\* Model representation

Example

**Housing Prices**
**(Portland, OR)**



Price
230
(in 1000s
of dollars)

1250 → predict
Size (feet²)

Supervised Learning | Regression problem
given "right answer" | Predict real-value output

training set

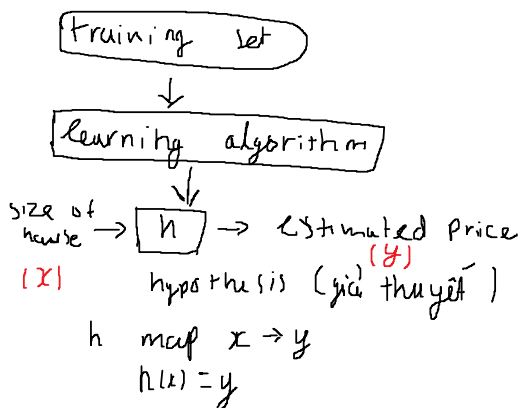| Size of feet² (x) | price (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1536 | 315 |
| 852 | 178 |

$m$ = number of training example
$x$'s = Input features
$y$'s = Output features

$x^{(i)}_{1}, y^{(i)}_{1}$ - training example $i$ th

How ML Work

$\boxed{\text{training set}}$
↓
$\boxed{\text{learning algorithm}}$
↓
size of house → $\boxed{h}$ → Estimated Price
[x]     hypothesis (giả thuyết)    (y)
    h map $x \to y$
    $h(x) = y$

represent h

$h_\theta(x) = \theta_0 + \theta_1 x$     $(f(x) = ax + b)$
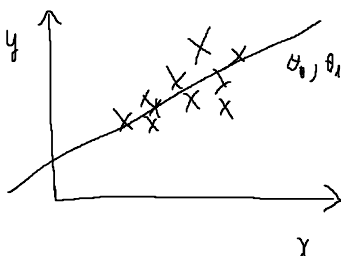


$h_{(x)} = \theta_0 + \theta_1 x$
linear

linear regression with one variable
univariable linear regression

---

\* Cost function

$h_\theta(x) = \theta_0 + \theta_1 x$
    \     /
    Parameter
    (tham số)

$\underset{\theta_0, \theta_1}{minimize} \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Kết quả            Thực tế
dự đoán



$J(\theta_0, \theta_1) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left( h_{(\theta)}(x^{(i)}) - y^{(i)} \right)^2$

Cost function

Square error function

$x$

ý tưởng : tìm tham số $\theta_0$, $\theta_1$ cho
$h_\theta(x)$ gần với $y$ của các mẫu huấn
luyện $(x, y)$

Square error function

Simplified
$$h_\theta(x) = \theta_1 x \qquad \theta_0 \equiv 0$$
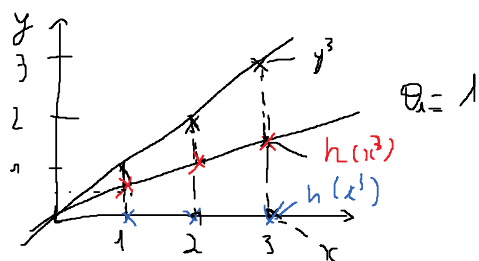
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^i) - y^i \right)^2$$
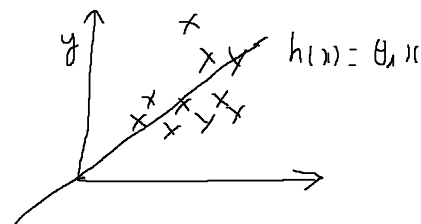
minimize $\quad J(\theta_1)$
$\theta_1$

$y$

$h(x) = \theta_1 x$

$h_\theta(x)$
cố định $\theta_1$, $x$ tham số

$\theta_1 = 1$

$h(x^3)$
$h(x^1)$

Với $\theta_1 = 1$
$$J(1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h(x^i) - y^i \right)^2$$
$$= \frac{1}{6} \cdot \left( (1-1)^2 + (2-2)^2 + (3-3)^2 \right) = 0$$

$J(\theta_1)$
$\theta_1$ là tham số

$J(\theta_1)$

$\theta_1 = 1$

Parabol

$J(1) = 0$

$J(0,5) = 0,58$

$J(0) = 2,167$

minimize $J(\theta_1)$
$\theta_1$

for $h_\theta(x) = \theta_0 + \theta_1 x$
$h_\theta(x)$ $x$ tham số

Price

300
200
100

$h_\theta(x)$

$\theta_0 = 50$
$\theta_1 = 0,06$

$J(\theta_0, \theta_1)$. $\theta_0$, $\theta_1$ tham số

100
75
50
25
0

$J(\theta_0, \theta_1)$

10

10

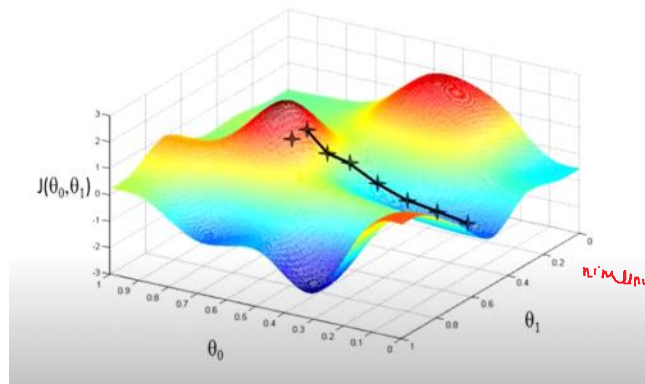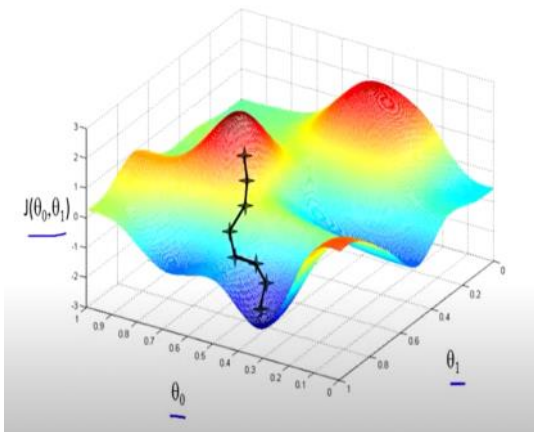$$h_{\theta(x)} = 50 + \theta_{1}.06 x$$

### $h_\theta(x)$
(for fixed $\theta_0, \theta_1$, this is a function of x)



### $J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)
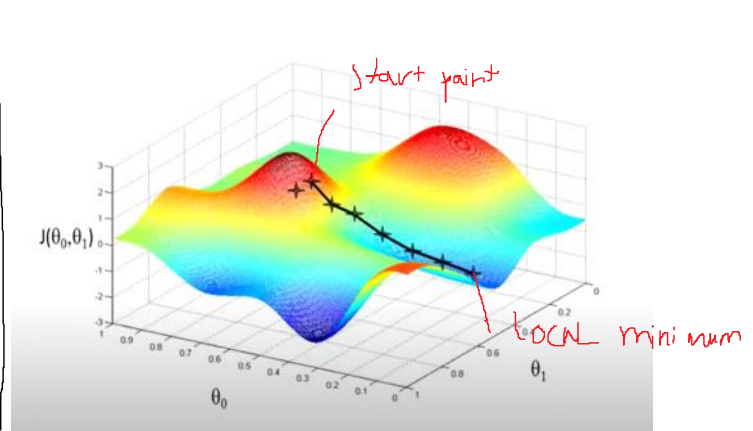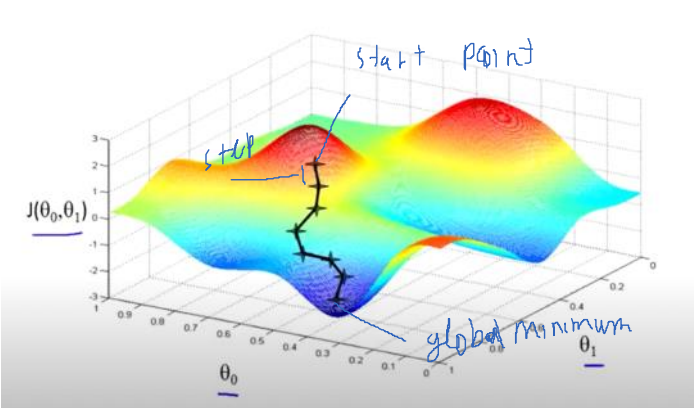


same value of $J(\theta_0, \theta_1)$

minimum





minimum

\* Gradient Descent

ý tưởng : có hàm $J(\theta_0, \theta_1)$ , muốn $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

— khởi đầu với $\theta_0, \theta_1$ bất kì

— thay đổi $\theta_0, \theta_1$ để giảm $J(\theta_0, \theta_1)$ tới khi đạt được minimum



start point

step

$J(\theta_0, \theta_1)$

$\theta_0$   $\theta_1$   global minimum



start point

$J(\theta_0, \theta_1)$

$\theta_0$   $\theta_1$   local minimum

⟩ phụ thuộc vào start point (tính chất của GD)

thuật toán

$$\theta_j = \theta_1 - \textcolor{red}{\boxed{\alpha}} \frac{\partial}{\partial \theta_j} . J(\theta_0, \theta_1) \qquad \text{for } j = 0 \text{ and } j = 1$$

learning rate

— Lặp lại cho tới khi $\theta_1$ không đổi ⟨⟩ $\boxed{\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} = 0}$ ⟨⟩ cực tiểu

— Cập nhật đồng thời các tham số

Simplifiend

$\min\limits_{\theta_1} J(\theta_1)$   ⟨⟩   $\theta_0 = 0$

for derivate



$J$

$J(\theta_1)$   $\theta_1 \in \mathbb{R}$

$\theta_1$

$$\theta_1 = \theta_1 - \alpha \frac{\partial J(\theta_1)}{\partial \theta_1}$$

⟨⟩ $\theta_1 - \alpha ( \text{số dương} ) \Rightarrow \theta_1$ giảm



$J$

$$\theta_1 = \theta_1 - \alpha \frac{\partial J(\theta_1)}{\partial \theta_1}$$

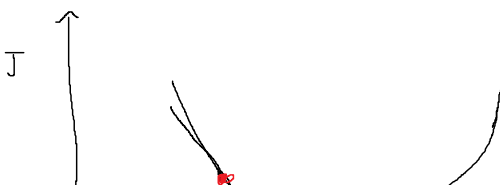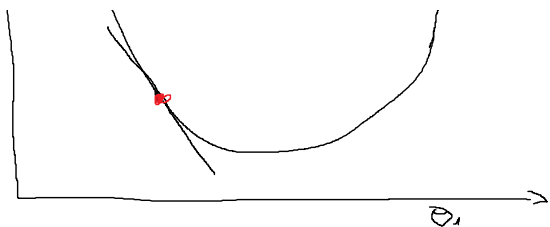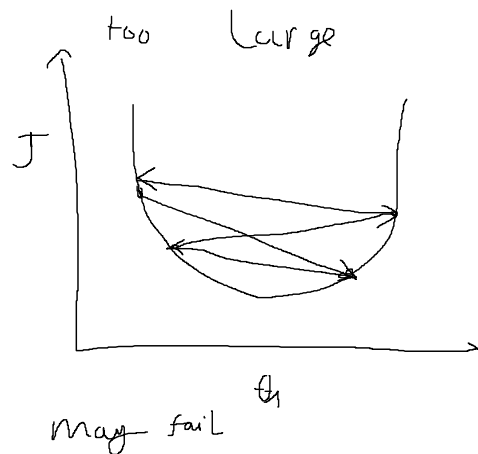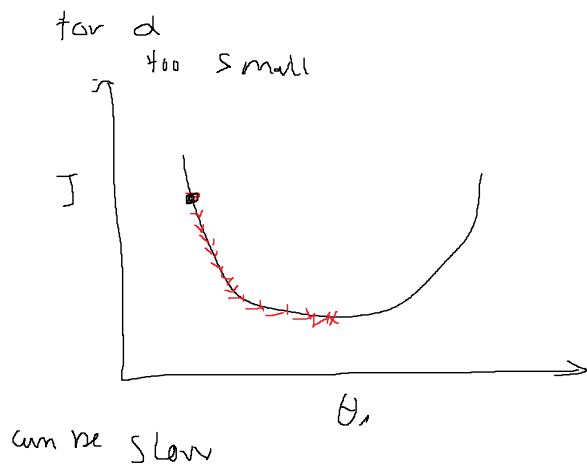$$\theta_1 = \theta_1 - \alpha \frac{\partial \ J(\theta_1)}{\partial \theta_1}$$

$$\Leftrightarrow \theta_1 - \alpha \cdot (\text{số âm})$$

$$\Rightarrow \theta_1 \ \text{tăng}$$

for $\alpha$ too small



cần be slow

too Large



may fail

Khi gần điểm uc tiểu thì các bước càng nhỏ nên không cần cập nhật lại $\alpha$

---

Gradient descent for Linear regression

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x_j^i) - y^i \right)^2$$

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

$\theta_0 \quad \frac{\partial}{\partial \theta_1} \quad = \frac{1}{m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 x - y^i \right)$

$\theta_1 \quad \frac{\partial}{\partial \theta_1} J \quad = \frac{1}{\cancel{2}m} \sum_{j}^{m} \cancel{2} \cdot x^i \left( \theta_0 + \theta_1 x^i - y^i \right)$

Gradient descent Algorithms

$$\theta_0 = \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1(x^i) - y^i \right)$$

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

Lặp Lại tới khi $\theta_0, \theta_1$ không đổi

Normal   Gradient Linear descent

Linear regression   Gradient descent

Convex

$J(\theta_0, \theta_1)$

$h_\theta(x)$
(for fixed $\theta_0, \theta_1$, this is a function of x)

Price $ (in 1000s)

1250

Size (feet²)

Training data
Current hypothesis

$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

other name : batch Gradient descent

## one variable

| Size x | Price y |
|--------|---------|
| 2014   | 460     |
| 1416   | 232     |
| 1552   | 301     |

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

## Multiple Var

| Size $x_1$ | No bedroom $x_2$ | No floor $x_3$ | age of home $x_4$ | price y |
|------------|------------------|----------------|-------------------|---------|
| 2864       | - - -5           | -1 --          | 115 - -           | 460     |
| 1416       | 3                | 2              | 40                | 232     |
| 1534       | 3                | 1              | 30                | 315     |

Note:

$n$ = number of feature

$x^i$ = Input of $i^{th}$ training exampl

$x_j^i$ = value of feature $j$ in $i^{th}$

Ex; $x^2 = \begin{bmatrix} 1414 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^4$

$x_4^2 = 40$

## hypothesis.

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

define $x_0 = 1 \iff x_0^i = 1$

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_1 \\ -- \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ --- \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

multivariable Linear regression

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

tham số: $\theta_0, \theta_1, \theta_2, \dots$ $\Rightarrow \theta$

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

Gradient descent

Lặp lại tối khi tham số không đổi

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

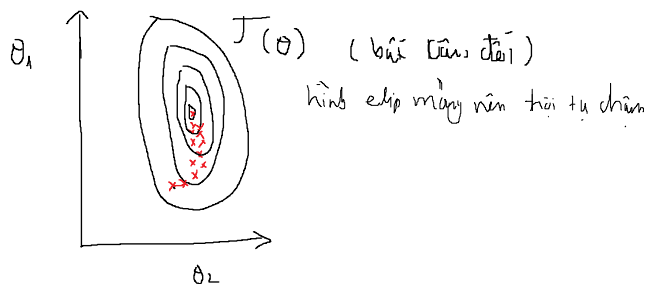$$= \theta_j - \alpha . \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i) . x_j$$

trick feature scaling
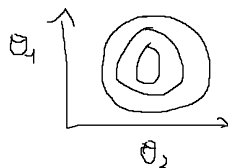
ý tưởng: biết miền giá trị của các đặc trưng

VD:
$x_1$ = size $(0 - 2000)$
$x_2$ = No bedroom $(1 - 5)$



$\theta_1$ $J(\theta)$ (bất đâu đôi)

hình elip mỏng nên tới tụ chậm

$$x_1 = \frac{size}{2000}$$

$$x_2 = \frac{No\ bedroom}{5}$$

(cân đôi)

hình tròn nên đi đến điểm cực tiểu nhanh hơn

làm cho các đặc trưng có giá trị trong khoảng $-1 \le x_i \le 1$

Mean normalization

thường là giá trị trung bình

ý tưởng: thay $x_i$ bằng $x_i - u_i$ để làm đặc trưng xấp xỉ 0

VP:
$$x_1 = \frac{x_1 - 1000}{2000} \qquad -0,5 \le x_1 \le 0,5$$

$$x_2 = \frac{x_2 - 2}{5} \qquad -0,5 \le x_2 \le 0,5$$

$$\boxed{x_j = \frac{x_i - avg}{Max - Min}}$$

Debugging

$$\min_\theta J(\theta)$$

$\min_{\theta} J(\theta)$

(graph with x-axis: 100, 50 -100 Số vòng lặp, 300, 400)

$J_{(\theta)}$ hội tụ nếu $J(\theta)$ giảm $< 10^{-3}$ trong 1 lần lặp

nếu Gradient descent không hoạt động tốt:
- giảm $\alpha$
- $\alpha$ quá nhỏ làm Gradient descent bị chậm
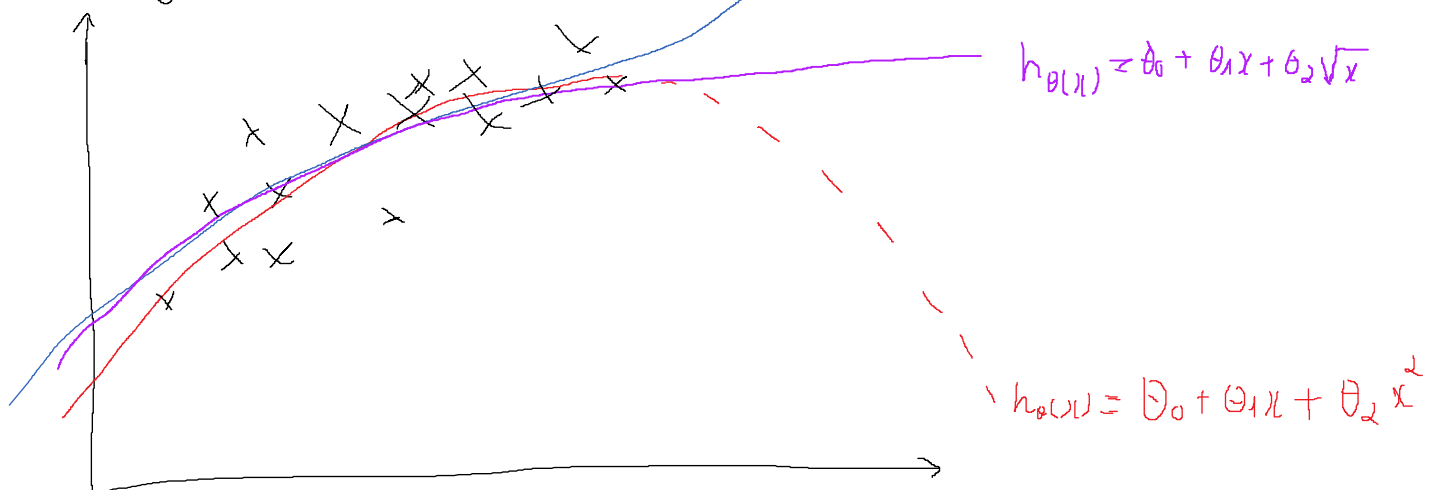
Lựa chọn feature và hồi quy đa thức

Example

dự đoán giá nhà

$h_\theta(x) = \theta_0 + \theta_1 \cdot front + \theta_2 \cdot depth$

$Area = front \cdot depth$

$\Rightarrow h_\theta(x) = \theta_0 + \theta_1 \cdot Area$

hồi quy đa biến



$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

áp dụng vào mutiple variable

$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$

$$= \theta_0 + \theta_1 \, size + \theta_2 \, size^2 + \theta_3 \, size^3$$

$x_1 = size$
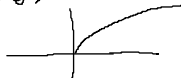
$x_2 = size^2$

$x_3 = size^3$

chú ý Fearture scaling để tối ưu

thay vì bậc ba thì có thể chọn

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$

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

đặt $X = [1, x_1, x_2, \dots, x_n] \in R^{m \times (n+1)}$

$1 \times (n+1)$

$\theta = [\theta_0, \theta_1, \theta_2, \dots, \theta_n] \in R^{1 \times (n+1)}$

$1 \times (n+1)$

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

$$J_\theta = \frac{1}{2m} \sum_{j=1}^m (h(x^i - y^j))^2 = \frac{1}{2m} \|X\theta^T - y\|_2^2$$

$$\frac{\partial}{\partial \theta_n} J_\theta = \frac{1}{2m} \sum_{j=1}^m (h_\theta(x^i) - y^i) x_n$$

$$\frac{\partial J_\theta}{\partial \theta} = X^T(X\theta^T - y) = 0$$

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

$\Leftrightarrow \quad X^T X \theta^T = X^T y$

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

ma trận nghịch đảo của $X^T X$

đk $X^T X$ khả nghịch $\Leftrightarrow$ det $\neq 0$

so sánh với Gradient descent

**G D**
- cần chọn $\alpha$
- nhiều vòng lặp
- hoạt động tốt ngay cả khi có nhiều feature
- $O(Kn^2)$

**N E**
- không cần $\alpha$
- không cần vòng lặp
- chậm nếu có nhiều feature
- $X^T X \quad O(n^3)$

$*$ nếu $X^T X$ không khả nghịch

- có feature bị thừa $\Rightarrow$ hai cột, hàng tỉ lệ $\Rightarrow$ det $= 0$
- có quá nhiều feature và ít train set

$X$ Like $\begin{bmatrix} x_0^1 & x_1^1 & x_2^1 & \cdots & x_n^1 \\ x_0^2 & x_1^2 & x_3^1 & \cdots & x_n^2 \\ \cdots & & & & \\ x_0^m & \cdots & & & x_n^m \end{bmatrix}$

$m \times (n+1)$

$\theta$ like $[\theta_0 \; \theta_1 \; \theta_2 \; \cdots \; \theta_n]$

$1 \times (n+1)$

$Y$ like $\begin{bmatrix} y^1 \\ y^2 \\ y^3 \\ y^n \end{bmatrix}$ $m \times 1$