

# Linear regression with multiple variables

## 4.1) Multiple features

Notation:

$n$  = Number of features

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

$x_j^{(i)}$  = Value of feature  $j$  in  $i^{\text{th}}$  training example

$\swarrow \searrow$  (Will be a number)

(Will be  $n$  dimensional vector)

$$x^{(i)} \in \mathbb{R}^n$$

$\Rightarrow$  New hypothesis function:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

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

$$\text{So } x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \& \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

$$h_{\theta}(x) = \theta^T x \quad \left\{ x = \begin{bmatrix} x_0 \\ x^{(i)} \end{bmatrix} \right\}$$

$\Rightarrow$  Multivariate linear regression

## 4.2) Gradient descent for multiple variables

### Cost function

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

Vector

Vector

### Gradient descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \left\{ j = 0, 1, \dots, n \right\}$$

}

⇓

Repeat {

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

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

}

## 4.37 Feature Scaling

If different features take different range of values then it becomes very difficult for algorithm to find minimum or it takes way more time for algorithm to find minimum.

→ To avoid this we scale every feature so that they are in similar range.

⇒ Get every feature into approximately a  $-1 \leq x_i \leq 1$  range.

## Mean Normalization

⇒ Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean.

General rule  $x_i \leftrightarrow \frac{x_i - \mu_i}{s_i}$

Annotations:

- $\mu_i$ : avg value of  $x_i$  in training set
- $s_i$ : range of  $x_i$  (i.e.  $\max x_i - \min x_i$ )

⇒ To make sure gradient descent is working correctly.

→ Plot  $J(\theta)$  vs (Number of iteration)

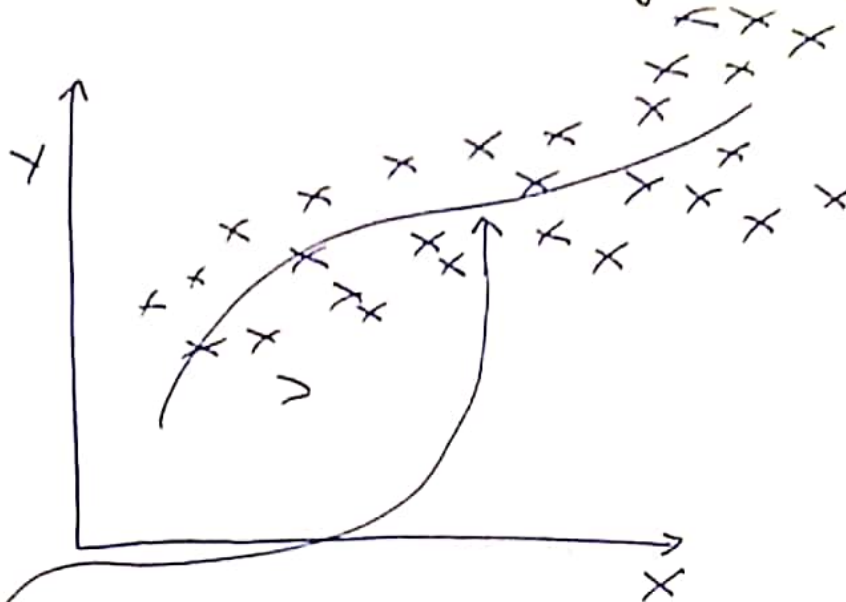
↓  
{ If it is continuously decreasing }  
→ then it's going well.

↓  
{ If it is almost flat then you }  
→ can assume that it has converged

$\Rightarrow$  To choose  $\alpha$ , try

... 0.001, 0.01, 0.1, 1 ...

#### 4.5 > Polynomial regression



$h_0(x) = \theta_0 + \theta_1 x$   $\left\{ \begin{array}{l} \text{X} \\ \text{doesn't give good} \\ \text{fit} \end{array} \right\}$

So,  $h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

$\left\{ \begin{array}{l} \text{Very good fit} \end{array} \right\}$

So  $\begin{array}{l} x_1 = x \\ x_2 = x^2 \\ x_3 = x^3 \end{array} \left\{ \begin{array}{l} \text{define three new} \\ \text{features} \end{array} \right\}$



## 4.6 Normal Equations

Method to solve for  $\theta$  analytically.

Let  $X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ \vdots & x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{bmatrix}_{m \times (n+1)}$

number of features

Matrix that contains all the training examples

number of training set

Similarly  $y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix}_{m \times 1}$  { Collection of all the target variables in training examples }

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

This gives you the value of  $\theta$  that minimizes the cost function

Note

$$J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$$

## Gradient Descent

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

$$n > 10000$$

## Normal Equation

- ⇒ No need to choose  $\alpha$ .
- ⇒ Don't need to iterate.
- ⇒ Slow if  $n$  is very large.

$$O(n^3)$$

$$n = 100 \quad \text{OK}$$

$$n = 1000 \quad \text{OK}$$

$$n = 10000 \leftarrow \text{Not greater than this}$$