

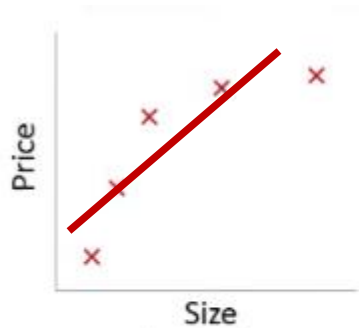


Machine Learning

# Regularization

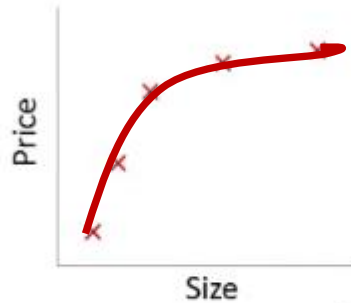
The problem of  
overfitting

## Example : Linear regression (housing prices)



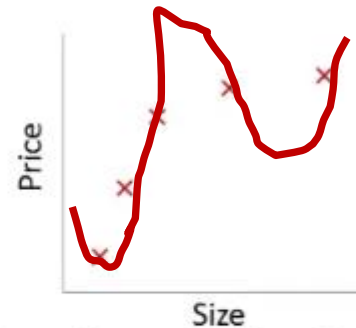
$$\theta_0 + \theta_1 x$$

**‘underfit’ ‘high bias’**  
(红点的纵坐标偏离直线)



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

**‘just right’**

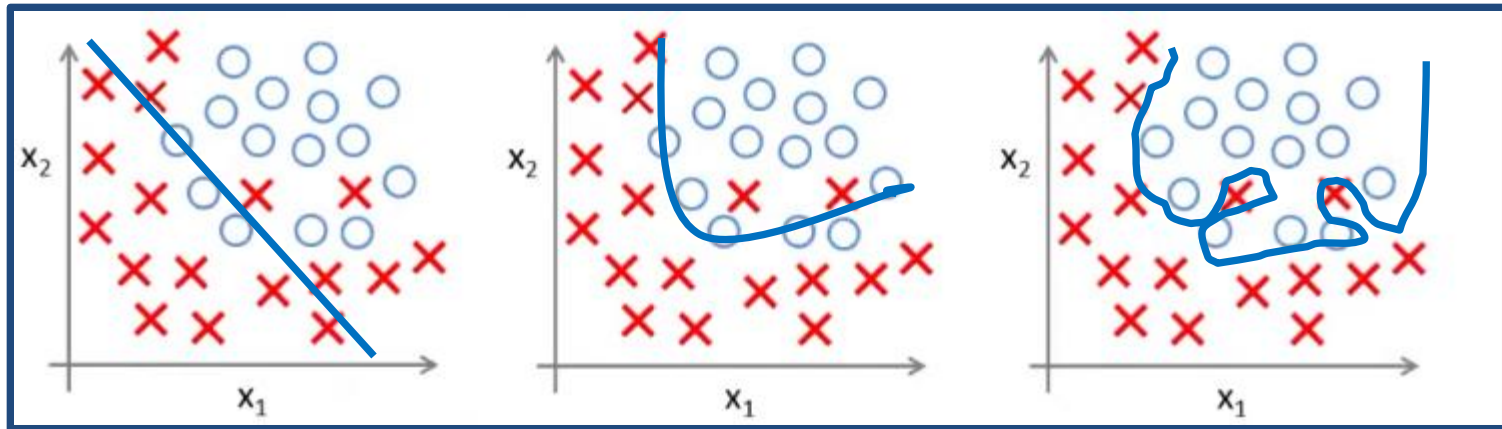


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

**‘overfit’ ‘high variance’**  
(横坐标改变一点点，纵坐标可能发生很大改变)

**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



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

(  $g$  = sigmoid function)

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

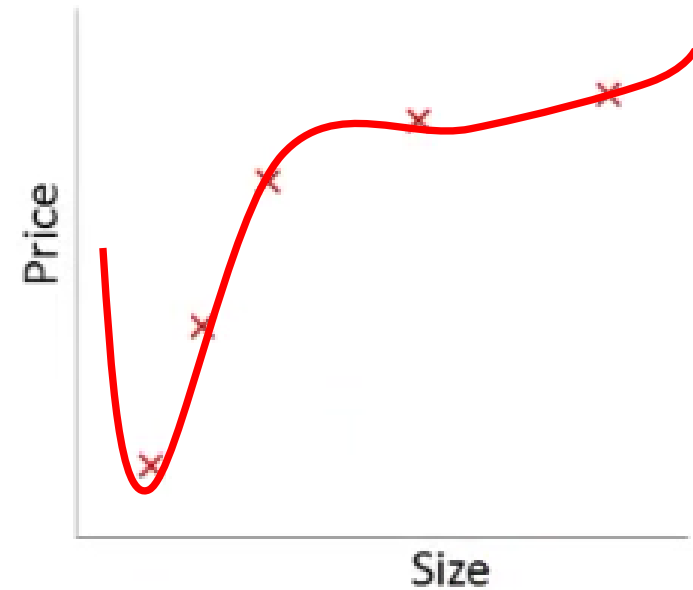
$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

‘underfit’

‘overfit’

# 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  
 $\vdots$   
 $x_{100}$



# Addressing overfitting

## Options

1. Reduce number of features.
  - Manually select which features to keep.
  - Model selection algorithm (later in course).
2. Regularization.
  - Keep all the features, but reduce magnitude/values of parameters  $\theta_j$ .
  - Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

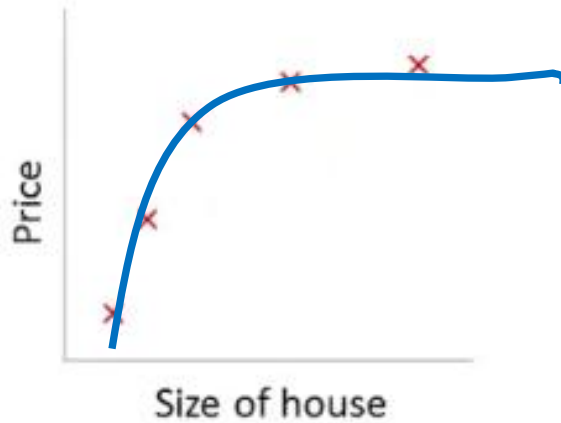


Machine Learning

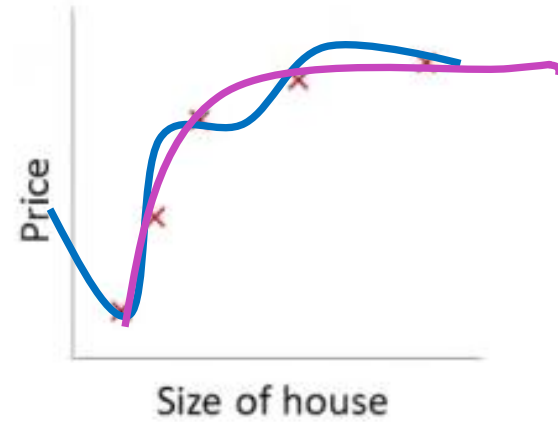
# Regularization

## Cost function

## Intuition



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



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

Suppose we penalize and make  $\theta_3, \theta_4$  really small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 10000\theta_3^2 + 10000\theta_4^2$$

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

## Regularization.

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$

- “Simpler” hypothesis
- Less prone to overfitting

$$\underline{\theta_3 \approx 0, \theta_4 \approx 0}$$

Housing:

- Features:  $\underline{x_1}, \underline{x_2}, \dots, x_{100}$
- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \underline{\theta_{100}}$

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

Regularization term



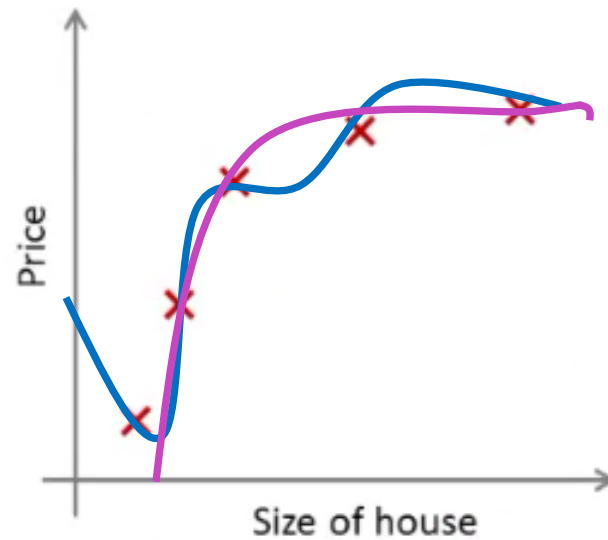
## Regularization.

Regularization parameter

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

←

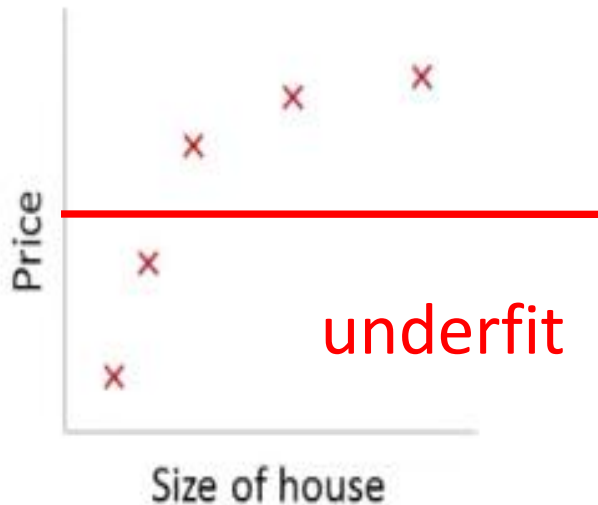
$$\min_{\theta} J(\theta)$$



In regularized linear regression, we choose  $\theta$  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  $\lambda$  is set to an extremely large value (perhaps far too large for our problem, say  $\lambda = 10^{10}$ )?



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

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

$$h(x) = \theta_0$$

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

## Gradient descent

Repeat {

→  $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$

→  $\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right]$   
 $(j = 1, 2, 3, \dots, n)$

}

$\frac{\partial}{\partial \theta_j} J(\theta)$

$\theta_j := \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

< 10.99

$\frac{\partial}{\partial \theta_0} J(\theta)$

## Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

**m\*(n+1)**

$$\min_{\theta} J(\theta)$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

**$R^m$**

$$\theta = (X^T X + \underbrace{\dots \begin{bmatrix} 0 & & \\ & 1 & \\ \dots & & \dots \\ & & & 1 \end{bmatrix} \dots)}_{(n+1)*(n+1)})^{-1} X^T y$$

**$(n+1)*(n+1)$**

## Non-invertibility (optional/advanced).

Suppose  $m \leq n$ ,  
(#examples) (#features)

$$\theta = \underline{(X^T X)^{-1}} X^T y$$

**Non-invertibility**

If  $\lambda > 0$ ,

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

---

**invertibility**



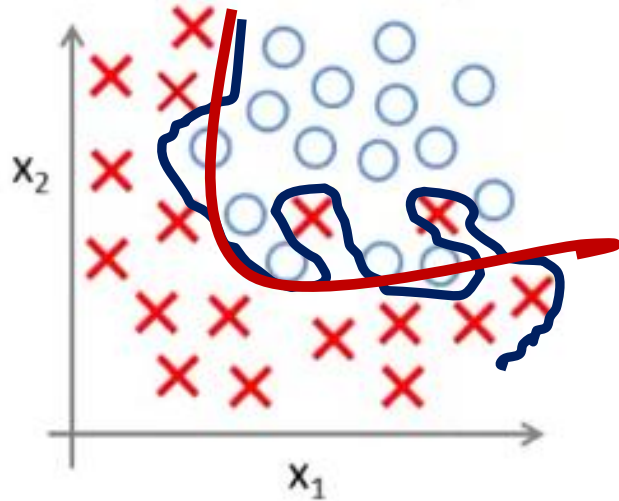
Machine Learning

# Regularization

---

Regularized  
logistic regression

## Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

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

$$+ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

# Gradient descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$(j = 0, 1, 2, 3, \dots, n)$

}



# Gradient descent

Repeat {

→  $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$

→  $\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right]$   
(j = ~~0~~, 1, 2, 3, ..., n)

}

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

$$\theta_1, \theta_2, \dots, \theta_n$$