

## Week 2

### Multivariate Linear Regression:

- Multiple Features:

$x_j^{(i)}$  = value of feature  $j$  in the  $i^{th}$  training example

$x^{(i)}$  = the input (features) of the  $i^{th}$  training example

$m$  = the number of training examples

$n$  = the number of features

$$\rightarrow h_{\theta}(x) = \underline{\theta_0} + \underline{\theta_1 x_1} + \underline{\theta_2 x_2} + \cdots + \underline{\theta_n x_n}$$

For convenience of notation, define  $\underline{x_0} = 1$  ( $x_0^{(i)} = 1$ )

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \vdots \\ \Theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \underline{\Theta_0 x_0 + \Theta_1 x_1 + \cdots + \Theta_n x_n} \\ = \underline{\Theta^T x}$$

$$\Theta^T = \begin{bmatrix} \Theta_0 & \Theta_1 & \cdots & \Theta_n \end{bmatrix}$$

$$(n+1) \times 1 \text{ matrix}$$

$$\Theta^T x$$

- Gradient Descent for Multiple Variables:

↗ New algorithm ( $n \geq 1$ ):

$$\text{Repeat } \left\{ \frac{\partial}{\partial \Theta_j} J(\Theta) \right\}$$

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

(simultaneously update  $\theta_j$  for  
 $j = 0, \dots, n$ )

$$\}$$

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

$$\rightarrow \underline{\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}}$$

$$\rightarrow \underline{\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}}$$

...

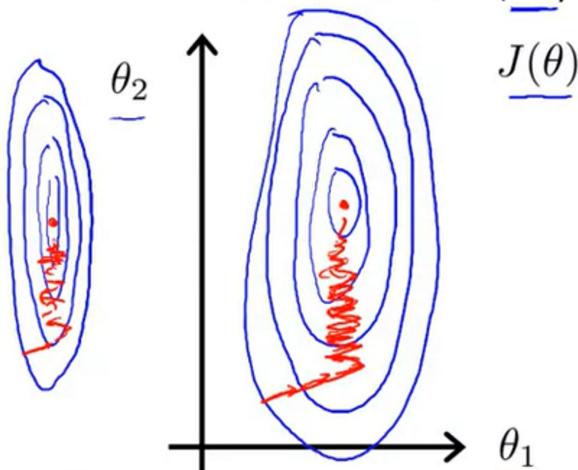
- Feature scaling: (make gradient descent run much faster and converge)

## Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet}^2)$  ↪

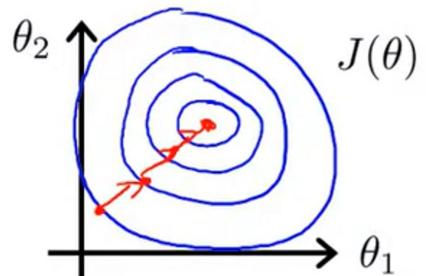
$x_2 = \text{number of bedrooms (1-5)}$  ↪



$$\rightarrow x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$\rightarrow x_2 = \frac{\text{number of bedrooms}}{5}$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



The scope of features are different, as a result, It's difficult to converge. So we're going to scale the feature into [-1,1] range.

(range isn't important, we just scale it into a small range.):

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq 100 \quad \times$$

$$-0.0001 \leq x_4 \leq 0.0001 \quad \times$$

Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean  
(Do not apply to  $x_0 = 1$ ).

E.g.  $\rightarrow x_1 = \frac{\text{size}-1000}{2000}$

Average size = 100

$$x_2 = \frac{\#bedrooms-2}{5}$$

1-5 bedrooms

$$[-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5]$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1}$$

avg value of  $x_1$  in training set

range (max-min) (or standard deviation)

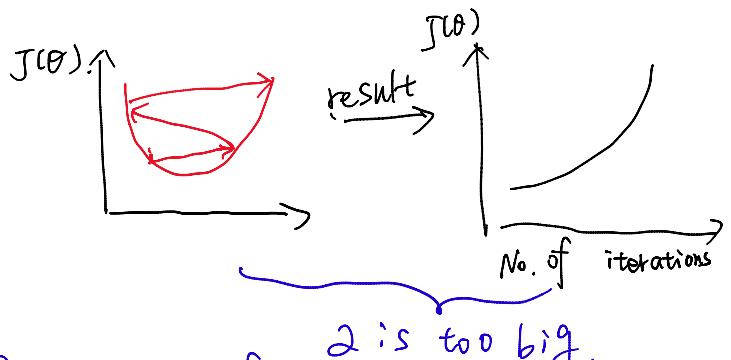
$$x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}$$

- Learning rate

$J(\theta)$  should decrease after every iteration

Example automatic convergence test:

Declare convergence if  $\frac{J(\theta)}{\sum}$  decreases by less than  $10^{-3}$  in one iteration.



plot  $J(\theta)$  as a function of the number of iterations. to figure out what's going on.

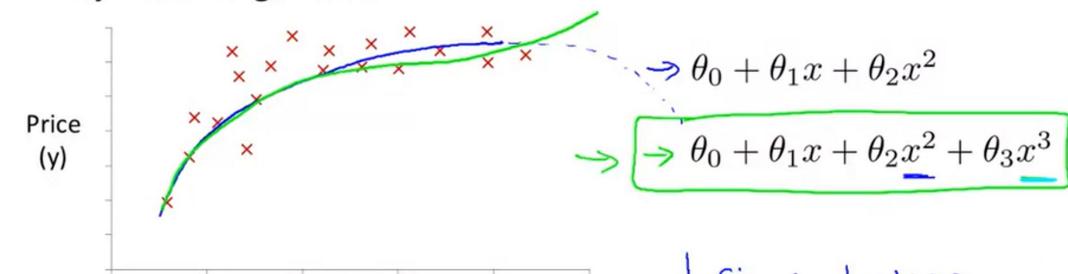
try a range of values for  $\alpha$ . and plot the plotting.

- Features and polynomial regression

Create a new features by yourself

### Polynomial regression

Transfer the problem into linear regression



$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

$\rightarrow x_1 = (\text{size})$

$\rightarrow x_2 = (\text{size})^2$

$\rightarrow x_3 = (\text{size})^3$

create new features

Size: 1 - 1600  
 $\text{Size}^2$ : 1 - 1000,000  
 $\text{Size}^3$ : 1 - 10<sup>9</sup>

Feature scaling is very important

## Computing Parameters Analytically:

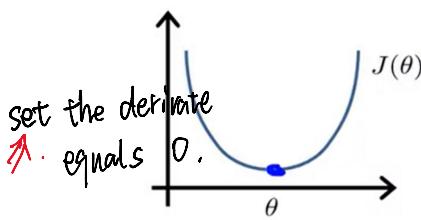
- Normal equation

Intuition: If 1D ( $\theta \in \mathbb{R}$ )

$$\rightarrow J(\theta) = a\theta^2 + b\theta + c$$

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

Solve for  $\theta$



$$\theta \in \mathbb{R}^{n+1}$$

$$J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

Solve for  $\theta_0, \theta_1, \dots, \theta_n$  set

we want..  $\theta = (X^T X)^{-1} X^T y$ .  
↓ inverse of  $X^T X$ .

$$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \xrightarrow{\text{design matrix}} m \times (n+1).$$

m training examples, n features.

### Gradient Descent

- Need to choose  $\alpha$ .
- Needs many iterations.
- Works well even when n is large.

$$n = 10^6$$

### Normal Equation

- No need to choose  $\alpha$ .
- Don't need to iterate.
- Need to compute  $(X^T X)^{-1}$   $n \times n$   $O(n^3)$
- Slow if n is very large.

$$n = 100$$

$$n = 1000$$

$$\dots \underset{n=10000}{\dots}$$

- Normal Equation Non-invertibility

Like above we use normal equation to compute theta, at that computing we use inverse operate. But sometimes, singular or degenerate matrices is non-invertible.

At the learning time, the reason for  $X$  transpose  $X$  to be non-invertible is redundant features and smaller  $m$  than  $n$  ( $m \leq n$ , the number of samples are less or equal to the number of features)

Tip : vectorization:

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

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

(n = 2)

Vectorized implementation:

$$\Theta := \Theta - \alpha \delta \quad \text{where } \delta = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) X^{(i)}$$

$$\delta_0 = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\delta_1 = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\delta_2 = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

$u(j) = 2v(j) + 5w(j)$  (for all j)

$$u = 2v + 5w$$

$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$