

definition of 2 form:

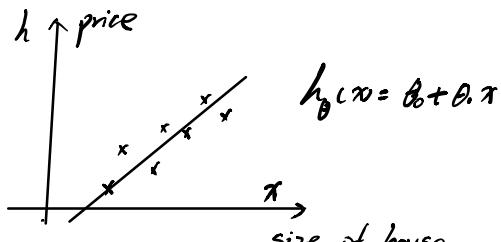
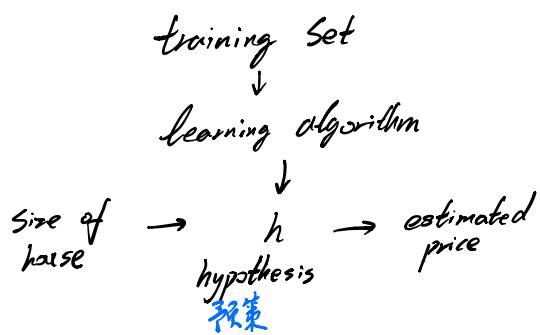
Arthur Samuel: (ML) field of study that gives computer ability to learn without being explicitly programmed.

Tom Mitchell: A computer program is said to learn from experience E with respect to some task T and some performance measured P , if its performance on T , as measured by P , improves with experience E .

ML algorithms } supervised learning
 } unsupervised learning

others: reinforcement learning, recommender system

Linear Regression Problem:



house price predict problem, supervised learning
linear regression with one variable
univariate linear regression
单变量

cost function:
损失函数

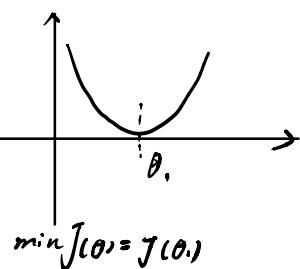
$$h_{\theta}(x) = \theta_0 + \theta_1 x, \text{ choose } \theta_0, \theta_1 \text{ to minimize } J(\theta)$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2, \text{ square error function, } m \text{ is the number of training set}$$

when θ is only one dimension.

$$h_{\theta}(x) = \theta_0 x$$

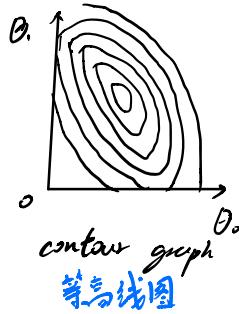
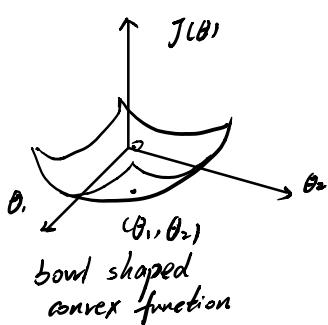
by plotting the graph of $J(\theta)$, is :



when θ is 2 dimension:

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

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



Gradient descent:

梯度下降

a more general method for calculating $\min J(\theta)$, $\theta = \langle \theta_0, \theta_1 \rangle$

outline: - start with some θ_0, θ_1

- keep changing θ_0, θ_1 to reduce $J(\theta)$, until end up at a minimum

repeat until convergence {

收敛

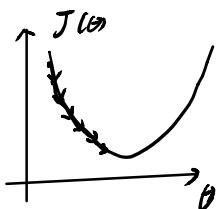
$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad j = 0, 1 \quad \alpha: \text{learning rate}$$

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$: partial derivative
同时更新 θ_0, θ_1 偏导数

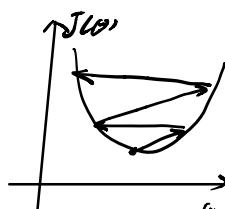
$$\left\{ \begin{array}{l} \theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 = t_0 \\ \theta_1 = t_1 \end{array} \right.$$

$$\left\{ \begin{array}{l} \theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{array} \right.$$

To choose the learning rate:

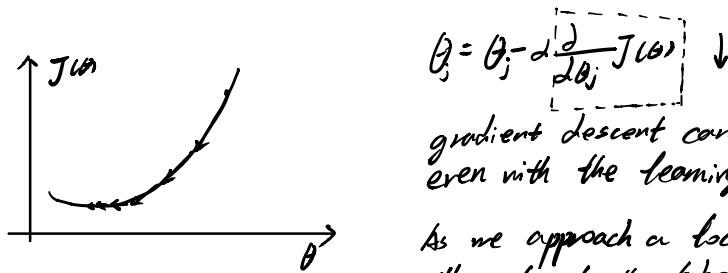
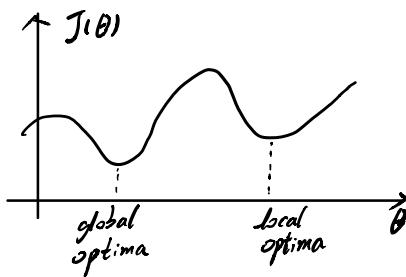


α is too small, gradient descent can be very slow



α too large, gradient descent can overshoot the minimum, it may fail to converge, or even diverge.

发散



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

gradient descent can converge to a local minimum even with the learning rate α fixed.

As we approach a local minimum, gradient descent will automatically take a smaller step, hence no need to decrease α over time.

gradient descent for linear regression:

for a 2 parameter $h_{\theta}(x)$:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

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

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

repeat until convergence :

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

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

}

2 variable linear regression's cost function $J(\theta_0, \theta_1)$ is a bowl shaped convex function, it always has a global minimum which gradient descent can converge to.

↓ is a "batch" gradient descent : $\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$
use all the training examples each step iterate

Review of linear algebra:

(only terminologies)

Matrix: rectangular array of numbers

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad 4 \times 2 \text{ matrix}$$

$$\mathbb{R}^{4 \times 2}$$

vector: $n \times 1$ matrix

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad 3 \text{ dimensional vector}$$

$$\mathbb{R}^3$$

Matrix Addition / subtraction

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 6 & 8 \\ 10 & 12 \end{bmatrix} \quad \begin{bmatrix} 3 & 4 \\ 5 & 6 \end{bmatrix} - \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

Scalar Multiplication
real number

$$3 \times \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 3 & 6 \\ 9 & 12 \end{bmatrix}$$

Matrix vector multiplication:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$$

Matrix Matrix Multiplication:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 3 & 4 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 13 & 16 \\ 29 & 36 \end{bmatrix}$$

Matrix Multiplication Properties

$A \times B \neq B \times A$ (not commutative)

Identity Matrix 单位矩阵

denoted as I or $I_{n \times n}$

$$I_{2 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$A = AI = IA$ for any matrix A

Matrix Inverse:

$$AA^{-1} = I = A^{-1}A$$

matrix that can't be inverse is "singular" or "degenerate"

matrix transpose:

$$B = A^T$$

$$b_{ij} = a_{ji}$$

Linear Regression with multiple variables (multiple features)

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

for convenience denote $x_0 = 1$

n : number of features
 x_i : feature i

$$h_{\theta}(x) = [\theta_0, \theta_1, \dots, \theta_n] \cdot \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \theta^T x \quad (\text{multivariate linear regression})$$

gradient descent:

repeat until convergence {

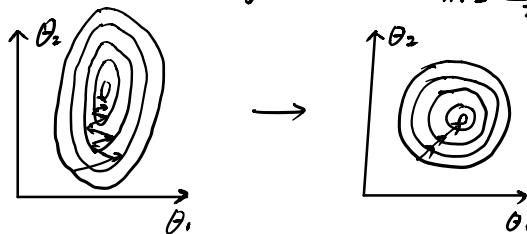
$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$j = 0, \dots, n$$

to improve the performance of gradient descent.

Feature Scaling: make sure features are on the similar scale

特征缩放
 e.g. $x_1 = \text{size} (0 \sim 2000 \text{ feet}^2) \rightarrow x_1 = \frac{\text{size}}{2000}$



get features 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 (not apply to $x_0 = 1$)

e.g. $x_1 = \text{size} (0 \sim 2000 \text{ feet}^2) \rightarrow x_1 = \frac{\text{size} - 1000}{2000} (-0.5 \sim 0.5)$

$$x \leftarrow \frac{x - \mu}{s} \quad \mu: \text{average value of } x \text{ in training set}$$

$s: \text{range, max-min}$

standardization:

标准化

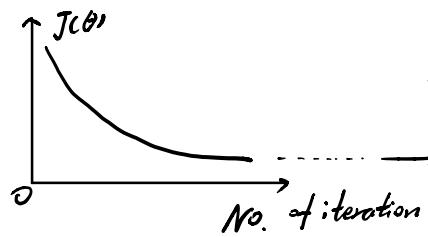
$$x \leftarrow \frac{x - \mu}{\text{std}(x)}$$

has zero-mean and unit-variance

零均值 单位方差

Debug Gradient Descent

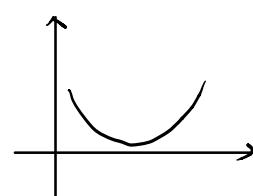
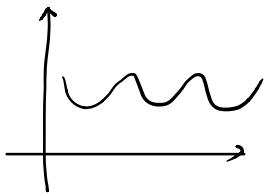
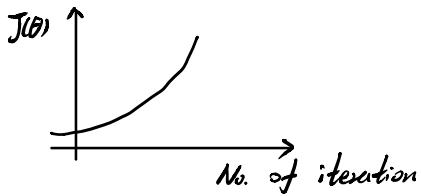
to choose the proper α to make it work correctly.



$J(\theta)$ should decrease every iteration
← example of automatic convergence test

decide convergence if $J(\theta)$ decrease by less than
② in one iteration.
 $\downarrow 10^{-3}$? hard to determine

in the following case, gradient descent works incorrectly for α is too large.



if α is small enough, $J(\theta)$ should decrease every iteration
choose α : usually from 0.001 to 1

$$\dots \underbrace{0.001}_{x^3}, \underbrace{0.003}_{x^3}, \underbrace{0.01}_{x^3}, \underbrace{0.03}_{x^3}, \underbrace{0.1}_{x^3}, \underbrace{0.3}_{x^3}, \underbrace{1}_{x^3}, \dots$$

Features And Polynomial Regression:

e.g.1: house pricing prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \underbrace{x_{frontage}}_{x_1} + \theta_2 \underbrace{x_{depth}}_{x_2} \quad \text{is linear function}$$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 \underbrace{x_{frontage}}_{x_1} \underbrace{x_{depth}}_{x_2} \\ &= \theta_0 + \theta_1 x_1 x_2 \end{aligned} \quad \text{is not linear function}$$

but if let $x = x_1 x_2$, $h_{\theta}(x) = \theta_0 + \theta_1 x$ it turns linear again

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

$$\leftarrow x_1 = x \quad 1 \sim 100$$

$$x_2 = x^2 \quad 1 \sim 10000$$

$$x_3 = x^3 \quad 1 \sim 1000000$$

it must use feature scaling or mean normalization

Normal Equation

- method to solve for θ analytically.

$$\theta \in \mathbb{R}^{n+1}, J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

solve for $\theta_0, \theta_1, \dots, \theta_n$

e.g.

$m=4$

	size	number of bedrooms	number of floors	age of home	price
x_0	x_1	x_2	x_3	x_4	x_5
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad m \times (n+1)$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \quad m \times 1$$

$$h_\theta(x) = \theta^T x$$

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

no need to choose α
 no need to iterate
 no need to do feature scaling
 $O(n^3)$, be slow if n is large

when has many features ($> 10^6$), use gradient descent

when $X^T X$ is non-invertible, two reason

1. redundant features (linearly dependent) e.g. $x_1 = kx_2 \dots$

2. too many features ($n > m$), just delete some features or use regularization.

Normal Equation Proof:

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

$$h_\theta(x) = \theta^T x = \sum_{i=0}^n \theta_i x_i$$

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

assume that we have n features and m training items
training set is:

$$X = \begin{bmatrix} x^{(1)} & \cdots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_0^{(m)} & & x_n^{(m)} \end{bmatrix}_{m \times (n+1)}$$

feature parameters:

$$\theta = [\theta_0, \theta_1, \dots, \theta_n]^T$$

output y :

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]^T$$

$$\begin{aligned} X\theta - Y &= \begin{bmatrix} (\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)}) - y^{(1)} \\ \vdots \\ (\theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)}) - y^{(m)} \end{bmatrix}_{m \times 1} = \begin{bmatrix} \theta^T x^{(1)} - y^{(1)} \\ \theta^T x^{(2)} - y^{(2)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{bmatrix}_{m \times 1} \\ &= \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}_{m \times 1} \end{aligned}$$

$$\therefore \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = (X\theta - Y)^T (X\theta - Y)$$

cost function.

$$\begin{aligned} J(\theta_0, \theta_1, \dots, \theta_n) &= \frac{1}{2m} (X\theta - Y)^T (X\theta - Y) \\ &= \frac{1}{2m} (Y^T Y + \theta^T X^T X \theta - Y^T X \theta - \theta^T X^T Y) \end{aligned}$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{2m} \left(\underset{\textcircled{1}}{\frac{\partial Y^T Y}{\partial \theta}} + \underset{\textcircled{2}}{\frac{\partial \theta^T X^T X \theta}{\partial \theta}} - \underset{\textcircled{3}}{\frac{\partial Y^T X \theta}{\partial \theta}} - \underset{\textcircled{4}}{\frac{\partial \theta^T X^T Y}{\partial \theta}} \right)$$

$$\therefore \frac{dAB}{d\theta} = A^T \quad \frac{dX^TAX}{d\theta} = 2AX$$

$$\textcircled{1} \quad \frac{\partial Y^TY}{\partial \theta} = 0$$

$$\begin{aligned}\textcircled{2} \quad \therefore Y^T X \theta &= [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \begin{bmatrix} x_0^{(1)} & \dots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_0^{(m)} & \dots & x_n^{(m)} \end{bmatrix} [\theta_0, \theta_1, \dots, \theta_n]^T \\ &= (x_0^{(1)}y^{(1)} + \dots + x_0^{(m)}y^{(m)})\theta_0 + (x_1^{(1)}y^{(1)} + \dots + x_1^{(m)}y^{(m)})\theta_1 + \dots \\ &\quad + (x_n^{(1)}y^{(1)} + \dots + x_n^{(m)}y^{(m)})\theta_n\end{aligned}$$

$$\therefore \frac{\partial Y^T X \theta}{\partial \theta} = \begin{bmatrix} \frac{\partial Y^T X \theta}{\partial \theta_0} \\ \vdots \\ \frac{\partial Y^T X \theta}{\partial \theta_n} \end{bmatrix} = \begin{bmatrix} x_0^{(1)}y^{(1)} + \dots + x_0^{(m)}y^{(m)} \\ \vdots \\ x_n^{(1)}y^{(1)} + \dots + x_n^{(m)}y^{(m)} \end{bmatrix} = X^T Y$$

$$\begin{aligned}\textcircled{3} \quad \theta^T X^T Y &= [\theta_0, \theta_1, \dots, \theta_n] \begin{bmatrix} x_0^{(1)} & \dots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_0^{(m)} & \dots & x_n^{(m)} \end{bmatrix}^T [y^{(1)}, y^{(2)}, \dots, y^{(m)}]^T \\ &= (x_0^{(1)}\theta_0 + \dots + x_n^{(1)}\theta_n)y^{(1)} + (x_0^{(2)}\theta_0 + \dots + x_n^{(2)}\theta_n)y^{(2)} \\ &\quad + \dots + (x_0^{(m)}\theta_0 + \dots + x_n^{(m)}\theta_n)y^{(m)}\end{aligned}$$

$$\frac{\partial \theta^T X^T Y}{\partial \theta} = \begin{bmatrix} \frac{\partial \theta^T X^T Y}{\partial \theta_0} \\ \vdots \\ \frac{\partial \theta^T X^T Y}{\partial \theta_n} \end{bmatrix} = \begin{bmatrix} x_0^{(1)}y^{(1)} + \dots + x_0^{(m)}y^{(m)} \\ \vdots \\ x_n^{(1)}y^{(1)} + \dots + x_n^{(m)}y^{(m)} \end{bmatrix} = X^T Y$$

$$\textcircled{4} \quad \theta^T X^T X \theta = (X^T X)(\theta_0^2 + \theta_1^2 + \dots + \theta_n^2)$$

$$\frac{\partial \theta^T X^T X \theta}{\partial \theta} = \begin{bmatrix} \frac{\partial \theta^T X^T X \theta}{\partial \theta_0} \\ \vdots \\ \frac{\partial \theta^T X^T X \theta}{\partial \theta_n} \end{bmatrix} = 2(X^T X) \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} = 2X^T X \theta$$

$$\therefore \frac{\partial J(\theta)}{\partial \theta} = \frac{1}{2m} (-2X^T Y + 2X^T X \theta) = 0$$

$$\therefore X^T X \theta = X^T Y$$

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

Another Proof Not Using Matrix Calculus

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

$$\left\{ \begin{array}{l} \frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} = 0 \quad \textcircled{i} \\ \frac{\partial J(\theta)}{\partial \theta_n} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_n^{(i)} = 0 \quad \textcircled{n} \end{array} \right.$$

$$\textcircled{i}: \frac{1}{m} [x_0^{(1)} \ x_0^{(2)} \ \dots \ x_0^{(m)}] \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix} = 0$$

$$[x_0^{(1)} \ x_0^{(2)} \ \dots \ x_0^{(m)}] \left(\begin{bmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} \theta - y \right) = 0$$

$$[x_0^{(1)} \ x_0^{(2)} \ \dots \ x_0^{(m)}] (X\theta - y) = 0$$

$$\therefore [x_1^{(1)} \ x_1^{(2)} \ \dots \ x_1^{(m)}] (X\theta - y) = 0$$

⋮

$$[x_n^{(1)} \ x_n^{(2)} \ \dots \ x_n^{(m)}] (X\theta - y) = 0$$

$$\therefore \begin{bmatrix} x_0^{(1)} & x_0^{(2)} & \dots & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(m)} \end{bmatrix} (X\theta - y) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\| X^\top (X\theta - y) = 0$$

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

$$X^\top X\theta = X^\top y$$

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

vectorized :

e.g.

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i$$

formula :



$$h_{\theta}(x) = \theta^T x$$

Octave code.

```
prediction = 0.0  
for i = 1: n+1  
    prediction = prediction + theta(i) * x(i),  
end .
```



```
prediction = theta' * x
```

vectorize can use Gpu to accelerate running. and shorten the code.