



Machine Learning

Model Performance: Key Concepts in Machine Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



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



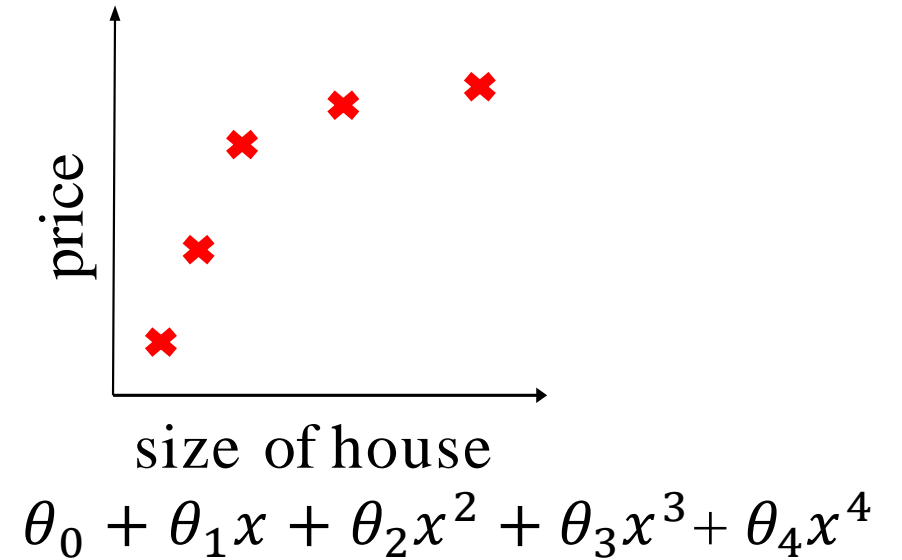
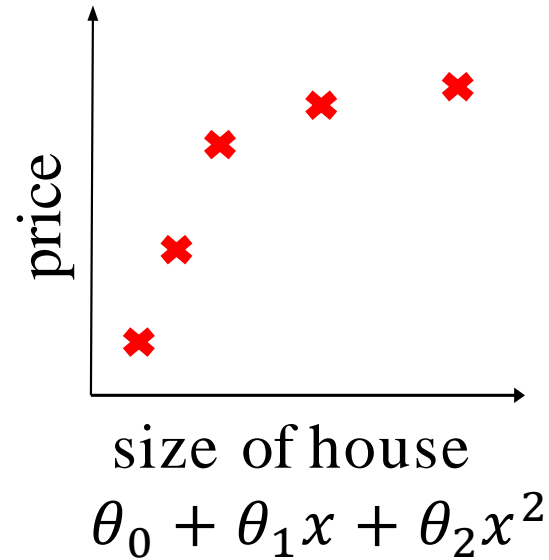
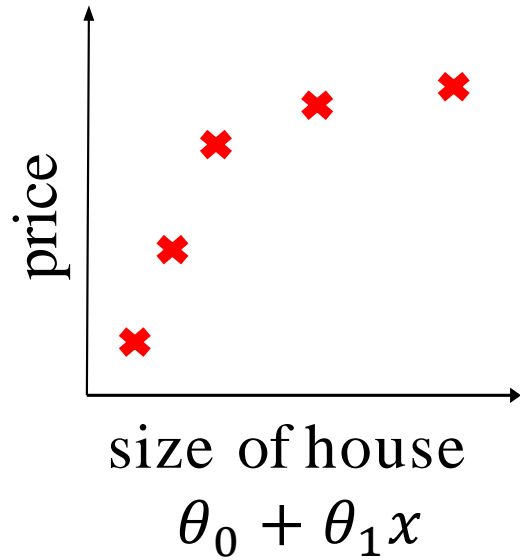
https://github.com/safayani/machine_learning_course



Contents

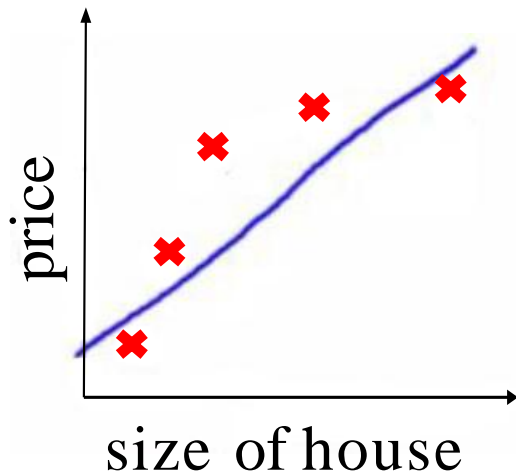
- Underfitting, Overfitting
- Model Selection
- Train-validation test split
- Regularization: ridge and lasso regression
- Bias-Variance Tradeoff
- K-fold Cross validation

Example: Linear regression (housing prices)

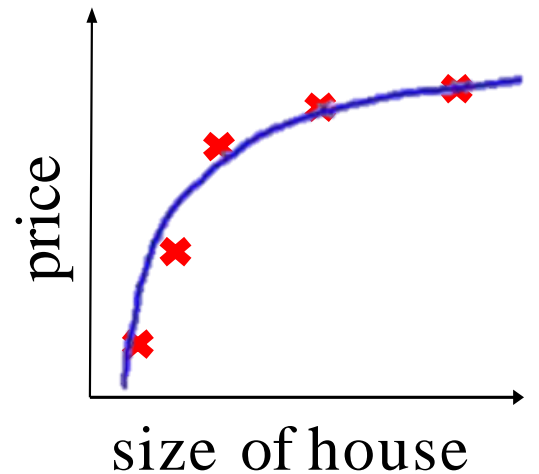


The slides are modified, based on original slides by [Andrew NG, Stanford university]

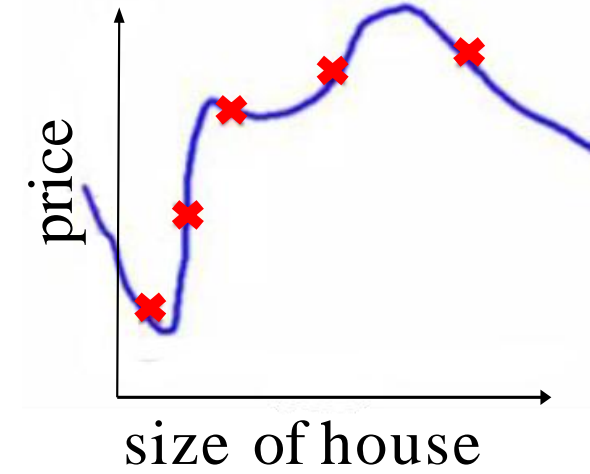
Example: Linear regression (housing prices)



$\theta_0 + \theta_1 x$
"Underfit" "High bias"



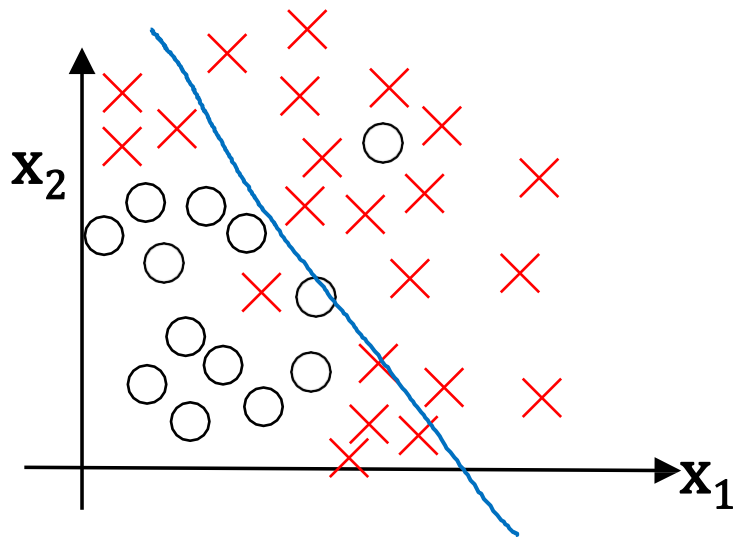
$\theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"



$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"overfit" "High variance"

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

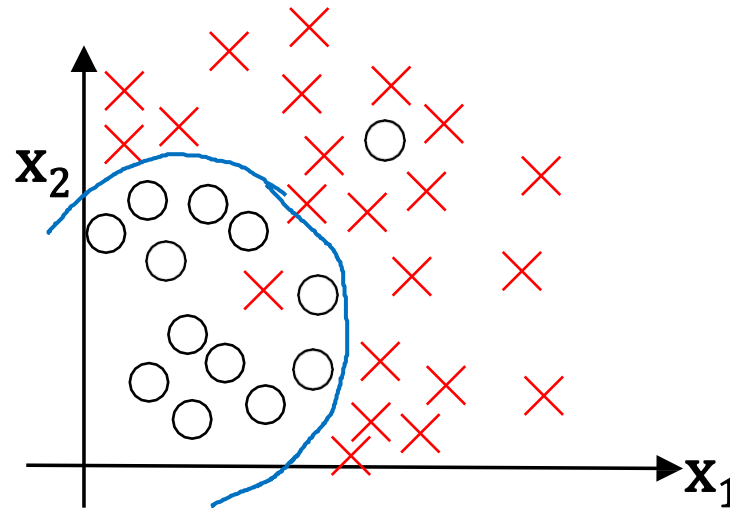
Example: Logistic regression



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

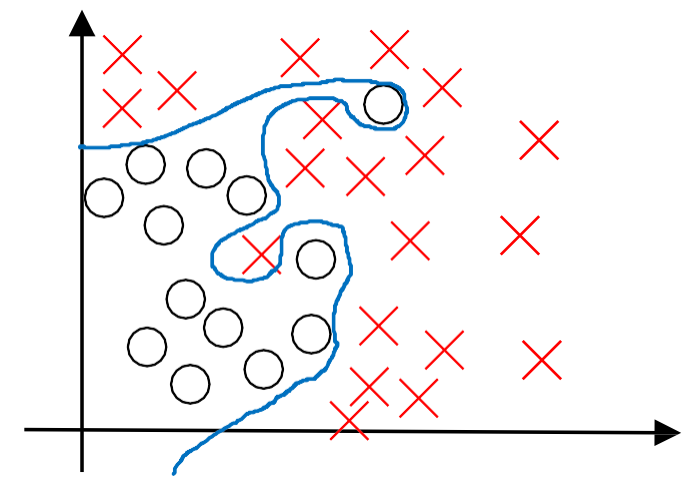
high bias

underfitting



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2) \quad g\left(\begin{array}{l} \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 \\ + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots \end{array}\right)$$

“just right”



high variance

overfitting

Addressing overfitting:

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

x_4 = age of house

x_5 = average income in neighborhood

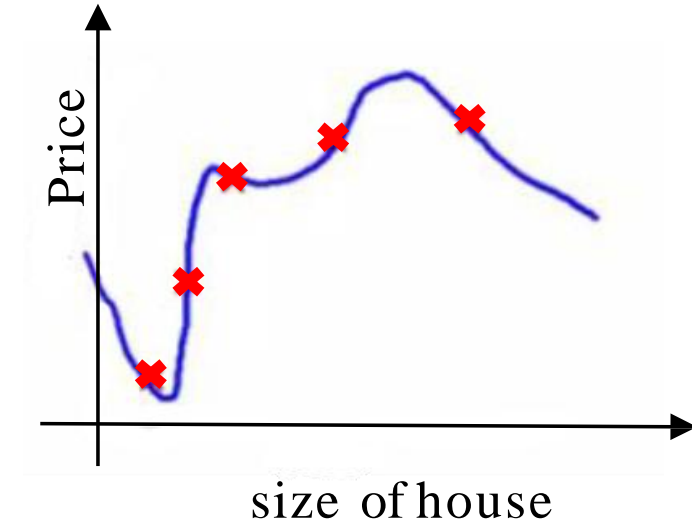
x_6 = kitchen size

.

.

.

x_{100}



Evaluating your hypothesis

- Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243

60%
Training set

20%
Cross validation set (cv)

20%
Test set

$(x^{(1)}, y^{(1)})$

$(x^{(2)}, y^{(2)})$

.

.

.

$(x^{(m)}, y^{(m)})$

$(x_{cv}^{(1)}, y_{cv}^{(1)})$

$(x_{cv}^{(2)}, y_{cv}^{(2)})$

.

.

.

$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$(x_{test}^{(1)}, y_{test}^{(1)})$

$(x_{test}^{(2)}, y_{test}^{(2)})$

.

.

.

$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Train/validation/test error

- Training error:

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

- Cross Validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

- Test error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Model Selection

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

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

.

.

.

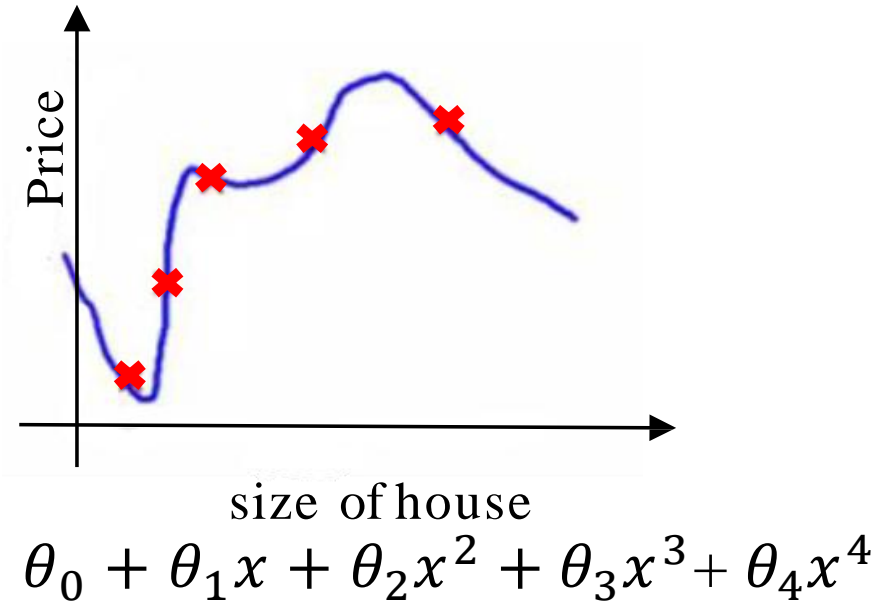
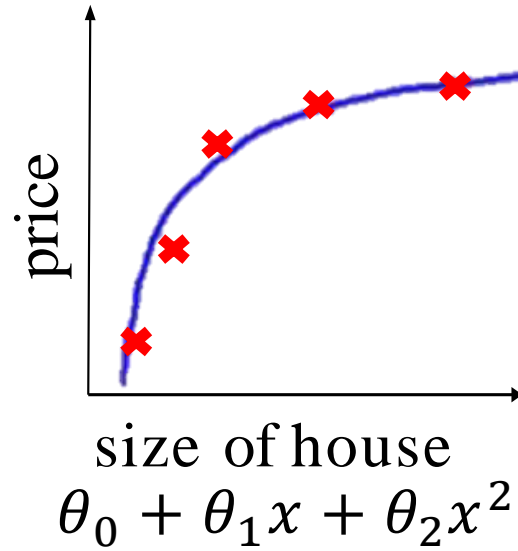
$$h_{\theta_{10}}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \quad J_{cv}(\theta^{10})$$

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



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

Regularization Intuition



- Suppose we penalize and make θ_3, θ_4 really small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \underbrace{1000\theta_3^2}_{\theta_3 \approx 0} + \underbrace{1000\theta_4^2}_{\theta_4 \approx 0}$$

Regularization

- Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$
 - "Simpler" hypothesis
 - Less prone to overfitting
- Housing:
 - Features: x_1, x_2, \dots, x_{100}
 - Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

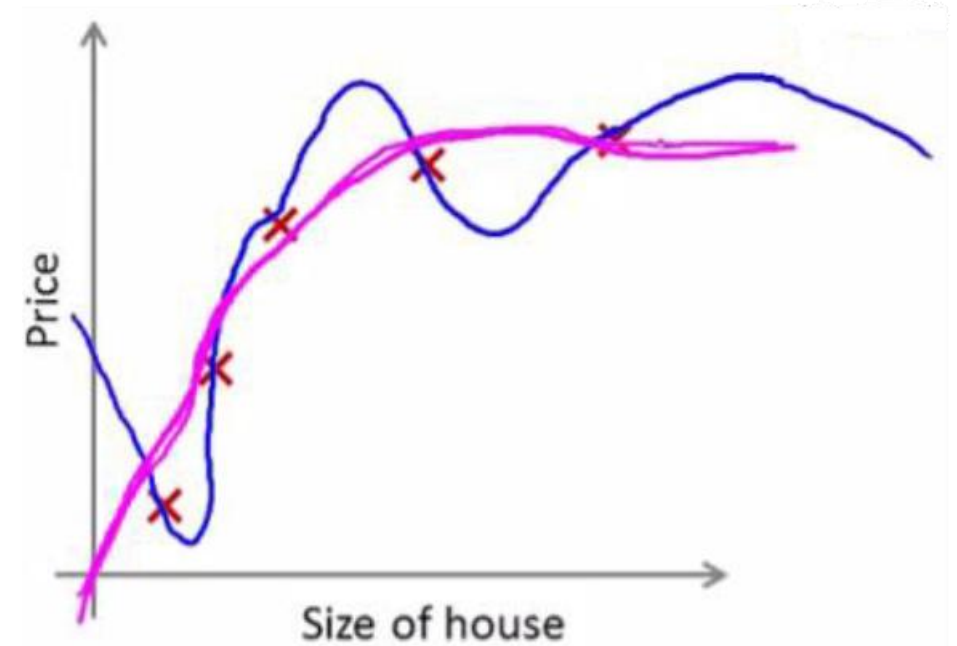
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

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

λ = regularization parameter
lambda lambda

Regularization

- $J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$
 $\min_{\theta} J(\theta)$

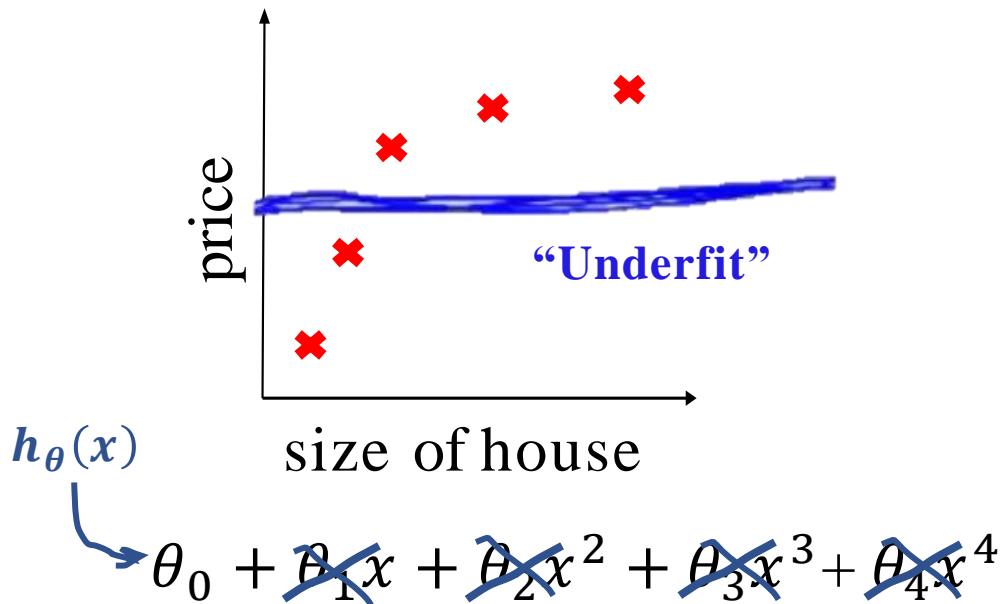


Regularization: Ridge Regression

- In regularized linear regression, we choose θ to minimize

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

- What if λ is set to an extremely large value (perhaps far too large for our problem, say $\lambda = 1010$)?



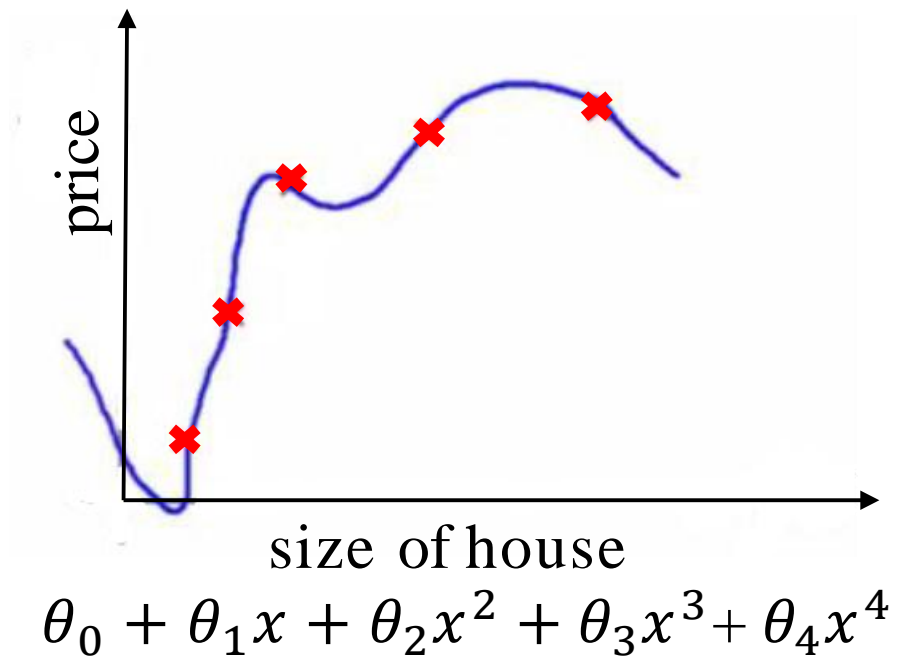
$$\theta_1, \theta_2, \theta_3, \theta_4$$

$$\theta_1 \approx 0, \theta_2 \approx 0$$

$$\theta_3 \approx 0, \theta_4 \approx 0$$

$$\boxed{h_{\theta}(x) = \theta_0}$$

Evaluating your hypothesis



- Fails to generalize to new examples not in training set.

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

.

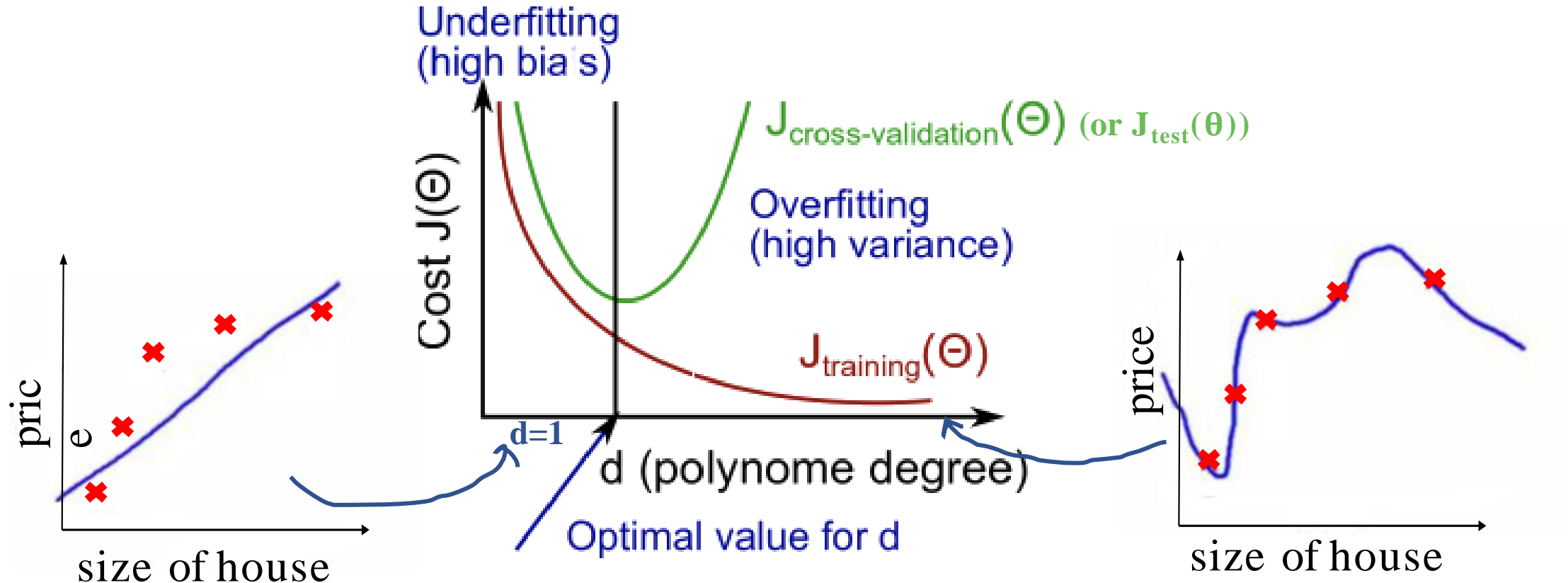
.

.

x_{100}

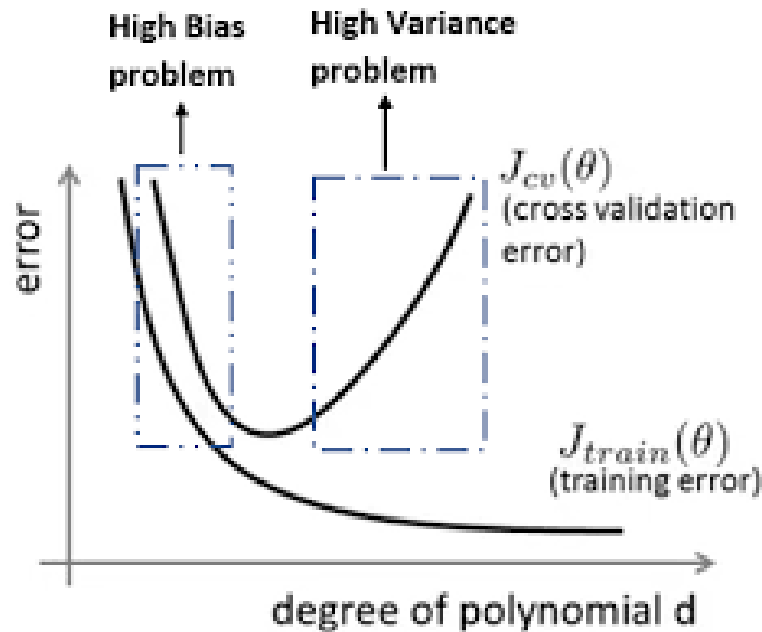
Bias/variance

- Training error: $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (or $J_{test}(\theta)$)



Diagnosing bias vs. variance

- Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$) is high.) Is it a bias problem or a variance problem?



<https://towardsdatascience.com/>

Bias (underfit):

$J_{train}(\theta)$ will be high

$$J_{cv}(\theta) \approx J_{train}(\theta)$$

Variance (overfit):

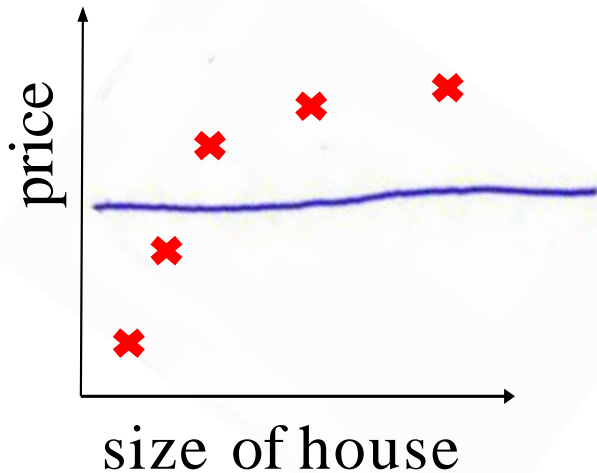
$J_{train}(\theta)$ will be low

$$J_{cv}(\theta) \gg J_{train}(\theta)$$

Linear regression with regularization

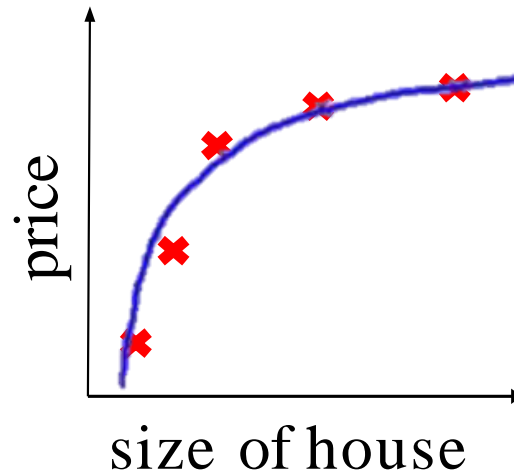
Model: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

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

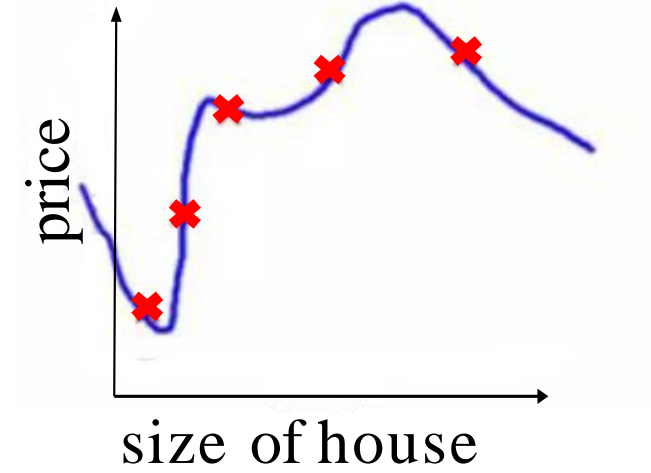


Large λ
High bias(underfit)

$$\lambda = 10000. \theta_1 \approx 0, \theta_2 \approx 0, \dots$$
$$h_{\theta}(x) \approx \theta_0$$



Intermediate λ
"Just right"



Small λ
High variance (overfit)

$$\lambda = 0$$

Choosing the regularization parameter λ

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

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

1. Try $\lambda = 0 \rightarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. Try $\lambda = 0.01 \rightarrow \min J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. Try $\lambda = 0.02 \rightarrow \min J(\theta) \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
4. Try $\lambda = 0.04$
5. Try $\lambda = 0.08 \rightarrow \min J(\theta) \rightarrow \theta^{(5)} \rightarrow J_{cv}(\theta^{(5)})$
- .
- .
- .
12. Try $\lambda = 10$ $\rightarrow \min J(\theta) \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
10.24 Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$

Bias/variance as a function of the regularization parameter λ

- **Training error:**

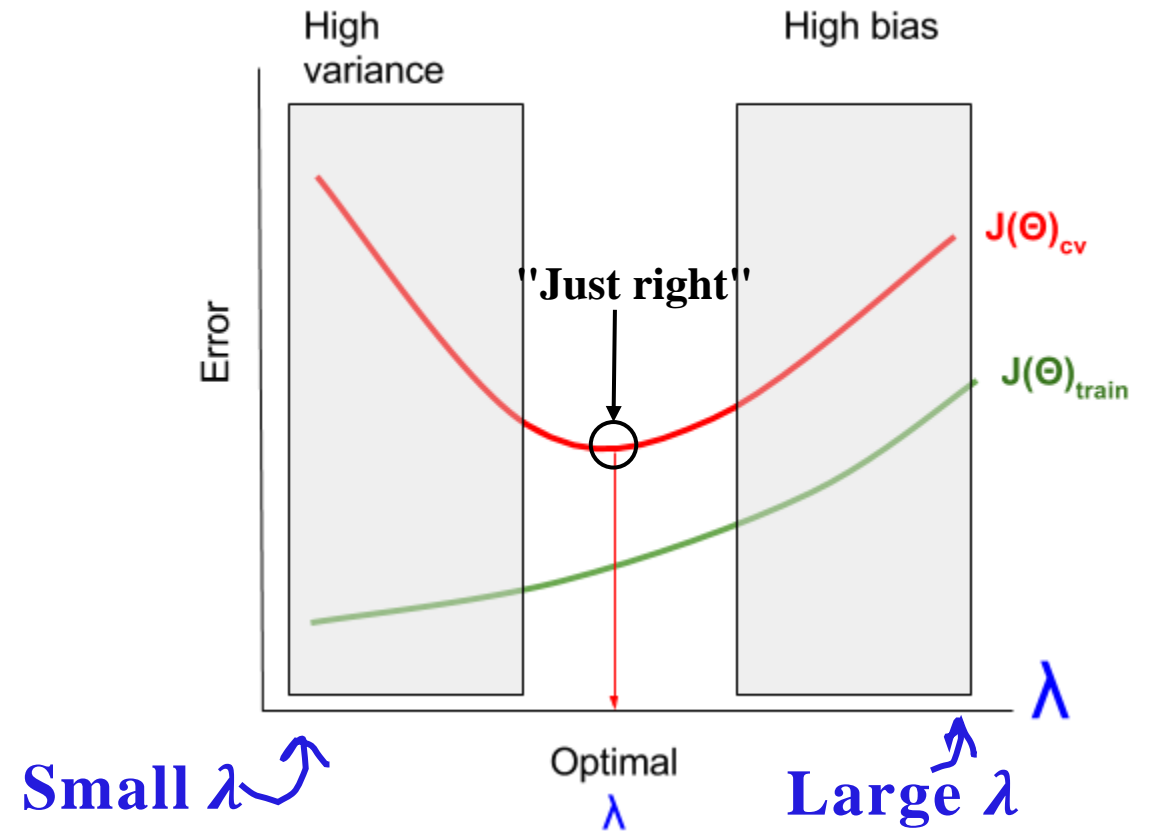
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

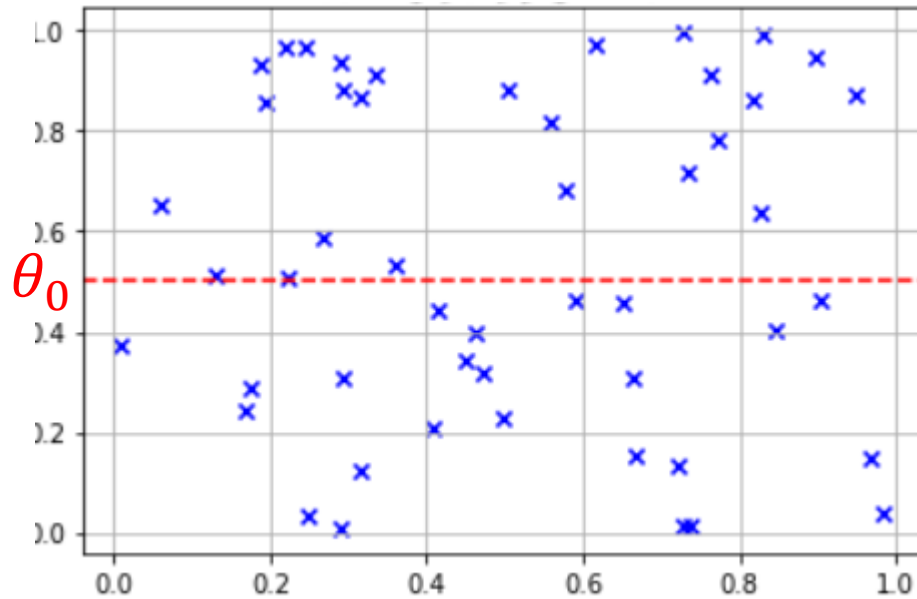
- **Cross Validation error:**

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

- **Test error:**

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$





GD:

Repeat until convergence{

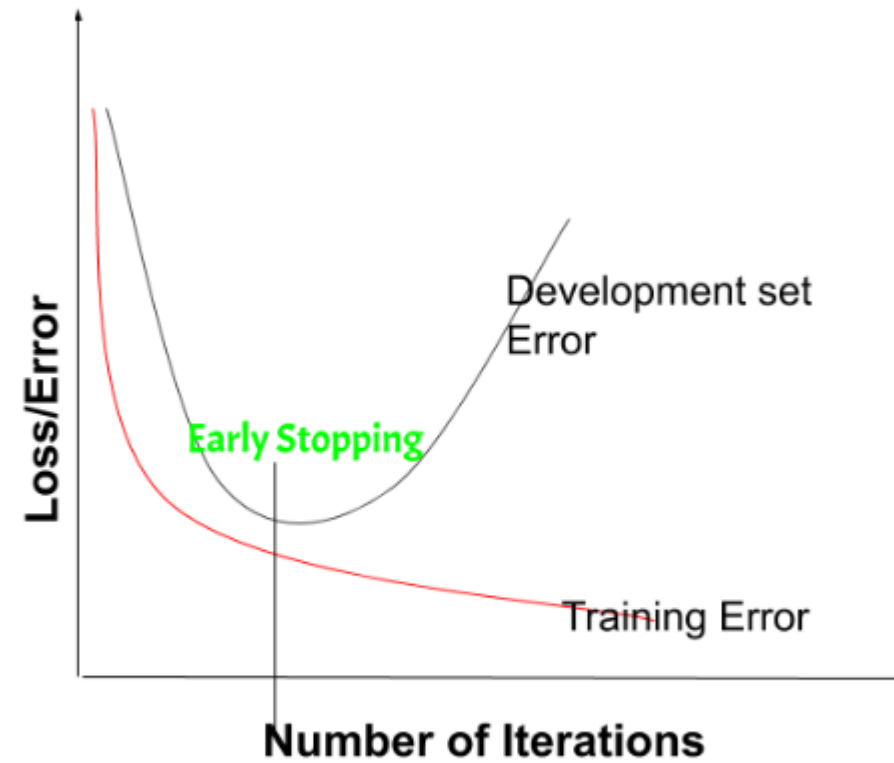
$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \quad (x_0^i = 1)$$

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

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

}

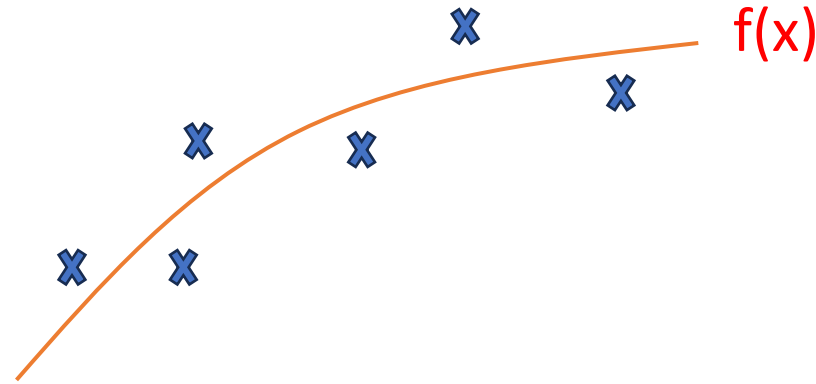
Early stopping



تعریف تئوری بایاس – واریانس

مدل مولد داده:

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



ε : نویز با توزیع D_ε که مستقل از داده ها است.

S_{train} : داده های آموزشی

D : فضای داده ها

محاسبه رابطه خطا

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

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

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

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:

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

ادامه محاسبه رابطه خطا

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

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

ادامه محاسبه رابطه خطا

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

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

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

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

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.

Examples

- $f_S(x) = k$

$$\underbrace{\left(f(\mathbf{x}_0) - \mathbb{E}_{s'_{\text{train}} \sim \mathcal{D}} \left[f_{s'_{\text{train}}}(\mathbf{x}_0)\right]\right)^2}_{\text{bias}}$$

$$\underbrace{\mathbb{E}_{s_{\text{train}} \sim \mathcal{D}} \left[\left(\mathbb{E}_{s'_{\text{train}} \sim \mathcal{D}} \left[f_{s'_{\text{train}}}(\mathbf{x}_0) \right] - f_{s_{\text{train}}}(\mathbf{x}_0) \right)^2 \right]}_{\text{variance}}.$$

Prove it

Bias is high

Variance is 0

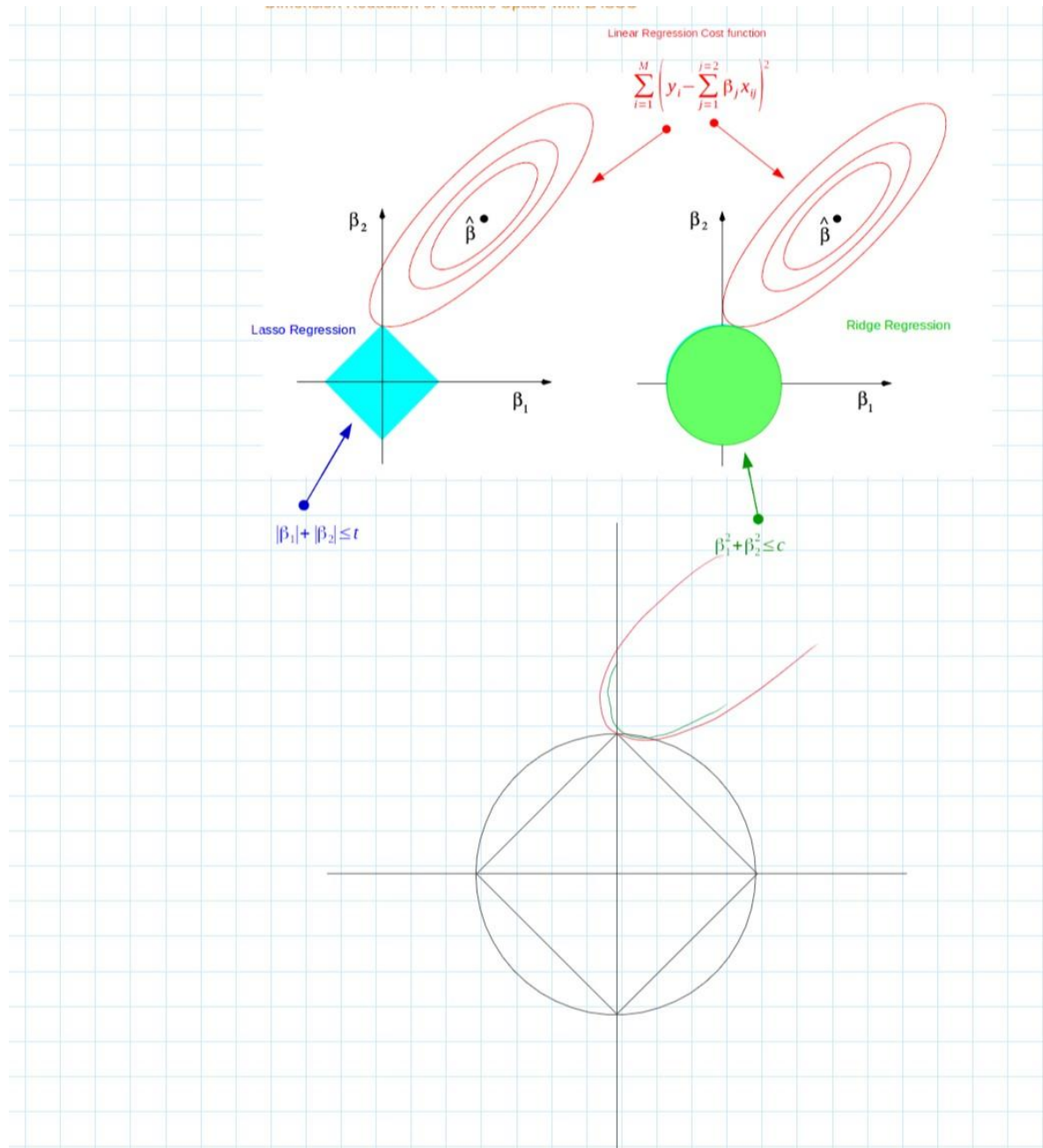
$$f_S(\mathbf{x}) = f(\mathbf{x}) + \varepsilon$$

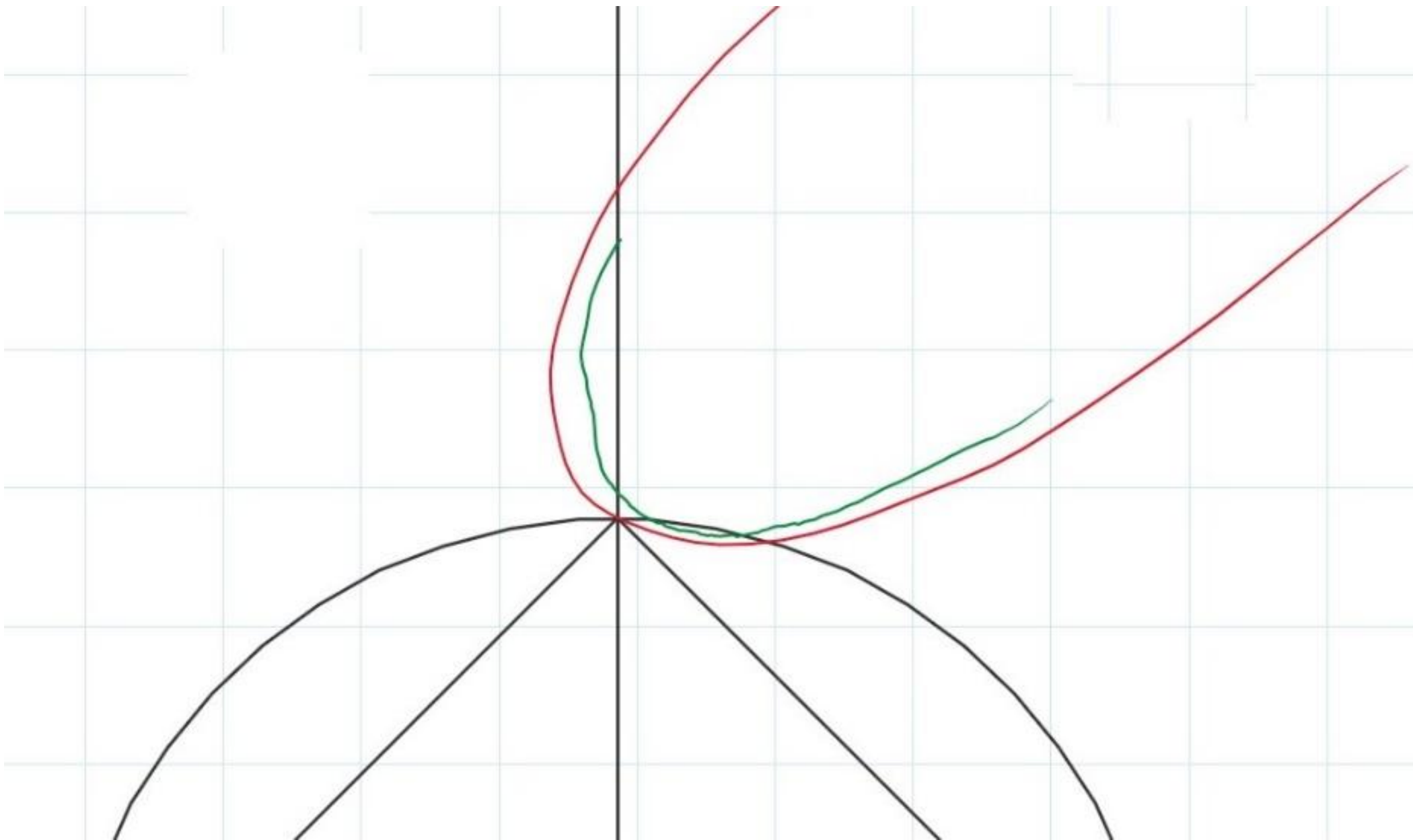
Bias is zero

Variance is $\text{var}_{\varepsilon \sim D_\varepsilon}[\varepsilon]$

Lasso Regression

- $J(\theta) = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j|$





چند مثال

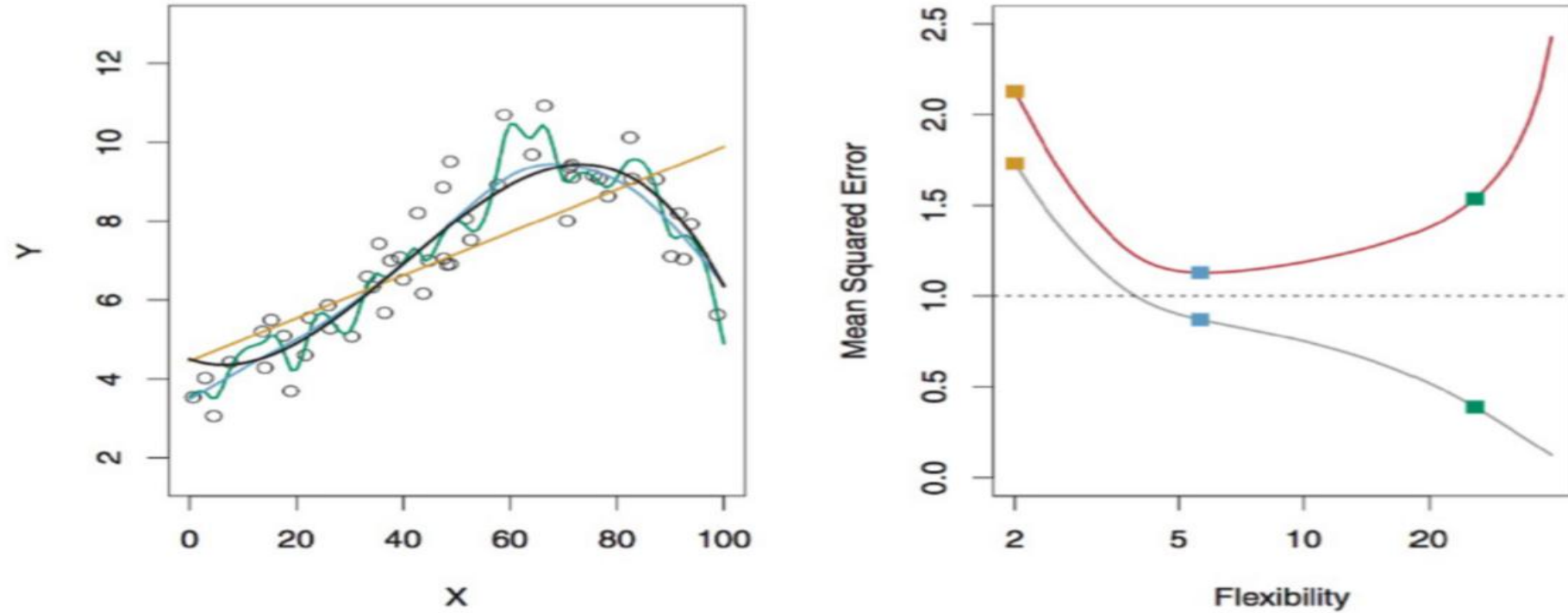


FIGURE 2.9. Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line).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

چند مثال

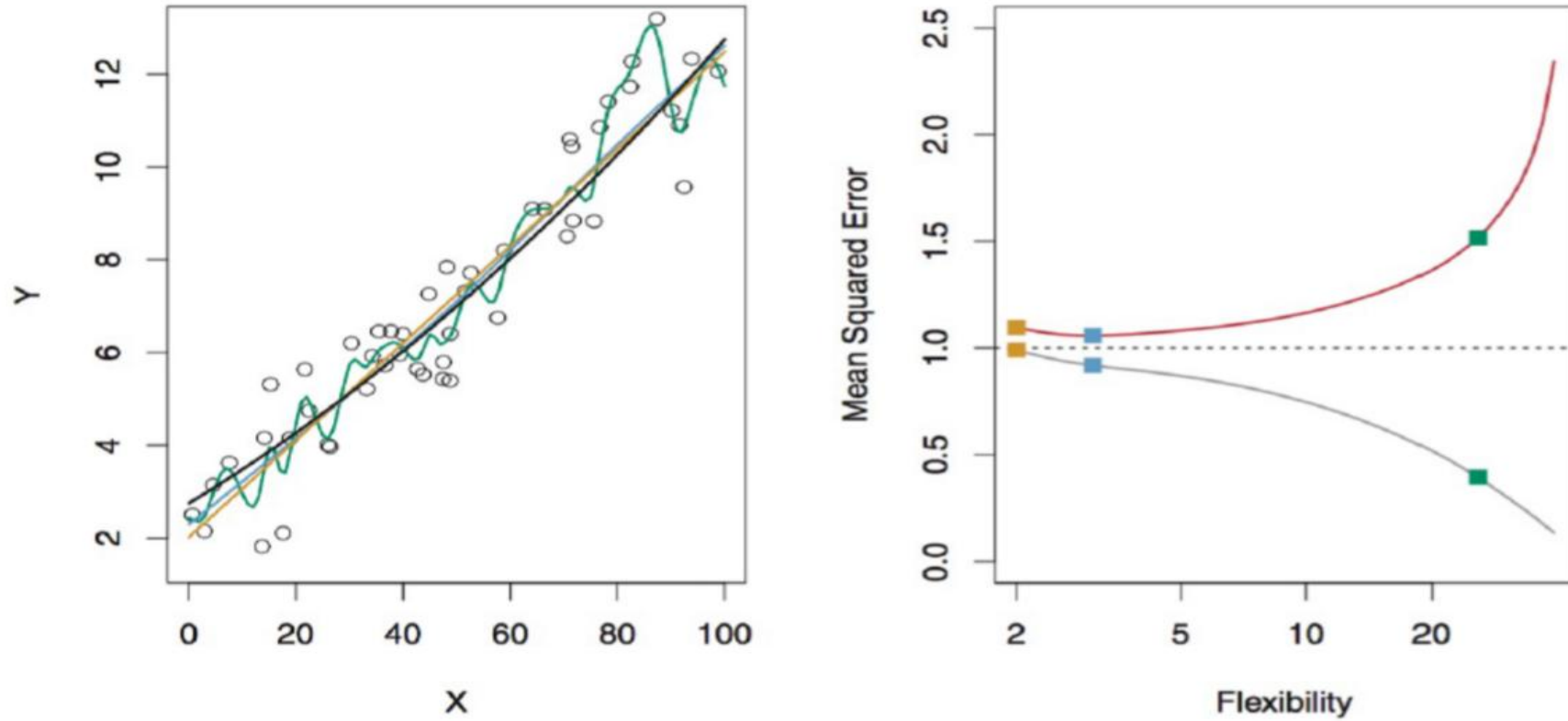


FIGURE 2.10. Details are as in Figure 2.9, using a different true f that is much closer to linear. In this setting, linear regression provides a very good fit to the data.

چند مثال

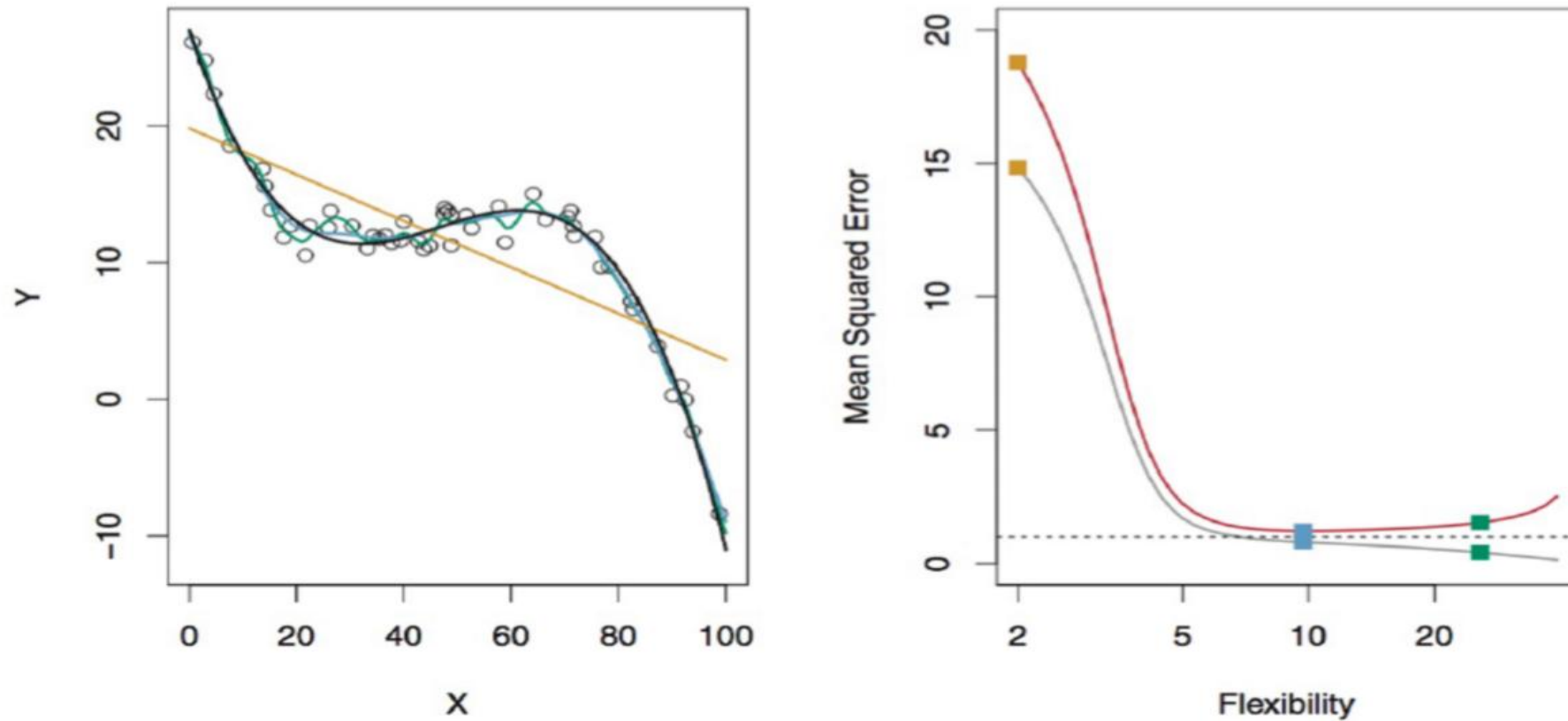


FIGURE 2.11. Details are as in Figure 2.9, using a different f that is far from linear. In this setting, linear regression provides a very poor fit to the data.

چند مثال

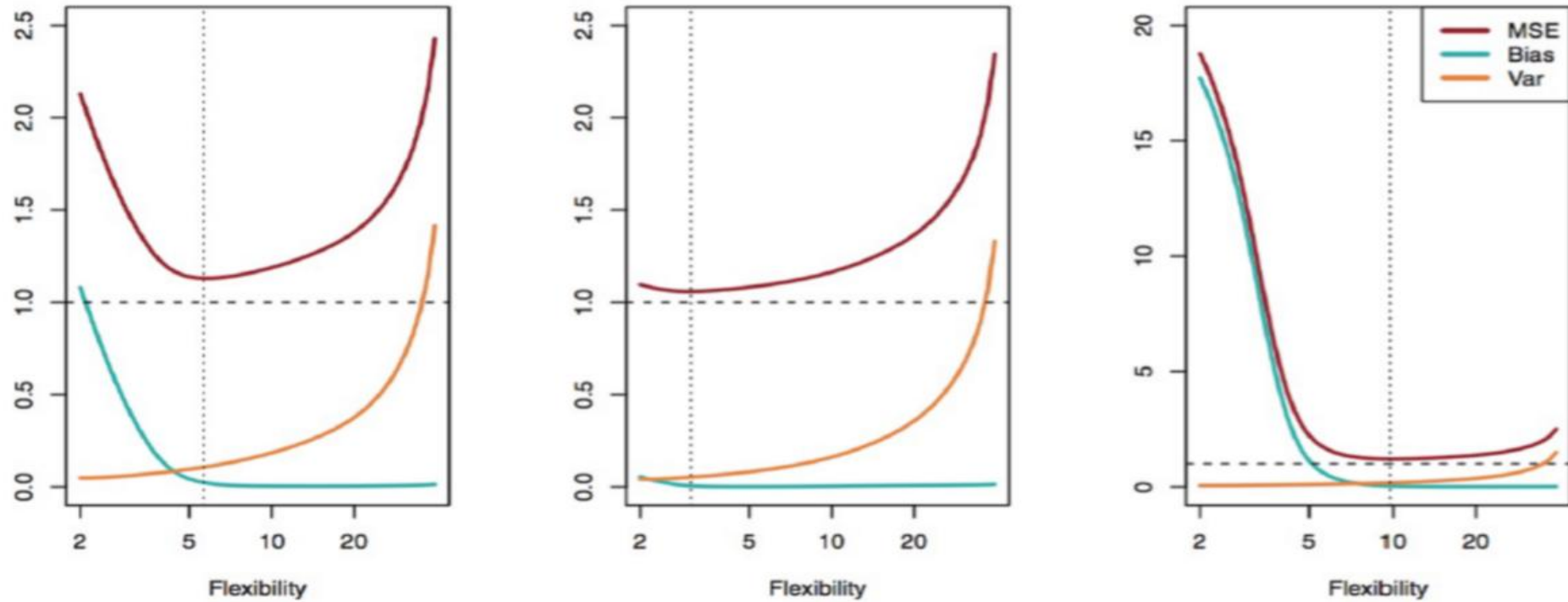


FIGURE 2.12. Squared bias (blue curve), variance (orange curve), $\text{Var}(\epsilon)$ (dashed line), and test MSE (red curve) for the three data sets in Figures 2.9–2.11. The vertical dotted line indicates the flexibility level corresponding to the smallest test MSE.

Old way of splitting data

- Deep learning

train	dev	test
-------	-----	------

60% 20% 20%
1000,000 data 1% 10,000
98% 1% 1%
99.5% 0.4% 0.1%

- K-fold cv

