



# Supervised Learning & Linear Regression

---

Deep learning



Woohwan Jung

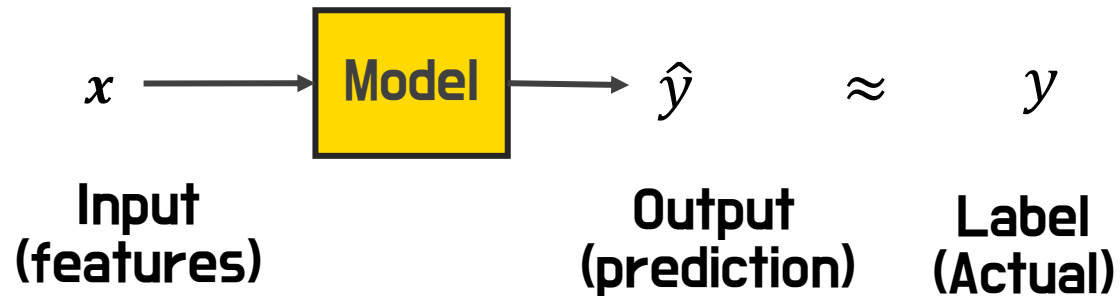
# Topics

- Supervised learning
- Linear regression
  - Model
  - Cost/Loss function
  - Optimization

# Supervised learning

# Supervised learning

Goal: generalize the input-output relationship



Training?

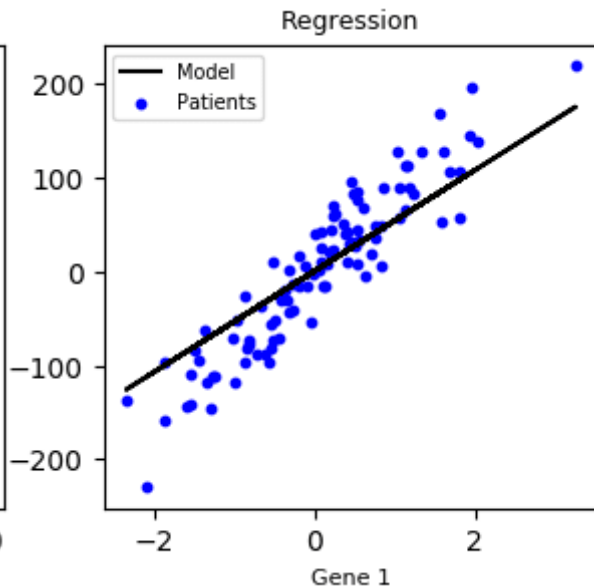
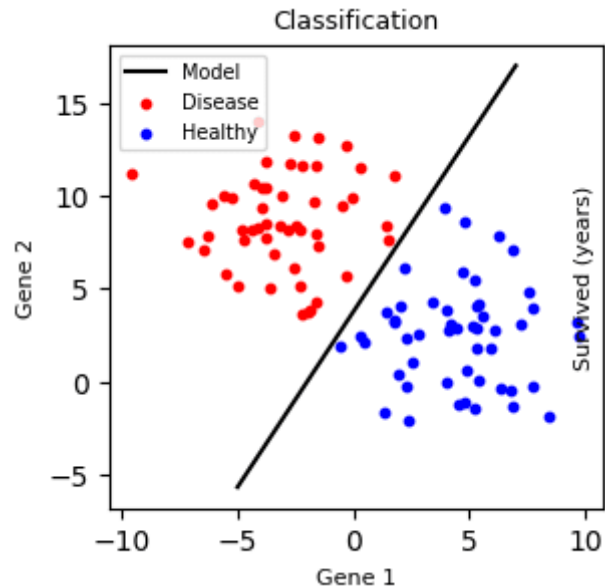
**Building a model** to make the model can **predict the labels** by using train data

## Dataset

Features				Label
Size	Beds	Baths	Zip	Price
1100	1	1	64576	1.29
1900	3	1.5	78321	2.14
2800	3	3	98712	3.10
3400	4	3.5	25721	3.75

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})\}$$

# Classification vs Regression



Output      Categorical value  
(Class)

Numeric value

Q1. Classification? Regression?



Q2. Classification? Regression?



---> Cat

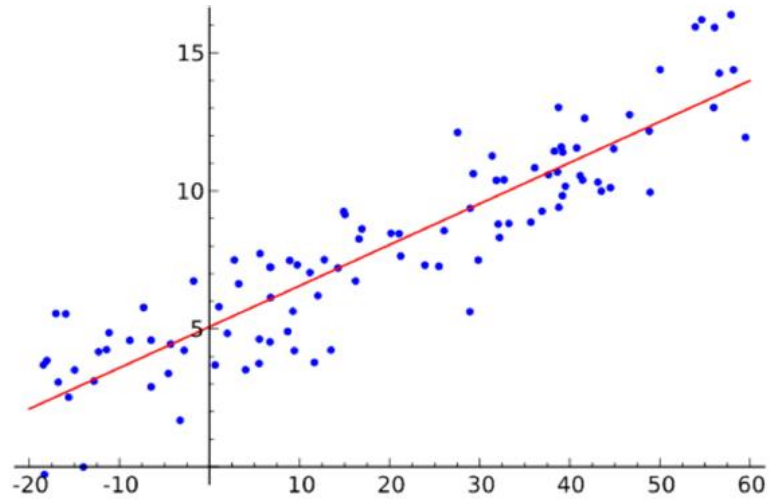


---> Dog

# Linear Regression

---

- Modelling the **linear** relationship between a scalar response (label) and one or more explanatory variables (features)



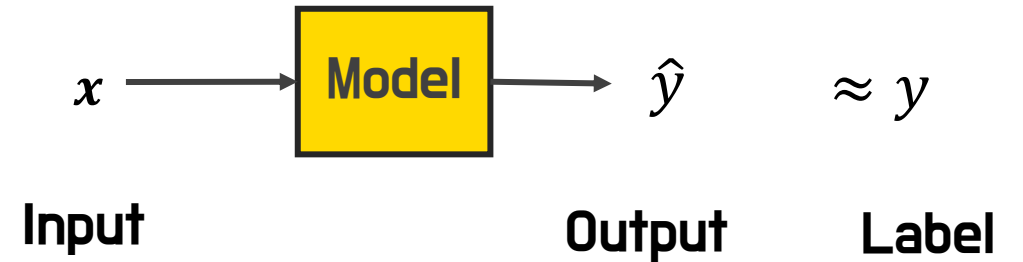
Examples)

Area of house -> house price

# of iPhones sold -> Apple's sales

# Linear Regression

- Data  $D = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots (\mathbf{x}^{(m)}, y^{(m)})\}$
- Model
  - Input:  $\mathbf{x}^{(i)} \in \mathbb{R}^d$
  - Output  $\hat{y}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$ 
    - Parameters:  $\mathbf{w} \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$



## Training a linear regression model?

Finding the model parameters  $\mathbf{w}$  and  $b$  which make  $\hat{y} \approx y$

## Measuring the distance between $\hat{y}$ and $y$

Squared error:  $(\hat{y} - y)^2$

Loss function  $L(\hat{y}^{(i)}, y^{(i)}) = (y^{(i)} - \hat{y}^{(i)})^2$

Cost function  $J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$

# Training a linear regression model

- **Given**

- Training data  $D = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

- **Our goal**

- Find  $\mathbf{w}, b$  that minimizes  $J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$

Q. How?

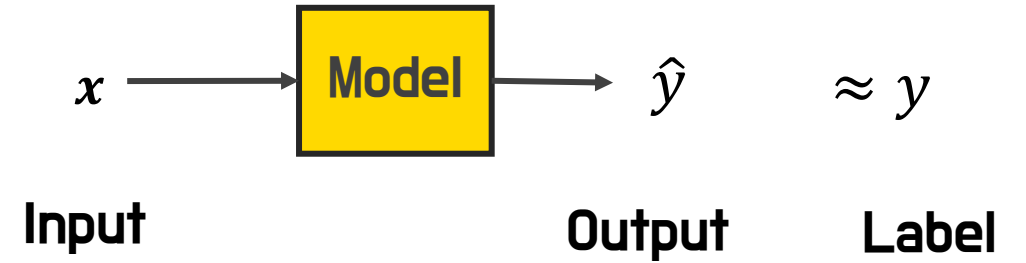
Applicable methods: gradient descent, linear least squares, ...

We are going to use gradient descent!



# **Gradient Descent :Linear Regression**

# Optimization



- Linear regression model

- $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$

- Loss function  $L(\hat{y}^{(i)}, y^{(i)}) = (y^{(i)} - \hat{y}^{(i)})^2$

- Cost function

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

- Our goal

- Find parameters  $\mathbf{w} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$  that minimize  $J(\mathbf{w}, b)$

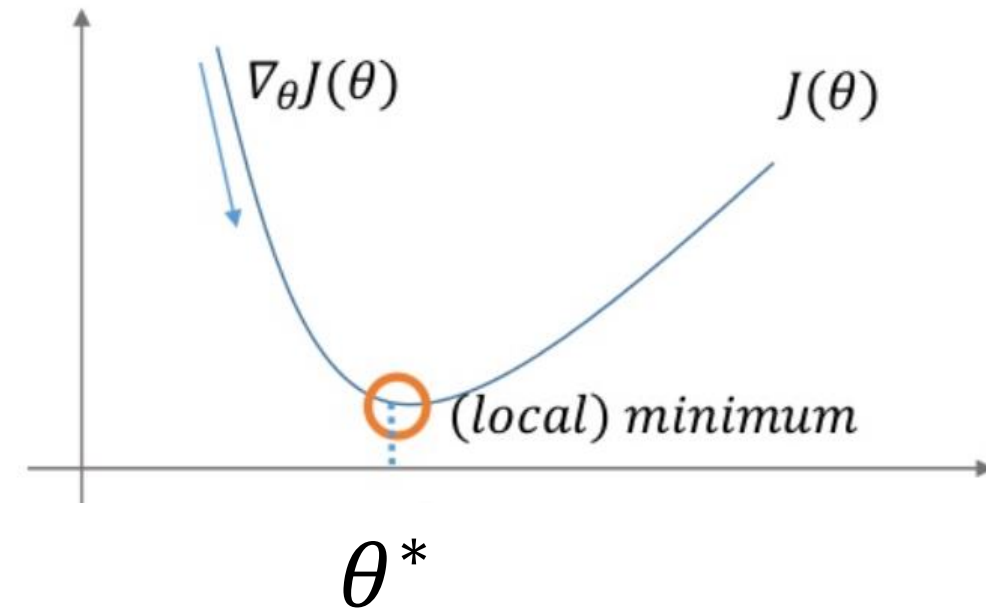
- Gradient Descent!

# Gradient Descent

- A popular optimization algorithm to minimize a cost function  $J(\theta)$ 
  - $J(\theta)$ : cost function
  - $\theta$ : model parameters
- Iteratively update model parameters to find the values that result in the lowest  $J(\theta)$

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

- $\eta$ : Learning rate



$$\nabla_{\theta} J(\theta) = \left[ \frac{\partial J(\theta)}{\partial \theta_1} \quad \frac{\partial J(\theta)}{\partial \theta_2} \quad \dots \quad \frac{\partial J(\theta)}{\partial \theta_k} \right]^T \text{ where } \theta = [\theta_1 \quad \theta_2 \quad \dots \quad \theta_k]^T$$

# Gradient of L

- Output  $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$
- Loss function  $L(\hat{y}, y) = (y - \hat{y})^2$

$$\frac{\partial L}{\partial w_k} = -2(y - \hat{y})x_k$$

$$\frac{\partial L}{\partial b} = -2(y - \hat{y})$$

# Gradient descent on m examples

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Gradient descent  
 $\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \cdot \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

$$\frac{\partial}{\partial w_k} J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_k} L(\hat{y}^{(i)}, y^{(i)}) = \frac{2}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y) x_k$$

$$\frac{\partial}{\partial b} J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b} L(\hat{y}^{(i)}, y^{(i)}) = \frac{2}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y)$$

# Gradient descent for training a Linear Regression model (n = 2)

- Randomly Initialize  $w, b$
- $lr = 0.1$
- For  $e = 1$  to  $n_{epoch}$ :
  - $d\_w1 = 0; d\_w2 = 0; d\_b = 0$
  - For  $i = 1$  to  $m$ :
    - $a = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$
    - $d\_w1 += 2(a - y)x_1^{(i)}$
    - $d\_w2 += 2(a - y)x_2^{(i)}$
    - $d\_b += 2(a - y)$
  - $w_1 -= lr * d\_w1 / m$
  - $w_2 -= lr * d\_w2 / m$
  - $b -= lr * d\_b / m$

$$\text{Gradient descent}$$
$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \cdot \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

$$\frac{\partial}{\partial w_k} J(\mathbf{w}, b) = \frac{2}{m} \sum_{i=1}^m (a - y) x_k$$

$$\frac{\partial}{\partial b} J(\mathbf{w}, b) = \frac{2}{m} \sum_{i=1}^m (a - y)$$