

4. Linear Regression with multiple variables



4.1 Multiple features

Multiple features (variables).

Size (feet ²) x_1	Number of bedrooms x_2	Number of floors x_3	Age of home (years) x_4	Price (\$1000) y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$m=47$
example 4

Notation

n = feature 4

$x^{(i)}$ = i th training example

$x_j^{(i)}$ = i th training example j th feature

ex) $x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$, $x_3^{(2)} = 2$

hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (x_0^{(i)} = 1)$$

$$= \theta^T \cdot x$$

$$\hookrightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\Rightarrow \underbrace{[\theta_0 \ \theta_1 \ \dots \ \theta_n]}_{\theta^T} \cdot \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}}_x = \theta^T \cdot x$$

4. Linear Regression with multiple variables



4.2 Gradient descent for multiple variables

Cost function:

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

Gradient descent:

$$\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta) \quad (\text{simultaneously updated})$$

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

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \end{cases}$$

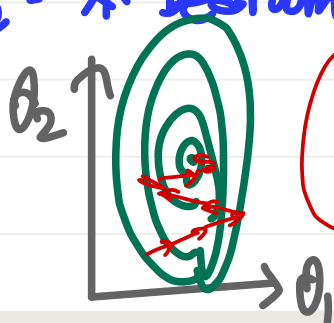
4.3 Gradient descent in practice I: Feature Scaling

make sure features are on a similar scale.

ex) $x_1 = \text{size (0 ~ 2000 feet}^2)$ $\rightarrow x_1 = \frac{\text{size (feet}^2)}{2000}$

$x_2 = \text{# bedrooms (1 ~ 5)}$ $\rightarrow x_2 = \frac{\text{\# bedrooms}}{5}$

$$-1 \leq x_i \leq 1$$



보통, 높은 차원차량
이런 시간 동안
앞뒤로 진동 \Rightarrow 최솟값 도달



더 균형적인 원형
빠르게 수렴,
안정적인 경로

4. Linear Regression with multiple variables



Feature scaling 함으로써 Gradient descent 더 빨리 되고
적은 수의 반복으로 수렴!

- 정확하게 같은 값이 나오지 않아도 괜찮! (충분히 근접하면 OK)

ex) $-3 \leq x_1 \leq 3$ (OK)

$-100 \leq x_3 \leq 100$ (X)

$-\frac{1}{3} \leq x_2 \leq \frac{1}{3}$ (OK)

$-0.001 \leq x_4 \leq 0.001$ (X)

• Mean normalization

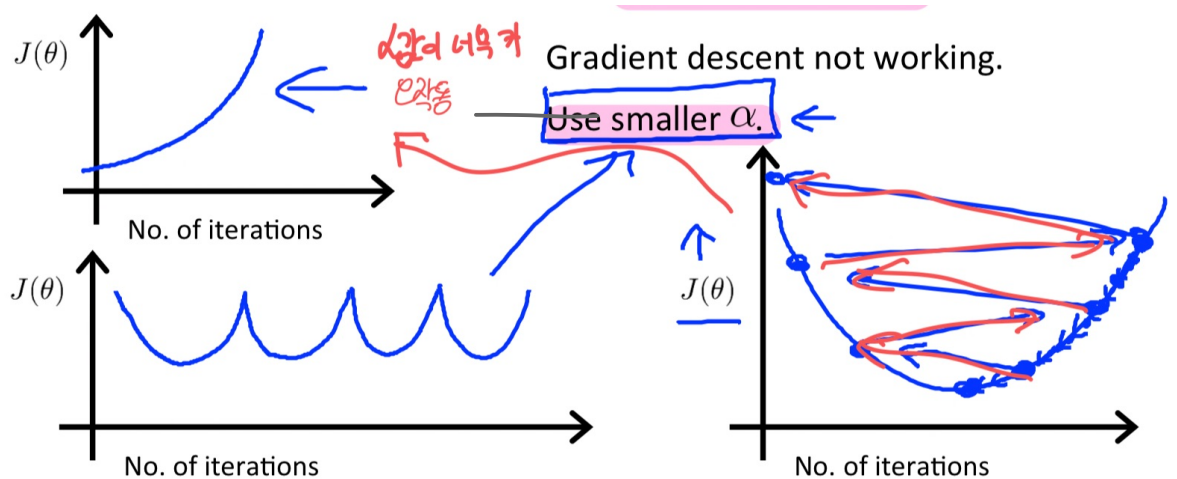
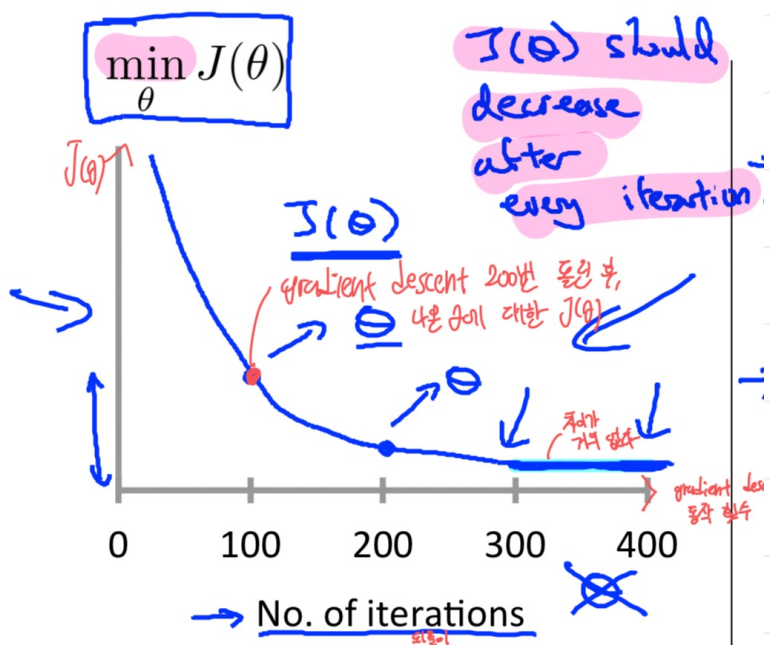
ex) $x_1 = \frac{\text{size} - 1000}{2000 - 0}$ (average)

$x_2 = \frac{\text{rooms} - 2}{5 - 1}$

range (max-min)
or S (표준편차)

$$\frac{x_i - \mu_i}{s_i}$$

4.4 Gradient descent in practice II: Learning rate (α)



α 가 너무 작으면 \rightarrow slow convergence

α 가 너무 크면 \rightarrow not decrease on every iteration, may not converge

교수님 $\rightarrow \alpha$ 3배씩 늘리며 $J(\theta)$ 값도 관찰 \rightarrow 적절한 α 찾기

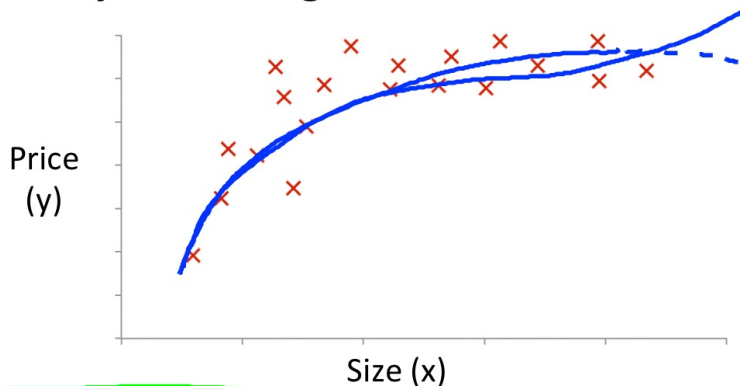
4. Linear Regression with multiple variables



4.5 Features and polynomial regression

(다항 회귀, 비선형)

Polynomial regression



$\theta_0 + \theta_1 x + \theta_2 x^2$ (2차식)
 $\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$ (3차식)

자동으로 feature
선택해주는 알고리즘 존재



선택한 feature 선택하는
다양한 방법이 있음

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

$$= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

$\rightarrow x_1 = (\text{size})$
 $\rightarrow x_2 = (\text{size})^2$
 $\rightarrow x_3 = (\text{size})^3$

이렇게 직용하면
 feature scaling의 과정이 중요해짐
 (매우 다른 범위를 가지게 되므로 → 비슷한 범위로 바꿔주는 것 필요)

Size: $1 \sim 1000$ (orange)
 Size²: $1 \sim 1,000,000$
 Size³: $1 \sim 10^9$

Andrew Ng

4.6 Normal equation (정규 방정식)

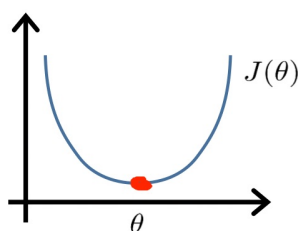
Gradient Descent에서 θ 의 최적값 (반복없이) 한 번에 구할 수 있다!

Intuition: If 1D ($\theta \in \mathbb{R}$) scalar (not vector)

$\rightarrow J(\theta) = a\theta^2 + b\theta + c$

$\frac{d}{d\theta} J(\theta) = \dots$ set to 0

Solve for θ



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

↑ cost fn 최적화 하는 값

* θ 가 vector 일 경우

$\theta \in \mathbb{R}^{n+1}$

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

Examples: $m = 4$.

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

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

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

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

cost fn 최적화 하는 값

m -dimensional vector

Andrew Ng

4. Linear Regression with multiple variables



m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n features.

$$\underline{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

(design matrix)

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$n \times (n+1)$

E.g. If $\underline{x}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times 2$

Gradient Descent

- θ 선택해야 함
- 반복 필요
- n 이 커도 작동함 $O(n^2)$

* feature 수 1000보다 작은 경우 효과적!

Normal Equation

- θ 선택할 필요 X
- 반복 필요 X
- 계산 필요
- n 이 커지면 느려짐 $O(n^3)$

4.7 Normal equation and non-invertibility

$X^T X$: 역행렬 없는 경우 (매우 드물다)

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

- ① Redundant features \rightarrow 불필요한 feature 없애기
한번 후, delete
- ② Too many features \rightarrow Delete some features
 \rightarrow regularization