

# multiple linear regression (multiple variables/features)

(3-1)

$n$ : number of all features.  $m$ : the number of training examples. (行)

$x_{i,j}^{(i)}$ : input (features) of  $i^{th}$  training example  $x_{i,j}^{(i)} \rightarrow$  index of training sets (examples).

$x_j$ : Value of feature  $j$  in  $i^{th}$  training examples.

$x_n^{(i)}$  第  $i$  个样本的第  $n$  个特征.

hypothesis:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \quad n: \text{the number of features.}$$

For convenience of notation: define  $\theta_0 = 1$ .  $(x_0^{(i)}) = 1$

(the zero feature)

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

feature vector.

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

parameter vector.

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots, \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\theta^T \cdot x$$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$= \underline{\theta^T \cdot x}$$

$$h_\theta(x) = \theta_0 + \theta_1 x.$$

## Multivariate linear regression.

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

$\theta$  is a  $n+1$ -dimensional vector.

$$\text{cost function: } J(\theta_0, \theta_1, \dots, \theta_n) = \underbrace{\frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2}_{J(\theta)}$$

Gradient Descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$$

$$= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

j.

simultaneously update for every  $j = 0, 1, \dots, n$

$n > 1$

Repeat {

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

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

},

$n \geq 1$

$$h_\theta(x) = \theta_0 + \theta_1 x.$$

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

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

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

e.g.

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_0^{(i)}) - y_0^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_1^{(i)}) - y_1^{(i)}) \cdot x_1^{(i)}$$

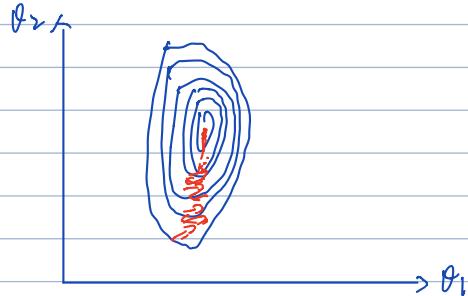
$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_2^{(i)}) - y_2^{(i)}) \cdot x_2^{(i)}$$

assume  $x_0^{(i)} = 1$

## Feature Scaling

Idea: make sure features are on a similar scale, then the GD will be converged faster.

e.g.  $\theta_1 = \text{size } (\theta \sim 2000 \text{ feet}^2)$   
 $\theta_2 = \text{number of bedrooms } (1-5)$ .

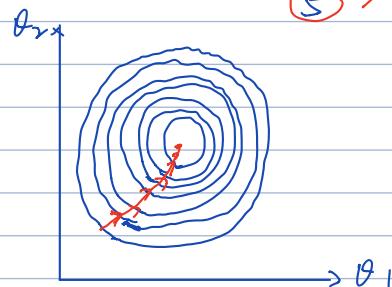


Ignore  $\theta_0$ .

converge slowly

$$\theta_1 = \frac{\text{size.}}{2000} \quad 0 \leq \theta_1 \leq 1$$

$$\theta_2 = \frac{\text{number}}{5} \quad 0 \leq \theta_2 \leq 1$$



converge faster and more directly.

normally, we get every feature into approximately a  $-1 \leq \theta_i \leq 1$  range or near it.

$$\theta_0 = 1 \checkmark$$

$$0 \leq \theta_1 \leq 3. \checkmark$$

$$-2 \leq \theta_2 \leq 0.5 \checkmark$$

$$-100 \leq \theta_3 \leq 100 \times \text{too large.}$$

$$-0.0001 \leq \theta_4 \leq 0.0001 \times \text{too small.}$$

people have different notation on the range of feature scaling according to different applications.

## Mean Normalization: $\bar{x}_i - \bar{x}_i$ .

Replace  $x_i$  with  $\frac{x_i - \bar{x}_i}{s_i}$  to make features have approximately zero mean (do not apply to  $\theta_0 = 1$ ).

mean normalization involves subtracting the average value for an input variables from the values for that input variable resulting in a new average value for the input variables of just zero.

$$x_i = \frac{x_i - \bar{x}_i}{s_i} \rightarrow \text{average value of that variable}$$

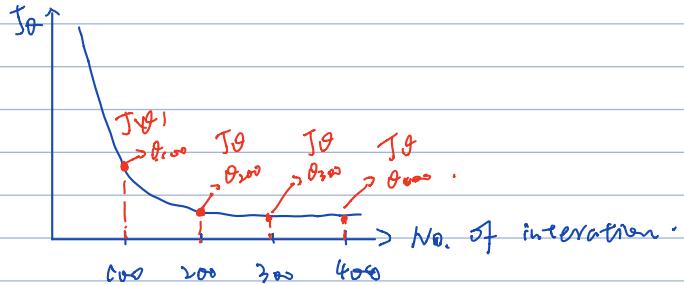
$s_i$   $\rightarrow$  the range of value (max - min) or standard deviation.  
 If range 和 布差有不同的效果.

## Debugging Gradient Descent

Make a plot with the number of iteration on the x-axis, the value of  $J(\theta)$  on the y-axis.

If  $J(\theta)$  is increasing, then you probably need to decrease the  $\alpha$ .

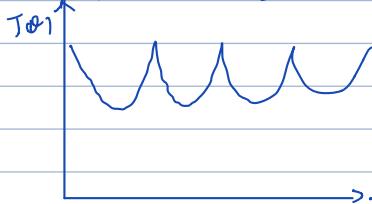
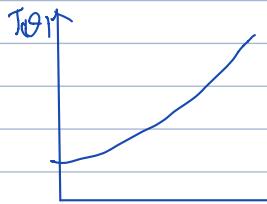
Normally,  $J(\theta)$  should decrease after every iteration,



Automatic convergence test:

Declare convergence if  $J(\theta)$  decrease by less than  $T\epsilon$  in one iteration where  $\epsilon$  is some small value such as  $10^{-3}$ . However, it's difficult to choose this threshold. value.

Gradient Descent is not working correctly:



May because the learning rate is too large.

Summary:

- $\lambda$  is too small : slow convergence
- $\lambda$  is too large : may not decrease on every iteration and thus may not converge.

## Feature and Polynomial Regression

Our hypothesis function need not be linear (a straight) line if that does not fit the data well.

We can use quadratic (= $x^2$ ) cubic ( $=x^3$ ) or square root function.

E.g.

original :  $h(x) = \theta_0 + \theta_1 x_1$ , we can create additional features based on  $x_1$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \leftarrow \text{quadratic}$$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 \leftarrow \text{cubic}$$

$x_1^2 = x_2 \quad x_1^3 = x_3$ .

Notice : if you choose your hypothesis function like this, feature scaling is increasingly important.

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$$

if  $x_1 \in [1, 1000] \quad x_1^2 \in [1, 10^6]$   
 $x_1^3 \in [1, 10^9]$

Normal Equation : Don't need to do feature scaling.  
 Method to solve for  $\theta$  analytically.

$$ID = (\theta \in \mathbb{R}^{k+1})$$

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

$$\frac{\partial}{\partial \theta} J(\theta) = 0$$

Solve for  $\theta$ .

$(n+1)$  features :

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

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

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

In training examples  $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ ,  $n$  features.

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

$(n+1) \times 1$

design matrix

$$X = \begin{bmatrix} \cdots (x^{(1)})^T \cdots \\ \cdots (x^{(2)})^T \cdots \\ \vdots \\ \cdots (x^{(m)})^T \cdots \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} y_0^{(1)} \\ y_1^{(2)} \\ \vdots \\ y_m^{(m)} \end{bmatrix}$$

$m \times 1$

$$\theta = (X^T X)^{-1} X^T y \Rightarrow \min J(\theta).$$

method to solve for  $\theta$  analytically. linear regression.  
Normal Equation ↴

Gradient Descent

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

(more than 10000)

- No need to choose  $\alpha$ .
- Don't need to integrate.
- Need to compute  $(X^T X)^{-1}$
- Show if  $n$  is very large.

Non-invertibility (singular / degenerate)

$\left\{ \begin{array}{l} \text{pinv} \rightarrow \text{even } (X^T X) \text{ is non-invertible.} \\ \text{inv} \end{array} \right.$

What if  $X^T X$  is non-invertible?

① Redundant features (linearly dependent).

E.g.,  $x_1 = \text{size in feet}^2$ .

$x_2 = \text{size in m}^2$ .

$x_1 - x_2$  always satisfy this equation =  $x_1 = (3.28)^2 x_2$

$1m = 3.28 \text{ feet}$

② Too many features (with less training examples). ( $m \leq n$ )

- Delete some features, OR use regularization.

cannot learn features from fewer training examples