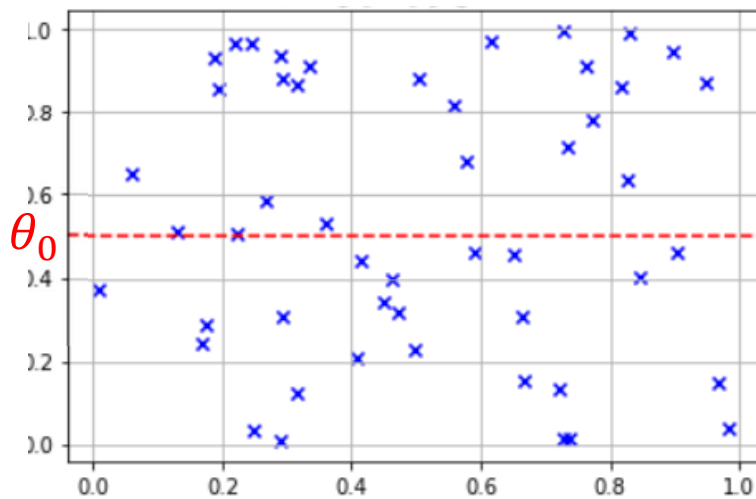
$$\text{Min}_{\theta} J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

ضرب Regularization

Regularization

If  $\lambda$  is very big ( $\lambda = 10^{10}$ )  $\longrightarrow$  underfit ,  $\theta_1, \theta_2, \dots, \theta_n = 0$  ,  $h_{\theta}(x_i) = \theta_0$

If  $\lambda$  is very small ( $\lambda = 0$ )  $\longrightarrow$  overfit



GD:

Repeat until convergence{

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \quad (x_0^i = 1)$$

$$\theta_j = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_j^i + \lambda \theta_j \right]$$

$$j = 1, 2, \dots, n$$

}

Train set: 60%

$$(x^1, y^1), (x^2, y^2), \dots, (x^M, y^M)$$

Cross validation set: 20%

$$(x_{cv}^1, y_{cv}^1), (x_{cv}^2, y_{cv}^2), \dots, (x_{cv}^M, y_{cv}^M)$$

Test set: 20%

$$(x_{test}^1, y_{test}^1), (x_{test}^2, y_{test}^2), \dots, (x_{test}^M, y_{test}^M)$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

$$J_{cv}(\theta) = \frac{1}{2m} \sum_{i=1}^{M_{cv}} (h_{\theta}(x_{cv}^i) - y_i)^2$$

$$J_{test}(\theta) = \frac{1}{2m} \sum_{i=1}^{M_{test}} (h_{\theta}(x_{test}^i) - y_i)^2$$

# Model Selection

$$h_{\theta_1}(x) = \theta_0 + \theta_1 x \quad J_{cv}(\theta^1)$$

$$h_{\theta_2}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \quad J_{cv}(\theta^2)$$

.

.

.

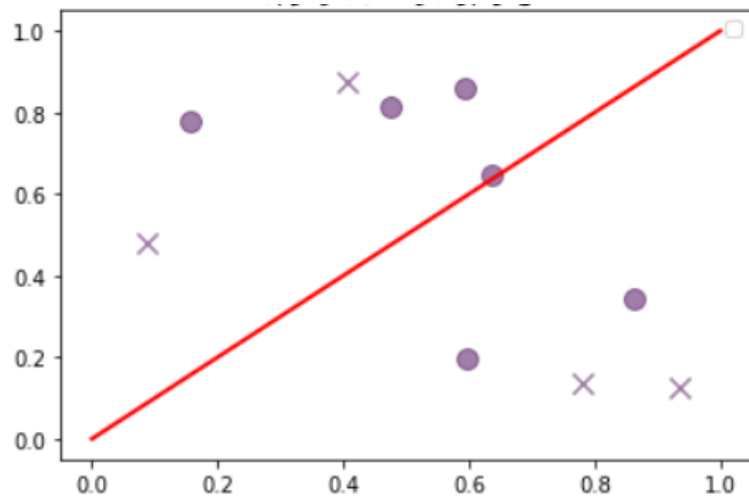
$$h_{\theta_{10}}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \quad J_{cv}(\theta^{10})$$

$$i^* = \underset{i}{\operatorname{argmin}} J_{cv}(\theta^i)$$



$$J_{test}(\theta^{i^*})$$

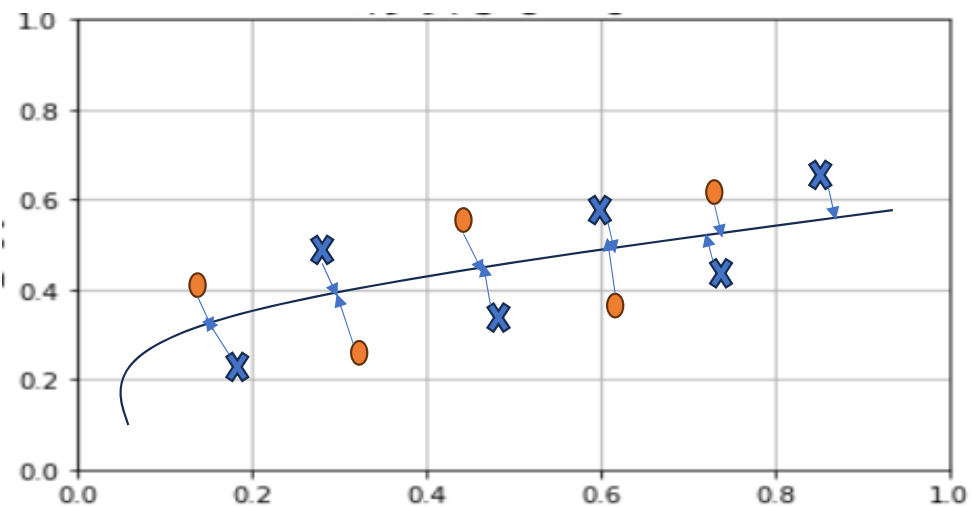




$$\theta_0 + \theta_1 x$$

Underfit  
High bias

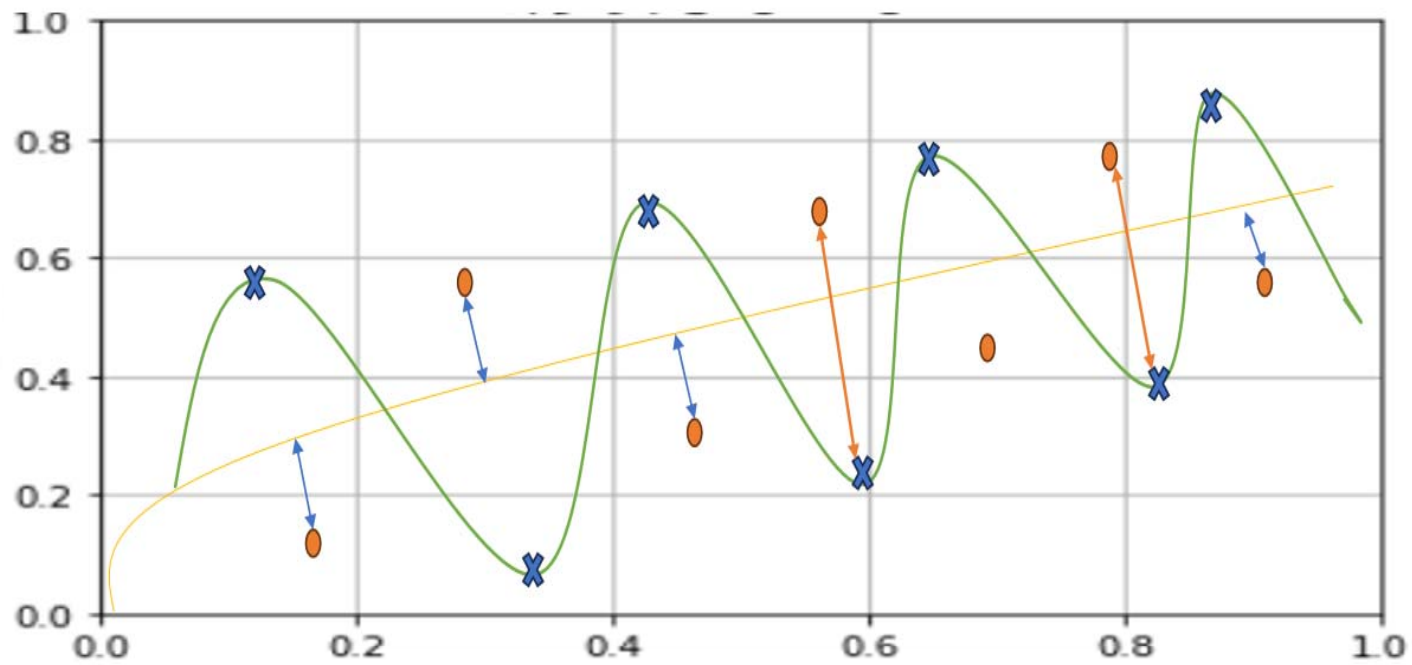
خطای زیاد آموزش  
خطای زیاد CV



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Right

خطای آموزش کم  
خطای CV کم  
نزدیک به یکدیگر

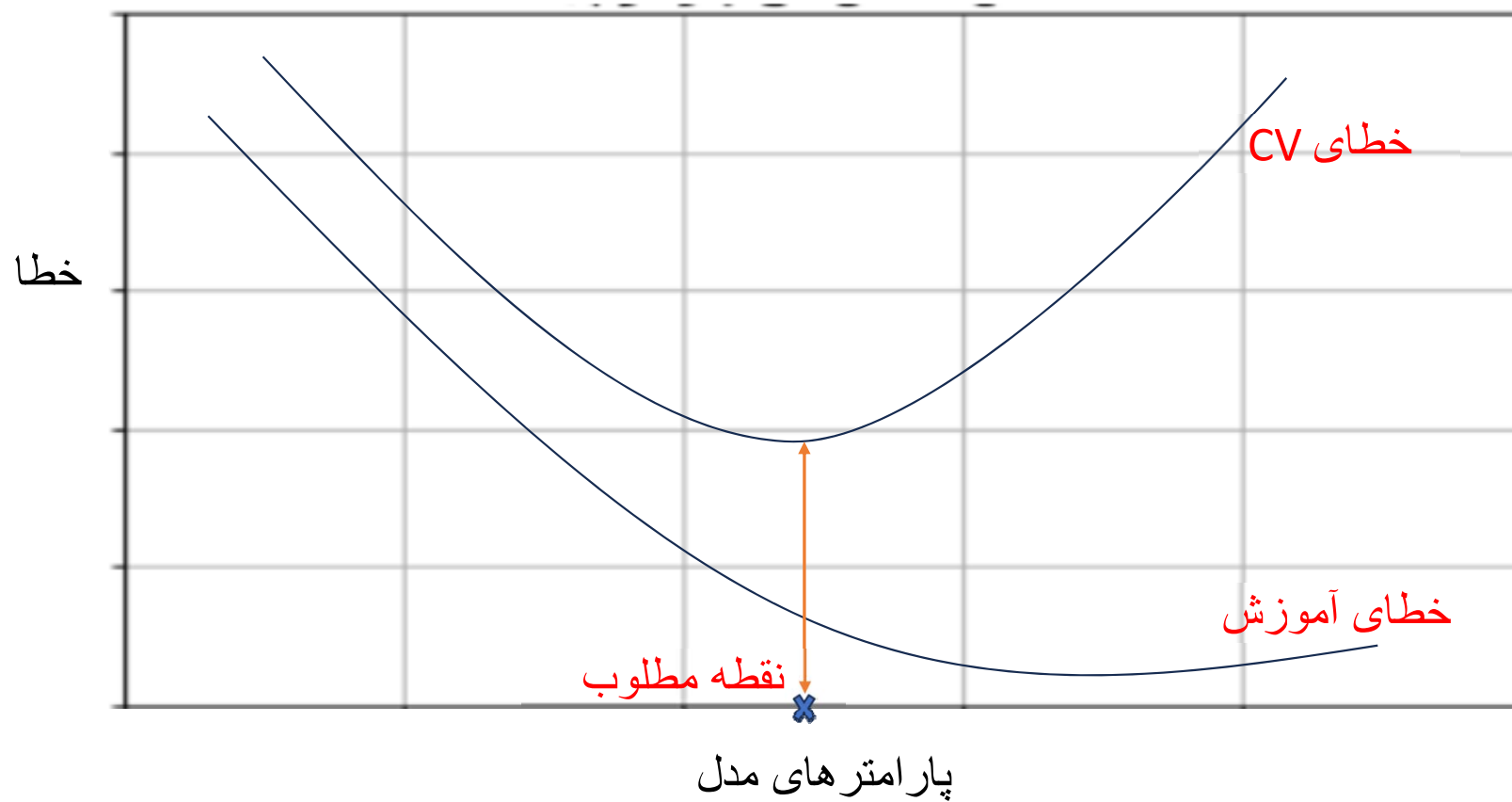


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfit

High variance

خطای آموزشی خیلی کم  
خطای CV زیاد



# Regularization

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

If  $\lambda$  is too large  $\rightarrow \lambda = 10000 \rightarrow$  underfit  $\rightarrow$  high bias  
 $\rightarrow \theta_1 = \theta_2 = \dots = \theta_5 = 0 \rightarrow h(x) = \theta_0$

Appropriate  $\lambda \rightarrow$  appropriate result

If  $\lambda$  is too small  $\rightarrow \lambda = 0 \rightarrow$  overfit  $\rightarrow$  high variance

# Regularization

$$\lambda = 0 \text{ ---> } J_{train}(m^1) \text{ ---> } J_{cv}(m^1)$$

$$\lambda = 0.01 \text{ ---> } J_{train}(m^2) \text{ ---> } J_{cv}(m^2)$$

$$\lambda = 0.02 \text{ ---> } J_{train}(m^3) \text{ ---> } J_{cv}(m^3)$$

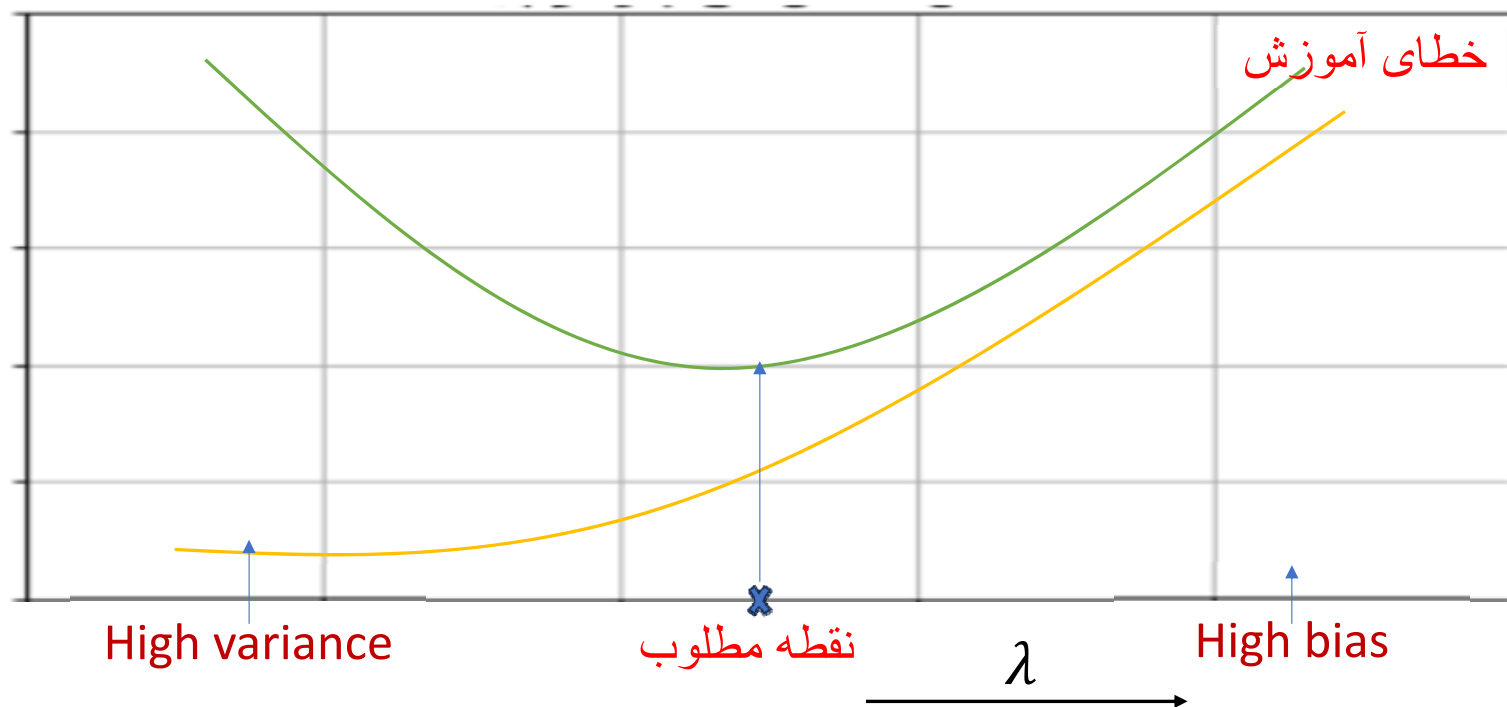
.

.

.

$$\lambda = 10 \text{ ---> } J_{train}(m^{10}) \text{ ---> } J_{cv}(m^{10})$$

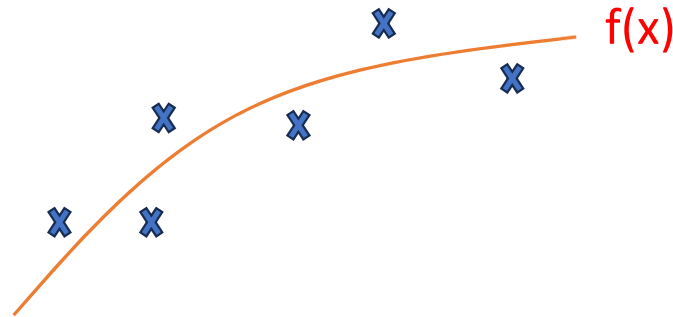
$$i^* = \underset{i}{\operatorname{argmin}} J_{cv}(m^i)$$



## تعریف تئوری بایاس – واریانس

مدل مولد داده:

$$Y = f(x) + \varepsilon$$



$\varepsilon$ : نویز با توزیع  $D_\varepsilon$  که مستقل از داده ها است.

$S_{\text{train}}$ : داده های آموزشی

$D$ : فضای داده ها

## محاسبه رابطه خطا

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[(f(\mathbf{x}) + \varepsilon - f_{S_{\text{train}}}(\mathbf{x}))^2]$$

برای یک نقطه  $\mathbf{x}_0$  خطا به صورت زیر است:

$$(f(\mathbf{x}_0) + \varepsilon - f_{S_{\text{train}}}(\mathbf{x}_0))^2.$$

فرض کنید که با داده های آموزشی مختلفی که از فضای داده  $\mathcal{D}$  نمونه گیری شده اند آزمایش را تکرار میکنیم. در این حالت خطای داده  $\mathbf{x}_0$  به صورت زیر محاسبه می شود:

$$\mathbb{E}_{S_{\text{train}} \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon}[(f(\mathbf{x}_0) + \varepsilon - f_{S_{\text{train}}}(\mathbf{x}_0))^2].$$



## ادامه محاسبه رابطه خطا

می توانیم رابطه بالا را به صورت زیر به دست آوریم:

$$\begin{aligned} & \mathbb{E}_{S_{\text{train}} \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon} [(f(\mathbf{x}_0) + \varepsilon - f_{S_{\text{train}}}(\mathbf{x}_0))^2] \\ & \stackrel{(a)}{=} \mathbb{E}_{\varepsilon \sim \mathcal{D}_\varepsilon} [\varepsilon^2] + \mathbb{E}_{S_{\text{train}} \sim \mathcal{D}} [(f(\mathbf{x}_0) - f_{S_{\text{train}}}(\mathbf{x}_0))^2] \\ & \stackrel{(b)}{=} \text{Var}_{\varepsilon \sim \mathcal{D}_\varepsilon} [\varepsilon] + \mathbb{E}_{S_{\text{train}} \sim \mathcal{D}} [(f(\mathbf{x}_0) - f_{S_{\text{train}}}(\mathbf{x}_0))^2] \\ & \stackrel{(c)}{=} \underbrace{\text{Var}_{\varepsilon \sim \mathcal{D}_\varepsilon} [\varepsilon]}_{\text{noise variance}} \\ & \quad + \underbrace{(f(\mathbf{x}_0) - \mathbb{E}_{S'_{\text{train}} \sim \mathcal{D}} [f_{S'_{\text{train}}}(\mathbf{x}_0)])^2}_{\text{bias}} \\ & \quad + \underbrace{\mathbb{E}_{S_{\text{train}} \sim \mathcal{D}} [(\mathbb{E}_{S'_{\text{train}} \sim \mathcal{D}} [f_{S'_{\text{train}}}(\mathbf{x}_0)] - f_{S_{\text{train}}}(\mathbf{x}_0))^2]}_{\text{variance}}. \end{aligned}$$

## ادامه محاسبه رابطه خطا

توجه کنید در بخش (a) عبارت زیر حذف شده است. چرا؟؟

$$\mathbb{E}_{S_{\text{train}} \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon} [2\varepsilon(f(\mathbf{x}_0) - f_{S_{\text{train}}}(\mathbf{x}_0))].$$

در بخش (b):

$$E_{\varepsilon \sim D_\varepsilon}[\varepsilon^2] = \text{var}_{\varepsilon \sim D_\varepsilon}[\varepsilon]$$

در بخش (c):

عبارت  $\mathbb{E}_{S'_{\text{train}} \sim \mathcal{D}}[f_{S'_{\text{train}}}(\mathbf{x}_0)]$  که  $S'$  یک مجموعه داده از  $\mathcal{D}$  است) را به رابطه اضافه و کم می کنیم و سپس توان ۲ را اعمال می کنیم. در این رابطه یک ترم سوم هم وجود دارد که نشان می دهیم که به صورت زیر برابر با صفر است:

ادامه محاسبه رابطه خطا

$$\begin{aligned} & \mathbb{E}_{S \sim \mathcal{D}} \left[ \left( f(\mathbf{x}_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] \right) \cdot \left( \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] - f_S(\mathbf{x}_0) \right) \right] \\ &= \left( f(\mathbf{x}_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] \right) \cdot \mathbb{E}_{S \sim \mathcal{D}} \left[ \left( \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] - f_S(\mathbf{x}_0) \right) \right] \\ &= \left( f(\mathbf{x}_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] \right) \cdot \left( \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(\mathbf{x}_0)] - \mathbb{E}_{S \sim \mathcal{D}}[f_S(\mathbf{x}_0)] \right) \\ &= 0. \end{aligned}$$

## ادامه محاسبه رابطه خطا

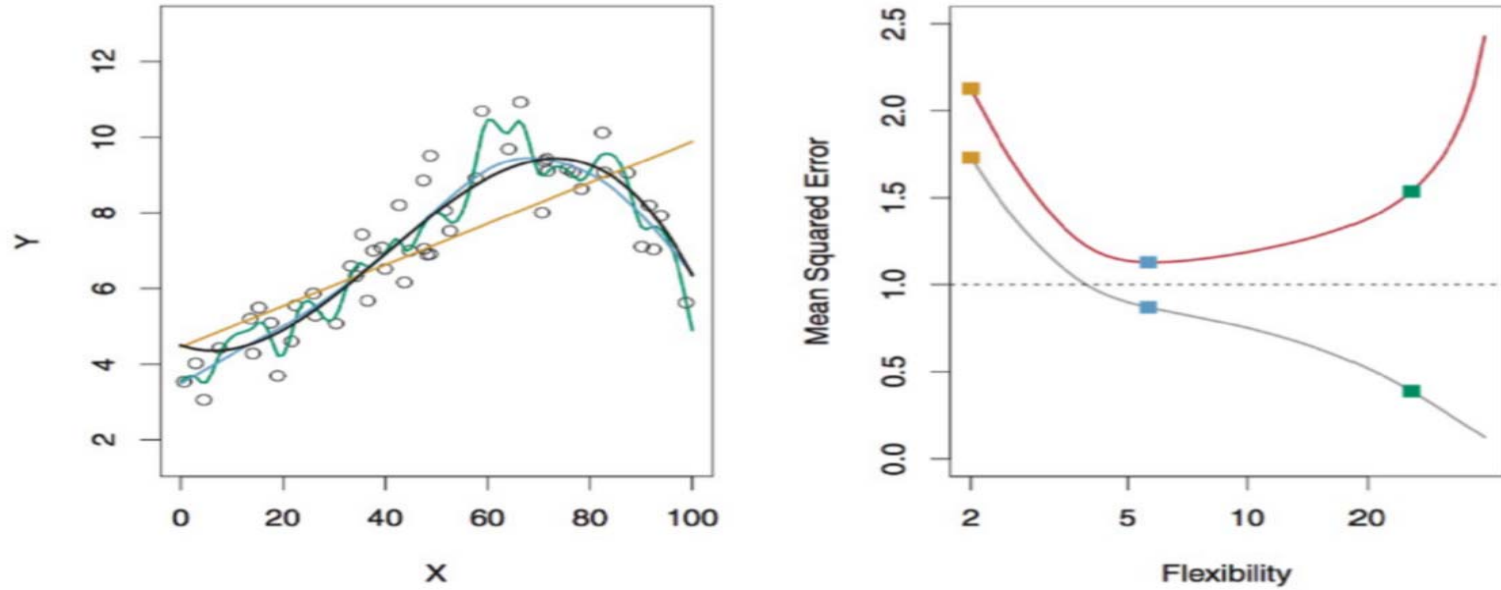
تعبیر رابطه (c):

از سه ترم مثبت تشکیل شده است. ترم اول ربطی به نحوه آموزش مدل ندارد و ناشی از عدم قطعیت ذاتی در داده ها است.

بایاس تفاضل مابین مقدار واقعی  $f(x_0)$  و متوسط مدل های مختلفی است که بر روی داده ها آموزش دیده اند. (مدل های ساده نمی توانند خوب بر روی داده ها تطبیق یابند. در نتیجه بایاس زیاد می شود.)

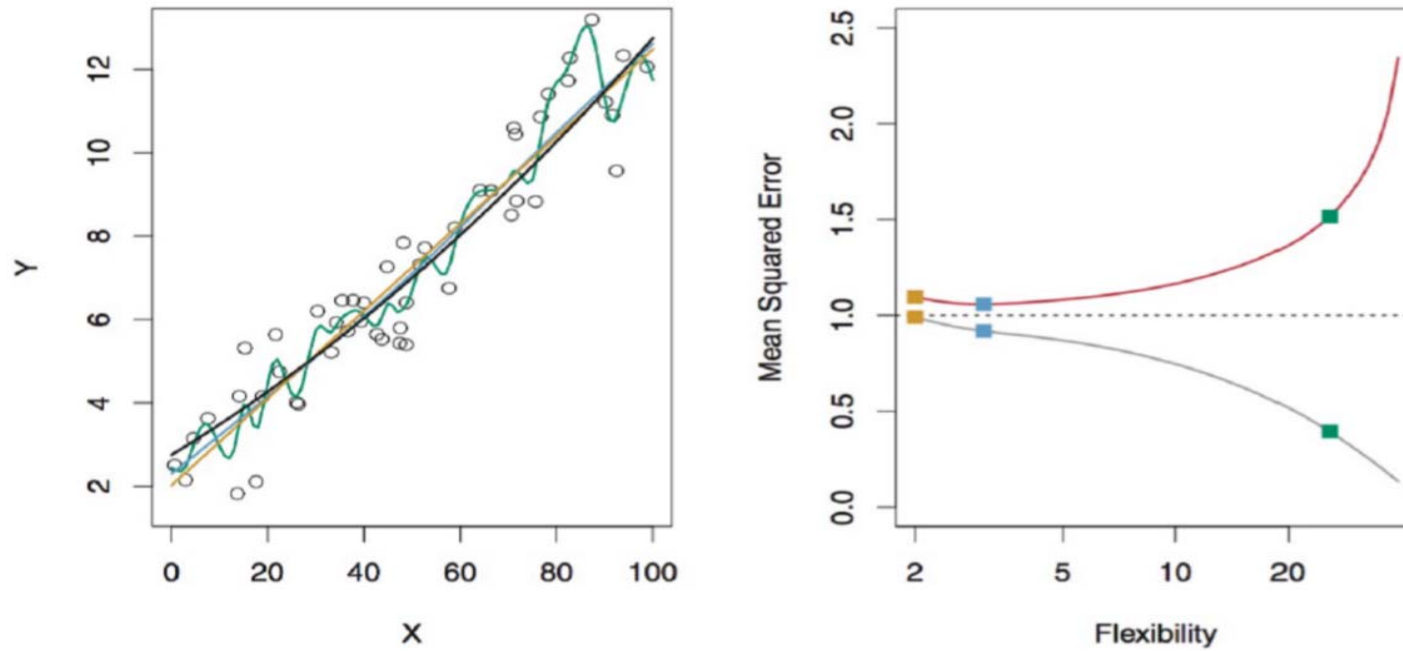
ترم واریانس در واقع واریانس مدل های مختلفی است که آموزش دیده اند. اگر مدل ما خیلی پیچیده باشد با تغییر اندکی در داده ها شکل مدل عوض می شود و پیش بینی بر روی  $x_0$  به میزان زیادی متغیر می شود.

## چند مثال



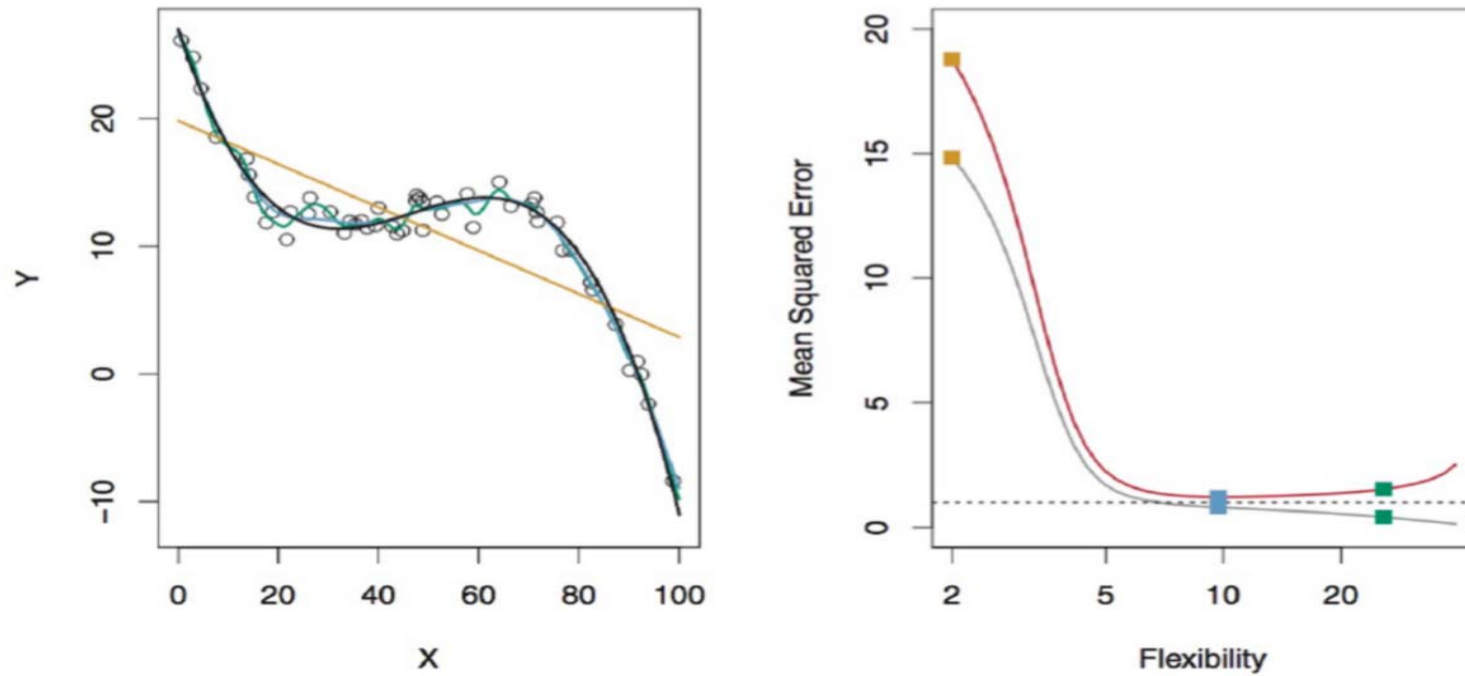
**FIGURE 2.9.** Left: Data simulated from  $f$ , shown in black. Three estimates of  $f$  are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

## چند مثال



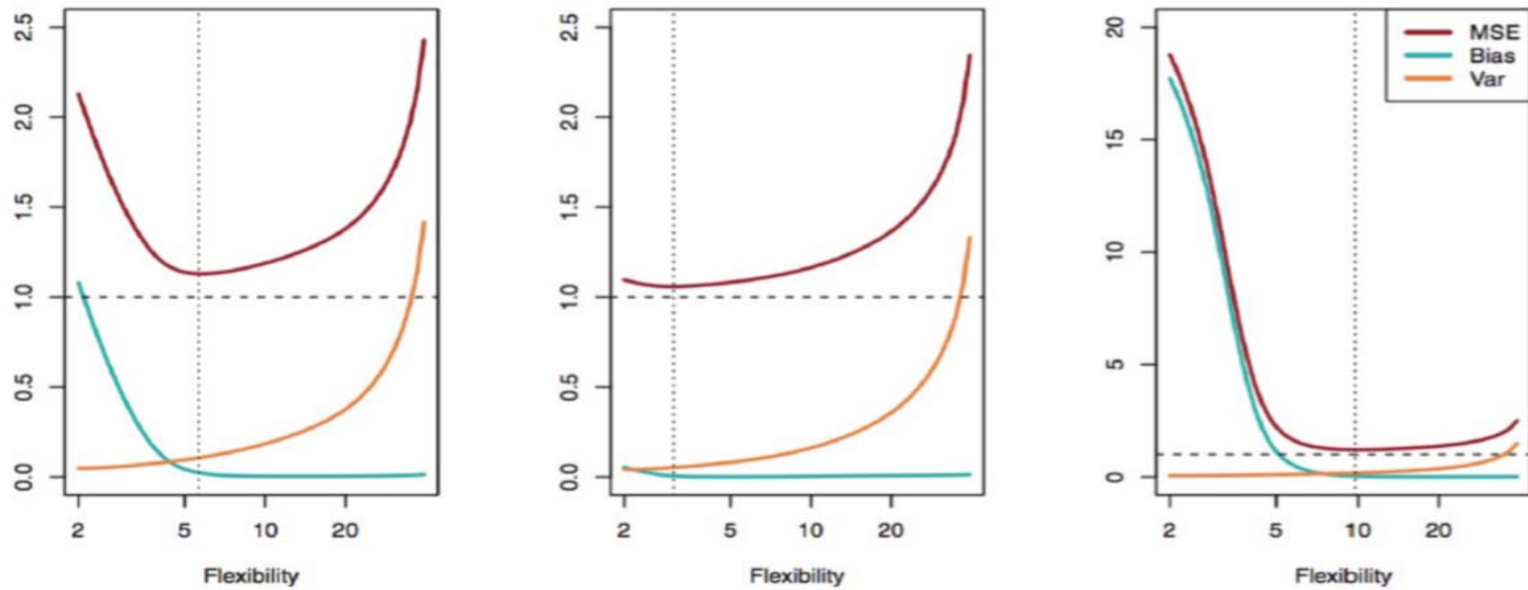
**FIGURE 2.10.** Details are as in Figure 2.9, using a different true  $f$  that is much closer to linear. In this setting, linear regression provides a very good fit to the data.

## چند مثال



**FIGURE 2.11.** Details are as in Figure 2.9, using a different  $f$  that is far from linear. In this setting, linear regression provides a very poor fit to the data.

## چند مثال



**FIGURE 2.12.** Squared bias (blue curve), variance (orange curve),  $\text{Var}(\epsilon)$  (dashed line), and test MSE (red curve) for the three data sets in Figures 2.9–2.11. The vertical dotted line indicates the flexibility level corresponding to the smallest test MSE.



# K fold CV

train

