# WEEK 2

## Multivariate Linear Regression

Notation:   $n$: # of features

$x^{(i)}$: input of $i^{th}$ training example

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

Hypothesis:   $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots \theta_n x_n$

define $x_0 = 1$, then $X = \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$

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

then: $h_\theta(x) = \theta^T X$

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

Gradient descent

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

Could be due to learning rate $\alpha$ two large, try use small $\alpha$

Mathematically, for sufficiently small $\alpha$, $J(\theta)$ should decrease on every iteration but if $\alpha$ is too small, gradient decent can be slow to converge

### Features and Polynomial Regression

eg: $h(\theta) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$

$\hookrightarrow x_1 = (size)$     $x_2 = (size)^2$   $x_3 = (size)^3$

One should take care of feature scaling, since $x_1 \sim x_3$ takes every different ranges

## Normal Equation:

Normal equation: method solve for $\theta$ analytically

$\dfrac{\partial J(\theta)}{\partial \theta_j} = 0 \implies \theta_0 \cdots \theta_n$

| $x_0$ | feature 1 $x_1$ | feature 2 $x_2$ | feature 3 $x_3$ | feature 4 $x_4$ | results $y$ |
|---|---|---|---|---|---|
| 1 |  |  |  |  |  |
| 1 |  |  |  |  |  |
| 1 |  |  |  |  |  |
| ⋮ |  |  |  |  |  |
| 1 |  |  |  |  |  |

$$X = \begin{bmatrix} 1 & & & \\ 1 & & & \\ 1 & & & \\ 1 & & & \end{bmatrix}_{m \times (n+1)} \quad y = \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{m \times 1}$$

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

Generally:

m examples. $(x^{(1)}, y^{(1)}), \cdots\cdots (x^{(m)}, y^{(m)})$

n features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \qquad X = \begin{bmatrix} \text{---} (X^{(1)})^T \text{---} \\ \vdots \quad \vdots \\ \text{---} (X^{(m)})^T \text{---} \end{bmatrix}$$

$$\underset{\text{design matrix}}{\uparrow} \qquad\qquad m \times (n+1)$$

It Does NOT matter whether using feature scaling or not in this case

Compare :

| Gradient decent | Normal equation |
|---|---|
| o need to chose $\alpha$ | o no need chose $\alpha$ |
| o needs many iterations | o No need of iteration |
| o Works well even feature $n$ is longer | o need to compute $[X^T X]^{-1}$ (slow when $n$ is large) |
| o more general | o linear regression only |
| | o $[X^T X]^{-1}$ may non-invertible |

what does $[X^T X]$ is non-invertible mean ?

     o Redundant features

     o Too many features ($m \leq n$)

Vectorization :

$$h_\theta(x) = \sum_{j=0}^{n} \theta_j x_j = \theta^T X$$