

2주차 강의 – Multiple features

박소정

TNT - ML팀



Contents



Multiple features

다중 선형회귀



Vectorization

벡터화



실습

Multiple features

-다중 선형회귀 -

Contents

✓ 2주차: 여러 입력 변수를 사용한 회귀 분석

▶ 동영상 49 분 남음 ④ 3개의 평가가 남음

이번 주에는 선형 회귀를 확장하여 여러 입력 특징을 처리하는 방법을 알아봅니다. 또한 벡터화, 특징 스케일링, 특징 엔지니어링, 다항식 회귀 등 모델의 훈련과 성능을 개선하기 위한 몇 가지 방법도 배웁니다. 마지막 주에는 코드에서 선형 회귀를 구현하는 연습을 하게 됩니다.

✓ 학습 목표 표시

다중 선형 회귀 - multiple features (variables)

- 여러 입력 변수를 사용한 회귀 분석
- ex) 주택 가격 예측 > 침대가 하나 추가 될때마다 4만달러

Model:

Previously: $f_{w,b}(x) = wx + b$

example

$$f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$
$$f_{w,b}(x) = 0.1 \underset{\substack{\uparrow \\ \text{size}}}{x_1} + 4 \underset{\substack{\uparrow \\ \text{\# bedrooms}}}{x_2} + 10 \underset{\substack{\uparrow \\ \text{\# floors}}}{x_3} + -2 \underset{\substack{\uparrow \\ \text{years}}}{x_4} + 80 \underset{\substack{\uparrow \\ \text{base price}}}{b}$$

단순(일차) 선형회귀 vs 다중 선형 회귀 vs 다항식 회귀

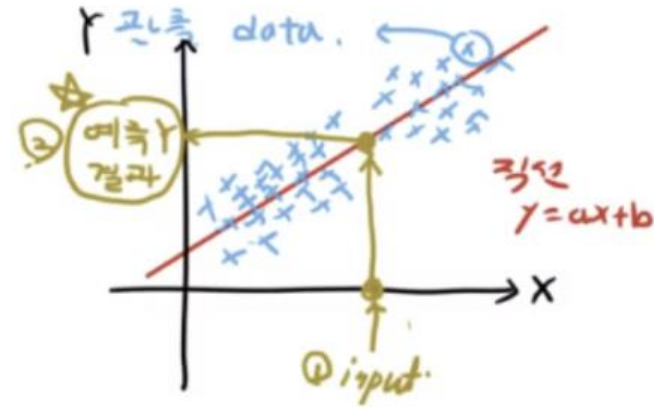
Simple + 직선 = 원인 1개
(Simple Linear Regression)



$$y = ax + b$$

종속변수 (Target, 결과) 독립변수 (Feature, 원인)

기울기 y절편



Multiple + 직선 = 원인 n개
(Multiple Linear Regression)

$$y = a_0 + a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_nx_n$$

종속변수 (Target) 독립변수 (Feature, 원인 n개)

y절편 기울기

비선형??? 원인변수의 차수가 1차가 아닌 항이 있으면...!!

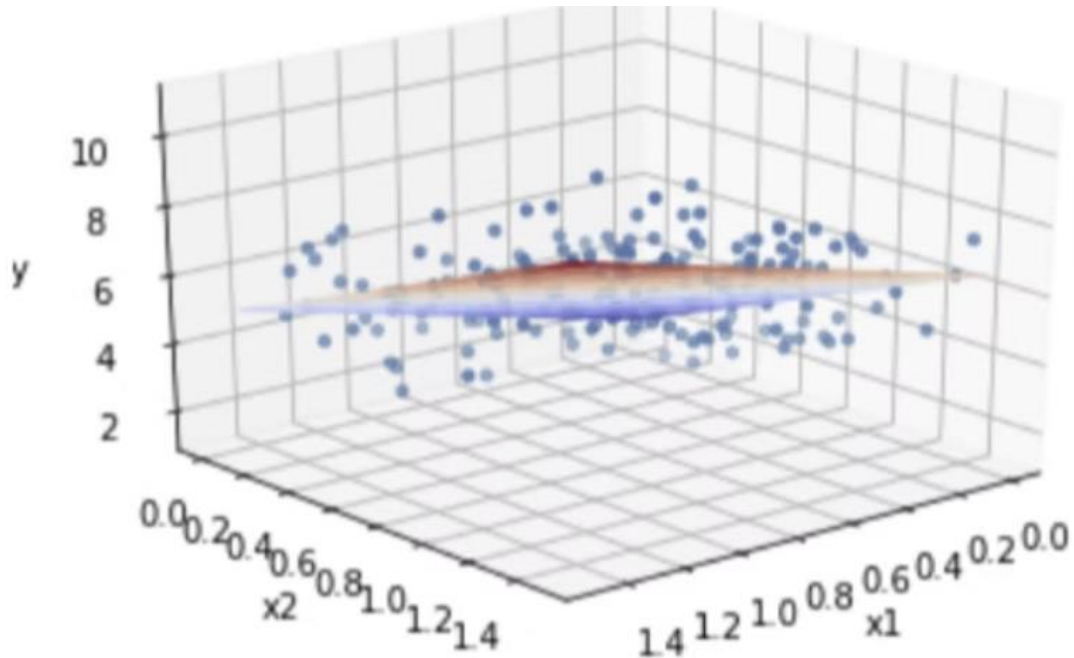
<비선형??>

$$y = a_0 + a_1x_1^2 + a_2x_2^3 + a_3x_3 + \dots + a_nx_n^k$$

★ 차수가 1차가 아닌 경우.

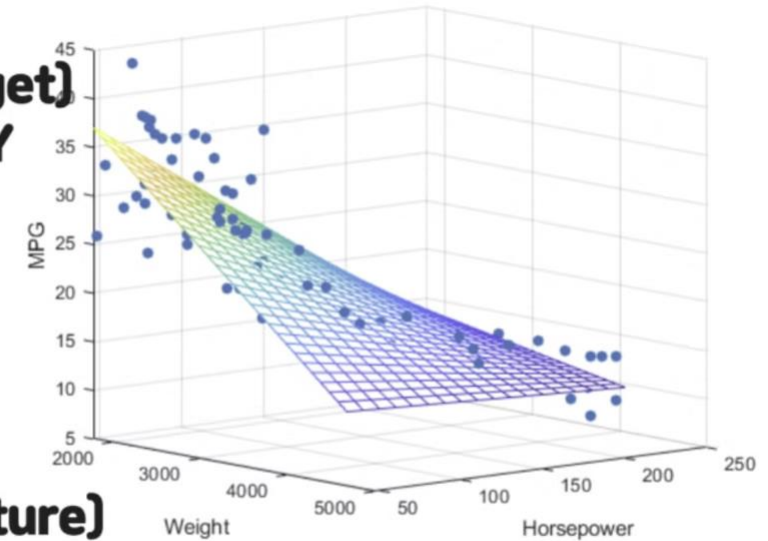
-> 독립변수의 개수, 차수 차이

다중 선형 회귀 vs 다항식 회귀



**결과(Target)
종속변수 Y**

**원인(Feature)
독립변수 X**



[https://medium.com/analytics-vidhya/](https://medium.com/analytics-vidhya/new-aspects-to-consider-while-moving-from-simple-linear-regression-to-multiple-linear-regression-dad06b3449ff)

new-aspects-to-consider-while-moving-from-simple-linear-regression-to-multiple-linear-regression-dad06b3449ff

- 독립변수가 2개일때 예시 (평면 vs 곡면)
- 다중 선형 회귀는 모든 독립 변수에 대해서 선형이다. 즉 일차 항 계수를 갖는다.

다중 선형 회귀 vs 다항식 회귀 예시 설명

대학생 스트레스 수준 $= a_1 \times \text{과제양} + a_2 \times \text{시험준비시간} +$
 $a_3 \times \text{수면시간} + a_4 \times \text{진로불확실성} + a_5 \times \text{경제적부담} + a_6 \times \text{인간관계}$

- 대학생의 스트레스 수준이 시험 준비 시간이 감소함에 따라 선형적으로 증가
- Vs
- 점점 더 빠르게 증가

변수의 중요도를 파악할때 - 회귀 계수 값을 비교x

회귀 계수는 각 독립 변수와 종속 변수와의 관계를 나타낸다.

결과적으로 계수가 큰 변수가 보다 큰 변화를 나타내므로 더 중요하다고 생각하기 쉽지만,

각 변수는 단위가 다르므로 직접적인 비교는 불가능하다.

- >

Model:

표준화 진행 (feature scaling)

Previously: $f_{w,b}(x) = wx + b$

example

$$f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$
$$f_{w,b}(x) = 0.1x_1 + 4x_2 + 10x_3 + -2x_4 + 80$$

↑ ↑ ↑ ↑ ↑
size #bedrooms #floors years base price

대학생 스트레스 수준 = $a_1 \times \text{과제양} + a_2 \times \text{시험준비시간} +$

$a_3 \times \text{수면시간} + a_4 \times \text{진로불확실성} + a_5 \times \text{경제적부담} + a_6 \times \text{인간관계}$

feature scaling
(1) mean - normalization

평균 정규화

$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_{1,scaled} = \frac{x_1}{2000}$$

max

$$x_{2,scaled} = \frac{x_2}{5}$$

max

$$0.15 \leq x_{1,scaled} \leq 1$$

$$0 \leq x_{2,scaled} \leq 1$$

feature scaling

(2) Z-score normalization

$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

$$-0.67 \leq x_1 \leq 3.1 \quad -1.6 \leq x_2 \leq 1.9$$

다중선형회귀분석 시 유의할 점

- 독립변수(설명변수)들끼리의 상관관계가 높으면 다중공선성 문제(**multicollinearity**)
- 다중 공선성이란, 사용해야 될 독립변수들끼리 서로 밀접한 상관관계가 있어서 다중선형회귀모델에서 각각의 요인들의 효과를 파악하기 어려워지는 것을 말합니다.

다중 공선성의 해결책

- 제거 또는 새로운 변수를 만들어서 해결한다.

Feature engineering

Feature engineering

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} \underbrace{x_1}_{\text{depth}} + w_2 \underbrace{x_2}_{\text{depth}} + b$$

$$\text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

new feature

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} x_1 + \underbrace{w_2}_{\text{depth}} x_2 + \underbrace{w_3}_{\text{area}} x_3 + b$$



벡터화

벡터화

vector $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$

$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b =$

Parameters and features

$$\vec{w} = [w_1 \ w_2 \ w_3] \quad n=3$$

b is a number

$$\vec{x} = [x_1 \ x_2 \ x_3]$$

linear algebra: count from 1

NumPy 

$w[0] \ w[1] \ w[2]$

```
w = np.array([1.0, 2.5, -3.3])
```

```
b = 4
```

$x[0] \ x[1] \ x[2]$

```
x = np.array([10, 20, 30])
```

code: count from 0

Without vectorization $n=100,000$

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + w_3x_3 + b$$

```
f = w[0] * x[0] +  
     w[1] * x[1] +  
     w[2] * x[2] + b
```



Without vectorization

$$f_{\vec{w},b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \rightarrow \begin{matrix} j=1 \dots n \\ 1, 2, 3 \end{matrix}$$

$\text{range}(0, n) \rightarrow j = 0 \dots n-1$

```
f = 0  
for j in range(0, n):  
    f = f + w[j] * x[j]  
f = f + b
```



Vectorization

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

```
f = np.dot(w, x) + b
```



-> 벡터화된 코드 작성시 더 빠르게 계산 가능함

Without vectorization

```
for j in range(0,16):  
    f = f + w[j] * x[j]
```

t_0
 $f + w[0] * x[0]$

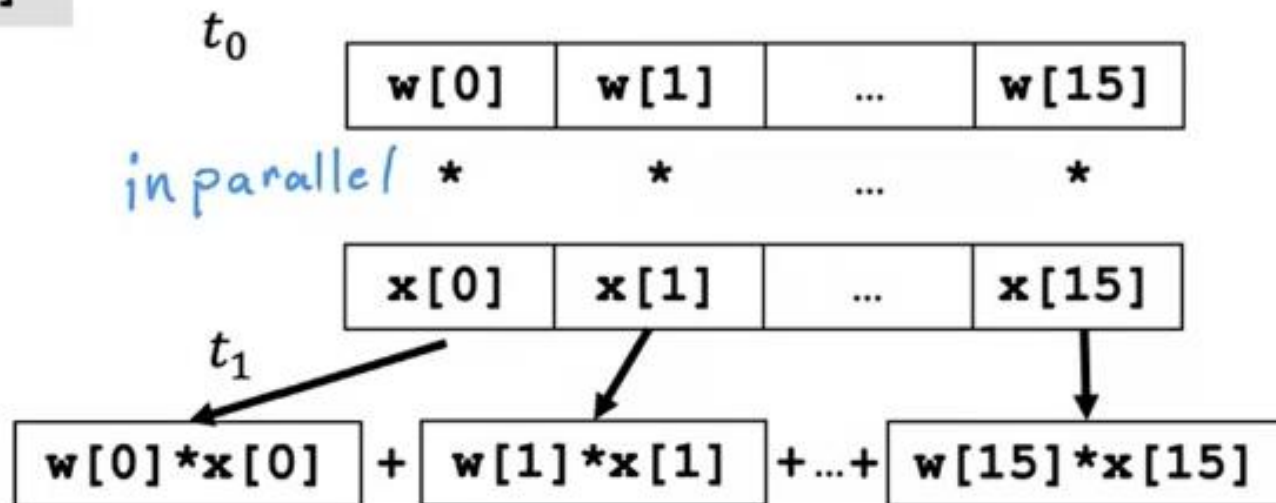
t_1
 $f + w[1] * x[1]$

...

t_{15}
 $f + w[15] * x[15]$

Vectorization

```
np.dot(w,x)
```



직관적으로 이게 왜 컴퓨터에서 계산이 훨 빠른지?

각 쌍을 동시에 병렬로 곱한다. > 차례로 별개의 덧셈을 수행하지 않아도 되서 시간이 훨빠름
대규모 학습 많기 때문에 필수적임

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \epsilon_i \quad (1)$$

여기서 아래첨자 y_i 는 i 번째 관측치를 의미하고, ϵ_i 는 이때의 오차항을 나타낸다. 우리가 추정하고자하는 값, 즉 회귀계수는 β_j ($0 \leq j \leq k$)이고, 독립변수 x_{ij} 는 known value이다. 식 (1)을 N 개의 샘플에 대하여 확장한 후, vector-matrix 형태로 표기하면 다음과 같다.

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \cdots & x_{Nk} \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}}_{\boldsymbol{\epsilon}} \quad (2)$$

여기서 $\mathbf{e} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_N)$.

- 경사 하강법에서 수렴과 학습률

learning rate

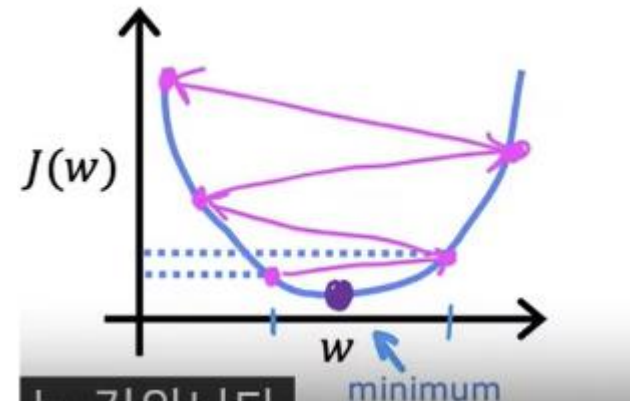
$$\underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$\underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

INL

- $w = w -$ (여기서 $=$ 은 같다가 아니라, 할당한다는 뜻. 업데이트한다는 뜻 !)

선형회귀 일때 알파는 학습률 (항상 양수)
 너무 작으면 느리고
 너무 크면 M/m 찾지 못하고,
 즉 수렴하지 못하고, 갈라질 수 있다.



적합한 학습률 찾는법

먼저 작은 숫자로 시도하고 3배, 10배 등으로 조금씩 조정한다.

- 0.001의 학습률을 시도하고
0.01과 0.1 등과 같이 10배 더 큰 학습률을 시도

learning rate

$$\underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b)$$
$$\underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

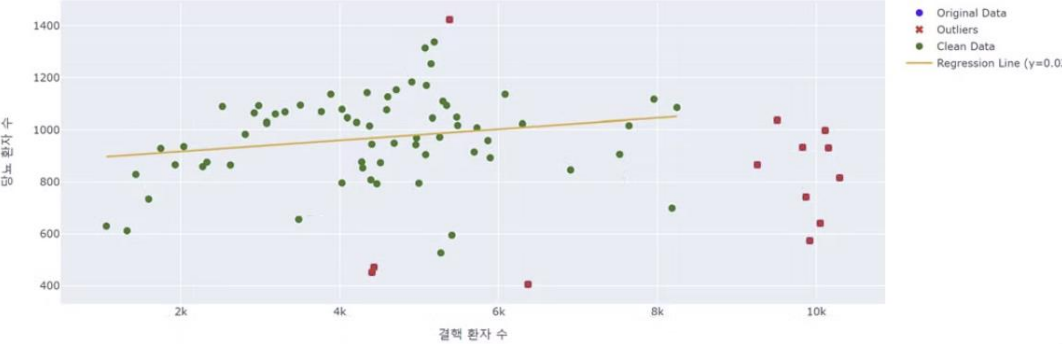
실습

- MSE (Mean Squared Error) = $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- MAE (Mean absolute error) = $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
- RMSE (Root Mean Squared Error) = \sqrt{MSE}
- R-squared (Coefficient of determination) = $1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SSE}{SST} = \frac{SSR}{SST}$

```
((y - y.mean())**2).sum()
```

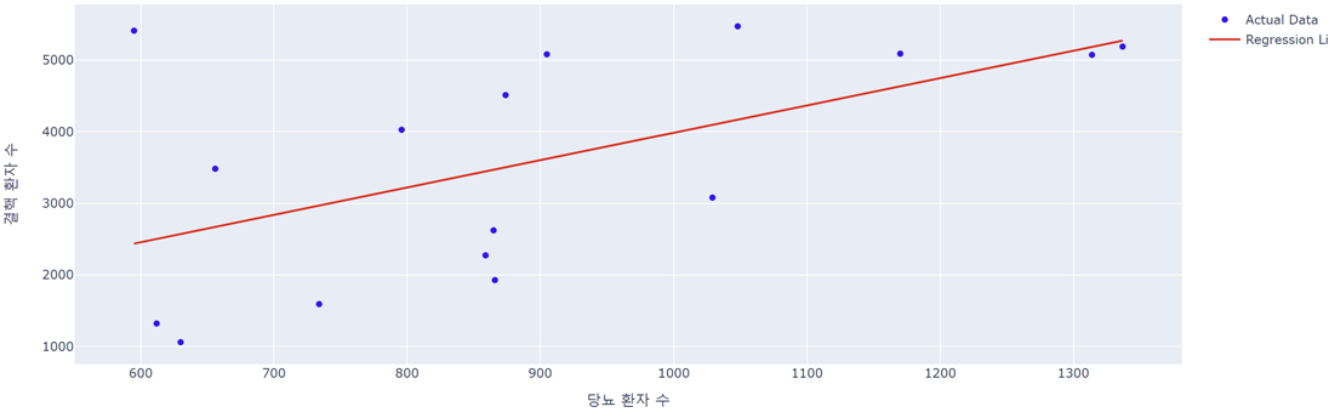
당뇨환자 수와 결핵환자 수 회귀 분석 실습

단순 선형 회귀 분석: 결핵 환자 수에 따른 당뇨 환자 수 (이상치 제거)



30~39세 연령대: 당뇨 환자 수에 따른 결핵 환자 수

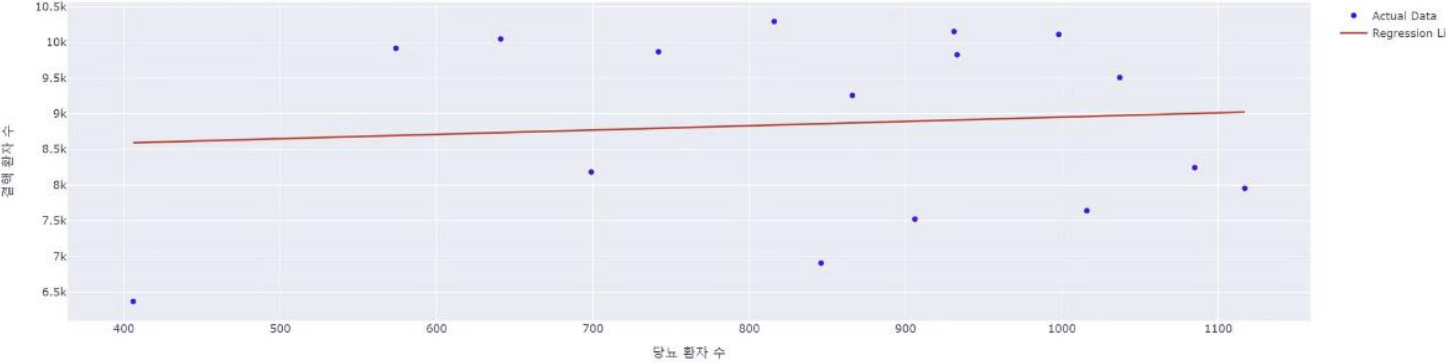
30~39세 연령대: 당뇨 환자 수에 따른 결핵 환자 수



ML

AI ML DL TF GS TP

70세 이상 연령대: 당뇨 환자 수에 따른 결핵 환자 수



뉴욕 맨해튼 의 주택 임대료 데이터 실습

manhattan.csv

git clone <https://github.com/Codecademy/datasets.git> head
[/home/ubuntu/workspace/datasets/streeteasy/manhattan.csv](https://github.com/Codecademy/datasets.git)

정답

<https://velog.io/@hyesoup/다준선형회귀Multiple-Linear-Regression-예제>
<https://recipesds.tistory.com/entry/맨하탄집값-분석-다중회귀-예측까지-실습>