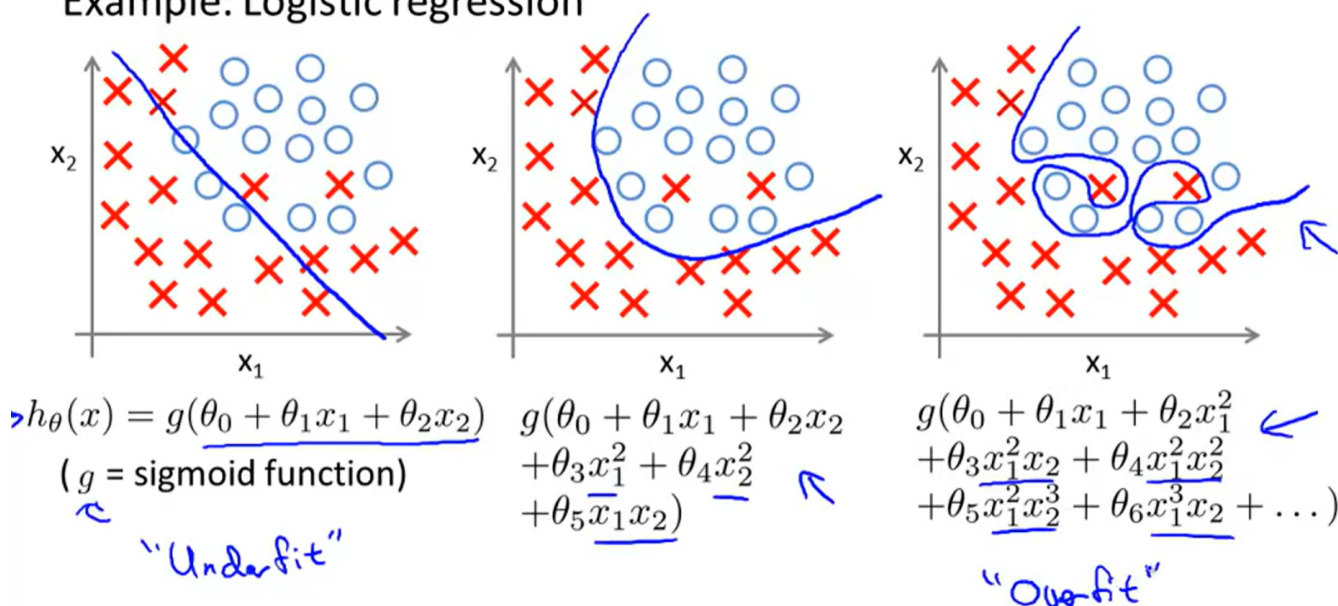# Week 3.4 Solving the Problem of Overfitting

## The Problem of Overfitting

> **Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well, but fail to generalize to new examples (predict prices on new examples).

### Example: Logistic regression



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

"Underfit"

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

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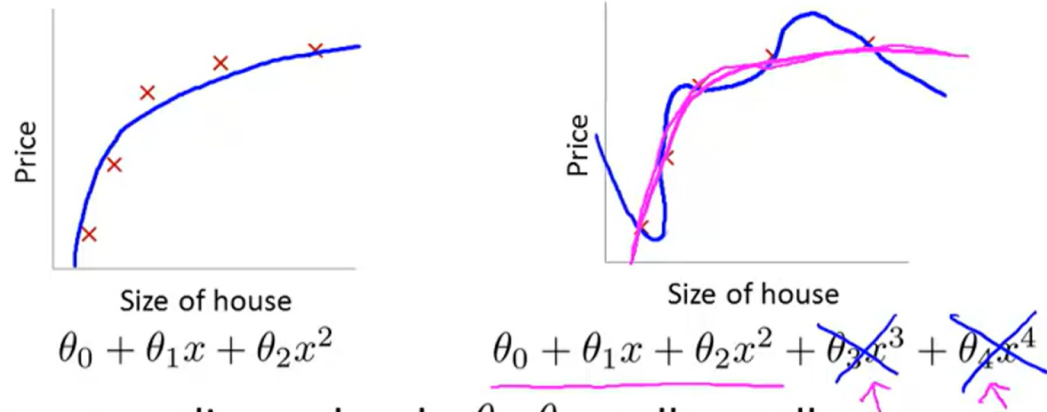
"Overfit"

Options:

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

## Cost Function

- penalize parameters (small values)

**Intuition**



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

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

$$\longrightarrow \quad \min_{\theta} \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\,\theta_3^2 + 1000\,\theta_4^2$$

$$\theta_3 \approx 0 \qquad\qquad \theta_4 \approx 0$$

- what if $\lambda$ is too large: **results in underfitting**

# Regularized Linear Regression

- For gradient descent, we repeat

- 
$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)} x_0^{(i)})$$
$$\theta_j = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

- We rewrite formula for $\theta_j$:

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

- $1 - \alpha \frac{\lambda}{m}$ is a bit smaller than 1

# Normal Equation Method

# Non-invertibility (optional/advanced).

Suppose $m \leq n$, $\leftarrow$

    (#examples)  (#features)

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

*non-invertible /singular*

*pinv*      *inv*

If $\lambda > 0$,

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

*invertible.*

# Regularized Logistic Regression

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

- Pseudocode:

**Advanced optimization**

*fminunc (@ costFunction)*

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ — *theta(1)*, *theta(2)*, *theta(n+1)*

```
function [jVal, gradient] = costFunction(theta)
    jVal = [code to compute J(θ)];
```
$$J(\theta) = \left[-\frac{1}{m} \sum_{i=1}^{m} y^{(i)} \log (h_\theta(x^{(i)}) + (1 - y^{(i)}) \log 1 - h_\theta(x^{(i)}))\right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$
```
    gradient(1) = [code to compute ∂/∂θ₀ J(θ)];
```
$$\frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)} \leftarrow$$
```
    gradient(2) = [code to compute ∂/∂θ₁ J(θ)];
```
$$\left(\frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}\right) + \frac{\lambda}{m}\theta_1 \leftarrow$$

*J(θ)*
```
    gradient(3) = [code to compute ∂/∂θ₂ J(θ)];
```
$$\left(\frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}\right) + \frac{\lambda}{m}\theta_2$$
```
        ⋮
    gradient(n+1) = [code to compute ∂/∂θₙ J(θ)];
```

-

推导过程：

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

$$y^{(i)}\log\left(h_\theta\left(x^{(i)}\right)\right) + \left(1-y^{(i)}\right)\log\left(1-h_\theta\left(x^{(i)}\right)\right) = y^{(i)}\log\left(\frac{1}{1+e^{\theta^T x^{(i)}}}\right) + \left(1-y^{(i)}\right)\log\left(1-\frac{1}{1+e^{\theta^T x^{(i)}}}\right)$$

$$= -y^{(i)}\log\left(1+e^{-\theta^T x^{(i)}}\right) - \left(1-y^{(i)}\right)\log\left(1+e^{\theta^T x^{(i)}}\right)$$

所以: $\frac{\partial}{\partial\theta_j}J\left(\theta\right) = \frac{\partial}{\partial\theta_j}\left[-\frac{1}{m}\sum_{i=1}^{m}\left[-y^{(i)}\log\left(1+e^{-\theta^T x^{(i)}}\right) - \left(1-y^{(i)}\right)\log\left(1+e^{\theta^T x^{(i)}}\right)\right]\right]$

$$= \frac{1}{m}\sum_{i=1}^{m}\left[y^{(i)}\frac{-x_j^{(i)}e^{-\theta^T x^{(i)}}}{1+e^{-\theta^T x^{(i)}}} - \left(1-y^{(i)}\right)\frac{x_j^{(i)}e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}}\right] = -\frac{1}{m}\sum_{i=1}^{m}y^{(i)}\frac{x_j^{(i)}}{1+e^{\theta^T x^{(i)}}} - \left(1-y^{(i)}\right)\frac{x_j^{(i)}e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}}]$$

$$= \frac{1}{m}\sum_{i=1}^{m}\frac{y^{(i)}x_j^{(i)}-x_j^{(i)}e^{\theta^T x^{(i)}}+y^{(i)}x_j^{(i)}e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}} = -\frac{1}{m}\sum_{i=1}^{m}\frac{y^{(i)}\left(1+e^{\theta^T x^{(i)}}\right)-e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}}x_j^{(i)} = -\frac{1}{m}\sum_{i=1}^{m}\left(y^{(i)}-\frac{e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}}\right)x_j^{(i)}$$

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

-