



[SWCON253] Machine Learning – Lec.06

Linear Regression + α

Fall 2025

김휘용

hykim.v@khu.ac.kr



경희대학교

KYUNG HEE UNIVERSITY

Contents

1. Linear Regression
2. Normal Equation
3. Polynomial Regression
4. For Better Results
 - Feature Normalization
 - Learning Rate Selection

References

- *Machine Learning* by Andrew Ng, Coursera (<https://www.coursera.org/learn/machine-learning>)

1. Linear Regression

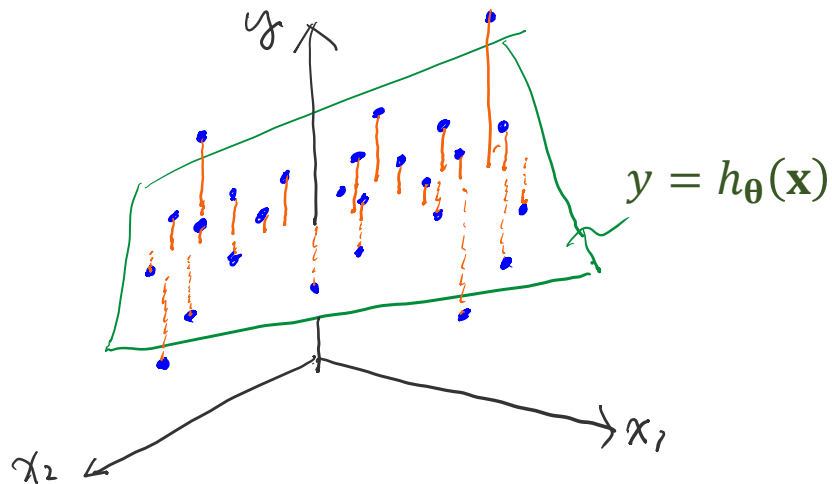
- ✓ Multivariate Linear Regression
- ✓ Data Representation
- ✓ Linear Model Representation
- ✓ Cost Function (& Gradient)
- ✓ Parameter Update (by Gradient Descent)

o. Multivariate Linear Regression

◆ Find the best linear function h_{θ} for the given training dataset \mathbb{D} with *multiple*(n)-features

- A feature vector: $\mathbf{x} = [x_1, \dots, x_n]^T$
- Training dataset: $\mathbb{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^M = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(M)}, y^{(M)})\}$
- **Linear model:**

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= h_{\theta}(x_1, \dots, x_n) \\ &= \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \end{aligned}$$



Ex) Housing Price Prediction

- *Multiple features*: size, # bedrooms, # floors, age
- *Single output*: the price of a house

→ Size (feet²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
→ 1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

- n = number of features $n = 4$
- $x^{(i)}$ = input (features) of i^{th} training example.
- $x_j^{(i)}$ = value of feature j in i^{th} training example.

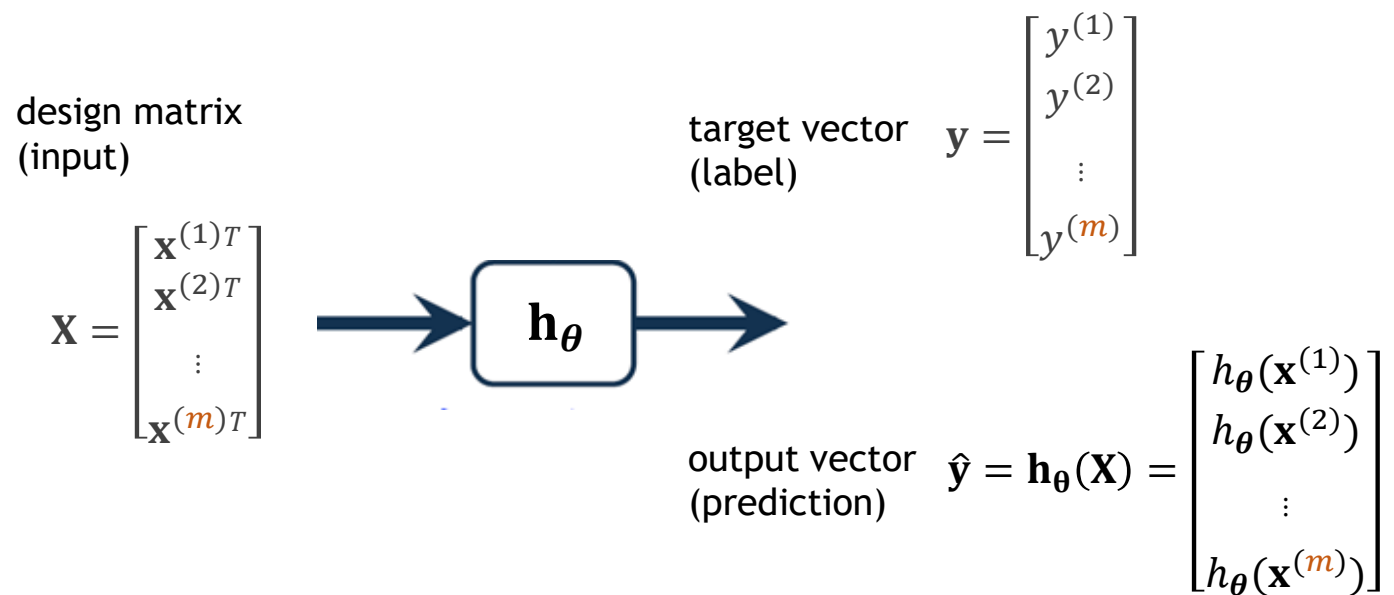
1. Data Representation

◆ Data with Multiple Features

- A feature vector: $\mathbf{x} = [x_1, \dots, x_n]^T$
- Training dataset: $\mathbb{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^M = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(M)}, y^{(M)})\}$

◆ For batch processing

- (mini-)batch size: $1 \leq m \leq M$



- $m=1$: online mode
- $m=M$: batch mode
- $1 < m < M$: mini-batch mode

- Shuffle the training example then apply minibatch mode
- Or randomly select training examples for each minibatch.

2. Linear Model Representation

◆ Representation 1

- Let $\mathbf{x} \triangleq [x_0, x_1, \dots, x_n]^T$, $\boldsymbol{\theta} \triangleq [\theta_0, \theta_1, \dots, \theta_n]^T$
★ where $x_0 = 1$.

- Then, for a **single** training example (i.e., $\mathbf{x}^{(i)}$):

$$\boxed{h_{\theta}(\underline{x})} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$
$$= [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \boxed{= \underline{\theta}^T \underline{x}}$$

- and for a **batch** of training examples (i.e., \mathbf{X}):

$$\boxed{\mathbf{h}_{\theta}(\mathbf{X}) = \begin{bmatrix} \boldsymbol{\theta}^T \mathbf{x}^{(1)} \\ \boldsymbol{\theta}^T \mathbf{x}^{(2)} \\ \vdots \\ \boldsymbol{\theta}^T \mathbf{x}^{(m)} \end{bmatrix} = \mathbf{X}\boldsymbol{\theta}}$$

◆ Representation 2 (weight & bias)

- Let $\mathbf{x} \triangleq [x_1, \dots, x_n]^T$, $\boldsymbol{\theta} \triangleq [\theta_1, \dots, \theta_n]^T$

- Then, for a **single** training example

$$h_{\theta}(\underline{x}) = \underline{\theta}^T \underline{x} + \theta_0$$

$$(h_w(\underline{x}) = \underline{w}^T \underline{x} + b)$$

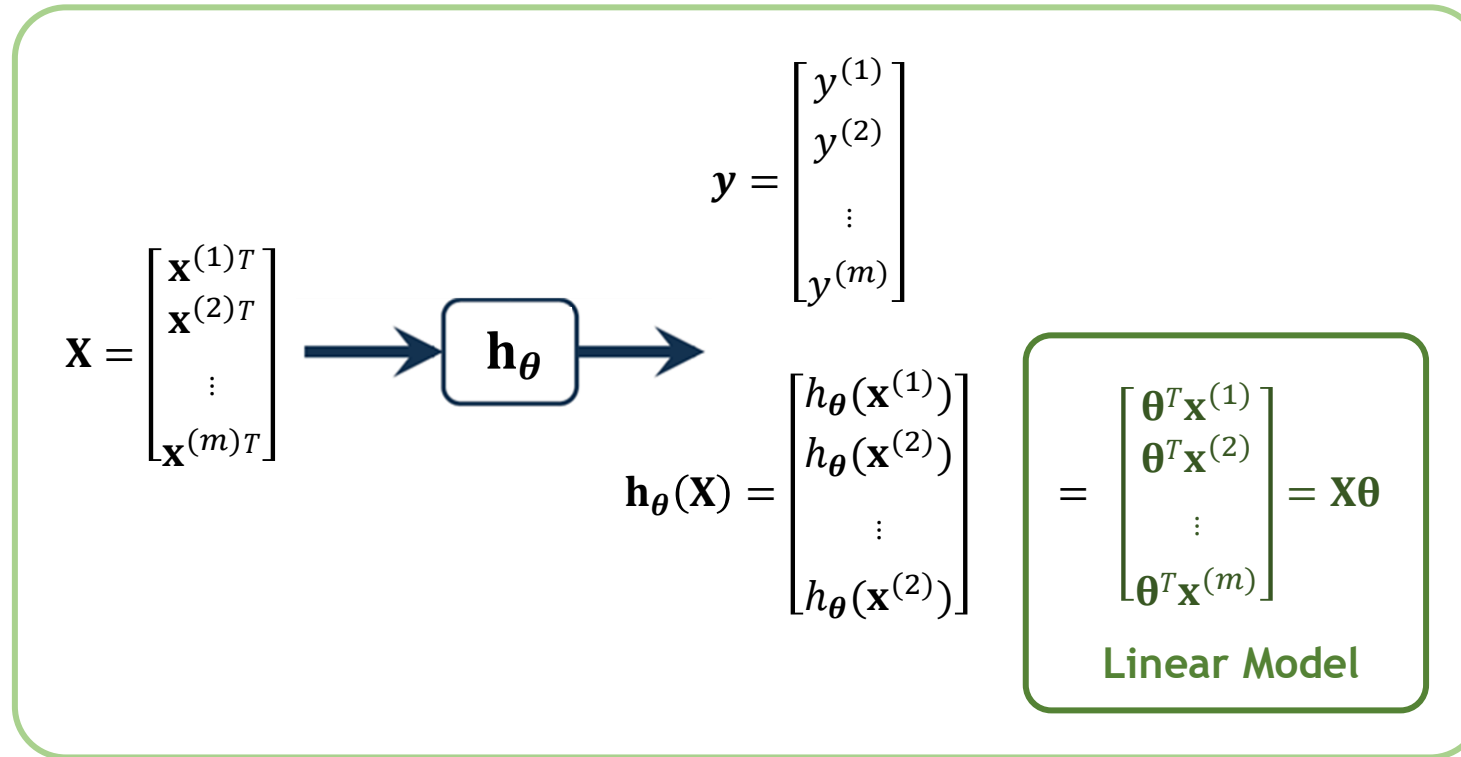
- and, for a **batch** of training example

$$h_{\theta}(\mathbf{X}) = \mathbf{X}\underline{\theta} + \theta_0 \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$(h_w(\mathbf{X}) = \mathbf{X}\underline{w} + \underline{b})$$

2. Linear Model Representation (cont'd)

◆ Summary of Input, Output, & Model



3. MSE Cost for Linear Model

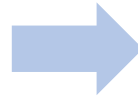
◆ MSE Cost

- Classic form:

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{2m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \end{aligned}$$

- Vector form

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{2m} \|\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{y}\|_2^2 \\ &= \frac{1}{2m} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 \\ &= \frac{1}{2m} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \end{aligned}$$



유도과정:
다음 슬라이드

◆ Gradient of the MSE Cost

- Classic form:

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

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_j^{(i)}$$

- Vector form

$$\begin{aligned} \nabla J(\boldsymbol{\theta}) &= \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \\ &= \frac{1}{m} \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \end{aligned}$$

3. MSE Cost for Linear Model (cont'd)

◆ $\nabla_{\underline{\theta}} \|x_{\underline{\theta}} - y\|_2^2$ 구하기 :

$$\begin{aligned} \|\underline{x}_0 - \underline{y}\|_2^2 &= (\underline{x}_0 - \underline{y})^T (\underline{x}_0 - \underline{y}) \stackrel{(1)}{=} (\underline{\theta}^T \underline{x}^T - \underline{y}^T) (\underline{x}_0 - \underline{y}) \\ &= \underline{\theta}^T \underline{x}^T \underline{x}_0 - \underline{y}^T \underline{x}_0 - \underline{\theta}^T \underline{x}^T \underline{y} + \underline{y}^T \underline{y} \\ &\stackrel{(2)}{=} \underline{\theta}^T \underline{x}^T \underline{x}_0 - 2 \underline{y}^T \underline{x}_0 + \underline{y}^T \underline{y} \end{aligned}$$

$$\textcircled{1} \quad \begin{cases} (\underline{a} - \underline{b})^T = \underline{a}^T - \underline{b}^T \\ (A\underline{x})^T = \underline{x}^T A^T \end{cases}$$

$$\textcircled{2} \quad \underline{\theta}^T \underline{X}^T \underline{y} = (\underline{x_0})^T \underline{y} = \underline{y}^T \underline{x_0}$$

$$\begin{aligned} \Rightarrow \nabla_{\underline{\theta}} \|\underline{X}\underline{\theta} - \underline{y}\|_2^2 &= \nabla_{\underline{\theta}} (\underline{\theta}^T \underline{X}^T \underline{X} \underline{\theta}) - 2 \nabla (\underline{y}^T \underline{X} \underline{\theta}) \\ &\stackrel{(3)}{=} 2 \underline{X}^T \underline{X} \underline{\theta} - 2 \underline{X}^T \underline{y} \\ &= \underline{2 \underline{X}^T (\underline{X} \underline{\theta} - \underline{y})} \end{aligned}$$

$$\textcircled{3} \quad \left\{ \begin{array}{l} \nabla_{\underline{x}} (\underline{x}^T A^T A \underline{x}) = 2 A^T A \underline{x} \\ \nabla_{\underline{x}} (\underline{b}^T A \underline{x}) = \nabla_{\underline{x}} \underline{c}^T \underline{x} = \underline{c} \end{array} \right.$$

4. Parameter Update by GD (for *Linear Model*)

◆ Gradient Descent for *Linear Regression with MSE Cost*

● Classic form:

Repeat until convergence {

Update $\nabla \theta_j$'s *simultaneously*:

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_j^{(i)}$$

}

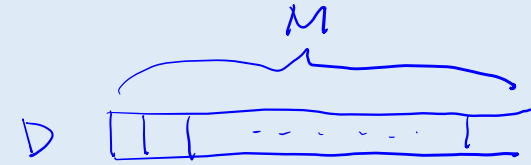
● Vector form

Repeat until convergence {

$$\theta := \theta - \alpha \nabla J(\theta) = \theta - \alpha \frac{1}{m} \mathbf{X}^T (\mathbf{X}\theta - \mathbf{y})$$

}

Learning Modes (Recap.)



- $m=1$: online mode
- $m=M$: batch mode
- $1 < m < M$: mini-batch mode

- Shuffle the training example then apply minibatch mode
- Or randomly select training examples for each minibatch.

2. Normal Equation

- ✓ Normal Equation
- ✓ Gradient Descent vs. Normal Equation

Normal Equations

◆ Analytic Solution to *Linear Regression with MSE Loss*

$$\nabla J(\theta) = \frac{1}{m} \mathbf{X}^T (\mathbf{X}\theta - \mathbf{y}) = 0$$



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

$\mathbf{X}\theta - \mathbf{y} = 0$ (즉, $\theta^* = \mathbf{X}^{-1}\mathbf{y}$)로 구하지 않은 이유는?

주어진 문제에서 \mathbf{X} 는 design matrix이고 차원이 대략 $m \times n$ (m 은 training example 개수, n 은 feature 차원)이 되므로 square 행렬이 아닐 수 있습니다.

역행렬은 square 행렬의 경우만 존재하므로 일반적인 $m \times n$ 행렬은 역행렬을 구할 수 없습니다.

반면에, $\mathbf{X}^T \mathbf{X}$ 나 $\mathbf{X} \mathbf{X}^T$ 형태로 만들면 square가 되어 역행렬을 구할 수 있습니다.

Examples: $m = 4$.

	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

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

$m \times (n+1)$

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

m -dimensional vector

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

Gradient Descent vs. Normal Equation

$$\nabla J(\theta) = \frac{1}{m} \mathbf{X}^T (\mathbf{X}\theta - \mathbf{y})$$

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

Gradient Descent	Normal Equation
Need to choose alpha <i>Bad</i>	No need to choose alpha <i>Good</i>
Needs many iterations <i>Bad</i>	No need to iterate <i>Good</i>
$O(mn^2)$ <i>Very Good!</i>	$O(n^3)$, need to calculate inverse of $\mathbf{X}^T \mathbf{X}$ <i>Very Bad!</i>
Works well when n is large	Slow if n is very large

$$\mathbf{X} = \begin{pmatrix} \text{---} n \text{---} \\ \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{pmatrix}_m$$

$$\mathbf{X}^T \mathbf{X} \rightarrow (n \times n)$$

$$(n \times m) \cdot (m \times n)$$

- With the normal equation, computing the inversion has complexity $O(n^3)$.
 - ★ So if we have a very large number of features (i.e., large n), the normal equation will be slow.
 - ★ In practice, **when n exceeds 10,000** it might be a good time **to go to an iterative process**.

3. Polynomial Regression

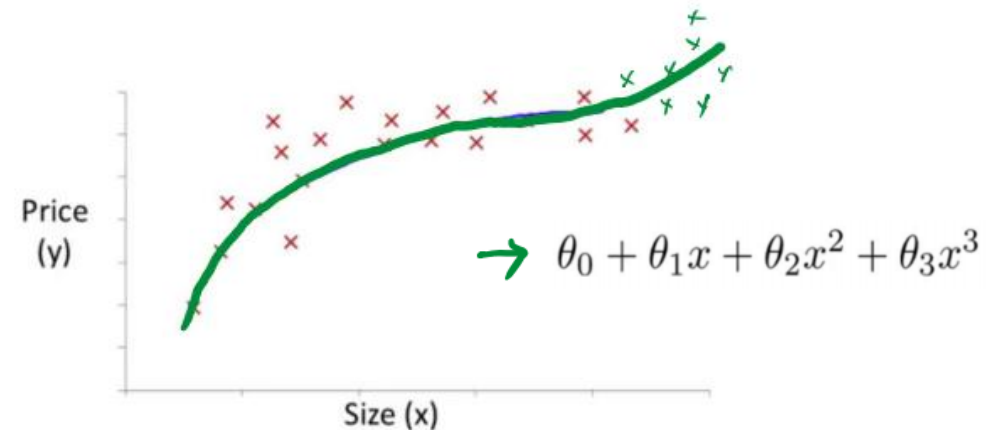
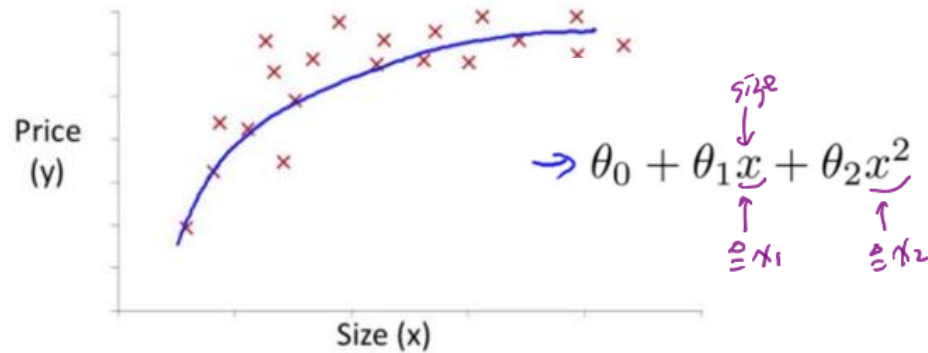
Polynomial Regression

◆ Polynomial Regression can be solved by Linear Regression

- Linear Regression 문제로 환원하여 풀 수 있음:
 $x_1 = (\text{size})$
 $x_2 = (\text{size})^2$
 $x_3 = (\text{size})^3$
- 이때 feature scaling이 중요해 짐

$$h_{\theta}(x) = \theta_0 + \theta_1 \underbrace{x}_{x_1} + \theta_2 \underbrace{x^2}_{x_2} + \theta_3 \underbrace{x^3}_{x_3}$$

● Ex) Housing Price Prediction



◆ Polynomial이외의 Nonlinear Function의 경우는?

- 마찬가지 방법(변수 치환)을 통해 선형회귀로 바꾸어 풀 수 있음!

$$h_{\theta}(x) = \theta_0 + \theta_1 \underbrace{\sqrt{x}}_{x_1} + \theta_2 \underbrace{\cos(x+\pi)}_{x_2}$$

4. For Better Results

- ✓ Feature Normalization
- ✓ Learning Rate Selection

Feature Normalization

◆ Feature Scaling (Range Normalization)

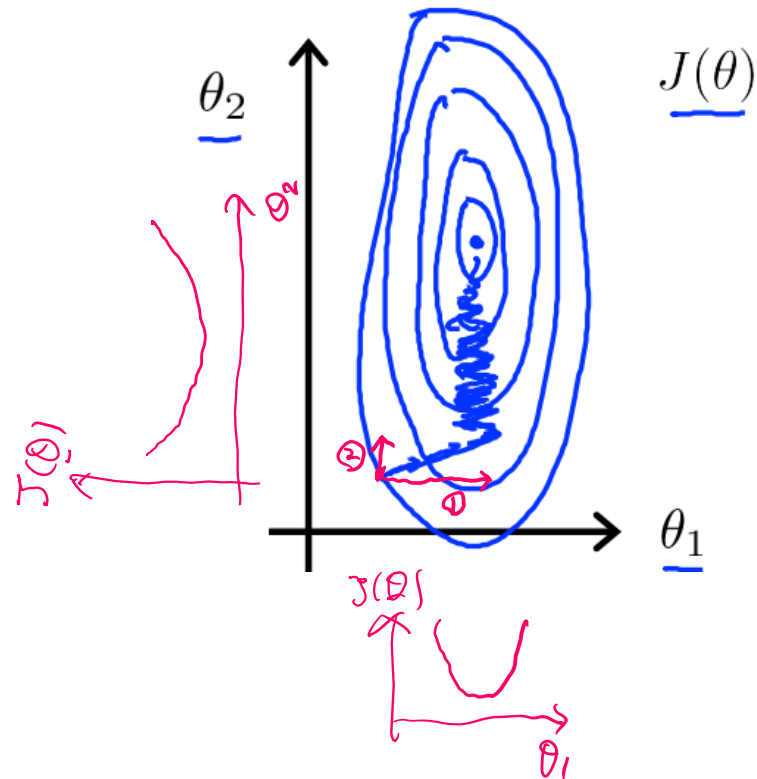
- Make sure features are on a similar scale

E.g. $x_1 = \text{size (0-2000 feet}^2\text{)}$ ←

$x_2 = \text{number of bedrooms (1-5)}$ ←

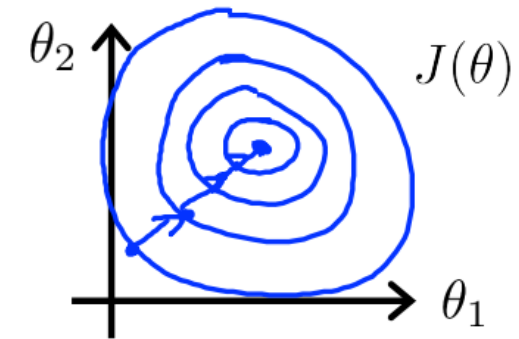


$$\begin{aligned} \rightarrow x_1 &= \frac{\text{size (feet}^2\text{)}}{2000} \\ \rightarrow x_2 &= \frac{\text{number of bedrooms}}{5} \end{aligned}$$



$$\begin{aligned} h_{\theta}(x) &= \theta_1 x_1 + \theta_2 x_2 \\ J(\theta) &= \frac{1}{2m} \sum_i (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ \textcircled{1} \frac{\partial J(\theta)}{\partial \theta_1} &= \frac{1}{m} \sum_i (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \textcircled{x_1^{(i)}} \\ \textcircled{2} \frac{\partial J(\theta)}{\partial \theta_2} &= \frac{1}{m} \sum_i (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \textcircled{x_2^{(i)}} \end{aligned}$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



∴ Converge faster!

Feature Normalization (cont'd)

◆ Mean Normalization (Mean Shifting)

- Replace x_j with $x_j - \mu_j$ to make features have approximately zero mean.

◆ Feature Normalization = Range Normalization + Mean Normalization

E.g. $\rightarrow x_1 = \frac{\text{size} - 1000}{2000}$

$x_2 = \frac{\#bedrooms - 2}{5}$

$\rightarrow -0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$

Feature Normalization (cont'd)

◆ Feature Normalization Summary

- Mean shifting:

$$x_j^{(i)} = x_j^{(i)} - \mu_j \quad (\text{i.e., } \mathbf{x}'^{(i)} = \mathbf{x}^{(i)} - \boldsymbol{\mu})$$

- With Scaling Method 1: Use **range** ($R = \max |x_j|$) to normalize

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\max |x_j|} \quad (\text{i.e., } \mathbf{x}'^{(i)} = \frac{\mathbf{x}^{(i)} - \boldsymbol{\mu}}{\max |\mathbf{x}^{(i)}|})$$

→ Normalized range: $-1 \leq x'_j \leq 1$ (strictly)

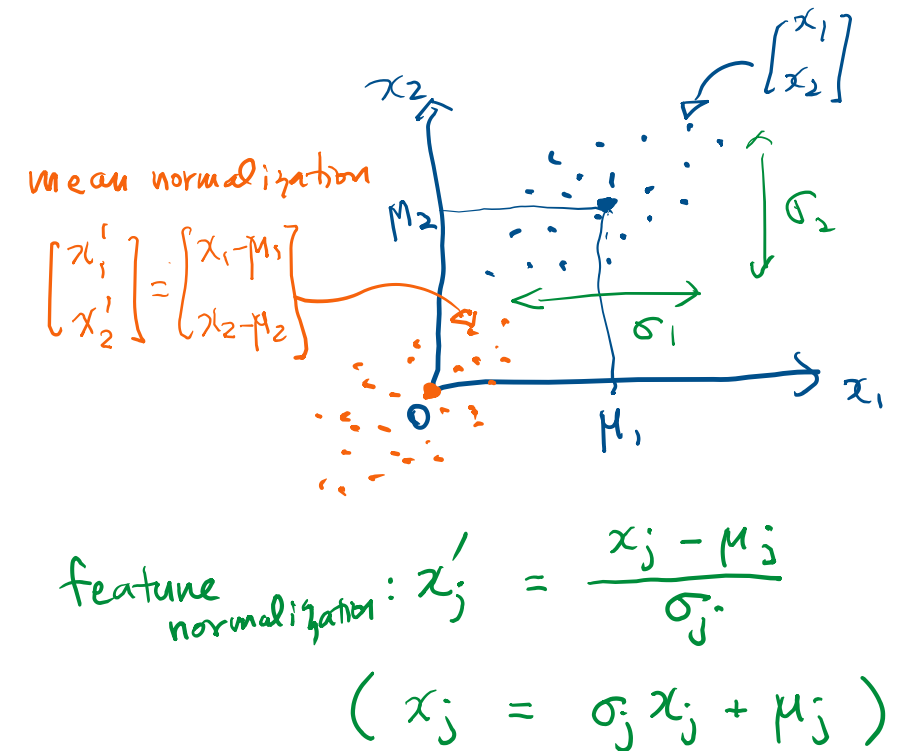
- With Scaling Method 2: Use **standard deviation** () to normalize

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad (\text{i.e., } \mathbf{x}'^{(i)} = \frac{\mathbf{x}^{(i)} - \boldsymbol{\mu}}{\boldsymbol{\sigma}})$$

→ Normalized range: $-1 \leq x'_j \leq 1$ (roughly)

◆ Caution:

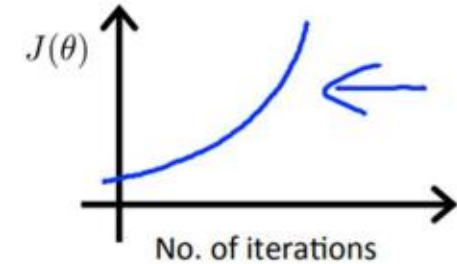
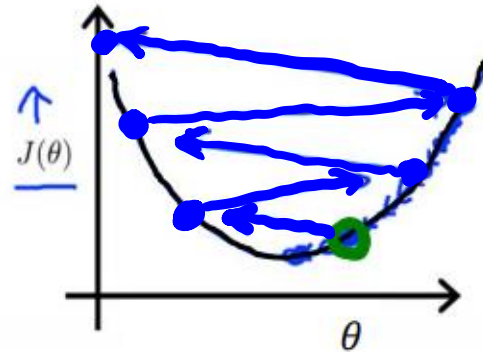
- Do not normalize $x_0 = 1$.
- There is **no** need to do feature normalization with the **normal equation**.



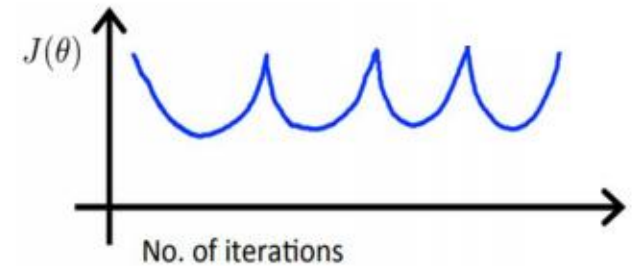
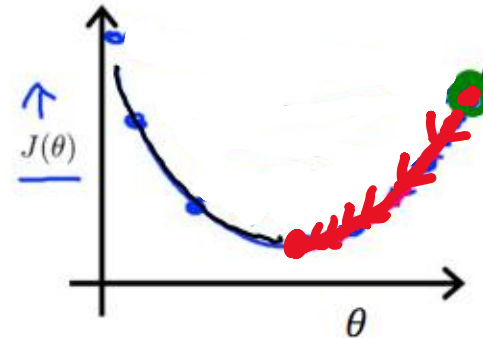
Learning Rate Selection

◆ How to Choose Learning Rate α ?

- Too large α : may not converge



- Too small α : slow convergence



- Rule of thumb:

★ Try ..., 0.001, ..., 0.01, ..., 0.1, ..., 1, ...
then try the in-betweens if not satisfactory

- Learning Rate Scheduling

★ Starts with some large α , and then decrease α according to a schedule

Q & A

본 강의 영상(자료)는 경희대학교 수업목적으로 제작·게시된 것이므로 수업목적 외 용도로 사용할 수 없으며, 무단으로 복제, 배포, 전송 또는 판매하는 행위를 금합니다. 이를 위반 시 민·형사상 법적 책임은 행위자 본인에게 있습니다.