$m$ : # training set

$\theta$ : parameter.　　$h_\theta(x)$ : hypothesis function.

loss function : $J(\theta) = \frac{1}{2} \sum_{j=1}^{m} (h(x) - y)^2$

min : $J(\theta)$

$\theta_i = \theta_i - \alpha \frac{d}{d\theta_i} J(\theta)$.

$\frac{d}{d\theta_i} J(\theta) = \frac{1}{2} \cdot 2 (h_\theta(x) - y) \cdot \frac{d}{d\theta_i} (h_\theta(x) - y)$

$= (h_\theta(x) - y) \cdot \frac{d}{d\theta_i} (\theta_0 x_0 + \dots + \theta_n x_n - y)$

$= (h_\theta(x) - y) \cdot x_i$

update : $\theta_i := \theta_i - \alpha \sum_{j=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_i^{(i)}$

　　　　　learning rate

repeat until convergence. (Batch Gradient Descent)

## # Stochastic Gradient Descent

Repeat {

　For $j = 1$ to $m$ {

　　$\theta_i := \theta_i - \alpha (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_i^{(i)}$ (for all $i$)

　}

}

Linear Algebra :

$\nabla_\theta J = \left[ \frac{\partial J}{\partial \theta_0}, \dots \frac{\partial J}{\partial \theta_n} \right]^T \in R^{n+1}$

Thus : Gradient descent :

$\theta := \theta - \alpha \nabla_\theta J$ 　 // $\theta, J \in R^{n+1}$

$f : R^{m \times n} \to R$ 　 $f(A)$

$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$

Fact :

$tr(AB) = tr(BA)$ 　 $tr(ABC) = tr(CAB) = tr(BCA)$

$\nabla_A tr AB = B^T$

$\nabla_A tr ABA^T C = CAB + C^T AB^T$

Design Matrix : $X\theta = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$ $\theta = \begin{bmatrix} x^{(1)T}\theta \\ x^{(2)T}\theta \\ \vdots \\ x^{(m)T}\theta \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) \\ \vdots \\ h_\theta(x^{(m)}) \end{bmatrix}$

$\bar{y} = [y^{(1)}, \dots, y^{(m)}]^T$

$X\theta - y = [h(x^{(1)}) - y^{(1)}, \dots, h(x^{(m)}) - y^{(m)}]^T$

Recall : $z^T z = \sum_i z_i^2$

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

Recall : min $J(\theta)$

$\nabla_\theta \frac{1}{2} (X\theta - y)^T (X\theta - y)$

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

$= \frac{1}{2} \nabla_\theta tr (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y)$

$= \frac{1}{2} [\nabla_\theta tr \theta\theta^T X^T X - \nabla_\theta tr y^T X\theta - \nabla_\theta tr y^T X\theta]$

$\nabla_\theta tr \underset{A}{\theta} \underset{B}{I} \underset{A^T}{\theta^T} \underset{C}{\underline{X^T X}} = X^T X\theta + X^T X\theta$

$\nabla_\theta tr \underset{B}{\underline{y^T X}} \underset{A}{\theta} = X^T y$

$\nabla_\theta J(\theta) = \frac{1}{2} [X^T X\theta + X^T X\theta - X^T y - X^T y] = 0$

$X^T X\theta - X^T y = 0$

$X^T X\theta = X^T y$ 　 Normal Equation.

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