

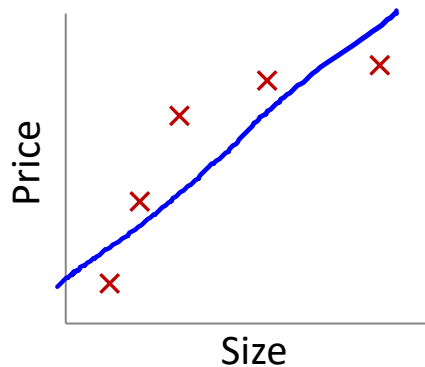
Machine Learning

# Regularization

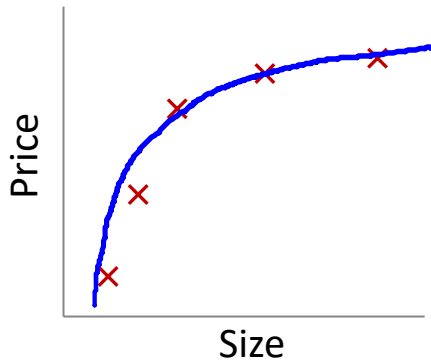
---

## The problem of overfitting

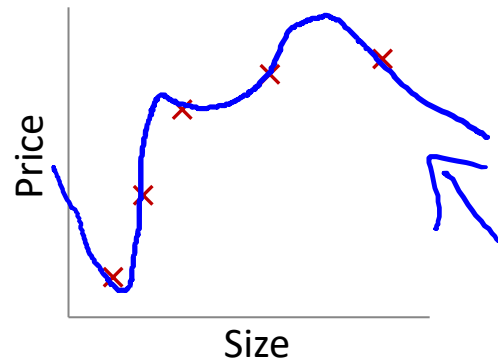
## Example: Linear regression (housing prices)



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



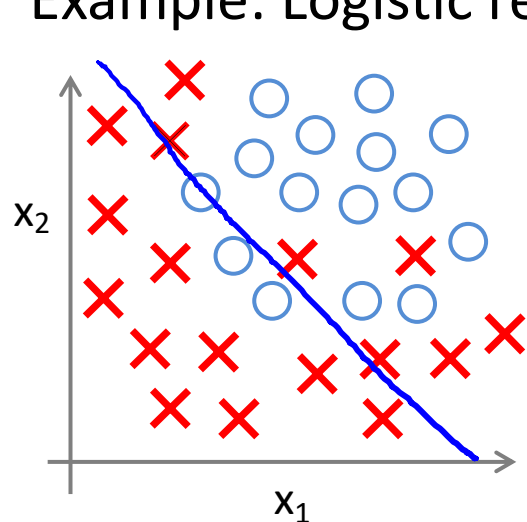
$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$   
"Just right"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$   
"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).

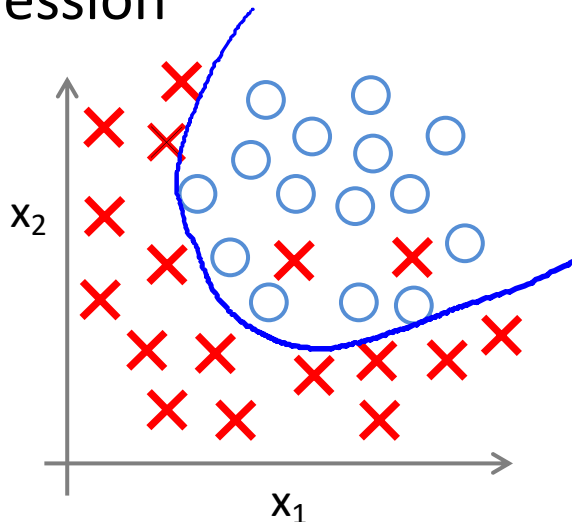
## Example: Logistic regression



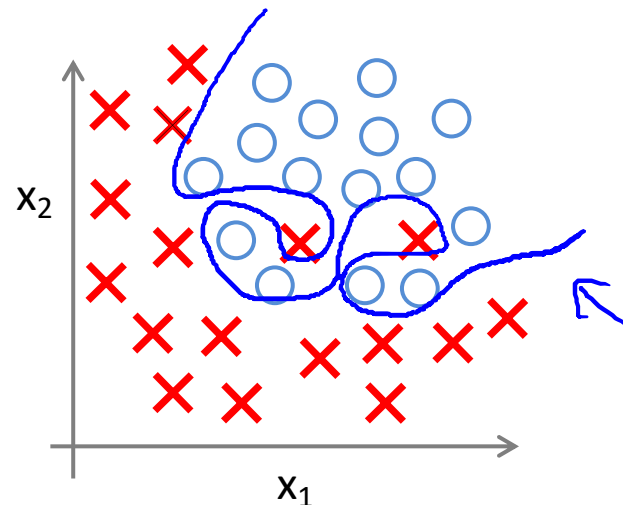
$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

( $g$  = sigmoid function)

"Underfit"



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



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

"Overfit"

## Addressing overfitting:

$x_1$  = size of house

$x_2$  = no. of bedrooms

$x_3$  = no. of floors

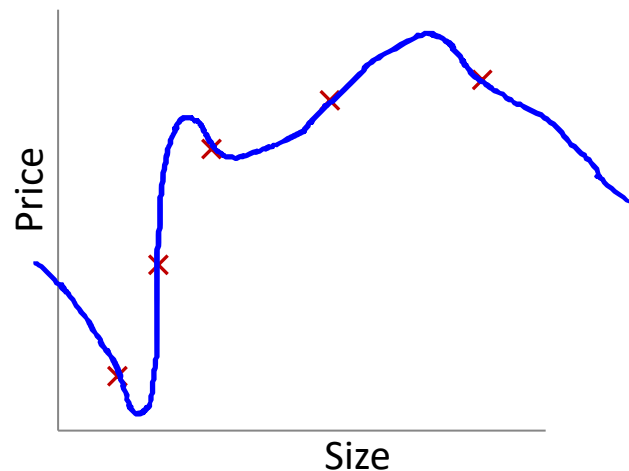
$x_4$  = age of house

$x_5$  = average income in neighborhood

$x_6$  = kitchen size

$\vdots$

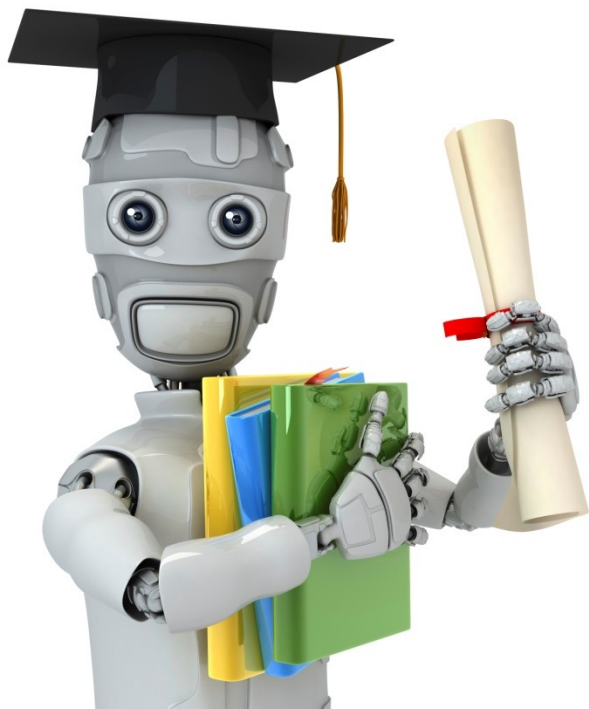
$x_{100}$



## 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  $y$ .



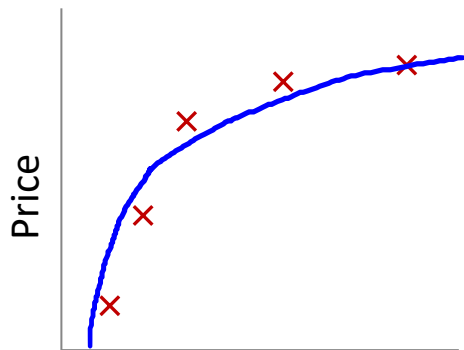
Machine Learning

# Regularization

---

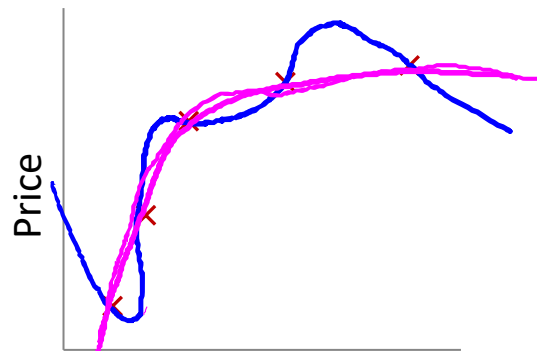
# Cost function

# Intuition



Size of house

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



Size of house

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

Two pink arrows point upwards from the pink underlines below the equation to the crossed-out terms  $\theta_3 x^3$  and  $\theta_4 x^4$ .

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

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

The equation is written in blue ink. The pink underlines and arrows from the previous block are also present. Below the equation, the terms  $\theta_3 \approx 0$  and  $\theta_4 \approx 0$  are written in blue ink, each with a pink underline.

## Regularization.

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

- “Simpler” hypothesis
- Less prone to overfitting

$$\rightarrow \boxed{\theta_3, \theta_4} \approx 0$$

Housing:

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

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

~~$\theta_0$~~   $\theta_1, \theta_2, \theta_3, \dots, \theta_{100}$   ~~$\theta_{100}$~~

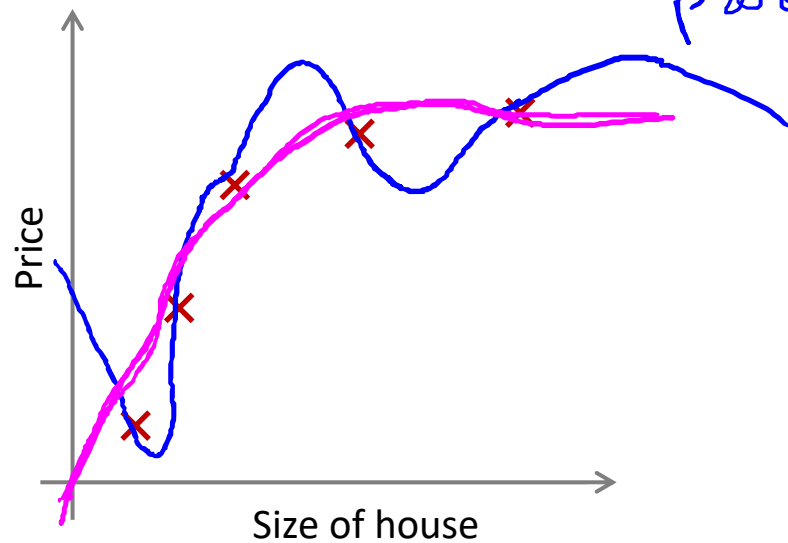


## Regularization.

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

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

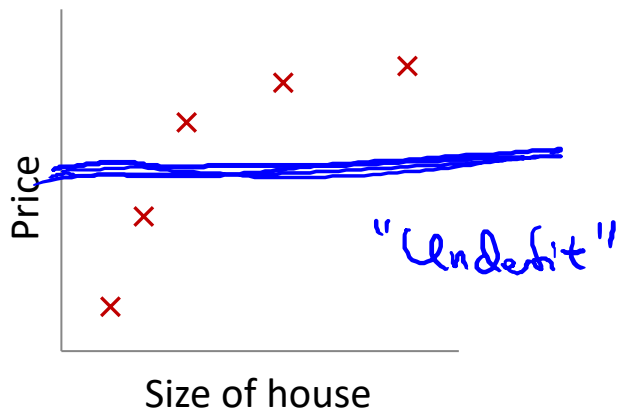
regularization parameter



In regularized linear regression, we choose  $\theta$  to minimize

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

What if  $\lambda$  is set to an extremely large value (perhaps far too large for our problem, say  $\lambda = 10^{10}$ )?



$h_{\theta}(x)$

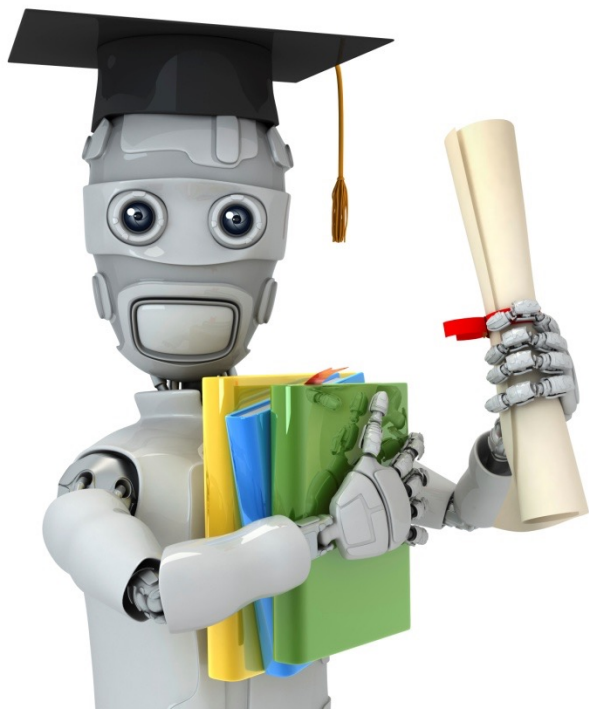
$$\theta_0 + \cancel{\theta_1 x} + \cancel{\theta_2 x^2} + \cancel{\theta_3 x^3} + \cancel{\theta_4 x^4}$$

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

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

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

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



Machine Learning

# Regularization

---

Regularized linear  
regression

# Regularized linear regression

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

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

# Gradient descent

$$\frac{\partial}{\partial \theta_0}$$

$$\theta_0, \theta_1, \theta_2, \dots, \theta_n$$

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

Repeat {

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

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

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

$$- \frac{\lambda}{m} \theta_j$$

}

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

$$\rightarrow J(\theta)$$

$$\theta_j^2$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

$$0.99$$

$$\theta_j \times 0.99$$

# Normal equation

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

$m \times (n+1)$

$$\underset{\uparrow}{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

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

$$\Rightarrow \Theta = \left( X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{(n+1) \times (n+1)} \right)^{-1} X^T y$$

$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0$

$\in \mathbb{R}^n \quad n=2$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

## Non-invertibility (optional/advanced).

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

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

$\underbrace{(X^T X)^{-1}}_{\text{non-invertible / singular}}$

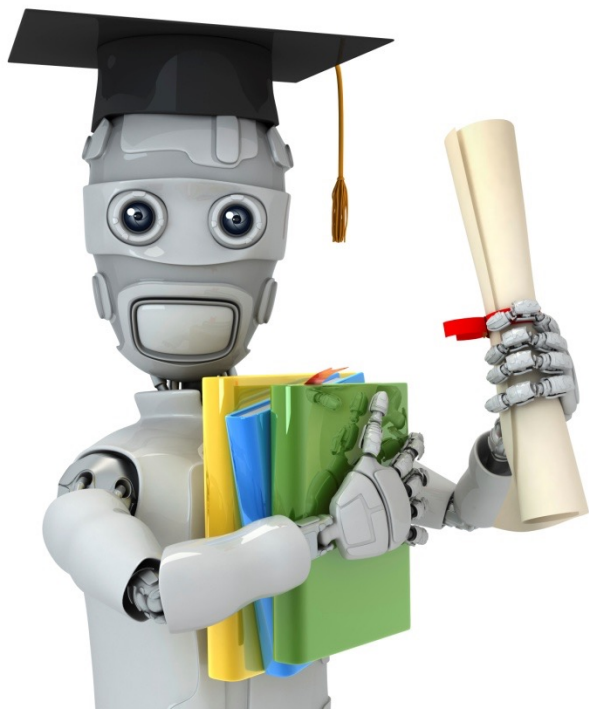
pinv

inv  
 $\nearrow$

If  $\lambda > 0$ ,

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

$\underbrace{\hspace{10em}}_{\text{invertible .}}$



Machine Learning

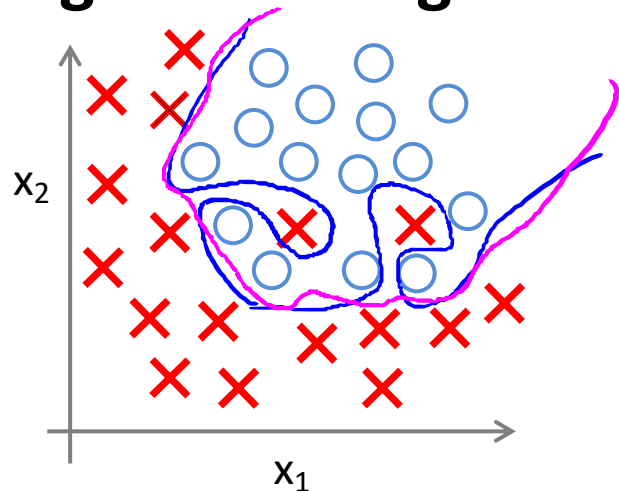
# Regularization

---

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:

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

$\boxed{\theta_1, \theta_2, \dots, \theta_n}$

# Gradient descent

Repeat {

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

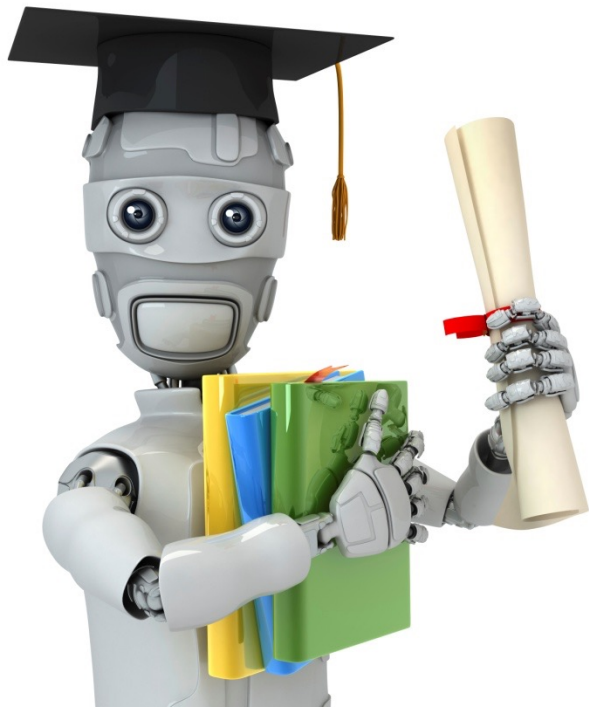
$$\rightarrow \theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{1}{n} \theta_j \right] \leftarrow$$

$(j = \text{red X}, 1, 2, 3, \dots, n)$

$\theta_1, \dots, \theta_n$

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

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



Machine Learning

# Regularization

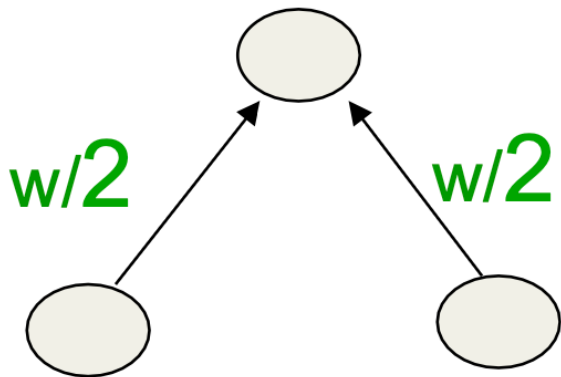
---

## Regularizer for sparsity

## L2-Norm vs L1-Norm

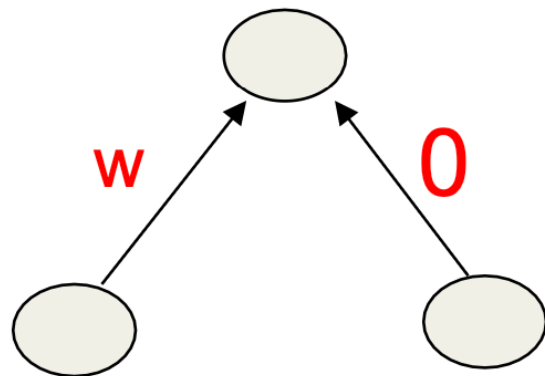
$$\frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x) - y)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

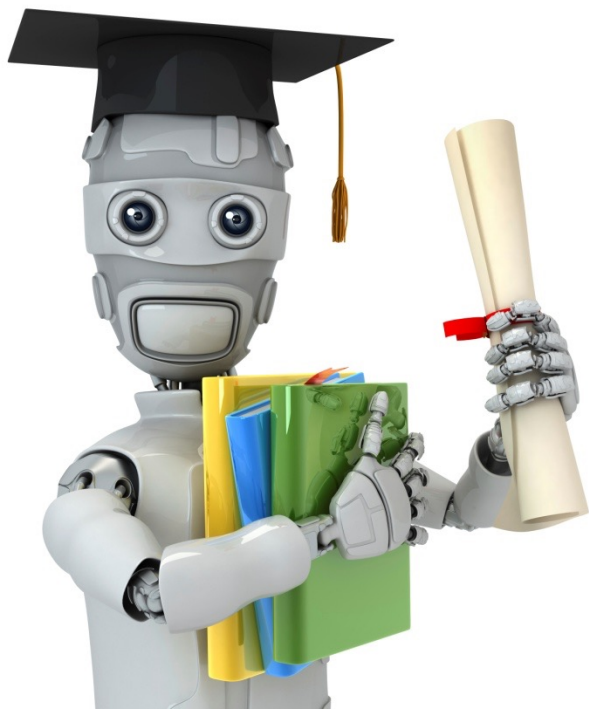
$$\theta_1^2 + \theta_2^2 \leq S$$



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

$$|\theta_1| + |\theta_2| \leq S$$





Machine Learning

# Validation

---

# Validation

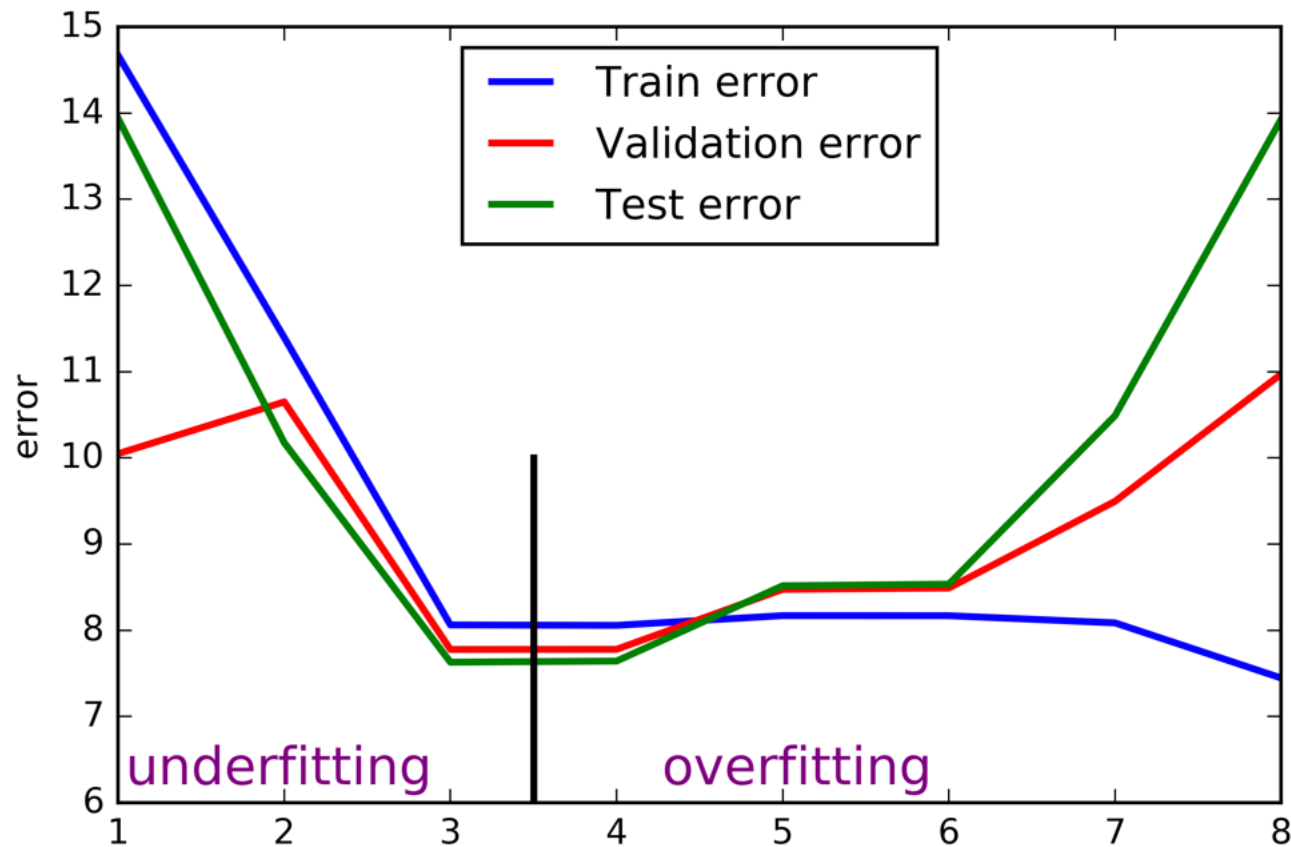
We often divide datasets into two sets:

- training data
- test data (not use in training phase)

How to know the quality of training model with unseen data?

- Extract a small set from training dataset → validation set
- Select the model that provides the small error in both training set and validation set

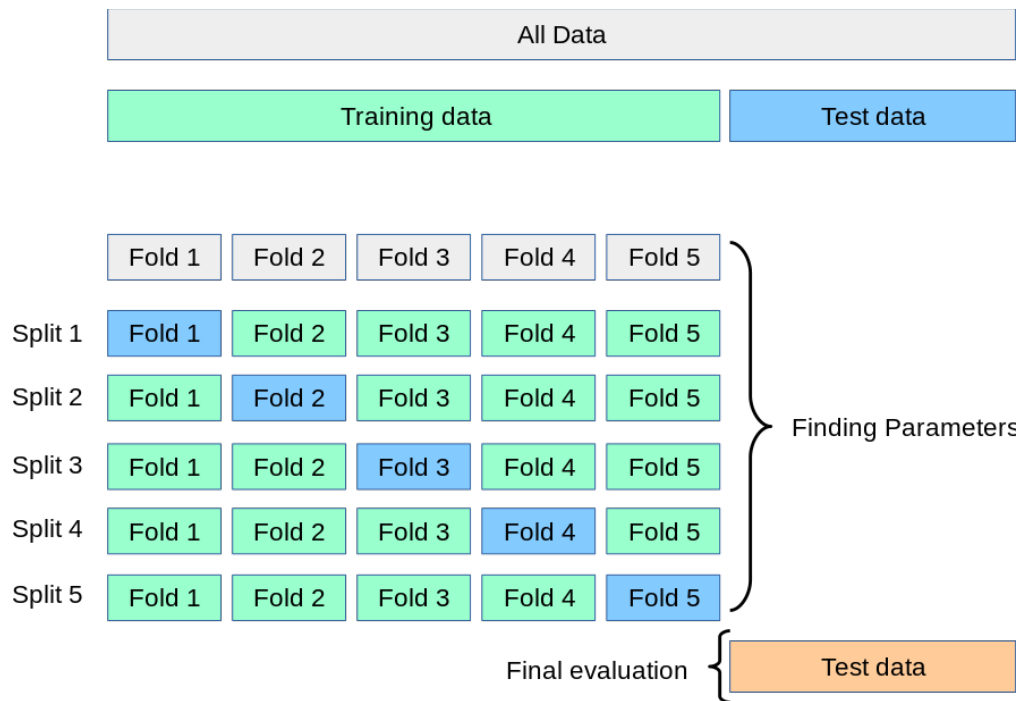
# Validation



# Cross-Validation

If the training dataset is very small, taking too much data from it to form validation set will affect the performance → **k-fold cross val.**

- Split the training set into  $k$  sub-sets
- At each run, one of  $k$  sub-sets is treated as validation set
- The rest  $(k-1)$  sub-sets are treated as training set





Các bạn sử dụng tập dữ liệu của bài tập 5, mô hình dữ liệu  $x$  theo dạng bậc 6 và làm các công việc sau:

1. Chia dữ liệu ra thành 2 tập: training (70%) và validation (30%). Phải đảm bảo việc chia dữ liệu là ngẫu nhiên và tỷ lệ positive và negative cân bằng.
2. Viết chương trình cho phép học các tham số của mô hình phân loại phi tuyến trên có sử dụng regularization L2 và L1.
3. Tính  $J$  ở mỗi vòng lặp cho cả hai tập, chọn điểm dừng phù hợp.
4. Thay đổi  $\lambda$  và tính  $J$  cho mỗi  $\lambda$  tương ứng cho cả hai tập. Vẽ biểu đồ quan hệ giữa  $J$  và  $\lambda$  từ đó lựa chọn  $\lambda$  phù hợp.