

Spring 2023

SWCON253: Machine Learning

Lecture 04

Linear Regression

Jinwoo Choi
Assistant Professor
CSE, Kyung Hee University



Contents

1. Linear Regression
2. Normal Equation
3. Polynomial Regression
4. For Better Results

- 
- Learning Rate Tuning

References

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



1. Linear Regression

1. Multivariate Linear Regression
2. Data Representation
3. Linear Model Representation
4. Cost Function (& Gradient)
5. Parameter Update (by Gradient Descent)

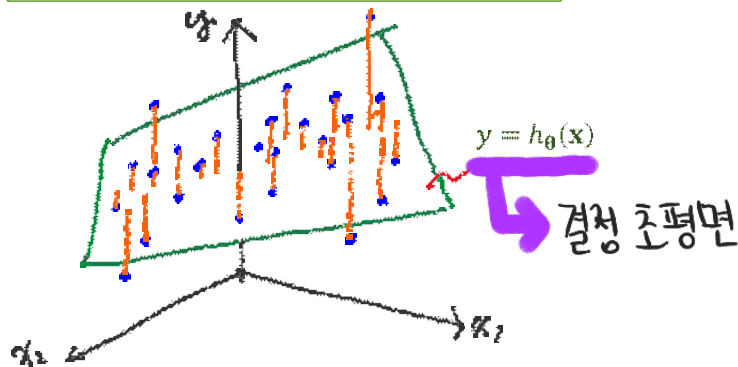


o. Multivariate Linear Regression

- ◆ 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 : y

Size (feet ²)	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 $n=4$ 4차원의 특징 벡터
- $x^{(i)}$ = input (features) of i^{th} training example. i 번째 훈련 샘플
- $x_j^{(i)}$ = value of feature j in i^{th} training example. // i 의 j 번째 특징 값



1. Data Representation

◆ Data with Multiple Features

● A feature vector: $\mathbf{x} = [x_1, \dots, x_n]^T$

● Training dataset: $\mathbb{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

$$\mathbf{X} = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^m & x_2^m & \dots & x_n^m \end{bmatrix}$$

◆ For a batch processing

design matrix
(input)

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



target vector
(label)

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

output vector
(prediction)

$$\hat{\mathbf{y}} = \mathbf{h}_{\theta}(\mathbf{X}) = \begin{bmatrix} h_{\theta}(\mathbf{x}^{(1)}) \\ h_{\theta}(\mathbf{x}^{(2)}) \\ \vdots \\ h_{\theta}(\mathbf{x}^{(m)}) \end{bmatrix}$$

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 $x_0 = 1$. 한 차원 추가 (상수항 θ_0 를 위해 추가)
- Then, for a single training example (i.e., $\mathbf{x}^{(i)}$):

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

$$= [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \underline{\boldsymbol{\theta}}^T \underline{\mathbf{x}}$$

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

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

같은 얘기
↔

◆ 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

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

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

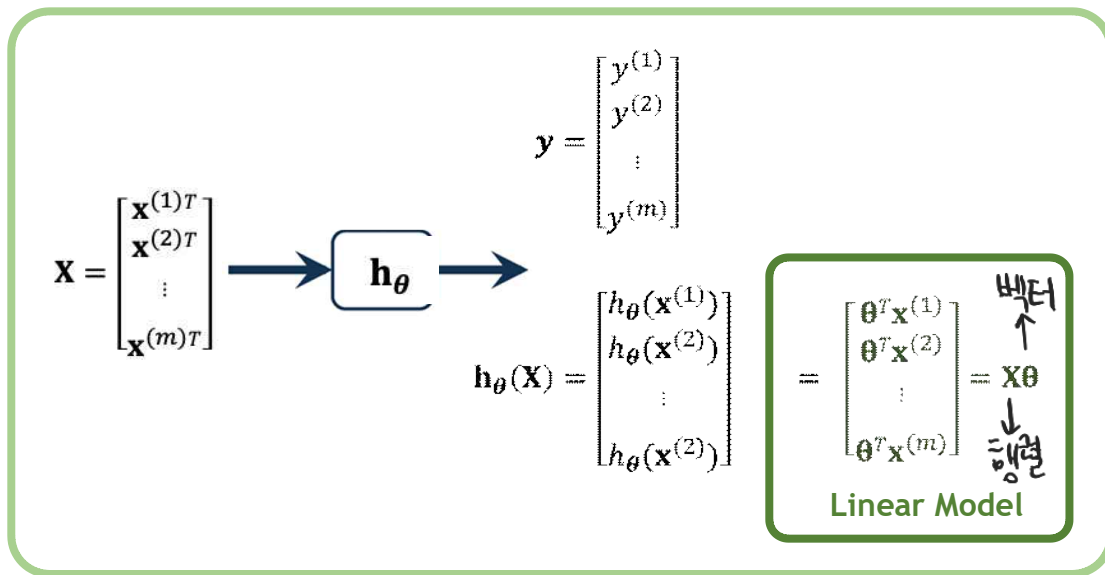
- and, for a batch of training example

$$\mathbf{h}_{\theta}(\mathbf{X}) = \mathbf{X} \boldsymbol{\theta} + \theta_0 \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

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

2. Linear Model Representation (cont'd)

◆ Summary of Input, Output, & Model



3. MSE Cost for Linear Model

◆ MSE Cost

- Classic form:

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

$\hat{y}^{(i)} = h_{\theta}(\mathbf{x}^{(i)})$ (Linear model)
 $\hat{y}^{(i)} - y^{(i)}$ is Error
 $\frac{\partial}{\partial \theta_j}$

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

- Vector form

$$J(\theta) = \frac{1}{2m} \|\mathbf{h}_{\theta}(\mathbf{X}) - \mathbf{y}\|_2^2$$

$\mathbf{h}_{\theta}(\mathbf{X})$ is Linear model
 $\mathbf{X}\theta$ is Linear model

$$= \frac{1}{2m} \|\mathbf{X}\theta - \mathbf{y}\|_2^2$$

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

cf) $\|\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{x}$ 이용

◆ Gradient of the MSE Cost

- Classic form:

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

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

- Vector form

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

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

유도과정:
다음 슬라이드



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

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

$$\|x\theta - y\|_2^2 = (x\theta - y)^T (x\theta - y) \stackrel{①}{=} (\theta^T x^T - y^T)(x\theta - y)$$

$$= \theta^T x^T x \theta - y^T x \theta - \underline{\theta^T x^T y} + y^T y$$

$$\stackrel{②}{=} \theta^T x^T x \theta - 2y^T x \theta + y^T y$$

Quadratic Form of θ

$$\Rightarrow \nabla_{\theta} \|x\theta - y\|_2^2 = \nabla_{\theta} (\theta^T x^T x \theta) - 2 \nabla (y^T x \theta)$$

$$\stackrel{③}{=} 2x^T x \theta - 2x^T y$$

$$= \underline{2x^T (x\theta - y)} //$$

$$\textcircled{1} \left\{ \begin{array}{l} (a-b)^T = a^T - b^T \\ (Ax)^T = x^T A^T \end{array} \right.$$

$$\textcircled{2} \theta^T x^T y = (x\theta)^T y = y^T x \theta$$

$$\textcircled{3} \left\{ \begin{array}{l} \nabla_x (x^T A^T A x) \xrightarrow{\text{symmetrize}} = 2A^T A x \\ \nabla_x (b^T A x) = \nabla_x e^T x = e \\ \quad \downarrow \\ \quad \text{a row vector} \\ \quad e = e^T \\ \quad \quad \quad = (A^T A)^T \\ \quad \quad \quad = A^T b \end{array} \right.$$



4. Parameter Update by GD (for Linear Model)

◆ Gradient Descent for Linear Regression with MSE Cost

- Classic form:

$$\theta = \theta + \Delta\theta = \theta - \alpha \nabla_{\theta}$$

Repeat until convergence {

Update $\forall \theta_j$'s simultaneously:

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

}

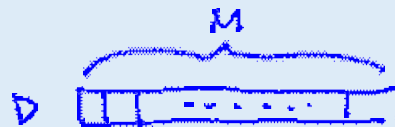
- Vector form

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

해석적인 방법

1. Normal Equation
2. Gradient Descent vs. Normal Equation



Normal Equations

◆ Analytic Solution to *Linear Regression with MSE Loss*

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

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

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

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

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

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

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$ m -dimensional vector

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

벡터

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

주어진 문제에서 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, $X^T X$ 나 XX^T 형태로 만들면 square가 되어 역행렬을 구할 수 있습니다.

일반적으로 역행렬이 없음



Gradient Descent vs. Normal Equation

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

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

→ 하이퍼 파라미터

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(m^2)$ 복잡도	$O(n^3)$, need to calculate inverse of $\mathbf{X}^T \mathbf{X}$
Works well when n is large	Slow if n is very large 특징 벡터의 차원이 커지면

$$\mathbf{X} = \begin{bmatrix} \text{---}^n \text{---} \\ \vdots \\ \text{---}^n \text{---} \end{bmatrix}$$

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

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

- 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

공식



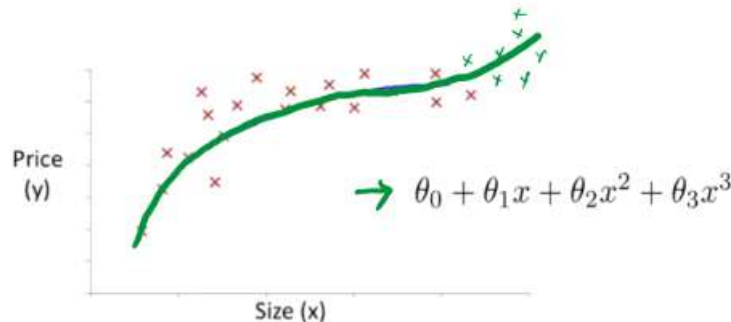
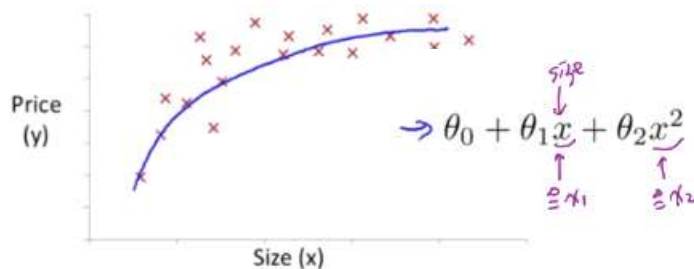
Polynomial Regression

◆ 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

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

- Ex) Housing Price Prediction
- 변수마다 가질 수 있는 범위가 상이하게 달라짐



◆ Polynomial이외의 Nonlinear Function의 경우는?

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

4. For Better Results

1. Feature Normalization
2. Learning Rate



Feature Normalization

◆ Feature Scaling (Range Normalization)

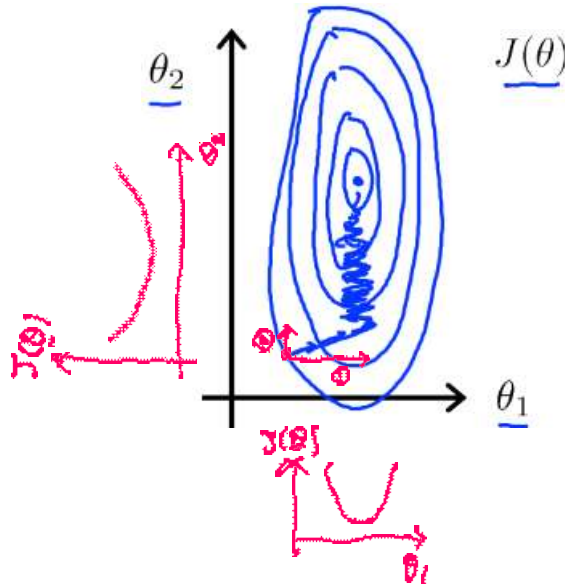
- Make sure features are on a similar scale

E.g. $x_1 = \text{size (0-2000 feet}^2\text{)}$ ← x_1 이 x_2 보다 10배에 큰 영향을 미침
 $x_2 = \text{number of bedrooms (1-5)}$ ←

Scaling

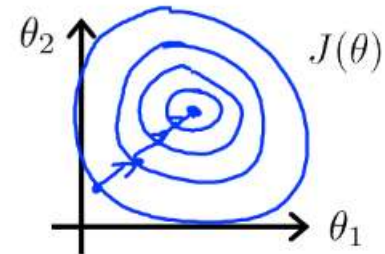
$$\begin{aligned} x_1 &= \frac{\text{size (feet}^2\text{)}}{2000} \\ x_2 &= \frac{\text{number of bedrooms}}{5} \end{aligned}$$

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



$$\begin{aligned} h_0(x) &= \theta_0 x_0 + \theta_1 x_1 \\ J(\theta) &= \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \\ \frac{\partial J(\theta)}{\partial \theta_1} &= \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \frac{\partial J(\theta)}{\partial \theta_2} &= \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_2^{(i)} \end{aligned}$$

→ 기울기 큰 곳



∴ Converge faster!

Feature Normalization (cont'd)

◆ Feature Scaling (cont'd)

- Eg., get every feature into approximately a $-1 \leq x_j \leq 1$ range.

$$\begin{array}{l|l} 0 \leq x_1 \leq 3 \quad \checkmark & -3 \text{ to } 3 \quad \checkmark \\ -2 \leq x_2 \leq 0.5 \quad \checkmark & -\frac{1}{2} \text{ to } \frac{1}{2} \quad \checkmark \\ -100 \leq x_3 \leq 100 \quad \times & \\ -0.0001 \leq x_4 \leq 0.0001 \quad \times & \end{array}$$

Feature Normalization (cont'd)

◆ Mean Normalization

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

E.g. →

$$x_1 = \frac{\text{size} - 1000}{2000}$$

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

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

→ 평균

모든 데이터에서 평균을 빼면 → 평균 0 됨

1.5 bed

◆ Caution:

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

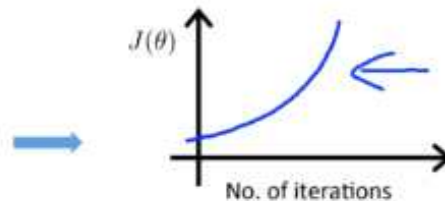
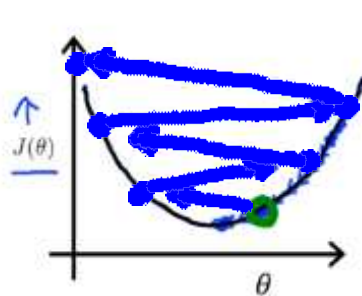


Learning Rate (α)

◆ How to Choose Learning Rate α ?

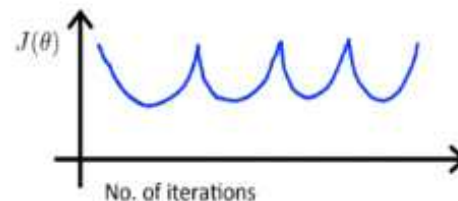
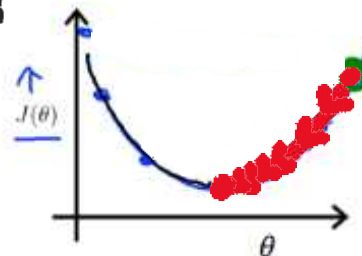
- Too large α : may not converge

바사
큰
→



- Too small α : slow convergence 오래걸림

→



- Rule of thumb:

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

$1e-1 \sim 1e-5$ or -6

- Learning Rate Scheduling

처음에 크게 했다가 점점 줄이는 방법

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

